

R 语言教程

李东风

2019-08-29

目录

前言	15
I 介绍与入门	17
1 R 语言介绍	19
1.1 R 的历史和特点	19
1.2 R 的下载与安装	21
1.3 基本 R 软件的使用	24
1.4 RStudio 软件	25
1.5 练习	31
2 R 语言入门运行样例	33
2.1 命令行界面	33
2.2 四则运算	33
2.3 数学函数	35
2.4 输出	37
2.5 向量计算与变量赋值	38
2.6 工作空间介绍	40
2.7 绘图示例	41
2.8 汇总统计示例	45
2.9 运行源程序文件	47
2.10 附录：数据	49

II	R 的数据类型与相应运算	51
3	常量与变量	53
3.1	常量	53
3.2	变量	54
3.3	R 数据类型	54
4	数值型向量及其运算	55
4.1	数值型向量	55
4.2	向量运算	56
4.3	向量函数	59
4.4	复数向量	62
4.5	练习	62
5	逻辑型向量及其运算	63
5.1	逻辑型向量与比较运算	63
5.2	逻辑运算	65
5.3	逻辑运算函数	66
6	字符型数据及其处理	69
6.1	字符型向量	69
6.2	paste() 函数	69
6.3	转换大小写	70
6.4	字符串长度	70
6.5	取子串	70
6.6	类型转换	71
6.7	字符串拆分	72
6.8	字符串替换功能	72
6.9	正则表达式	72
7	R 向量下标和子集	75
7.1	正整数下标	75
7.2	负整数下标	76
7.3	空下标与零下标	76
7.4	下标超界	77

目录	5
7.5 逻辑下标	78
7.6 <code>which()</code> 、 <code>which.min()</code> 、 <code>which.max()</code> 函数	79
7.7 元素名	79
7.8 用 R 向量下标作映射	81
7.9 集合运算	82
7.10 练习	83
8 R 数据类型的性质	85
8.1 存储模式与基本类型	85
8.2 类属	87
8.3 类型转换	88
8.4 属性	89
8.5 <code>str()</code> 函数	91
8.6 关于赋值	91
9 R 日期时间	93
9.1 R 日期和日期时间类型	93
9.2 从字符串生成日期数据	94
9.3 日期显示格式	96
9.4 访问日期时间的组成值	97
9.5 日期舍入计算	98
9.6 日期计算	99
9.7 基本 R 软件的日期功能	104
9.8 练习	108
10 R 因子类型	111
10.1 因子	111
10.2 <code>table()</code> 函数	113
10.3 <code>tapply()</code> 函数	114
10.4 <code>forcats</code> 包的因子函数	114
10.5 练习	117
11 R 矩阵和数组	119
11.1 R 矩阵	119
11.2 矩阵子集	120

11.3	<code>cbind()</code> 和 <code>rbind()</code> 函数	123
11.4	矩阵运算	124
11.5	逆矩阵与线性方程组求解	128
11.6	<code>apply()</code> 函数	129
11.7	多维数组	130
12	数据框	133
12.1	数据框	133
12.2	数据框内容访问	134
12.3	数据框与矩阵的区别	138
12.4	<code>gl()</code> 函数	138
12.5	<code>tibble</code> 类型	140
12.6	练习	143
13	列表类型	145
13.1	列表	145
13.2	列表元素访问	146
13.3	列表类型转换	149
13.4	返回列表的函数示例- <code>strsplit()</code>	150
14	工作空间	153
III	R 编程	155
15	R 输入输出	157
15.1	输入输出的简单方法	157
15.2	读取 CSV 文件	160
15.3	Excel 表访问	168
15.4	使用专用接口访问数据库	172
15.5	文件访问	175
15.6	中文编码问题	178
15.7	目录和文件管理	185
16	程序控制结构	187
16.1	表达式	187

目录	7
16.2 分支结构	187
16.3 循环结构	189
16.4 R 中判断条件	192
16.5 管道控制	192
17 函数	193
17.1 函数基础	193
17.2 变量作用域	200
17.3 函数进阶	206
17.4 程序调试	216
17.5 函数式编程介绍	221
18 R 程序效率	227
18.1 R 的运行效率	227
18.2 向量化编程	228
18.3 减少显式循环	233
18.4 R 的计算函数	242
18.5 并行计算	248
19 随机模拟	257
19.1 随机数	257
19.2 <code>sample()</code> 函数	258
19.3 随机模拟示例	259
IV 用 R 制作研究报告和图书	265
20 用 R 制作研究报告	267
21 Markdown 格式	271
21.1 介绍	271
21.2 Markdown 格式文件的应用	272
21.3 markdown 格式说明	272
21.4 附录: pandoc 软件介绍	287
22 R Markdown 文件格式	289

22.1 R Markdown 文件	289
22.2 R Markdown 文件的编译	291
22.3 在 R Markdown 文件中插入 R 代码	292
22.4 输出表格	294
22.5 利用 R 程序插图	296
22.6 代码段选项	297
22.7 章节目录链接问题	306
22.8 数学公式	306
22.9 其它编程语言引擎	310
22.10交互内容	310
22.11属性设置	310
22.12LaTeX 和 PDF 输出	317
22.13生成期刊文章	319
22.14附录：经验与问题	320
23 用 bookdown 制作图书	323
23.1 介绍	323
23.2 一本书的设置	324
23.3 章节结构	328
23.4 书的编译	329
23.5 交叉引用	330
23.6 数学公式和公式编号	331
23.7 定理类编号	332
23.8 文献引用	333
23.9 插图	334
23.10表格	336
23.11数学公式的设置	336
23.12使用经验	338
23.13bookdown 的一些使用问题	340
24 用 R Markdown 制作简易网站	343
24.1 介绍	343
24.2 简易网站制作	343
24.3 用 blogdown 制作网站	347

25 制作幻灯片	355
25.1 介绍	355
25.2 Slidy 幻灯片	355
25.3 MS PowerPoint 幻灯片	360
25.4 Beamer 幻灯片格式	361
25.5 R Presentation 格式	362
V R 数据处理	363
26 数据读取技巧	365
26.1 日期数据	365
26.2 缺失值处理	369
26.3 练习	371
27 数据整理	373
27.1 tidyverse 系统	373
27.2 用 <code>filter()</code> 选择行子集	376
27.3 用 <code>sample_n()</code> 对观测随机抽样	377
27.4 用 <code>distinct()</code> 去除重复行	378
27.5 用 <code>drop_na()</code> 去除指定的变量有缺失值的行	379
27.6 用 <code>select()</code> 选择列子集	379
27.7 用 <code>arrange()</code> 排序	384
27.8 用 <code>rename()</code> 修改变量名	386
27.9 用 <code>mutate()</code> 计算新变量	386
27.10 用管道连接多次操作	389
27.11 数据简单汇总	390
27.12 长宽表转换	395
27.13 拆分数据列	399
27.14 合并数据列	400
27.15 横向合并	401
27.16 利用第二个数据集筛选	404
27.17 数据集的集合操作	404
27.18 数据框纵向合并	404
27.19 标准化	406

27.20用 reshape 包做长宽表转换	409
28 数据汇总	419
28.1 用 summary() 函数作简单概括	419
28.2 连续型变量概括函数	422
28.3 分类变量概括	422
28.4 数据框概括	427
28.5 分类概括	428
28.6 练习	437
VI 绘图	439
29 绘图	441
29.1 常用高级图形	441
29.2 低级图形函数	469
29.3 图形参数	478
29.4 图形输出	484
29.5 包含多种中文字体的图形	485
29.6 其它图形	488
30 ggplot 作图入门	491
30.1 介绍	491
30.2 作图的一般原则	495
30.3 散点图: ggplot 入门	496
30.4 分组、小图、变换、条形图	511
30.5 更多图形种类	544
30.6 刻度 (scale)	577
30.7 如何使用颜色	585
30.8 标题、标注、指南、拼接	590
30.9 图形定制调整	598
30.10主题	602
30.11参考文献	603
31 ggplot 的各种图形	605

目录	11
31.1 介绍	605
31.2 表现数量	605
31.3 表现分布	624
31.4 表现比例	661
31.5 表现多个变量间的关系	673
31.6 时间序列图	685
31.7 拟合曲线图	696
31.8 表现不确定性	705
VII 统计分析	715
32 R 初等统计分析	717
32.1 概率分布	717
32.2 最大似然估计	718
32.3 假设检验和置信区间	724
33 R 相关与回归	755
33.1 相关分析	755
33.2 一元回归分析	765
33.3 多元线性回归	779
33.4 非参数回归	840
33.5 Logistic 回归	850
34 统计学习介绍	859
34.1 Hitters 数据分析	860
34.2 Heart 数据分析	891
34.3 汽车销量数据分析	903
34.4 波士顿郊区房价数据	913
34.5 附录	922
VIII 专题	931
35 R 语言的文本处理	933
35.1 简单的文本处理	933

35.2 文本文件读写	948
35.3 正则表达式	950
35.4 stringr 包的正则表达式函数	963
35.5 利用基本 R 函数进行正则表达式处理	968
35.6 正则表达式应用例子	983
IX 用 Rcpp 连接 C++ 代码	999
36 Rcpp 介绍	1001
36.1 Rcpp 的用途	1002
36.2 Rcpp 入门样例	1002
37 R 与 C++ 的类型转换	1009
37.1 用 wrap() 把 C++ 变量返回到 R 中	1009
37.2 用 as() 函数把 R 变量转换为 C++ 类型	1010
37.3 as() 和 wrap() 的隐含调用	1010
38 Rcpp 属性	1011
38.1 Rcpp 属性介绍	1011
38.2 在 C++ 源程序中指定要导出的 C++ 函数	1012
38.3 在 R 中编译链接 C++ 代码	1013
38.4 Rcpp 属性的其它功能	1015
39 Rcpp 提供的 C++ 数据类型	1017
39.1 RObject 类	1017
39.2 IntegerVector 类	1018
39.3 NumericVector 类	1021
39.4 Rcpp 的其它向量类	1025
39.5 Rcpp 提供的其它数据类型	1027
40 Rcpp 糖	1033
40.1 简单示例	1033
40.2 向量化的运算符	1034
40.3 用 Rcpp 访问数学函数	1035
40.4 返回单一逻辑值的函数	1038

目录	13
40.5 返回糖表达式的函数	1039
40.6 R 与 Rcpp 不同语法示例	1042
41 用 Rcpp 帮助制作 R 扩展包	1043
41.1 不用扩展包共享 C++ 代码的方法	1043
41.2 生成扩展包	1044
41.3 重新编译	1046
41.4 建立 C++ 用的接口界面	1047
X R 编程例子	1049
42 R 编程例子	1051
42.1 R 语言	1051
42.2 概率	1052
42.3 科学计算	1053
42.4 统计计算	1055
42.5 数据处理	1063
42.6 文本处理	1066

前言

李东风的《R 语言教程》的草稿。还在更新中。

相关下载:

- `Rbook-data.zip`: 一些配套数据的打包文件
- `Bookdown-template-v0-3.zip`: R Markdown 和 bookdown 的模板

书中的数学公式使用 MathJax 库显示, 下面是数学公式测试。如果数学公式显示的中文不正常, 在浏览器中用鼠标右键单击中文公式, 在弹出的菜单中选择 **Math Settings--Math Renderer** 选 HTML-CSS 或 SVG 即可。

公式测试:

$$f(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

中文公式测试:

$$\text{相对误差} = \frac{a - A}{A}$$

使用本教程必须安装的软件包:

- tidyverse
- bookdown
- xtable
- microbenchmark
- reshape

本教程中用到的软件包列表 (更新 R 软件后可以在一个基本 R 中运行如下命令):

```
pkgs <- c(
  "assertthat",
  "backports", "base64enc", "BH", "bindr", "bindrcpp", "bookdown", "broom",
  "callr", "cellranger", "cli", "clipr", "clorspace", "crayon", "curl",
  "DBI", "dbplyr", "dichromat", "digest", "dplyr",
  "evaluate",
  "forcats",
  "ggplot2", "glue", "gtable",
  "haven", "highr", "hms", "htmltools", "httr",
  "jsonlite",
  "knitr",
  "labeling", "lazyeval", "lubridate",
  "magrittr", "markdown", "microbenchmark", "mime", "mnormt", "modelr", "munsell",
  "openssl",
  "pillar", "pkgconfig", "plogr", "plyr", "psych", "purrr",
  "R6", "RColorBrewer", "Rcpp", "readr", "readxl",
  "rematch", "reprex", "reshape2", "rlang", "rmarkdown",
  "rprojroot", "rstudioapi", "rvest",
  "scales", "selectr", "stringi", "stringr",
  "tibble", "tidyr", "tidyselect", "tidyverse",
  "utf8",
  "viridisLite",
  "whisker",
  "xml2", "xtable",
  "yaml"
)
install.packages(unique(pkgs))
```


Part I

介绍与入门

Chapter 1

R 语言介绍

1.1 R 的历史和特点

1.1.1 R 的历史

R 语言来自 S 语言，是 S 语言的一个变种。S 语言由 Rick Becker, John Chambers 等人在贝尔实验室开发，著名的 C 语言、Unix 系统也是贝尔实验室开发的。

S 语言第一个版本开发于 1976-1980，基于 Fortran；于 1980 年移植到 Unix，并对外发布源代码。1984 年出版的“棕皮书” (Becker and Chambers, 1984) 总结了 1984 年为止的版本，并开始发布授权的源代码。这个版本叫做旧 S。与我们现在用的 S 语言有较大差别。

1989-1988 对 S 进行了较大更新，变成了我们现在使用的 S 语言，称为第二版。1988 年出版的“蓝皮书” (Becker et al., 1988) 做了总结。

1992 年出版的“白皮书” (Chambers and Hastie, 1992) 描述了在 S 语言中实现的统计建模功能，增强了面向对象的特性。软件称为第三版，这是我们现在用的多数版本。

1998 年出版的“绿皮书” (Chambers, 2008) 描述了第四版 S 语言，主要是编程功能的深层次改进。现行的 S 系统并没有都采用第四版，S-PLUS 的第 5 版

才采用了 S 语言第四版。

S 语言商业版本为 S-PLUS, 1988 年发布, 现在为 Tibco Software 拥有。命运多舛, 多次易主。

R 是一个自由软件, GPL 授权, 最初由新西兰 Auckland 大学的 Ross Ihaka 和 Robert Gentleman 于 1997 年发布, R 实现了与 S 语言基本相同的功能和统计功能。现在由 R 核心团队开发, 但全世界的用户都可以贡献软件包。R 的网站: <http://www.r-project.org/>

1.1.2 R 的特点

1.1.2.1 R 语言一般特点

- 自由软件, 免费、开放源代码, 支持各个主要计算机系统;
- 完整的程序设计语言, 基于函数和对象, 可以自定义函数, 调入 C、C++、Fortran 编译的代码;
- 具有完善的数据类型, 如向量、矩阵、因子、数据集、一般对象等, 支持缺失值, 代码像伪代码一样简洁、可读;
- 强调交互式数据分析, 支持复杂算法描述, 图形功能强;
- 实现了经典的、现代的统计方法, 如参数和非参数假设检验、线性回归、广义线性回归、非线性回归、可加模型、树回归、混合模型、方差分析、判别、聚类、时间序列分析等。
- 统计科研工作者广泛使用 R 进行计算和发表算法。R 有上万软件包 (截止 2019 年 7 月有一万四千多个)。

1.1.2.2 R 语言和 R 软件的技术特点

- 函数编程 (functional programming)。R 语言虽然不是严格的 functional programming 语言, 但可以遵照其原则编程, 得到可验证的可靠程序。
- 支持对象类和类方法。基于对象的程序设计。
- 是动态类型语言, 解释执行, 运行速度较慢。
- 数据框是基本的观测数据类型, 类似于数据库的表。
- 开源软件 (Open source software)。可深入探查, 开发者和用户交互。
- 可以用作 C 和 C++、FORTRAN 语言编写的算法库的接口。

- 主要数值算法采用已广泛测试和采纳的算法实现，如排序、随机数生成、线性代数（LAPACK 软件包）。

1.1.2.3 推荐参考书

- R.L. Kabacoff(2012)《R 语言实战》，人民邮电出版社。
- Hadley Wickham and Garrett Grolemund(2017)“R for Data Science”，<http://r4ds.had.co.nz/>, O'Reilly
- Hadley Wickham(2014)“Advanced R”，<http://adv-r.had.co.nz/>, Chapman & Hall/CRC The R Series
- R 网站上的初学者手册“An Introduction to R”和其它技术手册。
- John M. Chambers(2008), “Software for Data Analysis-Programming with R”, Springer.
- Venables, W. N. & Ripley, B. D.(2002) “Modern Applied Statistics with S”, Springer
- 薛毅、陈立萍（2007）《统计建模与 R 软件》，清华大学出版社。
- 汤银才（2008），《R 语言与统计分析》，高等教育出版社。
- 李东风（2006）《统计软件教程》，人民邮电出版社。

1.2 R 的下载与安装

1.2.1 R 的下载

以 MS Windows 操作系统为例。R 的主网站在<https://www.r-project.org/>。从 CRAN 的镜像网站下载软件，其中一个镜像如<http://mirror.bjtu.edu.cn/cran/>。选“Download R for Windows-base-Download R 3.4.1 for windows”（3.4.1 是版本号，应下载网站上给出的最新版本）链接进行下载。在“Download R for Windows”链接的页面，除了 base 为 R 的安装程序，还有 contrib 为 R 附加的扩展软件包下载链接（一般不需要从这里下载），以及 Rtools 链接，是在 R 中调用 C、C++ 和 Fortran 程序代码时需要用的编译工具。

RStudio (<https://www.rstudio.com/>) 是功能更强的一个 R 图形界面，在安装好 R 的官方版本后安装 RStudio 可以更方便地使用 R。

1.2.2 R 软件安装

下载官方的 R 软件后按提示安装。安装后获得一个桌面快捷方式，如“R i386 3.4.1”(这是 32 位版本)。如果是 64 位操作系统，可以同时安装 32 位版本和 64 位版本，对初学者这两种版本区别不大，尽量选用 64 位版本，这是将来的趋势。

安装官方的 R 软件后，可以安装 RStudio。平时使用可以使用 RStudio，其界面更方便，对 R Markdown 格式 (.Rmd) 文件支持更好。

如果使用 RStudio，每个分析项目需要单独建立一个“项目”(project)，每个项目也有一个工作文件夹。

1.2.3 辅助软件

R 可以把一段程序写在一个以 .r 或 .R 为扩展名的文本文件中，如“date.r”，称为一个 `_ 源程序 _` 文件，然后在 R 命令行用

```
source("date.r")
```

运行源程序。这样的文件可以用记事本生成和编辑。

在 MS Windows 操作系统中建议使用 notepad++ 软件，这是 MS Windows 下记事本程序的增强型软件。安装后，在 MS Windows 资源管理器中右键弹出菜单会有“edit with notepadpp”选项。notepad++ 可以方便地在不同的中文编码之间转换。

RStudio 则是一个集成环境，可以在 RStudio 内进行源程序文件编辑和运行。

1.2.4 R 扩展软件包的安装与管理

R 扩展软件包提供了特殊功能。以安装 sos 包为例。sos 包用来搜索某些函数的帮助文档。在 R 图形界面选菜单“程序包-安装程序包”，在弹出的“CRAN mirror”选择窗口中选择一个中国的镜像如“China (Beijing 2)”，然后在弹出的“Packages”选择窗口中选择要安装的扩展软件包名称，即可完成下载和安装。

还可以用如下程序制定镜像网站 (例子中是位于清华大学的镜像网站) 并安装指定的扩展包:

```
options(repos=c(CRAN="http://mirror.tuna.tsinghua.edu.cn/CRAN/"))
install.packages("sos")
```

还可以选择扩展包的安装路径, 如果权限允许, 可以选择安装在 R 软件的主目录内或者用户自己的私有目录位置。由于用户的对子目录的读写权限问题, 有时不允许一般用户安装扩展包到 R 的主目录中。用 `.libPaths()` 查看允许的扩展包安装位置, 在 `install.packages()` 中用 `lib=` 指定安装位置:

```
print(.libPaths())
## [1] "D:/R/R-3.3.1/library"
install.packages("sos", lib=.libPaths()[1])
```

在 RStudio 中用 “Tools” 菜单的 “Install Packages” 安装软件包。

在每一次 R 软件更新后, 需要重新安装原来的软件包, 这个过程很麻烦。如果仅仅是小的版本更新, 比如从 3.5.1 变成 3.5.2, 或者从 3.4.2 变成 3.5.0, 可以在安装新版本后, 将老版本的 library 子目录中所有内容复制到新版本的 library 子目录中, 同时尽量不要覆盖已有的内容, 然后在基本 R 中 (不要用 RStudio) 运行如下命令:

```
options(repos=c(CRAN="http://mirror.tuna.tsinghua.edu.cn/CRAN/"))
update.packages(checkBuilt=TRUE, ask=FALSE)
```

这个命令也可以用来成批地更新已安装的 R 扩展软件包。

如果版本改变比较大, 可以用如下方法批量地重新安装原有的软件包。首先, 在更新 R 软件前, 在原来的 R 中运行:

```
packages <- .packages(TRUE)
dump("packages", file="packages-20180704.R")
```

这样可以获得要安装的软件包的列表。在更新 R 软件后, 运行如下程序:

```
options(repos=c(CRAN="http://mirror.tuna.tsinghua.edu.cn/CRAN/"))
source("packages-20180704.R")
install.packages(packages)
```

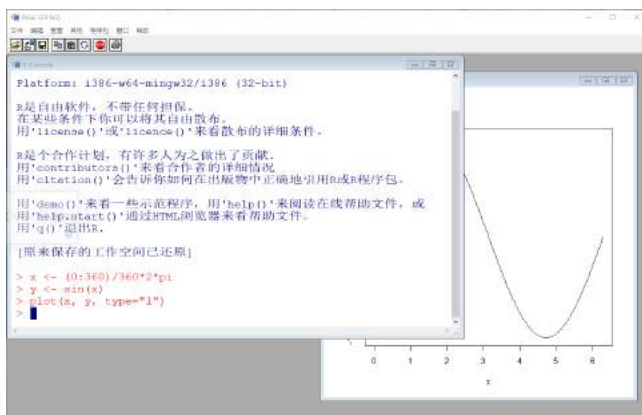


图 1.1: R GUI 截图

安装时如果提问是否安装需要编译的源代码包，最好选择否，因为安装源代码包速度很慢还有可能失败。

1.3 基本 R 软件的使用

1.3.1 基本运行

在 MS Windows 操作系统中的 R 软件有一个 R GUI 软件，即图形窗口模式的 R 软件，如图1.1。

R GUI 中有一个命令行窗口 (R Console)，以大于号为提示符，在提示符后面键入命令，命令的文字型结果马上显示在命令下方，命令的图形结果单独显示在一个图形窗口中。

在命令行可以通过左右光标键移动光标到适当位置进行修改。可以用上下光标在已经运行过的历史命令中切换，调回已经运行过的命令，修改后重新执行。

如果某个文件如 `myprog.R` 在当前工作目录中，保存的都是 R 程序，称这样的文件为源程序文件。可以在命令行用如下命令运行其中的程序：

```
source("myprog.R")
```

但是，在 MS Windows 操作系统中，默认的中文编码是 GB18030 编码。R 源

程序文件的中文编码可能是 GB18030 也可能是 UTF-8。UTF-8 是在世界范围更通用的编码。如果发现用如下命令运行时出现中文乱码，可能是因为源程序用了 UTF-8 编码，这时 `source()` 命令要加上编码选项如下：

```
source("myprog.R", encoding="UTF-8")
```

1.3.2 项目目录

用 R 进行数据分析,不同的分析问题需要放在不同的文件夹中。以 MS Windows 操作系统为例,设某个分析问题的数据与程序放在了 `c:\work` 文件夹中。把 R 的快捷方式从桌面复制入此文件夹,在 Windows 资源管理器中,右键单击此快捷方式,在弹出菜单中选“属性”,把“快捷方式”页面的“起始位置”的内容清除为空白,点击确定按钮。启动在 `work` 文件夹中的 R 快捷方式,出现命令行界面。这时, `C:\work` 称为**当前工作目录**。

在命令行运行如下命令可以显示当前工作目录位置：

```
getwd()  
## "C:/work"
```

显示结果中的目录、子目录、文件之间的分隔符用了 `/` 符号,在 Windows 操作系统中一般应该使用 `\` 符号,但是,在 R 的字符串中一个 `\` 需要写成两个,所以等价的写法是 `"C:\\work"`。

不同的分析项目需要存放在不同的文件夹中,每个文件夹都放置一个“起始位置”为空的 R 快捷方式图标,分析哪一个项目,就从对应的快捷图标启动,而不是从桌面上的 R 图标启动。这样做的好处是,用到源文件和数据文件时,只要文件在该项目的文件夹中,就不需要写完全路径而只需要用文件名即可。

1.4 RStudio 软件

1.4.1 介绍

RStudio 软件是 R 软件的应用界面与增强系统,可以在其中编辑、运行 R 的程序文件,可以跟踪运行,还可以构造文字、R 结果图表融合在一起的研究报

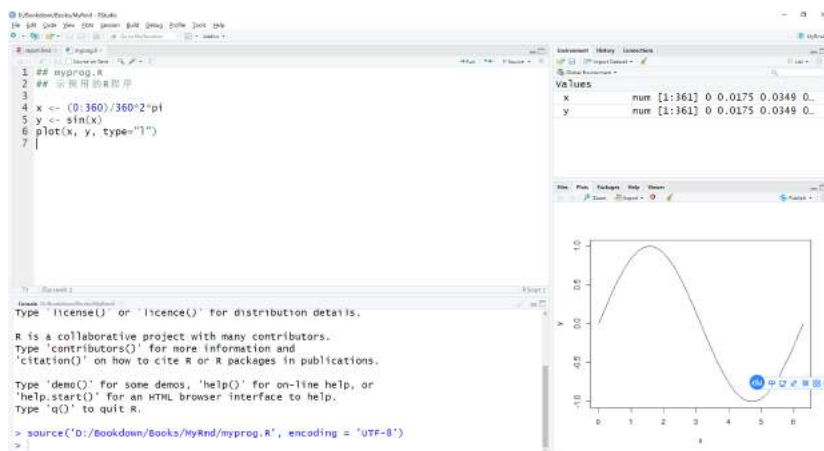


图 1.2: RStudio 截图

告、论文、图书、网站等。一个运行中的 RStudio 界面见图1.2。

界面一般分为四个窗格，其中编辑窗口与控制台（Console）是最重要的两个窗格。编辑窗格用来查看和编辑程序、文本型的数据文件、程序与文字融合在一起的 Rmd 文件等。控制台与基本 R 软件的命令行窗口基本相同，功能有所增强。

在编辑窗口中可以用操作系统中常用的编辑方法对源文件进行编辑，如复制、粘贴、查找、替换，还支持基于正则表达式的查找替换（关于正则表达式见35）。

其它的一些重要窗格包括：

- **Files:** 列出当前项目的目录（文件夹）内容。其中以 `.R` 或者 `.r` 为扩展名的是 R 源程序文件，单击某一源程序文件就可以在编辑窗格中打开该文件。
- **Plots:** 如果程序中有绘图结果，将会显示在这个窗格。因为绘图需要足够的空间，所以当屏幕分辨率过低或者 Plots 窗格太小的时候，可以点击“Zoom”图标将图形显示在一个单独的窗口中，或者将图形窗口作为唯一窗格显示。如何放大窗格见下面的使用技巧。
- **Help:** R 软件的文档与 RStudio 的文档都在这里。
- **Environment:** 已经有定义的变量、函数都显示在这里。
- **History:** 以前运行过的命令都显示在这里。不限于本次 RStudio 运行期间，也包括以前使用 RStudio 时运行过的命令。

- Packages: 显示已安装的 R 扩展包及其文档。
- Viewer, Connection, Build, Git 等窗格。

1.4.2 项目

用 R 和 RStudio 进行研究和数据分析，每个研究问题应该单独建立一个文件夹（目录）。该问题的所有数据、程序都放在对应的文件夹中。在 RStudio 中，用“File – New Project – Existing Directory”选中该问题的目录，建立一个新的“项目”（project）。

再次进入 RStudio 后，用菜单“File – Recent Projects”找到已有的项目打开，然后就可以针对该项目进行了。这样分项目进行研究的好处是，不同项目的可以使用同名的文件而不会有冲突，程序中用到某个文件时，只需要写文件名而不需要写文件所在的目录。

一个项目还可以有项目本身的一些特殊设置，用“Tools – Project Options”菜单打开设置。

1.4.3 帮助

在 RStudio 中有一个单独的 Help 窗格，如果需要，可以用菜单“View–Panes–Zoom help”将其放大到占据整个窗口空间。但是，这一功能目前不支持放大显示字体的功能，不如在浏览器中方便。

RStudio 的帮助窗格中包含 R 软件的官方文档，以及 RStudio 软件的文档。“Search engine and keywords”项下面有分类的帮助。有软件包列表。

在基本 R 软件而不是 RStudio 的命令行中运行命令 `help.start()` 或者用 RGUI 的帮助菜单中“html 帮助”可以打开系统默认的互联网浏览器，在其中查看帮助文档。

在命令行，用问号后面跟随函数名查询某函数的帮助。用 `example(函数名)` 的格式可以运行此函数的样例，如：

```
example(mean)
```

```
##
```

```
## mean> x <- c(0:10, 50)
##
## mean> xm <- mean(x)
##
## mean> c(xm, mean(x, trim = 0.10))
## [1] 8.75 5.50
```

有时仅知道一些方法的名字而不知道具体的扩展包和函数名称，可以安装 `sos` 扩展包 (package)，用 `findFn(" 函数名")` 查询某个函数，结果显示在互联网浏览器软件中。

1.4.4 使用技巧

RStudio 使用方法概要 PDF 下载：

[rstudio-ide.pdf](#)。

1.4.4.1 使用历史

在控制台（命令行窗格）中，除了可以用左右光标键移动光标位置，用上下光标键调回以前运行过的命令，还有一个重要的增强（以 MS Windows 操作系统为例）：键入要运行的命令的前几个字母，如 `book`，按“Ctrl+ 向上光标键”，就可以显示历史命令中以 `book` 开头的的所有命令，单击哪一个，哪一个就自动复制到命令行。这一技巧十分重要，我们需要反复允许同一命令时，这一方法让我们很容易从许多命令历史中找到所需的命令。

1.4.4.2 放大显示某一窗格

当屏幕分辨率较低时，将整个 RStudio 界面分为四个窗格会使得每个窗格都没有足够的显示精度。为此，可以将某个窗格放大到整个窗口区域，需要使用其它窗格时再恢复到四个窗格的状态或者直接放大其它窗格到整个窗口区域。

使用菜单“View – Panes – Zoom Source”可以将编辑窗格放到最大，在 MS Windows 下也可以使用快捷键“Ctrl+Alt+1”。其它操作系统也有类似的快捷键可用。使用菜单“View – Panes – Show All Panes”可以显示所有四个窗格。

放大其它窗格也可以用“Ctrl + Alt + 数字”，数字与窗格的对应关系为：

- 1: 编辑窗格；
- 2: 控制台 (Console)；
- 3: 帮助；
- 4: 历史；
- 5: 文件；
- 6: 图形；
- 7: 扩展包；
- 8: 已定义变量和函数；
- 9: 研究报告或网站结果显示。

1.4.4.3 运行程序

可以在命令行直接输入命令运行，文字结果会显示在命令行窗口，图形结果显示在“Plots”窗格中。在命令行窗口 (Console) 中可以用左右光标键移动光标，用上下光标键查找历史命令，输入命令的前几个字母后用“Ctrl+ 向上光标键”可以匹配地查找历史命令。

一般情况下，还是应该将 R 源程序保存在一个源程序文件中运行。RStudio 中“File – New File – R Script”可以打开一个新的无名的 R 源程序文件窗口供输入 R 源程序用。输入一些程序后，保存文件，然后点击“Source”快捷图标就可以运行整个文件中的所有源程序，并会自动加上关于编码的选项。

编写 R 程序的正常做法是一边写一遍试验运行，运行一般不是整体的运行而是写完一部分就运行一部分，运行没有错误才继续编写下一部分。在 R 源程序窗口中，当光标在某一程序上的时候，点击窗口的“Run”快捷图标或者用快捷键“Ctrl+Enter 键”可以运行该行；选中若干程序行后，点击窗口的“Run”快捷图标或者用快捷键“Ctrl+Enter 键”可以运行这些行。

1.4.4.4 中文编码问题

对于中文内容的 R 源程序、R Markdown 源文件 (.Rmd 文件)、文本型数据文件 (.txt, .csv)，其中的中文内容可能有不同的编码选择，在中国国内主要使

用 GB18030(基本兼容于 GB, GBK) 和 UTF-8, UTF-8 是国际上更普遍使用的统一文字编码, 涉及到计算机编程时应尽可能使用此编码系统。

在 RStudio 中新生成的 R 源程序、Rmd 源文件一般自动用 UTF-8 编码。点击 RStudio 的文件窗格中显示的源文件, 可以打开该源文件, 但是因为已有源文件的编码不一定与 RStudio 的默认编码一致, 可以会显示成乱码。为此, RStudio 提供了“File – Reopen with Encoding”命令, 我们主要试验其中 GB18030 和 UTF-8 两种选择一般就可以解决问题。如果选择 GB18030 显示就没有乱码了, 最好再用菜单“File – Save with Encoding”并选择 UTF-8 将其保存为 UTF-8 编码。

其它的文本格式的文件也可以类似地处理, 后面将会陆续提及。

1.4.5 Rmd 文件

在科学研究中, R 软件可以用来分析数据, 生成数据分析报表和图形。R Markdown(简称 Rmd) 是一种特殊的文件格式, 在这种文件中, 即有 R 程序, 又有说明文字, 通过 R 和 RStudio 软件, 可以运行其中的程序, 并将说明文字、程序、程序的文字结果、图形结果统一地转换为一个研究报告, 支持 Word、PDF、网页、网站、幻灯片等许多种输出格式。在打开的 Rmd 源文件中, 也可以选择其中的某一段 R 程序单独运行。所以, Rmd 文件也可以作为一种特殊的 R 源程序文件。

用 RStudio 的“File – New File – R Markdown”菜单就可以生成一个新的 Rmd 文件并显示在编辑窗格中, 其中已经有了一些样例内容, 可以修改这些样例内容为自己的文字和程序。

Rmd 文件中用 ````\{r}` 开头, 用 ````` 结尾的段落是 R 程序段, 在显示的程序段的右侧有一个向右箭头形状的小图标 (类似于媒体播放图标), 点击该图标就可以运行该程序段。

打开 Rmd 文件后, 用编辑窗口的 Knit 命令可以选择将文件整个地转换为 HTML(网页) 或者 MS Word 格式, 如果操作系统中安装有 LaTeX 软件, 还可以以 LaTeX 为中间格式转换为 PDF 文件。

详见20、21、22。

1.5 练习

1. 下载 R 安装程序，安装 R，建立 work 文件夹并在其中建立 R 的快捷方式。Windows 用户还需要下载 RTools 软件并安装。
2. 下载 RStudio 软件并安装。
3. 下载安装 notepad++ 软件 (仅 MS Windows 用户)。
4. 在 RStudio 中下载安装 sos 扩展软件包。

Chapter 2

R 语言入门运行样例

2.1 命令行界面

启动 R 软件后进入命令行界面，每输入一行命令，就在后面显示计算结果。可以用向上和向下箭头访问历史命令；可以从已经运行过的命令中用鼠标拖选加亮后，用 Ctrl+C 复制后用 Ctrl+V 粘贴，或用 Ctrl+X 一步完成复制粘贴，粘贴的目标都是当前命令行。

如果使用 RStudio 软件，有一个“Console 窗格”相当于命令行界面。在 RStudio 中，可以用 New File-Script file 功能建立一个源程序文件（脚本文件），在脚本文件中写程序，然后用 Run 图标或者 Ctrl+Enter 键运行当前行或者选定的部分。

2.2 四则运算

四则运算如：

```
5 + (2.3 - 1.125)*3.2/1.1 + 1.23E3
## [1] 1238.418
```

结果为 1238.418，前面显示的结果在行首加了井号，这在 R 语言中表示注释。本教程的输出前面一般都加了井号以区分于程序语句。输出前面的方括号和序

号 1 是在输出有多个值时提供的提示性序号，只有单个值时为了统一起见也显示出来了。这里 $1.23E3$ 是科学记数法，表示 1.23×10^3 。用星号 $*$ 表示乘法，用正斜杠 $/$ 表示除法。

用 \wedge 表示乘方运算，如

```
2^10
## [1] 1024
```

重要提示：关闭中文输入法，否则输入一些中文标点将导致程序错误。

2.2.1 计算例子

从 52 张扑克牌中任取 3 张，有多少种不同的组合可能？解答：有

$$C_{52}^3 = \frac{52!}{3!(52-3)!} = \frac{52 \times 51 \times 50}{3 \times 2 \times 1}$$

种，在 R 中计算如：

```
52*51*50/(3*2)
## [1] 22100
```

2.2.2 练习

- 某人存入 10000 元 1 年期定期存款，年利率 3%，约定到期自动转存（包括利息）。问：
 - 10 年后本息共多少元？
 - 需要存多少年这 10000 元才能增值到 20000 元？
- 成语说：“智者千虑，必有一失；愚者千虑，必有一得”。设智者作判断的准确率为 $p_1 = 0.99$ ，愚者作判断的准确率为 $p_2 = 0.01$ ，计算智者做 1000 次独立的判断至少犯一次错误的概率，与愚者做 1000 次独立判断至少对一次的概率。

2.3 数学函数

2.3.1 数学函数——平方根、指数、对数

例:

```
sqrt(6.25)
## [1] 2.5
exp(1)
## [1] 2.718282
log10(10000)
## [1] 4
```

`sqrt(6.25)` 表示 $\sqrt{6.25}$, 结果为 2.5。`exp(1)` 表示 e^1 , 结果为 $e = 2.718282$ 。`log10(10000)` 表示 $\lg 10000$, 结果为 4。`log` 为自然对数。

2.3.2 数学函数——取整

例:

```
round(1.1234, 2)
## [1] 1.12
round(-1.9876, 2)
## [1] -1.99
floor(1.1234)
## [1] 1
floor(-1.1234)
## [1] -2
ceiling(1.1234)
## [1] 2
ceiling(-1.1234)
## [1] -1
```

`round(1.1234, 2)` 表示把 1.1234 四舍五入到两位小数。`floor(1.1234)` 表示把 1.1234 向下取整, 结果为 1。`ceiling(1.1234)` 表示把 1.1234 向上取整,

结果为 2。

2.3.3 数学函数——三角函数

例:

```
pi
## [1] 3.141593
sin(pi/6)
## [1] 0.5
cos(pi/6)
## [1] 0.8660254
sqrt(3)/2
## [1] 0.8660254
tan(pi/6)
## [1] 0.5773503
```

`pi` 表示圆周率 π 。`sin` 正弦, `cos` 余弦, `tan` 正切, 自变量以弧度为单位。`pi/6` 是 30° 。

2.3.4 数学函数——反三角函数

例:

```
pi/6
## [1] 0.5235988
asin(0.5)
## [1] 0.5235988
acos(sqrt(3)/2)
## [1] 0.5235988
atan(sqrt(3)/3)
## [1] 0.5235988
```

`asin` 反正弦, `acos` 反余弦, `atan` 反正切, 结果以弧度为单位。

2.3.5 分布函数和分位数函数

例:

```
dnorm(1.98)
## [1] 0.05618314
pnorm(1.98)
## [1] 0.9761482
qnorm(0.975)
## [1] 1.959964
```

`dnorm(x)` 表示标准正态分布密度 $\phi(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}x^2}$. `pnorm(x)` 表示标准正态分布函数 $\Phi(x) = \int_{-\infty}^x \phi(t) dt$. `qnorm(y)` 表示标准正态分布分位数函数 $\Phi^{-1}(x)$. 还有其它许多分布的密度函数、分布函数和分位数函数。例如,

```
qt(1 - 0.05/2, 10)
## [1] 2.228139
```

求自由度为 10 的 t 检验的双侧临界值。其中 `qt(y,df)` 表示自由度为 `df` 的 t 分布的分位数函数。

2.4 输出

2.4.1 简单输出

命令行的计算结果直接显示在命令的后面。在用 `source()` 运行程序文件时, 需要用 `print()` 函数显示一个表达式的结果, 如:

```
print(sin(pi/2))
## [1] 1
```

用 `cat()` 函数显示多项内容, 包括数值和文本, 文本包在两个单撇号或两个双撇号中, 如:

```
cat("sin(pi/2)=", sin(pi/2), "\n")
## sin(pi/2)= 1
```

`cat()` 函数最后一项一般是 `"\n"`, 表示换行。忽略此项将不换行。

再次提示: 要避免打开中文输入法导致误使用中文标点。

2.4.2 用 `sink()` 函数作运行记录

R 使用经常是在命令行逐行输入命令(程序), 结果紧接着显示在命令后面。如何保存这些命令和显示结果? 在 R 命令行中运行过的命令会被保存在运行的工作文件夹中的一个名为 `.Rhistory` 的文件中。用 `sink()` 函数打开一个文本文件开始记录文本型输出结果。结束记录时用空的 `sink()` 即可关闭文件不再记录。如

```
sink("tmpres01.txt", split=TRUE)
print(sin(pi/6))
print(cos(pi/6))
cat("t(10) 的双侧 0.05 分位数(临界值)=", qt(1 - 0.05/2, 10), "\n")
sink()
```

`sink()` 用作输出记录主要是在测试运行中使用, 正常的输出应该使用 `cat()` 函数、`write.table()`、`write.csv()` 等函数。

2.4.3 练习

1. 用 `cat()` 函数显示

```
log10(2)=*** log10(5)=***
```

其中 `***` 应该代以实际函数值。

2. 用 `sink()` 函数开始把运行过程记录到文件“`log001.txt`”中, 在命令行试验几个命令, 然后关闭运行记录, 查看生成的“`log001.txt`”的内容。

2.5 向量计算与变量赋值

R 语言以向量为最小单位。用 `<-` 赋值。如

```
x1 <- 1:10
x1
## [1] 1 2 3 4 5 6 7 8 9 10
```

一般的向量可以用 `c()` 生成，如

```
marks <- c(3, 5, 10, 5, 6)
```

在程序语言中，变量用来保存输入的值或计算的结果。变量可以存放各种不同类型的值，如单个数值、多个数值（称为向量）、单个字符串、多个字符串（称为字符型向量），等等。单个数值称为**标量**。

技术秘诀：用程序设计语言的术语描述，R 语言是动态类型的，其变量的类型不需要预先声明，运行过程中允许变量类型改变，实际上变量赋值是一种“绑定”（binding），将一个变量的名称（变量名）与实际的一个存储位置联系在一起。在命令行定义的变量称为**全局变量**。

用 `print()` 函数显示向量或在命令行中显示向量时，每行显示的行首会有方括号和数字序号，代表该行显示的第一个向量元素的下标。如

```
12345678901:12345678920
## [1] 12345678901 12345678902 12345678903 12345678904 12345678905
## [6] 12345678906 12345678907 12345678908 12345678909 12345678910
## [11] 12345678911 12345678912 12345678913 12345678914 12345678915
## [16] 12345678916 12345678917 12345678918 12345678919 12345678920
```

向量可以和一个标量作四则运算，结果是每个元素都和这个标量作四则运算，如：

```
x1 + 200
## [1] 201 202 203 204 205 206 207 208 209 210
2*x1
## [1] 2 4 6 8 10 12 14 16 18 20
2520/x1
## [1] 2520 1260 840 630 504 420 360 315 280 252
```

两个等长的向量可以进行四则运算，相当于对应元素进行四则运算，如

```
x2 <- x1 * 3
x2
## [1] 3 6 9 12 15 18 21 24 27 30
x2 - x1
## [1] 2 4 6 8 10 12 14 16 18 20
```

R 的许多函数都可以用向量作为自变量，结果是自变量的每个元素各自的函数值。如

```
sqrt(x1)
## [1] 1.000000 1.414214 1.732051 2.000000 2.236068 2.449490 2.645751
## [8] 2.828427 3.000000 3.162278
```

结果是 1 到 10 的整数各自的平方根。

2.6 工作空间介绍

在命令行中定义的变量，在退出 R 时，会提问是否保存工作空间，初学时可选择保存，真正用 R 进行数据分析时往往不保存工作空间。再次启动 R 后，能够看到以前定义的几个变量的值。

在使用 R 的官方版本时，如果在 Windows 中使用，一般把不同的数据分析项目放在不同的文件夹中。将 R 的程序快捷图标复制到每一个项目的文件夹中，并用右键菜单讲快捷图标的“属性”中“起始位置”改为空白。要分析哪一个项目的数据，就在那个项目文件夹中的 R 快捷图标启动，这样可以保证不同的项目有不同的工作空间。

如果使用 RStudio 软件，也需要把不同项目放在不同文件夹，并且每个项目在 RStudio 中单独建立一个“项目”（project）。要分析那个项目的数据，就打开那个项目。不同项目使用不同的工作空间。

RStudio 中的“Environment”窗格会显示当前已定义的 R 变量与函数。

2.6.1 练习

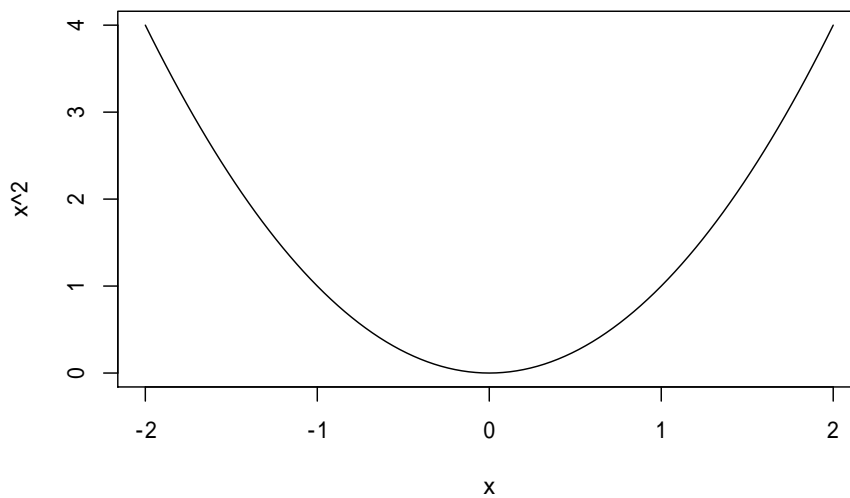
1. 某人存入 10000 元 1 年期定期存款，年利率 3%，约定到期自动转存（包括利息）。列出 1、2、.....、10 年后的本息金额。
2. 显示 2 的 1,2,....., 20 次方。
3. 定义 x1 为 1 到 10 的向量，定义 x2 为 x1 的 3 倍，然后退出 R，再次启动 R，查看 x1 和 x2 的值。

2.7 绘图示例

2.7.1 函数曲线示例

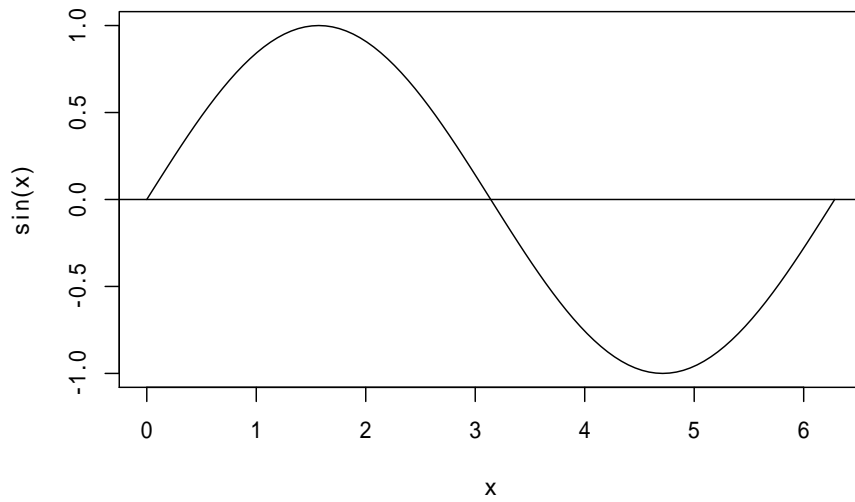
如下程序用 `curve()` 函数制作 $y = x^2$ 函数的曲线图，`curve()` 函数第二、第三自变量是绘图区间：

```
curve(x^2, -2, 2)
```



类似地， $\sin(x)$ 函数曲线图用如下程序可制作，用 `abline()` 函数添加参考线：

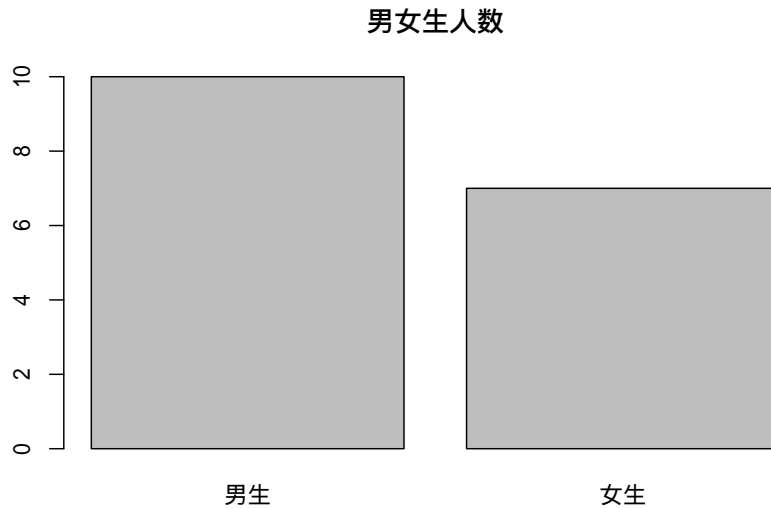
```
curve(sin(x), 0, 2*pi)
abline(h=0)
```



2.7.2 条形图示例

假设有 10 个男生，7 个女生，如下程序绘制男生、女生人数的条形图：

```
barplot(c(" 男生 "=10, " 女生 "=7),
        main=" 男女生人数")
```

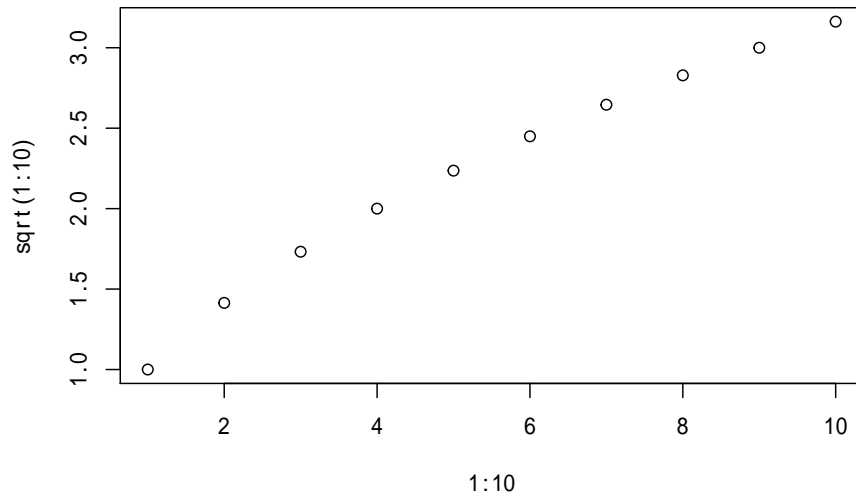


利用适当选项可以人为定制颜色、控制条形宽窄。实际问题中，个数（频数）一般是从数据中统计得到的。

2.7.3 散点图示例

下面的例子用 `plot()` 函数做了散点图, `plot()` 函数第一个自变量是各个点的横坐标值, 第二个自变量是对应的纵坐标值。

```
plot(1:10, sqrt(1:10))
```



2.7.4 R 软件自带的图形示例

R 软件中自带了一些演示图形。通过如下程序可以调用：

```
demo("graphics")  
demo("image")
```

2.7.5 练习

1. 画 $\exp(x)$ 在 $(-2, 2)$ 区间的函数图形。
2. 画 $\ln(x)$ 在 $(0.01, 10)$ 区间的函数图形。

2.8 汇总统计示例

2.8.1 表格数据

统计用的输入数据典型样式是 Excel 表那样的表格数据。表格数据特点：每一列应该是相同的类型（或者都是数值，或者都是文字，或者都是日期），每一列应该有一个名字。

这样的表格数据，一般可以保存为.csv 格式：数据项之间用逗号分开，文件本身是文本型的，可以用普通记事本程序查看和编辑。Excel 表可以用“另存为”命令保存为.csv 格式。常用的数据库管理系统一般也可以把表保存为.csv 格式。

2.8.2 读入表格数据

例如，`taxsamp.csv` 是这样一个 csv 格式表格数据文件，可以用 Excel 打开，也可以用记事本程序或 notepad++ 打开。内容见本页面后面的附录。

用如下程序可以把.csv 文件读入到 R 中：

```
tax.tab <- read.csv("taxsamp.csv", header=TRUE, as.is=TRUE)
print(head(tax.tab))
```

程序中的选项 `header=TRUE` 指明第一行作为变量名行，选项 `as.is=TRUE` 说明字符型列要原样读入而不是转换为因子 (factor)。读入的变量 `tax.tab` 称为一个数据框 (data.frame)。`head()` 函数返回数据框或向量的前几项。比较大的表最好不要显示整个表，会使得前面的运行过程难以查看。

技巧：`read.csv()` 的一个改进版本是 `readr` 扩展包的 `read_csv()` 函数，此函数读入较大表格速度要快得多，而且读入的转换设置更倾向于不做不必要的转换。但是，这两种输入方法的默认中文编码可能不一样。

2.8.3 练习

1. 用 Excel 软件查看“`taxsamp.csv`”的内容（双击即可）。
2. 用记事本程序或 notepad++ 软件查看“`taxsamp.csv`”的内容。

3. 读入 “taxsamp.csv” 到 R 数据框 tax.tab 中，查看 tax.tab 内容。

2.8.4 分类变量频数统计

在 tax.tab 中，“征收方式” 是一个分类变量。用 table() 函数计算每个不同值的个数，称为频数 (frequency):

```
table(tax.tab[[" 征收方式"]])

##
##      查帐征收  定期定额征收  定期定率征收
##           31           16           2
```

类似地可以统计“申报渠道”的取值频数:

```
table(tax.tab[[" 申报渠道"]])

##
##  大厅申报  网上申报
##       18       31
```

也可以用 table() 函数统计“征收方式”和“申报渠道”交叉分类频数，如:

```
table(tax.tab[[" 征收方式"]], tax.tab[[" 申报渠道"]])

##
##                大厅申报  网上申报
##  查帐征收           9       22
##  定期定额征收       9        7
##  定期定率征收       0        2
```

上述结果制表如下:

```
knitr::kable(table(tax.tab[[" 征收方式"]], tax.tab[[" 申报渠道"]]))
```

	大厅申报	网上申报
查帐征收	9	22
定期定额征收	9	7
定期定率征收	0	2

2.8.5 数值型变量的统计

数值型变量可以计算各种不同的统计量，如平均值、标准差和各个分位数。`summary()` 可以给出最小值、最大值、中位数、四分之一分位数、四分之三分位数和平均值。如

```
summary(tax.tab[["营业额"]])
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         0     650    2130 247327   9421 6048000
```

中位数是从小到大排序后排在中间的值。四分之一和四分之三分位数类似。

统计函数以一个数值型向量为自变量，包括 `sum`(求和), `mean`(平均值), `var`(样本方差), `sd`(样本标准差), `min`(最小值), `max`(最大值), `range`(最小值和最大值) 等。如

```
mean(tax.tab[["营业额"]])
```

```
## [1] 247327.4
```

```
sd(tax.tab[["营业额"]])
```

```
## [1] 1036453
```

如果数据中有缺失值，可以删去缺失值后计算统计量，这时在 `mean`, `sd` 等函数中加入 `na.rm=TRUE` 选项。

2.8.6 练习

用如下程序定义一个变量 `x`, 然后求 `x` 的平均值和最小值、最大值。

```
x <- c(3, 5, 10, 5, 6)
```

2.9 运行源程序文件

用 `source()` 函数可以运行保存在一个文本文件中的源程序。比如，如下内容保存在文件 `ssq.r` 中：

```
sum.of.squares <- function(x){  
  sum(x^2)  
}
```

用如下 `source()` 命令运行:

```
source("ssq.r")
```

运行后就可以调用自定义函数 `sum.of.squares()` 了。

2.9.1 源文件编码

源程序文件存在编码问题。对于源程序编码与系统默认编码不同的情况，在 `source()` 函数中可以添加 `encoding=` 选项。例如，保存为 UTF-8 编码的源程序在简体中文 MS Windows 系统的 R 中运行，可以在 `source()` 函数中可以添加 `encoding="UTF-8"` 选项。保存为 GBK 编码的源程序文件在 MAC 系统的 R 中运行，可以在 `source()` 函数中可以添加 `encoding="GBK"` 选项。

在 RStudio 中，可以打开一个源程序文件查看与编辑。用快捷键“Ctrl+Enter”或快捷图标“Run”可以运行当前行或者加亮选中行，快捷图标“Source”可以运行整个文件。如果发现中文乱码，可以用菜单“Reopen with encoding”选择合适的编码打开，用菜单“Save with encoding”选择需要的编码保存。

2.9.2 当前工作目录

在用 `source()` 调用源程序文件或者用 `read.csv()` 读入数据文件时，如果不写文件名的全路径，就认为文件位置是在所谓“当前工作目录”。用 `getwd()` 函数可以查询当前工作目录，用 `setwd()` 函数可以设置当前工作目录。在 RStudio 中用菜单“Session-Set working directory”设置当前工作目录。

在 MS Windows 操作系统中使用 R 软件时，一种好的做法是把某个研究项目所有数据和程序放在某个文件夹如 `c:\work` 中，把 R 的程序快捷图标复制到该目录中，在资源管理器中对该图标调用鼠标右键菜单“属性”，从弹出对话框中，把“起始位置”一栏清除。这样，每次从这个快捷图标启动 R，就可以自动以所在子目录为当前工作目录，工作空间和命令历史记录也默认存放在这里。

在 MS Windows 操作系统的 R 中使用文件路径时，要用正斜杠作为连接符，使用反斜杠则需要成对使用，如 `setwd("d:/work")` 或 `setwd("d:\\work")`。

如果使用 RStudio 软件，将某个研究项目所有数据和程序放在某个文件夹中，然后建立一个新项目（project）指向该文件夹。

2.9.3 练习

编辑生成 `ssq.r` 源程序文件并用 `source()` 函数运行，然后计算：

```
sum.of.squares(1:5)
```

2.10 附录：数据

2.10.1 公司纳税数据样例

本数据是某地区 2013 年 12 月所属税款申报信息的一个子集，仅含 20 个公司的数据。

数据文件显示

Part II

R 的数据类型与相应运算

Chapter 3

常量与变量

3.1 常量

R 语言基本的数据类型有数值型，逻辑型 (TRUE, FALSE)，文本 (字符串)。支持缺失值，有专门的复数类型。

常量是指直接写在程序中的值。

数值型常量包括整型、单精度、双精度等，一般不需要区分。写法如 123, 123.45, -123.45, -0.012, 1.23E2, -1.2E-2 等。为了表示 123 是整型，可以写成 123L。

字符型常量用两个双撇号或两个单撇号包围，如 "Li Ming" 或 'Li Ming'。字符型支持中文，如 "李明" 或 '李明'。国内的中文编码主要有 GBK 编码和 UTF-8 编码，有时会遇到编码错误造成乱码的问题，MS Windows 下 R 程序一般用 GBK 编码，但是 RStudio 软件采用 UTF-8 编码。在 R 软件内字符串一般用 UTF-8 编码保存。

逻辑型常量只有 TRUE 和 FALSE。

缺失值用 NA 表示。统计计算中经常会遇到缺失值，表示记录丢失、因为错误而不能用、节假日没有数据等。除了数值型，逻辑型和字符型也可以有缺失值，而且字符型的空白值不会自动识别为缺失值，需要自己规定。R 支持特殊的 Inf 值，这是实数型值，表示正无穷大，不算缺失值。

复数常量写法如 $2.2 + 3.5i$, $1i$ 等。

3.2 变量

程序语言中的变量用来保存输入的值或者计算得到的值。在 R 中，变量可以保存所有的数据类型，比如标量、向量、矩阵、数据框、函数等。

变量都有变量名，R 变量名必须以字母、数字、下划线和句点组成，变量名的第一个字符不能取为数字。在中文环境下，汉字也可以作为变量名的合法字符使用。变量名是区分大小写的， y 和 Y 是两个不同的变量名。

变量名举例: `x`, `x1`, `X`, `cancer.tab`, `clean_data`, `diseaseData`。

用 `<-` 赋值的方法定义变量。`<-` 也可以写成 `=`，但是 `<-` 更直观。如

```
x5 <- 6.25
x6 = sqrt(x5)
```

R 的变量没有固定的类型，给已有变量赋值为新的类型，该变量就变成新的类型，但一般应避免这样的行为。R 是“动态类型”语言，赋值实际上是“绑定” (binding)，即将一个变量名与一个存储地址联系在一起，同一个存储地址可以有多个变量名与其联系。

3.3 R 数据类型

R 语言基本的数据类型有数值，逻辑型 (TRUE, FALSE)，文本 (字符串)。支持缺失值，有专门的复数类型。

R 语言数据结构包括向量，矩阵和数据框，多维数组，列表，对象等。数据中元素、行、列还可以用名字访问。最基本的是向量类型。向量类型数据的访问方式也是其他数据类型访问方式的基础。

Chapter 4

数值型向量及其运算

4.1 数值型向量

向量是将若干个基础类型相同的值存储在一起，各个元素可以按序号访问。如果将若干个数值存储在一起可以用序号访问，就叫做一个数值型向量。

用 `c()` 函数把多个元素或向量组合成一个向量。如

```
marks <- c(10, 6, 4, 7, 8)
x <- c(1:3, 10:13)
x1 <- c(1, 2)
x2 <- c(3, 4)
x <- c(x1, x2)
x
```

```
## [1] 1 2 3 4
```

10:13 这样的写法表示从 10 到 13 的整数组成的向量。

用 `print()` 函数显示向量或在命令行中显示向量时，每行显示的行首会有方括号和数字序号，代表该行显示的第一个向量元素的下标。如

```
12345678901:12345678920
```

```
## [1] 12345678901 12345678902 12345678903 12345678904 12345678905
```

```
## [6] 12345678906 12345678907 12345678908 12345678909 12345678910
## [11] 12345678911 12345678912 12345678913 12345678914 12345678915
## [16] 12345678916 12345678917 12345678918 12345678919 12345678920
```

`length(x)` 可以求 `x` 的长度。长度为零的向量表示为 `numeric(0)`。`numeric()` 函数可以用来初始化一个指定元素个数而元素都等于零的数值型向量，如 `numeric(10)` 会生成元素为 10 个零的向量。

4.2 向量运算

4.2.1 标量和标量运算

单个数值称为标量，R 没有单独的标量类型，标量实际是长度为 1 的向量。

R 中四则运算用 `+` `-` `*` `/` `^` 表示 (加、减、乘、除、乘方)，如

```
1.5 + 2.3 - 0.6 + 2.1*1.2 - 1.5/0.5 + 2^3
```

```
## [1] 10.72
```

R 中四则运算仍遵从通常的优先级规则，可以用圆括号 `()` 改变运算的先后次序。如

```
1.5 + 2.3 - (0.6 + 2.1)*1.2 - 1.5/0.5 + 2^3
```

```
## [1] 5.56
```

除了加、减、乘、除、乘方，R 还支持整除运算和求余运算。用 `%%` 表示整除，用 `%%` 表示求余。如

```
5 %% 3
```

```
## [1] 1
```

```
5 %% 3
```

```
## [1] 2
```

```
5.1 %% 2.5
```

```
## [1] 2
```



```
5.1 %% 2.5
```

```
## [1] 0.1
```

4.2.2 向量与标量运算

向量与标量的运算为每个元素与标量的运算, 如

```
x <- c(1, 10)
```

```
x + 2
```

```
## [1] 3 12
```

```
x - 2
```

```
## [1] -1 8
```

```
x * 2
```

```
## [1] 2 20
```

```
x / 2
```

```
## [1] 0.5 5.0
```

```
x ^ 2
```

```
## [1] 1 100
```

```
2 / x
```

```
## [1] 2.0 0.2
```

```
2 ^ x
```

```
## [1] 2 1024
```

一个向量乘以一个标量, 就是线性代数中的数乘运算。

四则运算时如果有缺失值, 缺失元素参加的运算相应结果元素仍缺失。如

```
c(1, NA, 3) + 10
```

```
## [1] 11 NA 13
```

4.2.3 等长向量运算

等长向量的运算为对应元素两两运算。如

```
x1 <- c(1, 10)
x2 <- c(4, 2)
x1 + x2
```

```
## [1] 5 12
```

```
x1 - x2
```

```
## [1] -3 8
```

```
x1 * x2
```

```
## [1] 4 20
```

```
x1 / x2
```

```
## [1] 0.25 5.00
```

两个等长向量的加、减运算就是线性代数中两个向量的加、减运算。

4.2.4 不等长向量的运算

两个不等长向量的四则运算，如果其长度为倍数关系，规则是每次从头重复利用短的一个。如

```
x1 <- c(10, 20)
x2 <- c(1, 3, 5, 7)
x1 + x2
```

```
## [1] 11 23 15 27
```

```
x1 * x2
```

```
## [1] 10 60 50 140
```

不仅是四则运算，R 中有两个或多个向量按照元素一一对应参与某种运算或函数调用时，如果向量长度不同，一般都采用这样的规则。

如果两个向量的长度不是倍数关系，会给出警告信息。如

```
c(1,2) + c(1,2,3)
```

```
## Warning in c(1, 2) + c(1, 2, 3): 长的对象长度不是短的对象长度的整倍数
## [1] 2 4 4
```

4.3 向量函数

4.3.1 向量化的函数

R 中的函数一般都是向量化的: 在 R 中, 如果普通的一元函数以向量为自变量, 一般会对每个元素计算。这样的函数包括 `sqrt`, `log10`, `log`, `exp`, `sin`, `cos`, `tan` 等许多。如

```
sqrt(c(1, 4, 6.25))
```

```
## [1] 1.0 2.0 2.5
```

为了查看这些基础的数学函数的列表, 运行命令 `help.start()`, 点击链接“Search Engine and Keywords”, 找到“Mathematics”栏目, 浏览其中的“arith”和“math”链接中的说明。常用的数学函数有:

- 舍入: `ceiling`, `floor`, `round`, `signif`, `trunc`, `zapsmall`
- 符号函数 `sign`
- 绝对值 `abs`
- 平方根 `sqrt`
- 对数与指数函数 `log`, `exp`, `log10`, `log2`
- 三角函数 `sin`, `cos`, `tan`
- 反三角函数 `asin`, `acos`, `atan`, `atan2`
- 双曲函数 `sinh`, `cosh`, `tanh`
- 反双曲函数 `asinh`, `acosh`, `atanh`

有一些不太常用的数学函数:

- 贝塔函数 `beta`, `lbeta`
- 伽玛函数 `gamma`, `lgamma`, `digamma`, `trigamma`, `tetragamma`, `pentagamma`

- 组合数 `choose`, `lchoose`
- 富利叶变换和卷积 `fft`, `mvfft`, `convolve`
- 正交多项式 `poly`
- 求根 `polyroot`, `uniroot`
- 最优化 `optimize`, `optim`
- Bessel 函数 `besselI`, `besselK`, `besselJ`, `besselY`
- 样条插值 `spline`, `splinefun`
- 简单的微分 `deriv`

如果自己编写的函数没有考虑向量化问题，可以用 `Vectorize()` 函数将其转换成向量化版本。

4.3.2 排序函数

`sort(x)` 返回排序结果。`rev(x)` 返回把各元素排列次序反转后的结果。`order(x)` 返回排序用的下标。如

```
x <- c(33, 55, 11)
sort(x)
```

```
## [1] 11 33 55
```

```
rev(sort(x))
```

```
## [1] 55 33 11
```

```
order(x)
```

```
## [1] 3 1 2
```

```
x[order(x)]
```

```
## [1] 11 33 55
```

例子中，`order(x)` 结果中 3 是 `x` 的最小元素 11 所在的位置下标，1 是 `x` 的第二小元素 33 所在的位置下标，2 是 `x` 的最大元素 55 所在的位置下标。

4.3.3 统计函数

`sum`(求和), `mean`(求平均值), `var`(求样本方差), `sd`(求样本标准差), `min`(求最小值), `max`(求最大值), `range`(求最小值和最大值) 等函数称为统计函数, 把输入向量看作样本, 计算样本统计量。`prod` 求所有元素的乘积。

`cumsum` 和 `cumprod` 计算累加和累乘积。如

```
cumsum(1:5)
```

```
## [1] 1 3 6 10 15
```

```
cumprod(1:5)
```

```
## [1] 1 2 6 24 120
```

其它一些类似函数有 `pmax`, `pmin`, `cummax`, `cummin` 等。

4.3.4 生成规则序列的函数

`seq` 函数是冒号运算符的推广。比如, `seq(5)` 等同于 `1:5`。`seq(2,5)` 等同于 `2:5`。`seq(11, 15, by=2)` 产生 11,13,15。`seq(0, 2*pi, length.out=100)` 产生从 0 到 2π 的等间隔序列, 序列长度指定为 100。

从这些例子可以看出, S 函数可以带自变量名调用。每个函数的变量名和用法可以查询其帮助信息, 在命令行界面用“? 函数名”的方法查询。在使用变量名时次序可以颠倒, 比如 `seq(to=5, from=2)` 仍等同于 `2:5`。

`rep()` 函数用来产生重复数值。为了产生一个初值为零的长度为 `n` 的向量, 用 `x <- rep(0, n)`。`rep(c(1,3), 2)` 把第一个自变量重复两次, 结果相当于 `c(1,3,1,3)`。

`rep(c(1,3), c(2,4))` 则需要利用 R 的一般向量化规则, 把第一自变量的第一个元素 1 按照第二自变量中第一个元素 2 的次数重复, 把第一自变量中第二个元素 3 按照第二自变量中第二个元素 4 的次数重复, 结果相当于 `c(1,1,3,3,3,3)`。

如果希望重复完一个元素后再重复另一元素, 用 `each=` 选项, 比如 `rep(c(1,3), each=2)` 结果相当于 `c(1,1,3,3)`。

4.4 复数向量

复数常数表示如 $3.5+2.4i$, $1i$ 。用函数 `complex()` 生成复数向量, 指定实部和虚部。如 `complex(c(1,0,-1,0), c(0,1,0,-1))` 相当于 `c(1+0i, 1i, -1+0i, -1i)`。

在 `complex()` 中可以用 `mod` 和 `arg` 指定模和辐角, 如 `complex(mod=1, arg=(0:3)/2*pi)` 结果同上。

用 `Re(z)` 求 z 的实部, 用 `Im(z)` 求 z 的虚部, 用 `Mod(z)` 或 `abs(z)` 求 z 的模, 用 `Arg(z)` 求 z 的辐角, 用 `Conj(z)` 求 z 的共轭。

`sqrt`, `log`, `exp`, `sin` 等函数对复数也有定义, 但是函数定义域在自变量为实数时可能有限制而复数无限制, 这时需要区分自变量类型。如

```
sqrt(-1)
## [1] NaN
## Warning message:
## In sqrt(-1) : NaNs produced
sqrt(-1 + 0i)
## [1] 0+1i
```

4.5 练习

1. 显示 1 到 100 的整数的平方根和立方根 (提示: 立方根就是三分之一次方)。
2. 设有 10 个人的小测验成绩为:

77 60 91 73 85 82 35 100 66 75

- (1) 把这 10 个成绩存入变量 `x`;
 - (2) 从小到大排序;
 - (3) 计算 `order(x)`, 解释 `order(x)` 结果中第 3 项代表的意义。
 - (4) 计算这些成绩的平均值、标准差、最小值、最大值、中位数。
3. 生成 $[0, 1]$ 区间上等间隔的 100 个格子点存入变量 `x` 中。

Chapter 5

逻辑型向量及其运算

5.1 逻辑型向量与比较运算

逻辑型是 R 的基本数据类型之一，只有两个值 TRUE 和 FALSE, 缺失时为 NA。逻辑值一般产生自比较，如

```
sele <- (log10(15) < 2); print(sele)
```

```
## [1] TRUE
```

向量比较结果为逻辑型向量。如

```
c(1, 3, 5) > 2
```

```
## [1] FALSE TRUE TRUE
```

```
(1:4) >= (4:1)
```

```
## [1] FALSE FALSE TRUE TRUE
```

从例子可以看出，向量比较也遵从 R 的向量间运算的一般规则：向量与标量的运算是向量每个元素与标量都分别运算一次，等长向量的运算时对应元素的运算，不等长但长度为倍数关系的向量运算是把短的从头重复利用。

与 NA 比较产生 NA，如

```
c(1, NA, 3) > 2
```

```
## [1] FALSE NA TRUE
```

```
NA == NA
```

```
## [1] NA
```

为了判断向量每个元素是否 NA，用 `is.na()` 函数，如

```
is.na(c(1, NA, 3) > 2)
```

```
## [1] FALSE TRUE FALSE
```

用 `is.finite()` 判断向量每个元素是否 Inf 值。

比较运算符包括

```
< <= > >= == != %in%
```

分别表示小于、小于等于、大于、大于等于、等于、不等于、属于。要注意等于比较用了两个等号。

`%in%` 是比较特殊的比较，`x %in% y` 的运算把向量 `y` 看成集合，运算结果是一个逻辑型向量，第 i 个元素的值为 `x` 的第 i 元素是否属于 `y` 的逻辑型值。如

```
c(1,3) %in% c(2,3,4)
```

```
## [1] FALSE TRUE
```

```
c(NA,3) %in% c(2,3,4)
```

```
## [1] FALSE TRUE
```

```
c(1,3) %in% c(NA, 3, 4)
```

```
## [1] FALSE TRUE
```

```
c(NA,3) %in% c(NA, 3, 4)
```

```
## [1] TRUE TRUE
```

函数 `match(x, y)` 起到和 `x %in% y` 运算类似的作用，但是其返回结果不是找到与否，而是对 `x` 的每个元素，找到其在 `y` 中首次出现的下标，找不到时取缺失值，如


```
match(c(1, 3), c(2,3,4,3))
```

```
## [1] NA 2
```

5.2 逻辑运算

为了表达如“ $x > 0$ 而且 $x < 1$ ”，“ $x \leq 0$ 或者 $x \geq 1$ ”之类的复合比较，需要使用逻辑运算把两个比较连接起来。逻辑运算符为 `&`、`|` 和 `!`，分别表示“同时成立”、“两者至少其一成立”、“条件的反面”。比如，设 `age<=3` 表示婴儿，`sex=='女'` 表示女性，则 `age<=3 & sex=='女'` 表示女婴，`age<=3 | sex=='女'` 表示婴儿或妇女，`!(age<=3 | sex=='女')` 表示既非婴儿也非妇女。为了确定运算的先后次序可以用圆括号 `()` 指定。

用 `xor(x, y)` 表示 `x` 与 `y` 的异或运算，即值不相等时为真值，相等时为假值，有缺失值参加运算时为缺失值。

逻辑向量与逻辑标量之间的逻辑运算，两个逻辑向量之间的逻辑运算规则遵从一般 R 向量间运算规则。

在右运算符是缺失值时，如果左运算符能够确定结果真假，可以得到非缺失的结果。例如，`TRUE | NA` 为 `TRUE`，`FALSE & NA` 为 `FALSE`。不能确定结果时返回 `NA`，比如，`TRUE & NA` 为 `NA`，`FALSE | NA` 为 `NA`。

`&&` 和 `||` 分别为短路的标量逻辑与和短路的标量逻辑或，仅对两个标量进行运算，如果有向量也仅使用第一个元素。一般用在 `if` 语句、`while` 语句中，且只要第一个比较已经决定最终结果就不计算第二个比较。例如

```
if(TRUE || sqrt(-1)>0) next
```

其中的 `sqrt(-1)` 部分不会执行。

这里结果为 `TRUE`，第二部分没有参加计算，否则第二部分的计算会发生函数自变量范围错误。

5.3 逻辑运算函数

因为 R 中比较与逻辑运算都支持向量之间、向量与标量之间的运算，所以在需要一个标量结果时要特别注意，后面讲到的 `if` 结构、`while` 结构都需要逻辑标量而且不能是缺失值。这时，应该对缺失值结果单独考虑。

若 `cond` 是逻辑向量，用 `all(cond)` 测试 `cond` 的所有元素为真；用 `any(cond)` 测试 `cond` 至少一个元素为真。`cond` 中允许有缺失值，结果可能为缺失值。如

```
c(1, NA, 3) > 2
```

```
## [1] FALSE    NA    TRUE
```

```
all(c(1, NA, 3) > 2)
```

```
## [1] FALSE
```

```
any(c(1, NA, 3) > 2)
```

```
## [1] TRUE
```

```
all(NA)
```

```
## [1] NA
```

```
any(NA)
```

```
## [1] NA
```

函数 `which()` 返回真值对应的所有下标，如

```
which(c(FALSE, TRUE, TRUE, FALSE, NA))
```

```
## [1] 2 3
```

```
which((11:15) > 12)
```

```
## [1] 3 4 5
```

函数 `identical(x,y)` 比较两个 R 对象 `x` 与 `y` 的内容是否完全相同，结果只会取标量 `TRUE` 与 `FALSE` 两种。如

```
identical(c(1,2,3), c(1,2,NA))
```

```
## [1] FALSE
```

```
identical(c(1L,2L,3L), c(1,2,3))
```

```
## [1] FALSE
```

其中第二个结果假值是因为前一向量是整数型，后一向量是实数型。

函数 `all.equal()` 与 `identical()` 类似，但是在比较数值型时不区分整数型与实数型，而且相同时返回标量 `TRUE`，但是不同时会返回一个说明有何不同的字符串。如

```
all.equal(c(1,2,3), c(1,2,NA))
```

```
## [1] "'is.NA' value mismatch: 1 in current 0 in target"
```

```
all.equal(c(1L,2L,3L), c(1,2,3))
```

```
## [1] TRUE
```

函数 `duplicated()` 返回每个元素是否为重复值的结果，如：

```
duplicated(c(1,2,1,3,NA,4,NA))
```

```
## [1] FALSE FALSE TRUE FALSE FALSE FALSE TRUE
```

用函数 `unique()` 可以返回去掉重复值的结果。

Chapter 6

字符型数据及其处理

6.1 字符型向量

字符型向量是元素为字符串的向量。如

```
s1 <- c('abc', '', 'a cat', NA, ' 李明')
```

注意空字符串并不能自动认为是缺失值，字符型的缺失值仍用 `NA` 表示。

6.2 `paste()` 函数

针对字符型数据最常用的 R 函数是 `paste()` 函数。`paste()` 用来连接两个字符型向量，元素一一对应连接，默认用空格连接。如 `paste(c("ab", "cd"), c("ef", "gh"))` 结果相当于 `c("ab ef", "cd gh")`。

`paste()` 在连接两个字符型向量时采用 R 的一般向量间运算规则，而且可以自动把数值型向量转换为字符型向量。可以作一对多连接，如 `paste("x", 1:3)` 结果相当于 `c("x 1", "x 2", "x 3")`。

用 `sep=` 指定分隔符，如 `paste("x", 1:3, sep="")` 结果相当于 `c("x1", "x2", "x3")`。

使用 `collapse=` 参数可以把字符型向量的各个元素连接成一个单一的字符串, 如 `paste(c("a", "b", "c"), collapse="")` 结果相当于 "abc"。

6.3 转换大小写

`toupper()` 函数把字符型向量内容转为大写, `tolower()` 函数转为小写。比如, `toupper('aB cd')` 结果为 "AB CD", `tolower(c('aB', 'cd'))` 结果相当于 `c("ab" "cd")`。这两个函数可以用于不区分大小写的比较, 比如, 不论 `x` 的值是 'JAN', 'Jan' 还是 'jan', `toupper(x)=='JAN'` 的结果都为 TRUE。

6.4 字符串长度

用 `nchar(x, type='bytes')` 计算字符型向量 `x` 中每个字符串的以字节为单位的长度, 这一点对中英文是有差别的, 中文通常一个汉字占两个字节, 英文字母、数字、标点占一个字节。用 `nchar(x, type='chars')` 计算字符型向量 `x` 中每个字符串的以字符个数为单位的长度, 这时一个汉字算一个单位。

在画图时可以用 `strwidth()` 函数计算某个字符串或表达式占用的空间大小。

6.5 取子串

`substr(x, start, stop)` 从字符串 `x` 中取出从第 `start` 个到第 `stop` 个的子串, 如

```
substr('JAN07', 1, 3)
```

```
## [1] "JAN"
```

如果 `x` 是一个字符型向量, `substr` 将对每个元素取子串。如

```
substr(c('JAN07', 'MAR66'), 1, 3)
```

```
## [1] "JAN" "MAR"
```

用 `substring(x, start)` 可以从字符串 `x` 中取出从第 `start` 个到末尾的子串。如

```
substring(c('JAN07', 'MAR66'), 4)
```

```
## [1] "07" "66"
```

6.6 类型转换

用 `as.numeric()` 把内容是数字的字符型值转换为数值，如

```
substr('JAN07', 4, 5)
## [1] "07"
substr('JAN07', 4, 5) + 2000
## Error in substr("JAN07", 4, 5) + 2000 :
## non-numeric argument to binary operator
as.numeric(substr('JAN07', 4, 5)) + 2000
## [1] 2007
as.numeric(substr(c('JAN07', 'MAR66'), 4, 5))
## [1] 7 66
```

`as.numeric()` 是向量化的，可以转换一个向量的每个元素为数值型。

用 `as.character()` 函数把数值型转换为字符型，如

```
as.character((1:5)*5)
```

```
## [1] "5" "10" "15" "20" "25"
```

如果自变量本来已经是字符型则结果不变。

为了用指定的格式数值型转换成字符型，可以使用 `sprintf()` 函数，其用法与 C 语言的 `sprintf()` 函数相似，只不过是向量化的。例如

```
sprintf('file%03d.txt', c(1, 99, 100))
```

```
## [1] "file001.txt" "file099.txt" "file100.txt"
```

6.7 字符串拆分

用 `strsplit()` 函数可以把一个字符串按照某种分隔符拆分开，例如

```
x <- '10,8,7'
strsplit(x, ',', fixed=TRUE)[[1]]

## [1] "10" "8" "7"

sum(as.numeric(strsplit(x, ',', fixed=TRUE)[[1]]))

## [1] 25
```

因为 `strsplit()` 的结果是一个列表，这个函数延后再详细讲。

6.8 字符串替换功能

用 `gsub()` 可以替换字符串中的子串，这样的功能经常用在数据清理中。比如，把数据中的中文标点改为英文标点，去掉空格，等等。如

```
x <- '1, 3; 5'
gsub(';', ',', x, fixed=TRUE)

## [1] "1, 3, 5"

strsplit(gsub(';', ',', x, fixed=TRUE), ',')[[1]]

## [1] "1" " 3" " 5"
```

字符串 `x` 中分隔符既有逗号又有分号，上面的程序用 `gsub()` 把分号都换成逗号。

更多的文本数据（字符型数据）功能参见35。

6.9 正则表达式

正则表达式 (regular expression) 是一种匹配某种字符串模式的方法。用这样的方法，可以从字符串中查找某种模式的出现位置，替换某种模式，等等。这

样的技术可以用于文本数据的预处理，比如用网络爬虫下载的大量网页文本数据。R 中支持 perl 语言格式的正则表达式，`grep()` 和 `grep1()` 函数从字符串中查询某个模式，`sub()` 和 `gsub()` 替换某模式。比如，下面的程序把多于一个空格替换成一个空格

```
gsub('[:space:]+', ' ', 'a   cat   in a box', perl=TRUE)
```

```
## [1] "a cat in a box"
```

正则表达式功能强大但也不容易掌握。详见35。

Chapter 7

R 向量下标和子集

在 R 中下标与子集是极为强大的功能，需要一些练习才能熟练掌握，许多其它语言中需要多个语句才能完成的工作在 R 中都可以简单地通过下标和子集来完成。

7.1 正整数下标

对向量 x ，在后面加方括号和下标可以访问向量的元素和子集。

设 $x \leftarrow c(1, 4, 6.25)$ 。 $x[2]$ 取出第二个元素； $x[2] \leftarrow 99$ 修改第二个元素。 $x[c(1,3)]$ 取出第 1、3 号元素； $x[c(1,3)] \leftarrow c(11, 13)$ 修改第 1、3 号元素。下标可重复。例如

```
x <- c(1, 4, 6.25)
x[2]
```

```
## [1] 4
```

```
x[2] <- 99; x
```

```
## [1] 1.00 99.00 6.25
```

```
x[c(1,3)]
## [1] 1.00 6.25
x[c(1,3)] <- c(11, 13); x
## [1] 11 99 13
x[c(1,3,1)]
## [1] 11 13 11
```

7.2 负整数下标

负下标表示扣除相应的元素后的子集，如

```
x <- c(1,4,6.25)
x[-2]
## [1] 1.00 6.25
x[-c(1,3)]
## [1] 4
```

负整数下标不能与正整数下标同时用来从某一向量中取子集，比如，`x[c(1,-2)]` 没有意义。

7.3 空下标与零下标

`x[]` 表示取 `x` 的全部元素作为子集。这与 `x` 本身不同，比如

```
x <- c(1,4,6.25)
x[] <- 999
x
## [1] 999 999 999
```

```
x <- c(1,4,6.25)
x <- 999
x
```

```
## [1] 999
```

`x[0]` 是一种少见的做法,结果返回类型相同、长度为零的向量,如 `numeric(0)`。相当于空集。

当 0 与正整数下标一起使用时会被忽略。当 0 与负整数下标一起使用时也会被忽略。

7.4 下标超界

设向量 `x` 长度为 n , 则使用正整数下标时下标应在 $\{1, 2, \dots, n\}$ 中取值。如果使用大于 n 的下标,读取时返回缺失值,并不出错。给超出 n 的下标元素赋值,则向量自动变长,中间没有赋值的元素为缺失值。例如

```
x <- c(1,4,6.25)
x[5]
```

```
## [1] NA
```

```
x
```

```
## [1] 1.00 4.00 6.25
```

```
x[5] <- 9
```

```
x
```

```
## [1] 1.00 4.00 6.25 NA 9.00
```

虽然 R 的语法对下标超界不视作错误,但是这样的做法往往来自不良的程序思路,而且对程序效率有影响,所以实际编程中应避免下标超界。

7.5 逻辑下标

下标可以是与向量等长的逻辑表达式，一般是关于本向量或者与本向量等长的其它向量的比较结果，如

```
x <- c(1,4,6.25)
x[x > 3]
```

```
## [1] 4.00 6.25
```

取出 x 的大于 3 的元素组成的子集。

逻辑下标除了用来对向量取子集，还经常用来对数据框取子集，也用在向量化的运算中。例如，对如下示性函数

$$f(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

输入向量 x ，结果 y 需要也是一个向量，程序可以写成

```
f <- function(x){
  y <- numeric(length(x))
  y[x >= 0] <- 1
  y[x < 0] <- 0 # 此语句多余

  y
}
```

事实上，向量化的逻辑选择有一个 `ifelse()` 函数，比如，对上面的示性函数，如果 x 是一个向量，输出 y 向量可以写成 `y <- ifelse(x>=0, 1, 0)`。

要注意的是，如果逻辑下标中有缺失值，对应结果也是缺失值。所以，在用逻辑下标作子集选择时，一定要考虑到缺失值问题。正确的做法是加上 `!is.na` 前提，如

```
x <- c(1, 4, 6.25, NA)
x[x > 2]
```

```
## [1] 4.00 6.25 NA
```

```
x[!is.na(x) & x > 2]
```

```
## [1] 4.00 6.25
```

7.6 which()、which.min()、which.max() 函数

函数 `which()` 可以用来找到满足条件的下标，如

```
x <- c(3, 4, 3, 5, 7, 5, 9)
```

```
which(x > 5)
```

```
## [1] 5 7
```

```
seq(along=x)[x > 5]
```

```
## [1] 5 7
```

这里 `seq(along=x)` 会生成由 `x` 的下标组成的向量。用 `which.min()`、`which.max` 求最小值的下标和最大值的下标，不唯一时只取第一个。如

```
which.min(x)
```

```
## [1] 1
```

```
which.max(x)
```

```
## [1] 7
```

7.7 元素名

向量可以为每个元素命名。如

```
ages <- c("李明"=30, "张聪"=25, "刘颖"=28)
```

或

```
ages <- c(30, 25, 28)
```

```
names(ages) <- c("李明", "张聪", "刘颖")
```

或

```
ages <- setNames(c(30, 25, 28), c("李明", "张聪", "刘颖"))
```

这时可以用元素名或元素名向量作为向量的下标，如

```
ages["张聪"]
```

```
## 张聪  
##    25
```

```
ages[c("李明", "刘颖")]
```

```
## 李明 刘颖  
##    30  28
```

```
ages["张聪"] <- 26
```

这实际上建立了字符串到数值的映射表。

用字符串作为下标时，如果该字符串不在向量的元素名中，读取时返回缺失值结果，赋值时该向量会增加一个元素并以该字符串为元素名。

带有元素名的向量也可以是字符型或其它基本类型，如

```
sex <- c("李明"="男", "张聪"="男", "刘颖"="女")
```

除了给向量元素命名外，在矩阵和数据框中还可以给行、列命名，这会使得程序的扩展更为容易和安全。

R 允许仅给部分元素命名，这时其它元素名字为空字符串。不同元素的元素名一般应该是不同的，否则在使用元素作为下标时会发生误读，但是 R 语法允许存在重名。

用 `unname(x)` 返回去掉了元素名的 `x` 的副本，用 `names(x) <- NULL` 可以去掉 `x` 的元素名。

7.8 用 R 向量下标作映射

R 在使用整数作为向量下标时，允许使用重复下标，这样可以把数组 x 看成一个 $1:n$ 的整数到 $x[1], x[2], \dots, x[n]$ 的一个映射表，其中 n 是 x 的长度。比如，某商店有三种礼品，编号为 1,2,3，价格分别为 68, 88 和 168。令

```
price.map <- c(68, 88, 168)
```

设某个收银员在一天内分别售出礼品编号为 3,2,1,1,2,2,3，可以用如下的映射方式获得售出的这些礼品对应的价格：

```
items <- c(3,2,1,1,2,2,3)
y <- price.map[items]; print(y)
```

```
## [1] 168 88 68 68 88 88 168
```

R 向量可以用字符型向量作下标，字符型下标也允许重复，所以可以把带有元素名的 R 向量看成是元素名到元素值的映射表。比如，设 sex 为 10 个学生的性别（男、女）

```
sex <- c("男", "男", "女", "女", "男", "女", "女", "女", "女", "男")
```

希望把每个学生按照性别分别对应到蓝色和红色。首先建立一个 R 向量当作映射

```
sex.color <- c('男'='blue', '女'='red')
```

用 R 向量 $sex.color$ 当作映射，可以获得每个学生对应的颜色

```
cols <- sex.color[sex]; print(cols)
```

```
##      男      男      女      女      男      女      女      女      女      男
## "blue" "blue" "red" "red" "blue" "red" "red" "red" "red" "blue"
```

这样的映射结果中带有不必要的元素名，用 `unname()` 函数可以去掉元素名，如

```
unname(cols)
```

```
## [1] "blue" "blue" "red" "red" "blue" "red" "red" "red" "red" "blue"
```

7.9 集合运算

可以把向量 x 看成一个集合，但是其中的元素允许有重复。用 `unique(x)` 可以获得 x 的所有不同值。如

```
unique(c(1, 5, 2, 5))
```

```
## [1] 1 5 2
```

用 `a %in% x` 判断 a 的每个元素是否属于向量 x ，如

```
5 %in% c(1,5,2)
```

```
## [1] TRUE
```

```
c(5,6) %in% c(1,5,2)
```

```
## [1] TRUE FALSE
```

与 `%in%` 运算符类似，函数 `match(x, table)` 对向量 x 的每个元素，从向量 $table$ 中查找其首次出现位置并返回这些位置。没有匹配到的元素位置返回 `NA_integer_` (整数型缺失值)。如

```
match(5, c(1,5,2))
```

```
## [1] 2
```

```
match(5, c(1,5,2,5))
```

```
## [1] 2
```

```
match(c(2,5), c(1,5,2,5))
```

```
## [1] 3 2
```

```
match(c(2,5,0), c(1,5,2,5))
```

```
## [1] 3 2 NA
```

用 `intersect(x,y)` 求交集，结果中不含重复元素，如

```
intersect(c(5, 7), c(1, 5, 2, 5))
```

```
## [1] 5
```

用 `union(x,y)` 求并集，结果中不含重复元素，如

```
union(c(5, 7), c(1, 5, 2, 5))
```

```
## [1] 5 7 1 2
```

用 `setdiff(x,y)` 求差集，即 `x` 的元素中不属于 `y` 的元素组成的集合，结果中不含重复元素，如

```
setdiff(c(5, 7), c(1, 5, 2, 5))
```

```
## [1] 7
```

用 `setequal(x,y)` 判断两个集合是否相等，不受次序与重复元素的影响，如

```
setequal(c(1,5,2), c(2,5,1))
```

```
## [1] TRUE
```

```
setequal(c(1,5,2), c(2,5,1,5))
```

```
## [1] TRUE
```

7.10 练习

设文件 `class.csv` 内容如下:

```
name,sex,age,height,weight
Alice,F,13,56.5,84
Becka,F,13,65.3,98
Gail,F,14,64.3,90
Karen,F,12,56.3,77
Kathy,F,12,59.8,84.5
Mary,F,15,66.5,112
Sandy,F,11,51.3,50.5
Sharon,F,15,62.5,112.5
Tammy,F,14,62.8,102.5
Alfred,M,14,69,112.5
Duke,M,14,63.5,102.5
```

```
Guido,M,15,67,133
James,M,12,57.3,83
Jeffrey,M,13,62.5,84
John,M,12,59,99.5
Philip,M,16,72,150
Robert,M,12,64.8,128
Thomas,M,11,57.5,85
William,M,15,66.5,112
```

用如下程序可以把上述文件读入为 R 数据框 `d.class`, 并取出其中的 `name` 和 `age` 列到变量 `name` 和 `age` 中:

```
d.class <- read.csv('class.csv', header=TRUE, stringsAsFactors=FALSE)
name <- d.class[, 'name']
age <- d.class[, 'age']
```

- (1) 求出 `age` 中第 3, 5, 7 号的值;
- (2) 用变量 `age`, 求出达到 15 岁及以上的那些值;
- (3) 用变量 `name` 和 `age`, 求出 Mary 与 James 的年龄。
- (4) 求 `age` 中除 Mary 与 James 这两人之外的那些人的年龄值, 保存到变量 `age1` 中。
- (5) 假设向量 `x` 长度为 `n`, 其元素是 $\{1, 2, \dots, n\}$ 的一个重排。可以把 `x` 看成一个 `i` 到 `x[i]` 的映射 (`i` 在 $\{1, 2, \dots, n\}$ 中取值)。求向量 `y`, 保存了上述映射的逆映射, 即: 如果 `x[i]=j`, 则 `y[j]=i`。

Chapter 8

R 数据类型的性质

8.1 存储模式与基本类型

R 的变量可以存储多种不同的数据类型，可以用 `typeof()` 函数来返回一个变量或表达式的类型。比如

```
typeof(1:3)
```

```
## [1] "integer"
```

```
typeof(c(1,2,3))
```

```
## [1] "double"
```

```
typeof(c(1, 2.1, 3))
```

```
## [1] "double"
```

```
typeof(c(TRUE, NA, FALSE))
```

```
## [1] "logical"
```

```
typeof('Abc')
```

```
## [1] "character"
```

```
typeof(factor(c('F', 'M', 'M', 'F')))
```

```
## [1] "integer"
```

注意因子的结果是 `integer` 而不是因子。函数 `mode()` 和 `storage.mode()` 以及 `typeof()` 类似, 但返回结果有差别。这三个函数都是与存储类型有关, 不依赖于变量和表达式的实际作用。

R 中数据的最基本的类型包括 `logical`, `integer`, `double`, `character`, `complex`, `raw`, 其它数据类型都是由基本类型组合或转变得到的。`character` 类型就是字符串类型, `raw` 类型是直接使用其二进制内容的类型。为了判断某个向量 `x` 保存的基本类型, 可以用 `is.xxx()` 类函数, 如 `is.integer(x)`, `is.double(x)`, `is.numeric(x)`, `is.logical(x)`, `is.character(x)`, `is.complex(x)`, `is.raw(x)`。其中 `is.numeric(x)` 对 `integer` 和 `double` 内容都返回真值。

在 R 语言中数值一般看作 `double`, 如果需要明确表明某些数值是整数, 可以在数值后面附加字母 `L`, 如

```
is.integer(c(1, -3))
```

```
## [1] FALSE
```

```
is.integer(c(1L, -3L))
```

```
## [1] TRUE
```

整数型的缺失值是 `NA`, 而 `double` 型的特殊值除了 `NA` 外, 还包括 `Inf`, `-Inf` 和 `NaN`, 其中 `NaN` 也算是缺失值, `Inf` 和 `-Inf` 不算是缺失值。如:

```
c(-1, 0, 1)/0
```

```
## [1] -Inf NaN Inf
```

```
is.na(c(-1, 0, 1)/0)
```

```
## [1] FALSE TRUE FALSE
```

对 `double` 类型, 可以用 `is.finite()` 判断是否有限值, `NA`, `Inf`, `-Inf` 和 `NaN` 都不是有限值; 用 `is.infinite()` 判断是否 `Inf` 或 `-Inf`; `is.na()` 判断是否 `NA` 或 `NaN`; `is.nan()` 判断是否 `NaN`。

严格说来, NA 表示逻辑型缺失值, 但是当作其它类型缺失值时一般能自动识别。NA_integer_ 是整数型缺失值, NA_real_ 是 double 型缺失值, NA_character_ 是字符型缺失值。

在 R 的向量类型中, integer 类型、double 类型、logical 类型、character 类型、还有 complex 类型和 raw 类型称为原子类型 (atomic types), 原子类型的向量中元素都是同一基本类型的。比如, double 型向量的元素都是 double 或者缺失值。除了原子类型的向量, 在 R 语言的定义中, 向量还包括后面要讲到的列表 (list), 列表的元素不需要属于相同的基本类型, 而且列表的元素可以不是单一基本类型元素。用 typeof() 函数可以返回向量的类型, 列表返回结果为 "list":

```
typeof(list("a", 1L, 1.5))
```

```
## [1] "list"
```

原子类型的各个元素除了基本类型相同, 还不包含任何嵌套结构, 如:

```
c(1, c(2,3, c(4,5)))
```

```
## [1] 1 2 3 4 5
```

8.2 类属

R 具有一定的面向对象语言特征, 其数据类型有一个 class 属性, 函数 class() 可以返回变量类型的类属, 比如

```
typeof(factor(c('F', 'M', 'M', 'F')))
```

```
## [1] "integer"
```

```
mode(factor(c('F', 'M', 'M', 'F')))
```

```
## [1] "numeric"
```

```
storage.mode(factor(c('F', 'M', 'M', 'F')))
```

```
## [1] "integer"
```

```
class(factor(c('F', 'M', 'M', 'F')))
```

```
## [1] "factor"
```

```
class(as.numeric(factor(c('F', 'M', 'M', 'F'))))
```

```
## [1] "numeric"
```

R 有一个特殊的 NULL 类型，这个类型只有唯一的一个 NULL 值，表示不存在。要把 NULL 与 NA 区分开来，NA 是有类型的（integer、double、logical、character 等），NA 表示存在但是未知。用 `is.null()` 函数判断某个变量是否取 NULL。

8.3 类型转换

可以用 `as.xxx()` 类的函数在不同类型之间进行强制转换。如

```
as.numeric(c(FALSE, TRUE))
```

```
## [1] 0 1
```

```
as.character(sqrt(1:4))
```

```
## [1] "1" "1.4142135623731" "1.73205080756888"
```

```
## [4] "2"
```

类型转换也可能是隐含的，比如，四则运算中数值会被统一转换为 double 类型，逻辑运算中运算元素会被统一转换为 logical 类型。逻辑值转换成数值时，TRUE 转换成 1，FALSE 转换成 0。

在用 `c()` 函数合并若干元素时，如果元素基本类型不同，将统一转换成最复杂的一个，复杂程度从简单到复杂依次为：`logical<integer<double<character`。如

```
c(FALSE, 1L, 2.5, "3.6")
```

```
## [1] "FALSE" "1" "2.5" "3.6"
```

不同类型参与要求类型相同的运算时，也会统一转换为最复杂的类型，如：


```
TRUE + 10
```

```
## [1] 11
```

```
paste("abc", 1)
```

```
## [1] "abc 1"
```

8.4 属性

除了 NULL 以外，R 的变量都可以看成是对象，都可以有属性。在 R 语言中，属性是把变量看成对象后，除了其存储内容（如元素）之外的其它附加信息，如维数、类属等。对象 `x` 的所有属性可以用 `attributes()` 读取，如

```
x <- table(c(1,2,1,3,2,1)); print(x)
```

```
##
```

```
## 1 2 3
```

```
## 3 2 1
```

```
attributes(x)
```

```
## $dim
```

```
## [1] 3
```

```
##
```

```
## $dimnames
```

```
## $dimnames[[1]]
```

```
## [1] "1" "2" "3"
```

```
##
```

```
##
```

```
## $class
```

```
## [1] "table"
```

`table()` 函数用了输出其自变量中每个不同值的出现次数，称为频数。从上例可以看出，`table()` 函数的结果有三个属性，前两个是 `dim` 和 `dimnames`，这是数组 (array) 具有的属性；另一个是 `class` 属性，值为 `"table"`。因为 `x` 是数组，可以访问如

```
x[1]

## 1
## 3
x["3"]

## 3
## 1
```

也可以用 `attributes()` 函数修改属性，如

```
attributes(x) <- NULL
x
```

```
## [1] 3 2 1
```

如上修改后 `x` 不再是数组，也不是 `table`。

`class` 属性是特殊的。如果一个对象具有 `class` 属性，某些所谓“通用函数 (generic functions)”会针对这样的对象进行专门的操作，比如，`print()` 函数在显示向量和回归结果时采用完全不同的格式。这在其它程序设计语言中称为“重载”(overloading)。

可以用 `attr(对象, "属性名")` 读取和修改单个属性。向量的元素名是 `names` 属性，例如

```
ages <- c("李明"=30, "张聪"=25, "刘颖"=28)
names(ages)
```

```
## [1] "李明" "张聪" "刘颖"
```

```
attr(ages, "names")
```

```
## [1] "李明" "张聪" "刘颖"
```

```
attr(ages, "names") <- NULL
ages
```

```
## [1] 30 25 28
```

还可以用 `unname()` 函数返回一个去掉了 `names` 属性的副本。

8.5 str() 函数

用 `print()` 函数可以显示对象内容。如果内容很多，显示行数可能也很多。用 `str()` 函数可以显示对象的类型和主要结构及典型内容。例如

```
s <- 101:200
attr(s,'author') <- ' 李小明'
attr(s,'date') <- '2016-09-12'
str(s)

## int [1:100] 101 102 103 104 105 106 107 108 109 110 ...
## - attr(*, "author")= chr "李小明"
## - attr(*, "date")= chr "2016-09-12"
```

8.6 关于赋值

要注意的是，在 R 中赋值本质上是把一个存储的对象与一个变量名联系在一起 (binding)，多个变量名可以指向同一个对象。对于基本的数据类型如数值型向量，两个指向相同对象的变量当一个变量被修改时自动制作副本，如

```
x <- 1:5
y <- x
y[3] <- 0
x

## [1] 1 2 3 4 5
y

## [1] 1 2 0 4 5
```

这里如果 `y` 没有与其它变量指向同一对象，则修改时直接修改该对象而不制作副本。

但是对于有些比较复杂的类型，两个指向同一对象的变量是同步修改的。这样的类型的典型代表是闭包 (closure)，它带有一个环境，环境的内容是不自动制作副本的。

Chapter 9

R 日期时间

9.1 R 日期和日期时间类型

R 日期可以保存为 Date 类型，一般用整数保存，数值为从 1970-1-1 经过的天数。

R 中用一种叫做 POSIXct 和 POSIXlt 的特殊数据类型保存日期和时间，可以仅包含日期部分，也可以同时有日期和时间。技术上，POSIXct 把日期时间保存为从 1970 年 1 月 1 日零时到该日期时间的时间间隔秒数，所以数据框中需要保存日期时用 POSIXct 比较合适，需要显示时再转换成字符串形式；POSIXlt 把日期时间保存为一个包含年、月、日、星期、时、分、秒等成分列表，所以求这些成分可以从 POSIXlt 格式日期的列表变量中获得。日期时间会涉及到所在时区、夏时制等问题，比较复杂。

基础的 R 用 `as.Date()`、`as.POSIXct()` 等函数生成日期型和日期时间型，R 扩展包 `lubridate` 提供了多个方便函数，可以更容易地生成、转换、管理日期型和日期时间型数据。

```
library(lubridate)
```

```
##
```

```
## 载入程辑包: 'lubridate'
```

```
## The following object is masked from 'package:base':
```

```
##
##      date
```

9.2 从字符串生成日期数据

函数 `lubridate::today()` 返回当前日期:

```
today()
```

```
## [1] "2019-08-29"
```

函数 `lubridate::now()` 返回当前日期时间:

```
now()
```

```
## [1] "2019-08-29 21:52:01 CST"
```

结果显示中出现的 `CST` 是时区, 这里使用了操作系统提供的当前时区。`CST` 不是一个含义清晰的时区, 在不同国家对应不同的时区, 在中国代表中国标准时间 (北京时间)。

用 `lubridate::ymd()`, `lubridate::mdy()`, `lubridate::dmy()` 将字符型向量转换为日期型向量, 如:

```
ymd(c("1998-3-10", "2018-01-17", "18-1-17"))
```

```
## [1] "1998-03-10" "2018-01-17" "2018-01-17"
```

```
mdy(c("3-10-1998", "01-17-2018"))
```

```
## [1] "1998-03-10" "2018-01-17"
```

```
dmy(c("10-3-1998", "17-01-2018"))
```

```
## [1] "1998-03-10" "2018-01-17"
```

在年号只有两位数字时, 默认对应到 1969-2068 范围。

`lubridate::make_date(year, month, day)` 可以从三个数值构成日期向量。如

```
make_date(1998, 3, 10)
```

```
## [1] "1998-03-10"
```

lubridate 包的 `ymd`、`mdy`、`dmy` 等函数添加 `hms`、`hm`、`h` 等后缀，可以用于将字符串转换成日期时间。如

```
ymd_hms("1998-03-16 13:15:45")
```

```
## [1] "1998-03-16 13:15:45 UTC"
```

结果显示中 UTC 是时区，UTC 是协调世界时 (Universal Time Coordinated) 英文缩写，是由国际无线电咨询委员会规定和推荐，并由国际时间局 (BIH) 负责保持的以秒为基础的时间标度。UTC 相当于本初子午线 (即经度 0 度) 上的平均太阳时，过去曾用格林威治平均时 (GMT) 来表示。北京时间比 UTC 时间早 8 小时，以 1999 年 1 月 1 日 0000UTC 为例，UTC 时间是零点，北京时间为 1999 年 1 月 1 日早上 8 点整。

在 `Date()`、`as.DateTime()`、`ymd()` 等函数中，可以用 `tz=` 指定时区，比如北京时间可指定为 `tz="Etc/GMT+8"` 或 `tz="Asia/Shanghai"`。

`lubridate::make_datetime(year, month, day, hour, min, sec)` 可以从最多六个数值组成日期时间，其中时分秒缺省值都是 0。如

```
make_datetime(1998, 3, 16, 13, 15, 45.2)
```

```
## [1] "1998-03-16 13:15:45 UTC"
```

用 `lubridate::as_date()` 可以将日期时间型转换为日期型，如

```
as_date(as.POSIXct("1998-03-16 13:15:45"))
```

```
## [1] "1998-03-16"
```

用 `lubridate::as_datetime()` 可以将日期型数据转换为日期时间型，如

```
as_datetime(as.Date("1998-03-16"))
```

```
## [1] "1998-03-16 UTC"
```

9.3 日期显示格式

用 `as.character()` 函数把日期型数据转换为字符型, 如

```
x <- as.POSIXct(c('1998-03-16', '2015-11-22'))
as.character(x)
```

```
## [1] "1998-03-16" "2015-11-22"
```

在 `as.character()` 中可以用 `format` 选项指定显示格式, 如

```
as.character(x, format='%m/%d/%Y')
```

```
## [1] "03/16/1998" "11/22/2015"
```

格式中 “%Y” 代表四位的公元年号, “%m” 代表两位的月份数字, “%d” 代表两位的月内日期号。

“15Mar98” 这样的日期在英文环境中比较常见, 但是在 R 中的处理比较复杂。在下面的例子中, R 日期被转换成了类似 “Mar98” 这样的格式, 在 `format` 选项中用了 “%b” 代表三英文字母月份缩写, 但是因为月份缩写依赖于操作系统默认语言环境, 需要用 `Sys.setlocale()` 函数设置语言环境为 “C”。示例程序如下

```
x <- as.POSIXct(c('1998-03-16', '2015-11-22'))
old.lctime <- Sys.getlocale('LC_TIME')
Sys.setlocale('LC_TIME', 'C')
```

```
## [1] "C"
```

```
as.character(x, format='%b%y')
```

```
## [1] "Mar98" "Nov15"
```

```
Sys.setlocale('LC_TIME', old.lctime)
```

```
## [1] "Chinese (Simplified)_China.936"
```

`format` 选项中的 “%y” 表示两位数的年份, 应尽量避免使用两位数年份以避免混淆。

包含时间的转换如


```
x <- as.POSIXct('1998-03-16 13:15:45')
as.character(x)
```

```
## [1] "1998-03-16 13:15:45"
```

```
as.character(x, format='%H:%M:%S')
```

```
## [1] "13:15:45"
```

这里“%H”代表小时（按 24 小时制），“%M”代表两位的分钟数字，“%S”代表两位的秒数。

9.4 访问日期时间的组成值

lubridate 包的如下函数可以取出日期型或日期时间型数据中的组成部分：

- year() 取出年
- month() 取出月份数值
- mday() 取出日数值
- yday() 取出日期在一年中的序号，元旦为 1
- wday() 取出日期在一个星期内的序号，但是一个星期从星期天开始，星期日为 1, 星期一为 2, 星期六为 7。
- hour() 取出小时
- minute() 取出分钟
- second() 取出秒

比如, 2018-1-17 是星期三, 则

```
month(as.POSIXct("2018-1-17 13:15:40"))
```

```
## [1] 1
```

```
mday(as.POSIXct("2018-1-17 13:15:40"))
```

```
## [1] 17
```

```
wday(as.POSIXct("2018-1-17 13:15:40"))
```

```
## [1] 4
```

lubridate 的这些成分函数还允许被赋值，结果就修改了相应元素的值，如

```
x <- as.POSIXct("2018-1-17 13:15:40")
year(x) <- 2000
month(x) <- 1
mday(x) <- 1
x
```

```
## [1] "2000-01-01 13:15:40 CST"
```

update() 可以对一个日期或一个日期型向量统一修改其组成部分的值，如

```
x <- as.POSIXct("2018-1-17 13:15:40")
y <- update(x, year=2000)
y
```

```
## [1] "2000-01-17 13:15:40 CST"
```

update() 函数中可以用 year, month, mday, hour, minute, second 等参数修改日期的组成部分。

9.5 日期舍入计算

lubridate 包提供了 floor_date(), round_date(), ceiling_date() 等函数，对日期可以用 unit= 指定一个时间单位进行舍入。时间单位为字符串，如 seconds, 5 seconds, minutes, 2 minutes, hours, days, weeks, months, years 等。

比如，以 10 minutes 为单位，floor_date() 将时间向前归一化到 10 分钟的整数倍，ceiling_date() 将时间向后归一化到 10 分钟的整数倍，round_date() 将时间归一化到最近的 10 分钟的整数倍，时间恰好是 5 分钟倍数时按照类似四舍五入的原则向上取整。例如

```
x <- ymd_hms("2018-01-11 08:32:44")
floor_date(x, unit="10 minutes")
```

```
## [1] "2018-01-11 08:30:00 UTC"
```

```
ceiling_date(x, unit="10 minutes")
```

```
## [1] "2018-01-11 08:40:00 UTC"
```

```
round_date(x, unit="10 minutes")
```

```
## [1] "2018-01-11 08:30:00 UTC"
```

如果单位是星期，会涉及到一个星期周期的开始是星期日还是星期一的问题。用参数 `week_start=7` 指定开始是星期日，`week_start=1` 指定开始是星期一。

9.6 日期计算

在 `lubridate` 的支持下日期可以相减，可以进行加法、除法。`lubridate` 包提供了如下的三种与时间长短有关的数据类型：

- 时间长度 (`duration`)，按整秒计算
- 时间周期 (`period`)，如日、周
- 时间区间 (`interval`)，包括一个开始时间和一个结束时间

9.6.1 时间长度

R 的 `POSIXct` 日期时间之间可以相减，如

```
d1 <- ymd_hms("2000-01-01 0:0:0")
d2 <- ymd_hms("2000-01-02 12:0:5")
di <- d2 - d1; di
```

```
## Time difference of 1.500058 days
```

结果显示与日期之间差别大小有关系，结果是类型是 `difftime`。

`lubridate` 包提供了 `duration` 类型，处理更方便：

```
as.duration(di)
```

```
## [1] "129605s (~1.5 days)"
```

lubridate 的 `dseconds()`, `dminutes()`, `dhours()`, `ddays()`, `dweeks()`, `dyears()` 函数可以直接生成时间长度类型的数据, 如

```
dhours(1)
```

```
## [1] "3600s (~1 hours)"
```

lubridate 的时间长度类型总是以秒作为单位, 可以在时间长度之间相加, 也可以对时间长度乘以无量纲数, 如

```
dhours(1) + dseconds(5)
```

```
## [1] "3605s (~1 hours)"
```

```
dhours(1)*10
```

```
## [1] "36000s (~10 hours)"
```

可以给一个日期加或者减去一个时间长度, 结果严格按推移的秒数计算, 如

```
d2 <- ymd_hms("2000-01-02 12:0:5")
```

```
d2 - dhours(5)
```

```
## [1] "2000-01-02 07:00:05 UTC"
```

```
d2 + ddays(10)
```

```
## [1] "2000-01-12 12:00:05 UTC"
```

时间的前后推移在涉及到时区和夏时制时有可能出现未预料到的情况。

用 `unclass()` 函数将时间长度数据的类型转换为普通数值, 如:

```
unclass(dhours(1))
```

```
## [1] 3600
```

9.6.2 时间周期

时间长度的固定单位是秒, 但是像月、年这样的单位, 因为可能有不同的天数, 所以日历中的时间单位往往没有固定的时长。

lubridate 包的 `seconds()`, `minutes()`, `hours()`, `days()`, `weeks()`, `years()` 函数可以生成以日历中正常的周期为单位的时间长度, 不需要与秒数相联系, 可以用于时间的前后推移。这些时间周期的结果可以相加、乘以无量纲整数:

```
years(2) + 10*days(1)
```

```
## [1] "2y 0m 10d 0H 0M 0S"
```

lubridate 的月度周期因为与已有函数名冲突, 所以没有提供, 需要使用 `lubridate::period(num, units="month")` 的格式, 其中 `num` 是几个月的数值。

为了按照日历进行日期的前后平移, 而不是按照秒数进行日期的前后平移, 应该使用这些时间周期。例如, 因为 2016 年是闰年, 按秒数给 2016-01-01 加一年, 得到的并不是 2017-01-01:

```
ymd("2016-01-01") + dyears(1)
```

```
## [1] "2016-12-31"
```

使用时间周期函数则得到预期结果:

```
ymd("2016-01-01") + years(1)
```

```
## [1] "2017-01-01"
```

9.6.3 时间区间

lubridate 提供了 `%--%` 运算符构造一个时间期间。时间区间可以求交集、并集等。

构造如:

```
d1 <- ymd_hms("2000-01-01 0:0:0")
d2 <- ymd_hms("2000-01-02 12:0:5")
din <- (d1 %--% d2); din
```

```
## [1] 2000-01-01 UTC--2000-01-02 12:00:05 UTC
```

对一个时间区间可以用除法计算其时间长度, 如

```
din / ddays(1)
```

```
## [1] 1.500058
```

```
din / dseconds(1)
```

```
## [1] 129605
```

生成时间区间，也可以用 `lubridate::interval(start, end)` 函数，如

```
interval(ymd_hms("2000-01-01 0:0:0"), ymd_hms("2000-01-02 12:0:5"))
```

```
## [1] 2000-01-01 UTC--2000-01-02 12:00:05 UTC
```

可以指定时间长度和开始日期生成时间区间，如

```
d1 <- ymd("2018-01-15")
```

```
din <- as.interval(dweeks(1), start=d1); din
```

```
## [1] 2018-01-15 UTC--2018-01-22 UTC
```

注意这个时间区间表面上涉及到 8 个日期，但是实际长度还是只有 7 天，因为每一天的具体时间都是按零时计算，所以区间末尾的那一天实际不含在内。

用 `lubridate::int_start()` 和 `lubridate::int_end()` 函数访问时间区间的端点，如：

```
int_start(din)
```

```
## [1] "2018-01-15 UTC"
```

```
int_end(din)
```

```
## [1] "2018-01-22 UTC"
```

可以用 `as.duration()` 将一个时间区间转换成时间长度，用 `as.period()` 将一个时间区间转换为可变时长的时间周期个数。

用 `lubridate::int_shift()` 平移一个时间区间，如

```
din2 <- int_shift(din, by=ddays(3)); din2
```

```
## [1] 2018-01-18 UTC--2018-01-25 UTC
```

用 `lubridate::int_overlaps()` 判断两个时间区间是否有共同部分，如

```
int_overlaps(din, din2)
```

```
## [1] TRUE
```

时间区间允许开始时间比结束时间晚，用 `lubridate::int_standardize()` 可以将时间区间标准化成开始时间小于等于结束时间。`lubridate()` 现在没有提供求交集的功能，一个自定义求交集的函数如下：

```
int_intersect <- function(int1, int2){
  n <- length(int1)
  int1 <- lubridate::int_standardize(int1)
  int2 <- lubridate::int_standardize(int2)
  sele <- lubridate::int_overlaps(int1, int2)
  inter <- rep(lubridate::interval(NA, NA), n)
  if(any(sele)){
    inter[sele] <-
      lubridate::interval(pmax(lubridate::int_start(int1[sele]),
                               lubridate::int_start(int2[sele])),
                          pmin(lubridate::int_end(int1[sele]),
                               lubridate::int_end(int2[sele])))
  }
  inter
}
```

测试如：

```
d1 <- ymd(c("2018-01-15", "2018-01-18", "2018-01-25"))
d2 <- ymd(c("2018-01-21", "2018-01-23", "2018-01-30"))
din <- interval(d1, d2); din
```

```
## [1] 2018-01-15 UTC--2018-01-21 UTC 2018-01-18 UTC--2018-01-23 UTC
```

```
## [3] 2018-01-25 UTC--2018-01-30 UTC
```

```
int_intersect(rep(din[1], 2), din[2:3])
```

```
## [1] 2018-01-18 UTC--2018-01-21 UTC NA--NA
```

此自定义函数还可以进一步改成允许两个自变量长度不等的情形。

9.7 基本 R 软件的日期功能

9.7.1 生成日期和日期时间型数据

对 yyyy-mm-dd 或 yyyy/mm/dd 格式的数据，可以直接用 `as.Date()` 转换为 Date 类型，如：

```
x <- as.Date("1970-1-5"); x
```

```
## [1] "1970-01-05"
```

```
as.numeric(x)
```

```
## [1] 4
```

`as.Date()` 可以将多个日期字符串转换成 Date 类型，如

```
as.Date(c("1970-1-5", "2017-9-12"))
```

```
## [1] "1970-01-05" "2017-09-12"
```

对于非标准的格式，在 `as.Date()` 中可以增加一个 `format` 选项，其中用 %Y 表示四位数字的年，%m 表示月份数字，%d 表示日数字。如

```
as.Date("1/5/1970", format="%m/%d/%Y")
```

```
## [1] "1970-01-05"
```

用 `as.POSIXct()` 函数把年月日格式的日期转换为 R 的标准日期，没有时间部分就认为时间在午夜。如

```
as.POSIXct(c('1998-03-16'))
```

```
## [1] "1998-03-16 CST"
```

```
as.POSIXct(c('1998/03/16'))
```

```
## [1] "1998-03-16 CST"
```

年月日中间的分隔符可以用减号也可以用正斜杠，但不能同时有减号又有斜杠。

待转换的日期时间字符串，可以是年月日之后隔一个空格以“时:分:秒”格式带有时间。如

```
as.POSIXct('1998-03-16 13:15:45')
```

```
## [1] "1998-03-16 13:15:45 CST"
```

用 `as.POSIXct()` 可以同时转换多项日期时间，如

```
as.POSIXct(c('1998-03-16 13:15:45', '2015-11-22 9:45:3'))
```

```
## [1] "1998-03-16 13:15:45 CST" "2015-11-22 09:45:03 CST"
```

转换后的日期变量有 `class` 属性，取值为 `POSIXct` 与 `POSIXt`，并带有一个 `tzzone`（时区）属性。

```
x <- as.POSIXct(c('1998-03-16 13:15:45', '2015-11-22 9:45:3'))
attributes(x)
```

```
## $class
```

```
## [1] "POSIXct" "POSIXt"
```

```
##
```

```
## $tzzone
```

```
## [1] ""
```

在 `as.POSIXct()` 函数中用 `format` 参数指定一个日期格式。如

```
as.POSIXct('3/13/15', format='%m/%d/%y')
```

```
## [1] "2015-03-13 CST"
```

如果日期仅有年和月，必须添加日（添加 01 为日即可）才能读入。比如用 `'1991-12'` 表示 1991 年 12 月，则如下程序将其读入为 `'1991-12-01'`：

```
as.POSIXct(paste('1991-12', '-01', sep=''), format='%Y-%m-%d')
```

```
## [1] "1991-12-01 CST"
```

又如

```
old.lctime <- Sys.getlocale('LC_TIME')
```

```
Sys.setlocale('LC_TIME', 'C')
```

```
## [1] "C"
```

```
as.POSIXct(paste('01', 'DEC91', sep=' '), format='%d%b%y')
```

```
## [1] "1991-12-01 CST"
```

```
Sys.setlocale('LC_TIME', old.lctime)
```

```
## [1] "Chinese (Simplified)_China.936"
```

把 'DEC91' 转换成了 '1991-12-01'。

如果明确地知道时区，在 `as.POSIXct()` 和 `as.POSIXlt()` 中可以加选项 `tz=` 字符串。选项 `tz` 的缺省值为空字符串，这一般对应于当前操作系统的默认时区。但是，有些操作系统和 R 版本不能使用默认值，这时可以为 `tz` 指定时区，比如北京时间可指定为 `tz='Etc/GMT+8'`。如

```
as.POSIXct('1949-10-01', tz='Etc/GMT+8')
```

```
## [1] "1949-10-01 -08"
```

9.7.2 取出日期时间的组成值

把一个 R 日期时间值用 `as.POSIXlt()` 转换为 `POSIXlt` 类型，就可以用列表元素方法取出其组成的年、月、日、时、分、秒等数值。如

```
x <- as.POSIXct('1998-03-16 13:15:45')
y <- as.POSIXlt(x)
cat(1900+y$year, y$mon+1, y$mday, y$hour, y$min, y$sec, '\n')
```

```
## 1998 3 16 13 15 45
```

注意 `year` 要加 1900，`mon` 要加 1。另外，列表元素 `wday` 取值 1-6 时表示星期一到星期六，取值 0 时表示星期天。

对多个日期，`as.POSIXlt()` 会把它们转换成一个列表（列表类型稍后讲述），这时可以用列表元素 `year`, `mon`, `mday` 等取出日期成分。如

```
x <- as.POSIXct(c('1998-03-16', '2015-11-22'))
as.POSIXlt(x)$year + 1900
```

```
## [1] 1998 2015
```

9.7.3 日期计算

因为 Date 类型是用数值保存的，所以可以给日期加减一个整数，如：

```
x <- as.Date("1970-1-5")
x1 <- x + 10; x1
```

```
## [1] "1970-01-15"
```

```
x2 <- x - 5; x2
```

```
## [1] "1969-12-31"
```

所有的比较运算都适用于日期类型。

可以给一个日期加减一定的秒数，如

```
as.POSIXct(c('1998-03-16 13:15:45')) - 30
```

```
## [1] "1998-03-16 13:15:15 CST"
```

```
as.POSIXct(c('1998-03-16 13:15:45')) + 10
```

```
## [1] "1998-03-16 13:15:55 CST"
```

但是两个日期不能相加。

给一个日期加减一定天数，可以通过加减秒数实现，如

```
as.POSIXct(c('1998-03-16 13:15:45')) + 3600*24*2
```

```
## [1] "1998-03-18 13:15:45 CST"
```

这个例子把日期推后了两天。

用 `difftime(time1, time2, units='days')` 计算 time1 减去 time2 的天数，如

```
x <- as.POSIXct(c('1998-03-16', '2015-11-22'))
c(difftime(x[2], x[1], units='days'))
```

```
## Time difference of 6460 days
```

函数结果用 `c()` 包裹以转换为数值, 否则会带有单位。

调用 `difftime()` 时如果前两个自变量中含有时间部分, 则间隔天数也会带有小数部分。如

```
x <- as.POSIXct(c('1998-03-16 13:15:45', '2015-11-22 9:45:3'))
c(difftime(x[2], x[1], units='days'))
```

```
## Time difference of 6459.854 days
```

`difftime()` 中 `units` 选项还可以取为 `'secs'`, `'mins'`, `'hours'` 等。

9.8 练习

设文件 `dates.csv` 中包含如下内容:

```
"出生日期","发病日期"
"1941/3/8","2007/1/1"
"1972/1/24","2007/1/1"
"1932/6/1","2007/1/1"
"1947/5/17","2007/1/1"
"1943/3/10","2007/1/1"
"1940/1/8","2007/1/1"
"1947/8/5","2007/1/1"
"2005/4/14","2007/1/1"
"1961/6/23","2007/1/2"
"1949/1/10","2007/1/2"
```

把这个文件读入为 R 数据框 `dates.tab`, 运行如下程序定义 `date1` 和 `date2` 变量:

```
date1 <- dates.tab[, '出生日期']
date2 <- dates.tab[, '发病日期']
```

- (1) 把 `date1`、`date2` 转换为 R 的 `POSIXct` 日期型。
- (2) 求 `date1` 中的各个出生年。
- (3) 计算发病时的年龄, 以周岁论 (过生日才算)。

- (4) 把 `date2` 中发病年月转换为 `'monyy'` 格式，这里 `mon` 是如 `FEB` 这样英文三字母缩写，`yy` 是两数字的年份。
- (5) 对诸如 `'FEB91'`，`'OCT15'` 这样的年月数据，假设 `00—20` 表示 21 世纪年份，`21—99` 表示 20 实际年份。编写 R 函数，输入这样的字符型向量，返回相应的 `POSIXct` 格式日期，具体日期都取为相应月份的 1 号。这个习题和后两个习题可以预习函数部分来做。
- (6) 对 R 的 `POSIXct` 日期，写函数转换成 `'FEB91'`，`'OCT15'` 这样的年月表示，假设 `00—20` 表示 21 世纪年份，`21—99` 表示 20 实际年份。
- (7) 给定两个 `POSIXct` 日期向量 `birth` 和 `work`，`birth` 为生日，`work` 是入职日期，编写 R 函数，返回相应的入职周岁整数值（不到生日时周岁值要减一）。

Chapter 10

R 因子类型

10.1 因子

R 中用因子代表数据中分类变量, 如性别、省份、职业。有序因子代表有序量度, 如打分结果, 疾病严重程度等。

用 `factor()` 函数把字符型向量转换成因子, 如

```
x <- c("男", "女", "男", "男", "女")
sex <- factor(x)
sex
```

```
## [1] 男 女 男 男 女
## Levels: 男 女
```

```
attributes(sex)
```

```
## $levels
## [1] "男" "女"
##
## $class
## [1] "factor"
```

因子有 `class` 属性, 取值为"`factor`", 还有一个 `levels`(水平值) 属性, 此属

性可以用 `levels()` 函数访问，如

```
levels(sex)
```

```
## [1] "男" "女"
```

因子的 `levels` 属性可以看成是一个映射，把整数值 1,2,... 映射成这些水平值，因子在保存时会保存成整数值 1,2,... 等与水平值对应的编号。这样可以节省存储空间，在建模计算的程序中也比较有利于进行数学运算。

事实上，`read.csv()` 函数的默认操作会把输入文件的字符型列自动转换成因子，这对于性别、职业、地名这样的列是合适的，但是对于姓名、日期、详细地址这样的列则不合适。所以，在 `read.csv()` 调用中经常加选项 `stringsAsFactors=FALSE` 选项禁止这样的自动转换，还可以用 `colClasses` 选项逐个指定每列的类型。

用 `as.numeric()` 可以把因子转换为纯粹的整数值，如

```
as.numeric(sex)
```

```
## [1] 1 2 1 1 2
```

因为因子实际保存为整数值，所以对因子进行一些字符型操作可能导致错误。用 `as.character()` 可以把因子转换成原来的字符型，如

```
as.character(sex)
```

```
## [1] "男" "女" "男" "男" "女"
```

为了对因子执行字符型操作（如取子串），保险的做法是先用 `as.character()` 函数强制转换为字符型。

`factor()` 函数的一般形式为

```
factor(x, levels = sort(unique(x), na.last = TRUE),
      labels, exclude = NA, ordered = FALSE)
```

可以用选项 `levels` 自行指定各水平值，不指定时由 `x` 的不同值来求得。可以用选项 `labels` 指定各水平的标签，不指定时用各水平值的对应字符串。可以用 `exclude` 选项指定要转换为缺失值 (NA) 的元素值集合。如果指定了 `levels`，则当自变量 `x` 的某个元素等于第 `j` 个水平值时输出的因子对应元素

值取整数 j , 如果该元素值没有出现在 `levels` 中则输出的因子对应元素值取 NA。 `ordered` 取真值时表示因子水平是有次序的 (按编码次序)。

在使用 `factor()` 函数定义因子时, 如果知道自变量元素的所有可能取值, 应尽可能使用 `levels=` 参数指定这些不同可能取值, 这样, 即使某个取值没有出现, 此变量代表的含义和频数信息也是完整的。自己指定 `levels=` 的另一好处是可以按正确的次序显示因子的分类统计值。

因为一个因子的 `levels` 属性是该因子独有的, 所以合并两个因子有可能造成错误。如

```
li1 <- factor(c('男', '女'))
li2 <- factor(c('男', '男'))
c(li1, li2)
```

```
## [1] 1 2 1 1
```

结果不再是因子。正确的做法是

```
factor(c(as.character(li1), as.character(li2)))
```

```
## [1] 男 女 男 男
## Levels: 男 女
```

即恢复成字符型后合并, 然后再转换为因子。在合并两个数据框时也存在这样的问题。当然, 如果在定义 `li1` 和 `li2` 时都用了 `levels=c('男', '女')` 选项, `c(li1, li2)` 也能给出正确结果。

10.2 table() 函数

用 `table()` 函数统计因子各水平的出现次数 (称为频数或频率)。也可以对一般的向量统计每个不同元素的出现次数。如

```
table(sex)
```

```
## sex
## 男 女
## 3 2
```

对一个变量用 `table` 函数计数的结果是一个特殊的有元素名的向量，元素名是自变量的不同取值，结果的元素值是对应的频数。单个因子或单个向量的频数结果可以用向量的下标访问方法取出单个频数或若干个频数的子集。

10.3 `tapply()` 函数

可以按照因子分组然后每组计算另一变量的概括统计。如

```
h <- c(165, 170, 168, 172, 159)
tapply(h, sex, mean)
```

```
##      男      女
## 168.3333 164.5000
```

这里第一自变量 `h` 与第二自变量 `sex` 是等长的，对应元素分别为同一人的身高和性别，`tapply()` 函数分男女两组计算了身高平均值。

10.4 `forcats` 包的因子函数

```
library(forcats)
```

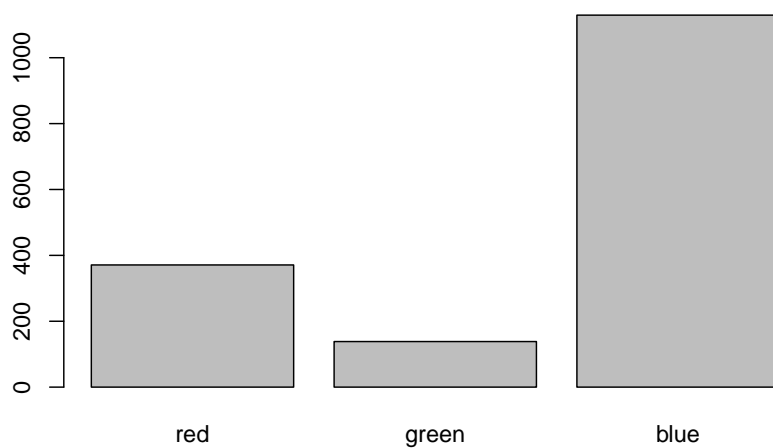
在分类变量类数较多时，往往需要对因子水平另外排序、合并等，`forcats` 包提供了一些针对因子的方便函数。

`forcats::fac_reorder()` 可以根据不同因子水平分成的组中另一数值型变量的统计量值排序。如：

```
set.seed(1)
fac <- sample(c("red", "green", "blue"), 30, replace=TRUE)
fac <- factor(fac, levels=c("red", "green", "blue"))
x <- round(100*(10+rt(30,2)))
res1 <- tapply(x, fac, sd); res1
```

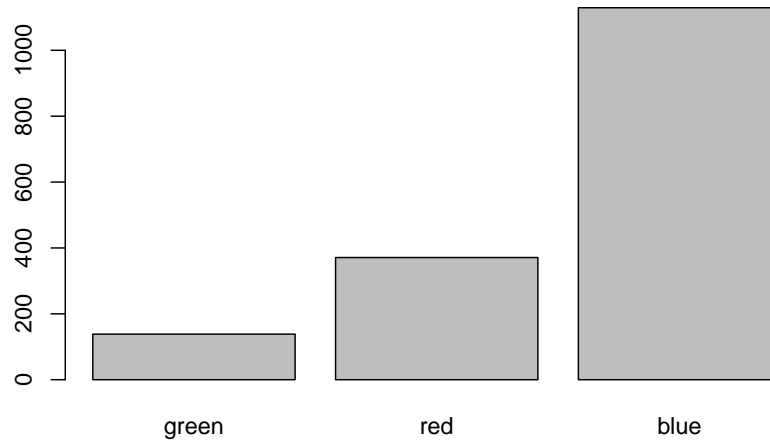
```
##      red      green      blue
## 370.9222 138.3185 1129.2587
```

```
barplot(res1)
```



如果希望按照统计量次序对因子排序，可以用 `forcats::fct_reorder()` 函数，如

```
fac2 <- fct_reorder(fac, x, sd)  
res2 <- tapply(x, fac2, sd)  
barplot(res2)
```



新的因子 `fac2` 的因子水平次序已经按照变量 `x` 的标准差从小到大排列。

有时在因子水平数较多时仅想将特定的一个或几个水平次序放到因子水平最前面，可以用 `forcats::fct_relevel()` 函数，如：

```
levels(fac)

## [1] "red" "green" "blue"

fac3 <- fct_relevel(fac, "blue"); levels(fac3)

## [1] "blue" "red" "green"
```

`fct_relevel()` 第一个参数是要修改次序的因子，后续可以有多个字符型参数表示要提前的水平。

`forcats::fct_reorder2(f, x, y)` 也调整因子 `f` 的水平次序，但是根据与每组中最大的 `x` 值相对应的 `y` 值大小调整次序，这样在作多个因子水平对应的曲线图时可以比较容易地区分多条曲线。

`forcats::fct_recode()` 可以修改每个水平的名称，如：

```
fac4 <- fct_recode(  
  fac,  
  " 红"="red", " 绿"="green", " 蓝"="blue")  
table(fac4)  
  
## fac4  
## 红 绿 蓝  
## 13 10 7
```

`fct_recode()` 在修改水平名时允许多个旧水平对应到一个新水平，从而合并原来的水平。如果合并很多，可以用 `fct_collapse()` 函数，如

```
compf <- fct_collapse(  
  comp,  
  " 其它"=c("", " 无名", " 无应答"),  
  " 联想"=c(" 联想", " 联想集团"),  
  " 百度"=c(" 百度", " 百度集团")  
)
```

如果某个因子频数少的水平很多，在统计时有过多水平不易展示主要的类别，可以用 `forcats::fct_lump(f)` 合并，缺省地从最少的类合并一直到“其它”类超过其它最小的类之前，可以用 `n=` 参数指定要保留多少个类。

10.5 练习

设文件 `class.csv` 中包含如下内容:

```
name,sex,age,height,weight  
Alice,F,13,56.5,84  
Becka,F,13,65.3,98  
Gail,F,14,64.3,90  
Karen,F,12,56.3,77  
Kathy,F,12,59.8,84.5  
Mary,F,15,66.5,112  
Sandy,F,11,51.3,50.5
```

```
Sharon,F,15,62.5,112.5
Tammy,F,14,62.8,102.5
Alfred,M,14,69,112.5
Duke,M,14,63.5,102.5
Guido,M,15,67,133
James,M,12,57.3,83
Jeffrey,M,13,62.5,84
John,M,12,59,99.5
Philip,M,16,72,150
Robert,M,12,64.8,128
Thomas,M,11,57.5,85
William,M,15,66.5,112
```

用如下程序把该文件读入为 R 数据框 `d.class`，其中的 `sex` 列已经自动转换为因子。取出其中的 `sex` 和 `age` 列到变量 `sex` 和 `age` 中

```
d.class <- read.csv('class.csv', header=TRUE)
sex <- d.class[, 'sex']
age <- d.class[, 'age']
```

- (1) 统计并显示列出 `sex` 的不同值频数；
- (2) 分男女两组分别求年龄最大值；
- (3) 把 `sex` 变量转换为一个新的因子，F 显示成 “Female”，M 显示成 “Male”。

Chapter 11

R 矩阵和数组

11.1 R 矩阵

矩阵用 `matrix` 函数定义，实际存储成一个向量，根据保存的行数和列数对应到矩阵的元素，存储次序为按列存储。定义如

```
A <- matrix(11:16, nrow=3, ncol=2); print(A)
```

```
##      [,1] [,2]
## [1,]  11  14
## [2,]  12  15
## [3,]  13  16
```

```
B <- matrix(c(1,-1, 1,1), nrow=2, ncol=2, byrow=TRUE); print(B)
```

```
##      [,1] [,2]
## [1,]   1  -1
## [2,]   1   1
```

`matrix()` 函数把矩阵元素以一个向量的形式输入，用 `nrow` 和 `ncol` 规定行数和列数，向量元素填入矩阵的缺省次序是按列填入，用 `byrow=TRUE` 选项可以转换成按行填入。

用 `nrow()` 和 `ncol()` 函数可以访问矩阵的行数和列数，如

```
nrow(A)
```

```
## [1] 3
```

```
ncol(A)
```

```
## [1] 2
```

矩阵有一个 `dim` 属性，内容是两个元素的向量，两个元素分别为矩阵的行数和列数。`dim` 属性可以用 `dim()` 函数访问。如

```
attributes(A)
```

```
## $dim
```

```
## [1] 3 2
```

```
dim(A)
```

```
## [1] 3 2
```

函数 `t(A)` 返回 `A` 的转置。

11.2 矩阵子集

用 `A[1,]` 取出 `A` 的第一行，变成一个普通向量。用 `A[,1]` 取出 `A` 的第一列，变成一个普通向量。用 `A[c(1,3),1:2]` 取出指定行、列对应的子矩阵。如

```
A
```

```
##      [,1] [,2]
```

```
## [1,]  11  14
```

```
## [2,]  12  15
```

```
## [3,]  13  16
```

```
A[1,]
```

```
## [1] 11 14
```

```
A[,1]
```

```
## [1] 11 12 13
```



```
A[c(1,3), 1:2]
```

```
##      [,1] [,2]
## [1,]  11  14
## [2,]  13  16
```

用 `colnames()` 函数可以给矩阵每列命名，也可以访问矩阵列名，用 `rownames()` 函数可以给矩阵每行命名，也可以访问矩阵行名。如

```
colnames(A) <- c('X', 'Y')
rownames(A) <- c('a', 'b', 'c')
```

```
A
```

```
##      X  Y
## a 11 14
## b 12 15
## c 13 16
```

矩阵可以有一个 `dimnames` 属性，此属性是两个元素的列表（列表见稍后部分的介绍），两个元素分别为矩阵的行名字符型向量与列名字符型向量。如果仅有其中之一，缺失的一个取为 `NULL`。

有了列名、行名后，矩阵下标可以用字符型向量，如

```
A[, 'Y']
```

```
## a b c
## 14 15 16
```

```
A['b', ]
```

```
## X Y
## 12 15
```

```
A[c('a', 'c'), 'Y']
```

```
## a c
## 14 16
```

注意在对矩阵取子集时，如果取出的子集仅有一行或仅有一列，结果就不再是

矩阵而是变成了 R 向量，R 向量既不是行向量也不是列向量。如果想避免这样的规则起作用，需要在方括号下标中加选项 `drop=FALSE`，如

```
A[,1,drop=FALSE]
```

```
##      X
## a 11
## b 12
## c 13
```

取出了 A 的第一列，作为列向量取出，所谓列向量实际是列数等于 1 的矩阵。如果用常量作为下标，其结果维数是确定的，不会出问题；如果用表达式作为下标，则表达式选出零个、一个、多个下标，结果维数会有不同，加 `drop=FALSE` 则是安全的做法。

矩阵也可以用逻辑下标取子集，比如

```
A
```

```
##      X Y
## a 11 14
## b 12 15
## c 13 16
```

```
A[A[,1]>=2,'Y']
```

```
##      a b c
## 14 15 16
```

矩阵本质上是一个向量添加了 `dim` 属性，实际保存还是保存成一个向量，其中元素的保存次序是按列填入，所以，也可以向对一个向量取子集那样，仅用一个正整数向量的矩阵取子集。如

```
A
```

```
##      X Y
## a 11 14
## b 12 15
## c 13 16
```

```
A[c(1,3,5)]
```

```
## [1] 11 13 15
```

为了挑选矩阵的任意元素组成的子集而不是子矩阵，可以用一个两列的矩阵作为下标，矩阵的每行的两个元素分别指定一个元素的行号和列号。如

```
ind <- matrix(c(1,1, 2,2, 3,2), ncol=2, byrow=TRUE)
```

```
A
```

```
##      X  Y
## a 11 14
## b 12 15
## c 13 16
```

```
ind
```

```
##      [,1] [,2]
## [1,]    1    1
## [2,]    2    2
## [3,]    3    2
```

```
A[ind]
```

```
## [1] 11 15 16
```

用 `c(A)` 或 `A[]` 返回矩阵 `A` 的所有元素。如果要修改矩阵 `A` 的所有元素，可以对 `A[]` 赋值。

对矩阵 `A`，`diag(A)` 访问 `A` 的主对角线元素组成的向量。另外，若 `x` 为正整数值标量，`diag(x)` 返回 `x` 阶单位阵；若 `x` 为长度大于 1 的向量，`diag(x)` 返回以 `x` 的元素为主对角线元素的对角矩阵。

11.3 `cbind()` 和 `rbind()` 函数

若 `x` 是向量，`cbind(x)` 把 `x` 变成列向量，即列数为 1 的矩阵，`rbind(x)` 把 `x` 变成行向量。

若 `x1`, `x2`, `x3` 是等长的向量，`cbind(x1, x2, x3)` 把它们看成列向量并在一

起组成一个矩阵。`cbind()` 的自变量可以同时包含向量与矩阵，向量的长度必须与矩阵行数相等。如

```
cbind(c(1,2), c(3,4), c(5,6))
```

```
##      [,1] [,2] [,3]
## [1,]    1    3    5
## [2,]    2    4    6
```

```
cbind(A, c(1,-1,10))
```

```
##      X Y
## a 11 14  1
## b 12 15 -1
## c 13 16 10
```

`cbind()` 的自变量中也允许有标量，这时此标量被重复使用。如

```
cbind(1, c(1,-1,10))
```

```
##      [,1] [,2]
## [1,]    1    1
## [2,]    1   -1
## [3,]    1   10
```

`rbind()` 用法类似，可以等长的向量看成行向量上下摞在一起，可以是矩阵与长度等于矩阵列数的向量上下摞在一起，向量长度为 1 也可以。

11.4 矩阵运算

11.4.1 四则运算

矩阵可以与标量作四则运算，结果为每个元素进行相应运算，如

```
A
```

```
##      X Y
## a 11 14
## b 12 15
```

```
## c 13 16  
C1 <- A + 2; C1
```

```
##      X  Y  
## a 13 16  
## b 14 17  
## c 15 18  
C2 <- A / 2; C2
```

```
##      X  Y  
## a 5.5 7.0  
## b 6.0 7.5  
## c 6.5 8.0
```

当运算为矩阵乘以一个标量时，就是线性代数中的矩阵的数乘运算。

两个同形状的矩阵进行加、减运算，即对应元素相加、相减，用 $A + B$, $A - B$ 表示，如

```
C1 + C2  
  
##      X  Y  
## a 18.5 23.0  
## b 20.0 24.5  
## c 21.5 26.0  
C1 - C2
```

```
##      X  Y  
## a 7.5  9.0  
## b 8.0  9.5  
## c 8.5 10.0
```

这就是线性代数中矩阵的加、减运算。

对两个同形状的矩阵，用 $*$ 表示两个矩阵对应元素相乘 (注意这不是线性代数中的矩阵乘法)，用 $/$ 表示两个矩阵对应元素相除。如

```
C1 * C2
```

```
##      X      Y
## a 71.5 112.0
## b 84.0 127.5
## c 97.5 144.0
```

```
C1 / C2
```

```
##           X           Y
## a 2.363636 2.285714
## b 2.333333 2.266667
## c 2.307692 2.250000
```

11.4.2 矩阵乘法

用`%%`表示矩阵乘法而不是用`*`表示，注意矩阵乘法要求左边的矩阵的列数等于右边的矩阵的行数。如

```
A
```

```
##      X  Y
## a 11 14
## b 12 15
## c 13 16
```

```
B
```

```
##      [,1] [,2]
## [1,]    1  -1
## [2,]    1   1
```

```
C3 <- A %% B; C3
```

```
##      [,1] [,2]
## a    25   3
## b    27   3
## c    29   3
```

11.4.3 向量与矩阵相乘

矩阵与向量进行乘法运算时，向量按需要解释成列向量或行向量。当向量左乘矩阵时，看成行向量；当向量右乘矩阵时，看成列向量。如

```
B
```

```
##      [,1] [,2]
## [1,]    1  -1
## [2,]    1   1
```

```
c(1,1) %*% B
```

```
##      [,1] [,2]
## [1,]    2   0
```

```
B %*% c(1,1)
```

```
##      [,1]
## [1,]    0
## [2,]    2
```

```
c(1,1) %*% B %*% c(1,1)
```

```
##      [,1]
## [1,]    2
```

注意矩阵乘法总是给出矩阵结果，即使此矩阵已经退化为行向量、列向量甚至于退化为标量也是一样。如果需要，可以用 `c()` 函数把一个矩阵转换成按列拉直的向量。

11.4.4 内积

设 x, y 是两个向量，计算向量内积，可以用 `sum(x*y)` 表示。

设 A, B 是两个矩阵， $A^T B$ 是广义的内积，也称为叉积 (crossprod)，结果是一个矩阵，元素为 A 的每列与 B 的每列计算内积的结果。 $A^T B$ 在 \mathbb{R} 中可以表示为 `crossprod(A, B)`， $A^T A$ 可以表示为 `crossprod(A)`。要注意的是，

`crossprod()` 的结果总是矩阵, 所以计算两个向量的内积用 `sum(x,y)` 而不用 `crossprod(x,y)`。

11.4.5 外积

R 向量支持外积运算, 记为`%o%`, 结果为矩阵。`x %o% y` 的第 i 行第 j 列元素等于 `x[i]` 乘以 `y[j]`。如

```
c(1,2,3) %o% c(1, -1)
```

```
##      [,1] [,2]
## [1,]   1  -1
## [2,]   2  -2
## [3,]   3  -3
```

这种运算还可以推广到 `x` 的每一元素与 `y` 的每一元素进行其它的某种运算, 而不仅限于乘积运算, 可以用 `outer(x,y,f)` 完成, 其中 `f` 是某种运算, 或者接受两个自变量的函数。

11.5 逆矩阵与线性方程组求解

用 `solve(A)` 求 `A` 的逆矩阵, 如

```
solve(B)
```

```
##      [,1] [,2]
## [1,]  0.5  0.5
## [2,] -0.5  0.5
```

用 `solve(A,b)` 求解线性方程组 $Ax = b$ 中的 x , 如

```
solve(B, c(1,2))
```

```
## [1] 1.5 0.5
```

求解了线性方程组

$$\begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} x = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

11.6 apply() 函数

`apply(A, 2, FUN)` 把矩阵 A 的每一列分别输入到函数 FUN 中，得到对应于每一列的结果，如

```
D <- matrix(c(6,2,3,5,4,1), nrow=3, ncol=2); D
```

```
##      [,1] [,2]
## [1,]    6    5
## [2,]    2    4
## [3,]    3    1
```

```
apply(D, 2, sum)
```

```
## [1] 11 10
```

`apply(A, 1, FUN)` 把矩阵 A 的每一行分别输入到函数 FUN 中，得到与每一行对应的结果，如

```
apply(D, 1, mean)
```

```
## [1] 5.5 3.0 2.0
```

如果函数 FUN 返回多个结果，则 `apply(A, 2, FUN)` 结果为矩阵，矩阵的每一列是输入矩阵相应列输入到 FUN 的结果，结果列数等于 A 的列数。如

```
apply(D, 2, range)
```

```
##      [,1] [,2]
## [1,]    2    1
## [2,]    6    5
```

如果函数 FUN 返回多个结果，为了对每行计算 FUN 的结果，结果存入一个与输入的矩阵行数相同的矩阵，应该用 `t(apply(A, 1, FUN))` 的形式，如

```
t(apply(D, 1, range))
```

```
##      [,1] [,2]
## [1,]    5    6
## [2,]    2    4
```

```
## [3,] 1 3
```

11.7 多维数组

矩阵是多维数组 (array) 的特例。矩阵是 $x_{ij}, i = 1, 2, \dots, n, j = 1, 2, \dots, m$ 这样的两下标数据的存储格式，三维数组是 $x_{ijk}, i = 1, 2, \dots, n, j = 1, 2, \dots, m, k = 1, 2, \dots, p$ 这样的三下标数据的存储格式， s 维数组则是有 s 个下标的数据的存储格式。实际上，给一个向量添加一个 `dim` 属性就可以把它变成多维数组。

多维数组的一般定义语法为

```
数组名 <- array(数组元素,
  dim=c(第一下标个数, 第二下标个数, ..., 第s下标个数))
```

其中数组元素的填入次序是第一下标变化最快，第二下标次之，最后一个下标是变化最慢的。这种次序称为 FORTRAN 次序。

下面是一个三维数组定义例子。

```
ara <- array(1:24, dim=c(2,3,4)); ara
```

```
## , , 1
##
##      [,1] [,2] [,3]
## [1,]  1   3   5
## [2,]  2   4   6
##
## , , 2
##
##      [,1] [,2] [,3]
## [1,]  7   9  11
## [2,]  8  10  12
##
## , , 3
##
```

```
##      [,1] [,2] [,3]
## [1,]  13  15  17
## [2,]  14  16  18
##
## , , 4
##
##      [,1] [,2] [,3]
## [1,]  19  21  23
## [2,]  20  22  24
```

这样的数组保存了 x_{ijk} , $i = 1, 2$, $j = 1, 2, 3$, $k = 1, 2, 3, 4$ 。三维数组 `ara` 可以看成是 4 个 2×3 矩阵。取出其中一个如 `ara[, , 2]` (取出第二个矩阵)

```
ara[, , 2]
```

```
##      [,1] [,2] [,3]
## [1,]    7    9   11
## [2,]    8   10   12
```

多维数组可以利用下标进行一般的子集操作，比如 `ara[, 2, 2:3]` 是 x_{ijk} , $i = 1, 2$, $j = 2$, $k = 2, 3$ 的值，结果是一个 2×2 矩阵：

```
ara[, 2, 2:3]
```

```
##      [,1] [,2]
## [1,]    9   15
## [2,]   10   16
```

多维数组在取子集时如果某一维下标是标量，则结果维数会减少，可以在方括号内用 `drop=FALSE` 选项避免这样的规则发生作用。

类似于矩阵，多维数组可以用一个矩阵作为下标，如果是三维数组，矩阵就需要有 3 列，四维数组需要用 4 列矩阵。下标矩阵的每行对应于一个数组元素。

Chapter 12

数据框

12.1 数据框

统计分析中最常见的原始数据形式是类似于数据库表或 Excel 数据表的形式。这样形式的数据在 R 中叫做数据框 (data.frame)。数据框类似于一个矩阵，有 n 行、 p 列，但各列允许有不同类型：数值型向量、因子、字符型向量、日期时间向量。同一列的数据类型相同。在 R 中数据框是一个特殊的列表，其每个列表元素都是一个长度相同的向量。事实上，数据框还允许一个元素是一个矩阵，但这样会使得某些读入数据框的函数发生错误。

函数 `data.frame()` 可以生成数据框，如

```
d <- data.frame(  
  name=c(" 李明", " 张聪", " 王建"),  
  age=c(30, 35, 28),  
  height=c(180, 162, 175),  
  stringsAsFactors=FALSE)  
print(d)
```

```
##   name age height  
## 1 李明  30   180  
## 2 张聪  35   162
```

```
## 3 王建 28 175
```

`data.frame()` 函数会将字符型列转换成因子,加选项 `stringsAsFactors=FALSE` 可以避免这样的转换。

数据框每列叫做一个变量,每列都有名字,称为列名或变量名,可以用 `names()` 函数和 `colnames()` 函数访问。如

```
names(d)
```

```
## [1] "name" "age" "height"
```

```
colnames(d)
```

```
## [1] "name" "age" "height"
```

给 `names(d)` 或 `colnames(d)` 赋值可以修改列名。

用 `as.data.frame(x)` 可以把 `x` 转换成数据框。如果 `x` 是一个向量,转换结果是以 `x` 为唯一一列的数据框。如果 `x` 是一个列表并且列表元素都是长度相同的向量,转换结果中每个列表变成数据框的一列。如果 `x` 是一个矩阵,转换结果把矩阵的每列变成数据框的一列。

数据框是一个随着 R 语言前身 S 语言继承下来的概念,现在已经有一些不足之处, `tibble` 包提供了 `tibble` 类,这是数据框的一个改进版本。

12.2 数据框内容访问

数据框可以用矩阵格式访问,如

```
d[2,3]
```

```
## [1] 162
```

访问单个元素。

```
d[[2]]
```

```
## [1] 30 35 28
```

访问第二列,结果为向量。

```
d[,2]
```

```
## [1] 30 35 28
```

也访问第二列，但是这种作法与 `tibble` 不兼容，所以应避免使用。

按列名访问列可用如

```
d[["age"]]
```

```
## [1] 30 35 28
```

```
d[, "age"]
```

```
## [1] 30 35 28
```

```
d$age
```

```
## [1] 30 35 28
```

其中第二种做法与 `tibble` 不兼容，应避免使用。

因为数据框的一行不一定是相同数据类型，所以数据框的一行作为子集，结果还是数据框，而不是向量。如

```
x <- d[2,]; x
```

```
##   name age height
```

```
## 2 张聪  35   162
```

```
is.data.frame(x)
```

```
## [1] TRUE
```

可以同时取行子集和列子集，如

```
d[1:2, 'age']
```

```
## [1] 30 35
```

```
d[1:2, c('age', 'height')]
```

```
##   age height
```

```
## 1  30   180
```

```
## 2 35 162
```

```
d[d[, 'age'] >= 30, ]
```

```
## name age height
## 1 李明 30 180
## 2 张聪 35 162
```

与矩阵类似地是，用如 `d[, 'age']`, `d[, 2]` 这样的方法取出的数据框的单个列是向量而不再是数据框。但是，如果取出两列或者两列以上，结果则是数据框。如果取列子集时不能预先知道取出的列个数，则子集结果有可能是向量也有可能是数据框，容易造成后续程序错误。对一般的数据框，可以在取子集的方括号内加上 `drop=FALSE` 选项，确保取列子集的结果总是数据框。数据框的改进类型 `tibble` 在取出列子集时保持为 `tibble` 格式。

对数据框变量名按照字符串与集合进行操作可以实现复杂的列子集筛选。

数据框每一行可以有行名，这在原始的 S 语言和传统的 R 语言中是重要的技术，但是在改进类型 `tibble` 中则取消了行名，需要用列名实现功能一般改用 `left_join()` 函数实现。

比如，每一行定义行名为身份证号，则可以唯一识别各行。下面的例子以姓名作为行名：

```
rownames(d) <- d$name
d$name <- NULL
d
```

```
## age height
## 李明 30 180
## 张聪 35 162
## 王建 28 175
```

用数据框的行名可以建立一个值到多个值的对应表。比如，有如下的数据框：

```
dm <- data.frame(
  '年级' = 1:6,
  '出游' = c(0, 2, 2, 2, 2, 1),
  '疫苗' = c(T, F, F, F, T, F)
```



```
)
```

其中“出游”是每个年级安排的出游次数，“疫苗”是该年级有全体无计划免疫注射。把年级变成行名，可以建立年级到出游次数与疫苗注射的对应表：

```
rownames(dm) <- dm[[' 年级']]
dm[[' 年级']] <- NULL
```

这样，假设某个社区的小学中抽取的 4 个班的年级为 `c(2,1,1,3)`，其对应的出游和疫苗注射信息可查询如下：

```
x <- c(2,1,1,3)
dm[as.character(x),]
```

```
##      出游  疫苗
## 2      2 FALSE
## 1      0  TRUE
## 1.1    0  TRUE
## 3      2 FALSE
```

结果中包含了不必要也不太合适的行名，可以去掉，以上程序改成：

```
x <- c(2,1,1,3)
xx <- dm[as.character(x),]
rownames(xx) <- NULL
xx
```

```
##      出游  疫苗
## 1      2 FALSE
## 2      0  TRUE
## 3      0  TRUE
## 4      2 FALSE
```

如果要从多个值建立映射，比如，从省名与县名映射到经度、纬度，可以预先用 `paste()` 函数把省名与县名合并在一起，中间以适当字符（如“-”）分隔，以这样的合并字符串为行名。

对于代替数据框的 `tibble` 类型，如果要实现行名的功能，可以将行名作为单独的一列，然后用 `dplyr` 包的 `inner_join()`、`left_join()`、`full_join()` 等

函数横向合并数据集。参见27.15。

12.3 数据框与矩阵的区别

数据框不能作为矩阵参加矩阵运算。需要时，可以用 `as.matrix()` 函数转换数据框或数据框的子集为矩阵。如

```
d2 <- as.matrix(d[,c("age", "height")])
d3 <- crossprod(d2); d3
```

```
##           age height
## age      2909 15970
## height 15970 89269
```

这里 `crossprod(A)` 表示 $A^T A$ 。

12.4 `gl()` 函数

可以用数据框保存试验结果，对有多个因素的试验，往往需要生成多个因素完全搭配并重复的表格。函数 `gl()` 可以生成这样的重复模式。比如，下面的例子：

```
d4 <- data.frame(
  group=gl(3, 10, length=30),
  subgroup=gl(5,2,length=30),
  obs=gl(2,1,length=30))
print(d4)
```

```
##   group subgroup obs
## 1     1         1   1
## 2     1         1   2
## 3     1         2   1
## 4     1         2   2
## 5     1         3   1
## 6     1         3   2
```

```
## 7      1      4      1
## 8      1      4      2
## 9      1      5      1
## 10     1      5      2
## 11     2      1      1
## 12     2      1      2
## 13     2      2      1
## 14     2      2      2
## 15     2      3      1
## 16     2      3      2
## 17     2      4      1
## 18     2      4      2
## 19     2      5      1
## 20     2      5      2
## 21     3      1      1
## 22     3      1      2
## 23     3      2      1
## 24     3      2      2
## 25     3      3      1
## 26     3      3      2
## 27     3      4      1
## 28     3      4      2
## 29     3      5      1
## 30     3      5      2
```

结果的数据框 `d` 有三个变量: `group` 是大组, 共分 3 个大组, 每组 10 个观测; `subgroup` 是子组, 在每个大组内分为 5 个子组, 每个子组 2 个观测。共有 $3 \times 5 \times 2 = 30$ 个观测 (行)。

`gl()` 第一个参数是因子水平个数, 第二个参数是同一因子水平连续重复次数, 第三个参数是总共需要的元素个数, 所有水平都出现后则重复整个模式直到长度满足要求。

12.5 tibble 类型

tibble 类型是一种改进的数据框。readr 包的 `read_csv()` 函数是 `read.csv()` 函数的一个改进版本, 它将 CSV 文件读入为 tibble 类型, 如文件 `class.csv` 的读入:

```
library(tibble)
library(readr)
t.class <- read_csv("class.csv")

## Parsed with column specification:
## cols(
##   name = col_character(),
##   sex = col_character(),
##   age = col_double(),
##   height = col_double(),
##   weight = col_double()
## )

t.class

## # A tibble: 19 x 5
##   name    sex    age height weight
##   <chr> <chr> <dbl> <dbl> <dbl>
## 1 Alice  F      13   56.5   84
## 2 Becka  F      13   65.3   98
## 3 Gail   F      14   64.3   90
## 4 Karen  F      12   56.3   77
## 5 Kathy  F      12   59.8  84.5
## 6 Mary   F      15   66.5  112
## 7 Sandy  F      11   51.3  50.5
## 8 Sharon F      15   62.5  112.
## 9 Tammy  F      14   62.8  102.
## 10 Alfred M      14    69   112.
## 11 Duke   M      14   63.5  102.
## 12 Guido  M      15    67   133
```

```
## 13 James M 12 57.3 83
## 14 Jeffrey M 13 62.5 84
## 15 John M 12 59 99.5
## 16 Philip M 16 72 150
## 17 Robert M 12 64.8 128
## 18 Thomas M 11 57.5 85
## 19 William M 15 66.5 112
```

tibble 类型的类属依次为 `tbl_df`, `tbl`, `data.frame`:

```
class(t.class)

## [1] "spec_tbl_df" "tbl_df"      "tbl"        "data.frame"
```

用 `as_tibble()` 可以将一个数据框转换为 tibble, `dplyr` 包提供了 `filter()`、`select()`、`arrange()`、`mutate()` 等函数用来对 tibble 选取行子集、列子集, 排序、修改或定义新变量, 等等。见27。

可以用 `tibble()` 函数生成小的 tibble, 如

```
t.bp <- tibble(
  `序号` = c(1,5,6,9,10,15),
  `收缩压` = c(145, 110, "未测", 150, "拒绝", 115))
t.bp

## # A tibble: 6 x 2
##   序号 收缩压
##   <dbl> <chr>
## 1     1 145
## 2     5 110
## 3     6 未测
## 4     9 150
## 5    10 拒绝
## 6    15 115
```

用 `tribble` 可以按类似于 CSV 格式输入一个 tibble, 如

```
t.bp <- tribble(
  ~`序号`, ~`收缩压`,
  1, 145,
  5, 110,
  6, "未测",
  9, 150,
  10, "拒绝",
  15, 115
)
```

```
## # A tibble: 6 x 2
##   序号 收缩压
##   <dbl> <chr>
## 1     1  145
## 2     5  110
## 3     6 未测
## 4     9  150
## 5    10 拒绝
## 6    15  115
```

注意 `tribble()` 中数据每行末尾也需要有逗号，最后一行末尾没有逗号。这比较适用于在程序中输入小的数据集。

`tibble` 与数据框的一大区别是在显示时不自动显示所有内容，这样可以避免显示很大的数据框将命令行的所有显示都充满。可以在 `print()` 用 `n=` 和 `width=` 选项指定要显示的行数和列数。

另外，用单重的方括号取列子集时，即使仅取一列，从 `tibble` 取出的一列结果仍是 `tibble` 而不是向量，这时应使用双方括号格式或 `$` 格式。因为这个原因有些原来的程序输入 `tibble` 会出错，这时可以用 `as.data.frame()` 转换成数据框。如：

```
t.bp[, "收缩压"]
```

```
## # A tibble: 6 x 1
```

```
## 收缩压
## <chr>
## 1 145
## 2 110
## 3 未测
## 4 150
## 5 拒绝
## 6 115
```

```
t.bp[["收缩压"]]
```

```
## [1] "145" "110" "未测" "150" "拒绝" "115"
```

tibble 在定义时不需要列名为合法变量名，但是作为变量名使用时需要用反单撇号包裹。tibble 不使用行名，需要行名时，将其保存为 tibble 的一列。原来用行名完成的功能，可以改用 dplyr 包的 `left_join()` 等函数，这些函数进行数据框的横向连接。

12.6 练习

假设 `class.csv` 已经读入为 R 数据框 `d.class`，其中的 `sex` 列已经自动转换为因子。

- (1) 显示 `d.class` 中年龄至少为 15 的行子集；
- (2) 显示女生且年龄至少为 15 的学生姓名和年龄；
- (3) 取出数据框中的 `age` 变量赋给变量 `x`。

Chapter 13

列表类型

13.1 列表

R 中列表 (list) 类型来保存不同类型的数据。一个主要目的是提供 R 分析结果输出包装：输出一个变量，这个变量包括回归系数、预测值、残差、检验结果等一系列不能放到规则形状数据结构中的内容。实际上，数据框也是列表的一种，但是数据框要求各列等长，而列表不要求。

列表可以有多个元素，但是与向量不同的是，列表的不同元素的类型可以不同，比如，一个元素是数值型向量，一个元素是字符串，一个元素是标量，一个元素是另一个列表。

定义列表用函数 `list()`，如

```
rec <- list(name=" 李明", age=30,  
           scores=c(85, 76, 90))  
rec
```

```
## $name  
## [1] "李明"  
##  
## $age  
## [1] 30
```

```
##  
## $scores  
## [1] 85 76 90
```

用 `typeof()` 函数判断一个列表，返回结果为 `list`。可以用 `is.list()` 函数判断某个对象是否是列表类型。

13.2 列表元素访问

列表的一个元素也可以称为列表的一个“变量”，单个列表元素必须用两重方括号格式访问，如

```
rec[[3]]  
  
## [1] 85 76 90
```

```
rec[[3]][2]  
  
## [1] 76
```

```
rec[["age"]]  
  
## [1] 30
```

如果使用单重方括号对列表取子集，结果还是列表而不是列表元素，如

```
rec[3]  
  
## $scores  
## [1] 85 76 90
```

列表的单个元素也可以用 `$` 格式访问，如

```
rec$age  
  
## [1] 30
```

列表一般都应该有元素名，元素名可以看成是变量名，列表中的每个元素看成一个变量。用 `names()` 函数查看和修改元素名。如

```
names(rec)

## [1] "name" "age" "scores"

names(rec)[names(rec)=='scores'] <- '三科分数'
names(rec)

## [1] "name" "age" "三科分数"

rec[["三科分数"]]

## [1] 85 76 90
```

可以修改列表元素内容。如

```
rec[["三科分数"]][2] <- 0
print(rec)
```

```
## $name
## [1] "李明"
##
## $age
## [1] 30
##
## $三科分数
## [1] 85 0 90
```

直接给列表不存在的元素名定义元素值就添加了新元素，而且不同于使用向量，对于列表而言这是很正常的做法，比如

```
rec[["身高"]] <- 178
print(rec)
```

```
## $name
## [1] "李明"
##
## $age
## [1] 30
##
```

```
## $三科分数
## [1] 85 0 90
##
## $身高
## [1] 178
```

把某个列表元素赋值为 `NULL` 就删掉这个元素。如

```
rec[['age']] <- NULL
print(rec)
```

```
## $name
## [1] "李明"
##
## $三科分数
## [1] 85 0 90
##
## $身高
## [1] 178
```

在 `list()` 函数中允许定义元素为 `NULL`，这样的元素是存在的，如：

```
li <- list(a=120, b='F', c=NULL); li
```

```
## $a
## [1] 120
##
## $b
## [1] "F"
##
## $c
## NULL
```

但是，要把已经存在的元素修改为 `NULL` 值而不是删除此元素，或者给列表增加一个取值为 `NULL` 的元素，这时需要用单重的方括号取子集，这样的子集会保持其列表类型，给这样的子列表赋值为 `list(NULL)`，如：

```
li['b'] <- list(NULL)
li['d'] <- list(NULL)
li
```

```
## $a
## [1] 120
##
## $b
## NULL
##
## $c
## NULL
##
## $d
## NULL
```

13.3 列表类型转换

用 `as.list()` 把一个其它类型的对象转换成列表；用 `unlist()` 函数把列表转换成基本向量。如

```
li1 <- as.list(1:3)
li1
```

```
## [[1]]
## [1] 1
##
## [[2]]
## [1] 2
##
## [[3]]
## [1] 3
```

```
li2 <- list(x=1, y=c(2,3))
unlist(li2)
```

```
## x y1 y2
## 1 2 3
```

13.4 返回列表的函数示例—`strsplit()`

`strsplit()` 输入一个字符型向量并指定一个分隔符，返回一个项数与字符型向量元素个数相同的列表，列表每项对应于字符型向量中一个元素的拆分结果。如

```
x <- c('10, 8, 7', '5, 2, 2', '3, 7, 8', '8, 8, 9')
res <- strsplit(x, ',');
```

```
## [[1]]
## [1] "10" " 8" " 7"
##
## [[2]]
## [1] "5" " 2" " 2"
##
## [[3]]
## [1] "3" " 7" " 8"
##
## [[4]]
## [1] "8" " 8" " 9"
```

为了把拆分结果进一步转换成一个数值型矩阵，可以使用 `sapply()` 函数如下：

```
t(sapply(res, as.numeric))
```

```
##      [,1] [,2] [,3]
## [1,]  10   8   7
## [2,]   5   2   2
## [3,]   3   7   8
## [4,]   8   8   9
```

sapply() 函数是 apply 类函数之一，稍后再详细进行讲解。

Chapter 14

工作空间

R 把在命令行定义的变量都保存到工作空间中，在退出 R 时可以选择是否保存工作空间。这也是 R 与其他如 C、Java 这样的语言的区别之一。

用 `ls()` 命令可以查看工作空间中的内容。

随着多次在命令行使用 R，工作空间的变量越来越多，使得重名的可能性越来越大，而且工作空间中变量太多也让我们不容易查看其内容。在命令行定义的变量称为“全局变量”，在编程实际中，全局变量是需要慎用的。

可以用 `rm()` 函数删除工作空间中的变量，格式如

```
rm(d, h, name, rec, sex, x)
```

要避免工作空间杂乱，最好的办法还是所有的运算都写到自定义函数中。自定义函数中定义的变量都是临时的，不会保存到工作空间中。这样，仅需要时才把变量值在命令行定义，这样的变量一般是读入的数据或自定义的函数（自定义函数也保存在工作空间中）。

可以定义如下的 `sandbox()` 函数：

```
sandbox <- function(){  
  cat(' 沙盘：接连的空行回车可以退出。\\n')  
  browser()  
}
```

运行 `sandbox()` 函数，将出现如下的 browser 命令行：

沙盘：接连的空行回车可以退出。

```
Called from: sandbox()
```

```
Browse[1]>
```

提示符变成了“Browser[n]”，其中 `n` 代表层次序号。在这样的 browser 命令行中随意定义变量，定义的变量不会保存到工作空间中。用“Q”命令可以退出这个沙盘环境，接连回车也可以退出。

Part III

R 编程

Chapter 15

R 输入输出

15.1 输入输出的简单方法

15.1.1 简单的输出

用 `print()` 函数显示某个变量或表达式的值，如

```
x <- 1.234  
print(x)
```

```
## [1] 1.234
```

```
y <- c(1,3,5)  
print(y[2:3])
```

```
## [1] 3 5
```

在命令行使用 R 时，直接以变量名或表达式作为命令可以起到用 `print()` 函数显示的相同效果。

用 `cat()` 函数把字符串、变量、表达式连接起来显示，其中变量和表达式的类型一般是标量或向量，不能是矩阵、列表等复杂数据。如

```
cat("x =", x, "\n")
```

```
## x = 1.234
```

```
cat("y =", y, "\n")
```

```
## y = 1 3 5
```

注意 `cat()` 显示中需要换行需要在自变量中包含字符串 `"\n"`，即换行符。

`cat()` 默认显示在命令行窗口，为了写入指定文件中，在 `cat()` 调用中用 `file=` 选项，这时如果已有文件会把原有内容覆盖，为了在已有文件时不覆盖原有内容而是在末尾添加，在 `cat()` 中使用 `append=TRUE` 选项。如：

```
cat("=== 结果文件 ===\n", file="res.txt")
cat("x =", x, "\n", file="res.txt", append=TRUE)
```

函数 `sink()` 可以用来把命令行窗口显示的运行结果转向保存到指定的文本文件中，如果希望保存到文件的同时也在命令行窗口显示，使用 `split=TRUE` 选项。如

```
sink("allres.txt", split=TRUE)
```

为了取消这样的输出文件记录，使用不带自变量的 `sink()` 调用，如

```
sink()
```

在 R 命令行环境中定义的变量、函数会保存在工作空间中，并在退出 R 会话时可以保存到硬盘文件中。用 `save()` 命令要求把指定的若干个变量（直接用名字，不需要表示成字符串）保存到用 `file=` 指定的文件中，随后可以用 `load()` 命令恢复到工作空间中。虽然允许保存多个变量到同一文件中，但尽可能仅保存一个变量，而且使用变量名作为文件名。用 `save()` 保存的 R 特殊格式的文件是通用的，不依赖于硬件和操作系统。如

```
save(scores, file="scores.RData")
load("scores.RData")
```

对于一个数据框，可以用 `write.csv()` 或 `readr::write_csv()` 将其保存为逗号分隔的文本文件，这样的文件可以很容易地被其它软件识别访问，如 Microsoft Excel 软件可以很容易地把这样的文件读成电子表格。用如

```
da <- tibble('name'=c(' 李明', ' 刘颖', ' 张浩'),
             'age'=c(15, 17, 16))
write_csv(da, path="mydata.csv")
```

结果生成的 mydata.csv 文件内容如下:

```
name,age
李明,15
刘颖,17
张浩,16
```

但是,在 Microsoft 的中文版 Windows 操作系统中,默认编码是 GB 编码,用 `write_csv()` 生成的 CSV 文件总是使用 UTF-8 编码,系统中的 MS Office 软件不能自动识别这样编码的 CSV 文件。`write.csv()` 函数不存在这个问题。

15.1.2 简单的输入

用 `scan()` 函数可以输入文本文件中的数值向量,文件名用 `file=` 选项给出。文件中数值之间以空格分开。如

```
cat(1:12, "\n", file="d:/work/x.txt")
x <- scan("d:/work/x.txt")
```

程序中用全路径给出了输入文件位置,注意路径中用了正斜杠/作为分隔符,如果在 MS Windows 环境下使用\作为分隔符,在 R 的字符串常量中\必须写成\\。

如果 `scan()` 中忽略输入文件参数,此函数将从命令行读入数据。可以在一行用空格分开多个数值,可以用多行输入直到空行结束输入。

这样的方法也可以用来读入矩阵。设文件 `mat.txt` 包含如下矩阵内容:

```
3 4 2
5 12 10
7 8 6
1 9 11
```

可以先把文件内容读入到一个 R 向量中,再利用 `matrix()` 函数转换成矩阵,

注意要使用 `byrow=TRUE` 选项，而且只要指定 `ncol` 选项，可以忽略 `nrow` 选项。如

```
M <- matrix(scan('mat.txt', quiet=TRUE), ncol=3, byrow=TRUE)
M
```

`scan()` 中的 `quite=TRUE` 选项使得读入时不自动显示读入的数值项数。

上面读入数值矩阵的方法在数据量较大的情形也可以使用，与之不同的是，`read.table()` 或 `readr::read_table()` 函数也可以读入这样的数据，但是会保存成数据框而不是矩阵，而且 `read.table()` 函数在读入大规模的矩阵时效率很低。

15.2 读取 CSV 文件

对于保存在文本文件中的电子表格数据，R 可以用 `read.csv()`, `read.table()`, `read.delim()`, `read.fwf()` 等函数读入，但是建议在 `readr` 包的支持下用 `read_csv()`, `read_table2()`, `read_delim()`, `read_fwf()` 等函数读入，这些将读入的数据框保存为 `tibble` 类型，`tibble` 是数据框的一个变种，改善了数据框的一些不适当的设计。`readr` 的读入速度比基本 R 软件的 `read.csv()` 等函数的速度快得多，速度可以相差 10 倍，也不自动将字符型列转换成因子，不自动修改变量名为合法变量名，不设置行名。

对于中小规模的数据，CSV 格式作为文件交换格式比较合适，兼容性强，各种数据管理软件与统计软件都可以很容易地读入和生成这样格式的文件，但是特别大型的数据读入效率很低。

CSV 格式的文件用逗号分隔开同一行的数据项，一般第一行是各列的列名（变量名）。对于数值型数据，只要表示成数值常量形式即可。对于字符型数据，可以用双撇号包围起来，也可以不用撇号包围。但是，如果数据项本身包含逗号，就需要用双撇号包围。例如，下面是一个名为 `testcsv.csv` 的文件内容，其中演示了内容中有逗号、有双撇号的情况。

```
id,words
1,"PhD"
2,Master's degree
```



```
3, "Bond, James"  
4, "A "special" gift"
```

为读入上面的内容，只要用如下程序：

```
d <- read_csv("testcsv.csv")
```

读入的数据框显示如下：

```
# A tibble: 4 × 2  
  id      words  
  <int>   <chr>  
1     1     PhD  
2     2 Master's degree  
3     3     Bond, James  
4     4 A "special" gift
```

`read_csv()` 还可以从字符串读入一个数据框，如

```
d.small <- read_csv("name,x,y  
John, 33, 95  
Kim, 21, 64  
Sandy, 49, 100  
")  
d.small
```

```
## # A tibble: 3 × 3  
##   name      x      y  
##   <chr> <dbl> <dbl>  
## 1 John      33      95  
## 2 Kim        21      64  
## 3 Sandy      49     100
```

`read_csv()` 的 `skip=` 选项跳过开头的若干行。当数据不包含列名时，只要指定 `col_names=FALSE`，变量将自动命名为 `X1`, `X2`, ...，也可以用 `col_names=` 指定各列的名字，如

```
d.small <- read_csv("John, 33, 95  
Kim, 21, 64
```

```
Sandy, 49, 100
", col_names=c("name", "x", "y") )
d.small
```

```
## # A tibble: 3 x 3
##   name      x      y
##   <chr> <dbl> <dbl>
## 1 John     33     95
## 2 Kim      21     64
## 3 Sandy    49    100
```

`read_csv()` 将空缺的值读入为缺失值, 将“NA”也读入为缺失值。可以用 `na=` 选项改变这样的设置。也可以将带有缺失值的列先按字符型原样读入, 然后再进行转换。

CSV 文件是文本文件, 是有编码问题的, 尤其是中文内容的文件。`readr` 包的默认编码是 UTF-8 编码。例如, 文件 `bp.csv` 以 GBK 编码 (有时称为 GB18030 编码, 这是中文 Windows 所用的中文编码) 保存了如下内容:

```
序号,收缩压
1,145
5,110
6, 未测
9,150
10, 拒绝
15,115
```

如果直接用 `read_csv()`:

```
d <- read_csv("bp.csv")
```

可能在读入时出错, 或者访问时出错。为了读入用 GBK 编码的中文 CSV 文件, 需要利用 `locale` 参数和 `locale()` 函数:

```
d <- read_csv("bp.csv", locale=locale(encoding="GBK"))
```

```
## Parsed with column specification:
## cols(
```

```
## 序号 = col_double(),  
## 收缩压 = col_character()  
## )
```

```
d
```

```
## # A tibble: 6 x 2  
##   序号 收缩压  
##   <dbl> <chr>  
## 1     1 145  
## 2     5 110  
## 3     6 未测  
## 4     9 150  
## 5    10 拒绝  
## 6    15 115
```

对每列的类型，`readr` 用前 1000 行猜测合理的类型，并在读取后显示猜测的每列类型。

但是有可能类型改变发生在 1000 行之后。`col_types` 选项可以指定每一列的类型，如"`col_double()`"、"`col_integer()`"、"`col_character()`"、"`col_factor()`"、"`col_date()`"、"`col_datetime()`" 等。`cols()` 函数可以用来规定各列类型，并且有一个 `.default` 参数指定缺省类型。对因子，需要在 `col_factor()` 中用 `levels=` 指定因子水平。

可以复制 `readr` 猜测的类型作为 `col_types` 的输入，这样当数据变化时不会因为偶尔猜测错误而使得程序出错。如

```
d <- read_csv("bp.csv", locale=locale(encoding="GBK"),  
              col_types=cols(  
                `序号` = col_integer(),  
                `收缩压` = col_character()  
              ))
```

```
d
```

```
## # A tibble: 6 x 2  
##   序号 收缩压  
##   <int> <chr>
```

```
## 1      1 145
## 2      5 110
## 3      6 未测
## 4      9 150
## 5     10 拒绝
## 6     15 115
```

当猜测的文件类型有问题的时候，可以先将所有列都读成字符型，然后用 `type_convert()` 函数转换，如：

```
d <- read_csv("filename.csv",
              col_types=cols(.default = col_character()))
d <- type_convert(d)
```

读入有错时，对特大文件可以先少读入一些行，用 `nmax=` 可以指定最多读入多少行。调试成功后再读入整个文件。

设文件 `class.csv` 内容如下：

```
name,sex,age,height,weight
Alice,F,13,56.5,84
Becka,F,13,65.3,98
Gail,F,14,64.3,90
Karen,F,12,56.3,77
Kathy,F,12,59.8,84.5
Mary,F,15,66.5,112
Sandy,F,11,51.3,50.5
Sharon,F,15,62.5,112.5
Tammy,F,14,62.8,102.5
Alfred,M,14,69,112.5
Duke,M,14,63.5,102.5
Guido,M,15,67,133
James,M,12,57.3,83
Jeffrey,M,13,62.5,84
John,M,12,59,99.5
Philip,M,16,72,150
```

```
Robert,M,12,64.8,128  
Thomas,M,11,57.5,85  
William,M,15,66.5,112
```

最简单地用 `read_csv()` 读入上述 CSV 文件，程序如：

```
d.class <- read_csv('class.csv')
```

```
## Parsed with column specification:  
## cols(  
##   name = col_character(),  
##   sex = col_character(),  
##   age = col_double(),  
##   height = col_double(),  
##   weight = col_double()  
## )
```

```
knitr::kable(d.class)
```

name	sex	age	height	weight
Alice	F	13	56.5	84.0
Becka	F	13	65.3	98.0
Gail	F	14	64.3	90.0
Karen	F	12	56.3	77.0
Kathy	F	12	59.8	84.5
Mary	F	15	66.5	112.0
Sandy	F	11	51.3	50.5
Sharon	F	15	62.5	112.5
Tammy	F	14	62.8	102.5
Alfred	M	14	69.0	112.5
Duke	M	14	63.5	102.5
Guido	M	15	67.0	133.0
James	M	12	57.3	83.0
Jeffrey	M	13	62.5	84.0
John	M	12	59.0	99.5
Philip	M	16	72.0	150.0
Robert	M	12	64.8	128.0
Thomas	M	11	57.5	85.0
William	M	15	66.5	112.0

从结果看出，读入后显示了每列的类型。对性别变量，没有自动转换成因子，而是保存为字符型。为了按自己的要求转换各列类型，用了 `read_csv()` 的 `coltypes=` 选项和 `cols()` 函数如下：

```
ct <- cols(
  .default = col_double(),
  name=col_character(),
  sex=col_factor(levels=c("M", "F"))
)
d.class <- read_csv('class.csv', col_types=ct)
str(d.class)
```

```
## Classes 'spec_tbl_df', 'tbl_df', 'tbl' and 'data.frame': 19 obs. of 5 variables
## $ name : chr "Alice" "Becka" "Gail" "Karen" ...
```

```
## $ sex : Factor w/ 2 levels "M","F": 2 2 2 2 2 2 2 2 2 1 ...
## $ age : num 13 13 14 12 12 15 11 15 14 14 ...
## $ height: num 56.5 65.3 64.3 56.3 59.8 66.5 51.3 62.5 62.8 69 ...
## $ weight: num 84 98 90 77 84.5 ...
## - attr(*, "spec")=
## .. cols(
## .. .default = col_double(),
## .. name = col_character(),
## .. sex = col_factor(levels = c("M", "F"), ordered = FALSE, include_na = FALSE),
## .. age = col_double(),
## .. height = col_double(),
## .. weight = col_double()
## .. )
```

其中 `str()` 函数可以显示数据框的行数 (obs.) 和变量数 (variables), 以及每个变量 (列) 的类属等信息。

除了 `read_csv()` 函数以外, R 扩展包 `readr` 还提供了其它的从文本数据读入数据框的函数, 如 `read_table2()`, `read_tsv()`, `read_fwf()` 等。这些函数读入的结果保存为 `tibble`。`read_table2()` 读入用空格作为间隔的文本文件, 同一行的两个数据项之间可以用一个或多个空格分隔, 不需要空格个数相同, 也不需要上下对齐。`read_tsv()` 读入用制表符分隔的文件。`read_fwf()` 读入上下对齐的文本文件。

另外, `read_lines()` 函数将文本文件各行读入为一个字符型向量。`read_file()` 将文件内容读入成一整个字符串, `read_file_raw()` 可以不管文件编码将文件读入为一个二进制字符串。

对特别大的文本格式数据, `data.table` 扩展包的 `fread()` 读入速度更快。

`readr` 包的 `write_excel_csv()` 函数将 `tibble` 保存为 `csv` 文件, 总是使用 UTF-8 编码, 结果可以被 MS Office 读取。

文本格式的文件都不适用于大型数据的读取与保存。大型数据可以通过数据库接口访问, 可以用 R 的 `save()` 和 `load()` 函数按照 R 的格式访问, 还有一些特殊的针对大数据集的 R 扩展包。

15.3 Excel 表访问

15.3.1 借助于文本格式

为了把 Microsoft Excel 格式的数据读入到 R 中，最容易的办法是在 Excel 软件中把数据表转存为 CSV 格式，然后用 `read.csv()` 读取。

为了把 R 的数据框保存为 Excel 格式，只要用 `write.csv()` 把数据框保存成 CSV 格式，然后在 Excel 中打开即可。例如，下面的程序演示了 `write.csv()` 的使用：

```
d1 <- tibble(" 学号"=c("101", "103", "104"),
             " 数学"=c(85, 60, 73),
             " 语文"=c(90, 78, 80))
write.csv(d1, file="tmp1.csv", row.names=FALSE)
```

保存在文件中的结果显示如下：

```
学号,数学,语文
101,85,90
103,60,78
104,73,80
```

15.3.2 使用剪贴板

为了把 Excel 软件中数据表的选中区域读入到 R 中，可以借助于剪贴板。在 Excel 中复制选中的区域，然后在 R 中用如

```
myDF <- read.delim("clipboard")
```

就可以把选中部分转换成一个 R 的数据框。如果复制的区域不含列名，应加上 `header=FALSE` 选项。

这种方法也可以从 R 中复制数据到在 Excel 中打开的电子表格中，例如

```
write.table(iris, file="clipboard", sep = "\t", col.names = NA)
```

首先把指定的数据框（这里是 `iris`）写入到了剪贴板，然后在用 Excel 软件打

开的工作簿中只要粘贴就可以。上述程序中 `write.table()` 函数把指定的数据框写入到指定的文件中, 其中的 `col.names=NA` 选项是一个特殊的约定, 这时保存的文件中第一行是列名, 如果有行名的话, 行名所在的列对应的列名是空白的 (但是存在此项)。

如果从 R 中复制数据框到打开的 Excel 文件中时不带行名, 但是带有列名, 可以写这样一个通用函数

```
write_clipboard <- function(df){  
  write.table(df, file="clipboard", sep='\t',  
             row.names=FALSE)  
}
```

15.3.3 利用 readxl 扩展包

readxl 扩展包的 `readxl()` 函数利用独立的 C 和 C++ 库函数读入.xls 和.xlsx 格式的 Excel 文件。一般格式为

```
read_excel(path, sheet = 1, col_names = TRUE,  
           col_types = NULL, na = "", skip = 0)
```

结果返回读入的表格为一个数据框。各个自变量为:

- **path**: 要读入的 Excel 文件名, 可以是全路径, 路径格式要符合所用操作系统要求。
- **sheet**: 要读入哪个工作簿 (sheet), 可以是整数序号, 也可以是工作簿名称的字符串。
- **col_names**: 是否用第一行内容作为列名, 缺省为是。
- **col_types**: 可以在读入时人为指定各列的数据类型, 缺省时从各列内容自动判断, 有可能会不够准确。人为指定时, 指定一个对应于各列的字符型向量, 元素可取值为:
 - **blank**: 自动判断该列;
 - **numeric**: 数值型;
 - **date**: 日期;
 - **text**: 字符型。

15.3.4 利用 RODBC 访问 Excel 文件

还可以用 RODBC 扩展包访问 Excel 文件。这样的方法不需要借助于 CSV 文件这个中间格式。RODBC 是一个通过 ODBC 协议访问数据文件与数据库的 R 扩展包。

先给出把 R 数据框保存为 Excel 文件的例子。如下的程序定义了两个数据框：

```
d1 <- data.frame("学号"=c("101", "103", "104"),
                 "数学"=c(85, 60, 73),
                 "语文"=c(90, 78, 80))
d2 <- data.frame("学号"=c("101", "103", "104"),
                 "性别"=c("女", "男", "男"))
```

在写入到 Excel 文件时，如果文件已经存在，会导致写入失败。比如，要写入到 `testwrite.xls` 中，可以用如下程序在文件已存在时先删除文件：

```
fname <- "testwrite.xls"
if(file.exists(fname)) file.remove(fname)
```

其中 `file.exists()` 检查文件是否已存在，`file.remove()` 删除指定文件。

使用 RODBC 比较麻烦，需要先用 `odbcConnectExcel()` 函数打开目的文件，然后可以用 `sqlSave()` 函数把数据框保存到目的文件中，保存完毕后需要用 `close()` 函数关闭打开的目的文件。目前 RODBC 的 `odbcConnectExcel()` 只能在 32 位版本的 R 软件中使用，而且操作系统中必须安装有 32 位的 ODBC 驱动程序。示例如下（需要使用 32 位 R 软件且需要操作系统中有 32 位版本的 ODBC 驱动程序）：

```
library(RODBC)
con <- odbcConnectExcel(fname, readOnly=FALSE)
res <- sqlSave(con, d1, tablename="成绩",
              rownames=F, colnames=F, safer=T)
res <- sqlSave(con, d2, tablename="性别",
              rownames=F, colnames=F, safer=T)
close(con)
```

用 `odbcConnectExcel2007()` 可以访问或生成 Excel 2007/2010 版本的 .xlsx

文件，此函数可以用在 64 位的 R 软件中，但是这时需要操作系统中安装有 64 位的 ODBC 驱动程序，而不能有 32 位的 ODBC 驱动程序。如果安装了 Office 软件，Office 软件是 32 位的，相应的 ODBC 驱动程序必须也是 32 位的；Office 软件是 64 位的，相应的 ODBC 驱动程序必须也是 64 位的。

RODBC 对 Excel 文件的支持还有一些其它的缺点，比如表名不规范，数据类型自动转换不一定合理等。在 Excel 中读入或者保存 CSV 格式会使得问题变得简单。大量数据或大量文件的问题就不应该使用 Excel 来管理了，一般会使用关系数据库系统，如 Oracle, MySQL 等。

为了读入 Excel 文件内容，先用 `odbcConnectExcel()` 函数打开文件，用 `sqlFetch()` 函数读入一个数据表为 R 数据框，读取完毕后用 `close()` 关闭打开的文件。如

```
require(RODBC)
con <- odbcConnectExcel('testwrite.xls')
rd1 <- sqlFetch(con, sqtable='成绩')
close(con)
```

读入的表显示如下：

	学号	数学	语文
1	101	85	90
2	103	60	78
3	104	73	80

15.3.5 用 RODBC 访问 Access 数据库

RODBC 还可以访问其他微机数据库软件的数据库。假设有 Access 数据库在文件 `c:/Friends/birthdays.mdb` 中，内有两个表 Men 和 Women，每个表包含域 Year, Month, Day, First Name, Last Name, Death。域名应尽量避免用空格。

下面的程序把女性记录的表读入为 R 数据框：

```
require(RODBC)
con <- odbcConnectAccess("c:/Friends/birthdays.mdb")
```

```
women <- sqlFetch(con, sqtable='Women')
close(con)
```

RODBC 还有许多与数据库访问有关的函数，比如，`sqlQuery()` 函数可以向打开的数据库提交任意符合标准的 SQL 查询。

15.4 使用专用接口访问数据库

15.4.1 访问 Oracle 数据库

Oracle 是最著名的数据库服务器软件。要访问的数据库，可以是安装在本机上的，也可以是安装在网络上某个服务器中的。如果是远程访问，需要在本机安装 Oracle 的客户端软件。

假设已经在本机安装了 Oracle 服务器软件，并设置 `orcl` 为本机安装的 Oracle 数据库软件或客户端软件定义的本地或远程 Oracle 数据库的标识，`test` 和 `oracle` 是此数据库的用户名和密码，`testtab` 是此数据库中的一个表。

为了在 R 中访问 Oracle 数据库服务器中的数据库，在 R 中需要安装 ROracle 包。这是一个源代码扩展包，需要用户自己编译安装。在 MS Windows 环境下，需要安装 R 软件和 RTools 软件包（在 CRAN 网站的 Windows 版本软件下载栏目中）。在 MS Windows 命令行窗口，用如下命令编译 R 的 ROracle 扩展包：

```
set OCI_LIB32=D:\oracle\product\10.2.0\db_1\bin
set OCI_INC=D:\oracle\product\10.2.0\db_1\oci\include
set PATH=D:\oracle\product\10.2.0\db_1\bin;C:\Rtools\bin;C:\Rtools\gcc-4.6.3\bin;"
C:\R\R-3.2.0\bin\i386\rcmd INSTALL ROracle_1.2-1.tar.gz
```

其中的前三个 `set` 命令设置了 Oracle 数据库程序或客户端程序链接库、头文件和可执行程序的位置，第三个 `set` 命令还设置了 RTools 编译器的路径。这些路径需要根据实际情况修改。这里的设置是在本机运行的 Oracle 10g 服务器软件的情况。最后一个命令编译 ROracle 扩展包，相应的 `rcmd` 程序路径需要改成自己的安装路径。

如果服务器在远程服务器上，设远程服务器的数据库标识名为 `ORCL`，本机

需要安装客户端 Oracle instant client 软件, 此客户端软件需要与服务器同版本号, 如 `instantclient-basic-win32-10.2.0.5.zip`, 这个软件不需要安装, 只需要解压到一个目录如 `C:\instantclient_10_2` 中。在本机 (以 MS Windows 操作系统为例) 中, 双击系统, 选择高级-环境变量, 增加如下三个环境变量:

```
NLS_LANG = SIMPLIFIED CHINESE_CHINA.ZHS16GBK
ORACLE_HOME = C:\instantclient_10_2
TNS_ADMIN = C:\instantclient_10_2
```

并在环境变量 PATH 的值的末尾增加 Oracle 客户端软件所在的目录 `verb|C:\instantclient_10_2`, 并与前面内容用分号分开。

然后, 在 client 所在的目录 `C:\instantclient_10_2` 中增加如下内容的 `tnsnames.ora` 文件

```
orcl =
(DESCRIPTION =
  (ADDRESS = (PROTOCOL = TCP)(HOST = 192.168.1.102 )
    (PORT = 1521))
  (CONNECT_DATA =
    (SERVER = DEDICATED)
    (SERVICE_NAME = orcl)
  )
)
```

其中 HOST 的值是安装 Oracle 服务器的服务器的 IP 地址, orcl 是一个服务器实例名, 能够在服务器端的 `tnsnames.ora` 文件中查到, 等号前面的 orcl 是对数据库给出的客户端别名, 这里就干脆用了和服务器端的数据库标识名相同的名字 orcl。

不论是在本机的数据框服务器还是在本机安装设置好客户端后, 在 R 中用如下的程序可以读入数据库中的表:

```
require(ROracle)
drv <- dbDriver("Oracle")

conn <- dbConnect(drv, username="test",
```

```
password="oracle", dbname="orcl")  
  
rs <- dbSendQuery(conn, "select * from testtab")  
d <- fetch(rs)
```

可以用 `dbGetTable()` 取出一个表并存入 R 数据框中。用 `dbSendQuery()` 发出一个 SQL 命令，用 `fetch()` 可以一次性取回或者分批取回，在表行数很多时这种方法更适用。

15.4.2 MySQL 数据库访问

MySQL 是高效、免费的数据库服务器软件，在很多行业尤其是互联网行业占有很大的市场。为了在 R 中访问 MySQL 数据库，只要安装 RMySQL 扩展包（有二进制版本）。假设服务器地址在 192.168.1.111，可访问的数据库名为 world，用户为 test，密码为 mysql。设 world 库中有表 country。

在 R 中要访问 MySQL 数据框，首先要建立与数据库服务器的连接：

```
con <- dbConnect(RMySQL::MySQL(),  
  dbname='world',  
  username='test', password='mysql',  
  host='192.168.1.111')
```

下列代码列出 world 库中的所有表，然后列出其中的 country 表的所有变量：

```
dbListTables(con)  
dbListFields(con, 'country')
```

下列代码取出 country 表并存入 R 数据框 d.country 中：

```
d.country <- dbReadTable(con, 'country')
```

下列代码把 R 中的示例数据框 USArrests 写入 MySQL 库 world 的表 arrests 中：

```
data(USArrests)  
dbWriteTable(con, 'arrests', USArrests,  
  overwrite=TRUE)
```

当然，这需要用户对该库有写权限。

可以用 `dbGetQuery()` 执行一个 SQL 查询并返回结果，如

```
dbGetQuery(con, 'select count(*) from arrests')
```

当表很大时，可以用 `dbSendQuery()` 发送一个 SQL 命令，返回一个查询结果指针对象，用 `dbFetch()` 从指针对象位置读取指定行数，用 `dbHasCompleted()` 判断是否已读取结束。如

```
res <- dbSendQuery(con, "SELECT * FROM country")
while(!dbHasCompleted(res)){
  chunk <- dbFetch(res, n = 5)
  print(chunk[,1:2])
}
dbClearResult(res)
```

数据库使用完毕时，需要关闭用 `dbConnect()` 打开的连接：

```
dbDisconnect(con)
```

15.5 文件访问

15.5.1 连接

输入输出可以针对命令行，针对文件，R 支持扩展的文件类型，称为“连接 (connection)”。

函数 `file()` 生成到一个普通文件的连接，函数 `url()` 生成一个到指定的 URL 的连接，函数 `gzfile`, `bzfile`, `xzfile`, `unz` 支持对压缩过的文件的访问（不是压缩包，只对一个文件压缩）。这些函数大概的用法如下：

```
file("path", open="", blocking=T,
     encoding = getOption("encoding"),
     raw = FALSE)

url(description, open = "", blocking = TRUE,
```

```
encoding = getOption("encoding"))

textConnection(description, open="r",
  local = FALSE,
  encoding = c("", "bytes", "UTF-8"))

gzfile(description, open = "",
  encoding = getOption("encoding"),
  compression = 6)

bzfile(description, open = "",
  encoding = getOption("encoding"),
  compression = 9)

xzfile(description, open = "",
  encoding = getOption("encoding"),
  compression = 6)

unz(description, filename, open = "",
  encoding = getOption("encoding"))
```

生成连接的函数不自动打开连接。给定一个未打开的连接，读取函数从中读取时会自动打开连接，函数结束时自动关闭连接。用 `open()` 函数打开连接，返回一个句柄；生成连接时可以用 `open` 参数要求打开连接。要多次从一个连接读取时就应该先打开连接，读取完毕用 `close` 函数关闭。

函数 `textConnection()` 打开一个字符串用于读写。

在生成连接与打开连接的函数中用 `open` 参数指定打开方式，取值为：

- `r`—文本型只读；
- `w`—文本型只写；
- `a`—文本型末尾添加；
- `rb`—二进制只读；
- `wb`—二进制只写；

- `ab`—二进制末尾添加;
- `r+` 或 `r+b`—允许读和写;
- `w+` 或 `w+b`—允许读和写, 但刚打开时清空文件;
- `a+` 或 `a+b`—末尾添加并允许读。

15.5.2 文本文件访问

函数 `readLines()`, `scan()` 可以从一个文本型连接读取。

给定一个打开的连接 `con`, 用 `readLines` 函数可以把文件各行读入为字符型向量的各个元素, 不包含文件中用来分开各行的换行标志。可以指定要读的行数。如

```
ll <- readLines(file('class.csv'))
print(ll)
```

用 `writeLines` 函数可以把一个字符型向量各元素作为不同行写入一个文本型连接。如

```
vnames <- strsplit(ll, ',')[[1]]
writeLines(vnames, con='class-names.txt')
```

其中的 `con` 参数应该是一个打开的文本型写入连接, 但是可以直接给出一个要写入的文件名。

15.5.3 二进制文件访问

函数 `save` 用来保存 R 变量到文件, 函数 `load` 用来从文件中读取保存的 R 变量。

函数 `readBin` 和 `writeBin` 对 R 变量进行二进制文件存取。

如果要访问其它软件系统的二进制文件, 请参考 R 手册中的“R Data Import/Export Manual”。

15.5.4 字符型连接

函数 `textConnection` 打开一个字符串用于读取或写入，是很好用的一个 R 功能。可以把一个小文件存放在一个长字符串中，然后用 `textConnection` 读取，如

```
fstr <-  
"name,score  
王芳,78  
孙莉,85  
张聪,80  
"  
d <- read.csv(textConnection(fstr), header=T)  
print(d)
```

读取用的 `textConnection` 的参数是一个字符型变量。

在整理输出结果时，经常可以向一个字符型变量连接写入，最后再输出整个字符串值。例如：

```
tc <- textConnection("sres", open="w")  
cat('Trial of text connection.\n', file=tc)  
cat(1:10, '\n', file=tc, append=T)  
close(tc)  
print(sres)
```

注意写入用的 `textConnection` 的第一个参数是保存了将要写入的字符型变量名的字符串，而不是变量名本身，第二个参数表明是写入操作，使用完毕需要用 `close` 关闭。

15.6 中文编码问题

读写文本格式的数据，或者用 `readLines()`、`readr::read_lines()` 读写文本文件，可能会遇到中文编码不匹配的问题。这里总结一些常用解决方法，所用的操作系统为中文 Windows10，在 RStudio 中运行，R 版本为 3.4.3。常见的中文编码有 GBK(或 GB18030, GB)，UTF-8，UTF-8 有 BOM 标志等。

可以用 `iconvlist()` 查看 R 支持的编码名称。

假设有如下的含有中文的文件：

```
序号,收缩压
1,145
5,110
6,未测
9,150
10,拒绝
15,115
```

这个文件是在中文版 MS Office 的 Excel 软件中输入后，用 Office 的“文件——另存为——.csv 格式”生成的，结果的编码是 GBK 编码，或 GB18030 编码。文件下载：[bp.csv](#)

我们用工具软件将其转换成 UTF-8 无 BOM 格式，下载链接：[bp-utf8nobom.csv](#)

转为 UTF-8 有 BOM 格式，下载链接：[bp-utf8bom.csv](#)

15.6.1 用基本 R 的读取函数读取

与所用操作系统默认编码相同的文本文件，R 基本软件的 `read.csv()`、`read.table()`、`readLines()` 函数都可以正常读取，所以 `bp.csv` 文件可以正常读取，如

```
read.csv("bp.csv")
```

```
##  序号 收缩压
## 1    1   145
## 2    5   110
## 3    6 未测
## 4    9   150
## 5   10 拒绝
## 6   15   115
```

```
readLines("bp.csv")
```

```
## [1] "序号,收缩压" "1,145"          "5,110"          "6, 未测"        "9,150"
## [6] "10, 拒绝"      "15,115"
```

但是另外两个以 UTF-8 编码的文件则不能正确读入:

```
read.csv("bp-utf8nobom.csv")
```

```
## Error in make.names(col.names, unique = TRUE) : invalid multibyte string 2
```

```
readLines("bp-utf8bom.csv")
```

```
## [1] "锃垮箬莖\xb7, 鎗剝緝劍\x8b" "1,145"
## [3] "5,110"                          "6, 鏈 械"
## [5] "9,150"                            "10, 鎬 擊 糲"
## [7] "15,115"
```

读取 UTF-8 编码无 BOM 的文件时, 在 `read.csv()` 和 `read.table()` 等函数中加 `fileEncoding="UTF-8"` 选项可以纠正编码问题:

```
read.csv("bp-utf8nobom.csv", fileEncoding="UTF-8")
```

```
##   序号 收缩压
## 1    1    145
## 2    5    110
## 3    6   未测
## 4    9    150
## 5   10   拒绝
## 6   15   115
```

读取 UTF-8 编码无 BOM 或者有 BOM 的文件时, 在 `readLines()` 函数中加 `encoding="UTF-8"` 选项可以纠正编码问题:

```
readLines("bp-utf8nobom.csv", encoding="UTF-8")
```

```
## [1] "序号,收缩压" "1,145"          "5,110"          "6, 未测"        "9,150"
## [6] "10, 拒绝"      "15,115"
```

```
readLines("bp-utf8bom.csv", encoding="UTF-8")
```

```
## [1] "<U+FEFF>序号,收缩压" "1,145"      "5,110"      "6, 未测"      "9,150"
## [6] "10, 拒绝"      "15,115"
```

但是, UTF-8 有 BOM 标志的文本文件不能被 `read.csv()` 识别:

```
read.csv("bp-utf8bom.csv", fileEncoding="UTF-8")
```

```
## invalid input found on input connection 'bp-utf8bom.csv'
## incomplete final line found by readTableHeader on 'bp-utf8bom.csv'
```

15.6.2 用 readr 包读取

readr 包的 `read_csv()`、`read_table2()`、`read_lines()` 函数默认从 UTF-8 编码的文件中读取, 无 BOM 或者有 BOM 都可以。如:

```
read_csv("bp-utf8nobom.csv")
```

```
## Parsed with column specification:
```

```
## cols(
##   序号 = col_double(),
##   收缩压 = col_character()
## )
```

```
## # A tibble: 6 x 2
```

```
##   序号 收缩压
##   <dbl> <chr>
## 1     1 145
## 2     5 110
## 3     6 未测
## 4     9 150
## 5    10 拒绝
## 6    15 115
```

```
read_csv("bp-utf8bom.csv")
```

```
## Parsed with column specification:
```

```
## cols(
##   序号 = col_double(),
##   收缩压 = col_character()
## )
```

```
## # A tibble: 6 x 2
##   序号 收缩压
##   <dbl> <chr>
## 1     1 145
## 2     5 110
## 3     6 未测
## 4     9 150
## 5    10 拒绝
## 6    15 115
```

```
read_lines("bp-utf8nobom.csv")
```

```
## [1] "序号,收缩压" "1,145"      "5,110"      "6, 未测"    "9,150"
## [6] "10, 拒绝"    "15,115"
```

```
read_lines("bp-utf8bom.csv")
```

```
## [1] "序号,收缩压" "1,145"      "5,110"      "6, 未测"    "9,150"
## [6] "10, 拒绝"    "15,115"
```

但是,对 GBK 编码的文件,不能直接读取:

```
read_csv("bp.csv")
```

```
read_lines("bp.csv")
```

```
## [1] "<d0><f2><U+00BA><c5>,<ca><d5><cb><f5><U+0479>" "1,145"
## [3] "5,110"                                     "6, <U+00B2><e2>"
## [5] "9,150"                                     "10, <U+00BE><U+073E><f8>"
## [7] "15,115"
```

为了读取 GBK(或 GB18030) 编码的文件,在 `read_csv()` 和 `read_lines()` 函数中加入 `locale=locale(encoding="GBK")` 选项:

```
read_csv("bp.csv", locale=locale(encoding="GBK"))
```

```
## Parsed with column specification:
## cols(
##   序号 = col_double(),
##   收缩压 = col_character()
## )
## # A tibble: 6 x 2
##   序号 收缩压
##   <dbl> <chr>
## 1     1 145
## 2     5 110
## 3     6 未测
## 4     9 150
## 5    10 拒绝
## 6    15 115
```

```
read_lines("bp.csv", locale=locale(encoding="GBK"))
```

```
## [1] "序号,收缩压" "1,145"          "5,110"          "6, 未测"        "9,150"
## [6] "10, 拒绝"     "15,115"
```

15.6.3 输出文件的编码

`write.csv()`、`writeLines()` 生成的含有中文的文件的编码默认为操作系统的默认中文编码，这里是 GB18030。

`readr` 的 `write_csv()`、`write_lines()` 函数生成的含有中文的文件的编码默认 UTF-8 无 BOM。如

```
write_csv(tibble("姓名"=c("张三", "李四")), "tmp.csv")
```

结果生成的文件编码为 UTF-8 无 BOM，这样的文件可以被 R 的 `readr::read_csv()` 正确读取，但是不能被 MS Excel 软件正确读取。

`write_lines()` 输出的文件也是编码为 UTF-8 无 BOM。

`write_excel_csv()` 可以生成带有 UTF-8 有 BOM 的 CSV 文件, 这样的文件可以被 MS Office 正确识别:

```
write_excel_csv(tibble("姓名"=c("张三", "李四")), "tmp2.csv")
```

15.6.4 分批读写

`readLines()`、`readr::read_lines()`、`writeLines()`、`readr::writeLines()` 支持分批读写。这需要预先打开要读取和写入的文件, 所有内容都处理一遍以后关闭读取和写入的文件。

使用 `file()` 函数打开文件用于读写, 使用 `close()` 函数关闭打开的文件。打开文件时可以用 `encoding=` 指定编码, 但是 `readr::read_lines()` 不支持分批读入。

下面的程序每次从 UTF-8 无 BOM 编码的 `bp-utf8nobom.csv` 读入两行, 不加处理第写入 `tmp.csv` 中, 使用 `readLines()` 和 `writeLines()`, 读入时用 `encoding="UTF-8"` 指定编码, 写出时不指定编码, 结果是操作系统默认的 GBK:

```
fin <- file("bp-utf8nobom.csv", "rt", encoding="UTF-8")
fout <- file("tmp.csv", "wt")
repeat{
  lines <- readLines(fin, n=2)
  print(lines)
  if(length(lines)==0) break
  writeLines(lines, fout)
}
close(fout)
close(fin)
## [1] " 序号, 收缩压" "1,145"
## [1] "5,110" "6, 未测"
## [1] "9,150" "10, 拒绝"
## [1] "15,115"
## character(0)
```


`file()` 中的 `encoding="UTF-8"` 特指 UTF-8 无 BOM 的格式, 有 BOM 的 UTF-8 编码文件无法用上述方法打开。

上面的例子生成的结果 `tmp.csv` 使用了中文 Windows 系统的默认编码 GBK 编码。为了生成 UTF-8 无 BOM 的结果, 可以在上述程序中打开输出文件时加选项 `encoding="UTF-8"`。即

```
fin <- file("bp-utf8nobom.csv", "rt", encoding="UTF-8")
fout <- file("tmp.csv", "wt", encoding="UTF-8")
...
```

`readr::read_lines()` 不支持从一个文件分批读入。`readr::write_lines()` 可以用 `append=TRUE` 选项向一个文件分批写出。

15.7 目录和文件管理

目录和文件管理函数:

- `getwd()`—返回当前工作目录。
- `setwd(path)`—设置当前工作目录。
- `list.files()` 或 `dir()`—查看目录中内容。`list.files(pattern='.*[.]r$')` 可以列出所有以“r”结尾的文件。
- `file.path()`—把目录和文件名组合得到文件路径。
- `file.info(filename)`—显示文件的详细信息。
- `file.exists()`—查看文件是否存在。
- `file.access()`—考察文件的访问权限。
- `create.dir()`—新建目录。
- `file.create()`—生成文件。
- `file.remove()` 或 `unlink()`—删除文件。`unlink()` 可以删除目录。
- `file.rename()`—为文件改名。
- `file.append()`—把两个文件相连。
- `file.copy()`—复制文件。
- `basename()` 和 `dirname()`— 从一个全路径文件名获取文件名和目录。

Chapter 16

程序控制结构

16.1 表达式

R 是一个表达式语言, 其任何一个语句都可以看成是一个表达式。表达式之间以分号分隔或用换行分隔。表达式可以续行, 只要前一行不是完整表达式 (比如末尾是加减乘除等运算符, 或有未配对的括号) 则下一行为上一行的继续。若干个表达式可以放在一起组成一个复合表达式, 作为一个表达式使用, 复合表达式的值为最后一个表达式的值, 组合用大括号表示, 如:

```
{  
  x <- 15  
  x  
}
```

16.2 分支结构

分支结构包括 if 结构:

if (条件) 表达式 1

或

if (条件) 表达式 1 else 表达式 2

其中的“条件”为一个标量的真或假值，表达式可以用大括号包围的复合表达式。如

```
if(is.na(lambda)) lambda <- 0.5
```

又如

```
if(x>1) {  
  y <- 2.5  
} else {  
  y <- -y  
}
```

多个分支，可以在中间增加 else if，如：

```
x <- c(0.05, 0.6, 0.3, 0.9)  
for(i in seq(along=x)){  
  if(x[i] <= 0.2){  
    cat("Small\n")  
  } else if(x[i] <= 0.8){  
    cat("Medium\n")  
  } else {  
    cat("Large\n")  
  }  
}
```

16.2.1 用逻辑下标代替分支结构

R 是向量化语言，尽可能少用标量运算。比如，x 为一个向量，要定义 y 与 x 等长，且 y 的每一个元素当且仅当 x 的对应元素为正数时等于 1，否则等于零。

这样是错误的：

```
if(x>0) y <- 1 else y <- 0
```

正解为：

```
y <- numeric(length(x))
y[x>0] <- 1
y
```

函数 `ifelse()` 可以根据一个逻辑向量中的多个条件，分别选择不同结果。如

```
x <- c(-2, 0, 1)
y <- ifelse(x >=0, 1, 0); print(y)
```

```
## [1] 0 1 1
```

函数 `switch()` 可以建立多分枝结构。

16.3 循环结构

16.3.1 计数循环

为了对向量每个元素、矩阵每行、矩阵每列循环处理，语法为

`for(循环变量 in 序列) 语句`

其中的语句一般是复合语句。如：

```
x <- rnorm(5)
y <- numeric(length(x))
for(i in 1:5){
  if(x[i]>=0) y[i] <- 1 else y[i] <- 0
}
print(y)
```

```
## [1] 1 0 0 0 0
```

其中 `rnorm(5)` 会生成 5 个标准正态分布随机数。`numeric(n)` 生成有 `n` 个 0 的数值型向量（基础类型为 `double`）。

如果需要对某个向量 `x` 按照下标循环，获得所有下标序列的标准写法是 `seq(along=x)`，而不用 `1:n` 的写法，因为在特殊情况下 `n` 可能等于零，这会导致错误下标，而 `seq(along=x)` 在 `x` 长度为零时返回零长度的下标。

例如，设序列 x_n 满足 $x_0 = 0$, $x_n = 2x_{n-1} + 1$, 求 $S_n = \sum_{i=1}^n x_n$:

```
x <- 0.0
s <- 0
n <- 5
for(i in 1:n){
  x <- 2*x + 1
  s <- s + x
}
print(s)
```

```
## [1] 57
```

在 R 中应尽量避免 for 循环：其速度比向量化版本慢一个数量级以上，而且写出的程序不够典雅。比如，前面那个示性函数例子实际上可以简单地写成

```
x <- rnorm(5)
y <- ifelse(x >= 0, 1, 0)
print(y)
```

```
## [1] 0 0 0 0 0
```

16.3.2 while 循环和 repeat 循环

用

```
while(循环继续条件) 语句
```

进行当型循环。其中的语句一般是复合语句。仅当条件成立时才继续循环，而且如果第一次条件就已经不成立就一次也不执行循环内的语句。

用

```
repeat 语句
```

进行无条件循环（一般在循环体内用 if 与 break 退出）。其中的语句一般是复合语句。如下的写法可以制作一个直到型循环：

```
repeat{
  ...
  if(循环退出条件) break
}
```

直到型循环至少执行一次，每次先执行... 代表的循环体语句，然后判断是否满足循环退出条件，满足条件就退出循环。

用 `break` 语句退出所在的循环。用 `next` 语句进入所在循环的下一轮。

例如，常量 e 的值可以用泰勒展开式表示为

$$e = 1 + \sum_{k=1}^{\infty} \frac{1}{k!}$$

R 函数 `exp(1)` 可以计算 e 的值，下面用泰勒展开逼近计算 e 的值：

```
e0 <- exp(1.0)
s <- 1.0
x <- 1
k <- 0
repeat{
  k <- k+1
  x <- x/k
  s <- s + x

  if(x < .Machine$double.eps) break
}
err <- s - e0
cat("k=", k, " s=", s, " e=", e0, " 误差 =", err, "\n")
```

```
## k= 18  s= 2.718282  e= 2.718282  误差= 4.440892e-16
```

其中 `.Machine$double.eps` 称为机器 ϵ ，是最小的加 1 之后可以使得结果大于 1 的正双精度数，小于此数的正双精度数加 1 结果还等于 1。用泰勒展开公式计算的结果与 `exp(1)` 得到的结果误差在 10^{-16} 左右。

16.4 R 中判断条件

if 语句和 while 语句中用到条件。条件必须是标量值，而且必须为 TRUE 或 FALSE，不能为 NA 或零长度。这是 R 编程时比较容易出错的地方。

16.5 管道控制

数据处理中经常会对同一个变量（特别是数据框）进行多个步骤的操作，比如，先筛选部分有用的变量，再定义若干新变量，再排序。R 的 `magrittr` 包提供了一个 `%>%` 运算符实现这样的操作流程。比如，变量 `x` 先用函数 `f(x)` 进行变换，再用函数 `g(x)` 进行变换，一般应该写成 `g(f(x))`，用 `%>%` 运算符，可以表示成 `x %>% f() %>% g()`。更多的处理，如 `h(g(f(x)))` 可以写成 `x %>% f() %>% g() %>% h()`。这样的表达更符合处理发生的次序，而且插入一个处理步骤也很容易。

处理用的函数也可以带有其它自变量，在管道控制中不要写第一个自变量。某个处理函数仅有一个自变量时，可以省略空的括号。

`tibble` 类型的数据框尤其适用于如此的管道操作。

将管道控制开始变量设置为 `.`，可以定义一个函数。

`magrittr` 包定义了 `%T%` 运算符，`x %T% f()` 返回 `x` 本身而不是用 `f()` 修改后的返回值 `f(x)`，这在中间步骤需要显示或者绘图但是需要进一步对输入数据进行处理时有用。

`magrittr` 包定义了 `%%$%` 运算符，此运算符的作用是将左运算元的各个变量（这时左运算元是数据框或列表）暴露出来，可以直接在右边调用其中的变量，类似于 `with()` 函数的作用。

`magrittr` 包定义了 `%<>%` 运算符，用在管道链的第一个连接，可以将处理结果存入最开始的变量中，类似于 C 语言的 `+=` 运算符。

如果一个操作是给变量加 `b`，可以写成 `add(b)`，给变量乘 `b`，可以写成 `multiply_by(b)`。

Chapter 17

函数

17.1 函数基础

17.1.1 介绍

在现代的编程语言中使用自定义函数，优点是代码复用、模块化设计。

在编程时，把编程任务分解成小的模块，每个模块用一个函数实现，可以降低复杂性，防止变量混杂。

函数的自变量是只读的，函数中定义的局部变量只在函数运行时起作用，不会与外部或其它函数中同名变量混杂。

函数返回一个对象作为输出，如果需要返回多个变量，可以用列表进行包装。

17.1.2 函数定义

函数定义使用 `function` 关键字，一般格式为

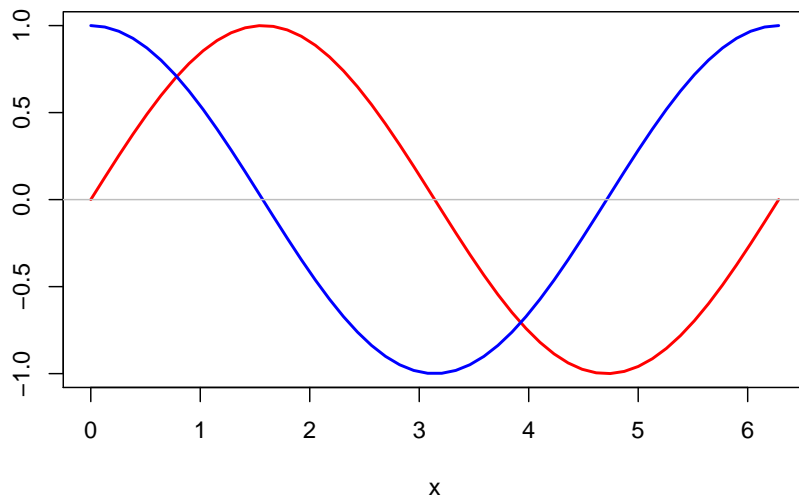
```
函数名 <- function(形式参数表) 函数体
```

函数体是一个表达式或复合表达式（复合语句），以复合表达式中最后一个表达式为返回值，也可以用 `return(x)` 返回 `x` 的值。如果函数需要返回多个结果，可以打包在一个列表（list）中返回。形式参数表相当于函数自变量，可以是空

的，形式参数可以有缺省值，R 的函数在调用时都可以用“形式参数名 = 实际参数”的格式输入自变量值。

下面的例子没有参数，仅画一个示例图：

```
f <- function() {  
  x <- seq(0, 2*pi, length=50)  
  y1 <- sin(x)  
  y2 <- cos(x)  
  plot(x, y1, type='l', lwd=2, col='red',  
        xlab='x', ylab='')  
  lines(x, y2, lwd=2, col='blue')  
  abline(h=0, col='gray')  
}  
f()
```



注意此自定义函数虽然没有参数，但是在定义与调用时都不能省略圆括号。

自定义函数也可以是简单的一元函数，与数学中一元函数基本相同，例如

```
f <- function(x) 1/sqrt(1 + x^2)
```

基本与数学函数 $f(x) = 1/\sqrt{1+x^2}$ 相对应。定义中的自变量 x 叫做形式参数或形参 (formal arguments)。函数调用时，形式参数得到实际值，叫做实参 (actual arguments)。R 函数有一个向量化的好处，在上述函数调用时，如果形式参数 x 的实参是一个向量，则结果也是向量，结果元素为实参向量中对应元素的变换值。如

```
f(0)
```

```
## [1] 1
```

```
f(c(-1, 0, 1, 2))
```

```
## [1] 0.7071068 1.0000000 0.7071068 0.4472136
```

第一次调用时，形式参数 x 得到实参 0，第二次调用时，形式参数 x 得到向量实参 $c(-1, 0, 1, 2)$ 。

函数实参是向量时，函数体中也可以计算对向量元素进行汇总统计的结果。例如，设 x_1, x_2, \dots, x_n 是一个总体的简单随机样本，其样本偏度统计量定义如下：

$$\hat{w} = \frac{n}{(n-1)(n-2)} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{S} \right)^3$$

其中 \bar{x} 与 S 分别是样本均值与样本标准差。如下的 R 函数可以把观测样本的值保存在一个向量中输入，计算并输出其样本偏度统计量值：

```
f <- function(x) {
  n <- length(x)
  xbar <- mean(x)
  S <- sd(x)
  n/(n-1)/(n-2)*sum( (x - xbar)^3 ) / S^3
}
```

函数体的最后一个表达式是函数返回值。

在函数体最后一个表达式中巧妙地利用了 R 的向量化运算 $(x - xbar)^3$ 与内建函数 (sum) 。这比用 `for` 循环计算效率高得多，计算速度相差几十倍。

请比较如下两个表达式：

```
n/(n-1)/(n-2)*sum( (x - xbar)^3 ) / S^3
n/(n-1)/(n-2)*sum( ((x - xbar)/S)^3 )
```

这两个表达式的值相同。表面上看，第二个表达式更贴近原始数学公式，但是在编程时，需要考虑计算效率问题，第一个表达式关于 S 只需要除一次，而第二个表达关于 S 除了 n 次，所以第一个表达式效率更高。

函数定义中的形式参数可以有多个，还可以指定缺省值。例如

```
fsub <- function(x, y=0){
  cat("x=", x, " y=", y, "\n")
  x - y
}
```

这里 x, y 是形式参数，其中 y 指定了缺省值为 0，有缺省值的形式参数在调用时可以省略对应的实参，省略时取缺省值。

实际上，“`function(参数表) 函数体`”这样的结构本身也是一个表达式，其结果是一个函数对象。在通常的函数定义中，函数名只不过是赋值给某个函数对象，或者说是“绑定”(bind)到某个函数对象上面。R 允许使用没有函数名的函数对象。

因为函数也是 R 对象，也可以拥有属性。所谓对象，就是 R 的变量所指向的各种不同类型的统称。

一个自定义 R 函数由三个部分组成：函数体 `body()`，即要函数定义内部要执行的代码；`formals()`，即函数的形式参数表以及可能存在的缺省值；`environment()`，是函数定义时所处的环境，这会影响到参数表中缺省值与函数体中非局部变量的查找。注意，函数名并不是函数对象的必要组成部分。如

```
body(fsub)
```

```
## {
##   cat("x=", x, " y=", y, "\n")
##   x - y
## }
```

```
formals(fsub)
```

```
## $x  
##  
##  
## $y  
## [1] 0
```

```
environment(fsub)
```

```
## <environment: R_GlobalEnv>
```

“环境”是 R 语言比较复杂的概念，后面再详细解释。

17.1.3 函数调用

函数调用时最基本的调用方式是把实参与形式参数按位置对准，这与我们在数学中使用多元函数的习惯类似。例如

```
fsub(3, 1)
```

```
## x= 3  y= 1  
## [1] 2
```

相当于以 $x=3, y=1$ 调用。

调用时可选参数可以省略实参，如

```
fsub(3)
```

```
## x= 3  y= 0  
## [1] 3
```

相当于以 $x=3, y=0$ 调用。

R 函数调用时全部或部分形参对应的实参可以用“形式参数名 = 实参”的格式给出，这样格式给出的实参不用考虑次序，不带形式参数名的则按先后位置对准。如

```
fsub(x=3, y=1)
## x= 3 y= 1
## [1] 2
fsub(y=1, x=3)
## x= 3 y= 1
## [1] 2
fsub(x=3)
## x= 3 y= 0
## [1] 3
fsub(3, y=1)
## x= 3 y= 1
## [1] 2
fsub(1, x=3)
## x= 3 y= 1
## [1] 2
fsub(x=3, 1)
## x= 3 y= 1
## [1] 2
```

注意作为好的程序习惯应该避免 `fsub(x=3, 1)` 这样的做法。虽然 R 的语法没有强行要求，调用 R 函数时，如果既有按位置对应的参数又有带名参数，按位置对应的参数都写在前面，带名参数写在后面，不遵守这样的约定容易使得程序被误读。

R 的形参、实参对应关系可以写成一个列表，如 `fsub(3, y=1)` 中的对应关系可以写成列表 `list(3, y=1)`，如果调用函数的形参、实参对应关系保存在列表中，可以用函数 `do.call()` 来表示函数调用，如

```
do.call(fsub, list(3, y=1))
```

与

```
fsub(3, y=1)
```

效果相同。

在自定义 R 函数的形参中，还允许有一个特殊的... 形参（三个小数点）。在

函数调用时，所有没有形参与之匹配的实参，不论是带有名字还是不带有名字的，都自动归入这个参数，这个参数的类型是一个列表。虽然很奇怪，这个语法在 R 里面是常用的，通常用来把函数内调用的其它函数的实参传递进来。

例如，`sapply(X, FUN, ...)` 中的形式参数 `FUN` 需要函数实参，此函数有可能需要更多的参数。例如，为了把 `1:5` 的每个元素都减去 2，可以写成

```
sapply(1:5, fsub, y=2)
```

```
## x= 1  y= 2
## x= 2  y= 2
## x= 3  y= 2
## x= 4  y= 2
## x= 5  y= 2

## [1] -1  0  1  2  3
```

或

```
sapply(1:5, fsub, 2)
```

```
## x= 1  y= 2
## x= 2  y= 2
## x= 3  y= 2
## x= 4  y= 2
## x= 5  y= 2

## [1] -1  0  1  2  3
```

实际上，R 语法中的大多数运算符如 `+`，`-`，`*`，`/`，`[`，`[[`，`(`，`{`等都是函数。这些特殊名字的函数要作为函数使用，需要使用反向单撇号‘包围，比如

```
1 + 2
```

```
## [1] 3
```

```
`+`(1, 2)
```

```
## [1] 3
```

效果相同。

这样，为了给 1:5 每个元素减去 2，还可以写成

```
sapply(1:5, `-`, 2)
```

```
## [1] -1  0  1  2  3
```

或

```
sapply(1:5, "-", 2)
```

```
## [1] -1  0  1  2  3
```

在上一写法中 `sapply` 的第二参数用了函数名字符串作为实参。

17.2 变量作用域

17.2.1 全局变量和工作空间

在所有函数外面（如 R 命令行）定义的变量是全局变量。在命令行定义的所有变量都保存在工作空间（workspace）中。用 `ls()` 查看工作空间内容。`ls()` 中加上 `pattern` 选项可以指定只显示符合一定命名模式的变量，如

```
ls(pattern='^tmp[.]')
```

显示所有以 `tmp.` 开头的变量。用 `object.size()` 函数查看变量占用存储大小。

因为 R 的函数调用时可以读取工作空间中的全局变量值，工作空间中过多的变量会引起莫名其妙的程序错误。用 `rm()` 函数删除指定的变量。`rm()` 中还可以用 `list` 参数指定一个要删除的变量名表。如

```
rm(list=ls(pattern='^tmp[.]'))
```

用 `save()` 函数保存工作空间中选择的某些变量；用 `load()` 函数载入保存在文件中的变量。如

```
save(my.large.data,  
      file='my-large-data.RData')  
load('my-large-data.RData')
```


实际上，R 的工作空间是 R 的变量搜索路径中的一层，大体相当于全局变量空间。R 的已启用的软件包中的变量以及用 `attach()` 命令引入的变量也在这个搜索路径中。

17.2.2 局部变量

在计算机语言中，“变量”实际是计算机内存中的一段存储空间。函数的参数（自变量）在定义时并没有对应的存储空间，所以也称函数定义中的参数为“形式参数”。

函数的形式参数在调用时被赋值为实参值（这是一般情形），形参变量和函数体内被赋值的变量都是局部的。这一点符合函数式编程 (functional programming) 的要求。所谓局部变量，就是仅在函数运行时才存在，一旦退出函数就不存在的变量。

17.2.2.1 自变量的局部性

在函数被调用时，形式参数（自变量）被赋值为实际的值（称为实参），如果实参是变量，形式参数实际变成了实参的一个副本，在函数内部对形式参数作任何修改在函数运行完成后都不影响原来的实参变量，而且函数运行完毕后形式参数对应的变量不再存在。

在下例中，在命令行定义了全局变量 `xv`, `x1`，然后作为函数 `f()` 的自变量值（实参）输入到函数中，函数中对两个形式参数作了修改，函数结束后实参变量 `xv`, `x1` 并未被修改，形参变量也消失了。例子程序如下：

```
xv <- c(1,2,3)
x1 <- list(a=11:15, b='James')
if(exists("x")) rm(x)
f <- function(x, y){
  cat(' 输入的 x=', x, '\n')
  x[2] <- -1
  cat(' 函数中修改后的 x=', x, '\n')
  cat(' 输入的 y 为:\n'); print(y)
  y[[2]] <- 'Mary'
```

```
    cat(' 函数中修改过的 y 为:\n'); print(y)
}
f(xv, xl)
## 输入的 x= 1 2 3
## 函数中修改后的 x= 1 -1 3
## 输入的 y 为:
## $a
## [1] 11 12 13 14 15
##
## $b
## [1] "James"
##
## 函数中修改过的 y 为:
## $a
## [1] 11 12 13 14 15
##
## $b
## [1] "Mary"
##
cat(' 函数运行完毕后原来变量 xv 不变: ', xv, '\n')
## 函数运行完毕后原来变量 xv 不变: 1 2 3
cat(' 函数运行完毕后原来变量 xl 不变: :\n'); print(xl)
## 函数运行完毕后原来变量 xl 不变: :
## $a
## [1] 11 12 13 14 15
##
## $b
## [1] "James"
##
cat(' 函数运行完毕后形式参数 x 不存在: :\n'); print(x)
## 函数运行完毕后形式参数 x 不存在: :
## Error in print(x) : object 'x' not found
```

R 语言的这种特点对于传递超大的数据是不利的，所以 R 中会容纳超大数据的类型往往涉及成修改副本时不占用不必要的额外存储空间，比如，tibble 类型就有这样的特点。

17.2.2.2 修改自变量

为了修改某个自变量，在函数内修改其值并将其作为函数返回值，赋值给原变量。

比如定义了如下函数：

```
f <- function(x, inc=1){  
  x <- x + inc  
  x  
}
```

调用如

```
x <- 100  
cat(' 原始 x=', x, '\n')
```

```
## 原始 x= 100
```

```
x <- f(x)  
cat(' 修改后 x=', x, '\n')
```

```
## 修改后 x= 101
```

17.2.2.3 函数内的局部变量

在函数内部用赋值定义的变量都是局部变量，即使在工作空间中有同名的变量，此变量在函数内部被赋值时就变成了局部变量，原来的全局变量不能被修改。局部变量在函数运行结束后就会消失。如

```
if('x' %in% ls()) rm(x)  
f <- function(){  
  x <- 123  
  cat(' 函数内: x = ', x, '\n')
```

```

}
f()
cat(' 函数运行完毕后: x=', x, '\n')
## 函数内: x = 123
> cat(' 函数运行完毕后: x=', x, '\n')
## Error in cat(" 函数运行完毕后: x=", x, "\n") : object 'x' not found

```

再比如，下面的函数试图知道自己被调用了多少次，但是因为每次函数调用完毕局部变量就消失，这样的程序不能达到目的：

```

f <- function(){
  if(!exists("runTimes")){
    runTimes <- 1
  } else {
    runTimes <- runTimes + 1
  }
  print(runTimes)
}
f()

```

```
## [1] 1
```

```
f()
```

```
## [1] 1
```

这个问题可以用 R 的 closure 来解决。

17.2.3 在函数内访问全局变量

函数内部可以读取全局变量的值，但一般不能修改全局变量的值。在现代编程指导思想中，全局变量容易造成不易察觉的错误，应谨慎使用，当然，也不是禁止使用，有些应用中不使用全局变量会使得程序更复杂且低效。

在下面的例子中，在命令行定义了全局变量 `x.g`，在函数 `f()` 读取了全局变量的值，但是在函数内给这样的变量赋值，结果得到的变量就变成了局部变量，全局变量本身不被修改：

```
x.g <- 9999
f <- function(x){
  cat(' 函数内读取：全局变量 x.g = ', x.g, '\n')
  x.g <- -1
  cat(' 函数内对与全局变量同名的变量赋值： x.g = ', x.g, '\n')
}
f()
```

```
## 函数内读取：全局变量 x.g = 9999
## 函数内对与全局变量同名的变量赋值： x.g = -1
cat(' 退出函数后原来的全局变量不变： x.g = ', x.g, '\n')
```

```
## 退出函数后原来的全局变量不变： x.g = 9999
```

在函数内部如果要修改全局变量的值，用 <<-代替 <-进行赋值。如

```
x.g <- 9999
f <- function(x){
  cat(' 函数内读取：全局变量 x.g = ', x.g, '\n')
  x.g <<- -1
  cat(' 函数内用"<<-" 对全局变量变量赋值： x.g = ', x.g, '\n')
}
f()
```

```
## 函数内读取：全局变量 x.g = 9999
## 函数内用"<<-"对全局变量变量赋值： x.g = -1
cat(' 退出函数后原来的全局变量被修改了： x.g = ', x.g, '\n')
```

```
## 退出函数后原来的全局变量被修改了： x.g = -1
```

后面将进一步解释函数在嵌套定义时 <<-的不同含义。

17.3 函数进阶

17.3.1 嵌套定义与句法作用域 (lexical scoping)

R 语言允许在函数体内定义函数。比如，

```
x <- -1
f0 <- function(x){
  f1 <- function(){
    x + 100
  }
  f1()
}
```

其中内嵌的函数 `f1()` 称为一个 closure(闭包)。

内嵌的函数体内在读取某个变量值时，如果此变量在函数体内还没有被赋值，它就不是局部的，会向定义的外面一层查找，外层一层找不到，就继续向外查找。上面例子 `f1()` 定义中的变量 `x` 不是局部变量，就向外一层查找，找到的会是 `f0` 的自变量 `x`，而不是全局空间中 `x`。如

```
f0(1)
```

```
## [1] 101
```

最后 `x+100` 中 `x` 取的是 `f0` 的实参值 `x=1`，而不是全局变量 `x=-1`。

这样的变量查找规则叫做句法作用域 (lexical scoping)，即函数运行时查找变量时，从其定义时的环境向外层逐层查找，而不是在运行时的环境中查找。句法作用域指的是可能有多个同名变量时查找变量按照定义时的环境查找，不是指查找变量值的规则。

例如，

```
f0 <- function(){
  f1 <- function(){
    x <- -1
    f2 <- function(){
      x + 100
    }
  }
}
```

```
    }  
    f2()  
  }  
  x <- 1000  
  f1()  
}  
f0()
```

```
## [1] 99
```

其中 `f2()` 运行时，用到的 `x` 是 `f1()` 函数体内的局部变量 `x=-1`，而不是被调用时 `f0()` 函数体内的局部变量 `x=1000`，所以结果是 $-1 + 100 = 99$ 。

“句法作用域”指的是函数调用时查找变量是查找其定义时的变量对应的存储空间，而不是定义时变量所取的历史值。函数运行时在找到某个变量对应的存储空间后，会使用该变量的当前值，而不是函数定义的时候该变量的历史值。例如

```
f0 <- function(){  
  x <- -1  
  f1 <- function(){  
    x + 100  
  }  
  x <- 1000  
  f1()  
}  
f0()
```

```
## [1] 1100
```

结果为什么不是 $-1 + 100 = 99$ 而是 $1000 + 100 = 1100$? 这是因为，`f1()` 在调用时，使用的 `x` 是 `f0` 函数体内局部变量 `x` 的值，但是要注意的是程序运行时会访问该变量的当前值，即 `1000`，而不是函数定义的时候 `x` 的历史值 `-1`。这个规则叫做“动态查找”(dynamic lookup)，句法作用域与动态查找一个说的是如何查找某个变量对应的存储空间，一个说的是使用该存储空间何时的存储值，程序运行时两个规则需要联合使用。

句法作用域不仅适用于查找变量，也适用于函数体内调用别的函数时查找函数。

查找函数的规则与查找变量规则相同。

17.3.1.1 辅助嵌套函数

有时内嵌函数仅仅是函数内用来实现模块化的一种工具，和正常的函数作用相同，没有任何特殊作用。例如，如下的程序在自变量 x 中输入一元二次方程 $ax^2 + bx + c = 0$ 的三个系数，输出解：

```
solve.sqe <- function(x){
  fd <- function(a, b, c) b^2 - 4*a*c
  d <- fd(x[1], x[2], x[3])
  if(d >= 0){
    return( (-x[2] + c(1,-1)*sqrt(d))/(2*x[1]) )
  } else {
    return( complex(real=-x[2], imag=c(1,-1)*sqrt(-d))/(2*x[1]) )
  }
}
```

在这个函数中内嵌的函数 `fd` 仅起到一个计算二次判别式公式的作用，没有用到任何的闭包特性，其中的形参变量 `a`, `b`, `c` 都是局部变量。运行如

```
solve.sqe(c(1, -2, 1))
```

```
## [1] 1 1
```

```
solve.sqe(c(1, -2, 0))
```

```
## [1] 2 0
```

```
solve.sqe(c(1, -2, 2))
```

```
## [1] 1+1i 1-1i
```

17.3.1.2 泛函

许多函数需要用函数作为参数，称这样的函数为泛函。比如，`apply` 类函数。这样的函数具有很好的通用性，因为需要进行的操作可以输入一个函数来规定，输入的函数规定什么样的操作，

用户可以自定义这样的函数。比如，希望对一个数据框中所有的数值型变量计算某种统计量，用来计算统计量的函数作为参数输入：

```
summary.df.numeric <- function(df, FUN, ...){
  vn <- names(df)
  vn <- vn[vapply(df, is.numeric, TRUE)]
  if(length(vn) > 0){
    sapply(df[,vn, drop=FALSE], FUN, ...)
  } else {
    numeric(0)
  }
}
```

这里参数 FUN 是用来计算统计量的函数。例如对 d.class 中每个数值型变量计算最小值：

```
d.class <- readr::read_csv("class.csv")
```

```
## Parsed with column specification:
## cols(
##   name = col_character(),
##   sex = col_character(),
##   age = col_double(),
##   height = col_double(),
##   weight = col_double()
## )
```

```
summary.df.numeric(d.class, min, na.rm=TRUE)
```

```
##   age height weight
##  11.0   51.3   50.5
```

17.3.1.3 函数工厂

利用嵌套定义在函数内的函数，可以解决上面的记录函数已运行次数的问题。如

```
f.gen <- function(){
  runTimes <- 0

  function(){
    runTimes <<- runTimes + 1
    print(runTimes)
  }
}
f <- f.gen()
f()
```

```
## [1] 1
```

```
f()
```

```
## [1] 2
```

在此例中, `f.gen` 中有局部变量 `runTimes`, `f.gen()` 的输出是一个函数, 输出结果保存到变量名 `f` 中, 所以 `f` 是一个函数, 调用 `f` 时, 查找变量 `runTimes` 时, 如果 `f` 的局部变量中没有 `runTimes`, 就从其定义的环境中逐层向外查找, 在 `f` 定义中用了 `<<-` 赋值, 这样赋值的含义是逐层向外查找变量是否存在, 在哪里找到变量就给那里的该变量赋值。 `f` 调用时向外查找到的变量在 `f.gen` 的局部空间中, 这是 `f` 函数的定义环境, 函数的定义环境是随函数本身一同保存的, 所以起到了把变量值 `runTimes` 与函数共同使用的效果。定义在函数内的函数称为一个 closure(闭包)。closure 最重要的作用就是定义能够保存历史运行状态的函数。

上面的 `f.gen` 这样的函数称为一个函数工厂, 因为它的结果是一个函数, 而且是一个闭包。闭包在 R 中的主要作用是带有历史状态的函数。

下面的例子也用了 closure, 可以显示从上次调用到下次调用之间经过的时间:

```
make_stop_watch <- function(){
  saved.time <- proc.time()[3]

  function(){
    t1 <- proc.time()[3]
```

```
    td <- t1 - saved.time
    saved.time <<- t1
    cat(" 流逝时间 (秒): ", td, "\n")
    invisible(td)
  }
}
ticker <- make_stop_watch()
ticker()
## 流逝时间 (秒):  0
for(i in 1:1000) sort(runif(10000))
ticker()
## 流逝时间 (秒):  1.53
```

其中 `proc.time()` 返回当前的 R 会话已运行的时间，结果在 MS Windows 系统中三个值，分别是用户时间、系统时间、流逝时间，其中流逝时间比较客观。

17.3.2 懒惰求值

R 函数在调用执行时，除非用到某个形式变量的值才求出其对应实参的值。这一点在实参是常数时无所谓，但是如果实参是表达式就不一样了。形参缺省值也是只有在函数运行时用到该形参的值时才求值。

例如，

```
f <- function(x, y=ifelse(x>0, TRUE, FALSE)){
  x <- -111
  if(y) x*2 else x*10
}
f(5)
```

```
## [1] -1110
```

可以看出，虽然形参 `x` 输入的实参值为 5，但是这时形参 `y` 并没按 `x=5` 被赋值为 `TRUE`，而是到函数体中第二个语句才被求值，这时 `x` 的值已经变成了 -111，故 `y` 的值是 `FALSE`。

17.3.3 递归调用

在函数内调用自己叫做递归调用。递归调用可以使得许多程序变得简单，但是往往导致程序效率很低，需谨慎使用。

R 中在递归调用时，最好用 `Recall` 代表调用自身，这样保证函数即使被改名（在 R 中函数是一个对象，改名后仍然有效）递归调用仍指向原来定义。

斐波那契数列是如下递推定义的数列：

$$\begin{aligned}x_0 &= 0, & x_1 &= 1 \\x_n &= x_{n-2} + x_{n-1}\end{aligned}$$

这个数列可以用如下递归程序自然地实现：

```
fib1 <- function(n){
  if(n == 0) return(0)
  else if(n == 1) return(1)
  else if(n >= 2) {
    Recall(n-1) + Recall(n-2)
  }
}
for(i in 0:10) cat('i =', i, ' x[i] =', fib1(i), '\n')
```

```
## i = 0 x[i] = 0
## i = 1 x[i] = 1
## i = 2 x[i] = 1
## i = 3 x[i] = 2
## i = 4 x[i] = 3
## i = 5 x[i] = 5
## i = 6 x[i] = 8
## i = 7 x[i] = 13
## i = 8 x[i] = 21
## i = 9 x[i] = 34
## i = 10 x[i] = 55
```

17.3.4 向量化

自定义的函数，如果其中的计算都是向量化的，那么函数自动地可以接受向量作为输入，结果输出向量。比如，将每个元素都变成原来的平方的函数：

```
f <- function(x){
  x^2
}
```

如果输入一个向量，结果也是向量，输出的每个元素是输入的对应元素的相应的平方值。

但是，如下的分段函数：

$$g(x) = \begin{cases} x^2, & |x| \leq 1, \\ 1, & |x| > 1 \end{cases}$$

其一元函数版本可以写成

```
g <- function(x){
  if(abs(x) <= 1) {
    y <- x^2
  } else {
    y <- 1
  }

  y
}
```

但是这个函数不能处理向量输入，因为 `if` 语句的条件必须是标量条件。一个容易想到的修改是

```
gv <- function(x){
  y <- numeric(length(x))
  sele <- abs(x) <= 1
  y[sele] <- x[sele]^2
  y[!sele] <- 1.0
}
```

```
y  
}
```

或者

```
gv <- function(x){  
  ifelse(abs(x) <= 1, x^2, 1)  
}
```

对于没有这样简单做法的问题，可以将原来的逻辑包在循环中，如

```
gv <- function(x){  
  y <- numeric(length(x))  
  for(i in seq(along=x)){  
    if(abs(x[i]) <= 1) {  
      y[i] <- x[i]^2  
    } else {  
      y[i] <- 1  
    }  
  }  
  y  
}
```

函数 `Vectorize` 可以将这样的操作自动化。如

```
g <- function(x){  
  if(abs(x) <= 1) {  
    y <- x^2  
  } else {  
    y <- 1  
  }  
  y  
}  
gv <- Vectorize(g)
```

```
gv(c(-2, -0.5, 0, 0.5, 1, 1.5))  
  
## [1] 1.00 0.25 0.00 0.25 1.00 1.00
```

17.3.5 纯函数与副作用

理想的自定义函数最好是像一般的数学函数那样，只要输入相同，输出也不变，而且除了利用输出值之外不能对程序环境做其它改变。这样的函数称为“纯函数”。R 的函数不能修改实参的值，这有助于实现纯函数的要求。

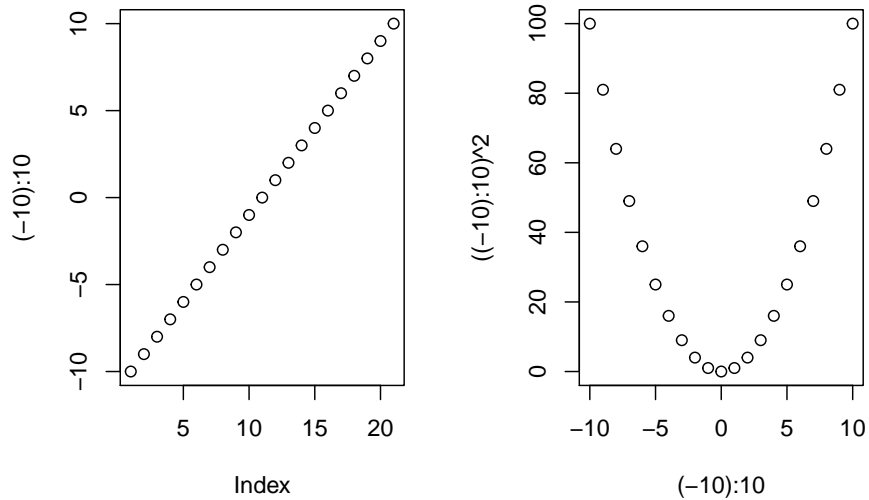
如果函数对相同的输入可以有不同的输出当然不是纯函数，例如 R 中的随机数函数 (`sample()`, `runif()`, `rnorm` 等)。

如果函数除了输出之外还在其它方面影响了运行环境，这样的函数就不是纯函数。所有画图函数 (`plot` 等)、输出函数 (`cat`, `print`, `save` 等) 都是这样的函数。这些对运行环境的改变叫做“副作用” (side effects)。又比如，`library()` 函数会引入新的函数和变量，`setwd()`, `Sys.setenv()`, `Sys.setlocale()` 会改变 R 运行环境，`options()`, `par()` 会改变 R 全局设置。自定义 R 函数中如果调用了非纯函数也就变成了非纯函数。编程中要尽量控制副作用而且要注意到副作用的影响，尤其是全局设置与全局变量的影响。

有些函数不可避免地要修改运行环境，如果可能的话，在函数结束运行前，应该恢复对运行环境的修改。为此，可以在函数体的前面部分调用 `on.exit()` 函数，此函数的参数是在函数退出前要执行的表达式或复合表达式。

例如，绘图的函数中经常需要用 `par()` 修改绘图参数，这会使得后续程序出错。为此，可以在函数开头保存原始的绘图参数，函数结束时恢复到原始的绘图参数。如

```
f <- function(){  
  opar <- par(mfrow=c(1,2))  
  on.exit(par(opar))  
  plot((-10):10)  
  plot((-10):10, ((-10):10)^2)  
}  
f()
```



如果函数中需要多次调用 `on.exit()` 指定多个恢复动作，除第一个调用的 `on.exit()` 以外都应该加上 `add=TRUE` 选项。

17.4 程序调试

17.4.1 跟踪调试

函数定义一般都包含多行，所以一般不在命令行定义函数，而是把函数定义写在源程序文件中，用 `source` 命令调入。用 `source` 命令调入运行的程序与在命令行运行的效果基本相同，这样定义的变量也是全局变量。

考虑如下函数定义：

```
f <- function(x){  
  for(i in 1:n){  
    s <- s + x[i]  
  }  
}
```


运行发现有错误:

```
f(1:5)
## Error in f(1:5) : object 'n' not found
```

简单的函数可以直接仔细检查发现错误, 用 `cat`, `print` 等输出中间结果查找错误。R 提供了一个 `browser()` 函数, 在程序中插入对 `browser()` 函数的调用, 可以进入跟踪调试状态, 可以实时地查看甚至修改运行时变量的值。

程序运行遇到 `browser()` 函数时程序进入 Browser 的调试命令行。在调试命令行, 用 `n` 命令逐句运行, 用 `s` 命令跟踪进调用的函数内部逐句运行, 用 `c` 命令恢复正常运行, 用 `q` 命令强制终止程序运行。可以如同在 R 命令行一样查看变量的值或修改变量的值。在 RStudio 中进入跟踪状态后有相应的运行控制图标, 可以用鼠标点击某行程序的行号设置断点, 重新 `source()` 之后就可以在断点处进入跟踪状态。

为调试如上函数 `f` 的程序, 在定义中插入对 `browser()` 的调用如:

```
f <- function(x){
  browser()
  for(i in 1:n){
    s <- s + x[i]
  }
}
```

调试运行过程如下:

```
f(1:5)
## Called from: f(1:5)
## Browse[1]> n
## debug at #3: for (i in 1:n) {
##     s <- s + x[i]
## }
## Browse[2]> n
## Error in f(1:5) : object 'n' not found
```

发现是在 `for(i in 1:n)` 行遇到未定义的变量 `n`。

在源文件中把出错行改为 `for(i in 1:length(x))`，再次运行，发现在运行 `s <- s + x[i]` 行时，遇到未定义的变量 `s`。这是忘记初始化引起的。在 `for` 语句前添加 `s <- 0` 语句，函数定义变成：

```
f <- function(x){
  browser()
  s <- 0
  for(i in 1:length(x)){
    s <- s + x[i]
  }
}
```

再次运行，在 `Browse[1]>` 提示下命令 `c` 表示恢复正常运行，程序不显示错误但是也没有显示求和结果。检查可以看出错误是忘记把函数返回值写在函数定义最后。

在函数定义最后添加 `s` 一行，再次运行，程序结果与手工验算结果一致。函数变成

```
f <- function(x){
  browser()
  n <- length(x)
  s <- 0
  for(i in 1:n){
    s <- s + x[i]
  }
  s
}
```

自定义函数应该用各种不同输入测试其正确性和稳定性。比如，上面的函数当自变量 `x` 为零长度向量时应该返回 `0` 才合适，但是上面的写法会返回一个 `numeric(0)` 结果，这个结果表示长度为零的向量：

```
f(numeric(0))
## Called from: f(numeric(0))
## Browse[1]> c
## numeric(0)
```

程序输入了零长度自变量, 我们期望其输出为零而不是 `numeric(0)`。在自变量 `x` 为零长度时, 函数中 `for(i in 1:length(x))` 应该一次都不进入循环, 跟踪运行可以发现实际对 `i=1` 和 `i=0` 共运行了两轮循环。把这里的 `1:length(x)` 改成 `seq(along=x)` 解决了问题, `seq(along=x)` 生成 `x` 的下标序列, 如果 `x` 是零长度的则下标序列为零长度向量。

函数不需要修改后, 可以把对 `browser()` 的调用删除或注释掉。函数最终修改为:

```
f <- function(x){
  s <- 0
  for(i in seq(along=x)){
    s <- s + x[i]
  }
  s
}
```

这里只是用这个简单函数演示如何调试程序, 求向量和本身是不需要我们去定义新函数的, `sum` 函数本来就是完成这样的功能。实际上, 许多我们认为需要自己编写程序作的事情, 在 R 网站都能找到别人已经完成的程序。

17.4.2 出错调试选项

比较长的程序在调试时如果从开头就跟踪, 比较耗时。可以设置成出错后自动进入跟踪模式, 检查出错时的变量值。只要进行如下设置:

```
options(error=recover)
```

则在出错后可以选择进入出错的某一层函数内部, 在 `browser` 环境中跟踪运行。

例如, 如上设置后, 前面那个求向量元素和的例子程序按最初的定义, 运行时出现如下的选择:

```
## Error in f(1:5) : object 'n' not found
##
## Enter a frame number, or 0 to exit
##
```

```
## 1: f(1:5)
##
## Selection: f(1:5)
##
## Selection: 1
## Called from: top level
## Browse[1]>
```

在 `Selection` 后面输入了 1，就进入了函数内部跟踪。用 `Q` 终止运行并退出整个 `browser` 跟踪。当函数调用函数时可以选择进入哪一个函数进行跟踪。

17.4.3 警告处理

有些警告信息实际是错误，用 `options()` 的 `warn` 参数可以设置警告级别，如设置 `warn=2` 则所有警告当作错误处理。设置如

```
options(warn=2)
```

17.4.4 `stop()`、`warning()`、`message()`

编写程序时应尽可能提前发现不合法的输入和错误的状态。发现错误时，可以用 `stop(s)` 使程序运行出错停止，其中 `s` 是一个字符型对象，用来作为显示的出错信息。

发现某些问题后如果不严重，可以不停止程序运行，但用 `warning(s)` 提交一个警告信息，其中 `s` 是字符型的警告信息。警告信息的显示可能比实际运行要滞后一些。

函数 `message()` 与 `stop()`、`warning()` 类似，不算是错误或者警告，但仍算是某种非正常的信息输出。

17.4.5 预防性设计

在编写自定义函数时，可以检查自变量输入以确保输入符合要求。函数 `stopifnot` 可以指定自变量的若干个条件，当自变量不符合条件时自动出错停

止。

例如，函数 `f()` 需要输入两个数值型向量 `x`, `y`, 需要长度相等，可以用如下的程序

```
f <- function(x, y){
  stopifnot(is.numeric(x),
            is.numeric(y),
            length(x)==length(y))
  ## 函数体程序语句...
}
```

17.5 函数式编程介绍

R 可以算是一个函数式语言 (functional language):

1. R 语言的设计主要用函数求值来进行运算;
2. R 的用户主要使用函数调用来访问 R 的功能。

按照函数式编程的要求，每个 R 函数必须功能清晰、定义确切。比较容易控制的函数是纯函数，纯函数必须像数学中单值函数那样给定自变量输入有唯一确定的输出。比如，多个函数用全局变量传递信息，就不能算是纯函数。

R 支持类 (class) 和方法 (method)，实际提供了适用于多种自变量的通用函数 (generic function)，不同自变量类型调用该类特有的方法，但函数名可以保持不变。

函数式编程语言提供了定义纯函数的功能。这样的函数不能有副作用 (side effects): 函数返回值包含了函数执行的所有效果。函数定义仅由对所有可能的自变量值确定返回值来确定，不依赖于任何外部信息 (也就不能依赖于全局变量与系统设置值)。函数定义返回值的方式是隐含地遍历所有可能的参数值给出返回值，而不是用过程式的计算来修改对象的值。

函数式编程的目的是提供可理解、可证明正确的软件。R 虽然带有函数式编程语言特点，但并不强求使用函数式编程规范。典型的函数式编程语言如 Haskell, Lisp 的运行与 R 的显式的、顺序的执行方式相差很大。

17.5.1 函数式编程的要求

- 没有副作用。调用一个函数对后续运算没有影响，不管是再次调用此函数还是调用其它函数。这样，用全局变量在函数之间传递信息就是不允许的。其它副作用包括写文件、打印、绘图等，这样的副作用对函数式要求破坏不大。
- 不受外部影响。函数返回值只依赖于其自变量及函数的定义。
- 不受赋值影响。函数定义不需要反复对内部对象（所谓“状态变量”）赋值或修改。

R 只能部分满足这些要求。一个 R 函数是否满足这些要求不仅要看函数本身，还要看函数内部调用的其它函数。

像 `options()` 函数这样修改全局运行环境的功能会破坏函数式要求。尽可能让自己的函数不依赖于 `options()` 中的参数。

与具体硬件、软件环境有关的一些因素也破坏纯函数要求，如不同的硬件常数、精度等。调用操作系统的功能对函数式要求破坏较大。减少赋值主要需要减少循环，可以用 R 的向量化方法解决。

17.5.2 Map、Reduce、Filter

R 提供了 `Map`, `Reduce`, `Filter`, `Find`, `Negate`, `Position` 等支持函数式编程的工具函数。这些函数包含对列表每一项进行变换，列表数据汇总，列表元素筛选等功能。

17.5.2.1 Map 函数

`Map()` 以一个函数作为参数，可以对其它参数的每一对应元素进行变换，结果为列表。

例如，对数据框 `d`，如下的程序可以计算每列的平方和：

```
d <- data.frame(  
  x = c(1, 7, 2),
```

```
y = c(3, 5, 9)
Map(function(x) sum(x^2), d)
```

```
## $x
## [1] 54
##
## $y
## [1] 115
```

实际上，这个例子也可以用 `lapply()` 改写成

```
lapply(d, function(x) sum(x^2))
```

```
## $x
## [1] 54
##
## $y
## [1] 115
```

`Map()` 比 `lapply()` 增强的地方在于它允许对多个列表的对应元素逐一处理。例如，为了求出 `d` 中每一行的最大值，可以用

```
Map(max, d$x, d$y)
```

```
## [[1]]
## [1] 3
##
## [[2]]
## [1] 7
##
## [[3]]
## [1] 9
```

可以用 `unlist()` 函数将列表结果转换为向量，如

```
unlist(Map(max, d$x, d$y))
```

```
## [1] 3 7 9
```

17.5.2.2 Reduce 函数

Reduce 函数把输入列表（或向量）的元素逐次地用给定的函数进行合并计算。例如，

```
Reduce(sum, 1:4)
```

```
## [1] 10
```

实际执行的是 $((1 + 2) + 3) + 4$ 。当然，求 1:4 的和只需要 `sum(1:4)`，但是 Reduce 可以对元素为复杂类型的列表进行逐项合并计算。

例如，`intersect` 函数可以计算两个集合的交集；对多个集合，如何计算交集？下面的例子产生了 4 个集合，然后反复调用 `intersect()` 求出了交集：

```
set.seed(2)
```

```
x <- replicate(4, sample(1:5, 5, replace=TRUE), simplify=FALSE); x
```

```
## [[1]]
```

```
## [1] 5 1 5 1 4
```

```
##
```

```
## [[2]]
```

```
## [1] 5 1 2 3 1
```

```
##
```

```
## [[3]]
```

```
## [1] 3 2 3 1 1
```

```
##
```

```
## [[4]]
```

```
## [1] 4 3 1 5 3
```

```
intersect(intersect(intersect(x[[1]], x[[2]]), x[[3]]), x[[4]])
```

```
## [1] 1
```

也可以用 `magrittr` 包的 `%>%` 符号写成：

```
library(magrittr)
```

```
x[[1]] %>% intersect(x[[2]]) %>% intersect(x[[3]]) %>% intersect(x[[4]])
```

```
## [1] 1
```


还可以写成循环:

```
y <- x[[1]]
for(i in 2:4) y <- intersect(y, x[[i]])
y
```

```
## [1] 1
```

都比较繁琐。

利用 Reduce 函数, 只要写成

```
Reduce(intersect, x)
```

```
## [1] 1
```

Reduce 函数还可以用 right 参数选择是否从右向左合并, 用参数 init 给出合并初值, 用参数 accumulate 要求保留每一步合并的结果 (累计)。这个函数可以把很多仅适用于两个运算元的运算推广到多个参数的情形。

17.5.2.3 Filter、Find、Position 函数

Filter(f, x) 用一个判断真假的一元函数 f 作为筛选规则, 从列表或向量 x 中筛选出用 f 作用后为真值的元素子集。f 必须返回标量的 TRUE 或者 FALSE, 这样的函数称为示性函数 (predicate functions)。例如

```
f <- function(x) x > 0 & x < 1
Filter(f, c(-0.5, 0.5, 0.8, 1))
```

```
## [1] 0.5 0.8
```

当然, 这样的简单例子完全可以改写成:

```
f <- function(x) x > 0 & x < 1
x <- c(-0.5, 0.5, 0.8, 1)
x[x>0 & x < 1]
```

```
## [1] 0.5 0.8
```

但是, 对于比较复杂的判断, 特别是需要用许多个语句计算后进行的判断, 就需要把判断写成一个函数, 然后可以用 Filter 比较简单地表达按照判断规则

取子集的操作。

`Find()` 作用与 `Filter()` 类似，但是仅返回满足条件的第一个，也可以用参数 `right=TRUE` 要求返回满足条件的最后一个。

`Position()` 作用与 `Find()` 类似，但不是返回满足条件的元素而是返回第一个满足条件的元素所在的下标位置。

Chapter 18

R 程序效率

18.1 R 的运行效率

R 是解释型语言，在执行单个运算时，效率与编译代码相近；在执行迭代循环时，效率较低，与编译代码的速度可能相差几十倍。R 以向量、矩阵为基础运算单元，在进行向量、矩阵运算时效率很高，应尽量采用向量化编程。

另外，R 语言的设计为了方便进行数据分析和统计建模，有意地使语言特别灵活，比如，变量为动态类型而且内容可修改，变量查找在当前作用域查找不到可以向上层以及扩展包中查找，函数调用时自变量仅在使用其值时才求值（懒惰求值），这样的设计都为运行带来了额外的负担，使得运行变慢。

在计算总和、元素乘积或者每个向量元素的函数变换时，应使用相应的函数，如 `sum`, `prod`, `sqrt`, `log` 等。

对于从其它编程语言转移到 R 语言的学生，如果不细究 R 特有的编程模式，编制的程序可能效率比正常 R 程序慢上几十倍，而且繁琐冗长。

为了提高 R 程序的运行效率，需要尽可能利用 R 的向量化特点，尽可能使用已有的高效函数，还可以把运行速度瓶颈部分改用 C++、FORTRAN 等编译语言实现，可以用 R 的 `profler` 工具查找运行瓶颈。对于大量数据的长时间计算，可以借助于现代的并行计算工具。

对已有的程序，仅在运行速度不满意时才需要进行改进，否则没必要花费宝贵的

时间用来节省几秒钟的计算机运行时间。要改善运行速度，首先要找到运行的瓶颈，这可以用专门的性能分析（profiling）功能实现。R 软件中的 `Rprof()` 函数可以执行性能分析的数据收集工作，收集到的性能数据用 `summaryRprof()` 函数可以显示运行最慢的函数。如果使用 RStudio 软件，可以用 `Profile` 菜单执行性能数据收集与分析，可以在图形界面中显示程序中哪些部分运行花费时间最多。

在改进已有程序的效率时，第一要注意的就是不要把原来的正确算法改成一个速度更快但是结果错误的算法。这个问题可以通过建立试验套装，用原算法与新算法同时试验看结果是否一致来避免。多种解决方案的正确性都可以这样保证，也可以比较多种解决方案的效率。

本章后面部分描述常用的改善性能的方法。对于涉及到大量迭代的算法，如果用 R 实现性能太差不能满足要求，可以改成 C++ 编码，用 `Rcpp` 扩展包连接到 R 中。`Rcpp` 扩展包的使用将单独讲授。

R 的运行效率也受到内存的影响，占用内存过多的算法有可能受到物理内存大小限制无法运行，过多复制也会影响效率。

如果要实现一个比较单纯的不需要利用 R 已有功能的算法，发现用 R 计算速度很慢的时候，也可以考虑先用 Julia 语言实现。Julia 语言设计比 R 更先进，运算速度快得多，运算速度经常能与编译代码相比，缺点是刚刚诞生几年的时间，可用的软件包还比较少。

18.2 向量化编程

18.2.1 示例 1

假设要计算如下的统计量：

$$w = \frac{1}{n} \sum_{i=1}^n |x_i - \hat{m}|,$$

其中 x_1, x_2, \dots, x_n 是某总体的样本， \hat{m} 是样本中位数。用传统的编程风格，把这个统计量的计算变成一个 R 函数，可能会写成：

```
f1 <- function(x){  
  n <- length(x)  
  mhat <- median(x)  
  s <- 0.0  
  for(i in 1:n){  
    s <- s + abs(x[i] - mhat)  
  }  
  s <- s/n  
  return(s)  
}
```

用 R 的向量化编程，函数体只需要一个表达式：

```
f2 <- function(x) mean( abs(x - median(x)) )
```

其中 `x - median(x)` 利用了向量与标量运算结果是向量每个元素与标量运算的规则，`abs(x - median(x))` 利用了 `abs()` 这样的一元函数如果以向量为输入就输出每个元素的函数值组成的向量的规则，`mean(...)` 避免了求和再除以 `n` 的循环也不需要定义多余的变量 `n`。

显然，第二种做法的程序比第一种做法简洁的多，如果多次重复调用，第二种做法的计算速度比第一种要快几十倍甚至上百倍。在 R 中，用 `system.time()` 函数可以求某个表达式的计算时间，返回结果的第 3 项是流逝时间。下面对 `x` 采用 10000 个随机数，并重复计算 1000 次，比较两个程序的效率：

```
nrep <- 1000  
x <- runif(10000)  
y1 <- numeric(nrep); y2 <- y1  
system.time(for(i in 1:nrep) y1[i] <- f1(x) )[3]  
## elapsed  
## 10.08  
system.time(for(i in 1:nrep) y1[i] <- f2(x) )[3]  
## elapsed  
## 0.48
```

速度相差二十倍以上。

有一个 R 扩展包 `microbenchmark` 可以用来测量比较两个表达式的运行时间。如:

```
x <- runif(10000)
microbenchmark::microbenchmark(
  f1(x),
  f2(x)
)

## Unit: microseconds
##      expr      min       lq      mean   median      uq      max neval
## f1(x) 1990.426 2129.2210 2738.4549 2348.723 2967.390 7911.838   100
## f2(x)  417.929  441.3185  557.4435  477.559  582.941 3796.306   100
```

就平均运行时间（单位：毫秒）来看，`f2()` 比 `f1()` 快大约 30 倍。

18.2.2 示例 2

假设要编写函数计算

$$f(x) = \begin{cases} 1 & x \geq 0, \\ 0 & \text{其它} \end{cases}$$

利用传统思维，程序写成

```
f1 <- function(x){
  n <- length(x)
  y <- numeric(n)

  for(i in seq(along=x)){
    if(x[i] >= 0) y[i] <- 1
    else y[i] <- 0
  }

  y
}
```

实际上, `y <- numeric(n)` 使得 `y` 的每个元素都初始化为 0, 所以程序中 `else y[i] <- 0` 可以去掉。

利用向量化与逻辑下标, 程序可以写成:

```
f2 <- function(x){
  n <- length(x)
  y <- numeric(n)
  y[x >= 0] <- 1

  y
}
```

但是, 利用 R 中内建函数 `ifelse()`, 可以把函数体压缩到仅用一个语句:

```
f2 <- function(x) ifelse(x >= 0, 1, 0)
```

18.2.3 示例 3

考虑一个班的学生存在生日相同的概率。假设一共有 365 个生日 (只考虑月、日)。设一个班有 n 个人, 当 n 大于 365 时 {至少两个人有生日相同} 是必然事件 (概率等于 1)。

当 n 小于等于 365 时:

$$\begin{aligned} & P(\text{至少有两人同生日}) \\ &= 1 - P(n \text{ 个人生日彼此不同}) \\ &= 1 - \frac{365 \times 364 \times \cdots \times (365 - (n - 1))}{365^n} \\ &= 1 - \frac{365 - 0}{365} \cdot \frac{365 - 1}{365} \cdots \frac{365 - (n - 1)}{365} \end{aligned}$$

对 $n = 1, 2, \dots, 365$ 来计算对应的概率。完全用循环 (两重循环), 程序写成:

```
f1 <- function(){
  ny <- 365
  x <- numeric(ny)
  for(n in 1:ny){
```

```
s <- 1
for(j in 0:(n-1)){
  s <- s * (365-j)/365
}
x[n] <- 1 - s
}
x
}
```

注意，不能先计算 $365 \times 364 \times \dots \times (365 - (n - 1))$ 和 365^n 再相除，这会造成数值溢出。

用 `prod()` 函数可以向量化内层循环：

```
f2 <- function(){
  ny <- 365
  x <- numeric(ny)
  for(n in 1:ny){
    x[n] <- 1 - prod((365:(365-n+1))/365)
  }
  x
}
```

程序利用了向量与标量的除法，以及内建函数 `prod()`。

把程序用 `cumprod()` 函数改写，可以完全避免循环：

```
f3 <- function(){
  ny <- 365
  x <- 1 - cumprod((ny:1)/ny)
  x
}
```

用 `microbenchmark` 比较：

```
microbenchmark::microbenchmark(
  f1(),
  f2(),
```



```

f3()
)
## Unit: microseconds
##      expr      min       lq      mean   median      uq      max neval
## f1() 2534.807 2577.730 2679.48244 2615.256 2712.9270 3408.187   100
## f2()  323.855  333.365  414.81692  344.160  369.3485 5868.456   100
## f3()   1.028   1.542   2.52415   2.056   2.5700  25.189   100

```

f2() 比 f1() 快大约 7 倍, f3() 比 f2() 又快了大约 160 倍, f3() 比 f1() 快了一千倍以上!

18.3 减少显式循环

显式循环是 R 运行速度较慢的部分, 有循环的程序也比较冗长, 与 R 的向量化简洁风格不太匹配。另外, 在循环内修改数据子集, 例如数据框子集, 可能会先制作副本再修改, 这当然会损失很多效率。R 3.1.0 版本以后列表元素在修改时不制作副本。

前面已经指出, 利用 R 的向量化运算可以减少很多循环程序。

R 中的有些运算可以用内建函数完成, 如 `sum`, `prod`, `cumsum`, `cumprod`, `mean`, `var`, `sd` 等。这些函数以编译程序的速度运行, 不存在效率损失。

R 的 `sin`, `sqrt`, `log` 等函数都是向量化的, 可以直接对输入向量的每个元素进行变换。

对矩阵, 用 `apply` 函数汇总矩阵每行或每列。`colMeans`, `rowMeans` 可以计算矩阵列平均和行平均, `colSums`, `rowSums` 可以计算矩阵列和与行和。

`apply` 类函数有多个, 包括 `apply`, `sapply`, `lapply`, `tapply`, `vapply`, `replicate` 等。这些函数不一定能提高程序运行速度, 但是使用这些函数更符合 R 的程序设计风格, 使程序变得简洁, 当然, 程序更简洁并不等同于程序更容易理解, 要理解这样的程序, 需要更多学习与实践。

18.3.1 lapply()、sapply() 和 vapply() 函数

对列表，`lapply` 函数操作列表每个元素，格式为

```
lapply(X, FUN)
```

其中 `X` 是一个列表或向量，`FUN` 是一个函数（可以是有名或无名函数），结果也总是一个列表，结果列表的第 i 个元素是将 `X` 的第 i 个元素输入到 `FUN` 中的返回结果。如果输入不是列表，就转换为列表再对每一元素做变换。

`sapply` 与 `lapply` 函数类似，但是 `sapply` 试图简化输出结果为向量或矩阵，在不可行时才和 `lapply` 返回列表结果。如果 `X` 长度为零，结果是长度为零的列表；如果 `FUN(X[i])` 都是长度为 1 的结果，`sapply()` 结果是一个向量；如果 `FUN(X[i])` 都是长度相同且长度大于 1 的向量，`sapply()` 结果是一个矩阵，矩阵的第 i 列保存 `FUN(X[i])` 的结果。因为 `sapply()` 的结果类型的不确定性，在自定义函数中应慎用。

`vapply()` 函数与 `sapply()` 函数类似，但是它需要第三个参数即函数返回值类型的例子，格式为

```
vapply(X, FUN, FUN.VALUE)
```

其中 `FUN.VALUE` 是每个 `FUN(X[i])` 的返回值的例子，要求所有 `FUN(X[i])` 结果类型和长度相同。

18.3.1.1 示例 1：数据框变量类型

`typeof()` 函数求变量的存储类型，如

```
d.class <- read.csv('class.csv', header=TRUE)
typeof(d.class[, 'age'])
```

```
## [1] "integer"
```

这里 `d.class` 是一个数据框，数据框也是一个列表，每个列表元素是数据框的一列。如下程序使用 `sapply()` 求每一列的存储类型：

```
sapply(d.class, typeof)
```

```
##      name      sex      age      height      weight
```

```
## "integer" "integer" "integer" "double" "double"
```

注意因为 CSV 文件读入为数据框时把姓名、性别都转换成了因子，所以这两个变量的存储类型也是整数。为了避免这样的转换，在 `read.csv()` 中使用选项 `stringsAsFactors=FALSE` 选项。

关于一个数据框的结构，用 `str()` 函数可以得到更为详细的信息：

```
str(d.class)
```

```
## 'data.frame': 19 obs. of 5 variables:
## $ name : Factor w/ 19 levels "Alfred","Alice",...: 2 3 5 10 11 12 15 16 17 1 ...
## $ sex : Factor w/ 2 levels "F","M": 1 1 1 1 1 1 1 1 1 2 ...
## $ age : int 13 13 14 12 12 15 11 15 14 14 ...
## $ height: num 56.5 65.3 64.3 56.3 59.8 66.5 51.3 62.5 62.8 69 ...
## $ weight: num 84 98 90 77 84.5 ...
```

18.3.1.2 示例 2: `strsplit()` 函数结果处理

假设有 4 个学生的 3 次小测验成绩，每个学生的成绩记录到了一个以逗号分隔的字符串中，如：

```
s <- c('10, 8, 7',
      '5, 2, 2',
      '3, 7, 8',
      '8, 8, 9')
```

对单个学生，可以用 `strsplit()` 函数把三个成绩拆分，如：

```
strsplit(s[1], ',', fixed=TRUE)[[1]]
```

```
## [1] "10" " 8" " 7"
```

注意这里 `strsplit()` 的结果是仅有一个元素的列表，用了 “[[...]]” 格式取出列表元素。拆分的结果可以用 `as.numeric()` 转换为有三个元素的数值型向量：

```
as.numeric(strsplit(s[1], ',', fixed=TRUE)[[1]])
```

```
## [1] 10 8 7
```

还可以求三次小测验的总分：

```
sum(as.numeric(strsplit(s[1], ',', fixed=TRUE)[[1]]))
```

```
## [1] 25
```

用 `strsplit()` 处理有 4 个字符串的字符型向量 `s`, 结果是长度为 4 的列表：

```
tmp.res <- strsplit(s, ',', fixed=TRUE); tmp.res
```

```
## [[1]]
## [1] "10" " 8" " 7"
##
## [[2]]
## [1] "5"  " 2" " 2"
##
## [[3]]
## [1] "3"  " 7" " 8"
##
## [[4]]
## [1] "8"  " 8" " 9"
```

用 `sapply()` 和 `as.numeric()` 可以把列表中所有字符型转为数值型, 并以矩阵格式输出：

```
sapply(tmp.res, as.numeric)
```

```
##      [,1] [,2] [,3] [,4]
## [1,]  10   5   3   8
## [2,]   8   2   7   8
## [3,]   7   2   8   9
```

但是, 在通用程序中使用 `sapply()` 有可能会发生结果类型可变的情况。为此, 上面可以用 `vapply()` 改写成

```
vapply(tmp.res, as.numeric, numeric(3))
```

```
##      [,1] [,2] [,3] [,4]
## [1,]  10   5   3   8
```

```
## [2,] 8 2 7 8
## [3,] 7 2 8 9
```

其中第三个参数 `numeric(3)` 给出了对 `tmp.res` 中的任意一项应用 `as.numeric` 函数的结果的例子。

如果要求每个学生的小测验总分，只要对结果矩阵每列求和：

```
colSums(vapply(tmp.res, as.numeric, numeric(3)))
```

```
## [1] 25 9 18 25
```

使用 `apply` 类函数的程序写法简洁，但是对于初学者需要比较长的时间来读懂，需要更长的时间用到自己的程序中。

上例中的嵌套调用用 `magrittr` 包的管道运算符更容易理解：

```
library(magrittr)
s %>%
  strsplit(",", fixed=TRUE) %>%
  sapply(as.numeric) %>%
  colSums()
```

```
## [1] 25 9 18 25
```

18.3.1.3 示例 3: 不等长结果

调用 `sapply(列表, 函数)` 时，如果“函数”结果长度有变化，结果只能以列表输出。这时，`sapply` 与 `lapply` 返回相同的结果。一般地，`sapply` 试图把结果简化为向量、矩阵、多维数组，在无法简化时就返回列表；`lapply` 总是返回列表。

设数据框 `d` 中有两列数，希望将每列变成没有重复值的。数据例子如下：

```
d1 <- data.frame(x1=c(1,3,3,2), x2=c(3,5,5,3)); d1
```

```
##   x1 x2
## 1  1  3
## 2  3  5
## 3  3  5
```

```
## 4 2 3
```

因为 x1 和 x2 两列的无重复值个数不同，结果只能是列表：

```
sapply(d1, unique)
```

```
## $x1
## [1] 1 3 2
##
## $x2
## [1] 3 5
```

与 `lapply(d, unique)` 效果相同。

如果使用了 `vapply`，在遇到结果长度变化时会明确报错，如：

```
vapply(d1, unique, numeric(3))
## Error in vapply(d1, unique, numeric(3)) : values must be length 3,
## but FUN(X[[2]]) result is length 2
```

在以上的例子中，输入的 `d1` 数据框若无重复个数相同，`sapply` 结果就是矩阵，而 `lapply` 结果仍然是列表：

```
d2 <- data.frame(x1=c(1,3,3,2), x2=c(3,5,5,7)); d2
```

```
##   x1 x2
## 1  1  3
## 2  3  5
## 3  3  5
## 4  2  7
```

```
sapply(d2, unique)
```

```
##      x1 x2
## [1,]  1  3
## [2,]  3  5
## [3,]  2  7
```

```
lapply(d2, unique)
```

```
## $x1
```

```
## [1] 1 3 2
##
## $x2
## [1] 3 5 7
```

`sapply()` 的这种特点会造成编程判断困难，所以在不能确定函数结果长度是否保持不变时，应该用 `lapply()` 代替 `sapply()`，`lapply()` 总是返回列表。如果能够确定函数结果长度保持不变，在通用程序中应该用 `vapply()` 取代 `sapply()`，使得程序结果总是一致的。

18.3.1.4 示例 4: 无名函数

`sapply`、`vapply` 和 `lapply` 中要做的函数变换可以当场定义，不需要函数名。

仍使用示例 2 的数据，任务是从输入的逗号分隔成绩中求每个学生的三科总分。

用 `strsplit()` 拆分可得列表，每个列表是由三个成绩字符串的字符型向量。如下代码可以求得总分：

```
s <- c('10, 8, 7',
      '5, 2, 2',
      '3, 7, 8',
      '8, 8, 9')
sapply( strsplit(s, ',', fixed=TRUE),
        function(ss) sum(as.numeric(ss)) )
```

```
## [1] 25 9 18 25
```

这里没有预先定义处理函数，也没有函数名，而是直接对 `sapply` 的第二自变量使用了一个无名函数。实际上，R 的函数定义也是函数名被赋值为一个函数对象。

但是，如果 `sapply()` 函数中的无名函数访问其它变量的话，容易产生作用域问题。

18.3.2 replicate() 函数

`replicate()` 函数用来执行某段程序若干次，类似于 `for()` 循环但是没有计数变量。常用于随机模拟。`replicate()` 的缺省设置会把重复结果尽可能整齐排列成一个多维数组输出。

语法为

```
replicate(重复次数, 要重复的表达式)
```

其中的表达式可以是复合语句，也可以是执行一次模拟的函数。

下面举一个简单模拟例子。设总体 X 为 $N(0,1)$ ，取样本量 $n = 5$ ，重复地生成模拟样本共 $B = 6$ 组，输出每组样本的样本均值和样本标准差。模拟可以用如下的 `replicate()` 实现：

```
set.seed(1)
replicate(6, {
  x <- rnorm(5, 0, 1);
  c(mean(x), sd(x)) })
```

```
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
## [1,] 0.1292699 0.1351357 0.03812297 0.4595670 0.08123054 -0.3485770
## [2,] 0.9610394 0.6688342 1.49887443 0.4648177 1.20109623 0.7046822
```

结果是一个矩阵，矩阵行数与每次模拟的结果（均值、标准差）个数相同，这里第一行都是均值，第二行都是标准差；矩阵每一列对应于一次模拟。此结果转置可能更合适。

18.3.3 Map() 和 mapply()

`lapply()`、`sapply()`、`vapply()` 只能针对单个列表 X 的每个元素重复处理。如果有两个列表 X 和 Y 要进行对应元素的处理，用这三个函数不容易做到，这时可以用 `Map()` 或 `mapply()`。`Map()` 的格式为

```
Map(f, ...)
```

其中 f 是一个函数，`Map` 的其它参数都是每次取出对应的元素作为 $f()$ 的输入。`Map()` 的结果总是列表。

例如，下面有两个向量：

```
x <- c(1, 7, 2)
y <- c(3, 5, 9)
```

为了求得 `x` 和 `y` 的每个对应元素的最大值，可以用

```
Map(max, x, y)
```

```
## [[1]]
## [1] 3
##
## [[2]]
## [1] 7
##
## [[3]]
## [1] 9
```

结果是一个列表。为了把列表转换为普通向量，可以用 `unlist()` 函数，如

```
unlist(Map(max, x, y))
```

```
## [1] 3 7 9
```

这个例子演示了 `Map` 的用法。实际上，为了求多个向量对应元素的最大值，可以用 `pmax` 函数，如

```
pmax(x, y)
```

```
## [1] 3 7 9
```

`mapply()` 函数与 `Map()` 类似，但是可以自动简化结果类型，可以看成是 `sapply()` 推广到了可以对多个输入的对应元素逐项处理。`mapply()` 可以用参数 `MoreArgs` 指定逐项处理时一些共同的参数。简单的调用格式为

```
mapply(FUN, ...)
```

如

```
mapply(max, x, y)
```

```
## [1] 3 7 9
```

18.4 R 的计算函数

R 中提供了大量的数学函数、统计函数和特殊函数，可以打开 R 的 HTML 帮助页面，进入“Search Enging & Keywords”链接，查看其中与算术、数学、优化、线性代数等有关的专题。

这里简单列出一部分常用函数，对函数 `filter`, `fft`, `convolve` 进行说明。

18.4.1 数学函数

常用数学函数包括 `abs`, `sign`, `log`, `log10`, `sqrt`, `exp`, `sin`, `cos`, `tan`, `asin`, `acos`, `atan`, `atan2`, `sinh`, `cosh`, `tanh`。还有 `gamma`, `lgamma`(伽玛函数的自然对数)。

用于取整的函数有 `ceiling`, `floor`, `round`, `trunc`, `signif`, `as.integer` 等。这些函数是向量化的一元函数。

`choose(n,k)` 返回从 n 中取 k 的组合数。`factorial(x)` 返回 $x!$ 结果。`combn(x,m)` 返回从集合 x 中每次取出 m 个的所有不同取法，结果为一个矩阵，矩阵每列为一种取法的 m 个元素值。

18.4.2 概括函数

`sum` 对向量求和, `prod` 求乘积。

`cumsum` 和 `cumprod` 计算累计，得到和输入等长的向量结果。

`diff` 计算前后两项的差分（后一项减去前一项）。

`mean` 计算均值, `var` 计算样本方差或协方差矩阵, `sd` 计算样本标准差, `median` 计算中位数, `quantile` 计算样本分位数。 `cor` 计算相关系数。

`colSums`, `colMeans`, `rowSums`, `rowMeans` 对矩阵的每列或每行计算总和或者平均值，效率比用 `apply` 函数要高。

`rle` 和 `inverse.rle` 用来计算数列中“连”长度及其逆向恢复，“连”经常用在统计学的随机性检验中。

18.4.3 最值

`max` 和 `min` 求最大和最小, `cummax` 和 `cummin` 累进计算。

`range` 返回最小值和最大值两个元素。

对于 `max`, `min`, `range`, 如果有多个自变量可以把这些自变量连接起来后计算。

`pmax(x1,x2,...)` 对若干个等长向量计算对应元素的最大值, 不等长时短的被重复使用。`pmin` 类似。比如, `pmax(0, pmin(1,x))` 把 `x` 限制到 `[0,1]` 内。

18.4.4 排序

`sort` 返回排序结果。可以用 `decreasing=TRUE` 选项进行降序排序。`sort` 可以有一个 `partial=` 选项, 这样只保证结果中 `partial=` 指定的下标位置是正确的。比如:

```
sort(c(3,1,4,2,5), partial=3)
```

```
## [1] 2 1 3 4 5
```

只保证结果的第三个元素正确。可以用来计算样本分位数估计。

在 `sort()` 中用选项 `na.last` 指定缺失值的处理, 取 `NA` 则删去缺失值, 取 `TRUE` 则把缺失值排在最后面, 取 `FALSE` 则把缺失值排在最前面。

`order` 返回排序用的下标序列, 它可以有多个自变量, 按这些自变量的字典序排序。可以用 `decreasing=TRUE` 选项进行降序排序。如果只有一个自变量, 可以使用 `sort.list` 函数。

`rank` 计算秩统计量, 可以用 `ties.method` 指定同名次处理方法, 如 `ties.method=min` 取最小秩。

`order`, `sort.list`, `rank` 也可以有 `na.last` 选项, 只能为 `TRUE` 或 `FALSE`。

`unique()` 返回去掉重复元素的结果, `duplicated()` 对每个元素用一个逻辑值表示是否与前面某个元素重复。如

```
unique(c(1,2,2,3,1))
```

```
## [1] 1 2 3
```

```
duplicated(c(1,2,2,3,1))
```

```
## [1] FALSE FALSE TRUE FALSE TRUE
```

rev 反转序列。

18.4.5 一元定积分 integrate

integrate(f, lower, upper) 对一元函数 f 计算从 lower 到 upper 的定积分。使用自适应算法保证精度。如：

```
integrate(sin, 0, pi)
```

```
## 2 with absolute error < 2.2e-14
```

函数的返回值不仅仅包含定积分数值，还包含精度等信息。

18.4.6 一元函数求根 uniroot

uniroot(f, interval) 对函数 f 在给定区间内求一个根，interval 为区间的两个端点。要求 f 在两个区间端点的值异号。即使有多个根也只能给出一个。如

```
uniroot(function(x) x*(x-1)*(x+1), c(-2, 2))
```

```
## $root
```

```
## [1] 0
```

```
##
```

```
## $f.root
```

```
## [1] 0
```

```
##
```

```
## $iter
```

```
## [1] 1
```

```
##
```

```
## $init.it
```

```
## [1] NA
```

```
##
```

```
## $estim.prec
## [1] 2
```

对于多项式，可以用 `polyroot` 函数求出所有的复根。

18.4.7 离散傅立叶变换 `fft`

R 中 `fft` 函数使用快速傅立叶变换算法计算离散傅立叶变换。设 x 为长度 n 的向量， $y = \text{fft}(x)$ ，则

```
y[k] = sum(x * complex(
  argument = -2*pi * (0:(n-1)) * (k-1)/n))
```

即

$$y_{k+1} = \sum_{j=0}^{n-1} x_{j+1} \exp\left(-i2\pi \frac{kj}{n}\right), \quad k = 0, 1, \dots, n-1.$$

注意没有除以 n ，结果是复数向量。

另外，若 $y = \text{fft}(x)$ ， $z = \text{fft}(y, \text{inverse}=T)$ ，则 $x == z/\text{length}(x)$ 。

快速傅立叶变换是数值计算中十分常用的工具，R 软件包 `fftw` 可以进行优化的快速傅立叶变换。

18.4.8 用 `filter` 函数作迭代

R 在遇到向量自身迭代时很难用向量化编程解决，`filter` 函数可以解决其中部分问题。`filter` 函数可以进行卷积型或自回归型的迭代。语法为

```
filter(x, filter,
  method = c("convolution", "recursive"),
  sides=2, circular =FALSE, init)
```

下面用例子演示此函数的用途。

18.4.8.1 示例 1: 双侧滤波

对输入序列 $x_t, t = 1, 2, \dots, n$, 希望进行如下滤波计算:

$$y_t = \sum_{j=-k}^k a_j x_{t-j}, \quad k+1 \leq t \leq n-k-1,$$

其中 $(a_{-k}, \dots, a_0, \dots, a_k)$ 是长度为 $2k+1$ 的向量。注意公式中 a_j 与 x_{t-j} 对应。

假设 x 保存在向量 x 中, $(a_{-k}, \dots, a_0, \dots, a_k)$ 保存在向量 f 中, y_{k+1}, \dots, y_{n-k} 保存在向量 y 中, 无定义部分取 NA, 程序可以写成

```
y <- filter(x, f, method="convolution", sides=2)
```

比如, 设 $x = (1, 3, 7, 12, 17, 23)$, $(a_{-1}, a_0, a_1) = (0.1, 0.5, 0.4)$, 则

$$y_t = 0.1 \times x_{t+1} + 0.5 \times x_t + 0.4 \times x_{t-1}, \quad t = 2, 3, \dots, 5$$

用 `filter()` 函数计算, 程序为:

```
y <- filter(c(1,3,7,12,17,23), c(0.1, 0.5, 0.4),
           method="convolution", sides=2)
```

```
y
## Time Series:
## Start = 1
## End = 6
## Frequency = 1
## [1] NA 2.6 5.9 10.5 15.6 NA
```

18.4.8.2 示例 2: 单侧滤波

对输入序列 $x_t, t = 1, 2, \dots, n$, 希望进行如下滤波计算:

$$y_t = \sum_{j=0}^k a_j x_{t-j}, \quad k+1 \leq t \leq n,$$

其中 (a_0, \dots, a_k) 是长度为 $k+1$ 的向量。注意公式中 a_j 与 x_{t-j} 对应。

假设 x 保存在向量 x 中, (a_0, \dots, a_k) 保存在向量 f 中, y_{k+1}, \dots, y_n 保存在向量 y 中, 无定义部分取 NA, 程序可以写成

```
y <- filter(x, f, method="convolution", sides=1)
```

比如, 设 $x = (1, 3, 7, 12, 17, 23)$, $(a_0, a_1, a_2) = (0.1, 0.5, 0.4)$, 则

$$y_t = 0.1 \times x_t + 0.5 \times x_{t-1} + 0.4 \times x_{t-2}, \quad t = 3, 4, \dots, 6$$

程序为

```
y <- filter(c(1,3,7,12,17,23), c(0.1, 0.5, 0.4),
           method="convolution", sides=1)
y
## Time Series:
## Start = 1
## End = 6
## Frequency = 1
## [1] NA NA 2.6 5.9 10.5 15.6
```

18.4.8.3 示例 3: 自回归迭代

设输入 $e_t, t = 1, 2, \dots, n$, 要计算

$$y_t = \sum_{j=1}^k a_j y_{t-j} + e_t, \quad t = 1, 2, \dots, n,$$

其中 (a_1, \dots, a_k) 是 k 个实数, (y_{-k+1}, \dots, y_0) 已知。

设 x 保存在向量 \mathbf{x} 中, (a_1, \dots, a_k) 保存在向量 \mathbf{a} 中, (y_1, \dots, y_n) 保存在向量 \mathbf{y} 中。

如果 (y_{-k+1}, \dots, y_0) 都等于零, 可以用如下程序计算 y_1, y_2, \dots, y_n :

```
filter(x, a, method="recursive")
```

如果 (y_0, \dots, y_{-k+1}) 保存在向量 \mathbf{b} 中 (注意与时间顺序相反), 可以用如下程序计算 y_1, y_2, \dots, y_n :

```
filter(x, a, method="recursive", init=b)
```

比如, 设 $e = (0.1, -0.2, -0.1, 0.2, 0.3, -0.2)$, $(a_1, a_2) = (0.9, 0.1)$, $y_{-1} = y_0 = 0$, 则

$$y_t = 0.9 \times y_{t-1} + 0.1 \times y_{t-2} + e_t, \quad t = 1, 2, \dots, 6$$

迭代程序和结果为

```
y <- filter(c(0.1, -0.2, -0.1, 0.2, 0.3, -0.2),
           c(0.9, 0.1), method="recursive")
print(y, digits=3)
## Time Series:
## Start = 1
## End = 6
## Frequency = 1
## [1] 0.1000 -0.1100 -0.1890 0.0189 0.2981 0.0702
```

这个例子中, 如果已知 $y_0 = 200, y_{-1} = 100$, 迭代程序和结果为:

```
y <- filter(c(0.1, -0.2, -0.1, 0.2, 0.3, -0.2),
           c(0.9, 0.1), init=c(200, 100),
           method="recursive")
print(y, digits=6)
## Time Series:
## Start = 1
## End = 6
## Frequency = 1
## [1] 190.100 190.890 190.711 190.929 191.207 190.979
```

18.5 并行计算

现代桌面电脑和笔记本电脑的 CPU 通常有多个核心或虚拟核心 (线程), 如 2 核心或 4 虚拟核心。通常 R 运行并不能利用全部的 CPU 能力, 仅能利用其中的一个虚拟核心。使用特制的 BLAS 库 (非 R 原有) 可以并发置信多个线程, 一些 R 扩展包也可以利用多个线程。利用多台计算机、多个 CPU、CPU 中的多核心和多线程同时完成一个计算任务称为并行计算。

想要充分利用多个电脑、多个 CPU 和 CPU 内的虚拟核心, 技术上比较复杂,

涉及到计算机之间与进程之间的通讯问题，在要交流的数据量比较大时会造成并行计算的瓶颈。

实际上，有些问题可以很容易地进行简单地并行计算。比如，在一个统计研究中，需要对 100 组参数组合进行模拟，评估不同参数组合下模型的性能。假设研究人员有两台安装了 R 软件的计算机，就可以在两台计算机上进行各自 50 组参数组合的模拟，最后汇总在一起就可以了。

R 的 `parallel` 包提供了一种比较简单的利用 CPU 多核心的功能，思路与上面的例子类似，如果有多个任务互相之间没有互相依赖，就可以分解到多个计算机、多个 CPU、多个虚拟核心中并行计算。最简单的情形是一台具有单个 CPU、多个虚拟核心的台式电脑或者笔记本电脑。但是，统计计算中最常见耗时计算任务是随机模拟，随机模拟要设法避免不同进程的随机数序列的重复可能，以及同一进程中不同线程的随机数序列的重复可能。

`parallel` 包提供了 `parLapply()`、`parSapply()`、`parApply()` 函数，作为 `lapply()`、`sapply()`、`apply()` 函数的并行版本，与非并行版本相比，需要用一个临时集群对象作为第一自变量。

18.5.1 例 1：完全不互相依赖的并行运算

考虑如下计算问题：

$$S_{k,n} = \sum_{i=1}^n \frac{1}{i^k}$$

下面的程序取 n 为一百万， k 为 2 到 21，循环地用单线程计算。

```
f10 <- function(k=2, n=1000){
  s <- 0.0
  for(i in seq(n)) s <- s + 1/i^k
  s
}
f11 <- function(n=1000000){
  nk <- 20
  v <- sapply(2:(nk+1), function(k) f10(k, n))
  v
}
```

```
system.time(f11())[3]
## elapsed
## 2.87
```

因为对不同的 k , $f0(k)$ 计算互相不依赖, 也不涉及到随机数序列, 所以可以简单地并行计算而没有任何风险。先查看本计算机的虚拟核心(线程)数:

```
library(parallel)
detectCores()
## [1] 8
```

用 `makeCluster()` 建立临时的有 8 个节点的单机集群:

```
nNodes <- 8
cpucl <- makeCluster(nNodes)
```

用 `parSapply()` 或者 `parLapply()` 关于 k 并行地循环:

```
f12 <- function(n=1000000){
  f10 <- function(k=2, n=1000){
    s <- 0.0
    for(i in seq(n)) s <- s + 1/i^k
    s
  }

  nk <- 20
  v <- parSapply(cpucl, 2:(nk+1), function(k) f10(k, n))
  v
}
system.time(f12())[3]
## elapsed
## 1.19
```

并行版本速度提高了 140% 左右。

并行执行结束后, 需要解散临时的集群, 否则可能会有内存泄漏:

```
stopCluster(cpucl)
```

注意并行版本的程序还需要一些在每个计算节点上的初始化，比如调入扩展包，定义函数，初始化不同的随机数序列等。parallel 包的并行执行用的是不同的进程，所以传送给每个节点的计算函数要包括所有的依赖内容。比如，f2() 中内嵌了 f0() 的定义，如果不将 f0() 定义在 f2() 内部，就需要预先将 f0() 的定义传递给每个节点。

parallel 包的 clusterExport() 函数可以用来把计算所依赖的对象预先传送到每个节点。比如，上面的 f2() 可以不包含 f0() 的定义，而是用 clusterExport() 预先传递：

```
cpucl <- makeCluster(nNodes)
clusterExport(cpucl, c("f10"))
f13 <- function(n=1000000){
  nk <- 20
  v <- parSapply(cpucl, 2:(nk+1), function(k) f10(k, n))
  v
}
system.time(f13())[3]
## elapsed
## 1.08
stopCluster(cpucl)
```

如果需要在每个节点预先执行一些语句，可以用 clusterEvalQ() 函数执行，如

```
clusterEvalQ(cpucl, library(dplyr))
```

18.5.2 例 2：使用相同随机数序列的并行计算

为了估计总体中某个比例 p 的置信区间，调查了一组样本，在 n 个受访者中选“是”的比例为 \hat{p} 。令 λ 为标准正态分布的双侧 α 分位数，参数 p 的近似 $1 - \alpha$

置信区间为

$$\frac{\hat{p} + \frac{\lambda^2}{2n}}{1 + \frac{\lambda^2}{n}} \pm \frac{\lambda}{\sqrt{n}} \frac{\sqrt{\hat{p}(1 - \hat{p}) + \frac{\lambda^2}{4n}}}{1 + \frac{\lambda^2}{n}}$$

称为 Wilson 置信区间。

假设要估计不同 $1 - \alpha, n, p$ 情况下，置信区间的覆盖率（即置信区间包含真实参数 p 的概率）。可以将这些参数组合定义成一个列表，列表中每一项是一种参数组合，对每一组合分别进行随机模拟，估计覆盖率。因为不同参数组合之间没有互相依赖的关系，随机数序列完全可以使用同一个序列。

不并行计算的程序示例：

```
wilson <- function(n, x, conf){
  hatp <- x/n
  lam <- qnorm((conf+1)/2)
  lam2 <- lam^2 / n
  p1 <- (hatp + lam2/2)/(1 + lam2)
  delta <- lam / sqrt(n) * sqrt(hatp*(1-hatp) + lam2/4) / (1 + lam2)
  c(p1-delta, p1+delta)
}

f20 <- function(cpar){
  set.seed(101)
  conf <- cpar[1]
  n <- cpar[2]
  p0 <- cpar[3]
  nsim <- 100000
  cover <- 0
  for(i in seq(nsim)){
    x <- rbinom(1, n, p0)
    cf <- wilson(n, x, conf)
    if(p0 >= cf[1] && p0 <= cf[2]) cover <- cover+1
  }
  cover/nsim
}

f21 <- function(){
```

```
dp <- rbind(rep(c(0.8, 0.9), each=4),
            rep(rep(c(30, 100), each=2), 2),
            rep(c(0.5, 0.1), 4))
lp <- as.list(as.data.frame(dp))
res <- sapply(lp, f20)
res
}
system.time(f21())[3]
## elapsed
## 4.3
```

约运行 4.3 秒。

改为并行版本：

```
library(parallel)
nNodes <- 8
cpucl <- makeCluster(nNodes)
clusterExport(cpucl, c("f20", "wilson"))
f22 <- function(){
  dp <- rbind(rep(c(0.8, 0.9), each=4),
              rep(rep(c(30, 100), each=2), 2),
              rep(c(0.5, 0.1), 4))
  lp <- as.list(as.data.frame(dp))
  res <- parSapply(cpucl, lp, f20)
  res
}
system.time(f22())[3]
## elapsed
## 1.25
stopCluster(cpucl)
```

运行约 1.25 秒，速度提高 240% 左右。这里模拟了 8 种参数组合，每种参数组合模拟了十万次，每种参数组合模拟所用的随机数序列是相同的。

18.5.3 例 3: 使用独立随机数序列的并行计算

大量的耗时的统计计算是随机模拟，有时需要并行计算的部分必须使用独立的随机数序列。比如，需要进行一千次重复模拟，每次使用不同的随机数序列，可以将其分解为 10 组模拟，每组模拟一百万次，这就要求这 10 组模拟使用的随机数序列不重复。

R 中实现了 L'Ecuyer 的多步迭代复合随机数发生器，此随机数发生器周期很长，而且很容易将发生器的状态前进指定的步数。parallel 包的 `nextRNGStream()` 函数可以将该发生器前进到下一段的开始，每一段都足够长，可以用于一个节点。

以 Wilson 置信区间的模拟为例。设 $n = 30$, $p = 0.01$, $1 - \alpha = 0.95$ ，取重复模拟次数为 1 千万次，估计 Wilson 置信区间的覆盖率。单线程版本为：

```
f31 <- function(nsim=1E7){
  set.seed(101)
  n <- 30; p0 <- 0.01; conf <- 0.95
  cover <- 0
  for(i in seq(nsim)){
    x <- rbinom(1, n, p0)
    cf <- wilson(n, x, conf)
    if(p0 >= cf[1] && p0 <= cf[2]) cover <- cover+1
  }
  cover/nsim
}
system.time(cvg1 <- f31())[3]
## elapsed
## 42.61
```

单线程版本运行了大约 43 秒。

改成并行版本。比例 2 多出的部分是为每个节点分别计算一个随机数种子将不同的种子传给不同节点。parallel 包的 `clusterApply()` 函数为临时集群的每个节点分别执行同一函数，但对每个节点分别使用列表的不同元素作为函数的自变量。

```
library(parallel)
nNodes <- 8
cpucl <- makeCluster(nNodes)
each.seed <- function(s){
  assign(".Random.seed", s, envir = .GlobalEnv)
}
RNGkind("L'Ecuyer-CMRG")
set.seed(101)
seed0 <- .Random.seed
seeds <- as.list(1:nNodes)
for(i in 1:nNodes){ # 给每个节点制作不同的种子
  seed0 <- nextRNGStream(seed0)
  seeds[[i]] <- seed0
}
## 给每个节点传送不同种子:
junk <- clusterApply(cpucl, seeds, each.seed)

f32 <- function(isim, nsimsub=10000){
  n <- 30; p0 <- 0.01; conf <- 0.95
  cover <- 0
  for(i in seq(nsimsub)){
    x <- rbinom(1, n, p0)
    cf <- wilson(n, x, conf)
    if(p0 >= cf[1] && p0 <= cf[2]) cover <- cover+1
  }
  cover
}

clusterExport(cpucl, c("f32", "wilson"))

f33 <- function(nsim=1E7){
  nbatch <- 40
  nsimsub <- nsim / nbatch
  cvs <- parSapply(cpucl, 1:nbatch, f32, nsimsub=nsimsub)
```

```
print(cvs)
sum(cvs)/(nsim*nbatch)
}

system.time(cvg2 <- f33())[3]
## [1] 963759 963660 963885 963739 963714 964171 963615 963822 963720 963939
## elapsed
## 13.63
stopCluster(cpucl)
```

并行版本运行了大约 14 秒，速度提高约 210%。从两个版本各自一千万次重复模拟结果来看，用随机模拟方法得到的覆盖率估计的精度大约为 3 位有效数字。

更大规模的随机模拟问题，可以考虑使用多 CPU 的计算工作站或者服务器，或用多台计算机通过高速局域网组成并行计算集群。

还有一种选择是租用云计算服务。

Chapter 19

随机模拟

19.1 随机数

随机模拟是统计研究的重要方法，另外许多现代统计计算方法（如 MCMC）也是基于随机模拟的。R 中提供了多种不同概率分布的随机数函数，可以批量地产生随机数。一些 R 扩展包利用了随机模拟方法，如 `boot` 包进行 bootstrap 估计。

所谓随机数，实际是“伪随机数”，是从一组起始值（称为种子），按照某种递推算法向前递推得到的。所以，从同一种子出发，得到的随机数序列是相同的。

为了得到可重现的结果，随机模拟应该从固定不变的种子开始模拟。用 `set.seed(k)` 指定一个编号为 `k` 的种子，这样每次从编号 `k` 种子运行相同的模拟程序可以得到相同的结果。

还可以用 `set.seed()` 加选项 `kind=` 指定后续程序要使用的随机数发生器名称，用 `normal.kind=` 指定要使用的正态分布随机数发生器名称。

R 提供了多种分布的随机数函数，如 `runif(n)` 产生 `n` 个标准均匀分布随机数，`rnorm(n)` 产生 `n` 个标准正态分布随机数。例如：

```
round(runif(5), 2)
## [1] 0.44 0.56 0.93 0.23 0.22
```

```
round(rnorm(5), 2)
## [1] -0.20  1.10 -0.02  0.16  2.02
```

注意因为没有指定种子，每次运行会得到不同的结果。

在 R 命令行运行

```
?Distributions
```

可以查看 R 中提供的不同概率分布。

19.2 sample() 函数

sample() 函数从一个有限集中无放回或有放回地随机抽取，产生随机结果。

例如，为了设随机变量 X 取值于 {正面,反面}，且 $P(X = \text{正面}) = 0.7 = 1 - P(X = \text{反面})$ ，如下程序产生 X 的 10 个随机抽样值：

```
sample(c('正面', '反面'), size=10,
       prob=c(0.7, 0.3), replace=TRUE)
## [1] "反面" "反面" "反面" "反面" "正面"
## [6] "正面" "正面" "正面" "反面" "反面"
```

sample() 的选项 size 指定抽样个数，prob 指定每个值的概率，replace=TRUE 说明是有放回抽样。

如果要做无放回等概率的随机抽样，可以不指定 prob 和 replace(缺省是 FALSE)。比如，下面的程序从 1:10 随机抽取 4 个：

```
sample(1:10, size=4)
## [1] 1 5 8 10
```

如果要从 $1:n$ 中等概率无放回随机抽样直到每一个都被抽过，只要用如：

```
sample(10)
## [1] 3 5 9 2 10 7 4 1 6 8
```

这实际上返回了 $1:10$ 的一个重排。

dplyr 包的 `sample_n()` 函数与 `sample()` 类似，但输入为数据框，输出为随机抽取的数据框行子集。

19.3 随机模拟示例

19.3.1 线性回归模拟

考虑如下线性回归模型

$$y = 10 + 2x + \varepsilon, \varepsilon \sim N(0, 0.5^2).$$

假设有样本量 $n = 10$ 的一组样本，R 函数 `lm()` 可以得到截距 a ，斜率 b 的估计 \hat{a}, \hat{b} ，以及相应的标准误差 $SE(\hat{a}), SE(\hat{b})$ 。样本可以模拟产生。

模型中的自变量 x 可以用随机数产生，比如，用 `sample()` 函数从 $1:10$ 中随机有放回地抽取 n 个。模型中的随机误差项 ε 可以用 `rnorm()` 产生。产生一组样本的程序如：

```
n <- 10; a <- 10; b <- 2
x <- sample(1:10, size=n, replace=TRUE)
eps <- rnorm(n, 0, 0.5)
y <- a + b * x + eps
```

如下程序计算线性回归：

```
lm(y ~ x)

##
## Call:
## lm(formula = y ~ x)
##
## Coefficients:
## (Intercept)          x
##      9.968         1.986
```

如下程序计算线性回归的多种统计量：

```
summary(lm(y ~ x))

##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.8143 -0.3718  0.1188  0.3165  0.5401
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  9.96827    0.55005   18.12 8.83e-08 ***
## x            1.98635    0.07882   25.20 6.58e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4865 on 8 degrees of freedom
## Multiple R-squared:  0.9876, Adjusted R-squared:  0.986
## F-statistic: 635.1 on 1 and 8 DF,  p-value: 6.581e-09
```

如下程序返回一个矩阵，包括 a, b 的估计值、标准误差、 t 检验统计量、检验 p 值：

```
summary(lm(y ~ x))$coefficients

##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept) 9.968269 0.55005071 18.12245 8.827682e-08
## x           1.986352 0.07882033 25.20101 6.581144e-09
```

如下程序把上述矩阵的前两列拉直成一个向量返回：

```
c(summary(lm(y ~ x))$coefficients[,1:2])

## [1] 9.96826910 1.98635176 0.55005071 0.07882033
```

这样得到 $\hat{a}, \hat{b}, SE(\hat{a}), SE(\hat{b})$ 这四个值。

用 `replicate()` (复合语句) 执行多次模拟, 返回向量或矩阵结果, 返回矩阵时, 每列是一次模拟的结果。下面是线性回归整个模拟程序, 写成了一个函数。

```
reg.sim <- function(
  a=10, b=2, sigma=0.5,
  n=10, B=1000){
  set.seed(1)
  resm <- replicate(B, {
    x <- sample(1:10, size=n, replace=TRUE)
    eps <- rnorm(n, 0, 0.5)
    y <- a + b * x + eps
    c(summary(lm(y ~ x))$coefficients[,1:2])
  })
  resm <- t(resm)
  colnames(resm) <- c('a', 'b', 'SE.a', 'SE.b')
  cat(B, ' 次模拟的平均值:\n')
  print( apply(resm, 2, mean) )
  cat(B, ' 次模拟的标准差:\n')
  print( apply(resm, 2, sd) )
}
```

运行测试:

```
set.seed(1)
reg.sim()
```

```
## 1000 次模拟的平均值:
##          a          b      SE.a      SE.b
## 9.9970476 1.9994490 0.3639505 0.0592510
## 1000 次模拟的标准差:
##          a          b      SE.a      SE.b
## 0.37974881 0.06297733 0.11992515 0.01795926
```

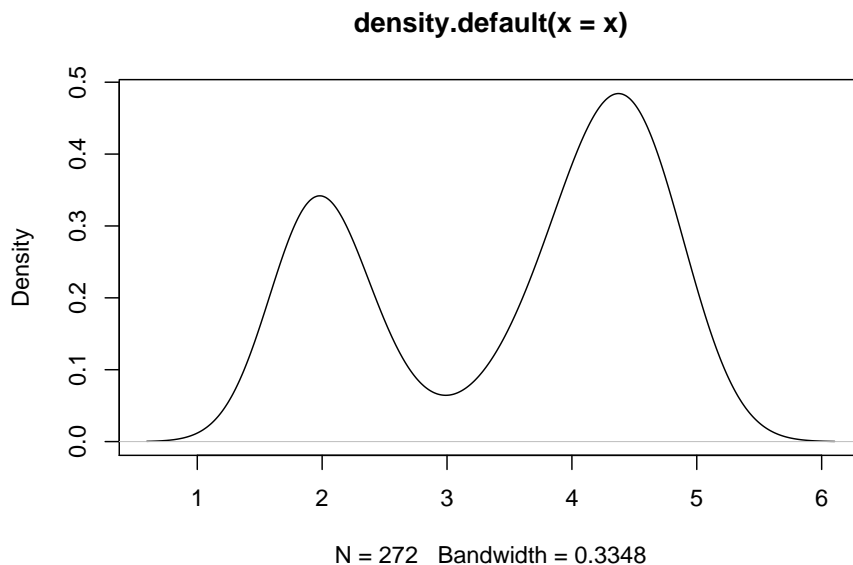
可以看出, 标准误差作为 \hat{a}, \hat{b} 的标准差估计, 与多次模拟得到多个 \hat{a}, \hat{b} 样本计算得到的标准差估计是比较接近的。结果中 $SE(\hat{a})$ 的平均值为 0.363, 1000 次模拟的 \hat{a} 的样本标准差为 0.393, 比较接近; $SE(\hat{b})$ 的平均值为 0.0594, 1000

次模拟的 \hat{b} 的样本标准差为 0.0637, 比较接近。

19.3.2 核密度的 bootstrap 置信区间

R 自带的数据库 `faithful` 内保存了美国黄石国家公园 Faithful 火山的 272 次爆发持续时间和间歇时间。为估计爆发持续时间的密度, 可以用核密度估计方法, R 函数 `density` 可以执行此估计, 返回 N 个格子点上的密度曲线坐标:

```
x <- faithful$eruptions
est0 <- density(x)
plot(est0)
```



这个密度估计明显呈现出双峰形态。

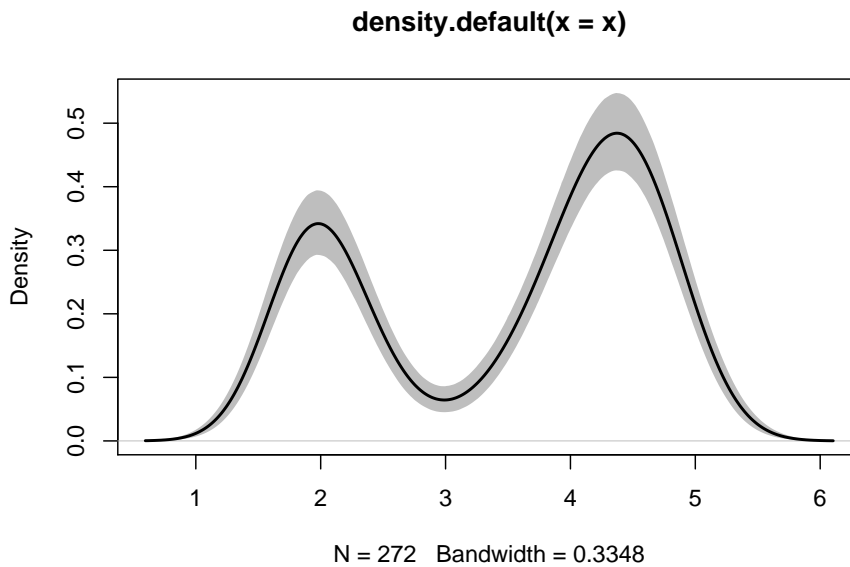
核密度估计是统计估计, 为了得到其置信区间 (给定每个 x 坐标, 真实密度 $f(x)$ 的单个的置信区间), 采用如下非参数 bootstrap 方法:

重复 $B = 10000$ 次, 每次从原始样本中有重复地抽取与原来大小相同的一组样本, 对这组样本计算核密度估计, 结果为 $(x_i, y_i^{(j)}), i = 1, 2, \dots, N, j = 1, 2, \dots, B$, 每组样本估计 N 个格子点的密度曲线坐标, 横坐标不随样本改变。

对每个横坐标 x_i , 取 bootstrap 得到的 B 个 $y_i^{(j)}, j = 1, 2, \dots, B$ 的 0.025 和 0.975 样本分位数, 作为真实密度 $f(x_i)$ 的 bootstrap 置信区间。

在 R 中利用 `replicate()` 函数实现:

```
set.seed(1)
resm <- replicate(10000, {
  x1 <- sample(x, replace=TRUE)
  density(x1, from=min(est0$x),
          to=max(est0$x))$y
})
CI <- apply(resm, 1, quantile, c(0.025, 0.975))
plot(est0, ylim=range(CI), type='n')
polygon(c(est0$x, rev(est0$x)),
        c(CI[1,], rev(CI[2,])),
        col='grey', border=FALSE)
lines(est0, lwd=2)
```



程序中用 `set.seed(1)` 保证每次运行得到的结果是不变的, `replicate()` 函数第一参数是重复模拟次数, 第二参数是复合语句, 这些语句是每次模拟要执

行的计算。在每次模拟中，用带有 `replace=TRUE` 选项的 `sample()` 函数从样本中有放回地抽样得到一组 bootstrap 样本，每次模拟的结果是在指定格子点上计算的核密度估计的纵坐标。`replicate()` 的结果为一个矩阵，每一列是一次模拟得到的纵坐标集合。对每个横坐标格子点，用 `quantile()` 函数计算 B 个 bootstrap 样本的 2.5% 和 97.5% 分位数，循环用 `apply()` 函数表示。`polygon()` 函数指定一个多边形的顺序的顶点坐标用 `col=` 指定的颜色填充，本程序中实现了置信下限与置信上限两条曲线之间的颜色填充。`lines()` 函数绘制了与原始样本对应的核密度估计曲线。

Part IV

用 R 制作研究报告和图书

Chapter 20

用 R 制作研究报告

一个统计或数据分析的科研项目，都会产生一个或多个研究报告。因为使用统计与数据分析不可避免地有很多计算涉及在内，这里假设使用 R 软件做了计算。科研是一个不断改进的过程，所以每一次重新做了计算，研究报告中的汇总表格、图形都要更新。这样的任务比较繁琐，也容易出错。

“文学式编程”(literate programming, (Knuth, 1984)) 是这样一种思想，把撰写报告与计算程序有机地结合在一起，用一个源文件文件既包含报告内容，又包含计算程序。每次产生研究报告时，先运行源文件中的计算程序得到计算结果，这些结果包括文字性内容与图形，然后利用适当软件自动地把这些原始文字、计算结果组合成最终的报告。利用这样的思想，可以自动生成重复的例行报告，还可以作为“可重复科学研究”的载体。

上述的源文件一般是文本文件，格式可以是 Markdown 格式，也可以是 LaTeX 格式等许多格式。R 的 knitr 软件包就是用来支持这样的编程思想的一个扩展包。

Markdown 是一种很简单的文本文件格式，通常保存为.md 扩展名，利用 R 程序进行扩展的版本保存为.Rmd 扩展名。Markdown 文件里面有一些简单的格式标注方法，比如两个星号之间的文字会转化为斜体，缩进四个空格或一个制表符的内容会看成代码。Markdown 仅适用于比较简单的文章、源程序说明等，不太适用于复杂的含有大量数学公式、图表的文章，从 markdown 格式比较适合转换为 html(网页) 格式，也可以转换为 MS Word 的 docx 格式，通过

LaTeX 编译器可以转换为 PDF 格式，也可以将 docx 转存为 PDF 格式，或者用输出到 PDF 文件的打印机将网页格式转换为 PDF 格式。

LaTeX 是一个文档排版系统，功能强大，结果美观，设计合理。缺点是需要学习类似于 HTML 的另一种语言。LaTeX 源文件主要是编译为 PDF。

R 扩展包 knitr 包支持在 Markdown 格式、LaTeX 格式等类型的文件中插入 R 代码，经过转换，文中的 R 代码可以变成代码的结果文字、表格、图形，与原有报告文字有机地结合在一起。

插入了 R 代码的 markdowng 格式的文件一般以 .Rmd 为扩展名，R 的 rmarkdown 扩展包和 knitr 包一起为 Rmd 格式提供了支持。

R 的 bookdown 包进一步增强了 R Markdown 格式的功能，支持生成 PDF、多文件互相链接的 HTML、Word 等输出，其中的表格、图形可以变成浮动表，公式、定理可以自动编号并支持文献、公式、定理、图表、章节的引用和链接，所以比较适用于编写一本书或论文、研究报告。

knitr 包也支持在 LaTeX 源文件中插入 R 代码，通过编译使得 R 代码运行结果、图形自动插入生成的研究报告中。原来有一个 Sweave 扩展包是支持在 LaTeX 中插入 R 代码的，现在 knitr 的功能已经涵盖并增强了此包功能。

下面几章分别介绍 markdown 格式、R markdown 格式、bookdown 扩展包、简易网站制作和幻灯片制作。参考：

- (Xie et al., 2019)
- (Xie, 2017)

利用 R Markdown 格式可以执行如下的任务：

- 将单一的 R Markdown 文件编译为不同的格式，如：
 - HTML;
 - PDF;
 - MS Word。
- 在 RStudio 软件内制作笔记本文档，可以在其中包含说明文字与 R 代码，可以在笔记本文档内交互地运行 R 代码并能将结果交互地显示在笔记本内。
- 生成演示幻灯片，可以是基于 HTML5 的，也可以是基于 LaTeX beamer 扩展包的，或者 MS Powerpoint。

- 编写多章节组成的书籍。
- 编写期刊论文。
- 制作网站或者博客。
- 制作商业智能仪表盘 (dashboards)，即数据可视化展示。

Chapter 21

Markdown 格式

21.1 介绍

Markdown 是一种很简单的文本文件格式，通常保存为.md 扩展名。Markdown 中文内容应该使用 UTF-8 编码。Markdown 文件里面有一些简单的格式标注方法，比如两个星号之间的文字会转化为斜体，缩进四个空格或一个制表符的内容会看成代码。

Markdown 适用于比较简单的文章、源程序说明等，不太适用于复杂的含有大量数学公式、图表的文章，从 markdown 格式比较适合转换为 html(网页) 格式，也可以转换为 MS Word 的 docx 格式，通过 docx 格式可以转存为 PDF 格式。R 扩展包 knitr、rmarkdown 和 bookdown 与 pandoc 软件一起大大扩展了 markdown 格式的适用范围。

如果需要作为一本书发布在网站上或者出版，可考虑使用 R 的 bookdown 扩展包。如果需要对出版格式进行精确控制，可考虑用 LaTeX 格式，LaTeX 格式很复杂，学习比较困难，但是表达能力强。

21.2 Markdown 格式文件的应用

Markdown 格式用在一些微博、论坛中作为缺省格式，用户在网络浏览器软件的输入框中按照 markdown 格式输入，网站自动将其转换为 html 富文本内容显示出来。

因为 markdown 格式就是纯文本，而且其格式十分简单，所以可以仅仅用普通文本编辑器编写 markdown 格式的文件，不需要转换为其它格式。

有一些独立运行的编辑软件，在其中输入了 markdown 格式文件后，可以并列地显示转化为 html 富文本后的结果，很多还支持自动备份到云服务中。但是，这些独立运行的编辑器大都有商业因素。比如，国内的印象笔记软件是提供了云存储的笔记软件，在其 PC 端就支持创建 markdown 格式的笔记，可以在输入时即时显示 html 效果。

RStudio 软件不是专用的 Markdown 编辑器，它是一个 R 程序的集成编辑、运行环境，对个人用户免费，在 R Studio 中可以编辑 Markdown 文件和含有 R 代码的 Markdown 文件，可以一键将其转化成 HTML、MS Word docx 文件，在单独安装的 LaTeX 编译软件支持下还可以直接编译成 PDF。RStudio 支持下的增强的 Markdown 格式，称为 RMarkdown 格式，以 `.Rmd` 为扩展名，支持大多数的 LaTeX 公式，在 bookdown 扩展包支持下还支持公式、定理、图表等自动编号和引用、链接。

21.3 markdown 格式说明

这部分内容参考了 markdown 手册和 Rstudio 的文档，以及 (Xie et al., 2019)。

21.3.1 概述

Markdown 格式是 John Gruber 于 2004 年创造的，Markdown 的目标是实现“易读易写”。Markdown 定义了一种简单好用的文本文件格式，作为单独的文本文件，此格式没有什么多余的标签，又可以转化为很多其它的格式。

Markdown 的语法全由一些符号所组成，这些符号经过精挑细选，其作用一目了然。比如：在文字两旁加上星号，看起来就像强调。Markdown 的列表看起

来就像我们平常在邮件中写一个列表的方法。Markdown 的区块引用看起来就真的像是引用一段文字，就像你曾在电子邮件中见过的那样。

需要时，可以直接在 markdown 中写 HTML 标记内容。markdown 能实现的功能是 HTML 的一部分，但是比 HTML 内容更干净，没有掺杂过多的与要表达的意思无关的标签。Markdown 的理念是，能让文档更容易读、写和随意改。

21.3.2 段落

一个段落由一行或连续的多行组成。段落之间以空行分隔。同一段落内的不同行在转换成 HTML 或 docx 等格式后会重新排列，原来的段内换行被当成了空格，这样的规定与 LaTeX 类似。普通段落不该用空格或制表符来缩进，不应在行尾留有空格。

为了在段内换行并且转化后仍保持段内换行，输入时在前面行的末尾输入两个或两个以上空格。例如：

白日依山尽，黄河入海流。
欲穷千里目，更上一层楼。

显示结果为：

白 日 依 山 尽， 黄 河 入 海 流。
欲穷千里目，更上一层楼。

这样做的缺点是末尾的空格时不可见的。可以使用 HTML 的 `
` 标签在段内换行，如：

白日依山尽，黄河入海流。
欲穷千里目，更上一层楼。

结果为

白日依山尽，黄河入海流。欲穷千里目，更上一层楼。

21.3.3 段内文字格式

在一段内，用星号或下划线包围的内容如 `* 强调 *` 是强调格式。用双星号或双下划线包围的内容如 `** 加重 **` 是加重格式。星号、下划线与要强调或加

重的内容之间不要空开，否则会当作普通星号或下划线解释，在行首还会当作列表。为了插入普通的星号或下划线，可以使用反斜杠保护，或者写成段内代码格式。

可以用一对 `~` 作为界定符给出下标，如 `H0~2~` 会变成 HO_2 。可以用一对 `^` 作为界定符给出上标，如 `Cu^2+^` 会变成 Cu^{2+} 。但是，数学公式一般还是应该使用 LaTeX 数学公式形式（见22.8）。

在普通段落内一部分内容希望显示成代码，对其中的特殊字符不进行解释，只要包在两个反向单撇号内。如 ``if(_x_ > 0) y=1;`` 会变成 `if(_x_ > 0) y=1;`。如果内容本身就包含反向单撇号，可以在两边使用更多个数的反向单撇号，如 ``` `x` ``` 会变成 ``x``。

在 Markdown 文件中，为了使得某些有特殊意义的字符不作特殊解释，可以在该字符前面加上反斜杠`\`，取消其特殊含义。

在 Rmarkdown 文件中，为了原样显示反向单撇号，可以在两边用双写的反向单撇号界定字符串。

21.3.4 标题和分隔线

以一个井号 `#` 开始的行是一级标题，以两个井号 `#` 开始的行是二级标题，.....，以六个井号 `#` 开始的行是六级标题。标题行前面应该空一行，否则可能把某些偶然出现在行首的 `#` 号误认为标题行的标志。

对一级标题，也可以用标题内容下面输入一行等于号 `=` 表示上一行内容是一级标题。对二级标题，可以用标题内容下面输入一行减号`-`表示上一行内容是二级标题。等于号和减号的个数不限。

用三个或三个以上连续的星号组成的行，可以转换成分隔线。下面是一个分隔线：

21.3.5 引用段落

可以用类似 Email 的回复包含原始邮件内容的办法输入引用段落，即，在段落的每行前面加一个大于号。比如下面的诗：

- > 白日依山尽，黄河入海流。
- > 欲穷千里目，更上一层楼。

转换成

白日依山尽，黄河入海流。欲穷千里目，更上一层楼。

注意引用也是段落模式，内容中的换行不起作用，空行导致分段。

引用段落也可以仅在段落第一行写大于号，其它行顶格写，例如下面的两段引用：

- > 远上寒山石径斜，
白云生处有人家。
- >
- > 停车坐爱枫林晚，
霜叶红于二月花。

转换成

远上寒山石径斜，白云生处有人家。

停车坐爱枫林晚，霜叶红于二月花。

引用也可以嵌套，如：

- > 张三说：李四这样说过
- >
- >> 不想当将军的木匠不是好厨子。
- >

转换成

张三说：李四这样说过

不想当将军的木匠不是好厨子。

注意嵌套内容前后都有空的引用行，否则不能实现嵌套引用。

引用内也可以嵌套其它的 Markdown 格式如标题、列表等。引用前后应该有空行把引用内容与其他内容分隔开。

21.3.6 列表

列表分为不编号的列表和编号的列表。不编号的列表转化后通常显示圆点开头的列表项。

在 Markdown 中，用星号表示一个不编号的列表项。星号也可以替换成加号或减号，后面必须有一个或多个空格。每个列表项可以输入多行，各行的内容最好左对齐，左对齐在使用文本格式时较易阅读，但不是必须的。两个列表项之间不要空行。例如：

- * 白日依山尽，
 黄河入海流。
- * 欲穷千里目，
 更上一层楼。

转换为

- 白日依山尽，黄河入海流。
- 欲穷千里目，更上一层楼。

段落顶头的数字加句点和空格表示编号列表，两个列表项之间尽量不要空行。例如：

1. 第一种解决方法，
 收买敌人的高官。
2. 第二种解决方法，
 尽可能拖延。

转换为

1. 第一种解决方法，收买敌人的高官。
2. 第二种解决方法，尽可能拖延。

标准的 markdown 编号列表不能自己定义数字的显示格式，不允许开始值不等于 1。pandoc 支持更自由的列表，允许输入时有括号或右括号，允许使用字母和罗马数字，但是括号会被去掉。如

- (1) 顶层一；
 - a) 内层二；
 - b) 内层三；

(2) 顶层二。

转换为

- (1) 顶层一;
 - a) 内层二;
 - b) 内层三;
- (2) 顶层二。

为了避免错误地产生非本意的编号列表，在行首写数字加句点和空格时，可以在句点前加反斜杠，或者在句点前加空格。例如，下面是一个年号：

2016\.

如上输入将不会错误地解释为有序列表。

如果列表项目中有多个段落，这时两个列表项之间应该以空行分隔，每个项目除了第一行外，输入的每行内容都应该缩进 4 个空格或者一个制表符。例如：

- * R语言第一个版本开发于1976-1980，基于Fortran；
于1980年移植到Unix，并对外发布源代码。
1984年出版的“棕皮书”
总结了1984年为止的版本，并开始发布授权的源代码。
这个版本叫做旧S。与我们现在用的S语言有较大差别。

1989--1988对S进行了较大更新，
变成了我们现在使用的S语言，称为第二版。
1988年出版的“蓝皮书”做了总结。

- * 1992年出版的“白皮书”描述了在S语言中实现的统计建模功能，
增强了面向对象的特性。软件称为第三版，这是我们现在用的多数版本。

1998年出版的“绿皮书”描述了第四版S语言，主要是编程功能的深层次改进。
现行的S系统并没有都采用第四版，S-PLUS的第5版才采用了S语言第四版。

转换为

- R语言第一个版本开发于1976-1980，基于Fortran；于1980年移植到Unix，并对外发布源代码。1984年出版的“棕皮书”总结了1984年为止

的版本, 并开始发布授权的源代码。这个版本叫做旧 S。与我们现在用的 S 语言有较大差别。

1989–1988 对 S 进行了较大更新, 变成了我们现在使用的 S 语言, 称为第二版。1988 年出版的“蓝皮书”做了总结。

- 1992 年出版的“白皮书”描述了在 S 语言中实现的统计建模功能, 增强了面向对象的特性。软件称为第三版, 这是我们现在用的多数版本。

1998 年出版的“绿皮书”描述了第四版 S 语言, 主要是编程功能的深层次改进。现行的 S 系统并没有都采用第四版, S-PLUS 的第 5 版才采用了 S 语言第四版。

列表项目内如果有引用段落, 需要都缩进 4 个空格。如果有程序代码, 需要缩进 4 个空格后用三个反单撇号表示开始与结束。

列表可以嵌套, 嵌套的列表需要缩进 4 个空格, 中间不需要空行。例如:

1. 第一类工作包括:
 - + 技术服务;
 - + 咨询服务。
2. 其它工作略。

转换为

1. 第一类工作包括:
 - 技术服务;
 - 咨询服务。
2. 其它工作略。

如果需要把每个列表项当作段落排版, 可以在每个列表项后空行。

21.3.7 源程序

为了让源程序能够自动显示成源程序的样式, 而不至于自动分行、特殊字符解释, 用空行把源程序与其它内容隔开, 并把源程序行都缩进 4 个空格(或以上)或者一个制表符。源程序格式持续到不缩进 4 个空格的地方为止。

例如, 下面的输入:

```
f <- function(x){  
  n <- length(x)  
  y <- numeric(n)  
  y[x >= 0] <- 1  
  ##y[x < 0] <- 0  
  
  y  
}
```

转换为

```
f <- function(x){  
  n <- length(x)  
  y <- numeric(n)  
  y[x >= 0] <- 1  
  ##y[x < 0] <- 0  
  
  y  
}
```

更适当的做法是用三个连续的反向单撇号表示代码开头与代码结束，中间就会当作源程序代码处理。例如下面的输入

```
...  
> x <- rnorm(100)  
> hist(x)  
...
```

转换为

```
> x <- rnorm(100)  
> hist(x)
```

R 的 knitr 包在 Markdown 格式的文件中插入 R 可执行代码时，就用了这样的方法。而且，R Markdown 格式的代码块不需要用空行与前后分隔开。

在 pandoc 程序的支持下，代码段还可以采用栅栏式代码段，在代码段开头前面一行加上至少三个连续 ~ 符号，在结尾后面一行加同样数目的 ~ 符号。这

样的代码段前后也必须空行以与其它内容分开。另外，如果代码内本身含有 ~ 行，只要使得开头与结尾标志中的 ~ 个数更多就可以了。例如下面的输入

```
~~~
#include <math.h>
double sqr(double x){
    return(x*x);
}
~~~
```

转换为

```
#include <math.h>
double sqr(double x){
    return(x*x);
}
```

使用栅栏式代码段时可以在开始行尾写大括号，在大括号内写选项。其中一种选项是要求按照某种编程语言对结果进行彩色语法显示，如 .cpp 表示 C++，.c 表示 C，.r 表示 R，.python 表示 python 等。选项 .numberLines 要求该代码行编号，选项 startFrom= 指定开始行号。如如：

```
~~~{.cpp .numberLines startFrom=101}
#include <math.h>
double sqr(double x){
    return(x*x);
}
~~~
```

转换为

```
101 #include <math.h>
102 double sqr(double x){
103     return(x*x);
104 }
```


21.3.8 链接

最简单的链接是原样显示的可点击的链接，只要把链接地址用小于号和大于号包在中间，两边用空格和其它内容隔开。如果是网页，需要加 `http://`，如果是邮箱，需要加 `mailto:`。例如，如下代码：

```
北京大学的网页地址是： <http://www.pku.edu.cn/> 。
```

显示为

北京大学的网页地址是： <http://www.pku.edu.cn/> 。

除此之外，Markdown 还支持两种形式的链接语法：行内式和引用式两种形式。不管是哪一种，链接的显示文字都是用方括号 [...] 来标记。

要建立一个行内式的链接，只要在方括号内写链接的显示文字，右方括号后面紧接着圆括号，并在圆括号中间插入网址链接即可。方括号部分与圆括号部分之间不能断开。例如：

```
请参考：[李东风的教学主页](http://www.math.pku.edu.cn/teachers/lidf/course/index.htm)
```

变成了

```
请参考：李东风的教学主页
```

在圆括号中，链接后面还可以包含用双撇号包围的标题文字，与链接之间用空格分开。

引用式的链接，需要在某处（比如文章结尾）定义一些链接的标识符，然后用方括号包围链接的显示文字，后面紧接着方括号包围着链接的标识符。

为了定义链接标识符，用方括号包围链接标识符，前面可以有至多 3 个空格缩进，右方括号后面紧接着冒号和一个空格，空格后写链接地址，然后空格，在两个双撇号中间写一个链接地址的标题。例如

```
[baidu]: http://baidu.com/ "百度"
```

```
[pku]: http://www.pku.edu.cn/ "北大"
```

这时，在文章中就可以用标识符调用链接，如 [北京大学主页][pku] 将变成链接北京大学主页。

使用引用式的链接，有些像论文中把所有参考文献排列在文章末尾，文中用到某一篇文献只要提及其序号。

还有一种链接是内部链接，用于文内跳转。在各级标题行的末尾，可以添加 `{# 自定义标签}` 这样的内容，其中“自定义标签”是自己写的一个标识符，标识符仅使用英文字母、数字、下划线、减号，用来区分不同的位置。比如，本文第一节“介绍”添加了 `markdown-intro` 为标签，就可以用“`[回到介绍](#markdown-intro)`”产生链接回到介绍。

21.3.9 插入图形

图形只能用链接形式，不可能保存到一个纯文本文件内。图形文件可以存在于远程服务器上，也可以是与生成的 HTML 文件在同一目录结构中的文件。语法仍使用行内式和参考式两种形式。转化成 HTML、PDF、Word 格式后可以把图形内嵌在输出文件内部。

行内式的图片链接，是普通行内链接格式前面添加了一个叹号，惊叹后面紧接着方括号，方括号内写图片的标题，标题可以空缺，在右方括号后面紧接着圆括号，圆括号内写图片的链接。

例如，下面的代码可以插入百度的一个 logo，使用的是网上的资源：

```

```

结果为



为了插入保存在本地的与 Markdown 源文件在同一子目录或下级子目录的图形，只要在圆括号中写图片文件名（如果与 Markdown 源文件在同一子目录）或相对路径。例如，在 `D:\work\figs` 下的图形文件 `baidu_logo.gif` 在本 Markdown 源文件所在目录的 `figs` 子目录中，可以用代码

```

```

结果为



这样含有网上图片和本地图片的 Markdown 源文件转化为 HTML 和 docx 格式，都可以正常显示插入的图片。

与链接类似，也可以在文章某处（比如末尾）定义图片的标识符，然后把行内图片引用中图片地址替换成图片标识符即可。

21.3.10 表格

Markdown 文本格式的表格就像是用减号、等号、竖线画的文本格式表格一样，转化为 HTML、docx 等格式后就变成了富文本的表格。有如下几种表格：

- 管道表
- 简单表
- 换行表
- 有格表

21.3.10.1 管道表

管道表在两列之间用竖线分开，在列标题下面用减号画横线，用如下方法指定各对齐方式：

- 在列标题下的横线开始加冒号，表示左对齐；
- 在列标题下的横线末尾加冒号，表示右对齐；
- 在列标题下的横线两端加冒号，表示居中对齐；

- 列标题下面仅有横线没有冒号，表示缺省对齐方式，一般是左对齐。

在表格内容后面空一行后写用 **Table:** 开头的表格说明。

这种方法不需要输入内容上下对齐，适用于中文内容。后面所讲的简单表、换行表、有格表需要能够输入内容对齐，对于中英文混合内容很难做到对齐，所以仅管道表比较适合中文内容。

例如:

```
| 姓名 | 收入 | 职业 | 颜色偏好 |
|:-----|:-----:|:-----:|:-----|
| 赵四海 | 123456 | 业务经理 | 红 |
| 刘英 | 50 | 无 | 蓝 |
| 钱德里 | 3200 | 保洁 | 灰 |
```

Table: 管道表示例

结果为

表 21.1: 管道表示例

姓名	收入	职业	颜色偏好
赵四海	123456	业务经理	红
刘英	50	无	蓝
钱德里	3200	保洁	灰

管道表不允许输入单元格换行，单元格内容太宽时转换结果可能自动换行，自动换行时列宽度与输入的列标题下横线宽度成比例。

21.3.10.2 简单表

简单表的格式是，第一行是各列标题，第二行是各标题下面用减号组成的表格线，同一行的不同列要用空格分开，从第三行开始是内容。

在表格前或表格后用空行隔开的以 **Table:** 开头的行是表格说明或标题。

为了确定表格每列单元格内容如何对齐，用列标题下的表格线给出提示:

- 表格线与列标题右对齐，表示该列右对齐；
- 表格线与列标题左对齐，表示该列左对齐；
- 列标题在表格线中间，表示该列居中对齐；
- 列标题左右都与表格线对齐，表示该列为缺省对齐方式，一般是左对齐。

一定要使用一个等宽字体来编辑这样的表格，否则对齐与否无法准确分辨。单元格内容不能超出表格线左端。经过试验发现，中文内容很难按这种方法对齐。

例如：

Name	Income	Job	Color
Jane	123456	Research Assistant	red
John	50	N/A	blue
William	3200	Cleaner	blue

Table: 一个简单表的例子

结果为：

表 21.2: 一个简单表的例子

Name	Income	Job	Color
Jane	123456	Research Assistant	red
John	50	N/A	blue
William	3200	Cleaner	blue

21.3.10.3 换行表

换行表在输入列标题和单元格内容时，允许输入内容拆分行，但是转化后并不拆分行。这样的表以一行减号开始，以一行减号结束，中间的表格用空行分开实际的不同行。例如：

```
-----
Name
of
```

```

Subject      Income      Job              color
-----
Jane         123456      Research
Ayer                    Assistant

John         50          N/A              blue
Tukey

William      3200       Cleaner         blue
Tale
-----

```

Table: 一个换行表的例子

结果为:

表 21.3: 一个换行表的例子

Name of Subject	Income	Job	color
Jane Ayer	123456	Research Assistant	red
John Tukey	50	N/A	blue
William Tale	3200	Cleaner	blue

换行表输入时各列的输入宽度是有作用的，输入较宽的列结果也较宽。

21.3.10.4 有格表

完全用减号、竖线、等于号、加号画出表格线。这样的表在文本格式下呈现出很好的表格形状。转化后不能指定对齐方式。

例如:

Table: 有格表示例

```

+-----+-----+-----+
| Fruit      | Price      | Advantages      |
+=====+=====+=====+
| Bananas    | $1.34      | - built-in wrapper |
|            |            | - bright color    |
+-----+-----+-----+
| Oranges    | $2.10      | - cures scurvy   |
|            |            | - tasty           |
+-----+-----+-----+

```

结果为

表 21.4: 有格表示例

Fruit	Price	Advantages
Bananas	\$1.34	<ul style="list-style-type: none"> • built-in wrapper • bright color
Oranges	\$2.10	<ul style="list-style-type: none"> • cures scurvy • tasty

21.4 附录：pandoc 软件介绍

pandoc是一个杰出的开源软件，作者为 John MacFarlane。RStudio 软件中已经包含了这个软件。软件在命令行运行，可以用来在多种不同的文件格式之间进行转换，输入格式一般是文本格式，比如 markdown、LaTeX、MediaWiki、reStructured Text、HTML、DocBook，但是也可以输入 MS Word docx、Open Office ODT、EPUB 格式。输出可以是这些文本格式，也可以是 MS Word docx、Open Office ODT、EPUB、LaTeX 等格式，在安装了 LaTeX 编译系

统如 MikTeX 时可以输出为 PDF。RStudio 软件中包括了一份 pandoc 软件程序。

可以用普通文本编辑器编写 markdown 格式的文件,用 pandoc 转换为 HTML 或 MS Word docx 等格式。因为 pandoc 需要 UTF8 编码的输入文件,所以应该把 markdown 文件保存为 UTF 格式,MS Windows 下的 Notepad++ 软件可以很容易地编辑文本文件并在各种编码之间转换。实际上,MS Windows 下有一个 markdown 编辑程序 Smarks 就是基于 pandoc 软件的。

首先,从 pandoc 网站下载并安装 pandoc。安装程序很奇怪地安装到了 C:\Users\登录用户名\AppData\Local\Pandoc 中,请将此子目录复制到一个合适的位置,比如 C:\Pandoc。如果把此路径加入到 Windows 系统的可执行文件搜索路径中(在“控制面板-系统和安全-系统-高级系统设置-环境变量”中,为系统变量的 Path 添加一个分号分开的路径 C:\Pandoc 即可以不用全路径访问 pandoc.exe 可执行文件。

为了用 pandoc 转化某个 markdown 文件,首先在该文件所在子目录打开一个 Windows 命令行窗口,在 MS Win10 系统中只要在文件管理器的“文件”菜单选择“打开命令提示符”即可。比如,文件名是 test.md,用如下 pandoc 命令可以转换为.docx 文件 (MS Word 文件的新版本):

```
C:\Pandoc\pandoc -o test.docx test.md
```

用如下 pandoc 命令可以转换为.html 文件:

```
C:\Pandoc\pandoc -s -o --mathjax test.html test.md
```

如果把 pandoc.exe 加入了 Windows 操作系统的 Path 环境变量中,上面的 C:\Pandoc\pandoc 可以简写为 pandoc。

pandoc 在其它操作系统中也有相应的版本。软件下载与文档见 Pandoc 的网站<http://pandoc.org>。

如果使用 RStudio 编辑 markdown 文件或者 R Markdown 文件,它会自动调用内置的 pandoc 程序。

Chapter 22

R Markdown 文件格式

22.1 R Markdown 文件

借助于 R 的 `knitr` 和 `rmarkdown` 扩展包的帮助，可以在 Markdown 格式的源文件中插入 R 代码，使得 R 代码的结果能够自动插入到最后生成的研究报告中。这种格式称为 R Markdown 格式，简称为 Rmd 格式，相应的源文件扩展名为 `.Rmd`。输出格式可以是 HTML、docx、pdf、beamer 等。R Markdown 的基础格式是 markdown 格式，严格说来是 Pandoc 软件支持的增强版的 markdown 格式，比如，支持 LaTeX 格式的数学公式，支持各种编程语言语法彩色加亮显示，等等。

`knitr` 的详细文档参见网站 `knitr` 文档。关于 R Markdown 可参考专著 (Xie et al., 2019)。RStudio 网站提供了一个 R Markdown 使用摘要下载: (`rmarkdown-2.0.pdf`)[`rmarkdown-2.0.pdf`]。Pandoc 的文档见 `pandoc` 网站。

一个 Rmd 文件中包含元数据 (metadata)、正文内容和 R 代码三种成分，比如，下面是一个简单的 Rmd 文件样例：

```
---  
title: "R 语言简介"  
author: " 李东风"  
date: "2019-08-29"  
output: html_document
```

这里是正文段落。

段落中仍可以利用一般的 Markdown 语法，
比如在两边用双星号表示 **** 粗体加重 ****，
在用两边反单撇号表示代码，如 `y <- sin(x)`。

下面是一个 R 代码段，有文字输出：

```
```{r}
set.seed(1)
x <- round(rnorm(10), 2)
print(x)
```
```

下面是一个 R 代码段，有图形输出：

```
```{r}
plot(x)
```
```

下面是一个 R 代码段，
有富文本表格输出：

```
```{r}
knitr::kable(as.data.frame(x))
```
```

在文字段落内部也可以有代码，
比如，x 的第一个元素值为 `r x[1]`。

```
```
```

在文件开头用三个减号组成的行包围的内容称为元数据，可以用来规定文章标题、作者、日期、输出格式、输出设置等属性。见22.11。

## 22.2 R Markdown 文件的编译

RStudio 是一个集成的 R 软件环境，可以用来编辑和执行 R 程序，这个软件也可以用来编辑和编译 R Markdown 格式的文件，使得 R Markdown 格式的文件变得容易使用。在 RStudio 中可以直接用一个快捷图标一次性地把 R 代码结果插入内容中并编译为 HTML 或 MS Word docx 格式，还支持 Markdown 中 LaTeX 格式的数学公式。建议使用 RStudio 软件作为 R Markdown 文件的编辑器。

在 RStudio 软件中，用菜单“File-New File-R Markdown”新建一个 R Markdown 文件，扩展名为 .Rmd。用快捷图标 Knit 可以将文件转换成 HTML 格式、PDF 格式（需要安装 LaTeX 编译软件）、MS Word 格式。

从 HTML 格式可以转换成 PDF 格式。为此，安装 Google 的 Chrome 浏览器，在 Chrome 中打开 HTML 文件后，从 Chrome 浏览器的菜单中找到打印菜单，从中选择打印机为“保存到 PDF”选项，就可以将 HTML 网页转换成 PDF，其中的数学公式、表格、图形都可以比较好地转换。但是，经试验，slidy 幻灯片用 Internet Explorer 浏览器转换 PDF 更合适。其它浏览器在转换数学公式时容易将公式右侧的编号裁切掉。

从 Word 文件也可以转换成 PDF 格式，用 MS Word 软件的“文件-导出”或者“文件-另存为”功能即可。

如果想将 R Markdown 文件借助于 LaTeX 格式转换为 PDF，需要在系统中安装一个 TeX 编译器，如 MS Windows 下有 MikTeX 软件包。如果仅在 R Markdown 文件中使用 LaTeX 功能，可以安装谢益辉的 TinyTeX 软件包。

如果不借助于 RStudio 软件，可以用 R 软件、knitr 包、rmarkdown 包、pandoc 软件来完成 R Markdown 源文件的编译。比如，假设 `test.Rmd` 是一个这样的 R Markdown 格式的文件，注意一定要使用 UTF-8 编码保存，可以在 R 或 RStudio 中运行如下命令以生成含有运行结果的 html 文件：

```
rmarkdown::render("myfile.Rmd", output_format = "html_document", encoding="UTF-8")
```

其中 `myfile.Rmd` 是源文件，产生的 HTML 文件带有图形、支持数学公式。

在 R 或 RStudio 中可以用如下命令把 .Rmd 文件转化为 MS Word docx 格式：

```
rmarkdown::render("myfile.Rmd", output_format = "word_document", encoding="UTF-8")
```

使用 RStudio 软件使得这些任务可以一键完成，而且有很好的数学公式支持，所以建议编辑 R Markdown 文件还是使用 RStudio 软件。

用 RStudio 的 Knit 图标一键编译与用 `rmarkdown::render()` 命令编译有一个重要差别：

- 用 Knit 图标编译，Rmd 文件中的程序会在一个崭新的会话中执行，当前会话中已经定义的函数、变量、导入的扩展包不会影响到编译结果；
- 用 `rmarkdown::render()` 编译，Rmd 文件中的程序是在当前会话中执行的，会带来一定的兼容性问题，有可能在别人的环境下就不能正确执行或者会给出不同结果。但是，`rmarkdown::render()` 可以通过程序调用，比如，循环地从同一个 Rmd 生成一系列不同的报告。为了不让当前会话环境干扰结果，可以人为地打开一个新会话。

## 22.3 在 R Markdown 文件中插入 R 代码

插入的 R 代码分为行内代码与代码块。

行内代码的结果插入到一个段落中间，代码以 ``r` 开头，以 ``` 结尾，如 `r sin(pi/2)` 在结果中会显示为 1。为了原样显示一个反向单撇号，可以在两边用双反向单撇号界定并用空格隔开内部的内容。

代码块则把结果当作单独的段落，按照 Markdown 格式的规定，代码块的前后需要有空行，但是 R Markdown 实际上放松了这个要求，允许前后不空行。R 代码段以单独的一行开头，此行以三个反单撇号开始，然后是 `{r}`，如 ````{r}`。代码段以占据单独一行的三个反单撇号 ````` 结尾。如

```
set.seed(1)
x <- round(rnorm(10), 2)
print(x)
```

结果将变成

```
set.seed(1)
x <- round(rnorm(10), 2)
```

```
print(x)
```

```
[1] -0.63 0.18 -0.84 1.60 0.33 -0.82 0.49 0.74 0.58 -0.31
```

可以看出，代码段程序会被插入到最终结果中，代码段的文本型输出会插入到程序的后面。

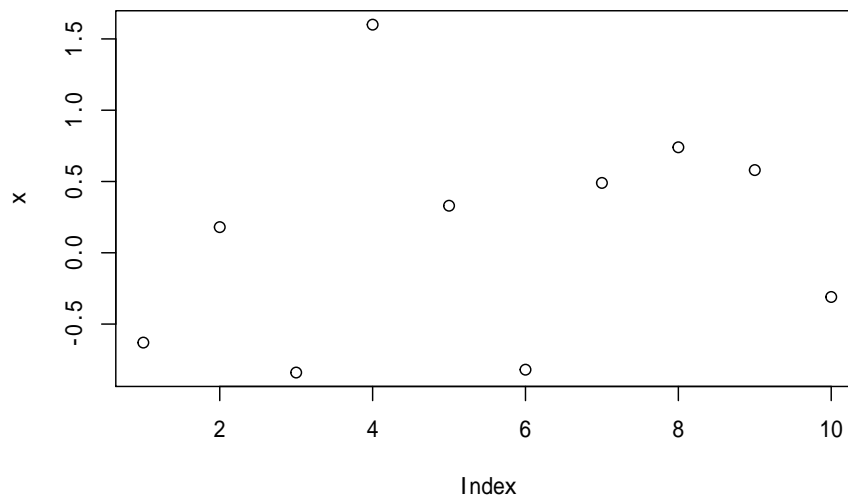
代码块也可以嵌入到引用、列表等环境中。

代码块中作的图将自动插入到当前位置。下面的程序：

```
plot(x)
```

结果将显示为：

```
plot(x)
```



在 RStudio 中，可以用 Insert 快捷图标插入代码段，还可以用 Ctrl+Alt+I 快捷键插入代码段。

## 22.4 输出表格

knitr 包提供了一个 `kable()` 函数可以用来把数据框或矩阵转化成有格式的表格，支持 HTML、docx、LaTeX 等格式。

例如，计算线性回归后，`summary()` 函数的输出中有 `coefficients` 一项，是一个矩阵，如果直接文本显示比较难看：

```
x <- 1:10; y <- x^2; lmr <- lm(y ~ x)
co <- summary(lmr)$coefficients
print(co)
```

```
Estimate Std. Error t value Pr(>|t|)
(Intercept) -22 5.5497748 -3.964125 4.152962e-03
x 11 0.8944272 12.298374 1.777539e-06
```

可以用 knitr 包的 `kable` 函数来显示：

```
knitr::kable(co)
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-22	5.5497748	-3.964125	0.0041530
x	11	0.8944272	12.298374	0.0000018

`kable()` 函数的 `digits=` 选项可以控制小数点后数字位数，`caption=` 选项可以指定表的标题内容。

R 扩展包 `xtable` 提供了一个 `xtable()` 函数，也可以用来生成 HTML 格式和 LaTeX 格式的表格，但是需要指定要输出的格式。`xtable` 对比较多的 R 数据类型和输出类型提供了表格式显示功能，包括矩阵、数据框、回归分析结果、方差分析结果、主成分分析结果、若干分析结果的 `summary` 结果等。例如，上面的回归结果用 `xtable()` 函数显示如：

```
print(xtable::xtable(lmr), type='html')
```

Estimate

Std. Error

t value

```

Pr(>|t|)
(Intercept)
-22.0000
5.5498
-3.96
0.0042
x
11.0000
0.8944
12.30
0.0000

```

这个代码段用了选项 `results='asis'`, 因为 `xtable` 生成的是直接用来插入到结果中的 `html` 代码。注意这里指定了输出为 `HTML` 类型。如果将本文件转化为 `docx`, `xtable` 的结果不可用。

R 扩展包 `pander` 提供了更好的表格能力, 也能与 `knitr` 包很好的合作输出。其 `pander()` 函数可以将多种 R 输出格式转换成 `knitr` 需要的表格形式。如

```
pander::pander(lmr)
```

表 22.1: Fitting linear model:  $y \sim x$

	Estimate	Std. Error	t value	Pr(> t )
<b>(Intercept)</b>	-22	5.55	-3.964	0.004153
<b>x</b>	11	0.8944	12.3	1.778e-06

但是, 经过试验发现, 表中中有中文时 `pander` 包会出错。

## 22.5 利用 R 程序插图

Rmd 文件的插图有两种，一种是已经保存为图形文件的，主要是 png 和 pdf 图片；另一种是文中的 R 代码生成的图形。

已经有图形文件的，可以用 markdown 格式原来的插图方法，见 markdown 格式介绍。但是，这样做不能给图形自动编号，另外因为制作图书是有网页和 PDF 书两种主要输出格式的，原有的插图方式在这两种输出格式上有细微的不一致。所以，最好是统一使用 Rmd 的插图方法。

Rmd 的插图方法就是写一段 R 代码段来插图，如果是用程序作图，则代码中写作图的代码；如果是已有的图形文件，可以在一个单独的 R 代码段中用类似下面的命令插图：

```
```{r, echo=FALSE}
knitr::include_graphics("figs/myfig01.png")
```
```

```
knitr::include_graphics("figs/myfig01.png")
```

其中 `figs` 是存放图形文件的子目录名，`myfig01.png` 是要插入的图形文件名。这样，如果同时还有 `myfig01.pdf` 的话，则 HTML 输出使用 png 图片而 PDF 输出自动选用 pdf 文件。另外，插图的选项在代码段的选项中规定：用代码段的 `fig.width` 和 `fig.height` 选项指定作图的宽和高（英寸），用 `out.width` 和 `out.height` 选项指定在输出中实际显示的宽和高，实际显示的宽和高如果使用如 "90%" 这样的百分数单位则可以自动适应输出的大小。

由于 PDF 中的中文编码不能自动识别，所以在每个 Rmd 源文件的开头应该加上如下的设置，使得生成 PDF 图时中文能够正确显示：

```
```{r setup-pdf, include=FALSE}
pdf.options(family="GB1")
```
```

其中 `include=FALSE` 表示要不显示代码段的代码，有运行结果也不插入到输出结果中，是否运行视缺省的 `eval=` 的值而定。



## 22.6 代码段选项

独立代码段以 ```{r}` 开头，在大括号内还可以写一些选项，选项之间以及与开始的 `r` 之间用逗号分隔，所有选项写在同一行内不要换行。选项都使用“选项名 = 选项值”的格式，选项值除了使用常量外也可以使用全局变量名或表达式。在大括号内开头的 `r` 空格后写一个由英文大小写字母、数字、减号组成的标识符，作为代码段的标签。标签中不要用其它类型的字符，下划线也不要。如

```
``{r firstCode}
cat('This is 第一段, 有标签.\n')
``
```

关于代码段选项，详见<https://yihui.name/knitr/options>。

### 22.6.1 代码和文本输出结果格式

R 代码块和 R 代码块的运行结果通常是代码块原样输出，运行结果用井号保护起来，这样有利于从文章中复制粘贴代码。如：

```
``{r}
s <- 0
for(x in 1:5) s <- s + x^x
s
``
```

结果为:

```
s <- 0
for(x in 1:5) s <- s + x^x
s
```

```
[1] 3413
```

### 22.6.1.1 highlight 选项

转化后的 R 代码块缺省显示为彩色加亮形式。用选项 `highlight=FALSE` 关闭彩色加亮功能。

### 22.6.1.2 prompt 和 comment 选项

如果希望代码用 R 的大于号提示符开始，用选项 `prompt=TRUE`。如果希望结果不用井号保护，使用选项 `comment=''`。例如：

```
```{r prompt=TRUE, comment=''}  
sum(1:5)  
```
```

结果为：

```
> sum(1:5)
```

```
[1] 15
```

### 22.6.1.3 echo 选项

如果希望不显示代码，加选项 `echo=FALSE`。如

```
```{r echo=FALSE}  
print(1:5)  
```
```

结果为：

```
[1] 1 2 3 4 5
```

### 22.6.1.4 tidy 选项

加选项 `tidy=TRUE` 可以自动重新排列代码段，使得代码段格式更符合规范。例如：

```
```{r tidy=TRUE}
```

```
s <- 0
for(x in 1:5) {s <- s + x^x; print(s)}
...
```

结果为:

```
s <- 0
for (x in 1:5) {
  s <- s + x^x
  print(s)
}
```

```
## [1] 1
## [1] 5
## [1] 32
## [1] 288
## [1] 3413
```

22.6.1.5 eval 选项和 include 选项

加选项 `eval=FALSE`, 可以使得代码仅显示而不实际运行。这样的代码段如果有标签, 可以在后续代码段中被引用。

加选项 `include=FALSE`, 则本代码段仅运行, 但是代码和结果都不写入到生成的文档中。

22.6.1.6 child 选项

加选项 `child=' 文件名.Rmd'` 可以调入另一个.Rmd 文件的内容。如果有多个.Rmd 文件依赖于相同的代码, 可以用这样的方法。

22.6.1.7 collapse 选项

一个代码块的代码、输出通常被分解为多个原样文本块中, 如果一个代码块希望所有的代码、输出都写到同一个原样文本块中, 加选项 `collapse=TRUE`。例如, 没有这个选项时:

```
```{r}
sin(pi/2)
cos(pi/2)
```
```

结果为:

```
sin(pi/2)
```

```
## [1] 1
```

```
cos(pi/2)
```

```
## [1] 6.123032e-17
```

代码和结果被分成了 4 个原样文本块。加上 `collapse=TRUE` 后, 结果为:

```
sin(pi/2)
```

```
## [1] 1
```

```
cos(pi/2)
```

```
## [1] 6.123032e-17
```

代码和结果都在一个原样文本块中。

22.6.1.8 results 选项

用选项 `results=` 选择文本型结果的类型。取值有:

- `markup`, 这是缺省选项, 会把文本型结果变成 HTML 的原样文本格式。
- `hide`, 运行了代码后不显示运行结果。
- `hold`, 一个代码块所有的代码都显示完, 才显示所有的结果。
- `asis`, 文本型输出直接进入到了 HTML 文件中, 这需要 R 代码直接生成 HTML 标签, knitr 包的 `kable()` 函数可以把数据框转换为 HTML 代码的表格。

例如: `results='hold'` 的示例:

```
```{r collapse=TRUE, results='hold'}
sin(pi/2)
cos(pi/2)
```

```
...
```

结果为:

```
sin(pi/2)
cos(pi/2)
[1] 1
[1] 6.123032e-17
```

### 22.6.1.9 错误信息选项

选项 `warning=FALSE` 使得代码段的警告信息不进入编译结果，而是在控制台 (console) 中显示。有一些扩展包的载入警告可以用这种办法屏蔽。

选项 `error=FALSE` 可以使得错误信息不进入编译结果，而是出错停止并将错误信息在控制台中显示。

选项 `message=FALSE` 可以使得 `message` 级别的信息不进入编译结果，而是在控制台中显示。

## 22.6.2 图形选项

### 22.6.2.1 图形大小

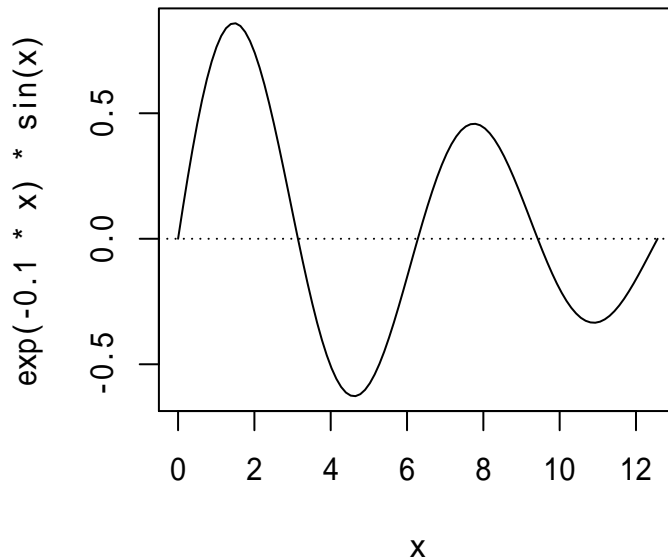
用 `fig.width=` 指定生成的图形的宽度，用 `fig.height=` 指定生成的图形的高度，单位是英寸（1 英寸等于 2.54 厘米）。

下面给出一个长宽都是 10 厘米的图例。

```
``{r fig.width=10/2.54, fig.height=10/2.54}
curve(exp(-0.1*x)*sin(x), 0, 4*pi)
abline(h=0, lty=3)
...`
```

结果为:

```
curve(exp(-0.1*x)*sin(x), 0, 4*pi)
abline(h=0, lty=3)
```

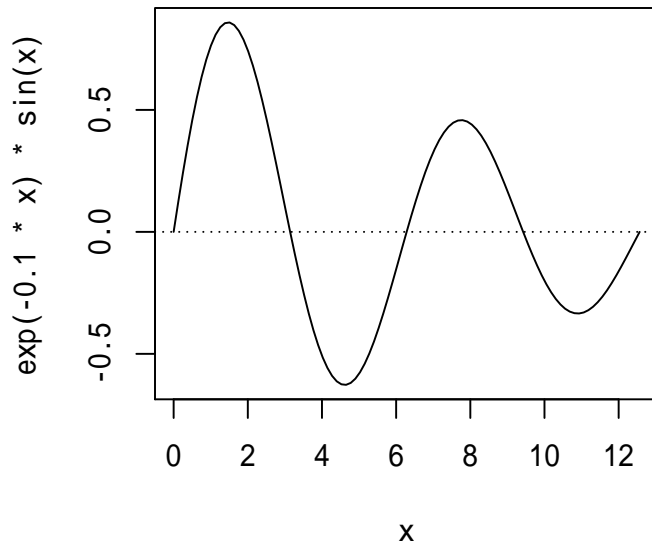


`fig.width=` 和 `fig.height=` 规定的是生成的图形大小，实际生成图形的显示大小会受到 `dpi`（分辨率）影响，默认 `dpi` 是 72（每英寸 72 个点），也可以用 `dpi=` 选择分辨率。转化后的 HTML 文件显示时不一定按原始大小显示。用 `out.width=` 和 `out.height=` 可以指定显示大小，但是必须是最终文件格式承认的单位，比如 HTML 的图形大小单位是点（平常说屏幕分辨率的单位）。或者写成如 90% 这样的百分比格式。例如在上面的例子中加上输出度为页面宽度 80% 的选项：

```
```{r fig.width=10/2.54, fig.height=10/2.54, out.width="80%"}
curve(exp(-0.1*x)*sin(x), 0, 4*pi)
abline(h=0, lty=3)
```
```

注意所有选项都要写在一行中，不能换行。结果为：

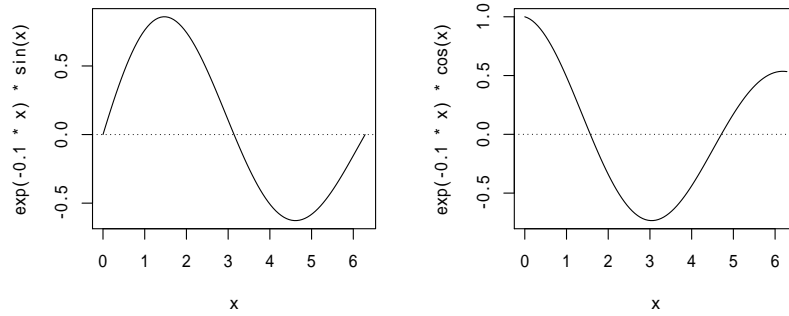
```
curve(exp(-0.1*x)*sin(x), 0, 4*pi)
abline(h=0, lty=3)
```



当输出结果为 HTML 结果或者 LaTeX 转换的 PDF 时，只要两个图形的总宽度小于页面宽度，两个连续的图形就可以左右并列放置。Word 不支持这种功能。例如：

```
```{r, echo=FALSE, fig.width=10/2.54, fig.height=10/2.54, out.width="45%"}  
curve(exp(-0.1*x)*sin(x), 0, 2*pi)  
abline(h=0, lty=3)  
curve(exp(-0.1*x)*cos(x), 0, 2*pi)  
abline(h=0, lty=3)  
```
```

结果为

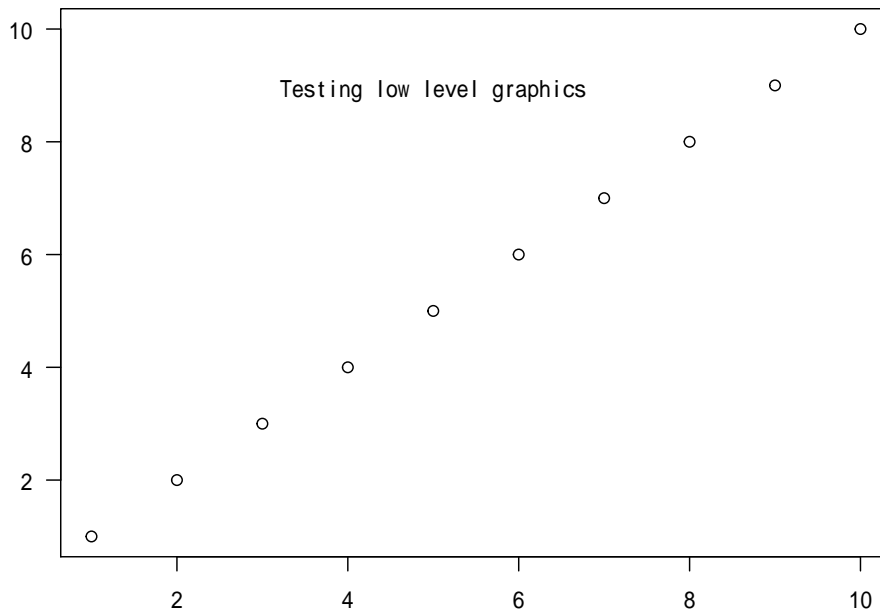


### 22.6.2.2 图形结果选择

用 `fig.keep=` 选项可以选择保留哪些 R 代码生成的图。缺省是 `fig.keep='high'`，即保留每个高级图形函数的结果图形，低级图形函数对高级图形函数的更改不单独保存而汇总到高级图形函数结果中。如

```
par(mar = c(3, 3, 0.1, 0.1))
plot(1:10, ann = FALSE, las = 1)
text(5, 9, "Testing low level graphics")
```





其中 `text()` 函数的结果与高级图形函数 `plot()` 的结果一起显示。

`fig.keep` 还可以取: `all`, 会把低级图形函数修改后的结果单独保存; `last`, 仅保留最后一个图形; `first`, 仅保留第一个图; `none`, 所有图都不显示出来。

### 22.6.3 缓存 (cache) 选项

当 R Markdown 文章比较长, 包含的 R 代码比较多, 或者代码段运行需要比较长时间时, 反复编译整篇文章会造成不必要的计算, 因为有些代码段并没有修改, 依赖的数据也没有改变。knitr 提供了缓存功能, 代码段选项 `cache=TRUE` 对代码段打开缓存, 允许暂存上次运行的结果 (包括文本结果和图形) 而不需要重复运行代码段。当代码段被修改时, 缓存被放弃, 编译时重新运行代码段。

缓存这种功能需要慎重使用, 免得错误地使用了旧的结果。当后面的代码段需要使用前面代码段结果时, 如果前面结果改了, 后面的代码段就不能使用缓存的结果而必须重新计算。为此, 在后面的代码段中应该加上 `dependson=` 选项, 比如 `dependson=c('codeA', 'codeB')`, 其中 `codeA` 和 `codeB` 是前面的缓存了的代码段的标签, 其结果会用在本地代码段中。也可以使用代码段选项 `autodep=TRUE`, knitr 试图自动确定前后代码段之间的依赖关系, 每当前面的

代码段改变时，后面的用到其结果代码段也自动重新计算而不使用缓存的旧结果。

建议仅对计算一次需要较长时间的代码段使用缓存功能，后面依赖于其结果的代码一定要加上 `dependson=` 选项。建议代码段尽可能有标签，这样在编译失败时能马上看出失败的地点。

为了保险起见，可以尽量不使用缓存功能，而是用 `save()` 功能将需要长时间结算的结果保存下来，用 `load()` 载入然后输出到结果文档中。

## 22.7 章节目录链接问题

对于英文文件，R Markdown 基于的 `pandoc` 软件可以自动从标题生成适当的链接标签，将 HTML 输出或者 PDF 输出设置属性 `toc: true` 后可以生成可点击的章节目录。但是，目前 `Pandoc` 对于中文文件的支持差一些，所以中文文件要自动生成可点击的目录，需要在每个标题行的末尾，空格后添加 `{#label}`，其中 `label` 是自己指定的标签内容，但是建议仅使用英文字母、数字、减号。如

```
第三章第一节标题 {#c3-s1-int}
```

如果希望某个标题不参与自动编号，可以在标题末尾空格之后加上 `{-}` 标记。

## 22.8 数学公式

### 22.8.1 在 Markdown 中输入数学公式

原始的 Markdown 格式并不支持数学公式。`Pandoc` 扩展的 `markdown` 格式提供了对数学公式的支持，可以在 Markdown 文件中插入 `LaTeX` 格式的数学公式。虽然不能提供所有的 `LaTeX` 公式能力，但是常用的数学公式还是能做得很好，转换到 HTML、docx 都可以得到正常显示的公式。

用 `RStudio` 软件编译 Markdown 文件，可以在其中插入 `LaTeX` 格式的数学公式，编译成 HTML 或者 docx 格式后都可以正常显示数学公式，在另外安装

的 LaTeX 编译器的支持下也可以将.Rmd 格式编译 LaTeX 格式然后再转换为 PDF 格式，这种方法对数学公式的支持会更完善。

关于数学公式的软件设置参见下面设置部分的内容。

### 22.8.2 数学公式类别

数学公式分为行内公式和独立公式。行内公式和段落文字混排，写在两个美元符号  $\$$  中间，或者  $\backslash$ (和 $\backslash$ ) 之间。例如  $\$f(x)=\frac{1}{2} \int_0^1 \sin^2(tx) dt\$$  变成  $f(x) = \frac{1}{2} \int_0^1 \sin^2(tx) dt$ 。开头的  $\$$  后面不能紧跟着空格，结尾的  $\$$  不能紧跟在空格后面。

独立公式写在成对的美分符号中间，或者  $\backslash[$ 和 $\backslash]$  之间。例如：

```


$$f(x) = \frac{1}{2} \sum_{j=1}^{\infty} \int_0^1 \sin^2(jtx) dt .$$


```

显示为

$$f(x) = \frac{1}{2} \sum_{j=1}^{\infty} \int_0^1 \sin^2(jtx) dt.$$

### 22.8.3 基本功能

在数学公式中，用下划线表示下标，比如  $\$x_1\$$  结果为  $x_1$ 。用  $\wedge$  表示上标，如  $\$x^2\$$  结果为  $x^2$ 。上下标都有，如  $\$x_1^2\$$  结果为  $x_1^2$ 。

公式中用大括号表示一个整体，比如  $\$x^{(1)}\$$  不能得到  $x^{(1)}$ ，需要用  $\$x^{\{(1)\}}\$$ 。

公式中用反斜杠开始一个命令，命令仅包含字母而不能包含数字，数字只能作为参数。如  $\$\frac{1}{2}\$$  和  $\$\frac{1}{2}\$$  都可以生成  $\frac{1}{2}$ 。类似的命令如  $\$\sqrt{2}\$$  为  $\sqrt{2}$ 。

希腊字母都有对应的命令，如  $\$\alpha\$$  为  $\alpha$ ， $\$\Sigma\$$  为  $\Sigma$ 。

常用数学函数有自己的命令，如  $\$\sin x\$$  变成  $\sin x$ ， $\$\exp(x)\$$  为  $\exp(x)$ 。

圆括号与方括号可以直接使用, 绝对值符号用 `|`, 如 `$|x|` 变成  $|x|$ 。但是, 大括号必须写成 `\{` 和 `\}` 的形式, 如 `$\exp\{-\frac{1}{2}x^2\}` 变成  $\exp\{-\frac{1}{2}x^2\}$ 。

为了使得括号能够与括号内部的内容等高, 可以用 `\left` 和 `\right` 命令进行修饰, 如

```


$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2\right\}$$


```

变成

$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2\right\}$$

注意 `\left` 与 `\right` 必须成对使用, 否则出错。

求和如 `$\sum_{i=1}^n x_i` 变成  $\sum_{i=1}^n x_i$ 。乘积如 `$\prod_{i=1}^n x_i` 变成  $\prod_{i=1}^n x_i$ 。积分如 `$\int_0^1 f(x) dx` 变成  $\int_0^1 f(x)dx$ 。

#### 22.8.4 修饰符

`$f'(x)` 显示为  $f'(x)$ , 表示导数; `$f''(x)` 显示为  $f''(x)$ , 表示二阶导数。顺便说一句, 偏导数写法如 `$\frac{\partial f(x,t)}{\partial x}`, 显示为  $\frac{\partial f(x,t)}{\partial x}$ 。

$\bar{x}$  的写法是 `$\bar{x}`。 $\overline{\text{span}}$  的写法是 `$\overline{\text{span}}`。 $\hat{x}$  的写法是 `$\hat{x}`。 $\tilde{x}$  的写法是 `$\tilde{x}`。 $\vec{x}$  的写法是 `$\vec{x}`。

#### 22.8.5 对齐与矩阵

为了产生对齐的公式, 在独立公式中使用 `aligned` 环境。公式中的环境以 `\begin{环境名}` 开始, 以 `\end{环境名}` 结束, 用 `\` 表示换行, 用 `&` 表示一个上下对齐位置。如

```


$$f(x) = \sum_{k=0}^{\infty} \frac{1}{k!} x^k$$

$$= e^x$$


```

`\end{aligned}`  
`$$`

转化成

$$f(x) = \sum_{k=0}^{\infty} \frac{1}{k!} x^k \\ = e^x$$

可以用 `pmatrix` 环境制作写在圆括号中的矩阵，用 `bmatrix` 环境制作写在方括号中的矩阵，用 `vmatrix` 制作写在绝对值号中的矩阵，如

$$\begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix}$$

列向量、行向量也可以用这种办法制作。

### 22.8.6 特殊字体

数学公式中的单词、文字必须用 `\text{...}` 保护，比如

`$$`  
`\text{CV} = \frac{S}{\bar{x}} \times 100 \%`  
`$$`

变成

$$CV = \frac{S}{\bar{x}} \times 100\%$$

上例中如果不用 `\text{}` 保护，就会显示为

$$CV = \frac{S}{\bar{x}} \times 100\%$$

这里的  $CV$  通常理解为  $C$  乘以  $V$ 。

有时用粗体表示向量或者矩阵，用 `\boldsymbol{...}` 说明。如

`$$`  
`\boldsymbol{v} = (v_1, v_2)^T`  
`$$`

变成

$$\mathbf{v} = (v_1, v_2)^T$$

美术体英文字母用`\mathcal{...}`，如  $\mathcal{A}, \mathcal{B}, \mathcal{C}$  写法为 `\mathcal{A}, \mathcal{B}, \mathcal{C}`。

手写花体字母用`\mathscr{...}`，如  $\mathscr{B}, \mathscr{C}, \mathscr{F}, \mathscr{G}$  写法为 `\mathscr{B}, \mathscr{C}, \mathscr{F}, \mathscr{G}`。

## 22.9 其它编程语言引擎

Rmd 文件中除了 R 代码段以外，还可以插入 Rcpp、Python、Julia、SQL 等许多编程语言的代码段，常用编程语言还可以与 R 代码段进行信息交换。

## 22.10 交互内容

在 Rmd 生成 HTML 结果时，可以在结果中包含允许读者交互操作的内容，比如，用户修改模型参数，使得模型报表、图形交互地改变。有两种方法可以实现交互：

- HTML 小程序 (widgets)，由 R 扩展包 `htmlwidgets` 实现，利用 JavaScript 函数库进行交互。在这一功能基础上有一些派生的扩展包，如 `DT`，`leaflet`，`dygraph`。
- 利用 `shiny` 框架，这也是一个 R 扩展包，交互功能由 R 的一个会话驱动。Rmd 文件的 YAML 元数据中需要设置 `runtime: shiny`。

## 22.11 属性设置

### 22.11.1 YAML 元数据

可以在 RStudio 软件中编辑 markdown 文件与 Rmarkdown 文件，RStudio 的 markdown 支持来自 `knitr`、`rmarkdown`、`bookdown` 扩展包和 `pandoc` 软件，.Rmd 文件支持一些特殊的文件设置，这些设置写在 markdown 文件和.Rmd

文件的开头位置，用三个减号组成的行作为开始标记与结束标记，称为 YAML 元数据块 (YAML meta data block)，块后面必须用空行分隔。元数据块中可以用“设置属性名: 设置属性值”的办法设置属性，用缩进表示属性内容，用上下对齐的行表示多项列表。常用属性有 `title`(标题)、`author`(作者)、`date`(日期)、`output_format`(输出格式) 等，如

```

title: " 临时测试"
author: " 李东风"
date: "2016 年 7 月 6 日"

```

这三个设置会出现在转换后结果的标题部分。RStudio 的中文支持还是有时出现问题，如果出现涉及到 YAML 的错误，先将中文内容替换为英文试一试。

因为冒号：在属性设置中有特殊意义，属性设置值如果含有冒号，需要把整个属性值两边用单撇号界定。

属性值可以是列表或多项，例如：

```

title: 'This is the title: it contains a colon'
author:
 - name: Author One
 affiliation: University of Somewhere
 - name: Author Two
 affiliation: University of Nowhere
tags: [nothing, nothingness]
abstract: |
 This is the abstract.

 It consists of two paragraphs.

```

此例中有 `title`、`author`、`tags`、`abstract` 四个属性。`author` 属性包含两个作者，每个作者又有 `name` 和 `affiliation` 两个属性。`tags` 属性是一个两项的列表。`abstract` 属性用管道符号 `|` 表示开头，不需要结束标志，内容缩进。

### 22.11.2 输出格式

.Rmd 文件开头的 YAML 元数据中，属性 `output` 选择输出的格式，如

- `html_document` 是 HTML 输出;
- `pdf_document` 是 PDF 输出，通过系统中另外安装 LaTeX 编译系统转换;
- `word_document` 是 docx 格式的 Word 文件输出，等等。

输出格式指定如：

```
output:
 html_document: default
 word_document: default
 pdf_document: default
```

其中 `default` 表示该输出格式完全使用默认设置。在输出格式没有定制设置时必须要有 `default` 指定。

### 22.11.3 输出格式设置

在每种输出格式后面还可以继续添加该输出特有属性。如：

```
output:
 html_document:
 toc: true
 number_sections: true
 word_document:
 toc: true
 pdf_document:
 toc: true
```

### 22.11.4 目录设置

输出格式的 `toc: true` 选项指定自动生成目录，HTML 输出、PDF 输出、Word 输出都支持此功能。见上例。



```
output:
 html_document:
 toc: true
 toc_depth: 2
```

上例中, `html_document` 是 `output` 的属性, 隶属关系用缩进表示; `html_document` 又有自己的属性 `toc` 和 `toc_depth`, 隶属关系用缩进表示。属性 `toc: true` 表示要自动生成可点击的目录。输出格式的属性 `toc_depth` 是目录包含的章节层级数。

对 `html_document` 输出, 属性 `toc_float: true` 使得生成的文档在左侧显示一个目录导航窗格。而 `toc_float` 属性又可以指定一些属性, 隶属关系用缩进表示。`toc_float` 的属性 `collapse: true` 表示仅展开章节第一层级, `smooth_scroll: true` 使得通过导航窗格跳跃时页面有平滑滚动过程显示。用如:

```
output:
 html_document:
 toc: true
 toc_float:
 collapse: true
 smooth_scroll: true
```

### 22.11.5 章节自动编号

输出格式的属性 `number_sections` 为 `true` 可以自动对章节编号。HTML 输出与 PDF 输出都支持此功能。

在章节自动编号时, 如果不是利用 `bookdown` 包而是基于基本的 R Markdown, 文中的分节最高层级最好用一个井号, 对应于节, 如果用两个井号, 编号会变成 0.1, 0.2 这样。

### 22.11.6 Word 输出章节自动编号及模板功能

Word 文件不能自动编号，解决办法是生成适当的模板文件，如 `wordstyle.docx`，然后在 Rmd 文件的元数据部分为 `word_document` 输出增加 `reference_docx` 选项，如：

```
output:
 html_document:
 toc: true
 number_sections: true
 word_document:
 toc: true
 reference_docx: wordstyle.docx
```

这样，新编译的 Rmd 文件就可以采用 `wordstyle.docx` 的格式设置，比如章节自动编号等。目前目录标题还是英文，需要在结果中人为地修改为中文。一个样例模板下载如下：`wordstyle.docx`

### 22.11.7 HTML 特有输出格式设置

对于 `html_document` 输出类型，可以用 `theme` 属性设置一个主题，取值如 `default`, `cerulean`, `journal`, `flatly`, `darkly`, `readable`, `spacelab`, `united`, `cosmo`, `lumen`, `paper`, `sandstone`, `simplex`, `yeti`。

用 `highlight` 属性设置程序语言语法高亮样式，可取值有 `default`, `tango`, `pygments`, `kate`, `monochrome`, `espresso`, `zenburn`, `haddock`, and `textmate`。用 `null` 表示取消语法高亮。

可以用 `css` 属性指定一个自定义的 CSS 样式表文件。如果希望完全用自己的样式代替原有样式，可以设置 `theme: null`。

`code_folding` 属性设置 HTML 结果中代码折叠选项，取值 `hide` 可以隐藏代码，但读者通过选项可见。取值 `show` 可以显示代码，但读者通过选项可隐藏代码显示。

### 22.11.8 关于数学公式支持的设置

Rmd 格式支持数学公式，在生成的 HTML、Word、PDF 中都可以正常显示数学公式。

以 HTML 为输出时，会使用 MathJax 库显示数学公式。这是一个用于在浏览器中显示数学公式的 Javascript 程序库，显示效果很好，还支持多种显示实现方式，支持包括 LaTeX 在内的多种数学公式输入方法。

MathJax 库很大，所以一般是从其网站按需远程调用的，但是远程调用 MathJax 库在网络不畅通时会使得公式显示极为缓慢甚至无法显示，所以 Rmd 允许将 MathJax 本地化。在.Rmd 文件头部 YAML 元数据的某种 HTML 输出格式的属性中，设置属性 `mathjax: local` 和就可以使得 MathJax 库被放在生成的 HTML 的一个下层目录中。设置如

```
output:
 html_document:
 toc: true
 self_contained: false
 mathjax: local
```

注意设置 `mathjax: local` 时必须设置 `self_contained: false`。在不使用本地 MathJax 副本时，可以设置 `self_contained: true`，这使得图形文件、JavaScript 代码等依赖项也打包在生成的单一 HTML 文件中，如果设为 `false`，这些依赖文件会存放在单独文件中。可以用 `lib_dir` 属性指定一个子目录用来存放这些依赖文件，多个 Rmd 文件可以共用同一个 `lib_dir` 值。

将 MathJax 用 `mathjax: local` 属性本地化有明显的缺点。这时，每个 Rmd 项目都需要一个 MathJax 库副本，使得文件备份变得很麻烦。下面的设置方法可以令多个 Rmd 项目共享一个共同的本地 MathJax 库。

### 22.11.9 输出设置文件

如果一个项目包含多个.Rmd 文件，每个文件单独用 YAML 元数据进行设置比较繁琐，修改时也很麻烦。可以在项目目录中增加一个 `_output.yml` 文件，其中包含所有.Rmd 文件共用的 `output` 属性的设置，各个.Rmd 文件还可以有

自己单独的 `output` 属性，`.Rmd` 文件 YAML 元数据中的属性优先级高于共同设置的属性。

例如，`_output.yml` 文件内容如下：

```
html_document:
 toc: yes
 number_sections: yes
 mathjax: "../../MathJax/MathJax.js"
 includes:
 in_header: "_header.html"
word_document:
 toc: yes
 reference_docx: wordstyle.docx
pdf_document:
 includes:
 in_header: preamble.tex
 latex_engine: xelatex
```

上面设置了使用本地硬盘的 MathJax 库，位于项目目录的上两层目录中的 MathJax 子目录中。这样可以使得多个项目共享一个 MathJax 库。将生成的结果当作网站发布时，只要将 MathJax 库也安装到项目在网站的目录的上两层的 MathJax 子目录中就可以。其中 `_header.html` 内容如下，是关于 MathJax 具体的一些设置：

```
<script type="text/x-mathjax-config">
MathJax.Hub.Config({
 jax: ["input/TeX","output/SVG"],
 extensions: ["tex2jax.js","MathMenu.js","MathZoom.js"],
 TeX: {
 extensions: ["AMSmath.js","AMSsymbols.js","noErrors.js","noUndefined.js"]
 }
});
</script>
```

注意，`_output.yml` 文件都是输出设置，所以不能再有 `output` 属性，所以下

面的 `_output.yml` 写法是不正确的:

```
output:
 html_document:
 toc: yes
```

## 22.12 LaTeX 和 PDF 输出

Rmd 文件的 `pdf_document` 输出类型是先转换为 LaTeX 文件,再借助计算机系统中另外安装的 LaTeX 编译程序转换为 PDF。需要在系统中安装一个 TeX 编译器,MS Windows 下有 MikTeX 软件包或者 CTEX 软件包。如果仅在 R Markdown 文件中使用 LaTeX 功能,可以安装谢益辉的 TinyTeX 软件包,优点是直接用 R 命令就可以安装,更新也由 R 自动进行,不需要用户干预。

许多 LaTeX 相关的 YAML 元数据不是设置在输出属性中,而是设置为文件的顶级元数据。如

```

title: "R 语言简介"
author: "李东风"
date: "2019-08-29"
output: pdf_document
fontsize: "11pt"
geometry: "margin=1in"

```

其中属性 `fontsize` 和 `geometry` 都是 LaTeX 选项,但是没有作为 `output--pdf_document` 的属性,而是作为 Rmd 文件的顶级元数据属性。

用属性 `documentclass` 指定 LaTeX 文档类型,比如 `article` 是论文格式, `book` 是图书格式。 `ctexart` 和 `ctexbook` 是适用于中文的格式。

`fontsize` 是正文字体大小,可取 `10pt`, `11pt`, `12pt`, 其中 `10pt` 是缺省大小。

`geometry` 用来指定 `geometry` 包的参数。

默认的 LaTeX 编译引擎是 pdfLaTeX。可以在 `pdf_document` 输出的属性中设置 `latex_engine: xelatex` 将引擎替换成 xeLaTeX。xeLaTeX 的优点是

支持 UTF8 编码，可以使用系统中现有字体，而我们在写中文的 R Markdown 文档时必须使用 UTF-8 编码。

Rmd 文件中的文献引用可以由 pandoc 软件执行，但是在用 LaTeX 转换 PDF 时，可以指定用一个 LaTeX 引用包，如：

```

output:
 pdf_document:
 latex_engine: xelatex
 citation_package: natbib

```

为了自动生成目录并自动为章节编号，仍在 `pdf_document` 的属性中指定 `toc: true` 和 `number_section: true`。

在使用 bookdown 管理 Rmd 文件时，文献引用、目录、章节编号方式有另外的指定办法以及默认规则。

为了能够检查转换过程中的错误，可以用 `pdf_document` 的属性 `keep_tex: true`，保留作为中间结果的 .tex 文件，人工检查其中的错误。

可以在 `pdf_document` 输出的属性中指定若干个 .tex 文件插入到 LaTeX 文档的导言、正文最前面、正文最后后面，如

```

output:
 pdf_document:
 latex_engine: xelatex
 keep_tex: true
 citation_package: natbib
 toc: true
 number_sections: true
 includes:
 in_header: preamble.tex
 before_body: doc-prefix.tex
 after_body: doc-suffix.tex

```

`preamble.tex` 主要用来载入额外的 LaTeX 包以及定义新的命令，一个 `preamble.tex` 的例子如下

```
\usepackage{ctex} % 支持中文的标点和章节模式

%\usepackage{xltextra} % XeLaTeX的一些额外符号
% 设置中文字体
\setCJKmainfont[BoldFont={黑体},ItalicFont={楷体}]{新宋体}

\usepackage{amsthm,mathrsfs} % 基本的ams数学公式
\usepackage{booktabs} % 增强的表格功能
\usepackage{longtable} % 支持超过一页的表格

% 下面是自定义的定理环境排版方法
\makeatletter
\def\thm@space@setup{%
 \thm@preskip=8pt plus 2pt minus 4pt
 \thm@postskip=\thm@preskip
}
\makeatother
```

## 22.13 生成期刊文章

节22.12的办法可以用来生成 LaTeX 然后转换为 PDF，可以用来撰写学位论文和期刊论文。但是学位论文和期刊论文都有比较严格的格式要求，用节22.12的一般方法需要对 LaTeX 比较熟悉才能达到格式要求。

用 R 扩展包 `rticles` 可以比较容易地将 Rmd 文件转换成期刊可接受的 LaTeX 和 PDF 格式。见<https://bookdown.org/yihui/rmarkdown/journals.html>。扩展包中对一些常见期刊提供了模板。

首先要安装 `rticles` 扩展包，命令为

```
install.packages("rticles")
```

为了生成新的论文源文件，在 RStudio 中，用“File – New File – R Markdown – from Template”菜单，从列出的文章模板中选一个。

杂志投稿的 Rmd 模板是用许多元数据来规定格式的。用户可以修改其中的标题、作者、作者单位等信息。

因为 rticles 的输出只有 PDF，所以必要时可以在 Rmd 源文件中直接使用 LaTeX 命令，但如果 Rmd 本身可以完成时尽可能用 Rmd 的功能。

因为定理、图表自动标号和索引功能是由 bookdown 包提供的，而 rticles 的输出类型实际是 `rmarkdown::pdf_document`，所以可以联合使用 bookdown 与 rticles，如

```
output:
 bookdown::pdf_book
 base_format: rticles::peerj_article
```

其中 `peerj_article` 可以替换成 rticles 提供的其它格式。

如果自己需要的模板 rticles 没有提供，可以自己设法定义模板，参见<https://bookdown.org/yihui/rmarkdown/document-templates.html>。

## 22.14 附录：经验与问题

### 22.14.1 Word 模板制作

如果有合适的 Word 格式样例，可以在 YAML 元数据的 `word_document` 的属性中，设置 `reference_docs` 属性指向该模板文件。

为了制作 Word 模板，不应该从头开始，而是在没有模板的情况下从 Rmd 文件转换生成一个 Word 文件，将其作为模板文件的基础进行修改。在修改过程中，应该制作过程的备份，以便修改错误时能退回上一个版本。

如果文章的大标题样式不是“标题”，将样式改为“标题”以免在后续的自动编号过程中参与编号。修改的方法是选中大标题，然后在 Word 软件主菜单的“开始-样式”选择框中选中“标题”，就可以设置其样式。

Word 的目录是与“标题 1”同级的，这样会在自动编号时参与编号；为此，选中



“Contents”，然后在主菜单的“开始-样式”选择框中，找到突出显示的“TOC 标题”，右键单击“修改”，在弹出的对话框中将“样式基准”从“标题 1”改为“标题”，并在同一对话框中将适当调整样式，如字体类型、大小、颜色、是否居中。改完后保存并保存一个副本。

如果需要修改“标题”，“作者”，“日期”，“标题 1”等等的字体、大小、颜色、对齐方式，可以选中该部分，并在“样式”中邮件单击相应的样式，选择“修改”进行修改，保存样式文件并试验 Rmd 文件是否正确地利用了修改的样式。

为了实现自动编号，先在正文中选中一个一级标题（不能是目录标题），从 Word 的主菜单栏找到“开始——段落——多级列表——定义新的多级列表”，在打开的对话框中点击“更多”选项，逐一地选择“单击要修改的级别”，将级别 1、2、3 等分别修改，主要是修改“将级别链接到样式”使得“级别 1”与“标题 1”对应，“级别 2”与“标题 2”对应，等等。注意这个修改要一次性完成，不能保存了文件后再次定义新的多级列表。

修改模板文件的页边距也会使得利用其为模板的 Rmd 文件转换结果被修改。为了修改模板文件的页边距，在主菜单“布局——页边距——自定义页边距”弹出的对话框中修改。

生成的目录的标题内容是“Table of Contents”，目前只能对 Rmd 输出的 docx 文件直接修改，修改模板文件的标题内容并不能使得 Rmd 输出被修改。

参考：Richard Layton 的网站文章 [https://rmarkdown.rstudio.com/articles\\_docx.html](https://rmarkdown.rstudio.com/articles_docx.html)。

### 22.14.2 数学公式设置补充

如果不使用 RStudio，R 扩展包 rmarkdown 的 `markdownToHTML()` 函数可以把含有数学公式的内容转换成可显示公式的 HTML 文件，R 扩展包 knitr 的 `knit2html()` 也可以实现此功能。用 rmarkdown 扩展包的 `markdownToHTML()` 函数生成的含有数学公式的内容。

单独使用 pandoc 软件时，也可以用普通文本编辑器在 Markdown 文件中输入 LaTeX 格式的数学公式，然后用 pandoc 软件转化为带有数学公式的 docx 文件，不需要额外选项。但是，为了能够在转换的 HTML 文件中正常显示数学公式，需要在运行 Pandoc 时增加运行选项 `--maxjax`，如

```
pandoc --mathjax -s -o test.html test.md
```

经试验，这样转化的 HTML 文件可以在 Firefox 以及 Microsoft 的 Edge 浏览器和 Internet Explorer 浏览器中正常显示数学公式。如果公式中的中文显示不正常，对显示的公式右键单击弹出选项菜单，选择“Math Settings–Math Renderer–SVG”或者“HTML–CSS”。

## Chapter 23

# 用 bookdown 制作图书

### 23.1 介绍

R 的 bookdown 扩展包 (<https://github.com/rstudio/bookdown>) 是继 knitr 和 rmarkdown 扩展包之后, 另一个增强 markdown 格式的扩展, 使得 Rmd 格式可以支持公式、定理、图表自动编号和引用、链接, 文献引用和链接等适用于编写书籍的功能。在 bookdown 的管理下一本书的内容可以分解成多个 Rmd 文件, 其中可以有可执行的 R 代码, R 代码生成的文字结果、表格、图形可以自动插入到生成的内容中, 表格和图形可以是浮动排版的。输出格式主要支持 gitbook 格式的网页图书, 这种图书在左侧显示目录, 右侧显示内容, 并可以自动链接到上一章和下一章; 通过单独安装的 LaTeX 编译器支持将书籍转换为一个 PDF 文件, 支持中文; 可以生成 ePub 等格式的电子书。

主要用于编写有多个章节的书籍, 也可以用来生成单一文件的研究报告。

建议使用 RStudio 集成环境制作这样的图书, 该软件内建了一键编译整本书的功能。需要安装 bookdown 扩展包的最新版本。bookdown 扩展包现在还比较新, 还有一些 BUG, 所以尽可能使用最新版的 bookdown 扩展包并且及时更新 RStudio 软件。查看编译的网站建议使用 Google Chrome 浏览器, 此浏览器对 gitbook 的支持较好。

为了新写一本书或者从已有的书转换, 最简单的做法是从 bookdown 的网

站下载 bookdown 配套的例书的 zip 文件 (见<https://github.com/rstudio/bookdown-demo>), 将其解压到本地硬盘某个子目录, 然后修改其中的内容适应自己的书的需要。

因为中文需要一些特殊的设置, 以及在网络条件不好的条件下支持数学公式显示, 本书作者提供了一个粗浅的中文书 bookdown 模板, 下载链接为:

- `Bookdown-template-v0-3.zip`。

## 23.2 一本书的设置

一本用 bookdown 管理的书, 一般放置在某个子目录下, 并作为一个 RStudio 项目 (project) 用 RStudio 管理。

也可以自己新建一个目录, 然后编辑生成必要的文件。注意, 所有的文本文件都要使用 UTF-8 编码。一本 bookdown 书, 一般都需要有一个 `index.Rmd` 文件, 这是最后生成的网站的主页的原始文件, 可以在这个文件中写一些书的说明, 并在开头的 YAML 元数据部分进行有关设置, 如标题、作者、日期等。`index.Rmd` 的一个例子如下:

```

title: " 统计计算"
author: " 李东风"
date: "2019-08-29"
site: bookdown::bookdown_site
output: bookdown::gitbook
documentclass: book
bibliography: [myrefs.bib]
biblio-style: apalike
link-citations: yes
description: " 本科生《统计计算》教材。采用 R 的 bookdown 制作, 输出格式为 bookdown:."

```

```
前言 {-}
```

统计计算研究如何将统计学的问题用计算机正确、高效地实现。

其中在三个减号组成的两行之间的内容叫做 YAML 元数据，是一本书的设置，上例中有书的标题、作者名、日期（用 R 程序自动生成）、描述。其中的 `site` 选项很重要，一定要有这个选项，`site: bookdown::bookdown_site` 使得 RStudio 软件能辨认这是一个 bookdown 图书项目，从而为其提供一键编译快捷方式。元数据中 `output` 项指定默认的输出格式。`documentclass` 项为借助 LaTeX 编译 PDF 格式指定 LaTeX 的模板，现在还不能支持 `ctexbook` 模板所以使用了 `book` 模板。`bibliography` 项指定一个或者几个 .bib 格式的文献数据库。

一个 bookdown 图书项目除了 `index.Rmd` 文件之外，一般还应该有一个 `_bookdown.yml` 文件存放与整本书有关的 YAML 元数据。例如

```
new_session: true
book_filename: 'statcompc'
language:
 label:
 thm: '定理'
 def: '定义'
 exm: '例'
 proof: '证明:'
 solution: '解:'
 fig: '图'
 tab: '表'
 ui:
 chapter_name: ''
delete_merged_file: true
```

其中 `new_session: true` 设置很重要，这使得每一个 Rmd 文件中的 R 程序都在一个单独的 R 会话中独立地运行，避免了不同 Rmd 文件之间同名变量和同名标签的互相干扰。`book_filename` 是最终生成的 LaTeX PDF 图书或者 ePub 电子书的主文件名。`language` 下可以定制一些与章节名、定理名等有关名称。

另外一个需要的设置文件是 `_output.yml` 文件，用于输出格式的设置。这部分内容也可以包含在 `index.Rmd` 的元数据部分。内容如

```

bookdown::gitbook:
 includes:
 in_header: mathjax-local.html
 config:
 toc:
 before: |
 <a href="http://www.math.pku.edu.cn/teachers/lidf/course/statcomp/_boo
 after: |
 编
 download: ["pdf"]
bookdown::pdf_book:
 includes:
 in_header: preamble.tex
 latex_engine: xelatex
 citation_package: natbib
 keep_tex: yes
bookdown::epub_book: default

```

其中的 `style.css` 是自定义的 CSS 显示格式，可以去掉这一行，使用默认的 CSS 格式。这个例子文件分为三部分，`gitbook`、`pdf_book` 和 `epub_book` 三种输出格式分别设置了一些输出选项。在 `gitbook` 部分，设置了目录上方显示的书的主页的链接（`before` 项）和目录下方显示的作者信息。在 `in_header` 部分插入了一部分个性化的 HTML 代码，这部分代码是使用本地的数学公式显示支持以免外网不同时数学公式不能显示，插入的内容将出现在每个生成的 HTML 文件的 `head` 部分。

在 `pdf_book` 部分，设置了通过 LaTeX 编译整本书为 PDF 的一些选项。指定了 `latex_engine` 为 `xelatex`，这对中文支持很重要。`in_header` 选项要求在 LaTeX 文件导言部分插入一个 `preamble.tex` 文件，内容如：

```

\usepackage{ctex}

%\usepackage{xltextra} % XeLaTeX的一些额外符号
% 设置中文字体
\setCJKmainfont[BoldFont={黑体},ItalicFont={楷体}]{新宋体}

```

```

\usepackage{amsthm,mathrsfs}
\usepackage{booktabs}
\usepackage{longtable}
\makeatletter
\def\thm@space@setup{%
 \thm@preskip=8pt plus 2pt minus 4pt
 \thm@postskip=\thm@preskip
}
\makeatother

```

其中很重要的是使用 `ctex` 包来支持中文。在元数据中也可以指定一些 LaTeX 选项，比如指定页边距等：

```

classoption: twoside
fontsize: 12pt
linestretch: 1.5
geometry: "left=4cm, right=3cm, top=2.5cm, bottom=2.5cm"
fontsize: 12pt
linestretch: 1.5
toc-depth: 1
lof: True
lot: True

```

在 `bookdown` 项目中与 `index.Rmd` 同级的所有 `.Rmd` 文件都自动作为书的一章，除非文件名以下划线开头。这样做的好处是作者可以任意地增删章节，编译整本书时章节编号会自动调整。但是，章节的顺序将按照文件名的字典序排列，所以，所有的包含一章内容的 `.Rmd` 文件，最好命名为类似 `0201-rng.Rmd` 这样的名字，文件名前面人为地加上排序用的序号，使得章节按照自己的次序排列。实际上，也可以设置不自动将每个 `.Rmd` 文件都作为一章，而是在 `_output.yml` 中设置一项 `rmd_files`，列出所有需要作为一章的文件，并以列出次序编译，如

```
rmd_files: ["index.Rmd", "rng.Rmd", "simulation.Rmd", "refs.Rmd"]
```

这时，应该添加一个 `_site.yml` 文件，内容如：

```
site: "bookdown::bookdown_site"
output: bookdown::gitbook
```

## 23.3 章节结构

除了 `index.Rmd` 文件，项目中每个 `.Rmd` 文件都作为一章。每个 `.Rmd` 文件第一行，应该是以一个井号和空格开头的一级标题，后面再加空格然后有大括号内以井号开头的章标签，如

```
随机数 {#rng}
```

这些章标签去掉井号后会作为生成的 HTML 文件的名称，所以一定要有章标签，而且章节标签在全书中都不要重复以免冲突。文件内可以用两个井号和一个空格开始的行表示节标题，最后也应该有大括号内以井号开头的节标签，如

```
随机数 {#rng}
```

```
均匀随机数发生器 {#rng-unif}
```

使用 `bookdown` 写书，一般每章不要太长，否则编译预览很慢，读者浏览网页格式也慢。

内容相近的章节可以作为一个“部分”。为此，在一个部分的第一个章节文件的章标题前面增加一行，以 `# (PART)` 开头，以 `{-}` 结尾，中间是部分的名称，如

```
(PART) 随机数和随机模拟 {-}
```

```
随机数 {#rng}
```

书的最后可以有附录，附录的章节将显示为 `A.1`, `B.1` 这样的格式。为此，在附录章节的第一个文件开头加如下前两行的标题行：

```
(APPENDIX) 附录 {-}
```

```
(PART) 附录 {-}
```



```
一些定理的证明 {#formula}
```

## 23.4 书的编译

建议使用 RStudio 软件编辑内容，管理和编译整本书。

在 `index.Rmd` 或者 `_bookdown.yml` 中设置 `site: bookdown::bookdown_site` 后，RStudio 就能识别这个项目是一个 bookdown 项目，这时 RStudio 会有一个 Build 窗格，其中有“Build book”快捷图标，从下拉菜单中选择一个输出格式（包括 `gitbook`、`pdf_book`、`epub_book`），就可以编译整本书。对 `gitbook` 格式，即 HTML 网页格式，编译完成后会弹出一个预览窗口，其中的“Open in Browser”按钮可以将内容在操作系统默认的网络浏览器中打开。

另一种办法是在命令窗口用如下命令编译（以输出 `gitbook` 为例），我个人认为这种办法更好用：

```
bookdown::render_book("index.Rmd",
 output_format="bookdown::gitbook", encoding="UTF-8")
```

编译结果默认保存在 `_book` 子目录中，可以在 `_bookdown.yml` 中设置 `output_dir` 项改为其它子目录。编译整本书为 `pdf_book` 格式时，如果成功编译，也会弹出一个 PDF 预览窗口。可以在 `_book` 子目录中找到这个 PDF 文件。

将书编译为 PDF 需要利用 LaTeX 编译器，这需要单独安装 LaTeX 编译软件，如 Windows 下的 CTEX 套装软件。LaTeX 编译器对输入要求十分严格，一丁点儿错误都会造成整本书的编译失败，所以对于不熟悉 LaTeX 的用户，不建议使用 bookdown 的 `pdf_book` 输出格式。如果仅在 R Markdown 中使用 LaTeX 编译器，可以安装谢益辉的 TinyTeX。

对于较短的书，做了一定修改后都可以重新编译 `gitbook` 结果和 `pdf_book` 结果。在书比较长了以后，每次编译都花费很长时间，所以可以仅编译 `gitbook` 格式的一章，修改满意后再编译整本书。仅编译一章也需要所有的 `.Rmd` 文件都是已经编译过一遍的，新增的 `Rmd` 文件和图形文件会使得编译单章出错，每次新增了 `Rmd` 文件和图形文件都应该重新编译整本书，但是内容修改后不必要重新编译整本书，可以仅编译单章。

编译单章现在没有快捷图标，只能在 RStudio 控制台（命令行）运行如下命令：

```
bookdown::preview_chapter("chap-name.Rmd",
 output_format="bookdown::gitbook", encoding="UTF-8")
```

其中 `chap-name.Rmd` 是要编译的单章的文件名。编译完成后在结果目录（默认是 `_book`）中找到相应的 HTML 文件打开查看，再次编译后仅需在浏览器中重新载入文件。建议使用 Google chrome 浏览器，用 MS IE 或者 Edge 浏览器对 gitbook 的 Javascript 支持不够好，使得目录的层级管理、自动滚动、单章编译后的目录更新不正常，而 chrome 则没有问题。

编译单章也不能解决所有的问题，有些问题还是需要编译整本书，而章节很多时整本书编译又太慢。为此，可以在项目中增加一个临时的部分内容子目录，如 `testing` 子目录，在子目录中存放相同的设置文件 `index.Rmd`、`_bookdown.yml`、`_output.yml`，以及图形文件、文献数据库文件，并将要检查的若干章节复制到 `testing` 子目录中，在 `testing` 中新建一个 bookdown 项目，然后编译其中的整本书。这在调试部分章节的 HTML 和 PDF 输出时很有效。解决问题后主要将修改过的章节复制回原始的书的目录中。

有时仅仅想验证某个长数学公式或者表格，用上述的编译单章或者单独一个小规模测试项目的办法也不经济。这时，单独开一个备用的普通 RStudio 项目，不能是 bookdown 项目，在其中的 Rmd 文件中验证数学公式和表格的编排，这样效率最快了。

## 23.5 交叉引用

在写作时，每个一级到三级标题都应该有自定义的标签，格式是在标题行末尾空格后添加 `{#label}`，其中 `label` 是自己指定的标签，使用英文、数字、减号，不要使用中文，而且整本书不要有重复的标签。为避免不同章节使用了重复标签，可以取 `label` 的前一部分为所在章节的文件名。

如果要引用某一章节，有如下的做法：

- `$\@ref(label)`，`label` 是某个标题对应的标签。结果显示为如 `$3.1.1` 这样的章节号，并可点击，点击时跳跃到相应的章节。

- [链接文本](#label) 其中 label 是某个标题对应的标签。结果产生一个链接，显示为链接文本，点击时跳到 label 对应的章节。

## 23.6 数学公式和公式编号

通过 R 的 knitr 和 rmarkdown 扩展包，.Rmd 格式文件已经支持数学公式，见 R Markdown 说明。

在用 \$\$ 符号在两端界定的公式后面，可以用 `\tag{标号}` 命令增加人为的公式编号，如

```
$$
y = f(x)
\tag{*}
$$
```

结果显示为

$$y = f(x) \tag{*}$$

要注意的是，在 \$\$ 界定的数学公式内用了 aligned 环境后，仅能在 `\end{aligned}` 之后加 `\tag{标号}` 命令，而不能写在 aligned 环境内。这样，多行的公式将不能为每行编号。

用 `\tag` 命令人为编号比较简单易用，但是在有大量公式需要编号时就很不方便，只要增加了一个公式就需要人为地重新编号并修改相应的引用。bookdown 包支持对公式自动编号，并可以按公式标签引用公式，引用带有超链接。

bookdown 的自动编号对 LaTeX 的 equation 环境、align 环境都可以使用，而且不需要在两端用 \$\$ 界定。在公式的末尾或者一行公式的 `\\` 换行符之前，写 `(\#eq:mylabel)`，其中 mylabel 是自己给公式的文字标签，文字标签可以使用英文字母、数字、减号、下划线。如

```
\begin{align}
f(x) &= \sum_{k=0}^{\infty} \frac{1}{k!} x^k \label{eq:efunc-sum} \\
&= e^x \label{eq:efunc-ex}
```

```
\end{align}
```

将会对两行公式自动编号。引用公式时，用如`\@ref(eq:mylabel)`，其中 `mylabel` 是公式的自定义标签，编译后这样的引用会变成带有链接的圆括号内的编号。

公式编号在全书中都不要有冲突（不同的公式定义了相同的编号）。一种办法是，自定义的公式标签的开头以章节文件名开头。

## 23.7 定理类编号

定理、引理、命题、例题等，使用特殊的 markdown 代码格式，以三个反单撇号开头，以三个反单撇号结尾，在开头的三个反单撇号后面写 `{theorem}` 表示定理。在 `theorem` 后面，用逗号分隔后写一个定理的自定义标签，因为现在 bookdown 的功能还不完善，所以所有的定理类都应该自定义一个标签。可以用 `name="定理名称"` 指定一个显示的定理名。标签也可以写成 `label="mythlabel"`，其中 `mythlabel` 是自定义的标签。

设某个定理的自定义标签是 `mythlabel`，则可以用如`\@ref(thm:mythlabel)` 引用此定理的编号，编号有自动生成的链接。

例如：

```
``{theorem norlim-weakconv, name="弱收敛"}
 ξ_n 依分布收敛到 ξ ，
当且仅当对任意 \mathbb{R} 上的一元实值连续函数 $f(\cdot)$ 都有

$$E f(\xi_n) \rightarrow E f(\xi), \quad n \rightarrow \infty$$

``
``
```

当定理或例子内有列表时，一定注意列表前后要空行，否则会导致嵌套错误。这种错误在编译 HTML 时无法发现，但是会造成结果莫名其妙地出错。

bookdown 提供了证明环境，但是不太实用。

对例题，将 `theorem` 替换成 `example`，在引用时将 `thm` 替换成 `exm`。如：

```
``{example noralim-aspr-ce, name="依概率收敛不a.s.收敛反例"}
a.s.收敛推出依概率收敛，
但是反之不然。给出反例。
...`
```

## 23.8 文献引用

bookdown 使用 .bib 格式的文献数据库, 关于 .bib 格式的文献数据库请参考 LaTeX 的有关说明。在 index.Rmd 的 YAML 元数据部分或者 \_bookdown.yml 中用 bibliography 可以设置使用的一个或者多个 .bib 格式的文献数据库文件。设某篇文章的 .bib 索引键是 Qin2007:comp, 用 @Qin2007:comp 可以引用此文献, 用 [Qin2007:comp] 可以生成带有括号的引用, 引用都有超链接。

指定 .bib 文件时可以用相对路径, 如 “../docs/mybib.bib”。

下面是一个样例 .bib 文件的内容 (主要要用 UTF-8 编码保存):

```
% Encoding: UTF-8

@Book{MWP06-HighStat,
 author = {茆诗松 and 王静龙 and 濮晓龙},
 title = {高等数理统计},
 year = {2006},
 edition = {第二版},
 publisher = {高等教育出版社}
}

@BOOK{Unwin-Visualize06,
 title = {Graphics of Large Datasets Visualizing a Million},
 publisher = {Springer},
 year = {2006},
 author = {Antony Unwin and Martin Theus and Heike Hofmann}
}
```

```

@Article{Wichmann1982:RNG,
 author = {Wichmann, B. A. and Hill, I. D.},
 title = {Algorithm as 183: An efficient and portable pseudo-random number gene
 journal = {Applied Statistics},
 year = {1982},
 volume = {31},
 pages = {188 - 190. Remarks: 34, 198 and 35, 89},
}

@Book{QGCZ2011:StochProc,
 author = {钱敏平 and 龚光鲁 and 陈大岳 and 章复熹},
 title = {应用随机过程},
 year = {2011},
 publisher = {高等教育出版社},
}

```

## 23.9 插图

bookdown 图书的插图有两种，一种是已经保存为图形文件的，主要是 png、jpg 和 pdf 图片；另一种是文中的 R 代码生成的图形。

已经有图形文件的，可以用 markdown 格式原来的插图方法，见 markdown 格式介绍。但是，这样做不能给图形自动编号，另外因为制作图书是有网页和 PDF 书两种主要输出格式的，原有的插图方式在这两种输出格式上有细微的不一致。所以，最好是统一使用 Rmd 的插图方法。

Rmd 的插图方法就是写一段 R 代码段来插图，如果是用程序作图，则代码中写作图的代码；如果是已有的图形文件，可以在一个单独的 R 代码段中用类似下面的命令插图：

```
knitr::include_graphics("figs/myfig01.png")
```

其中 `figs` 是存放图形文件的子目录名，`myfig01.png` 是要插入的图形文件名。这样，如果同时还有 `myfig01.pdf` 的话，则 HTML 输出使用 png 图片

而 PDF 输出自动选用 pdf 文件。另外，插图的选项在代码段的选项中规定：用代码段的 `fig.with` 和 `fig.height` 选项指定作图的宽和高（英寸），用 `out.width` 和 `out.height` 选项指定在输出中实际显示的宽和高，实际显示的宽和高如果使用如 "90%" 这样的百分数单位则可以自动适应输出的大小。

为了使得插图可以自动编号并可以被引用，为代码段指定标签并增加一个 `fig.cap="..."` 选项指定图形标题。代码段的标签变成浮动图形的标签，如 `myfiglabel`，则为了引用这个图只要用 `\@ref(fig:myfiglabel)`。注意，在整本书中这些标签都不能重复，否则编译 LaTeX 支持的 PDF 输出会失败。

有些插图会伴随很长的说明文字，这可以用代码段的 `fig.cap=` 选项指定，但是其中的 Markdown 特有的格式在转换 LaTeX 时不一定支持，而且在代码段选项中写太长的文字说明也是的程序难以辨认。所以，可以使用文字引用的方式：在单独的一段中，用如下格式定义一段可引用的文字内容：

其中 `\texttt{mylabel}` 是自己定义的仅由英文大小写字母、数字和减号组成的引用标志符。在需要使用这注意定义和引用都是用的 `\texttt{(ref:mylabel)}` 语法。

由于 PDF 中中文不能自动识别，所以在每个源文件的开头应该加上如下的设置，使得生成 PDF 图时中文能够正确显示：

```
\begin{verbatim}
```\r setup-pdf, include=FALSE}
pdf.options(family="GB1")
```\end{verbatim}
```

其中 `include=FALSE` 表示要不显示代码段的代码，有运行结果也不插入到输出结果中，是否允许运行视缺省的 `eval=` 的值而定。

## 23.10 表格

bookdown 书的表格也有两种，一种是原来 markdown 格式的表格，最好仅使用管道表，管道表对中文内容支持最好。为了对这样的表格自动编号，需要在表格的前面或者后面空开一行的位置，写

Table: `\label{tab:mylabel}` 表的说明

其中 `mylabel` 是自定义的表格标签。在引用这个表时用如 `\@ref(tab:mylabel)`。

另一种表格是 R 代码生成的表格，主要使用 `knitr::kable()` 函数。在 `knitr::kable()` 函数中用选项 `caption=` 指定表格的说明文字（标题），这时生成表格的 R 代码段的标签，如 `myfiglab`，就自动构成了表格的引用标签主干，实际引用如 `\@ref(tab:myfiglab)`。

为了适应较长的表，可以在 LaTeX 的 `preamble.tex` 中引入 `longtable` 包，并在 `knitr::kable()` 中加选项 `longtable=TRUE`。

## 23.11 数学公式的设置

bookdown 在生成 PDF 时使用 LaTeX 软件，所以 PDF 输出的数学公式的支持很好，但是 LaTeX 编译器也很挑剔，稍微一点错误也造成编译失败。比如，在行内公式内部如果紧邻 `$` 符号有多余的空格，如 `$ y $`，编译 PDF 时会出错。

bookdown 0.6 版生成的 gitbook 格式的网页书籍，在有数学公式时，使用 MathJax 库在浏览器中显示数学公式。MathJax 是用于网络浏览器中显示数学公式的优秀的 Javascript 程序库，可免费使用。但是，当数学公式中含有中文（用 `\text{}` 或 `\mbox{}` 命令）时，数学公式可能会显示不正常。另外，数学公式默认使用远程服务器上的 MathJax 程序库处理，在网络不通畅时显示很慢或者无法显示。

MathJax 有多种渲染输出选择，gitbook 的模板中已经固化了一种 Common-HTML，这种输出与 gitbook 的 `style.css` 配合，当中文公式在 `section` 标记内时，使得中文显示不正常。如果不修改 gitbook 的模板，通过在 `_output.yml`



设置文件中插入设置命令的办法不奏效,因为 gitbook 的模板是动态调入 MathJax 库的,不能在 HTML 文件头或尾插入静态设置命令修改输出格式。

为此,直接修改 bookdown 包中的 gitbook.html 模板,删去文件末尾动态调入 MathJax 的部分,改为在 `_output.yml` 中要求调入一个设置文件。但是,bookdown 的 0.9 版本已经没有这个问题,不需要人为修改 bookdown 包的 gitbook.html 模板了。

为了避免远程调用 MathJax 程序库的麻烦,改为本地使用。将 MathJax 安装在了书生成的网站主目录的上三层,用 `../../../../../MathJax/mathjax.js` 路径访问。假设下载整个 MathJax 库后解压放在了书的网页文件所在子目录的上三层的位置。

增加设置文件 `_header.html`:

```
<script type="text/x-mathjax-config">
MathJax.Hub.Config({
 jax: ["input/TeX","output/SVG"],
 extensions: ["tex2jax.js","MathMenu.js","MathZoom.js"],
 TeX: {
 extensions: ["AMSmath.js","AMSsymbols.js","noErrors.js","noUndefined.js"]
 }
});
</script>
<script type="text/javascript"
 src="../../../../../MathJax/MathJax.js">
</script>
```

如果希望用远程服务器上的 MathJax,可以使用如下的 `mathjax-cdnjs.html` 设置文件:

```
<script type="text/x-mathjax-config">
MathJax.Hub.Config({
 jax: ["input/TeX","output/HTML-CSS"],
 extensions: ["tex2jax.js","MathMenu.js","MathZoom.js"],
 TeX: {
 extensions: ["AMSmath.js","AMSsymbols.js","noErrors.js","noUndefined.js"]
 }
});
</script>
```

```

 }
 });
</script>
<script type="text/javascript"
 src="https://cdnjs.cloudflare.com/ajax/libs/mathjax/2.7.2/MathJax.js">
</script>

```

为了调用这样的设置文件，在 `_output.yml` 设置文件中如下设置：

```

bookdown::gitbook:
 css: style.css
 includes:
 in_header: _header.html

```

注意，MathJax 要求先设置再调入。由于浏览器对 MathJax 输出方式有记忆，所以如果不能正常显示中文公式，需要右键单击公式，选择“Math Settings—Math Renderer”中的“SVG”或者“HTML-CSS”。

## 23.12 使用经验

Ed Berry 在网上分享了用 bookdown 生成 PDF 学位论文的经验：<https://eddberry.netlify.com/post/writing-your-thesis-with-bookdown/>。Chester Ismay 提供了一个 R 扩展包 `thesisdown`，见 <https://github.com/ismayc/thesisdown>，例子见 <https://thesisdown.netlify.com/>。

有数学公式时，对行内公式边缘处多余的空格十分敏感，所以行内公式边缘不要有空格。

如果已经有用 LaTeX 写的书，要转换为 bookdown 的 Rmd 格式，可以用 RStudio 的支持 RegEx 的替换模式，如

- `\\keyword{([~]+?)}` 替换成 `**\1**`。
- `\\textbf{([~]+?)}` 替换成 `**\1**`。
- `\\texttt{([~]+?)}` 替换成 `\1`。

可以写一个函数对 LaTeX 文件进行转换生成 Rmd 文件，再手工修改。函数如

```

latex2rmd <- function(fname="_tmp/tomd.tex", encoding="UTF-8"){
 lines <- readLines(fname, encoding=encoding)
 if(encoding != "UTF-8")
 lines <- iconv(lines, from=encoding, to="UTF-8")
 print(head(lines, 10))
 lines <- gsub("\\\\keyword[{}](\\^+)?[{}]", "\\1", lines, perl=TRUE)
 lines <- gsub("\\\\textbf[{}](\\^+)?[{}]", "\\1", lines, perl=TRUE)
 lines <- gsub("\\\\texttt[{}](\\^+)?[{}]", "\\1", lines, perl=TRUE)
 lines <- gsub("\\\\label[{}](eq:\\^+)?[{}]", "\\#\\1", lines, perl=TRUE)
 lines <- gsub("\\\\eqref[{}](\\^+)?[{}]", "\\@ref(\\1)", lines, perl=TRUE)
 lines <- gsub("\\\\bs\\\\", "\\boldsymbol\\\\", lines, perl=TRUE)
 lines <- gsub("\\\\bs ", "\\boldsymbol ", lines, perl=TRUE)
 lines <- gsub("\\\\Rbb", "\\mathbb R", lines, perl=TRUE)
 lines <- gsub("\\\\Zbb", "\\mathbb Z", lines, perl=TRUE)
 lines <- gsub("\\\\Var[()]", "\\text\\{Var\\}(", lines, perl=TRUE)
 lines <- gsub("\\\\Cov[()]", "\\text\\{Cov\\}(", lines, perl=TRUE)
 lines <- gsub("\\\\S(\\d)", "\\$\\2", lines, perl=TRUE)
 lines <- gsub("\\\\S ", "\\$ ", lines, perl=TRUE)
 lines <- gsub("\\\\defeq", "\\stackrel{\\triangle}{=}", lines, perl=TRUE)
 lines <- gsub("\\\\argmin_", "\\mathop{\\text{argmin}}_", lines, perl=TRUE)
 lines <- gsub("^\\s*\\\\item", "*", lines, perl=TRUE)
 lines <- gsub("\\\\begin[{}]{frame}[{}]", "", lines, perl=TRUE)
 lines <- gsub("\\\\end[{}]{frame}[{}]", "", lines, perl=TRUE)
 lines <- gsub("\\\\begin[{}]{itemize}[{}]\\[<\\+>\\\\", "", lines, perl=TRUE)
 lines <- gsub("\\\\end[{}]{itemize}[{}]", "", lines, perl=TRUE)
 lines <- gsub("\\\\begin[{}]{itemize}[{}]", "", lines, perl=TRUE)
 lines <- gsub("\\\\bm ", "\\boldsymbol ", lines, perl=TRUE)
 lines <- gsub("\\\\bm\\\\", "\\boldsymbol\\\\", lines, perl=TRUE)
 lines <- gsub("\\\\tr[()]", "\\text{tr}(", lines, perl=TRUE)
 lines <- gsub("\\begin{align*}", "\\$\\begin{aligned}", lines, fixed=TRUE)
 lines <- gsub("\\end{align*}", "\\end{aligned}\\$\\$", lines, fixed=TRUE)

 print(head(lines, 20))
}

```

```
##writeLines(lines, "_tmp/fromtex-clean.Rmd")
con1 <- file("_tmp/fromtex-clean.Rmd", "wt", encoding="UTF-8")
writeLines(lines, con=con1)
close(con1)
}
```

Rmd 格式对算法编排的支持不够好。编排算法可以用表格来换行，仍用`\qqquad`和`\qquad`缩进，不要用空格缩进，因为在 LaTeX 转 PDF 时空格会损失。但是算法内容中有公式含有竖线时还是无法与表格线区分开来。

编排算法也可以用数学公式，写在 `$$` 界定范围内，用 `aligned` 环境分行，缩进使用`\quad`和`\qqquad`。只是无法对 `if` 这样的关键字用重体排印。

为了能够生成中文的 PDF，不要指定 `docclass` 为 `ctexbook`，而是指定为 `book`，然后在 `preamble.tex` 中引入 `ctex` 包。

为了 PDF 输出，不要引入太多的数学包，因为从 markdown 到 HTML 不支持复杂的数学。

用 R 作图时如果图形中有汉字，在代码块选项中加上 `dev="png"`，`dpi=300`。否则生成 PDF 时会有中文编码问题。另一办法是在每个.Rmd 文件开头的 `setup` 源代码段插入

```
pdf.options(family="GB1")
```

这样可以生成支持中文字的 PDF 图形。

在编译出错时，会在主目录留下编译的 `tex` 源文件及相关文件。但是，此 `tex` 源文件中使用的 R 生成的图片路径不对，需要将 `tex` 源文件复制到 `bookdown_files` 目录，将直接插入的图片也复制到这个目录，然后编译 `tex` 文件发现问题，逐个修复。建议建立小的测试项目专门调试有问题的文件。

连分数的加号是`\genfrac{}{}{0pt}{}{}{+}`。

### 23.13 bookdown 的一些使用问题

- 在数学公式中用`\text{\@ref(eq:label)}`引用公式，HTML 成功，LaTeX 版本会有 BUG，重复了`\`。如果使用文内的链接，则 LaTeX 成功，

HTML 不成功。

- YAML 部分有的中文文字（如作者名）会出错，代以 ASCII 则不出错。
- example 的自动编号有问题，不加 label 的 example，其 HTML 结果在不同章之间会编号混淆，避免问题的临时办法是 example 都加上 label。
- 用本地文件格式提供的下载文件不能用中文文件名。



## Chapter 24

# 用 R Markdown 制作简易网站

### 24.1 介绍

为了从多个.Rmd 源文件制作网站，可以使用 `blogdown` 扩展包或者 `bookdown` 扩展包。`bookdown` 扩展包可以生成 `gitbook` 格式的网站，带有左侧的章节目录和前后页面的导航链接，很适用于一本书的网站。但是，`bookdown` 使用时需要重新编译整个网站才能产生正确的链接。

也可以仅利用 `rmarkdown` 包的 `render.site()` 功能制作简易的网站。

### 24.2 简易网站制作

一种简易的制作网站的办法是仅仅利用 `rmarkdown` 包的 `render.site()` 功能。为了使得 `rmarkdown` 和 `RStudio` 将一个子目录（项目）看成一个网站项目，只要此项目中含有 `index.Rmd` 和 `_site.yml` 文件。所有网页.Rmd 源文件都必须在项目的顶级目录中。只要将编译所得的网站目录上传到任意的网站服务器就可以发布到网上。

### 24.2.1 网站结构

`index.Rmd` 是主页内容，可以在此处人工加入其它页面链接（参见 markdown 格式说明中链接写法），也可以制作单独的目录页面。内容如

```

title: " 概率论"

```

```
* [概率分布](dist.html)
* [期望](expect.html)
```

`_site.yml` 是一个 YAML 文件，其中包含站点的设定和输出设定。内容如

```
name: " 概率论"
output_dir: "_probbbook"
output:
 html_document:
 toc: true
 mathjax: "../../MathJax/MathJax.js"
 self_contained: false
 includes:
 in_header: "_header.html"
navbar:
 title: " 概率论"
 left:
 - text: "Home"
 href: index.html
 - text: "About"
 href: about.html
```

这里 `name` 设置站点名称，`output_dir` 给出生成的 html 及其他辅助文件的存放目录，缺省时用 `_site` 子目录。`navbar` 域设置了网站的菜单，这里包括 Home(主页), Contents(目录), About(关于) 三个页面的导航菜单。`output` 域设定输出的选项，其中 `html_document` 就是用 `rmarkdown::render_site()` 生成的网站的输出文件格式，其中 `toc: true` 说明每个页面都有自身内容的



目录, 关于 `mathjax`, `self_contained`, `includes` 都与数学公式设置有关, 这里设置数学公式使用与本网站在同一目录系统或同一网站地址存放的 MathJax 数学公式显示库, 该库放在生成的 HTML 文件所在目录向上两层的 MathJax 子目录中。 `in_header` 指定一个要插入到生成的每个 HTML 的 head 部分的内容文件, 这里内容文件 `_header.html` 的内容是

```
<script type="text/x-mathjax-config">
MathJax.Hub.Config({
 jax: ["input/TeX","output/SVG"],
 extensions: ["tex2jax.js","MathMenu.js","MathZoom.js"],
 TeX: {
 extensions: ["AMSmath.js","AMSsymbols.js","noErrors.js","noUndefined.js"]
 }
});
</script>
```

是对 MathJax 的一些设置, 主要使用 LaTeX 作为输入格式。

### 24.2.2 编译

具有 `index.Rmd` 和 `_site.yml` 文件的项目, 在 RStudio 软件中会显示一个 Build 窗格, 点击其中的 Build Website 可以自动调用 `rmarkdown::render_site()` 生成整个站点, 存放在 `output_dir` 指定的输出目录中, 默认目录名为 `_site`。为了将生成的站点发布到网站, 只要将输出目录的所有内容全部复制到服务器的某个子目录中就可以了。

`rmarkdown::render_site()` 生成整个站点时会自动逐个编译.Rmd 文件, 并将项目目录中其它的文件和子目录复制到输出目录中, 以句点或者下划线开头的文件和子目录不复制, `.R`, `.r`, `.Rmd` 文件不复制。

因为是自动依次编译所有的.Rmd 文件, 所以中间某个源文件编译错误就会使得整个编译失败, 修正错误后还是需要再次编译所有文件, 这一点不太方便。但是, `rmarkdown::render_site()` 提供了这样一种操作模式: 将有错的文件先移动到其它目录中, 编译整个站点, 这时生成的站点仅包含确认无误的各个页面; 然后, 再将可能有错的文件移动回到项目目录, 逐个地打开每个未成功编译的文件, 用 Knit 按钮单独编译, 编译成功后结果文件也自动进入输出目录。

RStudio 的 Knit 按钮或者 `rmarkdown::render_site()` 命令可以人工控制一个一个地编译.Rmd 文件, 编译成功一个就存放到输出目录中一个, 这是用 `rmarkdown::render_site()` 制作网站比用 bookdown 扩展包制作网站的一个优点。因为网站的目录是自己编写的链接, 所以这样逐个编译不会破坏网站的链接。编译单个文件的命令如

```
rmarkdown::render_site("sec1.Rmd", output_format = "html_document", encoding="UTF-
```

其中 `sec1.Rmd` 是某个网页内容源文件。

可以用 `rmarkdown::clean_site()` 删除生成的文件, 包括程序结果缓存。

如果用 bookdown 包, 只有编译所有文件才能生成正确的网站目录, bookdown 也能编译单个文件, 但是仅能作为调试, 调试成功后还是需要编译所有文件才能使得网站目录正确。

bookdown 站点的优势在于自动生成网站目录而不需要人工管理目录, 而且 bookdown 支持图书的公式、定理、图标自动编号和链接, 文献列表和文献引用, 所以 bookdown 是比 `rmarkdown::render_site()` 编译调试速度比较慢但是功能更强大的一种制作网站的工具, 而且还可以生成 PDF 版图书。

所以, 一本书如果不需要很多的公式自动编号、图表编号、文献, 在编写阶段可以使用 `rmarkdown::render_site()` 制作, 定稿后再改成 bookdown 图书格式, 实际上每个源文件内容部分不需要改变, 只需要修改一下头部就可以了。站点的设置需要将 `_site.yml` 中的 `site` 域的值改为 `bookdown::bookdown_site`, 并增加 `_bookdown.yml` 和 `_output.yml` 两个设置文件。详见23。

### 24.2.3 内容文件

内容文件只要用普通的 Rmd 文件即可, 如某一网页的源文件如

```

title: " 期望 "

```

随机变量 **\*\* 数学期望 \*\*** 可以看成是对随机变量取值的加权平均, 某个值的取值概率越大, 加权越大。

对离散随机变量，  
期望是加权和；  
对连续型随机变量，  
期望是用密度函数加权对  $x$  积分。

#### 24.2.4 网站设置

在 `_site.yml` 文件中进行网站设置，样例见24.2.1。顶级设置有：

- `name`: 网站名称；
- `output_dir`: 保存编译结果的子目录名，默认为 `_site`，如果设置为 `.`，则表示编译结果放在项目的顶级目录中而不是子目录中；
- `includes`: 指定输出的网站中要包含的文件，输出时默认不包含 `.R`，`.r`，`.Rmd`，`.RData`，`.rds` 文件和文件名以小数点或者下划线开头的文件，所以例外要写在这里，如 `includes: ["data1.R", "data2.R"]`；
- `excludes`: 指定项目中不希望输出的结果网站的文件，可以有通配符，如 `excludes: ["backup.txt", "*.xlsx"]`；
- `output`: 其中的 `html_document` 属性的设置是所有 Rmd 网页源文件的共同设置，参见22.11.3，每一个 Rmd 文件本身也可以有自己的属性设置；
- `navbar`: 用来设置导航菜单，属性 `left` 指定主菜单栏，根据所用的网站主题，可能显示在左侧或者顶部。属性 `right` 指定一个右侧导航栏。

`includes` 和 `excludes` 设置主要是为了改变默认的文件发布规则，所有的 `.Rmd` 文件只要不是文件名以下划线开头就都会被编译并且复制到结果子目录中。

### 24.3 用 blogdown 制作网站

R 扩展包 `blogdown` 可以与 Hugo 软件配合制作简单的静态网站。网站的所有文件都存在于一个目录中，只要上传到任意的网站服务器就可以发布，没有任何限制。

网站内容用 R Markdown 格式编写，在 R 或者 RStudio 中编译为网站。这样的做法使得网站内容容易再现，而且也可以将内容转换为 PDF、Word 等格式。

blogdown 包的 R Markdown 格式支持 bookdown 包的扩充，如公式、图表编号与引用。

blogdown 包的名称虽然包含了 blog，但是其建站并不限于博客类型，任何类型的静态网站都可以构建，比较适合的是个人网站。

参考：

- (Yihui Xie, 2017)
- (Xie et al., 2019)

### 24.3.1 生成新网站的框架

在 RStudio 中，菜单“File – New Project – blogdown site”可以生成一个网站框架，注意一定要新建目录而不是用已有目录，而且第一次新建网站时还会自动下载 Hugo 静态网站生成程序。

生成的网站框架包括如下文件与子目录：

- index.Rmd: 主页面的源文件。
- config.toml: 网站设置文件。
- ./content/: 存放作为网站内容的 md 或者 Rmd 文件。可以在下面再建立子目录。./content/post/用来保存博客类型的网页，但是不做要求。
- ./resources:
- ./static: 保存图片、CSS 等文件，在发布网站时会被复制到./public 子目录中。
- ./themes: 存放多个网站主题。
- ./public: 发布网站时所用到的所有文件。生成的内容被复制到这个目录中。只要将这个目录的内容全部复制到任何一个支持静态内容的网站服务器，就可以发布网站。

这样建站，建成的网站应该在 RStudio 的 Viewer 窗格自动显示，如果没有自动显示网站或者再次打开网站项目，点击图标栏的 `Addins -- Serve site` 可以启动网站服务器发布该网站。或者用命令

```
blogdown::serve_site()
```

启动网站服务器。每次在 RStudio 或者 R 中打开一个网站项目，都只需要启

动一次服务器。blogdown 的服务器会在后台运行，侦测到文件修改时会自动重建，这称为在线重建 (Live Reload)。在线重建功能基于 R 扩展包 `servr`。所以用户只要发布新的消息或者修改已有消息，不用关心后台的转换问题。

blogdown 提供了一些 RStudio 的 addin 程序，可以为网站功能提供方便，比如 `New Post` 可以生成新消息，`Update Metadata` 可以修改消息的 YAML 元数据，`Insert Image` 可以在消息中插入图片。

使用在线重建功能时，应关闭一些其他功能。在 RStudio 的 “Tools – Project Options – Build Tools” 对话框中，取消选定 “Preview site after building” 和 “Re-knit current preview when supporting files change”。

如果不使用在线重建功能，即不使用 `Serve site`，可以在 RStudio 的 Build 窗格中用 “Build Book” 功能手工控制网站的重建。

### 24.3.2 网页内容文件及其设置

网页内容文件可以用 Rmd 文件也可以用 md 文件，但是 Rmd 文件使用 Pandoc 和 bookdown 包支持的文件格式，功能更强，比如支持 LaTeX 数学公式，公式、图表自动标号与交叉索引，等等，所以最好是用 Rmd 文件。

如果有某个 Rmd 源文件暂时不想编译输出，可以将其扩展名临时改为一个不被编译的扩展名如 `.txt`。

blogdown 编译结果的网页是 HTML 格式，其 Rmd 元数据输出类型为 `blogdown::html_page`，此类型基于 `bookdown::html_document2`，又基于 `rmarkdown::html_document`。这样，可以在 Rmd 源文件的元数据部分对该输出类型设置相应的属性，如 `toc: true` 等。例子：

```

title: " 一个试验网页 "
author: " 李东风 "
date: "2019-07-08"
output:
 blogdown::html_page:
 toc: true
 fig_width: 6
```

```
dev: svg
```

为了给所有输出设置统一的属性，可以在项目目录中添加 `_output.yml` 文件，在其中设置输出文件的各种属性，如：

```
blogdown::html_page:
 toc: true
 fig_width: 6
 dev: svg
```

### 24.3.3 初学者的工作流程

Hugo 建站程序比较复杂，如果要挑选自己喜爱的网站主题 (theme)，移植网站也比较麻烦。对于了解网站技术不多的建站者，建议用如下的流程建立新站：

- 在 <https://themes.gohugo.io/> 仔细地挑选一个合适的网站主题 (theme)；
- 在 RStudio 中通过生成一个在已有目录的新项目；
- 在新项目的控制台中提交命令 `blogdown::new_site(theme = 'user/repo')`，其中 `user/repo` 是前面所选的主题的用户名与资源 (repository) 名；
- 试验新网站看是否满意。如果不满意可以重复以上步骤，如果满意，就可以修改 `config.toml` 中的设置。如果某个选项不明白，可以在该主题文档中查找，通常是资源的 README 文件。只需要修改必须需改的选项。

为修改已有网站内容，用如下的步骤：

- 点击 RStudio 的 addin “Serve Site”，可以在修改过程中动态地构建网站。每次打开 RStudio 的网站项目仅需执行一次这个命令。
- 用 “New Post” addin 添加新的网页。
- 在修改某个网页时，用 “Update Metadata” addin 修改网页元数据，主要是网页标签等。

为了发布网站，只要将项目内的 `public` 目录上传到任一支持静态内容的网站服务器。如果没有合适的服务提供商，可以考虑 <https://www.netlify.com/>，可

以用 Github 账户登录并有一定的免费建站服务。如果熟悉 Github, 还可以将自己的网站源文件作为一个 Github 资源 (repository), 将 Github 资源发布到 Netlify, 就不需将 public 的内容上传到网站, 而只要有源文件就可以了, Netlify 支持 Hugo。

#### 24.3.4 网站设置文件

在项目根目录中有一个 `config.toml` 文件, 是 Hugo 软件的建站设置文件, 可以设置网站标题、描述、菜单等。一个例子如下:

```
baseurl = "/"
languageCode = "en-us"
title = "试验Blogdown网站"
theme = "hugo-lithium"
googleAnalytics = ""
disqusShortname = ""
ignoreFiles = ["\\.Rmd$", "\\.Rmarkdown$", "_files$", "_cache$"]

[permalinks]
 post = "[:year/:month/:day/:slug/"

[[menu.main]]
 name = "About"
 url = "/about/"

[[menu.main]]
 name = "Li Dongfeng Homepage"
 url = "http://www.math.pku.edu.cn/teachers/lidf/"

[params]
 description = "用Hugo和blogdown制作的网站"

 # options for highlight.js (version, additional languages, and theme)
 highlightjsVersion = "9.12.0"
 highlightjsCDN = "//cdnjs.cloudflare.com/ajax/libs"
```

```
highlightjsLang = ["r", "yaml"]
highlightjsTheme = "github"

MathJaxCDN = "//cdnjs.cloudflare.com/ajax/libs"
MathJaxVersion = "2.7.5"

path to the favicon, under "static"
favicon = "favicon.ico"

[params.logo]
url = "logo.png"
width = 50
height = 50
alt = "Logo"
```

在设置文件中，字符串要写成用双撇号界定的形式，布尔值要写成没有双撇号的 `true` 或者 `false`。用方括号包围的项（如 `[permalinks]`）会定义一个列表类型的变量，下面缩进后的内容是列表的元素。用双方括号包围的项是列表数组，即每一项都是列表，但是合并在一起由组成一个元素为列表的数组。如上面 `[[menu.main]]` 有若干项，每一项是一个主菜单项，其设置是一个列表。主菜单位置、配色等在其它地方设置，主设置文件中仅设置其文字和对应链接。

不同的网站主题需要不同的 `config.toml` 文件。

下面给出一些选项的解释。

- **baseURL**: 这是你的网页发布到网站上时，该网站给你的主页面的地址。一般是需要设置为自己发布以后的 `divide`。
- **languageCode**: 中文是 `zh-cn`，英文是 `en-us`，可能会影响一些排版规则，但是编码都是 UTF-8。
- **permalinks**: Hugo 和 blogdown 在发布 HTML 与图片等内容时需要不变的链接，默认是 `content` 中的每个 Rmd 文件变成 `public` 或 `public/post` 下的同名子目录中的 `index.html` 文件。可以用 `permalinks` 指定命名规则。`slug` 是网页元数据可以设置的一项类似于文章标识符的东西，不设置则使用文章标题。
- **publishDir**: 生成的网站存放的文件，默认为 `public` 目录。



- `theme`: 在 `themes/` 目录中保存所选主题的目录的名字。
- `ignoreFiles`: 项目中那些文件不被复制到生成结果网站中。
- `hasCJKLanguage`: 如果有大量中文、日文、韩文内容, 设置为 `true` 可以在网页统计和单词书统计时更准确。
- `[params]`: 主题特有的设置放在这里。

### 24.3.5 静态文件

项目中在 `./static/` 下的文件称为静态文件, 可以保存自己的图片、Javascript 库等, 发布时自动复制到 `./public/` 中。比如, 文件 `.static/figs/mypic01.png`, 在 Rmd 文件中可以用 `` 引用。

每个主题中一般也有一个 `static` 目录, 自己的 `static` 目录中的同名文件在发布时可以覆盖主题中的同名文件, 这样可以修改主题的一些资源或者设置。从 Rmd 编译生成的文件也可能被 `static` 中的同名文件覆盖。

还可以将网页编译为 PDF 等格式存放在 `static` 目录中。`rmarkdown::build_dir()` 命令可以将目录中所有 Rmd 文件按照其元数据的输出格式编译输出。



## Chapter 25

# 制作幻灯片

### 25.1 介绍

R Markdown 文件 (.Rmd) 文件支持多种输出，如网页 (`html_document`)、MS Word(`word_document`)、PDF(`pdf_document`，需要 LaTeX 编译器支持) 等，还支持生成网页格式的幻灯片 (`slidy_presentation`，`ioslides_presentation`)，以及 LaTeX beamer 格式的 PDF 幻灯片 (`beamer_presentation`)。

### 25.2 Slidy 幻灯片

Rmd 文件选用输出格式 `slidy_presentation` 可以生成网页格式的幻灯片，并具有缩放字体大小、显示幻灯片目录等功能。只要在.Rmd 文件开头的 YAML 元数据部分指定 `output: slidy_presentation`。因为幻灯片的单位是帧 (frame)，与论文的结构有很大区别，所以幻灯片 Rmd 文件很难同时作为论文的源文件。

### 25.2.1 文件格式

幻灯片分为多个帧，每帧用二级标题作为标志并以其为标题。二级标题就是行首以两个井号和空格开始的行，或者在标题下面画由减号组成的线的行。用一级标题作为单独的分节帧，将单独显示在一帧中。一级标题是行首以两个井号和空格开始的行，或者在标题下面画由等于号组成的线的行。

帧也可以没有标题，比如仅有照片的帧，这时，用三个或三个以上的减号连在一起标识新帧的开始。

一个简单的 `slidy_presentation` 幻灯片源文件 `example-slidy.Rmd`，内容如：

```

title: "R 幻灯片演示样例"
author: "李东风"
date: "2017-11-16"
output: slidy_presentation

用 R Markdown 的 slidy 输出作幻灯片

幻灯片结构

- 用二级标题标志一个页面开始
- 用一级标题制作单独的分节页面
- 用三个或三个以上减号标志没有标题的页面开始
- 每个页面一般用 markdown 列表显示若干个项目

幻灯片编译

- 用 RStudio 编辑
- 用 RStudio 的 Knit 按钮，选`slidy_presentation`作为输出格式

[一个演示画面](figs/demoscreen.png)
```

### 25.2.2 幻灯片编译

幻灯片用 RStudio 的 knit 按钮编译，选择输出格式为 `slidy_presentation`，结果在浏览器中播放，最好使用外部浏览器而不使用 RStudio 自带的浏览器，自带的浏览器在显示数学公式时对网络环境条件有要求而且不能完美地支持关于数学公式的本地设置。

除了使用 knit 按钮，还可以用类似如下命令：

```
rmarkdown::render("mydemo.Rmd", output_format = "slidy_presentation", encoding="UTF-8")
```

其中 `mydemo.Rmd` 是源文件。

为了制作幻灯片，最好单独设置一个 RStudio 项目，并且此项目仅生成幻灯片，而不生成普通网页、Word、PDF 等输出，否则可能造成结果混乱。希望 `rmarkdown` 包的后续版本能取消这个限制。

### 25.2.3 播放控制

播放时，用如下方式控制：

- 鼠标左键单击、光标右移键、向下翻页键、空格键都可以翻到下一页；
- 光标左移键、向上翻页键回退一页；
- 单击下方的 Contents 或单击 C 键显示幻灯片目录列表，可单击转移到任意页面；
- Home 键回到幻灯片开头；
- 用 A 键切换是否将所有页面合并显示成一个长的网页；
- 用 S 键缩写字体，用 B 键放大字体；
- 通过选择保存为 PDF 的打印机进行打印，可以将幻灯片转换为 PDF，但是打印生成的 PDF 仍是每页仅有原来的一帧，所以转换成一个单页的 HTML 再打印可能更合适。

### 25.2.4 生成单页 HTML

对 `slidy` 幻灯片的 `Rmd` 源文件不加修改，也可以通过命令直接转换为普通的单页 HTML 文件，但是在 `slidy` 幻灯片源文件中二级标题用来分帧，而普

通 Rmd 文件中二级标题用来分小节，所以生成的单页 HTML 文件会有许多小节。这样的文件更适合打印以及转换为 PDF 文件。

命令如：

```
rmarkdown::render("mydemo.Rmd", output_file="handout.html", output_format="html_do
```

### 25.2.5 数学公式处理与输出设置文件

讲课用的幻灯片经常会有数学公式，比较关键的问题是数学公式如何处理。网页中的数学公式一般使用一个公开的自由 Javascript 库 MathJax 显示，但是这个库很大，如果使用远程的库，在网络不畅通时显示公式就不正常。更好的办法是使用局部的 MathJax 库或将 MathJax 库安装在临近的网站服务器上。

为了使用局部的 MathJax 库，简单的办法是在 YAML 的 `ioslides_presentation` 项目下面指定 `mathjax: local`，如：

```

title: "R 幻灯片演示样例"
output:
 slidy_presentation:
 mathjax: local

```

上述办法容易使用，缺点是多个不同的演示项目无法共用一个局部的 MathJax 库，生成的结果包含了许多小的支持文件。

为此，可以将 MathJax 库装在演示项目所在目录的上层（比如上两层的 MathJax 目录内），将设置 MathJax 的代码放在 `_header.html` 文件中，`_header.html` 中的内容如：

```
<script type="text/x-mathjax-config">
MathJax.Hub.Config({
 jax: ["input/TeX","output/SVG"],
 extensions: ["tex2jax.js","MathMenu.js","MathZoom.js"],
 TeX: {
 extensions: ["AMSmath.js","AMSsymbols.js","noErrors.js","noUndefined.js"]
```

```
 }
 });
</script>
```

在演示项目所在目录中增加一个 `_output.yml` 文件，这使得该项目所有输出的共用的输出设置，内容如：

```
slidy_presentation:
 toc: false
 mathjax: "../../MathJax/MathJax.js"
 self_contained: false
 includes:
 in_header: "_header.html"
html_document:
 toc: true
 number_sections: false
 mathjax: "../../MathJax/MathJax.js"
 self_contained: false
 includes:
 in_header: "_header.html"
```

这里关闭了 `self_contained` 选项，设置了 MathJax 库在本地目录，具体是演示项目所在目录上面两层的 MathJax 目录中。

公式中如果有中文，开始时可能显示不正常，这时右键点击公式弹出菜单选择“Math Settings-Math Renderer”，取为“HTML-CSS”或“SVG”应可解决问题。

上面的输出设置文件中也设置了输出 `html_document`，这是单页的 HTML 格式，取消了自动章节编号，因为源文件中每帧都是一个小节。

### 25.2.6 其它选项

用 `slidy_presentation` 的 `font_adjustment: -1` 可以缩小字体一号，类似可以设为 `+1`, `-2` 等。在播放时也可以用 S 和 B 键缩小或放大显示。

可以在播放时在浏览器状态栏显示倒数计时器，用 `slidy_presentation` 的 `duration: 5` 表示每帧显示 5 分钟。这是总的预计时间，适用于演讲有时间限制的情况。

可以用 `footer` 属性指定每帧都显示的状态栏脚注，如：

```

output:
 slidy_presentation:
 font_adjustment: -1
 duration: 5
 footer: " 北京大学"

```

`slidy_presentation` 输出属性 `incremental: true` 可以使得列表显示需要每点击一次才显示下一项。

### 25.2.7 slidy 幻灯片激光笔失效问题的修改

slidy 幻灯片翻页是用空格、左右光标、上下翻页键，而一般激光笔翻页是模拟上下光标键。为此，在安装的 R 软件目录的 `library/rmarkdown/rmd/slidy/Slidy2/scripts` 子目录中，找到 `slidy.js` 文件，用编辑器打开，用编辑器的搜索功能搜索 `key == 37`，将其替换成 `key == 37 || key == 38`，这里 37 是向左光标的编码，替换后就是向左或者向上光标。用编辑器的搜索功能搜索 `key == 39`，将其替换成 `key == 39 || key == 40`，这里 39 是向右光标的编码，替换后就是向右或者向上光标。修改完毕后保存，然后将 `slidy.js` 用文件压缩程序（如 7zip）压缩为 `slidy.js.gz`。这样就可以在用 `rmarkdown` 制作的 `slidy_presentation` 结果中支持激光笔翻页了。

这样修改后，如果某帧超高，需要滚动显示，就只能通过鼠标滚轮滚动了。

## 25.3 MS PowerPoint 幻灯片

`powerpoint_presentation` 输出格式可以生成 MS PowerPoint 文件。设置如：



```

output:
 powerpoint_presentation:
 slide_level: 2

```

可以用 `powerpoint_presentation` 的 `reference_doc` 属性指向一个 .pptx 的文件，作为模板，模板中的样式将被输出结果采用。

可以用 Rstudio 的 Knit 快捷图标实现转换（选其中的 knit to PowerPoint），或者用如下命令：

```
rmarkdown::render("slides.Rmd", output_format="powerpoint_presentation", encoding="UTF-8")
```

## 25.4 Bearmer 幻灯片格式

上述 Slidy 格式的幻灯片，也可以通过 LaTeX 编译器转换成 LaTeX beamer 格式的幻灯片，设置如：

```

output:
 beamer_presentation:
 latex_engine: xelatex
 slide_level: 2
 theme: CambridgeUS
 colortheme: dolphin
 includes:
 in_hearer: "preamble.tex"

```

其中 `slide_level` 用来规定几级标题开始新的一帧。`theme` 指定一种主题，`colortheme` 指定一种配色方案，`in_header` 在 LaTeX 导言部分插入 `preamble.tex`，内容见22.12。

编译命令如：

```
rmarkdown::render("mydemo.Rmd", output_format = "beamer_presentation", encoding="U
```

结果是一个 PDF 文件。

截止到 2019 年 7 月时，编译不成功。原因在于 Rmd 到 Beamer 幻灯片的转换是预期使用 `pdflatex` 引擎的，而中文更适合用 `xelatex` 引擎，但目前的 `rmarkdown` 包中 `beamer` 幻灯片功能对 `xelatex` 引擎支持不足。

如果是纯英文内容的幻灯片，可以在上面的设置中去掉关于 `latex_engine` 的设置并在 `preamble.tex` 中去掉与中文有关的内容，则可以编译成功。

## 25.5 R Presentation 格式

R Studio 软件单独提供了对一种 R Presentation 格式的源文件的支持，以 `.Rpres` 扩展名结尾，是一种特殊的 R Markdown 文件，与 `slidy` 的源文件也类似。

`Rpres` 文件编译为 HTML 格式的幻灯片，使用 `reveal.js` 控制显示。`reveal.js` 中也有对激光笔支持不好的问题，这是因为 `reveal.js` 中用向右光标键翻页，对向下光标另有定义，激光笔一般是模拟向下和向上光标键来翻页的。为了支持激光笔，找到 RStudio 的安装目录，在 `resources/presentation/revealjs/js` 中找到 `reveal.js` 文件，在文件编辑器中打开，通过搜索找到 `case 38:`，将其剪切到 `case 33:` 后面，变成 `case 33: case 38:`。找到 `case 40:`，将其剪切到 `case 34:` 后面，变成 `case 34: case 40:`。同一目录还有一个 `reveal.min.js` 文件，也进行上述修改。

## Part V

# R 数据处理



## Chapter 26

# 数据读取技巧

### 26.1 日期数据

设文件`dates.csv`中包含如下内容，并设其文件编码为 GBK:

序号,出生日期,发病日期

```
1,1941/3/8,2007/1/1
2,1972/1/24,2007/1/1
3,1932/6/1,2007/1/1
4,1947/5/17,2007/1/1
5,1943/3/10,2007/1/1
6,1940/1/8,2007/1/1
7,1947/8/5,2007/1/1
8,2005/4/14,2007/1/1
9,1961/6/23,2007/1/2
10,1949/1/10,2007/1/2
```

先把日期当作字符串读入:

```
d.dates <- read_csv('dates.csv', locale=locale(encoding="GBK"))

Parsed with column specification:
cols(
```

```
序号 = col_double(),
出生日期 = col_character(),
发病日期 = col_character()
)
```

然后用 `lubridate::ymd()` 函数转换为 R 日期类型:

```
d.dates[[" 出生日期 ct"]] <- lubridate::ymd(
 d.dates[[" 出生日期"]])
d.dates[[" 发病日期 ct"]] <- lubridate::ymd(
 d.dates[[" 发病日期"]])
```

也可以用 R 本身的 `as.POSIXct` 函数转换:

```
d.dates[[" 出生日期 ct"]] <- as.POSIXct(
 d.dates[[" 出生日期"]], format='%Y/%m/%d', tz='Etc/GMT+8')
d.dates[[" 发病日期 ct"]] <- as.POSIXct(
 d.dates[[" 发病日期"]], format='%Y/%m/%d', tz='Etc/GMT+8')
```

经过转换后的数据为:

```
knitr::kable(d.dates)
```

序号	出生日期	发病日期	出生日期 ct	发病日期 ct
1	1941/3/8	2007/1/1	1941-03-08	2007-01-01
2	1972/1/24	2007/1/1	1972-01-24	2007-01-01
3	1932/6/1	2007/1/1	1932-06-01	2007-01-01
4	1947/5/17	2007/1/1	1947-05-17	2007-01-01
5	1943/3/10	2007/1/1	1943-03-10	2007-01-01
6	1940/1/8	2007/1/1	1940-01-08	2007-01-01
7	1947/8/5	2007/1/1	1947-08-05	2007-01-01
8	2005/4/14	2007/1/1	2005-04-14	2007-01-01
9	1961/6/23	2007/1/2	1961-06-23	2007-01-02
10	1949/1/10	2007/1/2	1949-01-10	2007-01-02

以上读入日期是比较保险的做法。还可以直接在 `read_csv()` 函数中指定某列为 `col_date()`:

```
d.dates <- read_csv(
 'dates.csv', locale=locale(encoding="GBK"),
 col_types=cols(
 `序号`=col_integer(),
 `出生日期`=col_date(format="%Y/%m/%d"),
 `发病日期`=col_date(format="%Y/%m/%d")
)
)
print(d.dates)
```

```
A tibble: 10 x 3
序号 出生日期 发病日期
<int> <date> <date>
1 1 1941-03-08 2007-01-01
2 2 1972-01-24 2007-01-01
3 3 1932-06-01 2007-01-01
4 4 1947-05-17 2007-01-01
5 5 1943-03-10 2007-01-01
6 6 1940-01-08 2007-01-01
7 7 1947-08-05 2007-01-01
8 8 2005-04-14 2007-01-01
9 9 1961-06-23 2007-01-02
10 10 1949-01-10 2007-01-02
```

### 26.1.1 日期差计算

R 的日期可以用 `difftime` 计算差值。为了计算发病时的年龄，包括小数部分，可以这样计算：

```
d.dates[, '发病年龄 (带小数年)'] <- as.numeric(
 difftime(d.dates[["发病日期"]], d.dates[["出生日期"]], units='days')/365.25)
knitr::kable(d.dates, digits=2)
```

序号	出生日期	发病日期	发病年龄（带小数年）
1	1941-03-08	2007-01-01	65.82
2	1972-01-24	2007-01-01	34.94
3	1932-06-01	2007-01-01	74.58
4	1947-05-17	2007-01-01	59.63
5	1943-03-10	2007-01-01	63.81
6	1940-01-08	2007-01-01	66.98
7	1947-08-05	2007-01-01	59.41
8	2005-04-14	2007-01-01	1.72
9	1961-06-23	2007-01-02	45.53
10	1949-01-10	2007-01-02	57.98

### 26.1.2 计算周岁

如果按照我们通常计算周岁的方法计算年龄，算法就不仅包括年的差，还要判断是否到了本年的生日。

用 lubridate 包的功能计算周岁如下：

```
age.int <- function(birth, now){
 age <- year(now) - year(birth)
 sele <- (month(now) * 100 + mday(now)
 < month(birth) * 100 + mday(birth))
 ## sele 是那些没有到生日的人
 age[sele] <- age[sele] - 1

 age
}
```

用 R 本身的功能实现周岁计算如下：

```
age.int <- function(birth, now){
 date1 <- as.POSIXlt(birth)
 date2 <- as.POSIXlt(now)
 age <- date2$year - date1$year
}
```



```

sele <- (date2$mon * 100 + date2$mday
 < date1$mon * 100 + date1$mday)
sele 是那些没有到生日的人
age[sele] <- age[sele] - 1

age
}

```

用 `d.dates()` 计算发病时周岁年龄:

```

d.dates[[" 发病年龄"]] <- age.int(d.dates[[" 出生日期"]], d.dates[[" 发病日期"]])
knitr::kable(d.dates, digits=2)

```

序号	出生日期	发病日期	发病年龄（带小数年）	发病年龄
1	1941-03-08	2007-01-01	65.82	65
2	1972-01-24	2007-01-01	34.94	34
3	1932-06-01	2007-01-01	74.58	74
4	1947-05-17	2007-01-01	59.63	59
5	1943-03-10	2007-01-01	63.81	63
6	1940-01-08	2007-01-01	66.98	66
7	1947-08-05	2007-01-01	59.41	59
8	2005-04-14	2007-01-01	1.72	1
9	1961-06-23	2007-01-02	45.53	45
10	1949-01-10	2007-01-02	57.98	57

## 26.2 缺失值处理

设有如下的`bp.csv`数据，以 GBK 编码保存:

序号,收缩压

1,145

5,110

6,未测

9,150

10,拒绝

15,115

其中的血压有非数值内容。直接用 `read.csv` 读入:

```
d.bp <- read_csv('bp.csv', locale=locale(encoding="GBK"))
```

```
Parsed with column specification:
cols(
序号 = col_double(),
收缩压 = col_character()
)
```

```
print(d.bp)
```

```
A tibble: 6 x 2
序号 收缩压
<dbl> <chr>
1 1 145
2 5 110
3 6 未测
4 9 150
5 10 拒绝
6 15 115
```

读入的收缩压被当成了字符型列，无法进行计算。

把字符型的收缩压转换为数值型:

```
d.bp[["收缩压数值"]] <- as.numeric(d.bp[["收缩压"]])
```

```
Warning: 强制改变过程中产生了NA
```

```
knitr::kable(d.bp)
```

序号	收缩压	收缩压数值
1	145	145
5	110	110
6	未测	NA
9	150	150
10	拒绝	NA
15	115	115

收缩压中非数值的项被转换为数值型缺失值，并在转换时有警告信息。注意这里同时保存了原始输入的收缩压和转换为数值的收缩压，这样便于数据核对。

## 26.3 练习

- 读入 `patients.csv`，把其中的日期转换为 R 日期类型，计算发病年龄。保存为 tibble 数据框 `d.patients`。
- 读入 `cancer.csv`，保存为 `d.cancer`。v0 是放疗前肿瘤体积，v1 是放疗后肿瘤体积。计算放疗后肿瘤缩减率。



# Chapter 27

## 数据整理

### 27.1 tidyverse 系统

假设数据以 tibble 格式保存 (tibble 是数据框类型的改进, readr 包的 read\_csv() 会生成此类)。数据集经常需要选行子集、选列子集、排序、定义新变量、横向合并等操作,而且经常会用若干个连续的操作分步处理, magrittr 包的管道运算符 %>% 特别适用于这种分步处理。dplyr 包和 tidyr 包定义了一系列“动词”,可以用比较自然的方式进行数据整理。

为了使用这些功能,可以载入 tidyverse 包,则 magrittr 包, readr 包, dplyr 包和 tidyr 包都会被自动载入:

```
library(tidyverse)
```

下面的例子中用如下的一个班的学生数据作为例子,保存在如下 class.csv 文件中:

```
name,sex,age,height,weight
Alice,F,13,56.5,84
Becka,F,13,65.3,98
Gail,F,14,64.3,90
Karen,F,12,56.3,77
Kathy,F,12,59.8,84.5
```

```
Mary,F,15,66.5,112
Sandy,F,11,51.3,50.5
Sharon,F,15,62.5,112.5
Tammy,F,14,62.8,102.5
Alfred,M,14,69,112.5
Duke,M,14,63.5,102.5
Guido,M,15,67,133
James,M,12,57.3,83
Jeffrey,M,13,62.5,84
John,M,12,59,99.5
Philip,M,16,72,150
Robert,M,12,64.8,128
Thomas,M,11,57.5,85
William,M,15,66.5,112
```

读入为 tibble:

```
d.class <- read_csv(
 "class.csv",
 col_types=cols(
 .default = col_double(),
 name=col_character(),
 sex=col_factor(levels=c("M", "F"))
))
```

R 的 NHANES 扩展包提供了一个规模更大的示例数据框 NHANES，可以看作是美國扣除住院病人以外的人群的一个随机样本，有 10000 个观测，有 76 个变量，主题是个人健康与营养方面的信息。仅作为教学使用而不足以作为严谨的科研用数据。原始数据的情况详见<http://www.cdc.gov/nchs/nhanes.htm>。载入 NHANES 数据框：

```
library(NHANES)
data(NHANES)
```

```
print(dim(NHANES))
```

```
[1] 10000 76
```

```
print(names(NHANES))

[1] "ID" "SurveyYr" "Gender"
[4] "Age" "AgeDecade" "AgeMonths"
[7] "Race1" "Race3" "Education"
[10] "MaritalStatus" "HHIncome" "HHIncomeMid"
[13] "Poverty" "HomeRooms" "HomeOwn"
[16] "Work" "Weight" "Length"
[19] "HeadCirc" "Height" "BMI"
[22] "BMICatUnder20yrs" "BMI_WHO" "Pulse"
[25] "BPSysAve" "BPDiaAve" "BPSys1"
[28] "BPDia1" "BPSys2" "BPDia2"
[31] "BPSys3" "BPDia3" "Testosterone"
[34] "DirectChol" "TotChol" "UrineVol1"
[37] "UrineFlow1" "UrineVol2" "UrineFlow2"
[40] "Diabetes" "DiabetesAge" "HealthGen"
[43] "DaysPhysHlthBad" "DaysMentHlthBad" "LittleInterest"
[46] "Depressed" "nPregnancies" "nBabies"
[49] "Age1stBaby" "SleepHrsNight" "SleepTrouble"
[52] "PhysActive" "PhysActiveDays" "TVHrsDay"
[55] "CompHrsDay" "TVHrsDayChild" "CompHrsDayChild"
[58] "Alcohol12PlusYr" "AlcoholDay" "AlcoholYear"
[61] "SmokeNow" "Smoke100" "Smoke100n"
[64] "SmokeAge" "Marijuana" "AgeFirstMarij"
[67] "RegularMarij" "AgeRegMarij" "HardDrugs"
[70] "SexEver" "SexAge" "SexNumPartnLife"
[73] "SexNumPartYear" "SameSex" "SexOrientation"
[76] "PregnantNow"
```

变量 ID 是受试者编号，SurveyYr 是调查年份，同一受试者可能在多个调查年份中有数据。变量中包括性别、年龄、种族、收入等人口学数据，包括体重、身高、脉搏、血压等基本体检数据，以及是否糖尿病、是否抑郁、是否怀孕、已生产子女数等更详细的健康数据，运动习惯、饮酒、性生活等行为方面的数据。这个教学用数据集最初的使用者是 Cashmere 高中的 Michelle Dalrymple 和

新西兰奥克兰大学的 Chris Wild。

## 27.2 用 `filter()` 选择行子集

数据框的任何行子集仍为数据框，即使只有一行而且都是数值也是如此。行子集可以用行下标选取，如 `d.class[8:12,]`。函数 `head()` 取出数据框的前面若干行，`tail()` 取出数据框的最后若干行。

`dplyr` 包的 `filter()` 函数可以按条件选出符合条件的行组成的子集。下例从 `d.class` 中选出年龄在 13 岁和 13 岁以下的女生：

```
d.class %>%
 filter(sex=="F", age<=13)

A tibble: 5 x 5
name sex age height weight
<chr> <fct> <dbl> <dbl> <dbl>
1 Alice F 13 56.5 84
2 Becka F 13 65.3 98
3 Karen F 12 56.3 77
4 Kathy F 12 59.8 84.5
5 Sandy F 11 51.3 50.5
```

`filter()` 函数第一个参数是要选择的数据框，后续的参数是条件，这些条件是需要同时满足的，另外，条件中取缺失值的观测自动放弃。`filter()` 会自动舍弃行名，如果需要行名只能将其转换成数据框的一列。`filter()` 的结果为行子集数据框。用在管道操作当中的时候第一自变量省略（是管道传递下来的）。

函数 `head(x, n)` 可以用来选择数据框前面若干行，`tail(x, n)` 可以用来选择数据框后面若干行，如：

```
d.class %>%
 head(n=5)

A tibble: 5 x 5
name sex age height weight
<chr> <fct> <dbl> <dbl> <dbl>
```



```
1 Alice F 13 56.5 84
2 Becca F 13 65.3 98
3 Gail F 14 64.3 90
4 Karen F 12 56.3 77
5 Kathy F 12 59.8 84.5
```

`dplyr` 包的函数 `slice(.data, ...)` 可以用来选择指定序号的行子集，正的序号表示保留，负的序号表示排除。如：

```
d.class %>%
 slice(3:5)
```

```
A tibble: 3 x 5
name sex age height weight
<chr> <fct> <dbl> <dbl> <dbl>
1 Gail F 14 64.3 90
2 Karen F 12 56.3 77
3 Kathy F 12 59.8 84.5
```

## 27.3 用 `sample_n()` 对观测随机抽样

`dplyr` 包的 `sample_n(tbl, size)` 函数可以从数据集 `tbl` 中随机无放回抽取 `size` 行，如：

```
d.class %>%
 sample_n(size = 3)
```

```
A tibble: 3 x 5
name sex age height weight
<chr> <fct> <dbl> <dbl> <dbl>
1 Alfred M 14 69 112.
2 John M 12 59 99.5
3 Guido M 15 67 133
```

`sample_n()` 中加选项 `replace=TRUE` 可以变成有放回抽样。可以用 `weight` 选项指定数据框中的一列作为抽样权重，进行不等概抽样。

## 27.4 用 `distinct()` 去除重复行

有时我们希望得到一个或若干个变量组合的所有不同值。`dplyr` 包的 `distinct()` 函数可以对数据框指定若干变量，然后筛选出所有不同值，每组不同值仅保留一行。指定变量名时是写成字符串形式而是直接写变量名，这是 `dplyr` 和 `tidyr` 包的特点。例如，筛选出性别与年龄的所有不同组合：

```
d.class %>%
 distinct(sex, age)

A tibble: 11 x 2
sex age
<fct> <dbl>
1 F 13
2 F 14
3 F 12
4 F 15
5 F 11
6 M 14
7 M 15
8 M 12
9 M 13
10 M 16
11 M 11
```

如果希望保留数据框中其它变量，可以加选项 `.keep_all=TRUE`。

下面的程序查看 NHANES 数据框中 ID 与 SurveyYr 的组合的不同值的个数：

```
NHANES %>%
 distinct(ID, SurveyYr) %>%
 nrow()

[1] 6779
```

这个结果提示有些人在某一调查年中有多个观测。

## 27.5 用 `drop_na()` 去除指定的变量有缺失值的行

在进行统计建模时，通常需要用到的因变量和自变量都不包含缺失值。tidyr 包的 `drop_na()` 函数可以对数据框指定一到多个变量，删去指定的变量有缺失值的行。不指定变量时有任何变量缺失的行都会被删去。

用如

```
d.class %>%
 drop_na(age, height, weight)
```

又如，将 NHANES 中所有存在缺失值的行删去后数出保留的行数，原来有 10000 行：

```
NHANES %>%
 drop_na() %>%
 nrow()
```

```
[1] 0
```

可见所有行都有缺失值。下面仅剔除 AlcoholDay 缺失的观测并计数：

```
NHANES %>%
 drop_na(AlcoholDay) %>%
 nrow()
```

```
[1] 4914
```

## 27.6 用 `select()` 选择列子集

dplyr 包的 `select()` 选择列子集，并返回列子集结果。

可以指定变量名，如

```
d.class %>%
 select(name, age)
```

```
A tibble: 19 x 2
```

```
name age
```

```
<chr> <dbl>
1 Alice 13
2 Becca 13
3 Gail 14
4 Karen 12
5 Kathy 12
6 Mary 15
7 Sandy 11
8 Sharon 15
9 Tammy 14
10 Alfred 14
11 Duke 14
12 Guido 15
13 James 12
14 Jeffrey 13
15 John 12
16 Philip 16
17 Robert 12
18 Thomas 11
19 William 15
```

可以用冒号表示列范围，如

```
d.class %>%
 select(age:weight)
```

```
A tibble: 19 x 3
age height weight
<dbl> <dbl> <dbl>
1 13 56.5 84
2 13 65.3 98
3 14 64.3 90
4 12 56.3 77
5 12 59.8 84.5
6 15 66.5 112
```

```
7 11 51.3 50.5
8 15 62.5 112.
9 14 62.8 102.
10 14 69 112.
11 14 63.5 102.
12 15 67 133
13 12 57.3 83
14 13 62.5 84
15 12 59 99.5
16 16 72 150
17 12 64.8 128
18 11 57.5 85
19 15 66.5 112
```

可以用数字序号表示列范围，如

```
d.class %>%
 select(3:5)
```

```
A tibble: 19 x 3
age height weight
<dbl> <dbl> <dbl>
1 13 56.5 84
2 13 65.3 98
3 14 64.3 90
4 12 56.3 77
5 12 59.8 84.5
6 15 66.5 112
7 11 51.3 50.5
8 15 62.5 112.
9 14 62.8 102.
10 14 69 112.
11 14 63.5 102.
12 15 67 133
13 12 57.3 83
```

```
14 13 62.5 84
15 12 59 99.5
16 16 72 150
17 12 64.8 128
18 11 57.5 85
19 15 66.5 112
```

参数中前面写负号表示扣除，如

```
d.class %>%
 select(-name, -age)
```

```
A tibble: 19 x 3
sex height weight
<fct> <dbl> <dbl>
1 F 56.5 84
2 F 65.3 98
3 F 64.3 90
4 F 56.3 77
5 F 59.8 84.5
6 F 66.5 112
7 F 51.3 50.5
8 F 62.5 112.
9 F 62.8 102.
10 M 69 112.
11 M 63.5 102.
12 M 67 133
13 M 57.3 83
14 M 62.5 84
15 M 59 99.5
16 M 72 150
17 M 64.8 128
18 M 57.5 85
19 M 66.5 112
```

如果要选择的变量名已经保存为一个字符型向量，可以用 `one_of()` 函数引入，

如

```
vars <- c("name", "sex")
d.class %>%
 select(one_of(vars))
```

```
A tibble: 19 x 2
name sex
<chr> <fct>
1 Alice F
2 Becca F
3 Gail F
4 Karen F
5 Kathy F
6 Mary F
7 Sandy F
8 Sharon F
9 Tammy F
10 Alfred M
11 Duke M
12 Guido M
13 James M
14 Jeffrey M
15 John M
16 Philip M
17 Robert M
18 Thomas M
19 William M
```

`select()` 有若干个配套函数可以按名字的模式选择变量列，如

- `starts_with("se")`: 选择名字以“se”开头的变量列；
- `ends_with("ght")`: 选择名字以“ght”结尾的变量列；
- `contains("no")`: 选择名字中含有子串“no”的变量列；
- `matches("^[[:alpha:]]+[[:digit:]]+$")`，选择列名匹配某个正则表达式模式的变量列，这里匹配前一部分是字母，后一部分是数字的变量

名。

- `num_range("x", 1:3)`, 选择 `x1`, `x2`, `x3`。
- `everything()`: 代指所有选中的变量, 这可以用来将指定的变量次序提前, 其它变量排在后面。

R 的字符串函数 (如 `paste()`) 和正则表达式函数可以用来生成变量名子集。

R 函数 `subset` 也能对数据框选取列子集和行子集。

如果需要选择单个变量并使得结果为普通向量, 可以用 `dplyr` 包的 `pull()` 函数, 如:

```
d.class %>% pull(name)
```

```
[1] "Alice" "Becka" "Gail" "Karen" "Kathy" "Mary" "Sandy"
[8] "Sharon" "Tammy" "Alfred" "Duke" "Guido" "James" "Jeffrey"
[15] "John" "Philip" "Robert" "Thomas" "William"
```

`pull()` 可以指定单个变量名, 也可以指定变量序号, 负的变量序号从最后一个变量数起。缺省取出最后一个变量。

## 27.7 用 `arrange()` 排序

`dplyr` 包的 `arrange()` 按照数据框的某一行或某几列排序, 返回排序后的结果, 如

```
d.class %>%
 arrange(sex, age)
```

```
A tibble: 19 x 5
name sex age height weight
<chr> <fct> <dbl> <dbl> <dbl>
1 Thomas M 11 57.5 85
2 James M 12 57.3 83
3 John M 12 59 99.5
4 Robert M 12 64.8 128
5 Jeffrey M 13 62.5 84
```



```
6 Alfred M 14 69 112.
7 Duke M 14 63.5 102.
8 Guido M 15 67 133
9 William M 15 66.5 112
10 Philip M 16 72 150
11 Sandy F 11 51.3 50.5
12 Karen F 12 56.3 77
13 Kathy F 12 59.8 84.5
14 Alice F 13 56.5 84
15 Becca F 13 65.3 98
16 Gail F 14 64.3 90
17 Tammy F 14 62.8 102.
18 Mary F 15 66.5 112
19 Sharon F 15 62.5 112.
```

用 `desc()` 包裹想要降序排列的变量，如

```
d.class %>%
 arrange(sex, desc(age))
```

```
A tibble: 19 x 5
name sex age height weight
<chr> <fct> <dbl> <dbl> <dbl>
1 Philip M 16 72 150
2 Guido M 15 67 133
3 William M 15 66.5 112
4 Alfred M 14 69 112.
5 Duke M 14 63.5 102.
6 Jeffrey M 13 62.5 84
7 James M 12 57.3 83
8 John M 12 59 99.5
9 Robert M 12 64.8 128
10 Thomas M 11 57.5 85
11 Mary F 15 66.5 112
12 Sharon F 15 62.5 112.
```

```
13 Gail F 14 64.3 90
14 Tammy F 14 62.8 102.
15 Alice F 13 56.5 84
16 Becka F 13 65.3 98
17 Karen F 12 56.3 77
18 Kathy F 12 59.8 84.5
19 Sandy F 11 51.3 50.5
```

排序时不论升序还是降序，所有的缺失值都自动排到末尾。

R 函数 `order()` 可以用来给出数据框的排序次序从而将数据框排序。

## 27.8 用 `rename()` 修改变量名

在 `dplyr` 包的 `rename()` 中用“新名字 = 旧名字”格式修改变量名，如

```
d2.class <- d.class %>%
 rename(h=height, w=weight)
```

注意这样改名字不是对原始数据框修改而是返回改了名字后的新数据框。

## 27.9 用 `mutate()` 计算新变量

`dplyr` 包的 `mutate()` 可以为数据框计算新变量，返回含有新变量以及原变量的数据框。如

```
d.class %>%
 mutate(
 rwh=weight/height,
 sexc=ifelse(sex=="F", "女", "男"))

A tibble: 19 x 7
name sex age height weight rwh sexc
<chr> <fct> <dbl> <dbl> <dbl> <dbl> <chr>
1 Alice F 13 56.5 84 1.49 女
```

```
2 Becca F 13 65.3 98 1.50 女
3 Gail F 14 64.3 90 1.40 女
4 Karen F 12 56.3 77 1.37 女
5 Kathy F 12 59.8 84.5 1.41 女
6 Mary F 15 66.5 112 1.68 女
7 Sandy F 11 51.3 50.5 0.984 女
8 Sharon F 15 62.5 112. 1.8 女
9 Tammy F 14 62.8 102. 1.63 女
10 Alfred M 14 69 112. 1.63 男
11 Duke M 14 63.5 102. 1.61 男
12 Guido M 15 67 133 1.99 男
13 James M 12 57.3 83 1.45 男
14 Jeffrey M 13 62.5 84 1.34 男
15 John M 12 59 99.5 1.69 男
16 Philip M 16 72 150 2.08 男
17 Robert M 12 64.8 128 1.98 男
18 Thomas M 11 57.5 85 1.48 男
19 William M 15 66.5 112 1.68 男
```

用 `mutate()` 计算新变量时如果计算比较复杂，也可以用多个语句组成复合语句，如：

```
d.class %>%
 mutate(
 sexc = {
 x <- rep("男", length(sex))
 x[sex == "F"] <- "女"
 x
 }
)
```

```
A tibble: 19 x 6
name sex age height weight sexc
<chr> <fct> <dbl> <dbl> <dbl> <chr>
1 Alice F 13 56.5 84 女
```

```
2 Becca F 13 65.3 98 女
3 Gail F 14 64.3 90 女
4 Karen F 12 56.3 77 女
5 Kathy F 12 59.8 84.5 女
6 Mary F 15 66.5 112 女
7 Sandy F 11 51.3 50.5 女
8 Sharon F 15 62.5 112. 女
9 Tammy F 14 62.8 102. 女
10 Alfred M 14 69 112. 男
11 Duke M 14 63.5 102. 男
12 Guido M 15 67 133 男
13 James M 12 57.3 83 男
14 Jeffrey M 13 62.5 84 男
15 John M 12 59 99.5 男
16 Philip M 16 72 150 男
17 Robert M 12 64.8 128 男
18 Thomas M 11 57.5 85 男
19 William M 15 66.5 112 男
```

注意这样生成新变量不是在原来的数据框中添加，原来的数据框没有被修改，而是返回添加了新变量的新数据框。R 软件的巧妙设计保证了这样虽然是生成了新数据框，但是与原来数据框重复的列并不会重复保存。

计算公式中可以包含对数据框中变量的统计函数结果，如

```
d.class %>%
 mutate(
 cheight = height - mean(height))
```

新变量可以与老变量名相同，这样就在输出中修改了老变量。

函数 `transmute()` 用法与 `mutate()` 类似，但是仅保留新定义的变量，不保留原来的所有变量。如：

```
d.class %>%
 transmute(
 stdh = scale(height),
```

```
stdw = scale(weight))
```

可见结果中仅保留了新定义的变量。

定义新变量也可以直接为数据框的新变量赋值：

```
d.class[["rwh"]] <- d.class[["weight"]] / d.class[["height"]]
```

这样的做法与 `mutate()` 的区别是这样不会生成新数据框，新变量是在原数据框中增加的。

给数据框中某个变量赋值为 `NULL` 可以修改数据框，从数据框中删去该变量。

## 27.10 用管道连接多次操作

管道运算符特别适用于对同一数据集进行多次操作。例如，对 `t.class` 数据，先选出所有女生，再去掉性别和 `age` 变量：

```
d.class %>%
 filter(sex=="F") %>%
 select(-sex, -age)

A tibble: 9 x 4
name height weight rwh
<chr> <dbl> <dbl> <dbl>
1 Alice 56.5 84 1.49
2 Becka 65.3 98 1.50
3 Gail 64.3 90 1.40
4 Karen 56.3 77 1.37
5 Kathy 59.8 84.5 1.41
6 Mary 66.5 112 1.68
7 Sandy 51.3 50.5 0.984
8 Sharon 62.5 112. 1.8
9 Tammy 62.8 102. 1.63
```

管道操作的结果可以保存为新的 `tibble`，如：

```
class_F <- d.class %>%
 filter(sex=="F") %>%
 select(-sex, -age)
```

## 27.11 数据简单汇总

dplyr 包的 `summarise()` 函数可以对数据框计算统计量。这个函数针对少量变量时很方便，有大量变量需要对变量统一处理时不太方便。

以肺癌病人化疗数据 `cancer.csv` 为例，有 34 个肺癌病人的数据：

```
d.cancer <- read_csv(
 "cancer.csv", locale=locale(encoding="GBK"))
```

```
Parsed with column specification:
cols(
id = col_double(),
age = col_double(),
sex = col_character(),
type = col_character(),
v0 = col_double(),
v1 = col_double()
)
```

```
d.cancer
```

```
A tibble: 34 x 6
id age sex type v0 v1
<dbl> <dbl> <chr> <chr> <dbl> <dbl>
1 1 70 F 腺癌 26.5 2.91
2 2 70 F 腺癌 135. 35.1
3 3 69 F 腺癌 210. 74.4
4 4 68 M 腺癌 61 35.0
5 5 67 M 鳞癌 238. 128.
6 6 75 F 腺癌 330. 112.
```

```
7 7 52 M 鳞癌 105. 32.1
8 8 71 M 鳞癌 85.2 29.2
9 9 68 M 鳞癌 102. 22.2
10 10 79 M 鳞癌 65.5 21.9
... with 24 more rows
```

求年龄 (age) 的平均值、标准差:

```
d.cancer %>%
 summarise(mean.age=mean(age, na.rm=TRUE),
 sd.age=sd(age, na.rm=TRUE))
```

```
A tibble: 1 x 2
mean.age sd.age
<dbl> <dbl>
1 64.1 9.16
```

更重要的是分组汇总。dplyr 包的 `group_by()` 函数对数据框 (或 tibble) 分组, 随后的 `summarise()` 将按照分组汇总。比如, 按不同性别计算人数与年龄平均值:

```
d.cancer %>%
 group_by(sex) %>%
 summarise(count=n(), mean.age=mean(age, na.rm=TRUE))
```

```
A tibble: 2 x 3
sex count mean.age
<chr> <int> <dbl>
1 F 13 66.1
2 M 21 63.2
```

其中 `n()` 计算某类的观测数 (行数)。为了计算某列的非缺失值个数, 用 `sum(!is.na(x))`。

常用的汇总函数有:

- 位置度量: `mean()`, `median()`。
- 分散程度 (变异性) 度量: `sd()`, `IQR()`, `mad()`。
- 分位数: `min()`, `max()`, `quantile()`。

- 按下标查询，如 `first(x)` 取出 `x[1]`，`last(x)` 取出 `x` 的最后一个元素，`nth(x,2)` 取出 `x[2]`。可以提供一个缺省值以防某个下标位置不存在。
- 计数：`n()` 给出某个组的观测数，`sum(!is.na(x))` 统计 `x` 的非缺失值个数，`n_distinct(x)` 统计 `x` 的不同值个数（缺失值也算一个值）。`count(x)` 给出 `x` 的每个不同值的个数（类似于 `table()` 函数）。

这里有些函数是 `dplyr` 包提供的，仅适用于 `tibble` 类型。

下面的程序对 `d.cancer` 数据框分性别与病理类型分别统计人数：

```
d.cancer %>%
 group_by(sex, type) %>%
 summarise(freq=n())
```

```
A tibble: 4 x 3
Groups: sex [2]
sex type freq
<chr> <chr> <int>
1 F 鳞癌 4
2 F 腺癌 9
3 M 鳞癌 18
4 M 腺癌 3
```

在有交叉分类时，`n()`、`mean()` 这些汇总函数是针对最内层进行汇总，汇总后最内层的分类变量不再当作分类。比如，上述结果仍按照 `sex` 分类，但不再按照 `type` 分类。

事实上，不需要用 `group_by()`，交叉分类计算频数可以用 `dplyr` 的 `count()` 函数，如：

```
d.cancer %>%
 count(sex, type)
```

```
A tibble: 4 x 3
sex type n
<chr> <chr> <int>
1 F 鳞癌 4
```



```
2 F 腺癌 9
3 M 鳞癌 18
4 M 腺癌 3
```

又如，数出 NHANES 数据框中 ID 与 SurveyYr 每一对组合的出现次数，筛选出二次及以上者，并降序排列，仅显示前 100 行结果：

```
NHANES %>%
 count(ID, SurveyYr) %>%
 filter(n >=2) %>%
 arrange(desc(n)) %>%
 head(100)
```

```
A tibble: 100 x 3
ID SurveyYr n
<int> <fct> <int>
1 70324 2011_12 8
2 62927 2011_12 7
3 63297 2011_12 7
4 69626 2011_12 7
5 60566 2009_10 6
6 61442 2009_10 6
7 63163 2011_12 6
8 63330 2011_12 6
9 63390 2011_12 6
10 63744 2011_12 6
... with 90 more rows
```

用 `group_by()` 分组后除了可以分组汇总，还可以分组筛选：

```
d.cancer %>%
 group_by(sex) %>%
 filter(rank(desc(v0)) <= 2) %>%
 arrange(sex, desc(v0))
```

```
A tibble: 4 x 6
Groups: sex [2]
```

```
id age sex type v0 v1
<dbl> <dbl> <chr> <chr> <dbl> <dbl>
1 6 75 F 腺癌 330. 112.
2 25 NA F 鳞癌 223 25.6
3 5 67 M 鳞癌 238. 128.
4 16 76 M 鳞癌 231. 113.
```

以上程序按性别分组后，在每组中找出疗前体积排名在前两名的。

在分组后也可以根据每组的统计量用 `mutate()` 定义新变量。

用 `group_by()` 分组汇总后的结果不是普通的 tibble，总是带有分组信息。这在后续的使用中可能会产生问题，为此，可以用 `ungroup()` 函数取消分组。例如

```
d.cancer %>%
 group_by(sex, type) %>%
 summarise(freq=n()) %>%
 summarise(ntotal=sum(freq))
```

```
A tibble: 2 x 2
sex ntotal
<chr> <int>
1 F 13
2 M 21
```

可以看出并没有能够通过男、女分别的人数计算总人数。加入 `ungroup()`：

```
d.cancer %>%
 group_by(sex, type) %>%
 summarise(freq=n()) %>%
 ungroup() %>%
 summarise(ntotal=sum(freq))
```

```
A tibble: 1 x 1
ntotal
<int>
1 34
```

得到了需要的结果。

## 27.12 长宽表转换

考虑如下的宽表，保存在 CSV 文件 `widetab.csv` 中：

```
"subject","x_1","x_2","x_3","y_1","y_2","y_3"
1,5,7,8,9,7,6
2,8,2,10,1,1,9
3,7,2,5,10,8,3
4,1,5,6,10,1,1
5,9,7,10,8,8,10
```

这个数据是 5 名病人 3 次检查的记录，每次检查有 x 和 y 两个测量项目。每个病人的所有各次检查以及所有测量项目都在同一行，称这样的表为宽表。读入为

```
d.wide <- read_csv("widetab.csv")
```

```
Parsed with column specification:
cols(
subject = col_double(),
x_1 = col_double(),
x_2 = col_double(),
x_3 = col_double(),
y_1 = col_double(),
y_2 = col_double(),
y_3 = col_double()
)
```

现在希望将其中的时间分离出来，不同时间变成不同的观测，每个病人的三次检查转换成三个观测。tidyr 包的 `gather()` 函数将宽表变成长表，但是将不同测量项目也转换到了不同的行：

```
d.wide %>%
 gather(x_1, x_2, x_3, y_1, y_2, y_3, key="variable", value="value")
```

```
A tibble: 30 x 3
subject variable value
<dbl> <chr> <dbl>
1 1 1 x_1 5
2 2 2 x_1 8
3 3 3 x_1 7
4 4 4 x_1 1
5 5 5 x_1 9
6 1 1 x_2 7
7 2 2 x_2 2
8 3 3 x_2 2
9 4 4 x_2 5
10 5 5 x_2 7
... with 20 more rows
```

tidyr 包的 `separate()` 可以用来帮助拆分 `x_1`, `y_1` 这样的名字, 这样可以将测量项目名称与时间分离开来, 如:

```
d.wide %>%
 gather(x_1, x_2, x_3, y_1, y_2, y_3, key="variable", value="value") %>%
 separate(variable, into=c("variable", "time"), sep="_")
```

```
A tibble: 30 x 4
subject variable time value
<dbl> <chr> <chr> <dbl>
1 1 1 x 1 5
2 2 2 x 1 8
3 3 3 x 1 7
4 4 4 x 1 1
5 5 5 x 1 9
6 1 1 x 2 7
7 2 2 x 2 2
8 3 3 x 2 2
9 4 4 x 2 5
10 5 5 x 2 7
```

```
... with 20 more rows
```

tidyr 包的函数 `spread()` 可以将用变量名和变量值分别存储的变量，恢复为每个变量一列的形式，如：

```
d.wide %>%
 gather(x_1, x_2, x_3, y_1, y_2, y_3, key="variable", value="value") %>%
 separate(variable, into=c("variable", "time"), sep="_") %>%
 spread(key=variable, value=value)
```

```
A tibble: 15 x 4
subject time x y
<dbl> <chr> <dbl> <dbl>
1 1 1 1 5 9
2 2 1 2 7 7
3 3 1 3 8 6
4 4 2 1 8 1
5 5 2 2 2 1
6 6 2 3 10 9
7 7 3 1 7 10
8 8 3 2 2 8
9 9 3 3 5 3
10 4 1 1 10
11 4 2 5 1
12 4 3 6 1
13 5 1 9 8
14 5 2 7 8
15 5 3 10 10
```

这个结果已经变成每个病人每次检查为一个观测（记录）的形式。

如果变量名不像 `x_1`, `y_3` 这样整齐，为了拆分变量名与时间，也可以用字符串函数。比如，我们将 `d.wide` 中变量 `x_1` 改名为 `x1`, `x_3` 改名为 `x3`, `y_1` 改名为 `abc1`，删去其它变量：

```
d.wide %>%
 select(subject, x_1, x_3, y_1) %>%
```

```
rename(x1=x_1, x3=x_3, abc1=y_1)
```

```
A tibble: 5 x 4
subject x1 x3 abc1
<dbl> <dbl> <dbl> <dbl>
1 1 5 8 9
2 2 8 10 1
3 3 7 5 10
4 4 1 6 10
5 5 9 10 8
```

要将上述数据的变量与时间信息分离，因为没有变量名与时间之间没有明确的分隔符也没有固定宽度，`separate()` 函数难以分离这样的变量名与时间，可以利用正则表达式：

```
d.wide %>%
```

```
 select(subject, x_1, x_3, y_1) %>%
 rename(x1=x_1, x3=x_3, abc1=y_1) %>%
 gather(x1, x3, abc1, key="variable", value="value") %>%
 mutate(time=as.numeric(gsub("^([[:alpha:]]+)([[:digit:]]+)$", "\\2", variable)),
 variable=gsub("^([[:alpha:]]+)([[:digit:]]+)$", "\\1", variable)) %>%
 spread(key=variable, value=value)
```

```
A tibble: 10 x 4
subject time abc x
<dbl> <dbl> <dbl> <dbl>
1 1 1 9 5
2 1 3 NA 8
3 2 1 1 8
4 2 3 NA 10
5 3 1 10 7
6 3 3 NA 5
7 4 1 10 1
8 4 3 NA 6
9 5 1 8 9
```

```
10 5 3 NA 10
```

tidyr 包还有一些方便函数。长宽表转换问题中的某些数据是缺失的，这些缺失值在某些形式中可见，如上表中变量 abc 在 time=3 时的值，在某些形式中根本没有表现出来。在用 gather() 将宽表变长表时，可以加 na.rm=TRUE 选项将不必要的缺失值观测删去。函数 complete() 可以指定若干列，使得这些列的所有不同组合均出现。有时某个单元格缺失是表示该值等于其上一行的值，这时可以用函数 fill() 指定该列按照这样的规则填充值。

## 27.13 拆分数数据列

有时应该放在不同列的数据用分隔符分隔后放在同一列中了。比如，下面数据集中“succ/total”列存放了用“/”分隔开的成功数与试验数：

```
d.sep <- read_csv(
 "testdata, succ/total
 1, 1/10
 2, 3/5
 3, 2/8
 ")
d.sep
```

```
A tibble: 3 x 2
testid `succ/total`
<dbl> <chr>
1 1 1 1/10
2 2 2 3/5
3 3 3 2/8
```

用 tidyr::separate() 可以将这样的列拆分为各自的变量列，如

```
d.sep %>%
 separate(`succ/total`, into=c("succ", "total"),
 sep="/", convert=TRUE)
```

```
A tibble: 3 x 3
```

```
testid succ total
<dbl> <int> <int>
1 1 1 10
2 2 3 5
3 3 2 8
```

其中 `into` 指定拆分后新变量名，`sep` 指定分隔符，`convert=TRUE` 要求自动将分割后的值转换为适当的类型。`sep` 还可以指定取子串的字符位置，按位置拆分各个子串。

选项 `extra` 指出拆分时有多余内容的处理方法，选项 `fill` 指出有不足内容的处理方法。

函数 `extract()` 可以按照某种正则表达式表示的模式从指定列拆分出对应于正则表达式中捕获组的一列或多列内容。

## 27.14 合并数据列

`tidyr::unite()` 函数可以将同一行的两列或多列的内容合并成一列。这是 `separate()` 的反向操作，如：

```
d.sep %>%
 separate(`succ/total`, into=c("succ", "total"),
 sep="/", convert=TRUE) %>%
 unite(ratio, succ, total, sep=":")
```

```
A tibble: 3 x 2
testid ratio
<dbl> <chr>
1 1 1:10
2 2 3:5
3 3 2:8
```

`unite()` 的第一个参数是要修改的数据框，这里用管道`%>%`传递进来，第二个参数是合并后的变量名（`ratio` 变量），其它参数是要合并的变量名，`sep` 指定分隔符。实际上用 `mutate()`、`paste()` 或者 `sprintf()` 也能完成合并。



## 27.15 横向合并

实际数据往往没有存放在单一的表中，需要从多个表查找数据。多个表之间的连接，一般靠关键列（key）对准来连接。连接可以是一对一的，一对多的。多对多连接应用较少，因为多对多连接是所有两两组合。

在规范的数据库中，每个表都应该有主键，这可以是一列，也可以是多列的组合。为了确定某列是主键，可以用 `count()` 和 `filter()`，如

```
d.class %>%
 count(name) %>%
 filter(n>1)
```

```
A tibble: 0 x 2
... with 2 variables: name <chr>, n <int>
```

没有发现重复出现的 `name`，说明 `d.class` 中 `name` 可以作为主键。

为了演示一对一的横向连接，我们将 `d.class` 拆分为两个数据集 `d1.class` 和 `d2.class`，两个数据集都有主键 `name`，`d1.class` 包含变量 `name`, `sex`，`d2.class` 包含变量 `name`, `age`, `height`, `weight`，并删去某些观测：

```
d1.class <- d.class %>%
 select(name, sex) %>%
 filter(!(name %in% "Becka"))
d2.class <- d.class %>%
 select(name, age, height, weight)
```

用 `dplyr` 包的 `inner_join()` 函数将两个数据框按键值横向合并，仅保留能匹配的观测。因为 `d1.class` 中丢失了 `Becka` 的观测，所以合并后的数据框中也没有 `Becka` 的观测：

```
d1.class %>%
 inner_join(d2.class)

Joining, by = "name"

A tibble: 18 x 5
name sex age height weight
```

```
<chr> <fct> <dbl> <dbl> <dbl>
1 Alice F 13 56.5 84
2 Gail F 14 64.3 90
3 Karen F 12 56.3 77
4 Kathy F 12 59.8 84.5
5 Mary F 15 66.5 112
6 Sandy F 11 51.3 50.5
7 Sharon F 15 62.5 112.
8 Tammy F 14 62.8 102.
9 Alfred M 14 69 112.
10 Duke M 14 63.5 102.
11 Guido M 15 67 133
12 James M 12 57.3 83
13 Jeffrey M 13 62.5 84
14 John M 12 59 99.5
15 Philip M 16 72 150
16 Robert M 12 64.8 128
17 Thomas M 11 57.5 85
18 William M 15 66.5 112
```

横向连接自动找到了共同的变量 `name` 作为连接的键值,可以在 `inner_join()` 中用 `by=` 指定键值变量名,如果有不同的变量名,可以用 `by = c("a"="b")` 的格式指定左数据框的键值 `a` 与右数据框的键值 `b` 匹配进行连接。

两个表的横向连接,经常是多对一连接。例如, `d.stu` 中有学生学号、班级号、姓名、性别, `d.cl` 中有班级号、班主任名、年级,可以通过班级号将两个表连接起来:

```
d.stu <- tibble(
 sid=c(1,2,3,4,5,6),
 cid=c(1,2,1,2,1,2),
 sname=c("John", "Mary", "James", "Kitty", "Jasmine", "Kim"),
 sex=c("M", "F", "M", "F", "F", "M"))
d.stu

A tibble: 6 x 4
```

```
sid cid sname sex
<dbl> <dbl> <chr> <chr>
1 1 1 John M
2 2 2 Mary F
3 3 1 James M
4 4 2 Kitty F
5 5 1 Jasmine F
6 6 2 Kim M
```

```
d.cl <- tibble(
 cid=c(1,2),
 tname=c("Philip", "Joane"),
 grade=c("2017", "2016")
)
d.cl
```

```
A tibble: 2 x 3
cid tname grade
<dbl> <chr> <chr>
1 1 Philip 2017
2 2 Joane 2016
```

```
d.stu %>%
 left_join(d.cl, by="cid")
```

```
A tibble: 6 x 6
sid cid sname sex tname grade
<dbl> <dbl> <chr> <chr> <chr> <chr>
1 1 1 John M Philip 2017
2 2 2 Mary F Joane 2016
3 3 1 James M Philip 2017
4 4 2 Kitty F Joane 2016
5 5 1 Jasmine F Philip 2017
6 6 2 Kim M Joane 2016
```

`left_join()` 按照 `by` 变量指定的关键列匹配观测，左数据集所有观测不论匹

配与否全部保留，右数据集仅使用与左数据集能匹配的观测。不指定 `by` 变量时，使用左、右数据集的共同列作为关键列。如果左右数据集关键列变量名不同，可以用 `by=c("左名"="右名")` 的格式。

类似地，`right_join()` 保留右数据集的所有观测，而仅保留左数据集中能匹配的观测。`full_join()` 保留所有观测。`inner_join()` 仅保留能匹配的观测。

## 27.16 利用第二个数据集筛选

`left_join()` 将右表中与左表匹配的观测的额外的列添加到左表中。如果希望按照右表筛选左表的观测，可以用 `semi_join()`，函数 `anti_join()` 则是要求保留与右表不匹配的观测。

## 27.17 数据集的集合操作

R 的 `intersect()`，`union()`，`setdiff()` 本来是以向量作为集合进行集合操作。`dplyr` 包也提供了这些函数，但是将两个 `tibble` 的各行作为元素进行集合操作。

## 27.18 数据框纵向合并

矩阵或数据框要纵向合并，使用 `rbind` 函数即可。要求变量集合是相同的，变量次序可以不同。

比如，有如下两个分开男生、女生的数据框：

```
d3.class <- d.class %>%
 select(name, sex, age) %>%
 filter(sex=="M")
d4.class <- d.class %>%
 select(name, sex, age) %>%
 filter(sex=="F")
```

合并行如下：

```
rbind(d3.class, d4.class)
```

```
A tibble: 19 x 3
name sex age
<chr> <fct> <dbl>
1 Alfred M 14
2 Duke M 14
3 Guido M 15
4 James M 12
5 Jeffrey M 13
6 John M 12
7 Philip M 16
8 Robert M 12
9 Thomas M 11
10 William M 15
11 Alice F 13
12 Becka F 13
13 Gail F 14
14 Karen F 12
15 Kathy F 12
16 Mary F 15
17 Sandy F 11
18 Sharon F 15
19 Tammy F 14
```

将下面的数据框的变量列次序打乱，合并不受影响：

```
rbind(d3.class, d4.class[, c("age", "name", "sex")])
```

```
A tibble: 19 x 3
name sex age
<chr> <fct> <dbl>
1 Alfred M 14
2 Duke M 14
3 Guido M 15
```

```
4 James M 12
5 Jeffrey M 13
6 John M 12
7 Philip M 16
8 Robert M 12
9 Thomas M 11
10 William M 15
11 Alice F 13
12 Becca F 13
13 Gail F 14
14 Karen F 12
15 Kathy F 12
16 Mary F 15
17 Sandy F 11
18 Sharon F 15
19 Tammy F 14
```

## 27.19 标准化

设  $x$  是各列都为数值的列表（包括数据框和 tibble）或数值型矩阵，用 `scale(x)` 可以把每一列都标准化，即每一列都减去该列的平均值，然后除以该列的样本标准差。用 `scale(x, center=TRUE, scale=FALSE)` 仅中心化而不标准化。如

```
d.class %>%
 select(height, weight) %>%
 scale()
```

```
height weight
[1,] -1.13843504 -0.70371312
[2,] 0.57794313 -0.08897522
[3,] 0.38290015 -0.44025402
[4,] -1.17744363 -1.01108207
[5,] -0.49479323 -0.68175819
```

```
[6,] 0.81199469 0.52576268
[7,] -2.15265850 -2.17469309
[8,] 0.03182280 0.54771760
[9,] 0.09033569 0.10861910
[10,] 1.29960213 0.54771760
[11,] 0.22686577 0.10861910
[12,] 0.90951618 1.44786952
[13,] -0.98240066 -0.74762297
[14,] 0.03182280 -0.70371312
[15,] -0.65082761 -0.02311045
[16,] 1.88473105 2.19433697
[17,] 0.48042164 1.22832027
[18,] -0.94339207 -0.65980327
[19,] 0.81199469 0.52576268
attr("scaled:center")
height weight
62.33684 100.02632
attr("scaled:scale")
height weight
5.127075 22.773933
```

为了把  $x$  的每列变到  $[0, 1]$  内，可以用如下的方法：

```
d.class %>%
 select(height, weight) %>%
 scale(center=apply(., 2, min),
 scale=apply(., 2, max) - apply(., 2, min))
```

其中的 `.` 在管道操作中表示被传递处理的变量（一般是数据框）。也可以写一个自定义的进行零一标准化的函数：

```
scale01 <- function(x){
 mind <- apply(x, 2, min)
 maxd <- apply(x, 2, max)
 scale(x, center=mind, scale=maxd-mind)
}
```

```
d.class %>%
 select(height, weight) %>%
 scale01
```

```
height weight
[1,] 0.2512077 0.3366834
[2,] 0.6763285 0.4773869
[3,] 0.6280193 0.3969849
[4,] 0.2415459 0.2663317
[5,] 0.4106280 0.3417085
[6,] 0.7342995 0.6180905
[7,] 0.0000000 0.0000000
[8,] 0.5410628 0.6231156
[9,] 0.5555556 0.5226131
[10,] 0.8550725 0.6231156
[11,] 0.5893720 0.5226131
[12,] 0.7584541 0.8291457
[13,] 0.2898551 0.3266332
[14,] 0.5410628 0.3366834
[15,] 0.3719807 0.4924623
[16,] 1.0000000 1.0000000
[17,] 0.6521739 0.7788945
[18,] 0.2995169 0.3467337
[19,] 0.7342995 0.6180905
attr(,"scaled:center")
height weight
51.3 50.5
attr(,"scaled:scale")
height weight
20.7 99.5
```

注意在管道操作中某个操作除了被传递的第一自变量外没有其它自变量时，可以不写函数调用的空括号（）。

函数 `sweep()` 可以执行对每列更一般的变换。



## 27.20 用 reshape 包做长宽表转换

前面已经讲到，用 tidyr 的 `gather()`、`separate()`、`spread()` 等函数可以进行长宽表的转换。reshape 包是另一个可以进行长宽表转换的扩展包，这里列出其使用方法作为参考，建议主要使用 tidyr 的方法。

设数据框 `d.long` 变量为 `subject`(病人号, 有 5 个不同值), `time`(随访序号, 取 1,2,3), 变量 `x`, `y`(每个病人每次随访的两个测量指标值)。数据框有  $5 \times 3 = 15$  个观测 (行)。数据在如下的 `longtab.csv` 中:

```
subject,time,x,y
1,1,5,9
1,2,7,7
1,3,8,6
2,1,8,1
2,2,2,1
2,3,10,9
3,1,7,10
3,2,2,8
3,3,5,3
4,1,1,10
4,2,5,1
4,3,6,1
5,1,9,8
5,2,7,8
5,3,10,10
```

读入:

```
d.long <- as.data.frame(read_csv('longtab.csv'))
```

```
Parsed with column specification:
cols(
subject = col_double(),
time = col_double(),
x = col_double(),
```

```
y = col_double()
)
```

(reshape 包对 tibble 类型支持不好，所以转换成普通的数据框)

实际中，常常需要把每个病人的 3 次随访的 6 个测量指标合并到一行当中，这称为长表变宽表问题；反之则称为宽表变长表问题。

原始数据列表如下：

```
knitr::kable(d.long)
```

subject	time	x	y
1	1	5	9
1	2	7	7
1	3	8	6
2	1	8	1
2	2	2	1
2	3	10	9
3	1	7	10
3	2	2	8
3	3	5	3
4	1	1	10
4	2	5	1
4	3	6	1
5	1	9	8
5	2	7	8
5	3	10	10

### 27.20.1 用 melt() 融化

reshape 包用 melt() 函数把数据框转换为一个容易变形的格式，称为“融化”。melt() 函数把变量分为两种：分组用 (id.var)，测量用 (measure.var)。融化保持分组不变，但是将所有的测量变量合并到一列中，列名为 value；相应的变量名保存到一列中，列名为 variable。

如下程序把 d.long 转化为“融化”格式：

```
library(reshape)
```

```
##
```

```
载入程辑包: 'reshape'
```

```
The following object is masked from 'package:dplyr':
```

```
##
```

```
rename
```

```
The following objects are masked from 'package:tidyr':
```

```
##
```

```
expand, smiths
```

```
melt.long <- melt(as.data.frame(d.long),
 id.vars=c('subject', 'time'),
 measure.vars=c('x', 'y'))
```

融化后的数据框为:

```
knitr::kable(melt.long)
```

subject	time	variable	value
1	1	x	5
1	2	x	7
1	3	x	8
2	1	x	8
2	2	x	2
2	3	x	10
3	1	x	7
3	2	x	2
3	3	x	5
4	1	x	1
4	2	x	5
4	3	x	6
5	1	x	9
5	2	x	7
5	3	x	10
1	1	y	9
1	2	y	7
1	3	y	6
2	1	y	1
2	2	y	1
2	3	y	9
3	1	y	10
3	2	y	8
3	3	y	3
4	1	y	10
4	2	y	1
4	3	y	1
5	1	y	8
5	2	y	8
5	3	y	10

### 27.20.2 用 cast() 函数变形

用 `cast()` 函数把融化的数据框转换为要求的格式。例如，下面的程序把 `melt.long` 重新转化成了原来 `d.long` 的格式：

```
d1 <- cast(melt.long, subject + time ~ variable)
```

`cast()` 的第二自变量是公式，波折号左边是分组变量，右边是要从纵向转为横向的变量，这里把 `variable` 中的 `x` 和 `y` 转为横向，相应的变量值从 `value` 中取出。结果：

```
knitr::kable(d1)
```

subject	time	x	y
1	1	5	9
1	2	7	7
1	3	8	6
2	1	8	1
2	2	2	1
2	3	10	9
3	1	7	10
3	2	2	8
3	3	5	3
4	1	1	10
4	2	5	1
4	3	6	1
5	1	9	8
5	2	7	8
5	3	10	10

如下的程序把融化后的数据框转换为宽表，5 位病人每人的 3 次随访的 6 个测量值都合并到同一行中：

```
d2 <- cast(melt.long, subject ~ variable + time)
```

这里公式以病人作为仅有的分类，而变量 `x`, `y` 和 3 个 `time`(时间) 都变成了横向：

```
knitr::kable(d2)
```

subject	x_1	x_2	x_3	y_1	y_2	y_3
1	5	7	8	9	7	6
2	8	2	10	1	1	9
3	7	2	5	10	8	3
4	1	5	6	10	1	1
5	9	7	10	8	8	10

变量名与时间之间以下划线连接。

### 27.20.3 宽表变成长表

当数据框是宽表，变量名中用序号表示时间时，需要先把变量名与时间分离出来。比如，`d.wide` 数据框是如下的一个子集：

```
d.wide <- data.frame(
 subject=1:5,
 x1=c(5,8,7,1,9),
 x2=c(7,2,2,5,7),
 y1=c(9,1,10,10,8)
)
d.wide
```

```
subject x1 x2 y1
1 1 5 7 9
2 2 8 2 1
3 3 7 2 10
4 4 1 5 10
5 5 9 7 8
```

还是先融化为长表：

```
melt.wide <- melt(
 d.wide, id.vars=c("subject"), measure.vars=c("x1", "x2", "y1"))
```

宽表融化后结果：

```
knitr::kable(melt.wide)
```

subject	variable	value
1	x1	5
2	x1	8
3	x1	7
4	x1	1
5	x1	9
1	x2	7
2	x2	2
3	x2	2
4	x2	5
5	x2	7
1	y1	9
2	y1	1
3	y1	10
4	y1	10
5	y1	8

这个结果没有将时间分离出来，为此使用字符串处理：

```
melt.wide[, "time"] <- as.numeric(substr(melt.wide[, "variable"], 2))
melt.wide[, "variable"] <- substr(melt.wide[, "variable"], 1, 1)
```

```
knitr::kable(melt.wide)
```

subject	variable	value	time
1	x	5	1
2	x	8	1
3	x	7	1
4	x	1	1
5	x	9	1
1	x	7	2
2	x	2	2
3	x	2	2
4	x	5	2
5	x	7	2
1	y	9	1
2	y	1	1
3	y	10	1
4	y	10	1
5	y	8	1

现在的 `melt.wide` 已经有分类变量 `subject`，时间变量 `time`，变量名列 `variable` 和变量值列 `value`，可以用 `cast()` 转换了。例如，转换成每个病人两次随访的格式，这时因为 `y` 没有第二次随访值，会等于缺失值：

```
cast.wide <- cast(melt.wide, subject + time ~ variable)
knitr::kable(cast.wide)
```

subject	time	x	y
1	1	5	9
1	2	7	NA
2	1	8	1
2	2	2	NA
3	1	7	10
3	2	2	NA
4	1	1	10
4	2	5	NA
5	1	9	8
5	2	7	NA



变量名与时间之间更复杂的连接关系，如 `abc1`, `abc2`, `x1`, `y1`，不能简单用取子串完成，可以借助正则表达式拆开。

#### 27.20.4 变形同时进行概括

如果 `cast()` 后每组有不止一个值，可以指定一个统计函数进行汇总。下面的程序计算每个病人的所有 3 次随访的 `x` 平均值和 `y` 平均值：

```
cast(melt.long, subject ~ variable, mean)
```

```
subject x y
1 1 6.666667 7.333333
2 2 6.666667 3.666667
3 3 4.666667 7.000000
4 4 4.000000 4.000000
5 5 8.666667 8.666667
```

下面的程序对每个病人的每个随访时间，计算 `x` 和 `y` 的最大值：

```
cast(melt.long, subject + time ~ ., max)
```

```
subject time (all)
1 1 1 9
2 1 2 7
3 1 3 8
4 2 1 8
5 2 2 2
6 2 3 10
7 3 1 10
8 3 2 8
9 3 3 5
10 4 1 10
11 4 2 5
12 4 3 6
13 5 1 9
14 5 2 8
```

```
15 5 3 10
```

注意公式一端没有变量指定时用句点。

reshape 包的某些函数与 tidyverse 系统中函数冲突，所以使用完 reshape 后应卸载：

```
detach(package:reshape)
```

# Chapter 28

## 数据汇总

前面讲了用 `dplyr` 包的 `summarise()` 函数进行简单概括的方法。下面描述其它一些数据概括方法。

### 28.1 用 `summary()` 函数作简单概括

在 `d.cancer` 中保存了 34 个病人的代码 (`id`)、年龄 (`age`)、性别 (`sex`)、病理类型 (`type`)、放疗前肿瘤体积 (`v0`)、放疗后肿瘤体积 (`v1`)。其中, `sex`、`type` 用来作为分类, 年龄可以作为连续型取值变量, 也可以分组后作为分类。体积是连续型取值变量。

对数值型向量 `x`, 用 `summary(x)` 可以获得变量的平均值、中位数、最小值、最大值、四分之一和四分之三分位数。如

```
summary(d.cancer[["v0"]])
```

```
Min. 1st Qu. Median Mean 3rd Qu. Max.
12.58 43.77 93.40 110.08 157.18 330.24
```

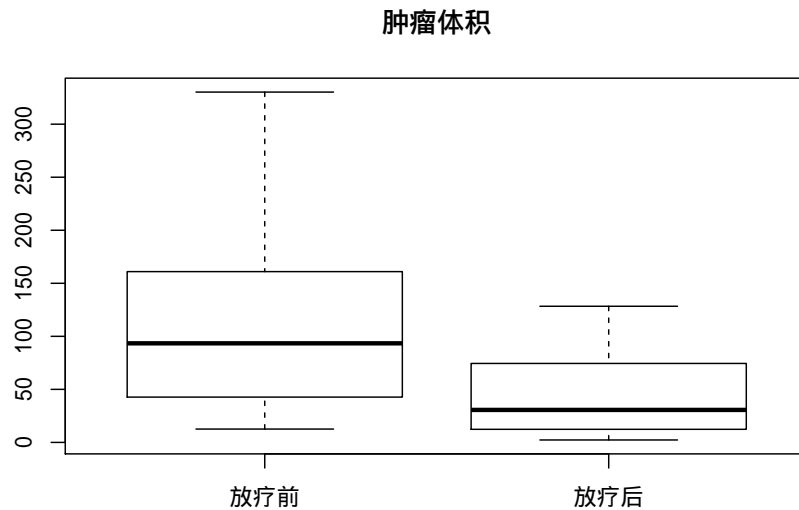
```
summary(d.cancer[["v1"]])
```

```
Min. 1st Qu. Median Mean 3rd Qu. Max.
2.30 12.73 30.62 44.69 72.94 128.34
```

可以看出放疗后体积减小了很多。

可以用盒形图表现类似的信息，如

```
boxplot(list('放疗前'=d.cancer[["v0"]],
 '放疗后'=d.cancer[["v1"]]), main='肿瘤体积')
```



对一个数据框 `d`，用 `summary(d)` 可以获得每个连续型变量的基本统计量，和每个离散取值变量的频率。如

```
summary(d.cancer)
```

```
id age sex type
Min. : 1.00 Min. :49.00 Length:34 Length:34
1st Qu.: 9.25 1st Qu.:55.00 Class :character Class :character
Median :17.50 Median :67.00 Mode :character Mode :character
Mean :17.50 Mean :64.13
3rd Qu.:25.75 3rd Qu.:70.00
Max. :34.00 Max. :79.00
##
NA's :11
v0 v1
```

```
Min. : 12.58 Min. : 2.30
1st Qu.: 43.77 1st Qu.: 12.73
Median : 93.40 Median : 30.62
Mean :110.08 Mean : 44.69
3rd Qu.:157.18 3rd Qu.: 72.94
Max. :330.24 Max. :128.34
##
```

对数据框 `d`，用 `str(d)` 可以获得各个变量的类型和取值样例。如

```
str(d.cancer)
```

```
Classes 'spec_tbl_df', 'tbl_df', 'tbl' and 'data.frame': 34 obs. of 6 variables:
$ id : num 1 2 3 4 5 6 7 8 9 10 ...
$ age : num 70 70 69 68 67 75 52 71 68 79 ...
$ sex : chr "F" "F" "F" "M" ...
$ type: chr "腺癌" "腺癌" "腺癌" "腺癌" ...
$ v0 : num 26.5 135.5 209.7 61 237.8 ...
$ v1 : num 2.91 35.08 74.44 34.97 128.34 ...
- attr(*, "spec")=
.. cols(
.. id = col_double(),
.. age = col_double(),
.. sex = col_character(),
.. type = col_character(),
.. v0 = col_double(),
.. v1 = col_double()
..)
```

用 `head(d)` 可以返回数据框（或向量、矩阵）的前几行，用 `tail(d)` 可以返回数据框的后几行。

## 28.2 连续型变量概括函数

对连续取值的变量  $x$ ，可以用 `mean`, `std`, `var`, `sum`, `prod`, `min`, `max` 等函数获取基本统计量。加 `na.rm=TRUE` 选项可以仅对非缺失值计算。

`sort(x)` 返回排序后的结果。`rev(x)` 把  $x$  所有元素次序颠倒后返回。`quantile(x, c(0.05, 0.95))` 可以求  $x$  的样本分位数。`rank(x)` 对  $x$  求秩得分（即名次，但从最小到最大排列）。

## 28.3 分类变量概括

分类变量一般输入为因子。对因子  $x$ ，`table(x)` 返回  $x$  的每个不同值的频率（出现次数），结果为一个类（class）为 `table` 的一维数组。每个元素有对应的元素名，为  $x$  的各水平值。如

```
res <- table(d.cancer[["sex"]]); res
```

```

F M
13 21
```

```
res['F']
```

```
F
13
```

对单个分类变量，`table` 结果是一个有元素名的向量。用 `as.data.frame()` 函数把 `table` 的结果转为数据框：

```
as.data.frame(res)
```

```
Var1 Freq
1 F 13
2 M 21
```

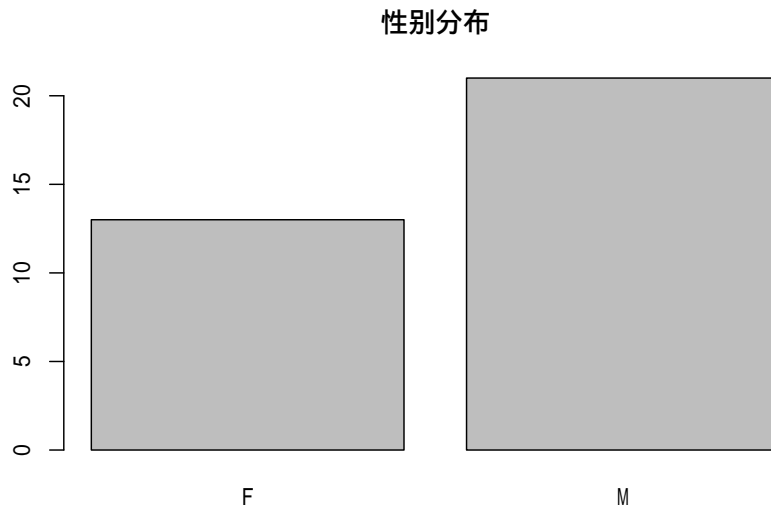
用 `prop.table()` 将频数转换成百分比：

```
prop.table(res)
```

```
##
F M
0.3823529 0.6176471
```

table 作的单变量频数表可以用 barplot 表现为图形，如：

```
barplot(res, main='性别分布')
```



对两个分类变量  $x_1$  和  $x_2$ ，其每个组合的出现次数可以用 `table(x1,x2)` 函数统计，结果叫做列联表。如

```
res2 <- with(d.cancer, table(sex, type)); res2
```

```
type
sex 鳞癌 腺癌
F 4 9
M 18 3
```

结果是一个类为 table 的二维数组（矩阵），每行以第一个变量  $x_1$  的各水平值为行名，每列以第二个变量  $x_2$  的各水平值为列名。这里用了 `with()` 函数引

入一个数据框，后续的参数中的表达式可以直接使用数据框的变量。

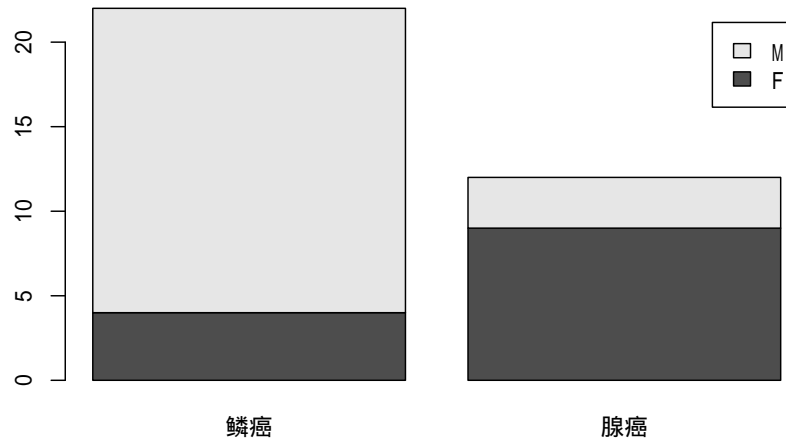
对两个分类变量, `table` 结果是一个矩阵。用 `as.data.frame` 函数把 `table` 的结果转为数据框:

```
as.data.frame(res2)
```

```
sex type Freq
1 F 鳞癌 4
2 M 鳞癌 18
3 F 腺癌 9
4 M 腺癌 3
```

列联表的结果可以用条形图表示。如

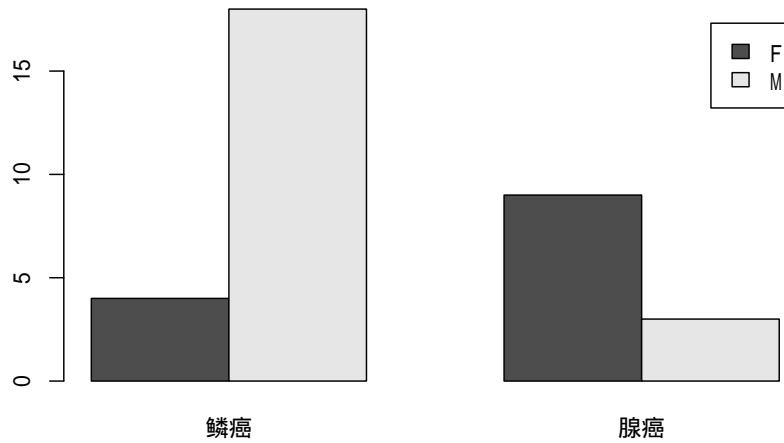
```
barplot(res2, legend=TRUE)
```



或



```
barplot(res2, legend=TRUE, beside=TRUE)
```



对于 `table()` 的结果列联表，可以用 `addmargins()` 函数增加行和与列和。  
如

```
addmargins(res2)
```

```
type
sex 鳞癌 腺癌 Sum
F 4 9 13
M 18 3 21
Sum 22 12 34
```

用 `margin.table()` 可以计算列联表行或列的和并返回，如

```
margin.table(res2, 1)
```

```
sex
F M
13 21
```

```
margin.table(res2, 2)
```

```
type
鳞癌 腺癌
22 12
```

用 `prop.table(r)` 把一个列联表 `r` 转换成百分比表。如

```
prop.table(res2)
```

```
type
sex 鳞癌 腺癌
F 0.11764706 0.26470588
M 0.52941176 0.08823529
```

用 `prop.table(res,1)` 把列联表 `res` 转换成行百分比表。用 `prop.table(res,2)` 把列联表 `res` 转换成列百分比表。如

```
prop.table(res2, 1)
```

```
type
sex 鳞癌 腺癌
F 0.3076923 0.6923077
M 0.8571429 0.1428571
```

```
prop.table(res2, 2)
```

```
type
sex 鳞癌 腺癌
F 0.1818182 0.7500000
M 0.8181818 0.2500000
```

在有多个分类变量时，用 `as.data.frame(table(x1, x2, x3))` 形成多个分类变量交叉分类的频数统计数据框。

`dplyr` 包的 `count()` 功能与 `table()` 类似。如

```
d.cancer %>%
 count(sex)
```

```
A tibble: 2 x 2
sex n
<chr> <int>
1 F 13
2 M 21
```

又如

```
d.cancer %>%
 count(sex, type)
```

```
A tibble: 4 x 3
sex type n
<chr> <chr> <int>
1 F 鳞癌 4
2 F 腺癌 9
3 M 鳞癌 18
4 M 腺癌 3
```

## 28.4 数据框概括

用 `colMeans()` 对数据框或矩阵的每列计算均值，用 `colSums()` 对数据框或矩阵的每列计算总和。用 `rowMeans()` 和 `rowSums()` 对矩阵的每行计算均值或总和。

数据框与矩阵有区别，某些适用于矩阵的计算对数据框不适用，例如矩阵乘法。用 `as.matrix()` 把数据框的数值子集转换成矩阵。

对矩阵，用 `apply(x, 1, FUN)` 对矩阵 `x` 的每一行使用函数 `FUN` 计算结果，用 `apply(x, 2, FUN)` 对矩阵 `x` 的每一列使用函数 `FUN` 计算结果。

如果 `apply(x,1,FUN)` 中的 `FUN` 对每个行变量得到多个 `m` 结果，结果将是一个矩阵，行数为 `m`，列数等于 `nrow(x)`。如果 `apply(x,2,FUN)` 中的 `FUN` 对每个列变量得到多个 `m` 结果，结果将是一个矩阵，行数为 `m`，列数等于 `ncol(x)`。例如：

```
apply(as.matrix(iris[,1:4]), 2,
 function(x)
 c(n=sum(!is.na(x)),
 mean=mean(x, na.rm=TRUE),
 sd=sd(x, na.rm=TRUE)))
```

```
Sepal.Length Sepal.Width Petal.Length Petal.Width
n 150.0000000 150.0000000 150.0000000 150.0000000
mean 5.8433333 3.0573333 3.7580000 1.1993333
sd 0.8280661 0.4358663 1.7652980 0.7622377
```

上面的例子如果改用 `dplyr` 包的 `summarise` 函数，需要对每个列分别写出汇总结果，会比较罗嗦。

## 28.5 分类概括

### 28.5.1 用 `dplyr` 包分类概括

用 `dplyr` 包的 `group_by()` 与 `summarise()` 配合可以比较简单地进行分类概括。适用于要概括的变量个数比较少的情形。

例如，按性别分组，计算 `v0` 的平均值：

```
d.cancer %>%
 group_by(sex) %>%
 summarise(mean.v0=mean(v0, na.rm=TRUE))
```

```
A tibble: 2 x 2
sex mean.v0
<chr> <dbl>
1 F 113.
2 M 108.
```

下面的程序按性别分组，分别计算 `v0` 与 `v1` 的平均值：

```
d.cancer %>%
 group_by(sex) %>%
 summarise(mean.v0=mean(v0, na.rm=TRUE),
 mean.v1=mean(v1, na.rm=TRUE))

A tibble: 2 x 3
sex mean.v0 mean.v1
<chr> <dbl> <dbl>
1 F 113. 42.7
2 M 108. 46.0
```

下面的程序按性别分组，计算 v0 和 v1 的平均值、标准差：

```
d.cancer %>%
 group_by(sex) %>%
 summarise(mean.v0=mean(v0, na.rm=TRUE),
 mean.v1=mean(v1, na.rm=TRUE),
 sd.v0=sd(v0, na.rm=TRUE),
 sd.v1=sd(v1, na.rm=TRUE))

A tibble: 2 x 5
sex mean.v0 mean.v1 sd.v0 sd.v1
<chr> <dbl> <dbl> <dbl> <dbl>
1 F 113. 42.7 100. 41.7
2 M 108. 46.0 66.5 37.3
```

以上结果如果每个变量的统计量分别占一行就好了，否则当需要分析的变量个数和统计量个数较多时结果表格可能过宽。从上面的例子还可以看出，当变量比较多、统计量比较多时，用 `summarise()` 写出的程序比较冗长。plyr 包功能更强，变量个数多、统计量多的时候能够统一处理，也能按用户需求排列结果，但是使用比 dplyr 包复杂一些。

dplyr 包也提供了一部分函数，可以解决对一批变量计算一批统计量的问题。`summarse_at()` 函数可以指定一批变量名与一批统计函数，自动命名结果变量，如：

```
d.cancer %>%
 group_by(sex) %>%
 summarise_at(c("v0", "v1"),
 list(avg = ~mean(.), std = ~sd(.)), na.rm=TRUE)

A tibble: 2 x 5
sex v0_avg v1_avg v0_std v1_std
<chr> <dbl> <dbl> <dbl> <dbl>
1 F 113. 42.7 100. 41.7
2 M 108. 46.0 66.5 37.3
```

其中的变量子集也可以用序号范围表示，或者用 `vars()` 函数写成不加撇号的格式，如：

```
d.cancer %>%
 group_by(sex) %>%
 summarise_at(vars(v0, v1),
 list(avg = ~mean(.), std = ~sd(.)), na.rm=TRUE)

A tibble: 2 x 5
sex v0_avg v1_avg v0_std v1_std
<chr> <dbl> <dbl> <dbl> <dbl>
1 F 113. 42.7 100. 41.7
2 M 108. 46.0 66.5 37.3
```

如果要对所有数值型变量计算某些统计量，可以用 `summarize_if(is.numeric, 变量名列表, list(变量后缀 =~ 统计函数名 <, ...>))`。如：

```
d.cancer %>%
 group_by(sex) %>%
 summarise_if(is_numeric,
 list(avg = ~mean(.), std = ~sd(.)), na.rm=TRUE)

Warning: Deprecated

Warning: Deprecated
```

```
Warning: Deprecated

Warning: Deprecated

Warning: Deprecated

A tibble: 2 x 9
sex id_avg age_avg v0_avg v1_avg id_std age_std v0_std v1_std
<chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 F 17.9 NA 113. 42.7 12.3 NaN 100. 41.7
2 M 17.2 NA 108. 46.0 8.50 NaN 66.5 37.3
```

### 28.5.2 用 `tapply()` 分组概括向量

用 `tapply()` 函数进行分组概括, 格式为:

```
tapply(X, INDEX, FUN)
```

其中 `X` 是一个向量, `INDEX` 是一个分类变量, `FUN` 是概括统计函数。

比如, 下面的程序分性别组计算疗前体积的均值:

```
with(d.cancer, tapply(v0, sex, mean))
```

```
F M
113.2354 108.1214
```

### 28.5.3 用 `aggregate()` 分组概括数据框

`aggregate` 函数对输入的数据框用指定的分组变量 (或交叉分组) 分组进行概括统计。例如, 下面的程序按性别分组计算年龄、疗前体积、疗后体积的平均值:

```
aggregate(d.cancer[,c("age", "v0", "v1")],
 list(sex=d.cancer[["sex"]]), mean, na.rm=TRUE)
```

```
sex age v0 v1
1 F 66.14286 113.2354 42.65538
```

```
2 M 63.25000 108.1214 45.95524
```

`aggregate()` 第一个参数是数据框，第二个参数是列表，列表元素是用来分组或交叉分组的变量，第三个参数是概括用的函数，概括用的函数的选项可以在后面给出。

可以同时计算多个概括统计量，如：

```
aggregate(d.cancer[,c('age', 'v0', 'v1')],
 list(sex=d.cancer[["sex"]]), summary)
```

上面的结果是两个观测、19 个变量的数据框，作为表格太宽了。后面讲的 `plyr` 包可以做出更合理的表格。

可以交叉分组后概括，如

```
with(d.cancer,
 aggregate(cbind(v0, v1), list(sex=sex, type=type), mean))
```

```
sex type v0 v1
1 F 鳞癌 126.99250 45.54750
2 M 鳞癌 113.55722 49.65556
3 F 腺癌 107.12111 41.37000
4 M 腺癌 75.50667 23.75333
```

#### 28.5.4 用 `split()` 函数分组后概括

`split` 函数可以把数据框的各行按照一个或几个分组变量分为子集的列表，然后可以用 `sapply()` 或 `vapply()` 对每组进行概括。如

```
sp <- split(d.cancer[,c('v0', 'v1')], d.cancer[, 'sex'])
sapply(sp, colMeans)
```

```
F M
v0 113.23538 108.12143
v1 42.65538 45.95524
```

返回矩阵，行为变量 `v0`, `v1`，列为不同性别，值为相应的变量在各性别组的平均值。当 `sapply()` 对列表每项的操作返回一个向量时，总是列表每项的输出



保存为结果的一列。colMeans 函数计算分组后数据框子集每列的平均值。

### 28.5.5 用 plyr 包进行分类概括

plyr 则是一个专注于分组后分别分析然后将分析结果尽可能合理地合并的扩展包, 功能强大, dplyr 包仅针对数据框, 使用更方便, 但是对于复杂情况功能不如 plyr 包强。

plyr 的输入支持数组、数据框、列表, 输出支持数组、数据框、列表或无输出。分组分析的函数输出格式需要与指定的输出格式一致。

这里主要介绍从数据框分组概括并将结果保存为数据框的方法, 使用 plyr 包的 ddply() 函数。实际上, dplyr 包的这种功能更方便。plyr 包的优点是可以自定义概括函数, 使得结果表格符合用户的预期, 处理多个变量时程序更简洁。

plyr 包与 dplyr 包的函数名冲突比较大, 所以需要先卸载 dplyr 包再调用 plyr 包:

```
if("dplyr" %in% .packages()) detach("package:dplyr")
library(plyr)
```

```
##
载入程辑包: 'plyr'
The following object is masked from 'package:purrr':
##
compact
```

ddply() 函数第一自变量是要分组的数据框, 第二自变量是分组用的变量名, 第三自变量是一个概括函数, 此概括函数以一个数据框子集(数据类型是数据框)为输入, 输出是一个数值、一个数值型向量或者一个数据框, 但最好是数据框。例如, 按性别分组, 计算 v0 的平均值:

```
ddply(d.cancer, 'sex',
 function(d) c(mean.v0 = mean(d[["v0"]], na.rm=TRUE)))
```

```
sex mean.v0
1 F 113.2354
2 M 108.1214
```

下面的程序按性别分组，分别计算 v0 与 v1 的平均值：

```
ddply(d.cancer, 'sex',
 function(d) colMeans(d[,c('v0', 'v1')]))

sex v0 v1
1 F 113.2354 42.65538
2 M 108.1214 45.95524
```

下面的程序按性别分组，计算 v0 和 v1 的平均值、标准差：

```
f1 <- function(dsub){
 tab <- tibble(
 '变量'=c('v0', 'v1'),
 '均值'=c(mean(dsub[, 'v0'], na.rm=TRUE),
 mean(dsub[, 'v1'], na.rm=TRUE)),
 '标准差'=c(sd(dsub[, 'v0'], na.rm=TRUE),
 sd(dsub[, 'v1'], na.rm=TRUE))
)
 tab
}
ddply(d.cancer, 'sex', f1)
```

```
sex 变量 均值 标准差
1 F v0 113.23538 100.06621
2 F v1 42.65538 41.72226
3 M v0 108.12143 66.45374
4 M v1 45.95524 37.27592
```

注意 `f1()` 结果是一个数据框。程序有些重复内容，对每个变量和每个统计量都需要分别写出，如果这样用 `plyr` 包就不如直接用 `dplyr::summarise()` 了。下面用 `vapply()` 简化程序。

按性别分组，然后 v0、v1 各自一行结果，计算非缺失值个数、均值、标准差、中位数：

```
f2 <- function(d, variables){
 d1 <- d[,variables,drop=FALSE]
 nnotmiss <- vapply(d1, function(x) sum(!is.na(x)), 1L)
```

```

xm <- vapply(d1, mean, 0.0, na.rm=TRUE)
xsd <- vapply(d1, sd, 1.0, na.rm=TRUE)
xmed <- vapply(d1, median, 0.0, na.rm=TRUE)
data.frame(variable=variables,
 n=nnotmiss,
 mean=xm,
 sd=xsd,
 median=xmed)
}
ddply(d.cancer, 'sex', f2, variables = c("v0", "v1"))

```

##	sex	variable	n	mean	sd	median
## 1	F	v0	13	113.23538	100.06621	67.37
## 2	F	v1	13	42.65538	41.72226	27.32
## 3	M	v0	21	108.12143	66.45374	101.65
## 4	M	v1	21	45.95524	37.27592	32.10

f2() 函数针对分组后的数据框子集，这样的函数可以先在一个子集上试验。在 f2() 函数中，设输入的数据子集为 d，要分析的变量组成的数据框为 d1，用 vapply() 函数对 d1 的每一列计算一种统计量，然后将每种统计量作为结果数据框的一列。vapply() 函数类似于 lapply() 和 sapply()，但是用第三个自变量表示要应用的变换函数的返回值类型和个数，用举例的方法给出。

ddply() 也可以对交叉分类后每个类分别汇总，例如按照性别与病理类型交叉分组后汇总 v0、v1:

```
ddply(d.cancer, c('sex', 'type'), f2, variables=c("v0", "v1"))
```

##	sex	type	variable	n	mean	sd	median
## 1	F	鳞癌	v0	4	126.99250	83.82544	119.000
## 2	F	鳞癌	v1	4	45.54750	23.55433	40.920
## 3	F	腺癌	v0	9	107.12111	110.67144	42.700
## 4	F	腺癌	v1	9	41.37000	48.95945	9.450
## 5	M	鳞癌	v0	18	113.55722	68.88281	103.275
## 6	M	鳞癌	v1	18	49.65556	38.96325	33.730
## 7	M	腺癌	v0	3	75.50667	44.36592	61.000

```
8 M 腺癌 v1 3 23.75333 11.32141 23.960
```

上面的程序写法适用于已知要分析的变量名的情况。如果想对每个数值型变量都分析，而且想把要计算的统计量用统一的格式调用，可以写成：

```
f3 <- function(d){
 ff <- function(x){
 c(n=sum(!is.na(x)),
 each(mean, sd, median)(x, na.rm=TRUE))
 }
 ldply(Filter(is.numeric, d), ff)
}
ddply(d.cancer, 'sex', f3)
```

```
sex .id n mean sd median
1 F id 13 17.92308 12.325188 19.00
2 F age 7 66.14286 6.792853 69.00
3 F v0 13 113.23538 100.066207 67.37
4 F v1 13 42.65538 41.722263 27.32
5 M id 21 17.23810 8.502381 17.00
6 M age 16 63.25000 10.096204 66.50
7 M v0 21 108.12143 66.453742 101.65
8 M v1 21 45.95524 37.275917 32.10
```

`ff()` 函数对输入的一个数值型向量计算 4 种统计量，返回一个长度为 4 的数值型向量，用来对分组后的数据子集中的一列计算 4 种统计量。`plyr` 包的 `each()` 函数接受多个函数，返回一个函数可以同时得到这几个函数的结果，结果中各元素用输入的函数名命名。`f3()` 函数中的 `Filter` 函数用于从列表或数据框中取出满足条件的项，这里取出输入的数据子集 `d` 中所有的数值型列。`f3()` 函数中的 `ldply()` 函数接受一个列表或看成列表的一个数据框，对数据框的每列应用 `ff()` 函数计算 4 种统计量，然后合并所有各列的统计量为一个数据框，结果数据框的每行对应于 `d` 中的一列。程序中的 `ddply()` 函数接受一个数据框，第二自变量指定用来将数据框分组的变量，第三自变量 `f3()` 是对分组后的数据框子集进行分析的函数，此函数接受一个数据框，输出一个数据框。

上面的程序也可以利用无名函数写成：

```
f4 <- function(x){
 c(n=sum(!is.na(x)),
 each(mean, sd, median)(x, na.rm=TRUE))
}
ddply(d.cancer, 'sex',
 function(d)
 ldply(Filter(is.numeric, d), f4))
```

```
sex .id n mean sd median
1 F id 13 17.92308 12.325188 19.00
2 F age 7 66.14286 6.792853 69.00
3 F v0 13 113.23538 100.066207 67.37
4 F v1 13 42.65538 41.722263 27.32
5 M id 21 17.23810 8.502381 17.00
6 M age 16 63.25000 10.096204 66.50
7 M v0 21 108.12143 66.453742 101.65
8 M v1 21 45.95524 37.275917 32.10
```

## 28.6 练习

- 把patients.csv读入“d.patients”中，并计算发病年龄、发病年、发病月、发病年月（格式如“200702”表示2007年2月份）。
- 把“现住地址国标”作为字符型，去掉最后两位，仅保留前6位数字，保存到变量“地址编码”中。
- 按照地址编码和发病年月交叉分类汇总发病人数，保存到数据框d.pas1中，然后保存为CSV文件“分区分年月统计.csv”中。要求结果有三列：“地址编码”、“发病年月”、“发病人数”。
- 按照地址编码和发病月分类汇总发病人数，保存到数据框d.pas2中，然后保存为CSV文件“分区分月统计.csv”中。要求每个地址编码占一行，各列为地址编码以及1、2、.....、12各月份，每行为同一地址编码各月份的发病数。

- 按发病年月和性别汇总发病人数，并计算同年月不分性别的发病总人数。结果保存到数据框 `d.pas3` 中，然后保存到 CSV 文件“分年月分性别统计.csv”中。要求每个不同年月占一行，变量包括年月、男性发病数、女性发病数、总计。
- 分析病人的职业分布，保存到数据框 `d.pas4` 中，然后保存到 CSV 文件“职业构成.csv”中。要求各列为职业、发病人数、百分比（结果乘以 100 并保留一位小数）。
- 把年龄分成 0—9, 11—19, ....., 70 以上各段，保存为“年龄段”变量。用年龄段和性别交叉汇总发病人数和百分比（结果乘以 100 并保留一位小数），保存到“年龄性别分布.csv”中。要求将每个年龄段的男性发病人数、发病率、女性发病人数、发病率存为一行。

## Part VI

### 绘图





# Chapter 29

## 绘图

R 语言的前身是 S 语言，S 语言的设计目的就是交互式数据分析、绘图。所以绘图是 R 的重要功能。

R 有最初的基本绘图，这是从 S 语言继承过来的，还有一些功能更易用、更强大的绘图系统，如 `lattice`、`ggplot2`。基本绘图使用简单，灵活性强，但是为了做出满意的图形需要比较多的调整。这里先讲解 R 语言的基本绘图功能。

R 的基本绘图功能有两类图形函数：高级图形函数，直接针对某一绘图任务作出完整图形；低级图形函数，在已有图形上添加内容。具备有限的与图形交互的能力（函数 `locator` 和 `identify`）。

### 29.1 常用高级图形

#### 29.1.1 条形图

`d.cancer` 数据框包含了肺癌病人放疗的一些数据，从 `cancer.csv` 读入：

```
d.cancer <- readr::read_csv("cancer.csv", locale=locale(encoding="GBK"))

Parsed with column specification:
cols(
id = col_double(),
```

```
age = col_double(),
sex = col_character(),
type = col_character(),
v0 = col_double(),
v1 = col_double()
)
```

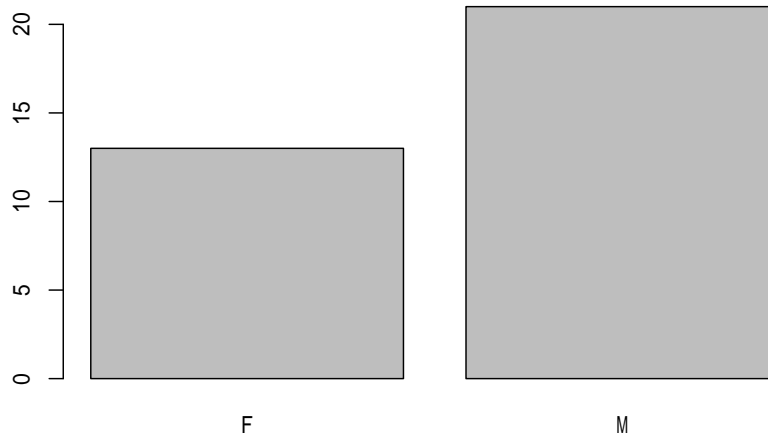
统计男女个数并用条形图表示:

```
res1 <- table(d.cancer[, 'sex']); print(res1)
```

```

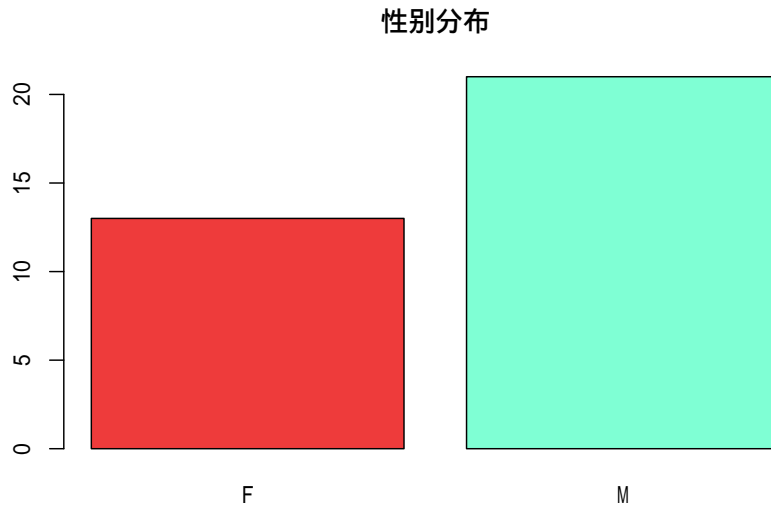
F M
13 21
```

```
barplot(res1)
```



可以增加标题, 采用不同的颜色:

```
barplot(res1, main=" 性别分布",
 col=c("brown2", "aquamarine1"))
```



R 函数 `colors()` 可以返回 R 中定义的用字符串表示的六百多种颜色名字。如

```
head(colors(), 6)
```

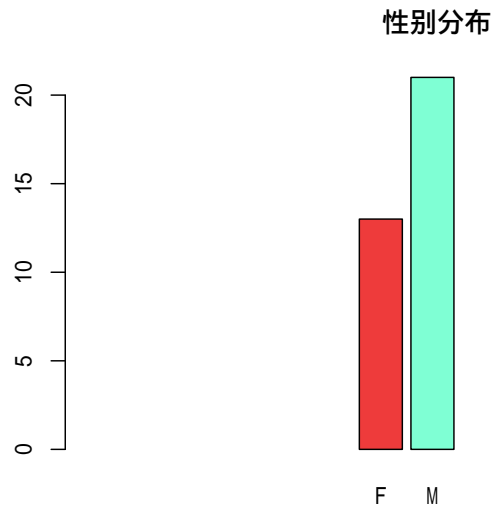
```
[1] "white" "aliceblue" "antiquewhite" "antiquewhite1"
[5] "antiquewhite2" "antiquewhite3"
```

下面的函数可以用来挑选颜色，鼠标点击画出的颜色就可以挑选，结果返回挑选出的颜色名：

```
select.colors <- function(){
 nc <- length(colors())
 x <- rep(seq(26), 26)[1:nc]
 y <- rep(seq(26), each=26)[1:nc]
 cols <- colors()
 plot(x, y, type="p", pch=16, cex=2,
 col=cols)
 res <- cols[identify(x,y, labels=cols)]
 res
}
```

用 `width` 选项与 `xlim` 选项配合可以调整条形宽度，如

```
barplot(res1, width=0.5, xlim=c(-3, 5),
 main=" 性别分布",
 col=c("brown2", "aquamarine1"))
```



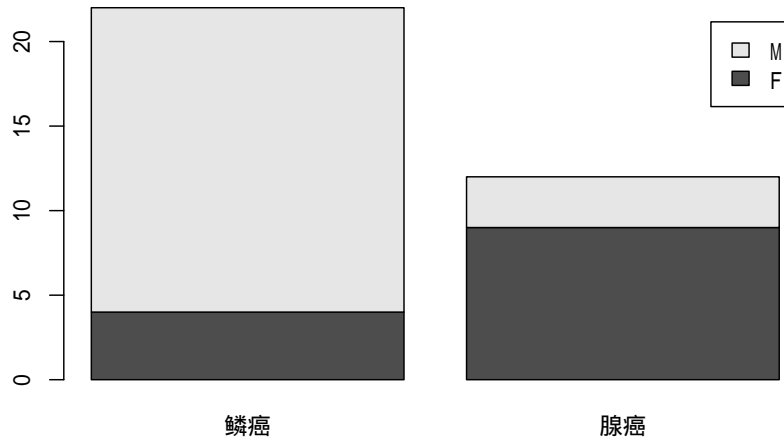
按性别与病理类型交叉分组后统计频数，结果称为列联表：

```
res2 <- with(d.cancer, table(sex, type)); res2
```

```
type
sex 鳞癌 腺癌
F 4 9
M 18 3
```

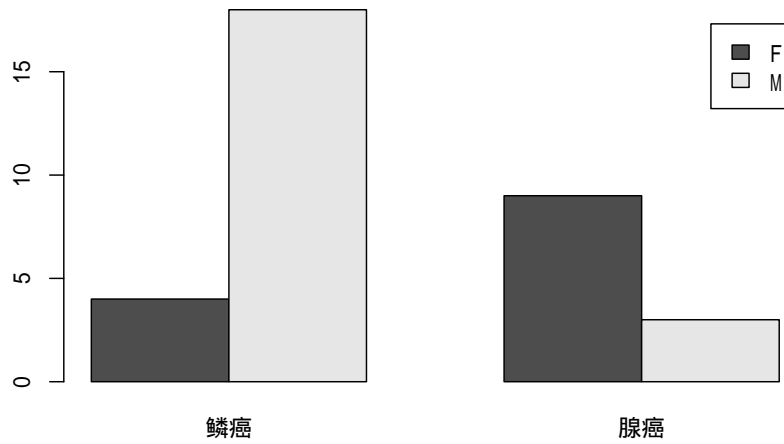
用分段条形图表现交叉分组频数，交叉频数表每列为一条：

```
barplot(res2, legend=TRUE)
```



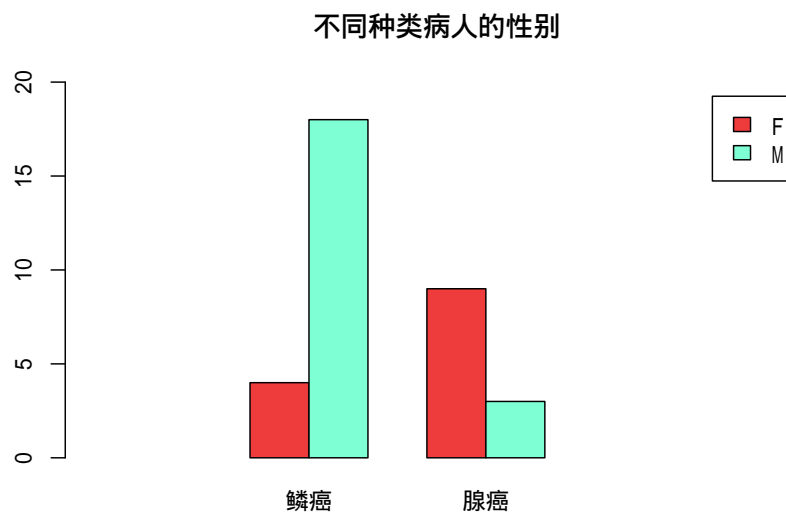
用并排条形图表现交叉分组频数，交叉频数表每列为一组：

```
barplot(res2, beside=TRUE, legend=TRUE)
```



增加标题，指定颜色，调整图例位置，调整条形宽度：

```
barplot(res2, beside=TRUE, legend=TRUE,
 main=' 不同种类病人的性别',
 ylim=c(0, 20),
 xlim=c(-1, 6), width=0.6,
 col=c("brown2", "aquamarine1"))
```



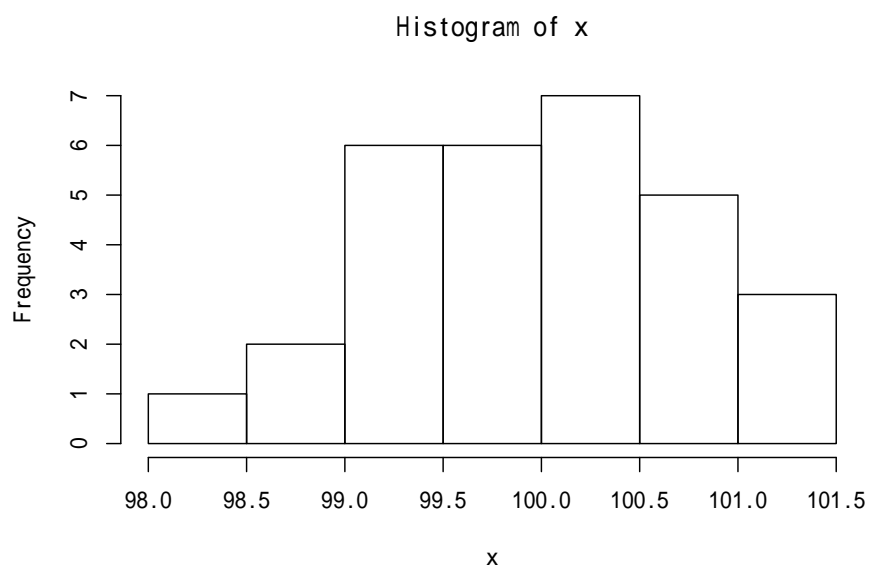
### 29.1.2 直方图和密度估计图

用 `hist` 作直方图以了解连续取值变量分布情况。例如，下面的程序模拟正态分布数据并做直方图：

```
x <- rnorm(30, mean=100, sd=1)
print(round(x,2))
```

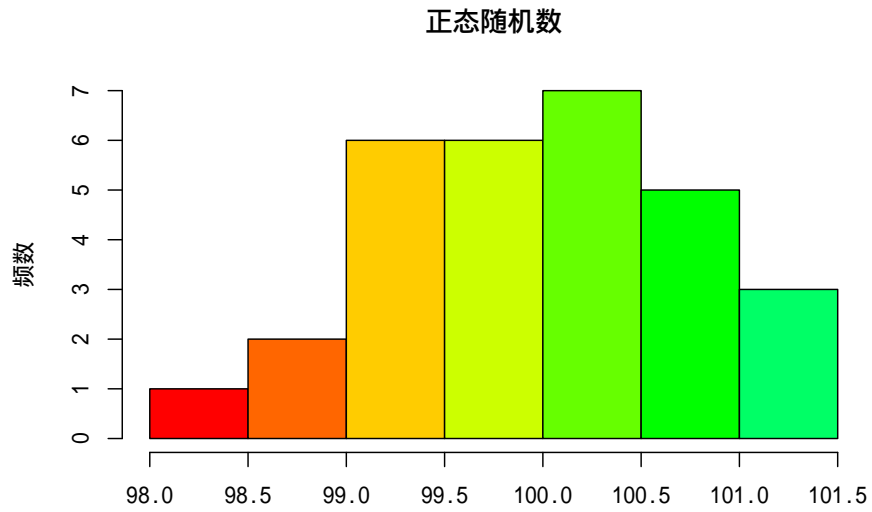
```
[1] 100.58 101.19 100.42 99.18 99.44 100.76 99.00 99.44 100.11 100.32
[11] 99.36 99.49 98.39 99.38 98.61 99.55 99.89 100.21 100.18 99.99
[21] 100.71 101.17 100.25 101.23 100.94 99.56 99.62 99.81 100.27 100.57
```

```
hist(x)
```



可以用 `main=`、`xlab=`、`ylab=` 等选项，可以用 `col=` 指定各个条形的颜色，如：

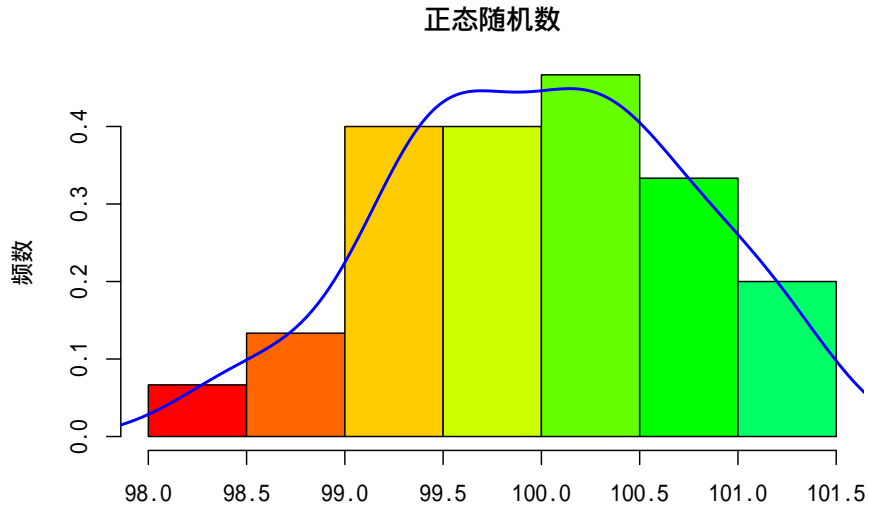
```
hist(x, col=rainbow(15),
 main=' 正态随机数', xlab='', ylab=' 频数')
```



函数 `density()` 估计核密度。下面的程序作直方图，并添加核密度曲线：

```
tmp.dens <- density(x)
hist(x, freq=FALSE,
 ylim=c(0,max(tmp.dens$y)),
 col=rainbow(15),
 main=' 正态随机数',
 xlab='', ylab=' 频数!')
lines(tmp.dens, lwd=2, col='blue')
```

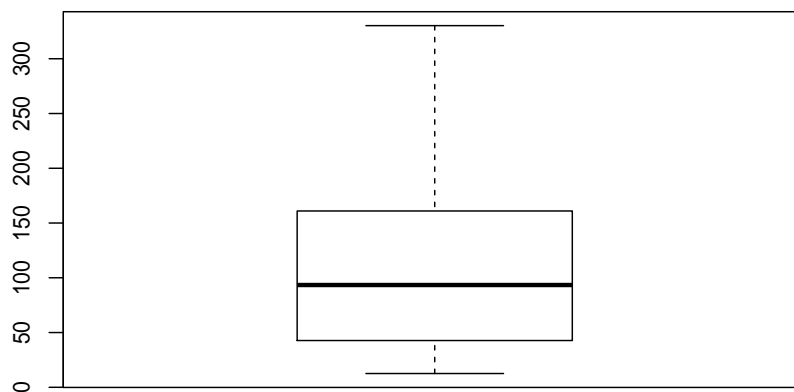




### 29.1.3 盒形图

盒形图可以简洁地表现变量分布，如

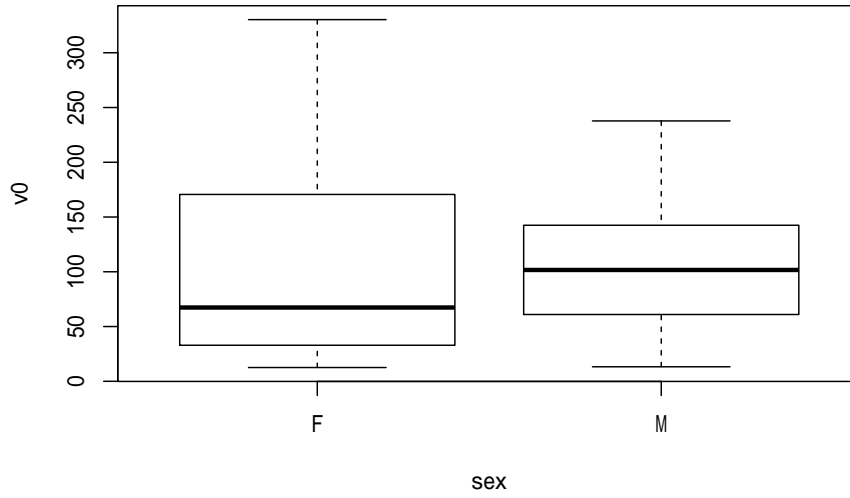
```
with(d.cancer, boxplot(v0))
```



其中中间粗线是中位数，盒子上下边缘是  $\frac{3}{4}$  和  $\frac{1}{4}$  分位数，两条触须线延伸到取值区域的边缘。

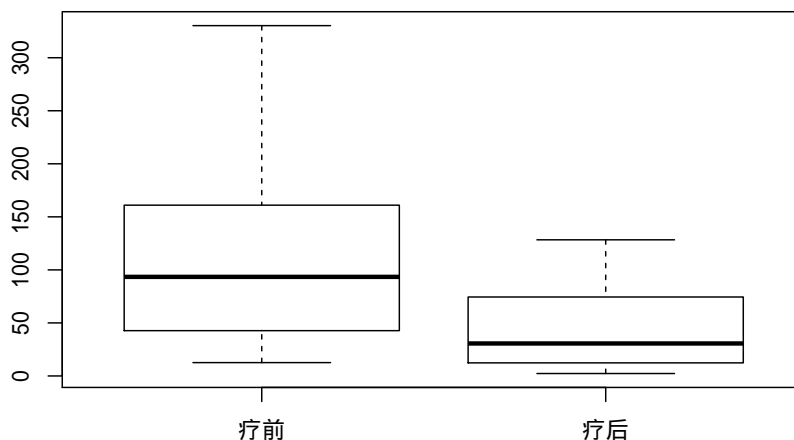
盒形图可以很容易地比较两组或多组，如

```
with(d.cancer, boxplot(v0 ~ sex))
```



也可以画若干个变量的并排盒形图，如

```
with(d.cancer,
 boxplot(list(' 疗前'=v0, ' 疗后'=v1)))
```



### 29.1.4 正态 QQ 图

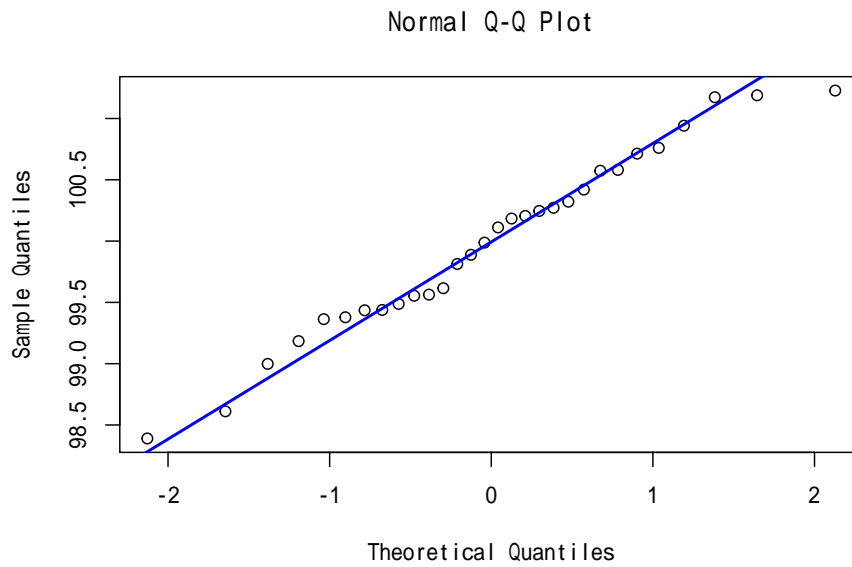
用 `qqnorm` 和 `qqline` 作正态 QQ 图。当变量样本来自正态分布总体时，正态 QQ 图的散点近似在一条直线周围。

QQ 图一般作法如下：设有  $n$  个观测  $y_1, y_2, \dots, y_n$ ，并已从小到大排列，则  $y_i$  是总体的  $i/n$  分位数的估计，设  $x_i$  是标准正态分布的  $i/n$  分位数，则在样本来自正态  $N(\mu, \sigma^2)$  的情况下，记  $N(\mu, \sigma^2)$  的分布函数为  $F(x)$ ，有  $y_i \approx F^{-1}(\frac{i}{n})$ ， $F(y_i) = \Phi(\frac{y_i - \mu}{\sigma}) \approx \frac{i}{n}$ ， $\frac{y_i - \mu}{\sigma} \approx \Phi^{-1}(\frac{i}{n}) = x_i$ ， $y_i \approx \mu + \sigma x_i$ ，于是用  $(x_i, y_i)$  ( $i = 1, 2, \dots, n$ ) 作为坐标画散点图应该近似呈现为截距  $\mu$ 、斜率  $\sigma$  的一条直线。

但是，上述的近似有一个小缺点： $y_1$  是观测到的最小值，对应于  $1/n$  分位数， $y_n$  是观测到的最大值，却对应于  $n/n = 100\%$  分位数，对最小和最大值的处理不对称，相当于说总体分布不能超过  $y_n$ ，这是不合理的。所以在实际画正态 QQ 图时， $y_i$  不是对应于标准正态的  $i/n$  分位数而是对应于略调整的如  $(i - 0.375)/(n + 0.25)$  这样的分位数，这种做法叫做连续性修正。这时， $y_1$  对应于  $\frac{0.625}{n+0.25}$  分位数而  $y_n$  对应于  $1 - \frac{0.625}{n+0.25}$  分位数，两边各留了  $\frac{0.625}{n+0.25}$ 。

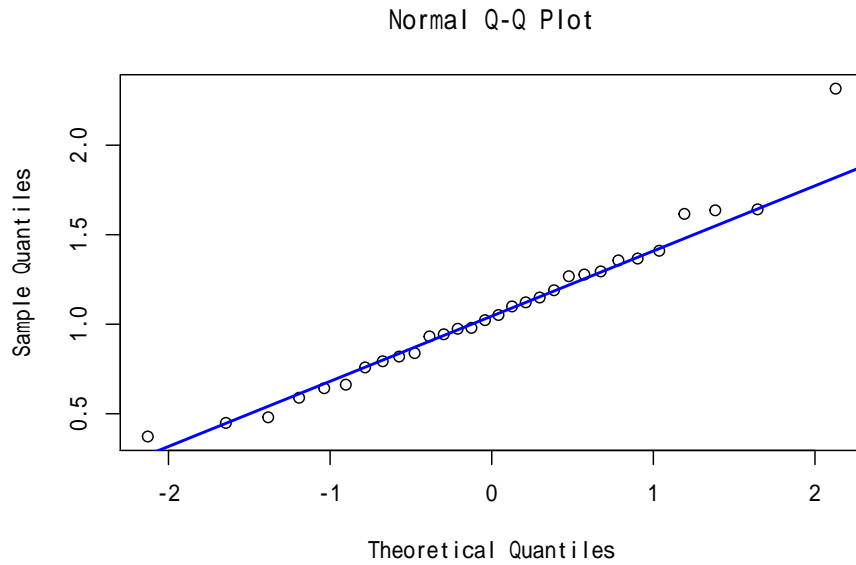
下面的程序模拟正态分布随机数，并作正态 QQ 图：

```
qqnorm(x)
qqline(x, lwd=2, col='blue')
```



下面的程序模拟对数正态数据，并作正态 QQ 图：

```
z <- 10^rnorm(30, mean=0, sd=0.2)
qqnorm(z)
qqline(z, lwd=2, col='blue')
```



### 29.1.5 散点图

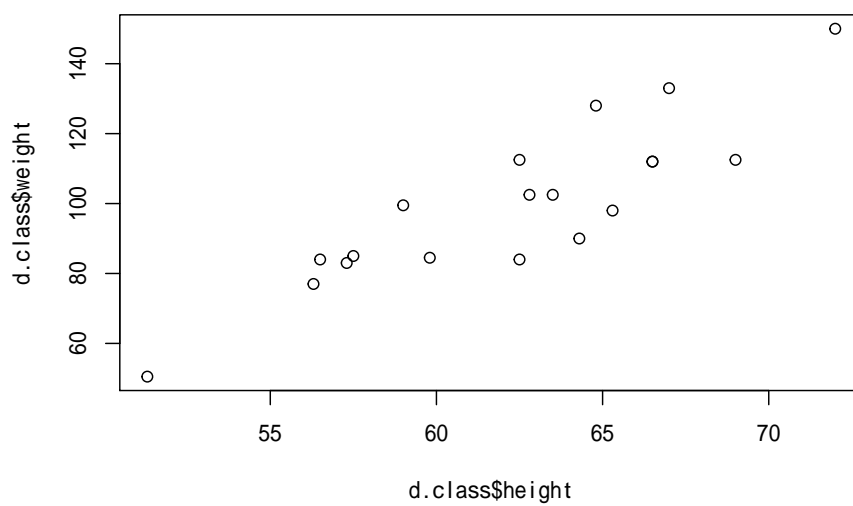
以 `d.class` 数据为例, 有 `name`, `sex`, `age`, `height`, `weight` 等变量。从 `class.csv` 读入:

```
d.class <- read_csv("class.csv")
```

```
Parsed with column specification:
cols(
name = col_character(),
sex = col_character(),
age = col_double(),
height = col_double(),
weight = col_double()
)
```

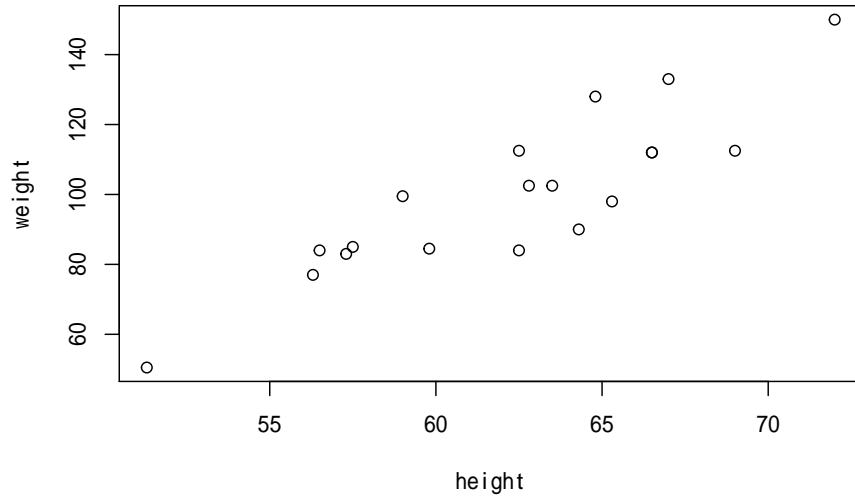
体重对身高的散点图:

```
plot(d.class$height, d.class$weight)
```



用 `with()` 函数简化数据框变量访问格式:

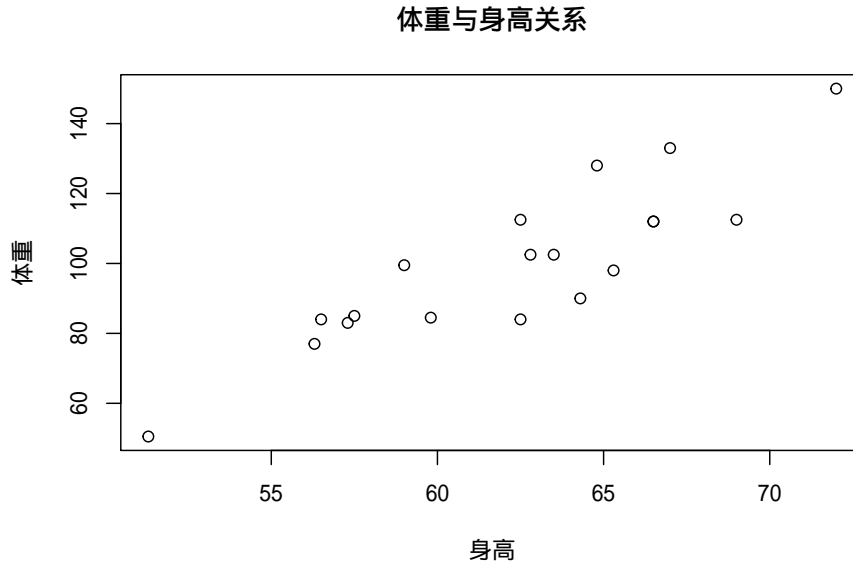
```
with(d.class,
 plot(height, weight))
```



在 `plot()` 函数内用 `main` 参数增加标题, 用 `xlab` 参数指定横轴标注, 用 `ylab` 参数指定纵轴标注, 如

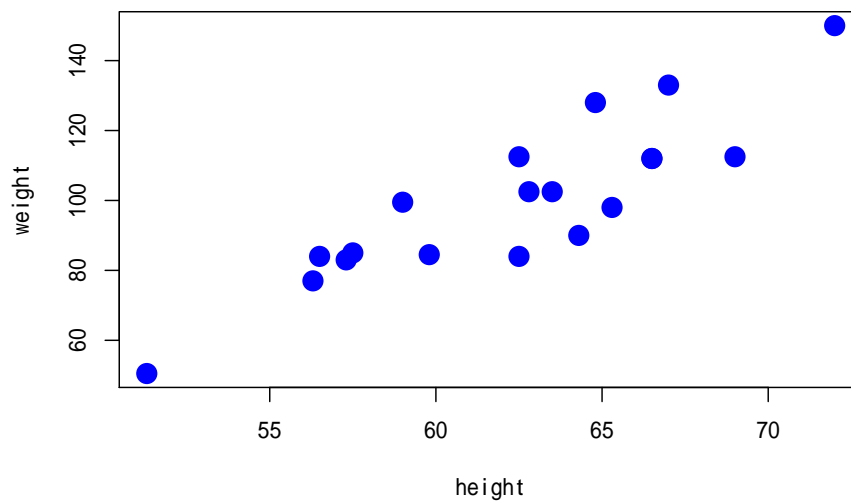
```
with(d.class,
 plot(height, weight,
 main=' 体重与身高关系',
 xlab=' 身高', ylab=' 体重'))
```





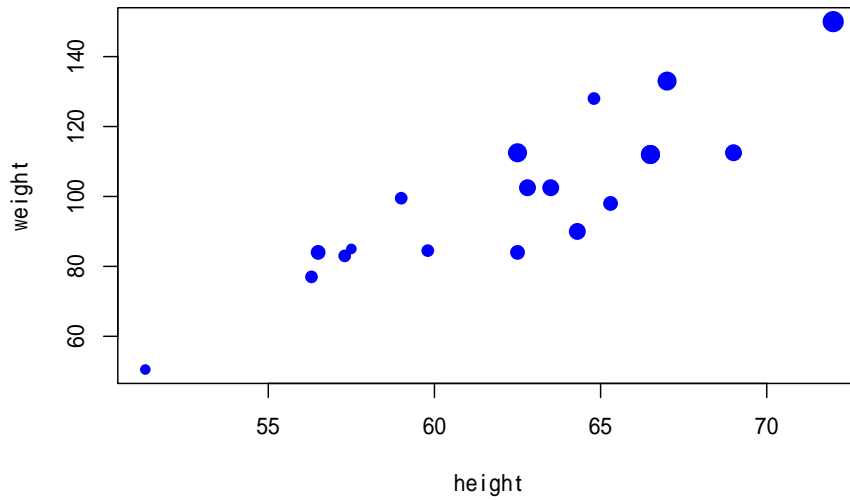
用 `pch` 参数指定不同散点形状, 用 `col` 参数指定颜色, 用 `cex` 参数指定大小, 如:

```
with(d.class,
 plot(height, weight,
 pch=16, col='blue',
 cex=2))
```



用气泡大小表现第三维（年龄）：

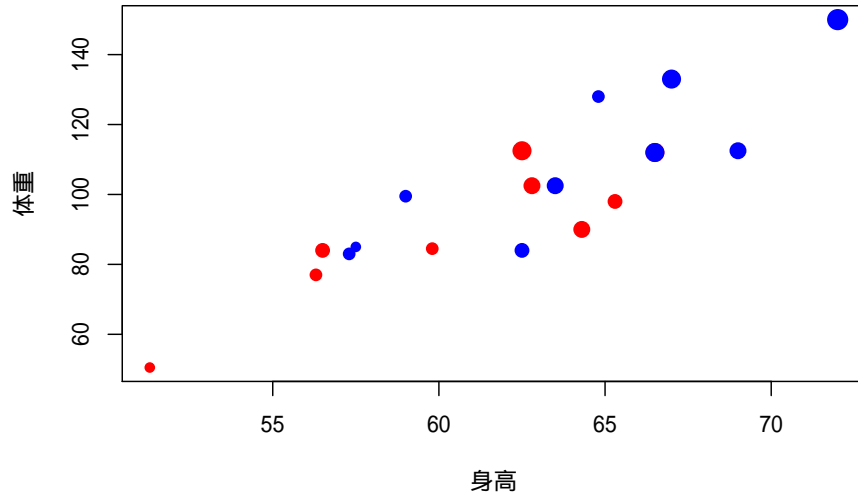
```
with(d.class,
 plot(height, weight,
 pch=16, col='blue',
 cex=1 + (age - min(age))/(max(age)-min(age))))
```



用气泡大小表现年龄，用颜色区分性别：

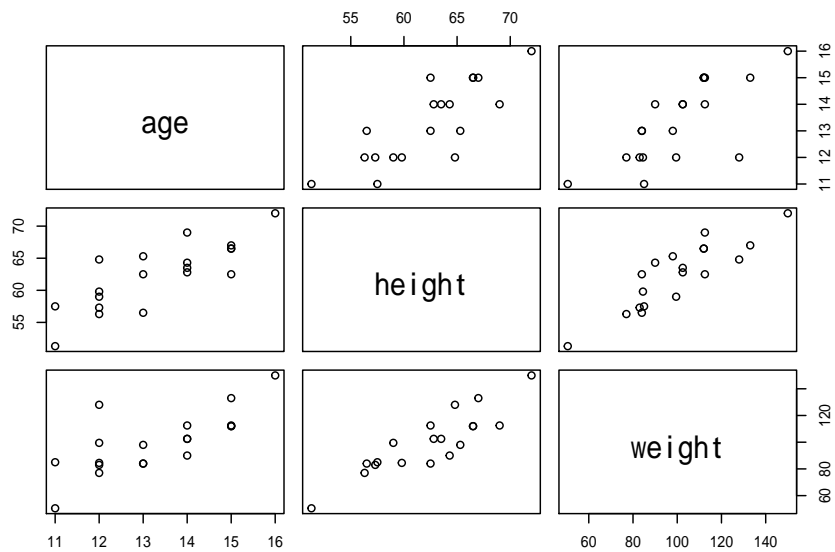
```
with(d.class,
 plot(height, weight,
 main=' 体重与身高关系',
 xlab=' 身高', ylab=' 体重',
 pch=16,
 col=ifelse(sex=='M', 'blue', 'red'),
 cex=1 + (age - min(age))
 / (max(age) - min(age)))
```

体重与身高关系



用 `pairs()` 函数可以做散点图矩阵:

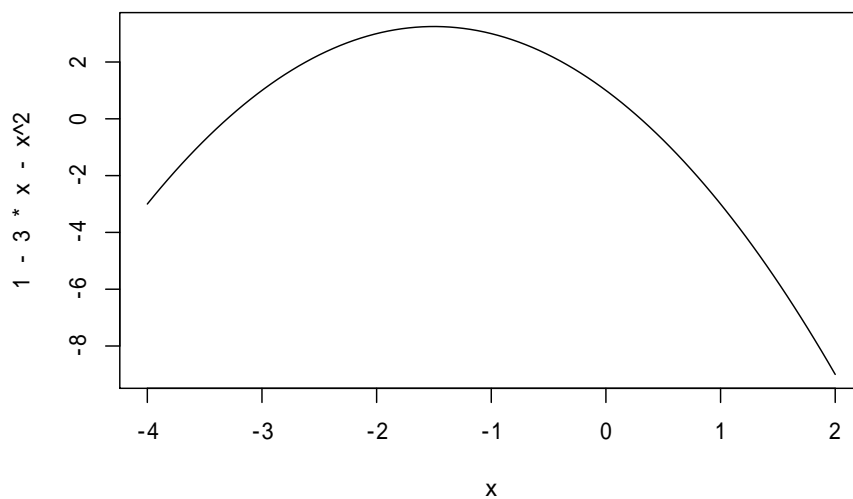
```
pairs(d.class[, c('age', 'height', 'weight')])
```



### 29.1.6 曲线图

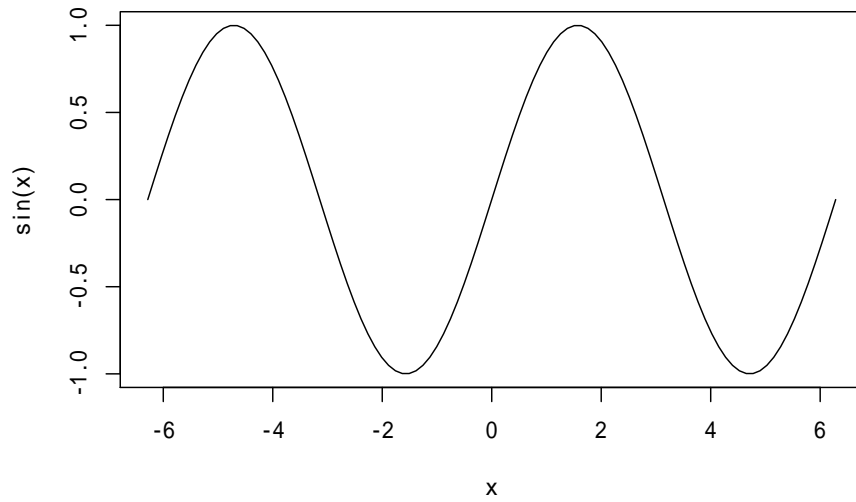
`curve()` 函数接受一个函数，或者一个以 `x` 为变量的表达式，以及曲线的自变量的左、右端点，绘制函数或者表达式的曲线图，如：

```
curve(1 - 3*x - x^2, -4, 2)
```



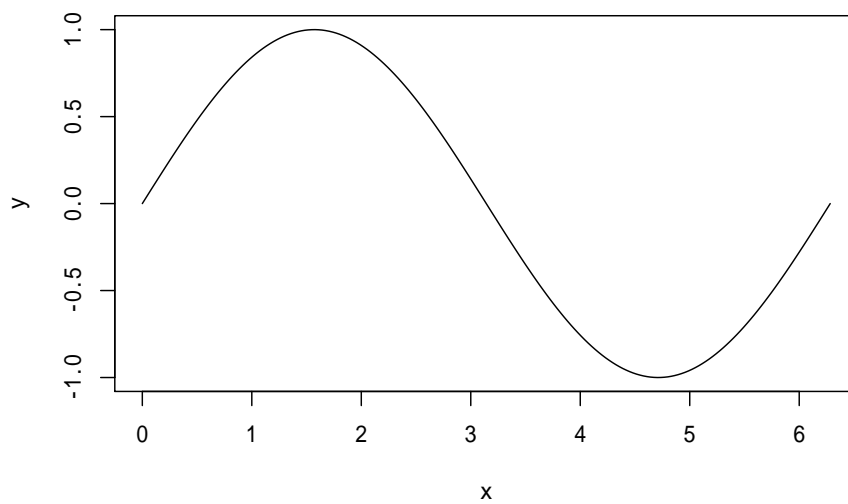
又如：

```
curve(sin, -2*pi, 2*pi)
```



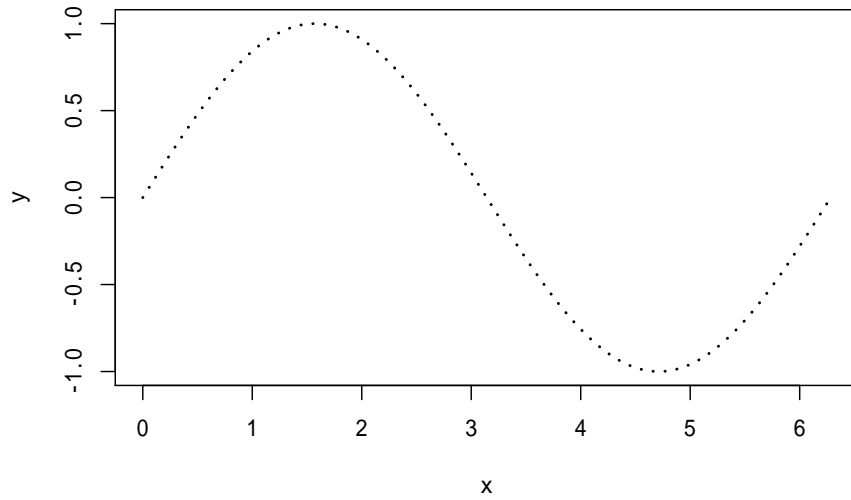
在 `plot` 函数中使用 `type='l'` 参数可以作曲线图，如

```
x <- seq(0, 2*pi, length=200)
y <- sin(x)
plot(x,y, type='l')
```



除了仍可以用 `main`, `xlab`, `ylab`, `col` 等参数外, 还可以用 `lwd` 指定线宽度, `lty` 指定虚线, 如

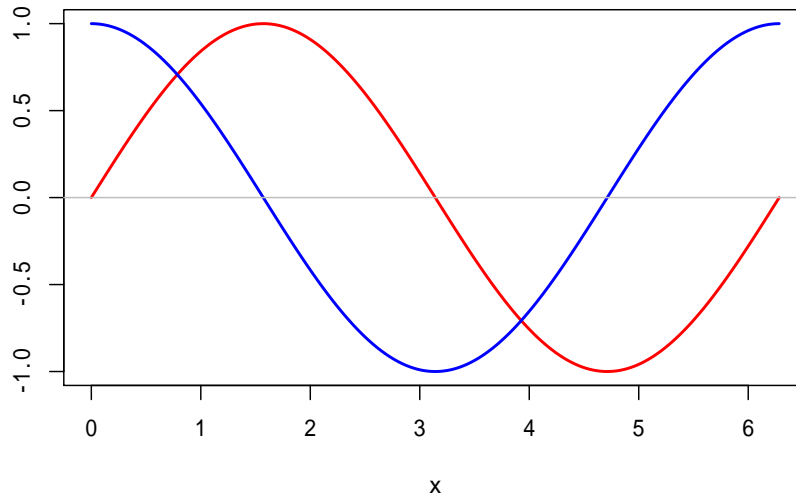
```
plot(x,y, type='l', lwd=2, lty=3)
```



多条曲线，可以用 `matplot()` 函数。例如

```
x <- seq(0, 2*pi, length=200)
y1 <- sin(x)
y2 <- cos(x)
matplot(x, cbind(y1, y2), type='l',
 lty=1, lwd=2, col=c("red", "blue"),
 xlab="x", ylab="")
abline(h=0, col='gray')
```





### 29.1.7 三维图

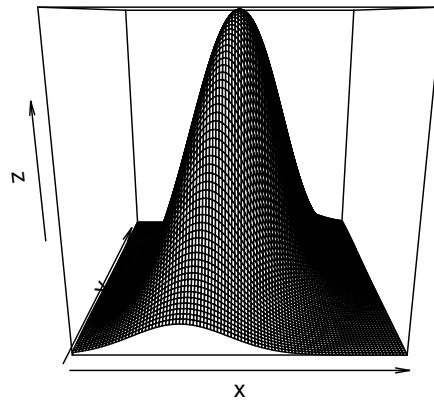
用 `persp` 函数作三维曲面图, `contour` 作等值线图, `image` 作色块图。坐标 `x` 和 `y` 构成一张平面网格, 数据 `z` 是包含 `z` 坐标的矩阵, 每行对应一个横坐标, 每列对应一个纵坐标。

下面的程序生成二元正态分布密度曲面数据:

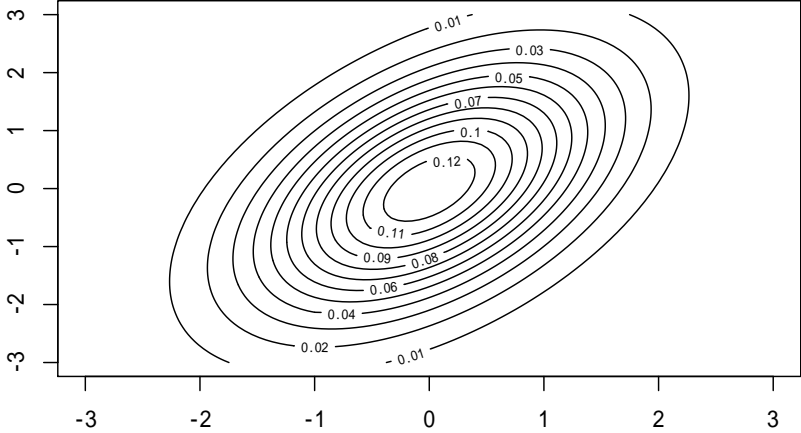
```
x <- seq(-3,3, length=100)
y <- x
f <- function(x,y,ssq1=1, ssq2=2, rho=0.5){
 det1 <- ssq1*ssq2*(1 - rho^2)
 s1 <- sqrt(ssq1)
 s2 <- sqrt(ssq2)
 1/(2*pi*sqrt(det1)) * exp(-0.5 / det1 * (
 ssq2*x^2 + ssq1*y^2 - 2*rho*s1*s2*x*y))
}
z <- outer(x, y, f)
```

作二元正态密度函数的三维曲面图、等高线图、色块图:

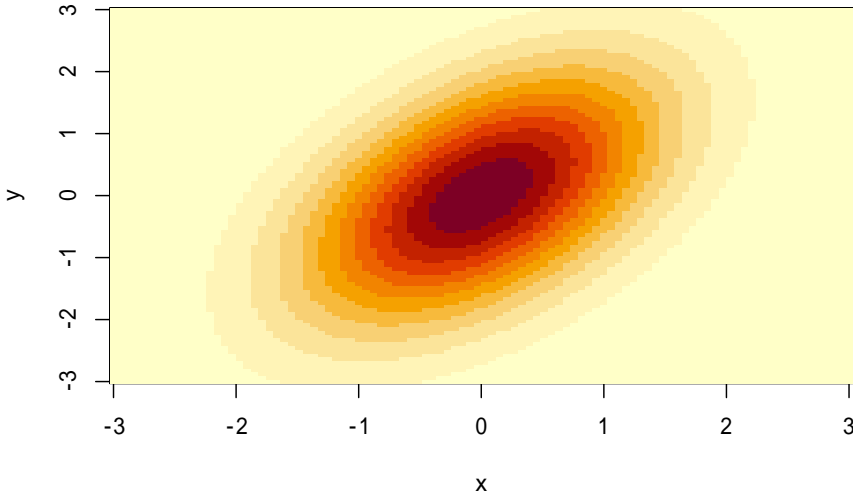
```
persp(x, y, z)
```



```
contour(x, y, z)
```



```
image(x, y, z)
```



### 29.1.8 动态三维图

rgl 包能制作动态的三维散点图与曲面图。

```
library(rgl)
```

iris 数据框包含了 3 种鸢尾花的各 50 个样品的测量值，测量值包括花萼长、宽，花瓣长、宽。用 rgl 的 plot3d() 作动态三维散点图如下：

```
with(iris, plot3d(
 Sepal.Length, Sepal.Width, Petal.Length,
 type="s", col=as.numeric(Species)))
```

这个图可以用鼠标拖动旋转。其中 type="s" 表示绘点符号是球体形状。还可选 "p"(点)、"l"(连线)、"h"(向 z=0 连线)。可以用 size= 指定大小倍数（缺省值为 3）。

用 rgl 的 persp3d() 函数作曲面图。如二元正态分布密度曲面：

```
x <- seq(-3,3, length=100)
y <- x
f <- function(x,y,ssq1=1, ssq2=2, rho=0.5){
 det1 <- ssq1*ssq2*(1 - rho^2)
 s1 <- sqrt(ssq1)
 s2 <- sqrt(ssq2)
 1/(2*pi*sqrt(det1)) * exp(-0.5 / det1 * (
 ssq2*x^2 + ssq1*y^2 - 2*rho*s1*s2*x*y))
}
z <- outer(x, y, f)
persp3d(x=x, y=y, z=z, col='red')
```

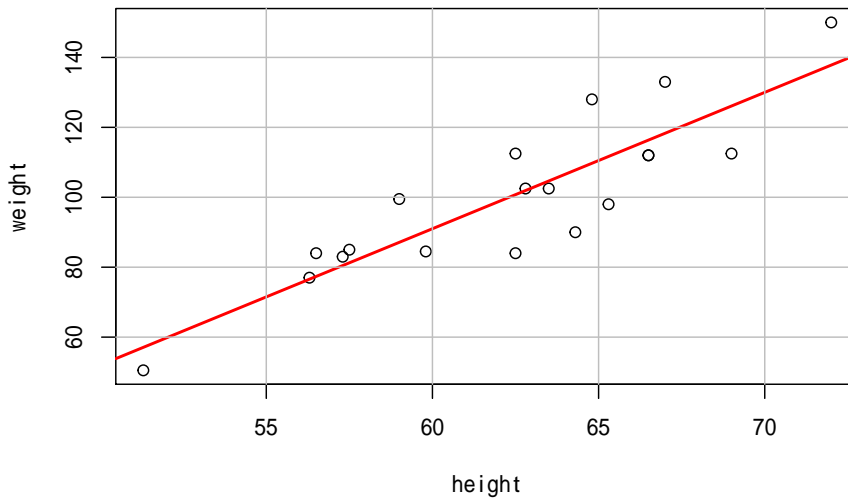
rgl 也有低级图形函数支持向已有图形添加物体、文字等，也支持并列多图。适当设置可以在 R Markdown 生成的 HTML 结果中动态显示三维图。

## 29.2 低级图形函数

### 29.2.1 abline()

用 `abline` 函数在图中增加直线。可以指定截距和斜率，或为竖线指定横坐标 (用参数 `v`)，为水平线指定纵坐标 (用参数 `h`)。如

```
with(d.class, plot(height, weight))
abline(-143, 3.9, col="red", lwd=2)
abline(v=c(55,60,65,70), col="gray")
abline(h=c(60,80,100,120,140), col="gray")
```

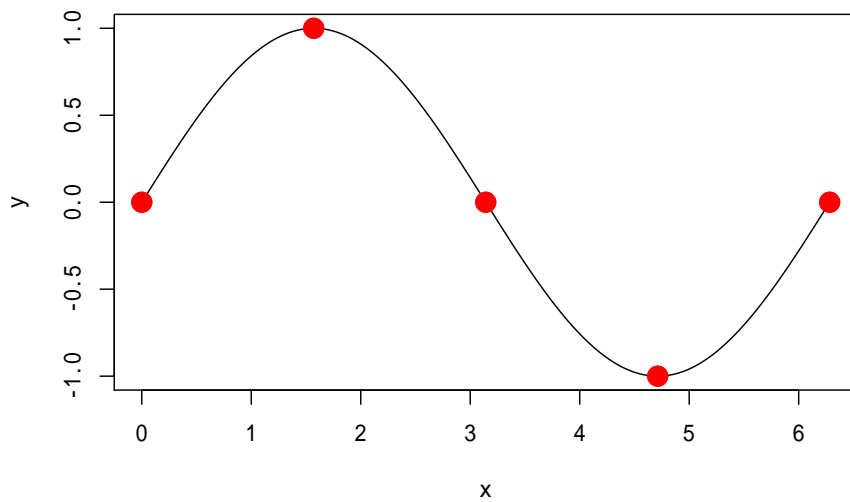


### 29.2.2 points()

用 `points` 函数增加散点，如：

```
x <- seq(0, 2*pi, length=200)
y <- sin(x)
special <- list(x=(0:4)*pi/2, y=sin((0:4)*pi/2))
```

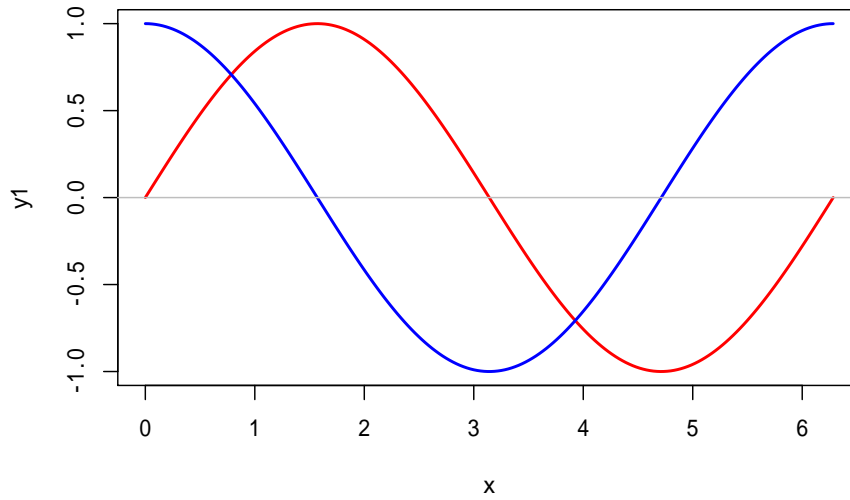
```
plot(x, y, type='l')
points(special$x, special$y,
 col="red", pch=16, cex=2)
points(special, col="red", pch=16, cex=2)
```



### 29.2.3 lines()

用 `lines` 函数增加曲线，如：

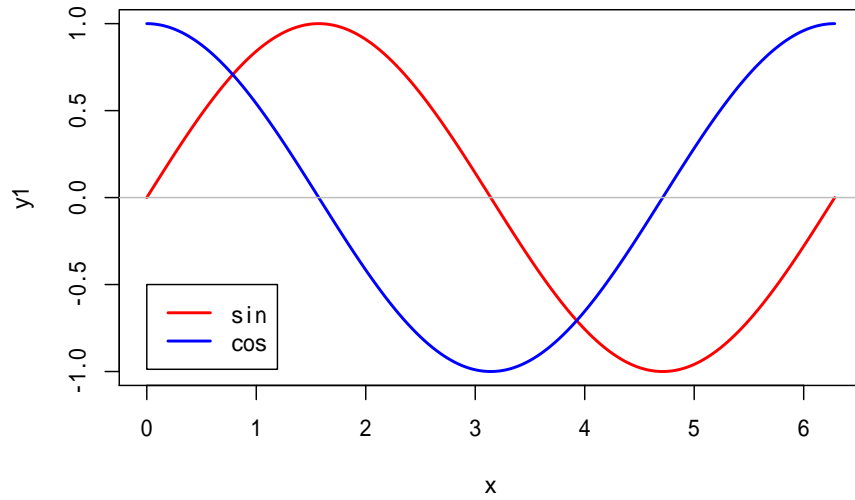
```
x <- seq(0, 2*pi, length=200)
y1 <- sin(x)
y2 <- cos(x)
plot(x, y1, type='l', lwd=2, col="red")
lines(x, y2, lwd=2, col="blue")
abline(h=0, col='gray')
```



#### 29.2.4 图例

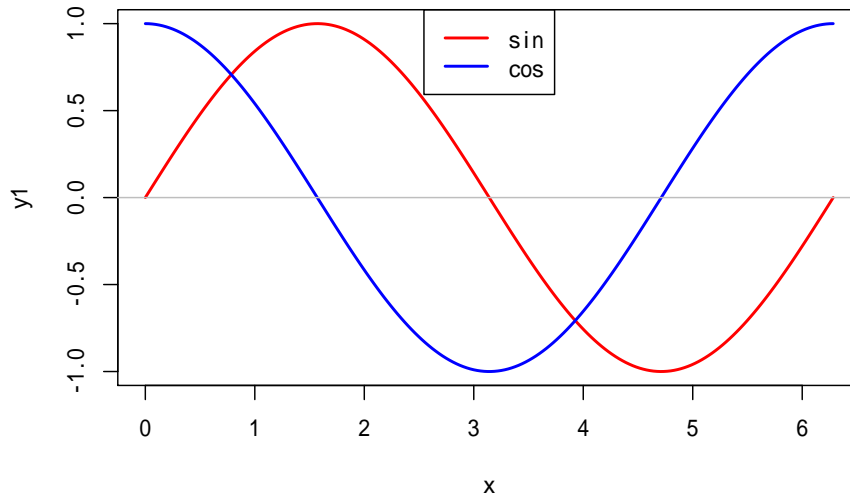
可以用 `legend` 函数增加标注, 如

```
x <- seq(0, 2*pi, length=200)
y1 <- sin(x)
y2 <- cos(x)
plot(x, y1, type='l', lwd=2, col="red")
lines(x, y2, lwd=2, col="blue")
abline(h=0, col='gray')
legend(0, -0.5, col=c("red", "blue"),
 lty=c(1,1), lwd=c(2,2),
 legend=c("sin", "cos"))
```



```
x <- seq(0, 2*pi, length=200)
y1 <- sin(x)
y2 <- cos(x)
plot(x, y1, type='l', lwd=2, col="red")
lines(x, y2, lwd=2, col="blue")
abline(h=0, col='gray')
legend('top', col=c("red", "blue"),
 lty=c(1,1), lwd=c(2,2),
 legend=c("sin", "cos"))
```

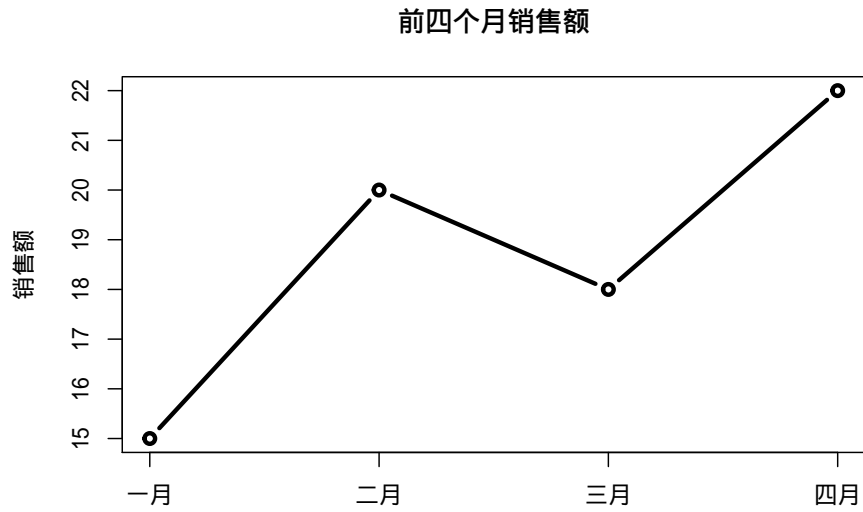




### 29.2.5 axis()

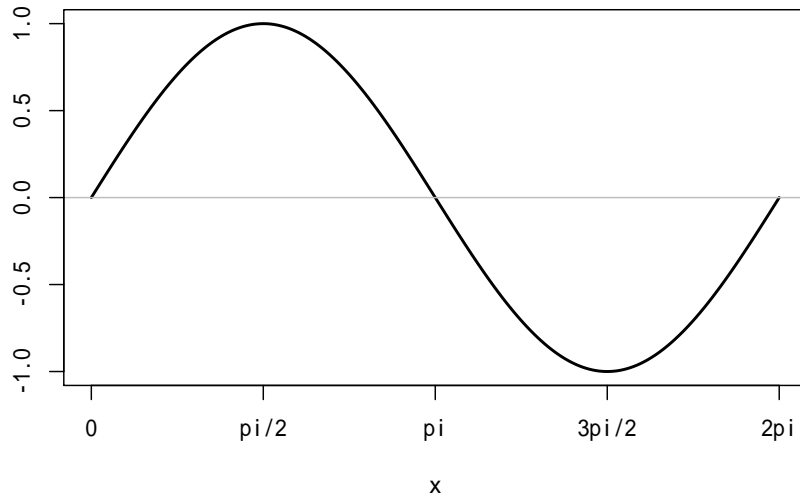
在 `plot()` 函数中用 `axes=FALSE` 可以取消自动的坐标轴。用 `box()` 函数画坐标边框。用 `axis` 函数单独绘制坐标轴。`axis` 的第一个参数取 1, 2, 3, 4, 分别表示横轴、纵轴、上方和右方。`axis` 的参数 `at` 为刻度线位置, `labels` 为标签。如

```
x <- c(' 一月'=15, ' 二月'=20,
 ' 三月'=18, ' 四月'=22)
plot(seq(along=x), x, axes=FALSE,
 type='b', lwd=3,
 main=' 前四个月销售额',
 xlab='', ylab=' 销售额')
box(); axis(2)
axis(1, at=seq(along=x), labels=names(x))
```



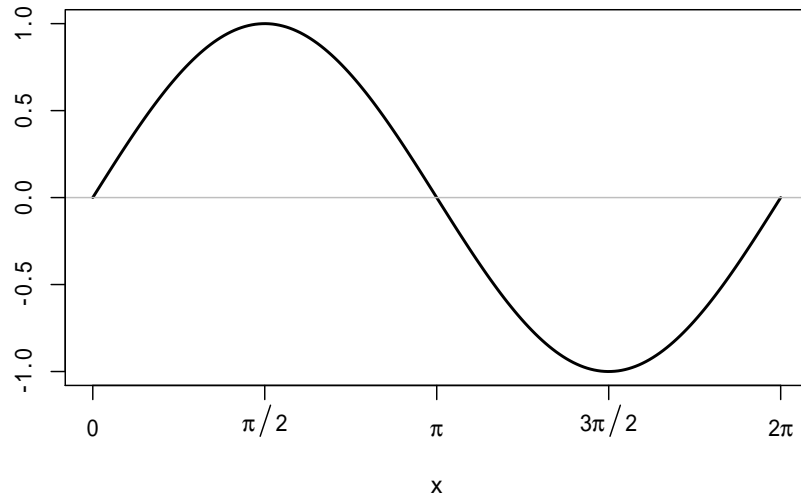
R 基本绘图支持少量的数学公式显示功能，如不用数学符号时：

```
x <- seq(0, 2*pi, length=200)
y1 <- sin(x)
plot(x, y1, type='l', lwd=2,
 axes=FALSE,
 xlab='x', ylab='')
abline(h=0, col='gray')
box()
axis(2)
axis(1, at=(0:4)/2*pi,
 labels=c('0', 'pi/2', 'pi', '3pi/2', '2pi'))
```



使用数学符号时:

```
x <- seq(0, 2*pi, length=200)
y1 <- sin(x)
plot(x, y1, type='l', lwd=2,
 axes=FALSE,
 xlab='x', ylab='')
abline(h=0, col='gray')
box()
axis(2)
axis(1, at=(0:4)/2*pi,
 labels=c(0, expression(pi/2),
 expression(pi), expression(3*pi/2),
 expression(2*pi)))
```



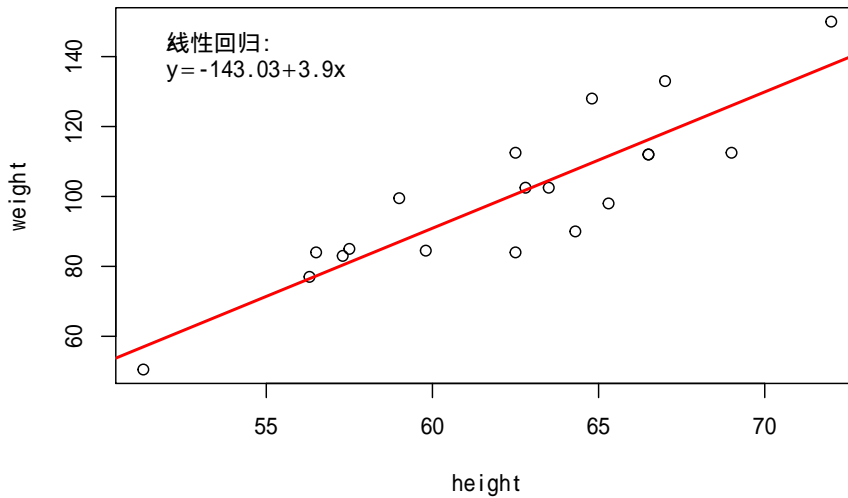
绘图中使用数学符号的演示：

```
demo(plotmath)
```

### 29.2.6 text()

`text()` 在坐标区域内添加文字。`mtext()` 在边空处添加文字。如

```
with(d.class, plot(height, weight))
lm1 <- lm(weight ~ height, data=d.class)
abline(lm1, col='red', lwd=2)
a <- coef(lm1)[1]
b <- coef(lm1)[2]
text(52, 145, adj=0, '线性回归:')
text(52, 140, adj=0,
 substitute(hat(y) == a + b*x,
 list(a=round(coef(lm1)[1], 2),
 b=round(coef(lm1)[2], 2))))
```



### 29.2.7 locator() 和 identify()

`locator()` 函数在执行时等待用户在图形的坐标区域内点击并返回点击处的坐标。可以用参数 `n` 指定要点击的点的个数。不指定个数则需要用右键菜单退出。这个函数也可以用来要求用户点击以进行到下一图形。如

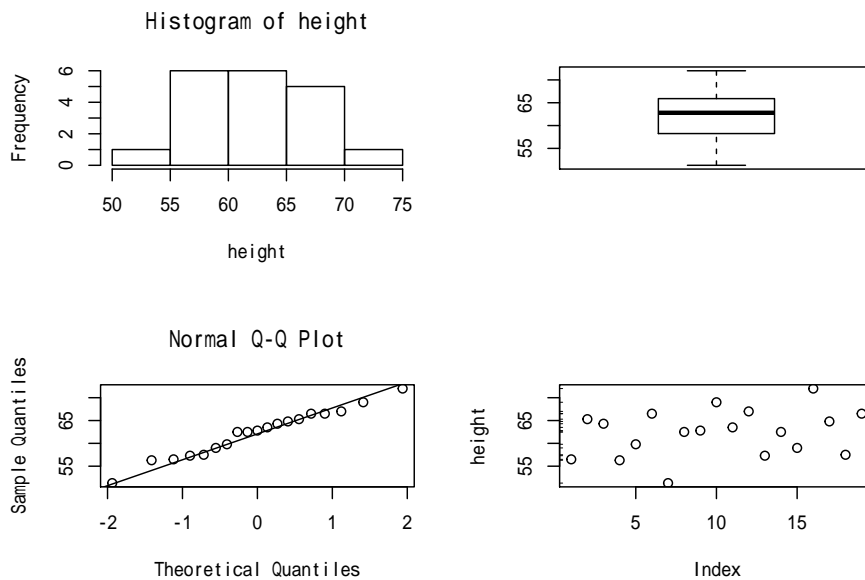
```
x <- seq(0, 2*pi, length=200)
y1 <- sin(x); y2 <- cos(x)
plot(x, y1, type='l',
 col="red")
lines(x, y2, col="blue")
legend(locator(1), col=c("red", "blue"),
 lty=c(1,1), legend=c("sin", "cos"))
```

`identify()` 可以识别点击处的点并标注标签，格式为 `identify(x, y, labels)`，其中 `(x,y)` 给出可点击的点的坐标，`labels` 是每个点对应的标签，点击那个点就在那个点旁边标对应的标签。

### 29.3 图形参数

用图形参数可以选择点的形状、颜色、线型、粗细、坐标轴做法、边空、一页多图等。有些参数直接用在绘图函数内，如 `plot` 函数可以用 `pch`、`col`、`cex`、`lty`、`lwd` 等参数。有些图形参数必须使用 `par()` 函数指定。`par` 函数指定图形参数并返回原来的参数值，所以在修改参数值作图后通常应该恢复原始参数值，做法如

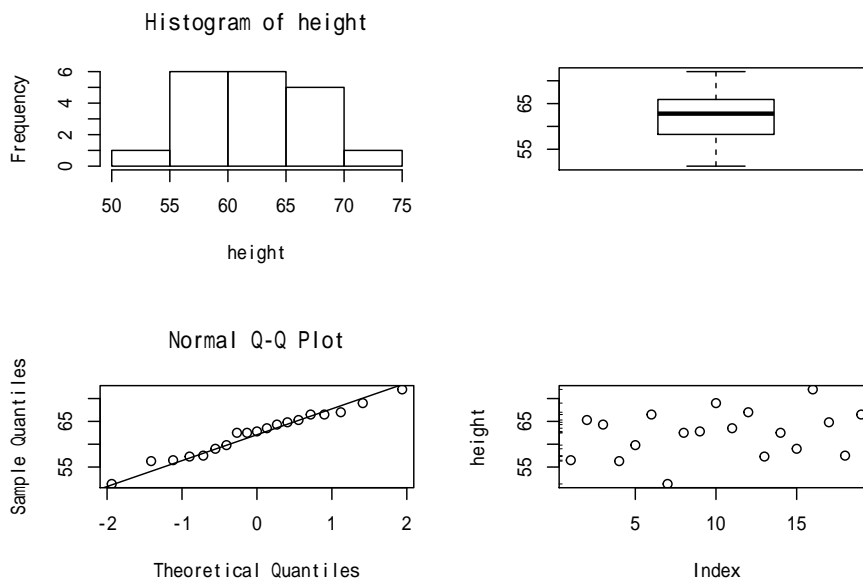
```
opar <- par(mfrow=c(2,2))
with(d.class, {hist(height);
 boxplot(height);
 qqnorm(height); qqline(height);
 plot(height); rug(height,side=2)})
```



```
par(opar)
```

在函数内，可以在函数开头修改了图形参数后，用 `on.exit()` 函数将恢复原始图形参数作为函数退出前必须完成的任务，如

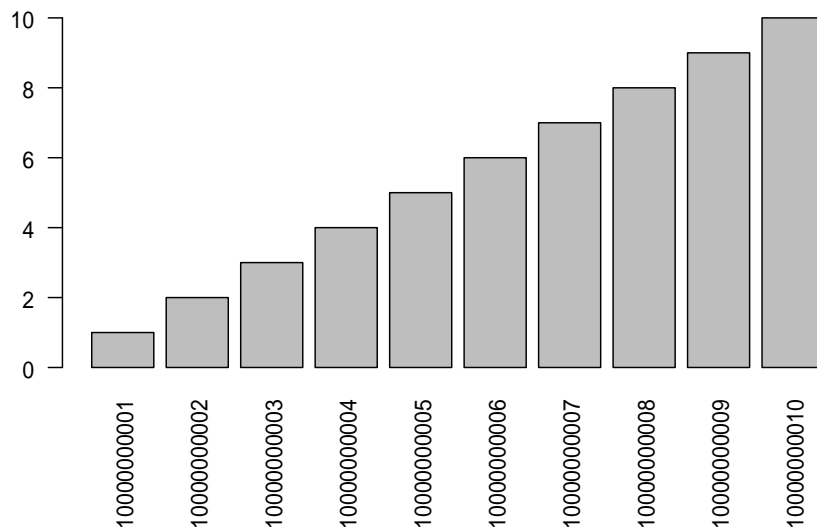
```
f <- function(){
 opar <- par(mfrow=c(2,2)); on.exit(par(opar))
 with(
 d.class,
 {hist(height);
 boxplot(height);
 qqnorm(height); qqline(height);
 plot(height); rug(height,side=2)
 })
}
```



### 29.3.1 例子：用图形参数解决 barplot 图形横坐标值过宽

barplot 的横坐标标注太宽时，自动将某些标注省略。用 `las=2` 指定坐标轴刻度标签垂直于坐标轴，这样 x 轴的刻度值就变成了纵向的。注意使用 `mar` 参数增加横坐标边空大小。例如

```
f <- function(){
 opar <- par(mar=c(8, 4, 2, 0.5)); on.exit(par(opar))
 x <- 1:10
 names(x) <- paste(10000000000 + (1:10))
 barplot(x, las=2)
}
f()
```



图形参数可以分为如下四类

- 图形元素控制;
- 坐标轴与坐标刻度;
- 图形边空;
- 一页多图。

### 29.3.2 图形元素控制

- `pch=16` 参数。散点符号，取 0 ~ 18 的数。
- `lty=2` 参数。线型，1 为实线，从 2 开始为各种虚线。



- `lwd=2` 参数, 线的粗细, 标准粗细为 1。
- `col="red"` 参数, 颜色, 可以是数字 1 ~ 8, 或颜色名字符串如"red", "blue" 等。用 `colors()` 函数查询有名字的颜色。用 `rainbow(n)` 函数产生连续变化的颜色。
- `font=2` 参数, 字体, 一般 `font=1` 是正体, 2 是粗体, 3 是斜体, 4 是粗斜体。
- `adj=-0.1` 指定文本相对于给定坐标的对齐方式。取 0 表示左对齐, 取 1 表示右对齐, 取 0.5 表示居中。此参数的值实际代表的是出现在给定坐标左边的文本的比例。
- `cex=1.5` 绘点符号大小倍数, 基本值为 1。

### 29.3.3 坐标轴与坐标刻度

- `mgp=c(3,1,0)` 坐标轴的标签、刻度值、坐标轴线到实际的坐标轴位置的距离, 以行高为单位。经常用来缩小坐标轴所占的空间, 如 `mgp=c(1.5, 0.5, 0)`。
- `lab=c(5,7,12)` 提供刻度线多少的建议, 第一个数为 x 轴刻度线个数, 第二个数为 y 轴刻度线个数, 第三个数是坐标刻度标签的字符宽度。
- `las=1` 坐标刻度标签的方向。0 表示总是平行于坐标轴, 1 表示总是水平, 2 表示总是垂直于坐标轴。
- `tck=0.01` 坐标轴刻度线长度, 以绘图区域大小为单位 1。
- `xaxs="s", yaxs="e"`: 控制 x 轴和 y 轴标刻度的方法。

取" s "(即 standard) 或" e "(即 extended) 的时候数据范围控制在最小刻度和最大刻度之间。取" e "时如果有数据点十分靠近边缘轴的范围会略微扩大。

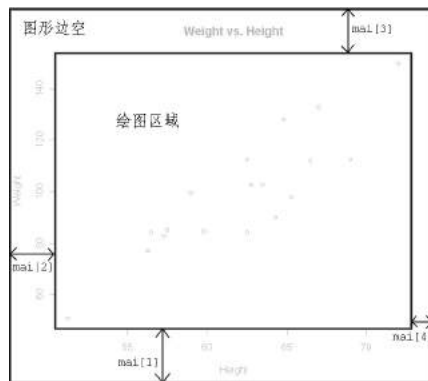
取值为" i "(即 internal) 或" r "(此为缺省) 使得刻度线都落在数据范围内部, 而" r "方式所留的边空较小。

取值设为" d "时会锁定此坐标轴, 后续的图形都使用与它完全相同的坐标轴, 这在要生成一系列可比较的图形的时候是有用的。要解除锁定需要把

这个图形参数设为其它值。

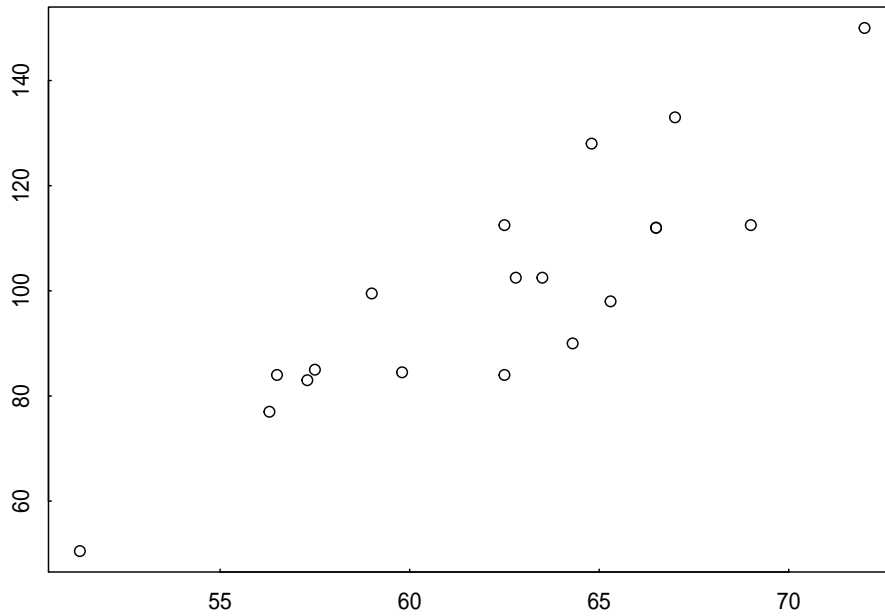
### 29.3.4 图形边空

一个单独的图由绘图区域 (绘图的点、线等画在这个区域中) 和包围绘图区域的边空组成, 边空中可以包含坐标轴标签、坐标轴刻度标签、标题、小标题等, 绘图区域一般被坐标轴包围。



边空的大小由 `mai` 参数或 `mar` 参数控制, 它们都是四个元素的向量, 分别规定下方、左方、上方、右方的边空大小, 其中 `mai` 取值的单位是英寸, 而 `mar` 的取值单位是文本行高度。例如:

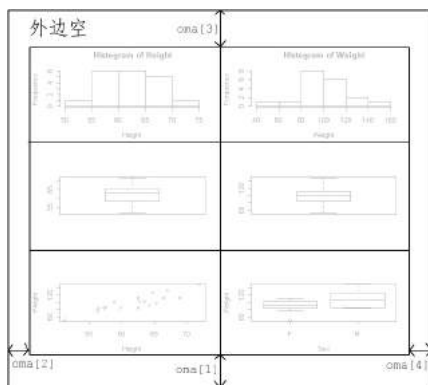
```
opar <- par(mar=c(2,2,0.5,0.5),
 mgp=c(0.5, 0.5, 0), tck=0.005)
with(d.class, plot(height, weight,
 xlab='', ylab=''))
```



`par(opar)`

### 29.3.5 一页多图

R 可以在同一页面开若干个按行、列排列的窗格, 在每个窗格中可以作一幅图。每个图有自己的内边空, 而所有图的外面可以包一个“外边空”。

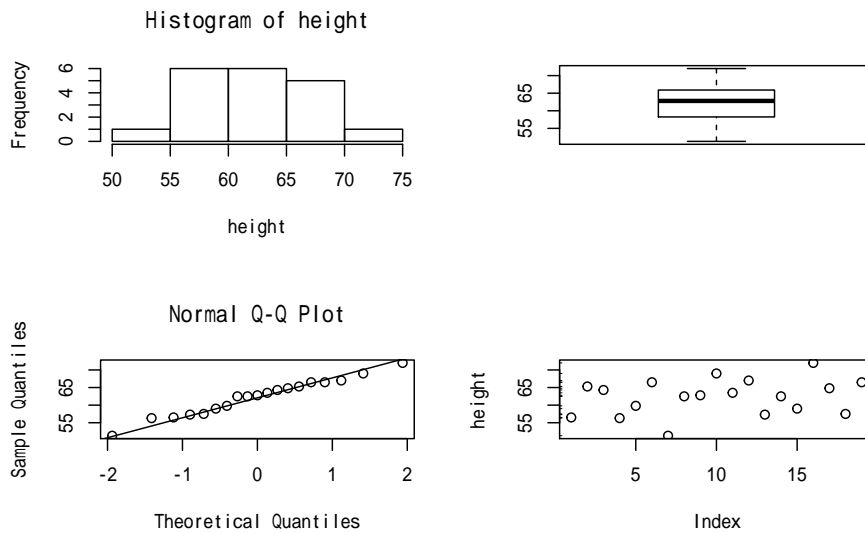


一页多图用 `mfrow` 参数或 `mfcol` 参数规定。用 `oma` 指定四个外边空的行数。用 `mtext` 加 `outer=T` 指定在外边空添加文本。如果没有 `outer=T` 则在内边

空添加文本。如

```
opar <- par(mfrow=c(2,2),
 oma=c(0,0,2,0))
with(d.class, {hist(height);
 boxplot(height);
 qqnorm(height); qqline(height);
 plot(height); rug(height,side=2)})
mtext(side=3, text=' 身高分布', cex=2, outer=T)
```

## 身高分布



```
par(opar)
```

## 29.4 图形输出

只要启用了高级绘图函数会自动选用当前绘图设备，缺省为屏幕窗口。

### 29.4.1 PDF 输出

用 `pdf` 函数可以指定输出到 PDF 文件。如

```
pdf(file='fig-hw.pdf', height=10/2.54,
 width=10/2.54, family='GB1')
with(d.class, plot(height, weight,
 main=' 体重与身高关系'))
dev.off()
```

用 `dev.off()` 关闭当前设备并生成输出文件（如果是屏幕窗口则没有保存结果）。

### 29.4.2 PNG 输出

```
png(file='fig-hw.png', height=1000, width=1000)
with(d.class, plot(height, weight,
 main=' 体重与身高关系'))
dev.off()
```

类似地，用 `jpeg()` 函数启用 JPEG 图形设备，用 `bmp()` 函数启用 BMP 图形设备，用 `postscript()` 函数启用 PostScript 图形设备。

## 29.5 包含多种中文字体的图形

为了使用 MS Windows 系统字体，一种办法是安装 `showtext` 包。该包替换画图时的添加文本函数命令，把文本内容替换成多边形（PDF 或 PS 图）或点阵（点阵图）。

需要的工作：

- 查看 Windows 的 `font` 目录内容，看文件名与字体名的对应关系。下面的程序中的列表是我的中文 Windows 10 的部分中文字体。
- 找到自己希望使用的中文字体的文件名。

- 用 `font.add()` 命令, 增加一套自定义字体 `family`, 一套中可以指定四种: 常规 (`regular`), 粗体 (`bold`), 斜体 (`italic`), 粗斜体 (`bolditalic`)
- 程序中调入 `showtext` 包并运行 `showtext.auto()` 命令, 这个命令使得文本命令采用 `showtext` 包
- 用 `par(family=)` 指定自定义的字体 `family`。
- 作图 (主要是 PDF)。关闭图形设备。

```
test.chinese <- function(){
 require(showtext); showtext.auto()

 ## 建立文件名到字体名对照表
 fmap <- c(
 'msyh'=' 微软雅黑常规',
 'msyhbd'=' 微软雅黑粗体',
 'msyhl'=' 微软雅黑细体',
 'simsun'=' 宋体',
 'simfang'=' 仿宋',
 'simkai'=' 楷体',
 'simhei'=' 黑体',
 'SIMLI'=' 隶书',
 'SIMYOU'=' 幼圆',
 'STSONG'=' 华文宋体',
 'STZHONGS'=' 华文中宋',
 'STFANGSO'=' 华文仿宋',
 'STKAITI'=' 华文楷体',
 'STXIHEI'=' 华文细黑',
 'STLITI'=' 华文隶书',
 'STXINGKA'=' 华文行楷',
 'STXINWEI'=' 华文新魏',
 'STCAIYUN'=' 华文彩云',
 'STHUPO'=' 华文琥珀'
)

 fmapr <- names(fmap); names(fmapr) <- unname(fmap)
 cat('==== 字体文件名与字体名称对应:\n')
```

```
print(fmap)
cat('==== 字体名与字体文件名对应:\n')
print(fmapr)

找到某个字体的字体文件
font.name 是字体名称
find.font <- function(font.name){
 fname <- fmapr[font.name]
 flist0 <- font.files()
 flist1 <- sapply(strsplit(flist0, '[.]'), function(it) it[1])
 flist0[flist1==fname]
}
ff1 <- find.font(' 宋体')
ff2 <- find.font(' 黑体')
ff3 <- find.font(' 仿宋')
ff4 <- find.font(' 隶书');

font.add('cjk4',
 regular=ff1,
 bold=ff2,
 italic=ff3,
 bolditalic=ff4)

##browser()

pdf('test-chinese.pdf'); on.exit(dev.off())
par(family='cjk4')

plot(c(0,1), c(0,1), type='n',
 axes=FALSE, xlab='', ylab='')
text(0.1, 0.9, ' 正体', font=1)
text(0.1, 0.8, ' 粗体', font=2)
text(0.1, 0.7, ' 斜体', font=3)
text(0.1, 0.6, ' 粗斜体', font=4)
```

```
}
test.chinese()
```

注意：图形参数 font=1 表示正体, font=2 表示粗体, font=3 表示斜体, font=4 表示粗斜体。

## 29.6 其它图形

### 29.6.1 相关系数图

用 `cor(x)` 可以计算数据框 `x` 的各列的相关系数阵。`corrgram` 包的 `corrgram()` 函数可以将相关系数阵用图形表示，系数绝对值大小用色块颜色深浅表示，正负用两种颜色区分。

例如，计算 `iris` 数据框中四个测量值的相关系数并用矩阵表示：

```
library(corrgram)

Registered S3 method overwritten by 'seriation':
method from
reorder.hclust gclus

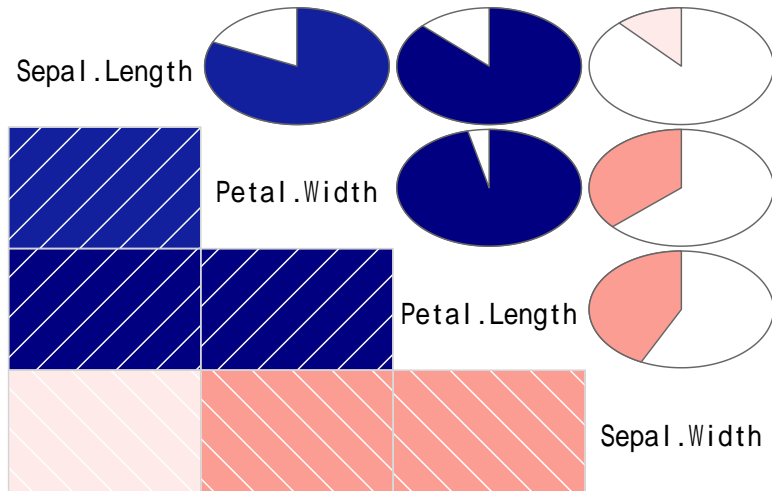
R.iris <- cor(iris[,1:4])
print(round(R.iris, 2))

Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length 1.00 -0.12 0.87 0.82
Sepal.Width -0.12 1.00 -0.43 -0.37
Petal.Length 0.87 -0.43 1.00 0.96
Petal.Width 0.82 -0.37 0.96 1.00

corrgram(
 R.iris, order=TRUE,
 lower.panel=panel.shade,
 upper.panel = panel.pie,
 text.panel = panel.txt
```



)



左下方用颜色代表正负相关，蓝色为正相关，红色为负相关，可以看出花萼长 (Sepal.Length)、花瓣宽 (Petal.Width) 和花瓣长 (Petal.Length) 相互为较强的正相关，但是花萼宽与其它三个变量为负相关。这个相关结果实际是虚假的，因为样本不是单一总体而是来自三个总体。图形右上方用阴影部分大小和颜色深度代表相关系数绝对值，用颜色区分正负。



# Chapter 30

## ggplot 作图入门

### 30.1 介绍

Hadley Wickem 的 `ggplot2` 包是 R 的一个作图用的扩展包，它实现了“图形的语法”，将一个作图任务分解为若干个子任务，只要完成各个子任务就可以完成作图。在作常用的图形时，只需要两个步骤：首先将图形所展现的数据输入到 `ggplot()` 函数中，然后调用某个 `geom_xxx()` 函数，指定图形类型，如散点图、曲线图、盒形图等。

如果需要进一步控制图形细节，只要继续调用其它函数，就可以控制变量值的表现方式 (`scale`)、图例、配色等。这使得我们很容易做出基本的图形，在有需要时再深入学习，做出更为满意的图形。

`ggplot2` 的作图一般步骤为：

- 准备数据，一般为数据框，且一般为长表，即每个观测时间占一行，每个观测变量占一列。
- 将数据输入到 `ggplot()` 函数中，并指定参与作图的每个变量分别映射到哪些图形特性，比如映射为 x 坐标、y 坐标、颜色、形状等。这些映射称为 `aesthetic mappings` 或 `aesthetics`。
- 选择一个合适的图形类型，函数名以 `geom_` 开头，如 `geom_point()` 表示散点图。图形类型简称为 `geom`。将 `ggplot()` 部分与 `geom_xxx()` 部

分用加号连接。到此已经可以作图，下面的步骤是进一步的细化设定。

- 设定适当的坐标系统，如 `coord_cartesian()`, `scale_x_log10()` 等。仍用加号连接。
- 设定标题和图例位置等，如 `labs()`。仍用加号连接。

这个流程的一个大致的模板为：

```
p <- ggplot(data=<输入数据框>,
 mapping=aes(<维度>=<变量名>,
 <维度>=<变量名>, <...>))
p + geom_<图形类型>(<...>) +
 scale_<映射>_<类型>(<...>) +
 coord_<类型>(<...>) +
 labs(<...>)
```

其中 `<...>` 表示额外的选项。变量 `p` 包含做出的图形的所有数据与设定，变量名可以任意取。

本章内容主要来自：

- Healy, Kieran (2018). *Data Visualization: A Practical Introduction*. Princeton University Press. <https://socviz.co/index.html> 这本书讲了 R 的 ggplot 的使用，也讲了一些可视化的一般性原则。
- Claus O. Wilke(2019). *Fundamentals of Data Visualization*. O'Reilly Media. <https://serialmentor.com/dataviz/> 这本书虽然也使用 R 的 ggplot2 包，但正文中没有代码，主要讲作图有哪些考虑、各种图形类型。代码在 github 上可下载。
- Winston Chang(2018). *R Graphics Cookbook*. O'Reilly Media. 网站：<https://r-graphics.org/> 为第二版。讲了各种图的 R 程序。
- Wickham, Hadley (2016). *Ggplot2: Elegant graphics for data analysis*. New York: Springer.

Wickham 的书主要需要安装 tidyverse 扩展包，安装时会自动安装其它一些有关扩展包。Healy 的书需要通过如下程序安装 socviz 软件包：

```
devtools::install_github("kjhealy/socviz")
```

后续的例子中用到一些数据集：

- 来自 `gapminder` 扩展包的 `gapminder` 数据集，有若干个国家不同年份的一些数据，包括所属洲、期望寿命、人口数、人均 GDP。有 1704 个观测和 6 个变量。
- `socviz` 包的 `gss_sm` 数据集，是 2016 年美国一般社会调查数据的部分内容。有 2867 个观测，32 个变量。社会调查数据的变量主要取属性值，比如无序分类、有序分类、分组的数值、整数值等。
- `socviz` 包的 `organdata` 数据集，是 17 个 OECD 国家历年的器官捐献情况以及一些其它记录。
- `socviz` 扩展包的 `elections_historic` 数据集。包括美国历次总统大选当选人、所属党派、支持比例等。
- `socviz` 扩展包的 `asasec` 数据集。这是美国社会学学会 (ASA) 的各分会 2005 年到 2015 年的一些数据。
- `ggplot2` 包中的 `midwest` 数据集包含了美国中西部的一些县的统计数据，如面积等。
- 来自 `ggplot2` 包的钻石数据集。

`gapminder` 的头部:

```
library(gapminder)
head(gapminder, 20)

A tibble: 20 x 6
country continent year lifeExp pop gdpPercap
<fct> <fct> <int> <dbl> <int> <dbl>
1 Afghanistan Asia 1952 28.8 8425333 779.
2 Afghanistan Asia 1957 30.3 9240934 821.
3 Afghanistan Asia 1962 32.0 10267083 853.
4 Afghanistan Asia 1967 34.0 11537966 836.
5 Afghanistan Asia 1972 36.1 13079460 740.
6 Afghanistan Asia 1977 38.4 14880372 786.
7 Afghanistan Asia 1982 39.9 12881816 978.
8 Afghanistan Asia 1987 40.8 13867957 852.
9 Afghanistan Asia 1992 41.7 16317921 649.
10 Afghanistan Asia 1997 41.8 22227415 635.
11 Afghanistan Asia 2002 42.1 25268405 727.
```

```
12 Afghanistan Asia 2007 43.8 31889923 975.
13 Albania Europe 1952 55.2 1282697 1601.
14 Albania Europe 1957 59.3 1476505 1942.
15 Albania Europe 1962 64.8 1728137 2313.
16 Albania Europe 1967 66.2 1984060 2760.
17 Albania Europe 1972 67.7 2263554 3313.
18 Albania Europe 1977 68.9 2509048 3533.
19 Albania Europe 1982 70.4 2780097 3631.
20 Albania Europe 1987 72 3075321 3739.
```

gss\_sm 的头部:

```
head(gss_sm, 20)
```

```
A tibble: 20 x 32
year id ballot age childs sibs degree race sex region income16
<dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <fct> <fct> <fct> <fct> <fct>
1 2016 1 1 47 3 2 Bache~ White Male New E~ $170000~
2 2016 2 2 61 0 3 High ~ White Male New E~ $50000 ~
3 2016 3 3 72 2 3 Bache~ White Male New E~ $75000 ~
4 2016 4 1 43 4 3 High ~ White Fema~ New E~ $170000~
5 2016 5 3 55 2 2 Gradu~ White Fema~ New E~ $170000~
6 2016 6 2 53 2 2 Junio~ White Fema~ New E~ $60000 ~
7 2016 7 1 50 2 2 High ~ White Male New E~ $170000~
8 2016 8 3 23 3 6 High ~ Other Fema~ Middl~ $30000 ~
9 2016 9 1 45 3 5 High ~ Black Male Middl~ $60000 ~
10 2016 10 3 71 4 1 Junio~ White Male Middl~ $60000 ~
11 2016 11 2 33 5 4 High ~ Black Fema~ Middl~ under $~
12 2016 12 1 86 4 4 High ~ White Fema~ Middl~ under $~
13 2016 13 2 32 3 3 High ~ Black Male Middl~ $8 000 ~
14 2016 14 3 60 5 6 High ~ Black Fema~ Middl~ $12500 ~
15 2016 15 2 76 7 0 Lt Hi~ White Male New E~ $40000 ~
16 2016 16 3 33 2 1 High ~ White Fema~ New E~ $50000 ~
17 2016 17 3 56 6 3 High ~ White Male New E~ $50000 ~
18 2016 18 2 62 5 8 Lt Hi~ Other Fema~ New E~ $5 000 ~
```

```
19 2016 19 2 31 0 2 Gradu~ Black Male New E~ $35000 ~
20 2016 20 1 43 2 0 High ~ Black Male New E~ $25000 ~
... with 21 more variables: relig <fct>, marital <fct>, padeg <fct>,
madeg <fct>, partyid <fct>, polviews <fct>, happy <fct>,
partners <fct>, grass <fct>, zodiac <fct>, pres12 <dbl>,
wtssall <dbl>, income_rc <fct>, agegrp <fct>, ageq <fct>,
siblings <fct>, kids <fct>, religion <fct>, bigregion <fct>,
partners_rc <fct>, obama <dbl>
```

## 30.2 作图的一般原则

关于什么是好的图形和坏的图形，William S. Cleveland, Edward R. Tufte 等人有很多的研究。

坏的图形可能有如下缺点：

- 坏的品味。统计图形应该用尽可能少的图形元素表示尽可能多的数据，从打印图形而言，即数据量与所用墨水比例越大越好。没有必要的颜色、三维形态经常会影响读者对图形的认读。这是 Edward R. Tufte 的观点，但是过于极端也不好。
- 坏的数据。即使图形本身的做法没有问题，选择了错误的或者不合适的数据也会误导读者，甚至于用错误数据做的很专业的图形会比粗陋的图形更能误导读者。
- 坏的感知。不好的颜色选择、三维形状、坐标轴范围、宽高比都有可能对读者的认知有影响。

作图时应考虑的一些因素：

- 数值型变量的不同值可以表示为：
  - 同一坐标轴上的不同位置、
  - 不同轴上的位置、
  - 不同长度、
  - 不同角度或者斜率、
  - 不同面积、
  - 三维空间中的不同位置、

- 颜色的不同明暗度、
- 不同颜色饱和度、
- 曲线的不同曲率、
- 三维体积、

这些表示的选择项越往后越难以被读者正确辨识。使用颜色时，应该使用渐变的明暗度或者渐变色。

- 分类变量的不同值可以表示为：
  - 不同分组、
  - 不同颜色、
  - 三维动态、
  - 不同符号。

这些表示的选择项越往后越难辨识。使用颜色时，应该使用明显不同的颜色而不应该使用渐变色。

- 对于最少是零的变量，是否应该以零作为坐标轴的最低值需要考虑，但没有一定的规则。同一组数据在不同的坐标范围或者长宽比下曲线的斜率会有很大差别。

## 30.3 散点图：ggplot 入门

### 30.3.1 基本的散点图

以 `gapminder` 数据集作为输入数据，做出简单的散点图，并逐步进行改善。这个数据集有多个国家在多个年份的期望寿命与人均 GDP 值，作期望寿命对人均 GDP 的散点图，每个国家的每个年份作为一个点。散点图最重要的映射是 x 轴与 y 轴两个维度。

首先调用 `ggplot()` 函数，指定数据集，将人均 GDP 映射到 x 轴，将期望寿命映射到 y 轴，结果保存为一个 R 变量：

```
p <- ggplot(data = gapminder,
 mapping = aes(
```

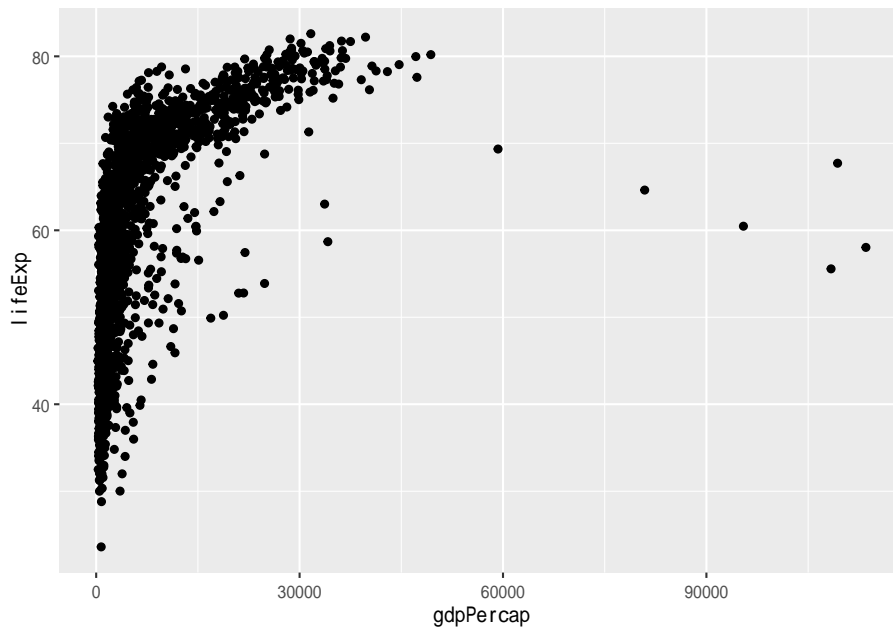


```
x = gdpPercap,
y = lifeExp))
```

x、y 轴是最常见的映射，也可以将变量映射为颜色、符号、线型等，这时不需要指定具体的颜色、符号、线型，而是将变量映射为这些图形元素类型。

在如上指定了数据和映射后，只要用 `geom_xxx()` 指定一个图形类型，并与 `ggplot()` 的结果用加号连接就可以作图了，如：

```
p + geom_point()
```



实际上，上面的程序等同于调用 `print(p + geom_point())`。在 R 函数中应该显示地调用 `print()`。

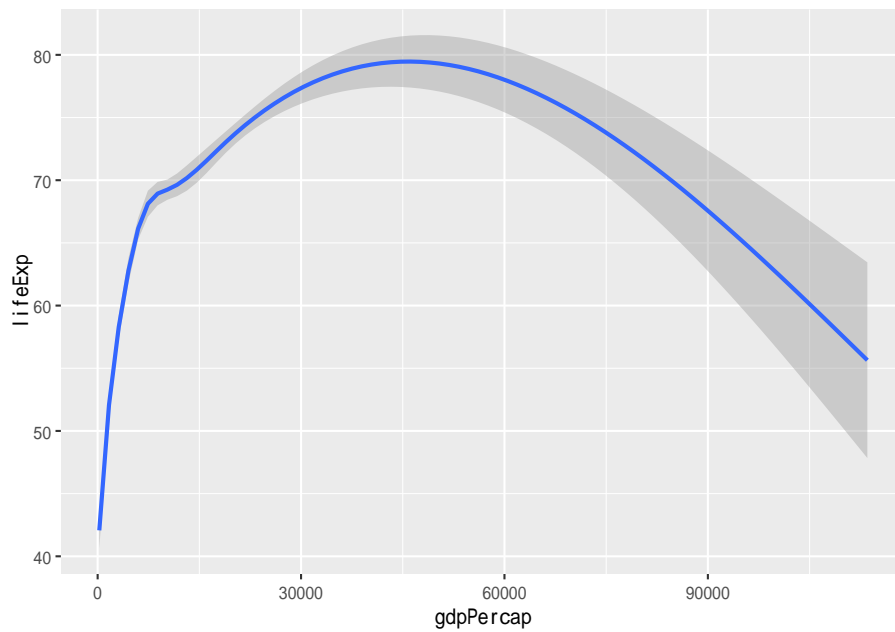
### 30.3.2 逐步改善

指定数据集、指定映射、选择适当的图形类型就可以做出基本的图形，随后可以逐步对坐标系、坐标系刻度、标签与图例、配色等进行改善。实际上，`ggplot2` 包已经提供了十分合理的预设值，用户只要进行一些必要的改动即可。

作图步骤之间用加号连接，这是 ggplot 包特有的语法。例如，用相同的映射做出拟合曲线图：

```
p + geom_smooth()
```

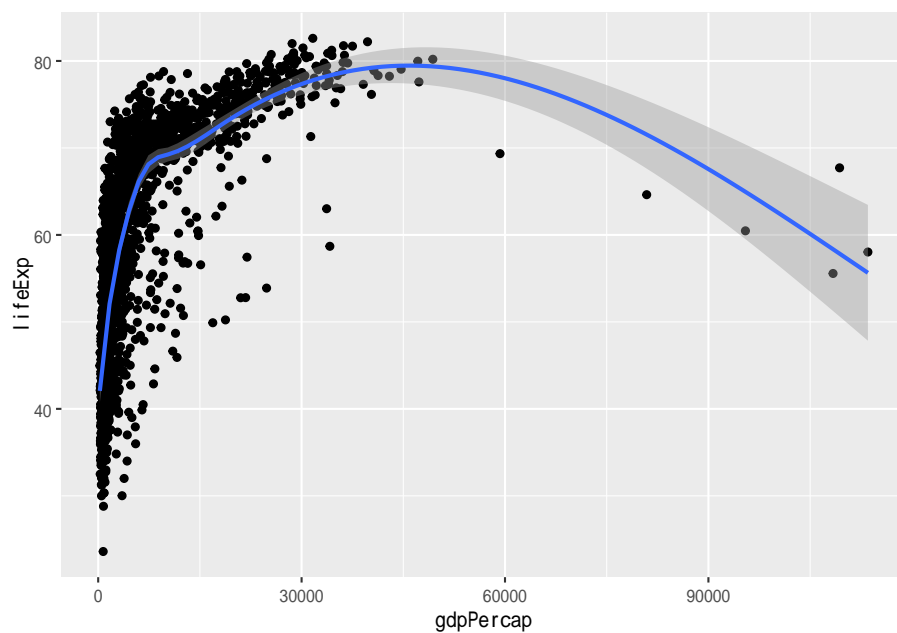
```
`geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```



用相同的映射做出散点图并叠加拟合曲线图：

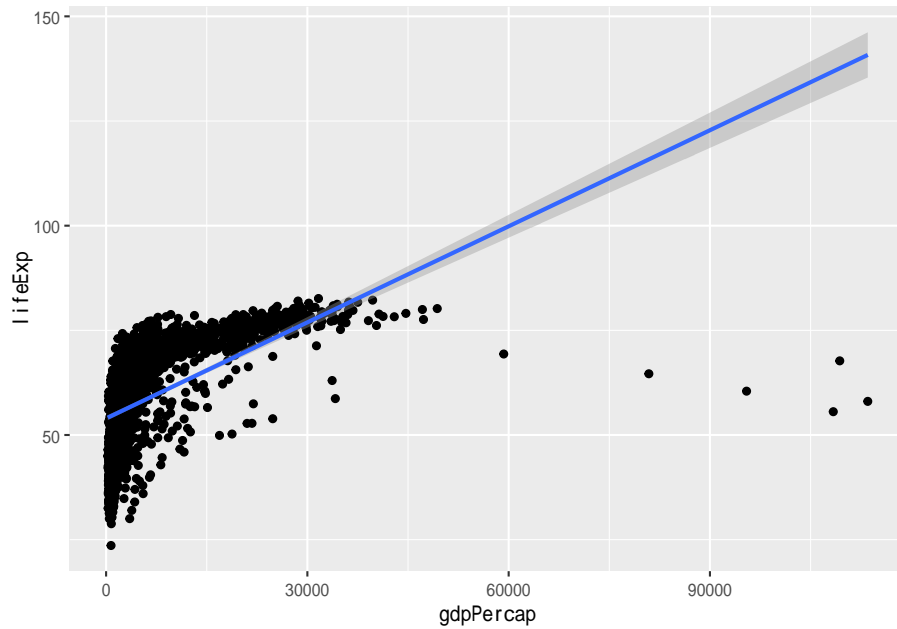
```
p + geom_point() + geom_smooth()
```

```
`geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```



`geom_smooth()` 的默认设置调用了 `gam()` 函数来拟合曲线, 可以用 `geom_smooth()` 的参数选择不同的拟合方法, 如直线拟合:

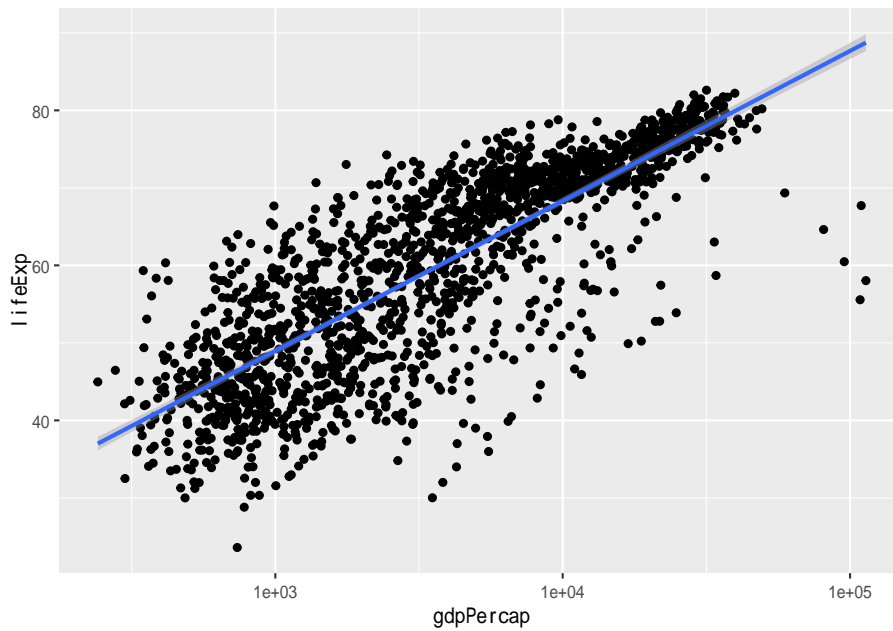
```
p + geom_point() + geom_smooth(method="lm")
```



注意 `geom_xxx()` 函数计算所需的变量值是从 `ggplot()` 函数保存在变量 `p` 中的信息提取的。

在以上的所有图形中，`x` 轴变量（人均 GDP）分布非正态，严重右偏，使得大多数散点重叠地分布在直角坐标系的左下角。将 `x` 轴用对数刻度可以改善，函数为 `scale_x_log10()`：

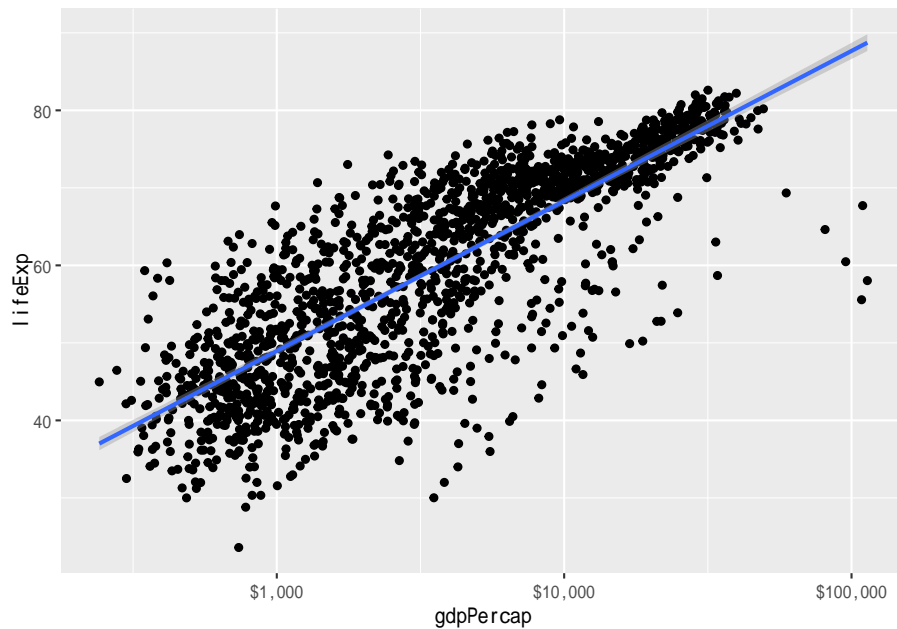
```
p + geom_point() +
 geom_smooth(method="gam") +
 scale_x_log10()
```



广义可加模型拟合的曲线基本是一条直线。注意，在进行拟合时，已经预先对人均 GDP 变量进行常用对数变换。

刚刚的图形的横坐标轴刻度不太友好，可以调用 `scales` 扩展包的适当函数进行改善，作为 `scale_x_log10()` 的 `labels` 选项：

```
p + geom_point() +
 geom_smooth(method="gam") +
 scale_x_log10(labels=scales::dollar)
```



`scale_xxx()` 的 `labels` 选项指定如何标出坐标刻度数字，参数值是一个函数对象，如果 `scales` 包中找不到适当的功能，可以自定义一个函数将数值转换为字符串。`scales` 包提供了 `comma`, `date`, `dollar`, `math`, `number`, `ordinal`, `pvalue`, `scientific`, `time` 等坐标刻度值转换函数。

### 30.3.3 颜色、符号、线型等映射

在 `ggplot()` 函数的 `mapping` 参数的 `aes()` 设定中将变量映射到 `x`、`y` 轴，颜色、符号、线型等图形元素类型，也可以作为图形设置将某些图形元素设置为固定值。

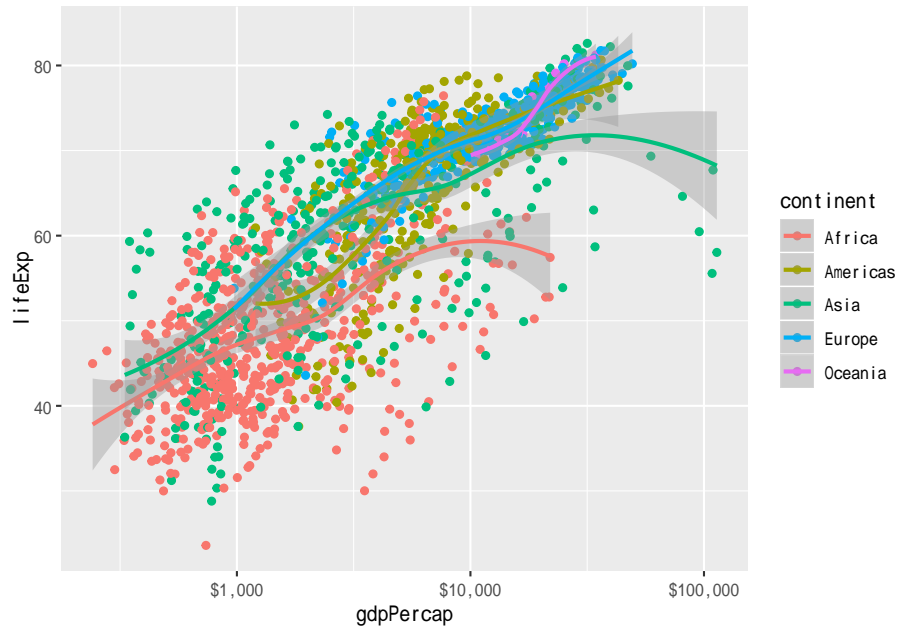
例如，用不同颜色表示不同大洲，就是将 `continent` 变量映射到 `color`:

```
p <- ggplot(data=gapminder,
 mapping = aes(
 x = gdpPerCap,
 y = lifeExp,
 color = continent))
```

程序中仅指定了将大洲映射到颜色维，并不具体指定所用的颜色。

作带有局部多项式曲线拟合的散点图:

```
p + geom_point() +
 geom_smooth(method="loess") +
 scale_x_log10(labels=scales::dollar)
```

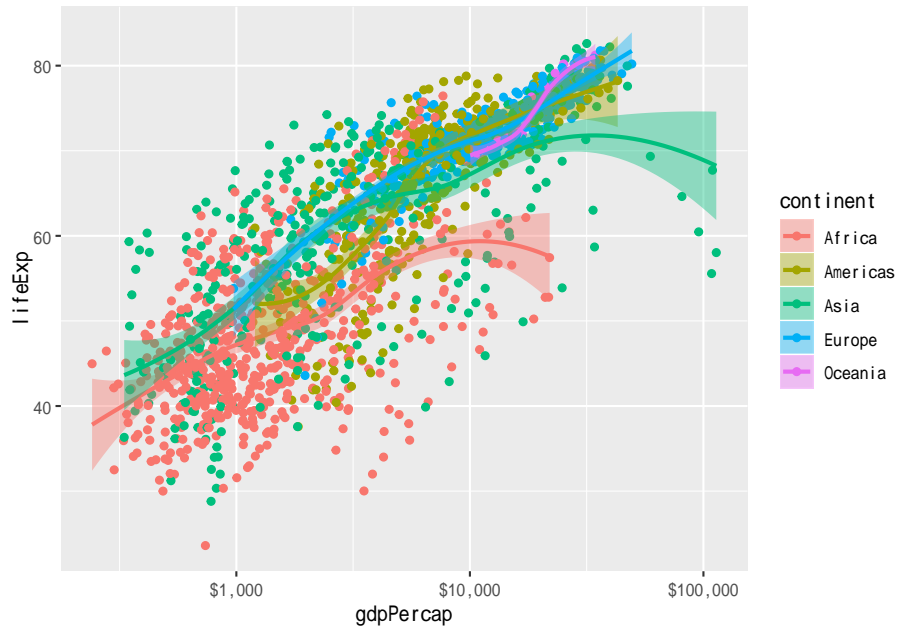


可以看出，不同散点用了不同颜色表示其 `continent` 变量的值，五个大洲分别进行了曲线拟合，使得图形比较复杂，难以认读。在图形右侧自动生成了颜色与 `continent` 变量值的对应关系图例。

下面的图形分不同大洲作曲线拟合，并将置信区间阴影的颜色也用不同大洲区分，方法是在 `aes()` 中将 `color` 和 `fill` 都指定为变量 `continent`:

```
p <- ggplot(data=gapminder,
 mapping = aes(
 x = gdpPercap,
 y = lifeExp,
 color = continent,
 fill = continent))
p + geom_point() +
 geom_smooth(method="loess") +
```

```
scale_x_log10(labels=scales::dollar)
```



尝试将颜色指定为一个固定值，如：

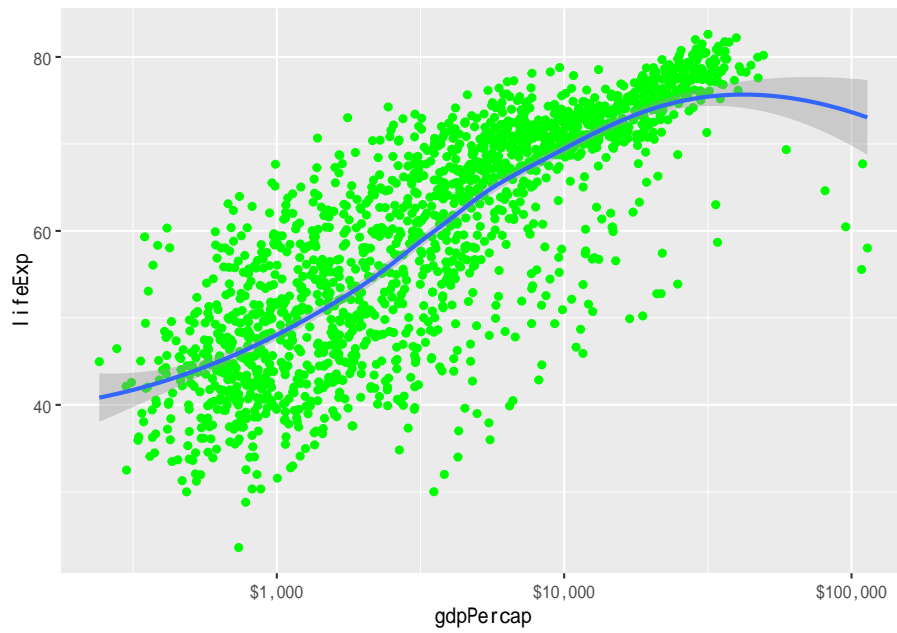
```
p <- ggplot(data=gapminder,
 mapping = aes(
 x = gdpPerCap,
 y = lifeExp,
 color = "green"))
p + geom_point() +
 geom_smooth(method="loess") +
 scale_x_log10(labels=scales::dollar)
```





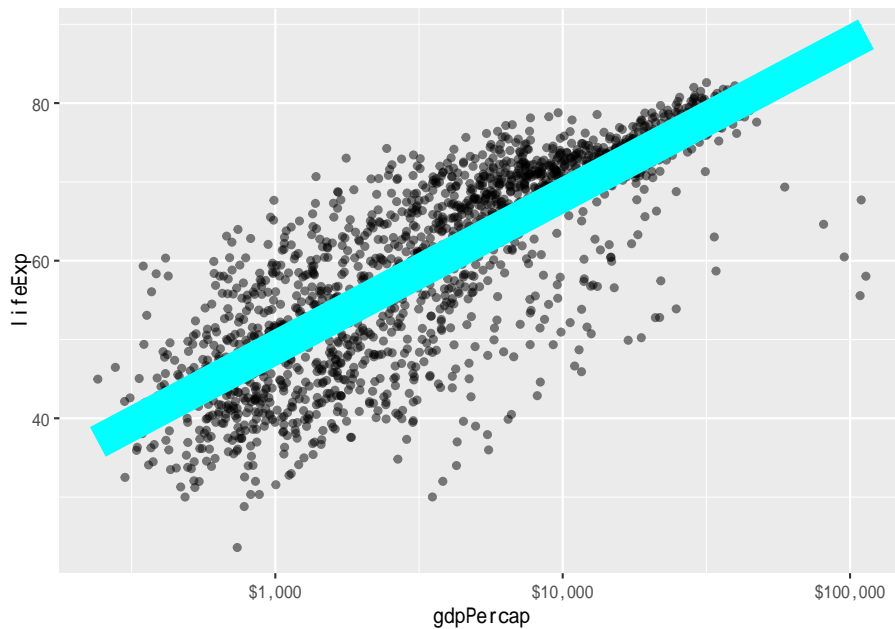
我们发现，散点并没有使用绿色，而且图形右侧有一个 `green` 图例。这是因为，`aes()` 仅用来指定变量与图形元素类型的映射，所以实际上是生成了一个仅有一个常数值 `"green"` 的新变量，用颜色表示这个新变量。为了指定固定颜色，可以作为 `geom_xxx()` 函数的选项，如：

```
p <- ggplot(data=gapminder,
 mapping = aes(
 x = gdpPercap,
 y = lifeExp))
p + geom_point(color="green") +
 geom_smooth(method="loess") +
 scale_x_log10(labels=scales::dollar)
```



`geom_xxx()` 函数接受许多关于颜色、透明度、符号、线型的设置参数。比如，下面的程序指定了散点的透明度，以及拟合直线的粗细：

```
p + geom_point(alpha=0.5) +
 geom_smooth(method="lm", color="cyan", se = FALSE, size = 8) +
 scale_x_log10(labels=scales::dollar)
```



程序中 `size` 指定了线的粗细倍数, `se = FALSE` 关闭了置信区间显示。用 `alpha =` 设置了透明度, 取 0 和 1 之间的值, 数值越小越透明。在有多个点时适当设置透明度可以比较好地显示出重叠的点, 重叠点越多点的颜色越深。虽然这里设置了固定的透明度, 也可以在 `aes()` 中将透明度 `alpha` 映射到某个变量, 使得该变量值大小用点的透明度表示。

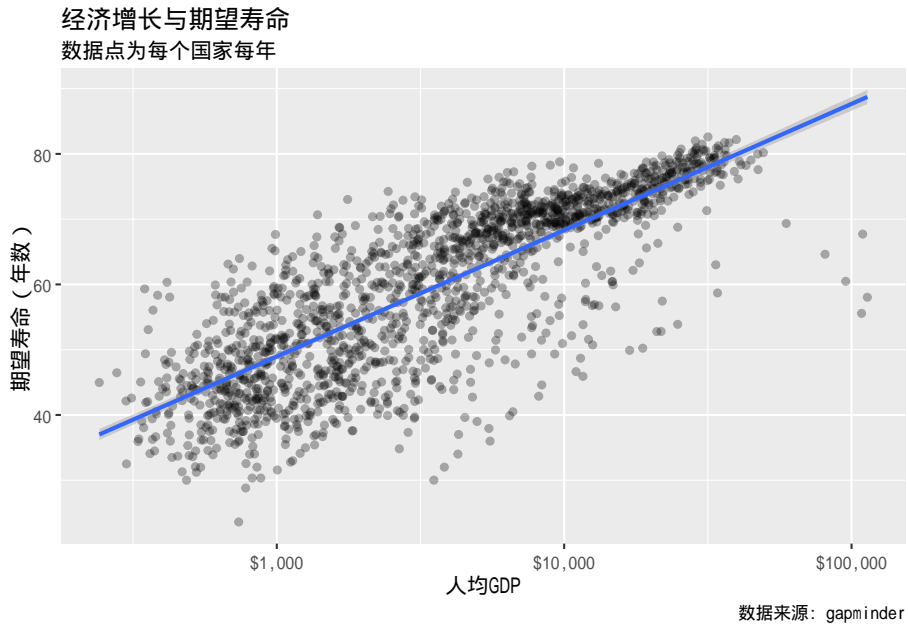
下面用 `labs()` 函数给图形加上适当的标题:

```
p <- ggplot(data=gapminder,
 mapping = aes(
 x = gdpPercap,
 y = lifeExp))
p + geom_point(alpha = 0.3) +
 geom_smooth(method="gam") +
 scale_x_log10(labels=scales::dollar) +
 labs(
 x = " 人均 GDP",
 y = " 期望寿命 (年数)",
 title = " 经济增长与期望寿命",
```

```

subtitle = " 数据点为每个国家每年",
caption = " 数据来源: gapminder")

```



可以看出, `labs()` 规定了上方的标题、小标题, x 轴、y 轴的标题, 右下方的标注 (`caption`)。坐标轴刻度数值的规定则需要在 `scale_xxx()` 函数中给出。

### 30.3.4 在 `geom` 函数中映射变量

在前面的一个例图中, 在 `ggplot()` 函数中将 `color` 和 `fill` 映射到了 `continent` 变量, 使得不仅散点颜色代表了不同大洲, 还使得每个大洲单独拟合了曲线。如果希望所有大洲拟合同一条曲线怎么办?

在必要时, 可以在 `geom_xxx()` 函数中用 `mapping = aes(<...>)` 单独指定变量映射。例如, 下面的程序在 `geom_point()` 中将不同大洲映射为不同颜色, 而不影响 `geom_smooth()` 中的颜色以及分组:

```

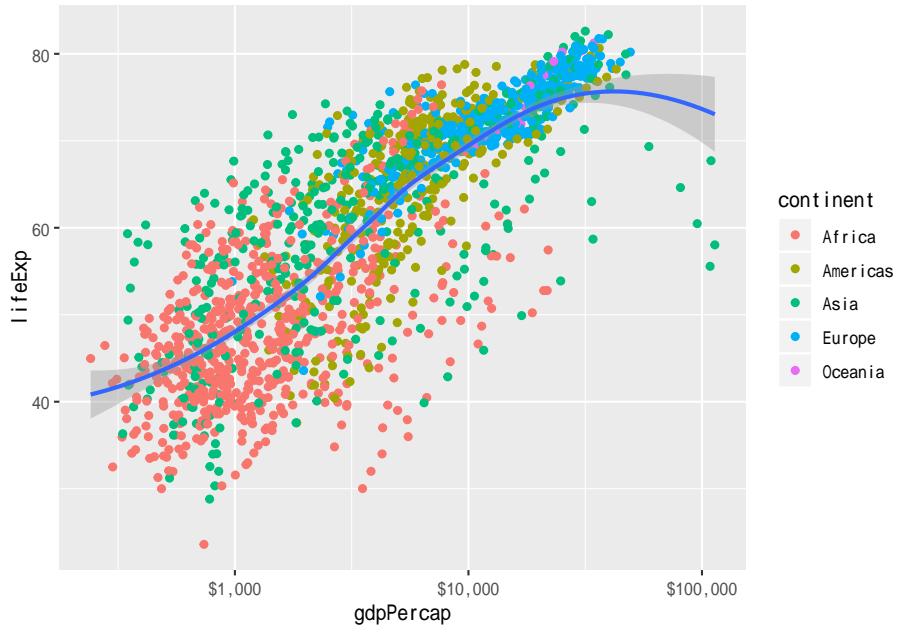
p <- ggplot(data=gapminder,
 mapping = aes(
 x = gdpPercap,

```

```

y = lifeExp))
p + geom_point(mapping = aes(color = continent)) +
 geom_smooth(method="loess") +
 scale_x_log10(labels=scales::dollar)

```



### 30.3.5 连续变量的颜色映射

也可以将连续变量映射为渐变色。除了表示二元函数的等值线图以外这种方法并不利于读者认读。

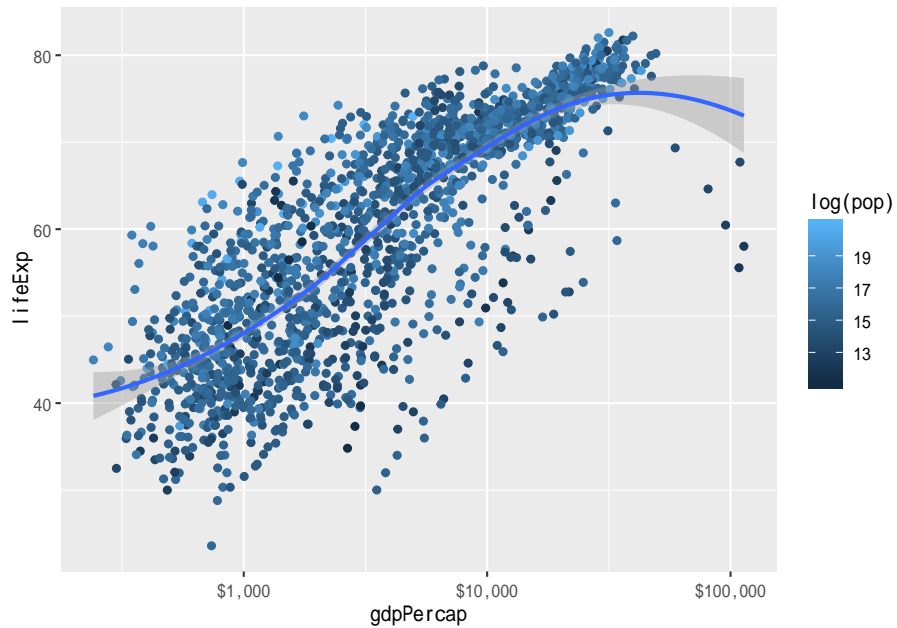
例如，将人口数取自然对数映射为渐变色：

```

p <- ggplot(data=gapminder,
 mapping = aes(
 x = gdpPercap,
 y = lifeExp,
 color = log(pop)))
p + geom_point() +
 geom_smooth(method="loess") +

```

```
scale_x_log10(labels=scales::dollar)
```



这里不同散点的颜色是连续变化的，右侧的图例仅显示了有限的一些代表值。

### 30.3.6 保存图像

如果使用 Rmarkdown 制作图文，图像会自动进入编译的结果（如 PDF、Word、HTML）中，图像大小、输出大小可以用 Rmarkdown 的设置调整。

为了将最近生成的图形保存为 PNG 格式，用命令如

```
ggsave(filename=" 文件名.png")
```

保存为 PDF 格式：

```
ggsave(filename=" 文件名.pdf")
```

如果将制作的图形保存到了一个 R 变量中，在 `ggsave()` 中可以用 `plot=` 参数指定，如

```
ggout01 <- p + geom_point()
ggsave(filename=" 文件名.pdf", plot=ggout01)
```

在 `ggsave()` 中可以用 `scale =` 指定放大比例，用 `height =` 指定高度，用 `width =` 指定宽度，用 `units =` 指定高度和宽度的单位，如：

```
ggsave(filename=" 文件名.pdf", plot=ggout01,
 height=12, width=8, units="cm")
```

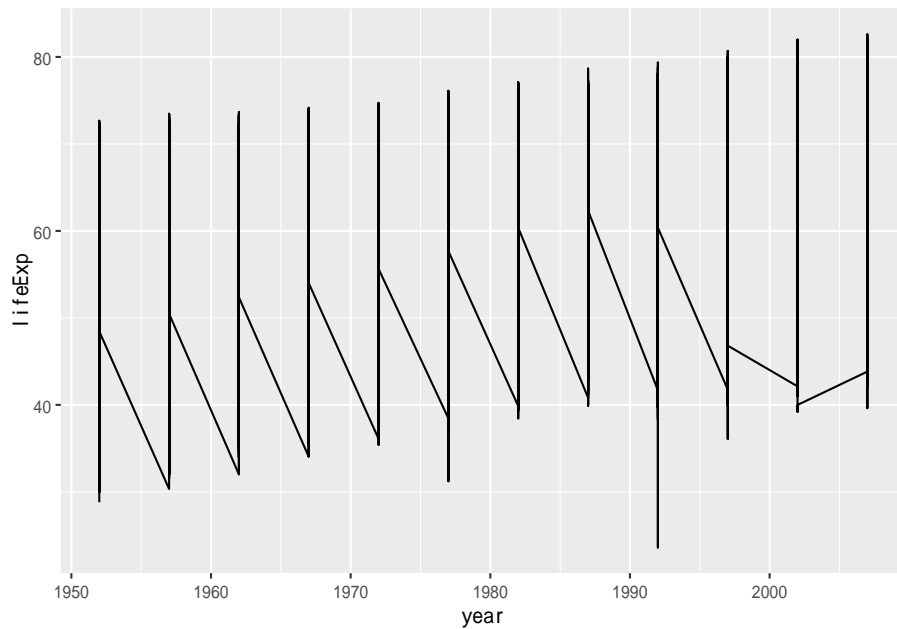
单位可以是 in, cm, mm。

## 30.4 分组、小图、变换、条形图

### 30.4.1 图形中的分组

考虑 `gapminder` 数据集中每个国家的期望寿命随时间（年）的变化。用 `geom_line()` 可以画折线图。因为有许多国家，所以无法直接作图，如：

```
p <- ggplot(data = gapminder,
 mapping = aes(
 x = year,
 y = lifeExp))
p + geom_line()
```



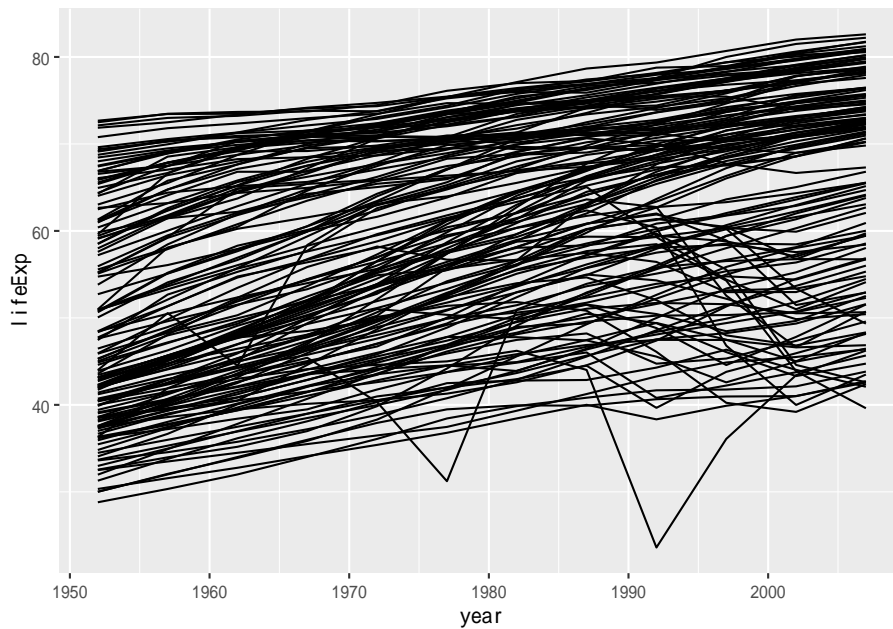
没有得出我们希望的每个国家一条曲线的效果。这是因为程序中没有指定需要按照国家分组，使得同一年的不同国家的坐标连成了一条竖线。

要注意的是，`geom_line()` 会自动将 `x` 坐标从小到大排序，然后再连接相邻的点。如果希望按输入数据的次序连接相邻的点，需要用 `geom_path()` 函数。

为了解决上图的问题，加入按照国家分组的设定。实际上，分组 (`group`) 与 `x`、`y`、`color`、`fill` 一样可以映射到一个变量，但仅能映射到分类变量。上述程序的改进如下：

```
p <- ggplot(data = gapminder,
 mapping = aes(
 x = year,
 y = lifeExp,
 group = country))
p + geom_line()
```





结果图形中每一条曲线对应一个国家。为了查探其中最下方的不稳定曲线是哪一个国家，使用筛选观测的功能：

```
gapminder %>%
 filter(lifeExp < 30, year >= 1990)

A tibble: 1 x 6
country continent year lifeExp pop gdpPercap
<fct> <fct> <int> <dbl> <int> <dbl>
1 Rwanda Africa 1992 23.6 7290203 737.
```

该国家为 Rwanda。

分组 (group) 映射仅在其它映射没有将分组考虑在内的情况。比如，分大洲用不同颜色，用了 `color=continent`，则不需要 `group=continent`。

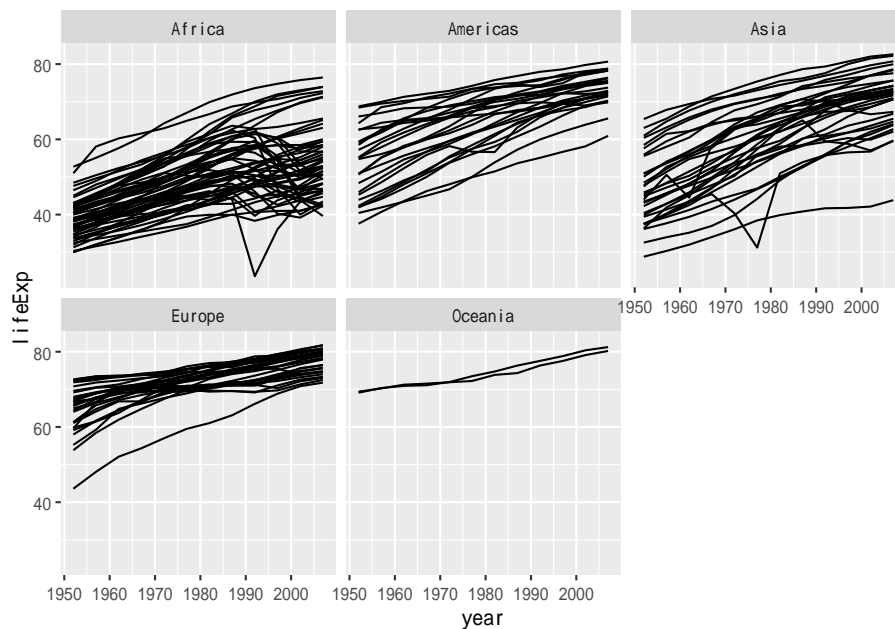
### 30.4.2 小图 (facet)

上面的图包含了过多的曲线，使得图形表现得很拥挤。可以将一个作图区域拆分成若干个小块，称为小图 (facet)，按照某一个或两个分类变量的不同值将数

据分为若干个子集，每个数据子集分别在小图上作图。

对于上面的例子，可以将每个大洲的图形分别放置在一个小图上。小图不是一种变量映射，而是一种图形摆放方法，所以不设置在 `aes()` 函数内，而是用 `facet_wrap()` 函数规定。这种功能与 `group` 映射的功能有些重复，所以有时需要与 `group` 映射配合使用，有时则不需要。程序如：

```
p <- ggplot(data = gapminder,
 mapping = aes(
 x = year,
 y = lifeExp,
 group = country))
p + geom_line() + facet_wrap(~ continent)
```



区分不同小图的标签写在每个小图的上方。可以用 `facet_wrap()` 参数 `strip_position` 和参数 `switch` 调整标签的上下左右。

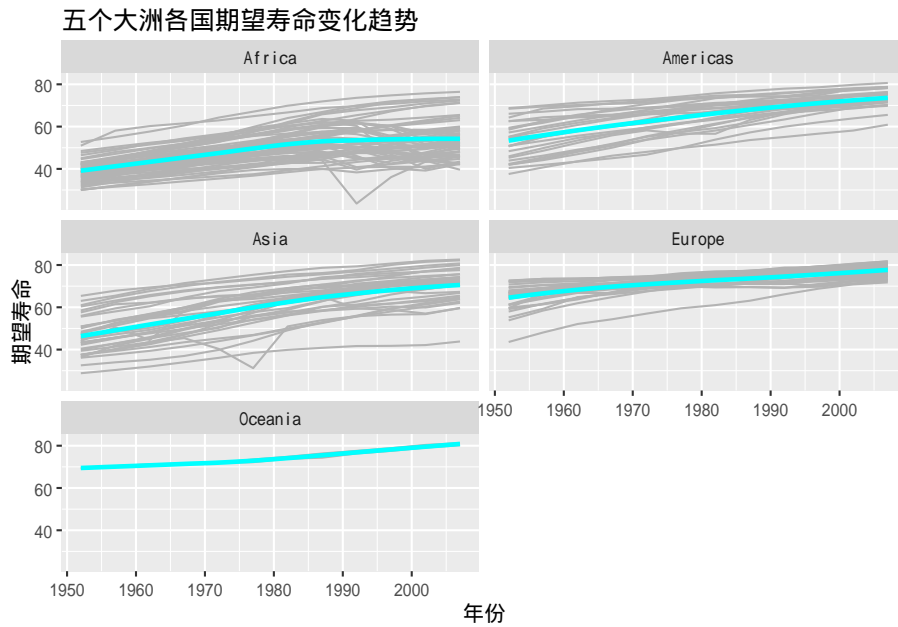
小图之间默认公用了横坐标和纵坐标且坐标范围保持一致。如果不保持一致，读者可能会有误解。但是 `x` 轴或 `y` 轴映射为分类变量且不同小图的分类型完全不同时，可以令各小图中该轴的取值不统一。`facet_wrap()` 选项 `scales` 默认为 "fixed"，即所有小图的 `x` 轴、`y` 轴都范围一致，取 "free\_x" 则允许各小图

的 x 轴不统一, "free\_y" 允许各小图的 y 轴不统一, "free" 允许各小图的 x 轴和 y 轴都不统一。

在 `facet_wrap()` 中可以用 `ncol` 参数指定小图的列数, 用 `nrow` 指定小图的行数。各个小图的次序应该设定为一定的合理次序, 比如用来分类的变量本身有序, 或者令各小图中的数据值有一定的增减次序。

下面的程序将曲线颜色变浅, 对每个大洲增加了拟合曲线, 增加了适当的标题和坐标轴标签:

```
p <- ggplot(data = gapminder,
 mapping = aes(
 x = year,
 y = lifeExp))
p + geom_line(mapping = aes(group = country), color = "gray70") +
 geom_smooth(method = "loess", color="cyan", se = FALSE, size = 1.1) +
 facet_wrap(~ continent, ncol = 2) +
 labs(
 x = " 年份",
 y = " 期望寿命",
 title = " 五个大洲各国期望寿命变化趋势"
)
```



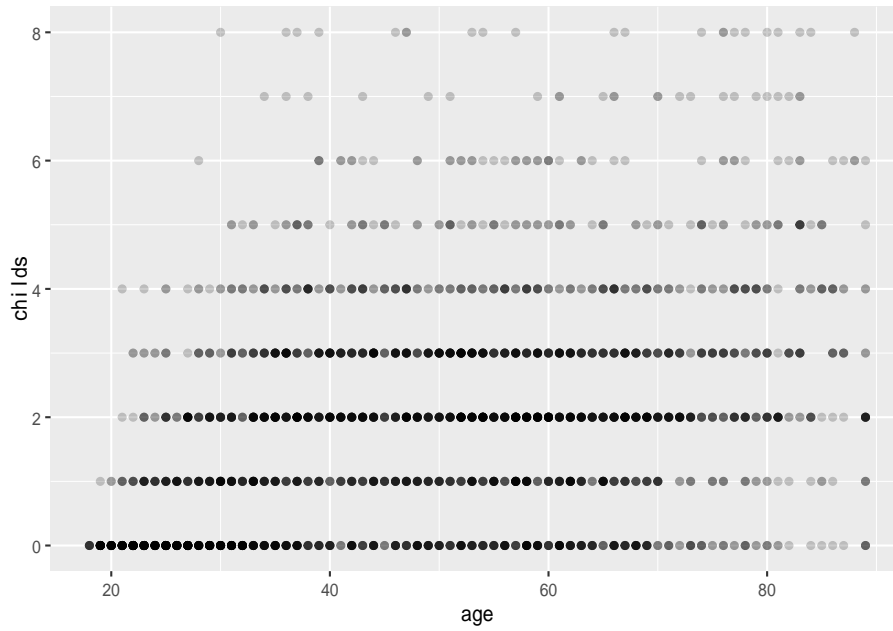
注意 `group = country` 的设置从 `ggplot()` 函数中转移到了 `geom_line()` 函数中，否则就意味着拟合线也需要按照国家分组，而不是按大洲分组。

`facet_wrap()` 主要适用于按照一个分类变量的值将不同观测在不同小图中表现，可以人为指定小图的行数和列数。如果需要按照两个分类变量交叉分组分配小图，可以用 `facet_grid()` 函数。

例如，对 `gss_sm` 数据集，作小孩个数对年龄的散点图：

```
p <- ggplot(data=gss_sm,
 mapping = aes(
 x = age,
 y = childs))
p + geom_point(alpha = 0.2)
```

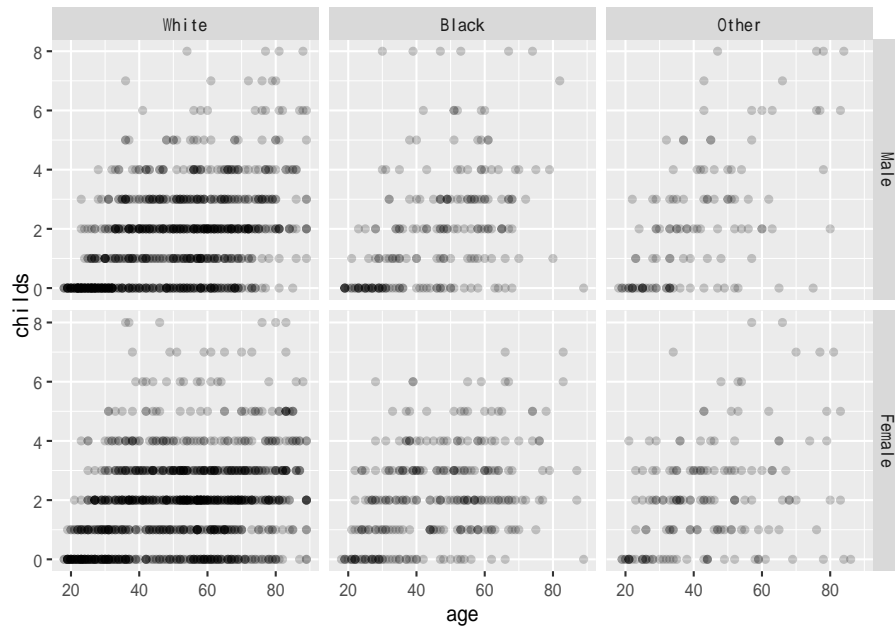
```
Warning: Removed 18 rows containing missing values (geom_point).
```



有过多的重叠点。将观测按照性别 (sex) 和种族 (race) 交叉分组，分配到不同的小图上：

```
p + geom_point(alpha = 0.2) +
 facet_grid(sex ~ race)
```

```
Warning: Removed 18 rows containing missing values (geom_point).
```



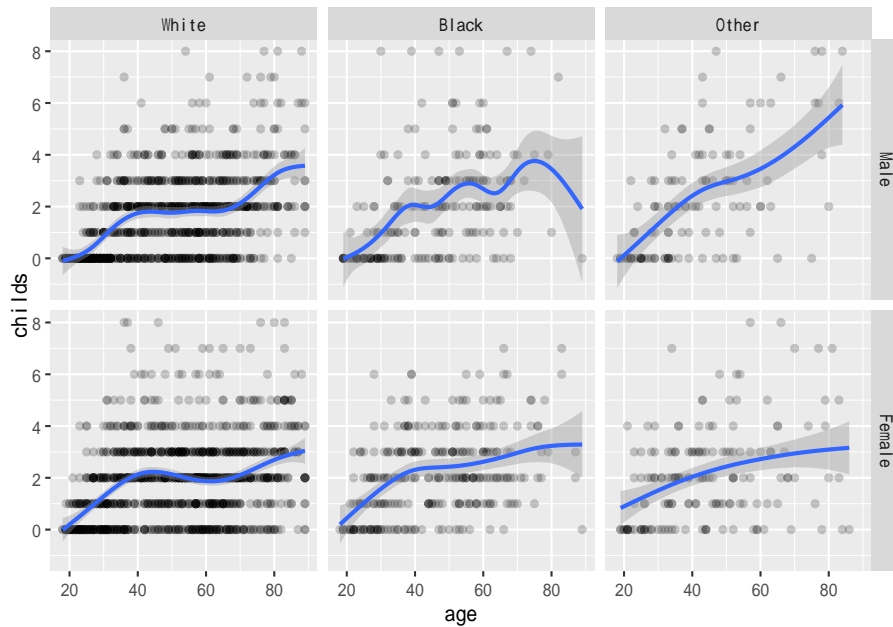
交叉分组时作小图时, `sex ~ race` 这种写法使得不同性别对应到不同行, 不同种族对应到不同列。在图形中增加拟合曲线:

```
p + geom_point(alpha = 0.2) +
 geom_smooth() +
 facet_grid(sex ~ race)
```

```
`geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```

```
Warning: Removed 18 rows containing non-finite values (stat_smooth).
```

```
Warning: Removed 18 rows containing missing values (geom_point).
```



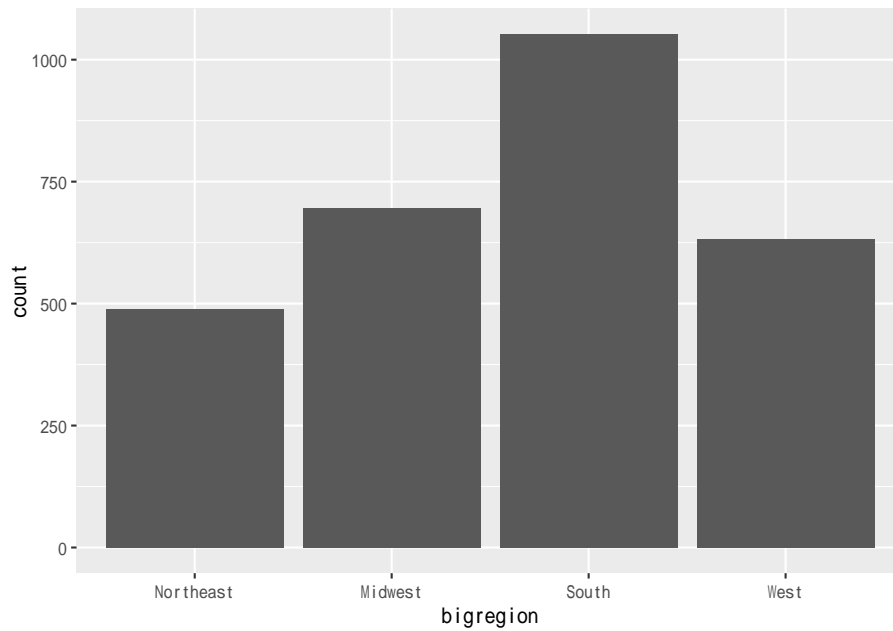
这里虽然没有映射 `group` 维，但还是按性别和种族对数据集分成了 6 个子集，每个小图中仅有一个自己的数据。

### 30.4.3 数据变换与条形图

有些 `geom_xxx()` 函数直接按照数据值作图，如 `geom_point()`、`geom_line()`，而 `geom_smooth()` 这样的函数则会按照某种算法计算并对计算结果作图。`geom_xxx()` 都有默认的 `stat_xxx()` 函数用来计算，也可以人为指定不同的统计规则。

考虑条形图的例子。ggplot2 中的条形图函数 `geom_bar()` 可以对一个分类变量自动统计频数，并作频数条形图。比如对 `gss_sm` 数据集的 `bigregion` 变量作频数条形图：

```
p <- ggplot(data = gss_sm,
 mapping = aes(x = bigregion))
p + geom_bar()
```

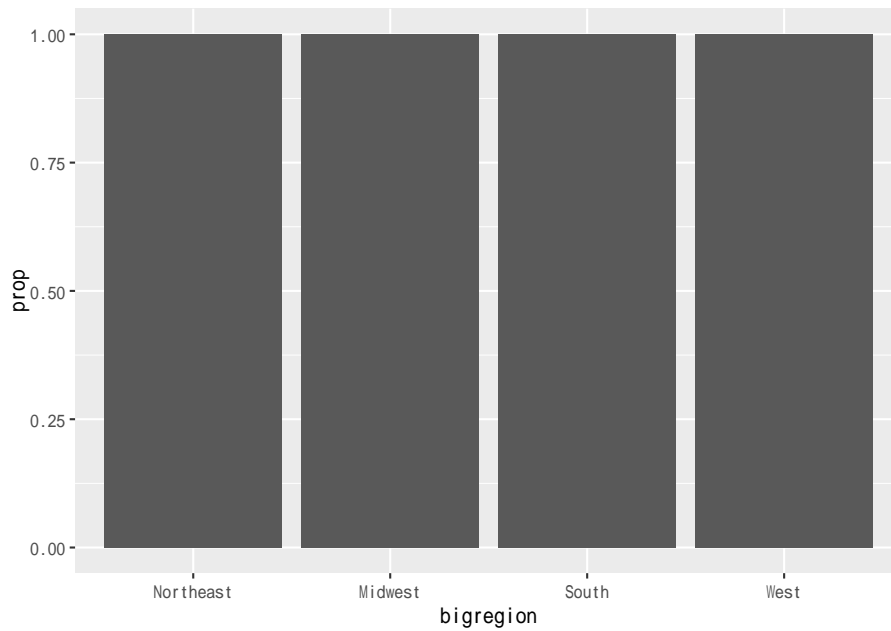


结果是每个大区的受访者人数的条形图。图形中 `x` 映射是用户指定的，而 `y` 轴则是自动计算的频数。实际上，`geom_bar()` 自动调用了统计函数 `stat_freq()` 对每个大区计算频数，生成新变量 `count` 和 `prop`。`geom_bar()` 默认使用 `count`(频数)。

下面的程序将纵坐标改成了比例：

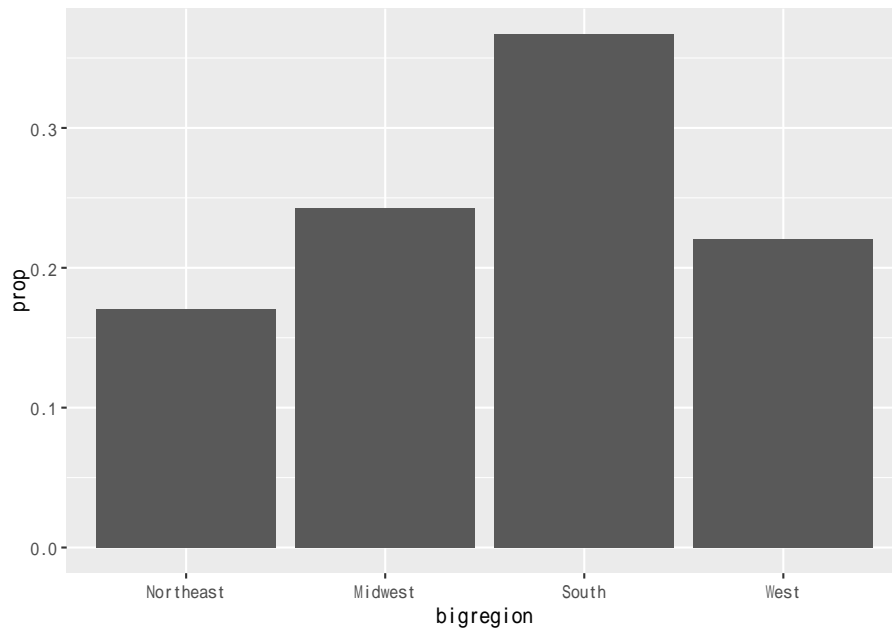
```
p + geom_bar(mapping = aes(y = ..prop..))
```





程序中的统计结果比例用特殊变量名 `..prop..` 表示。这个图形是错误的，原因是在 `ggplot()` 函数中指定 `x = bigregion` 后，因为条形图的主要变量变量 `bigregion` 是分组变量，这样在 `geom_bar()` 的执行中仍按照 `bigregion` 分组。解决方法是在 `geom_bar()` 的 `aes()` 中加入 `group = 1` 选项，这实际是定义了一个常数值的分组变量，也就仅有一组了：

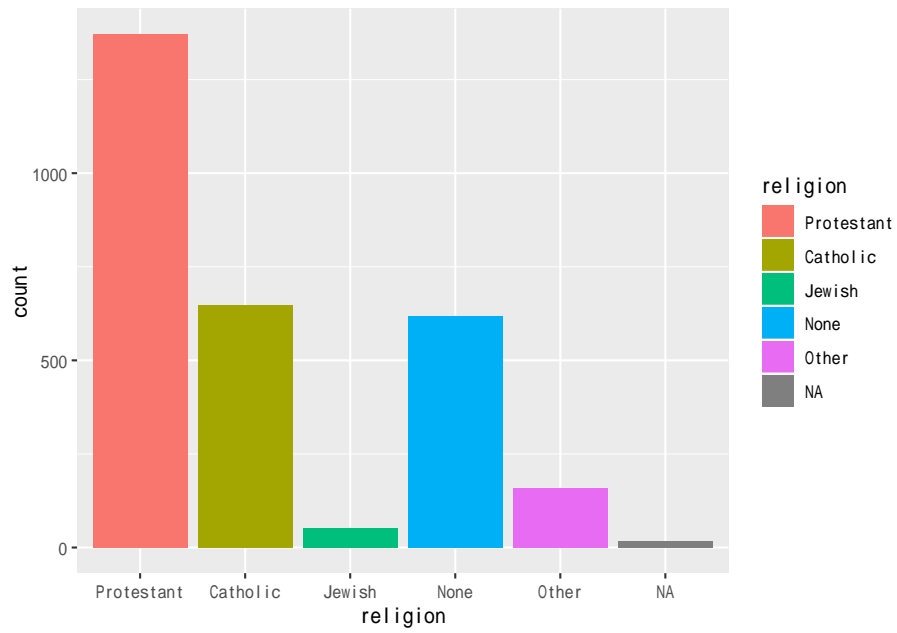
```
p + geom_bar(mapping = aes(y = ..prop.., group = 1))
```



虽然解决了问题，但是用 `geom_bar()` 函数直接从原始数据做比例的条形图还是很容易做错，所以建议从原始数据中用数据概括方法算出比例再用 `geom_col()` 函数作条形图，见30.4.6。

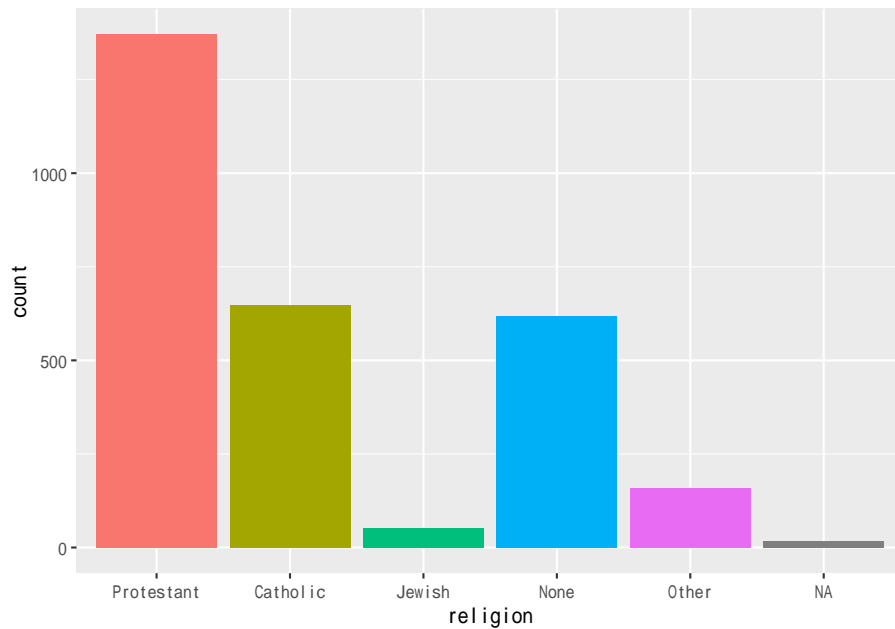
下面的例子作 `gss_sm` 数据集中 `religion` 变量的频数条形图，并给不同的条形自动分配不同的颜色，方法是指定 `fill = religion`:

```
p <- ggplot(data = gss_sm,
 mapping = aes(x = religion, fill = religion))
p + geom_bar()
```



因为将 `religion` 同时映射到 `x` 维与 `fill` 维，所以对应 `fill` 维在图形右侧出现了图例，这是多余的。调用 `guides(fill = FALSE)` 可以人为指定不做关于填充色的图例：

```
p + geom_bar() +
 guides(fill = FALSE)
```



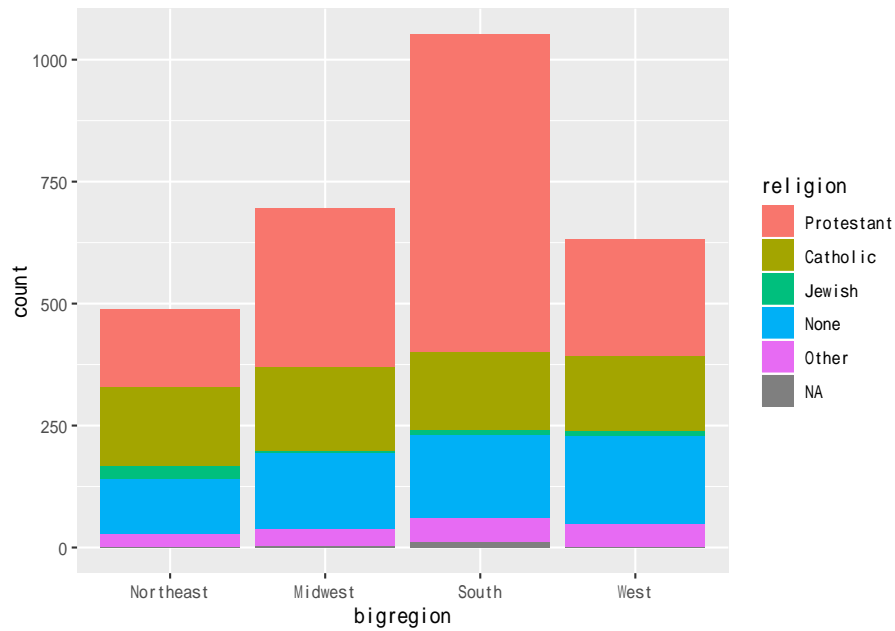
从可视化理论的角度看，上图中的不同颜色是多余的，用同一颜色更能强度数据本身。

#### 30.4.4 分段与并列条形图

上面的条形图展现了单个分类变量的频数分布。两个分类变量的交叉频数分布可以用分段条形图或者并列条形图表现。

例如，对 `gss_sm` 数据集，按照 `bigregion` 分组计算频数，每组内再按照 `religion` 计算频数，作图如下：

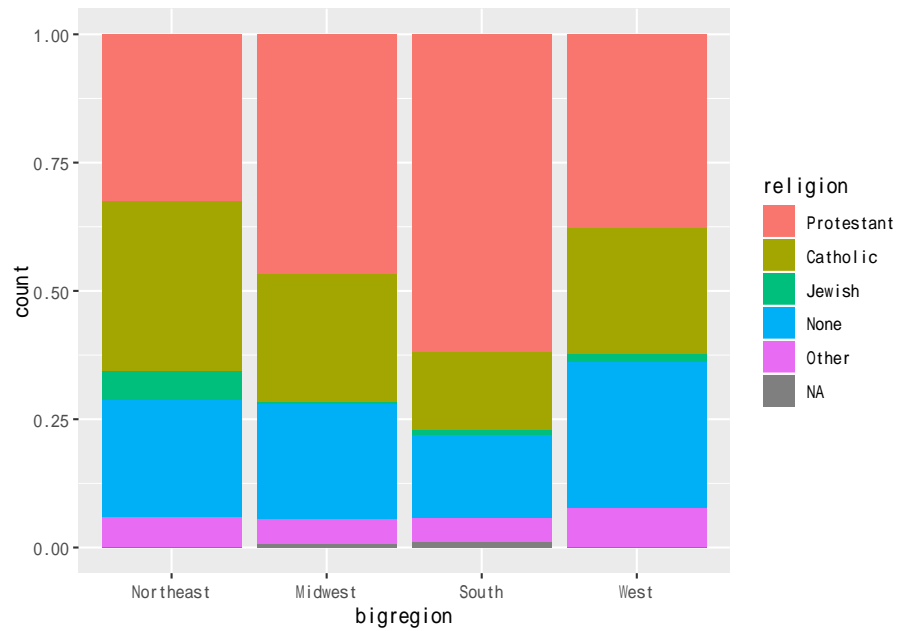
```
p <- ggplot(data = gss_sm,
 mapping = aes(
 x = bigregion,
 fill = religion))
p + geom_bar()
```



这样的图形可以很容易地比较大类 (这里是 bigregion) 的频数比例, 但大类内的小类 (这里是 religion) 可以比较容易地在同类内部比较, 但是在同类之间比较则较困难。如果仅有两个小类, 则小类在同类之间的比较也没有问题。

另一种做法是将同类的高度拉平, 图形仅表示每一同类内部小类的比例, 没有同类频数信息, 也不能比较两个同类之间的小类频数, 但可以大致地在同类之间比较小类的比例:

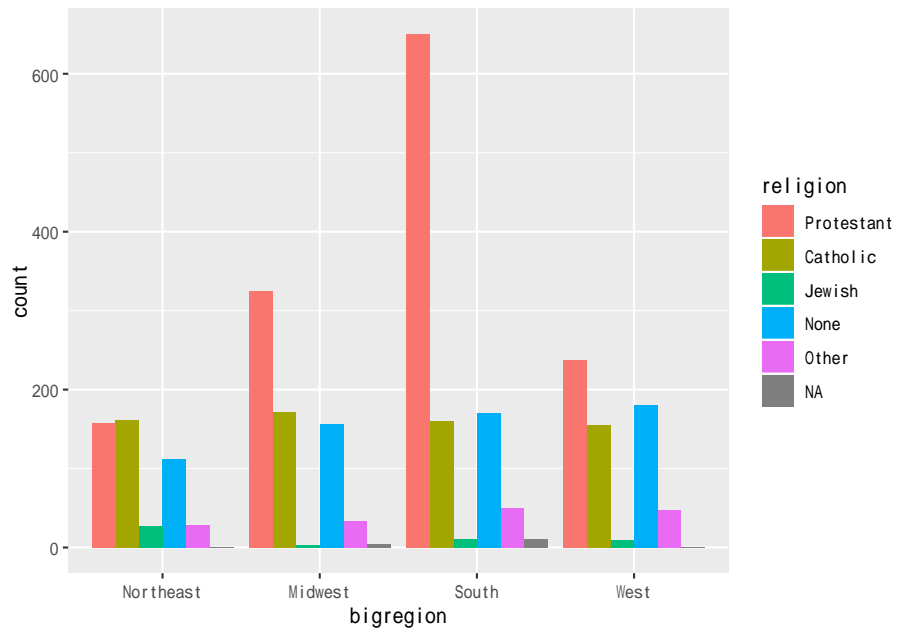
```
p + geom_bar(position = "fill")
```



上面的程序在 `geom_bar()` 中用了 `position = "fill"` 选项。

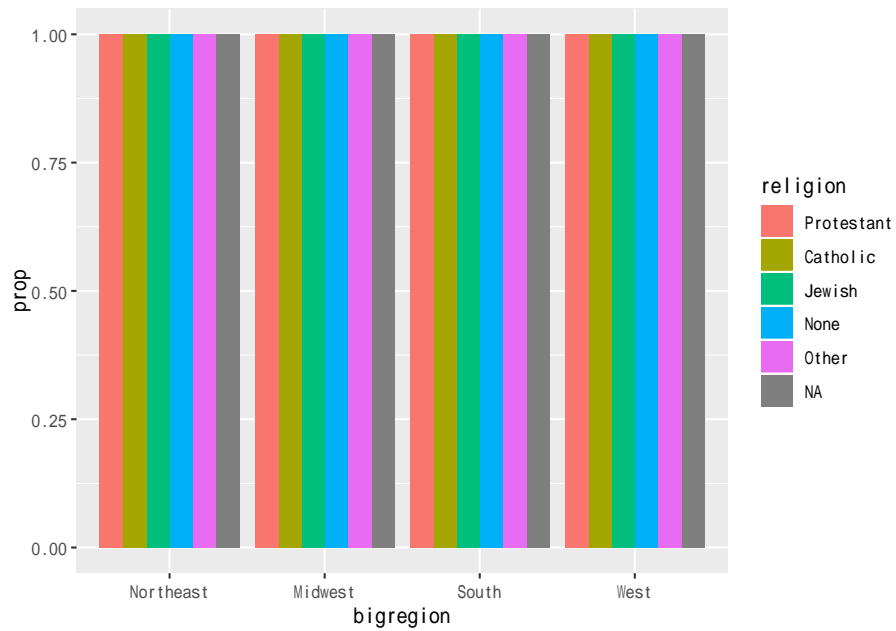
并排的条形图可以表现每个交叉类的频数，可以比较容易地比较每个大类内部的小类比例以及小类的频数，但是不容易比较大类的比例：

```
p + geom_bar(position = "dodge")
```



如果想将上图中的纵轴改为大类内的比例，下面的第一次尝试给出错误的结果：

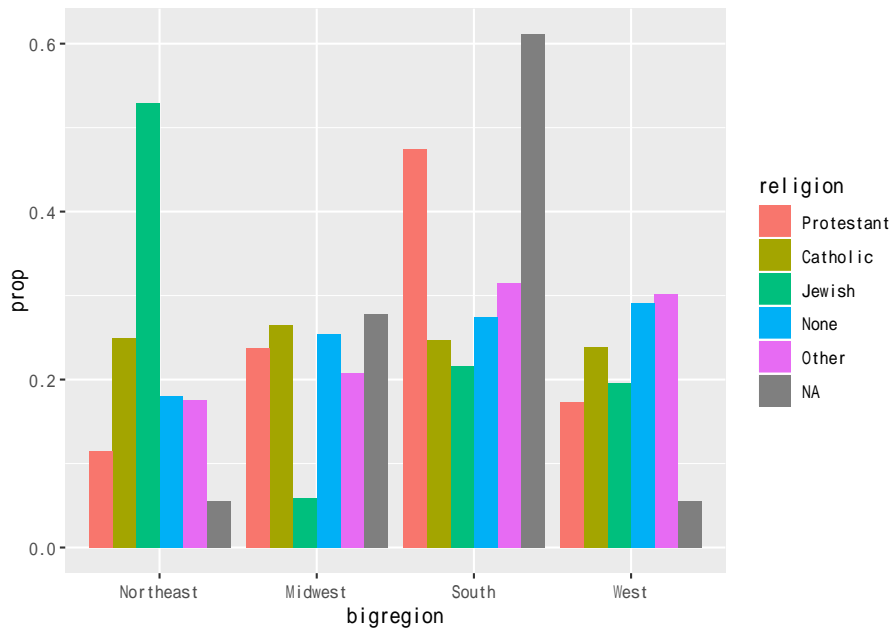
```
p + geom_bar(position = "dodge",
 mapping = aes(y = ..prop..))
```



出错的原因是，`ggplot()` 中的 `x =` 和 `fill =` 指定了两个分组变量，结果是交叉地分成了  $4 \times 6 = 24$  个组，每一组内的比例当然都是 100%。第二次尝试指定分组变量仅为 `religion`:

```
p + geom_bar(position = "dodge",
 mapping = aes(y = ..prop..., group=religion))
```



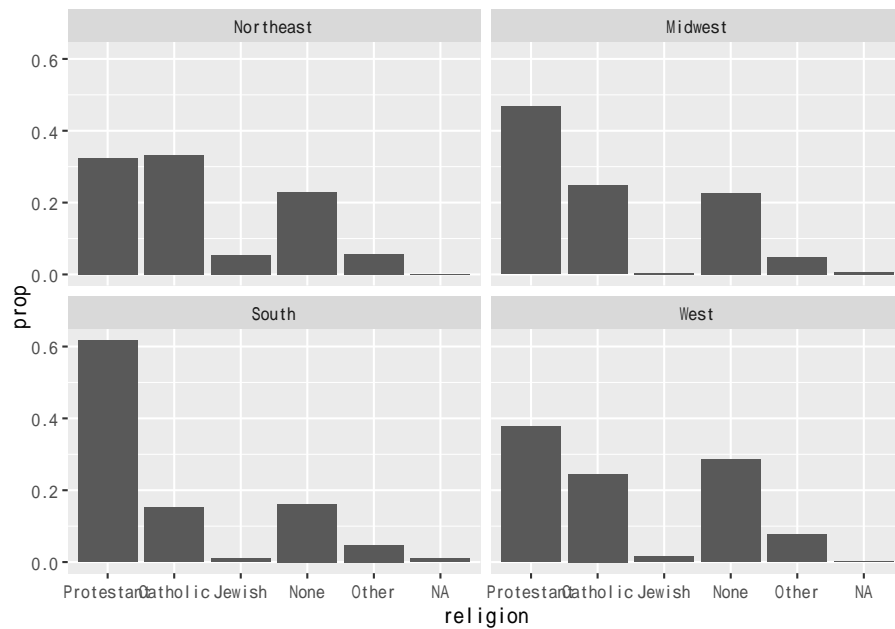


这个图形表面上没错，但是仔细查看发现，上图中每个大区内部的宗教比例之和不等于 1，而是每种宗教在不同大区之间的比例之和等于 1。

实际上，对于列联表，最好是预先统计出列联表结果，以列联表为输入用 `geom_col()` 函数作图，就可以避免许多错误图形。

为了在不同大区之间比较宗教比例分布，可以借助于小图，将每个大区分配到一个小图：

```
p <- ggplot(data = gss_sm,
 mapping = aes(x = religion, group=bigregion))
p + geom_bar(position = "dodge",
 mapping = aes(y = ..prop..)) +
 facet_wrap(~ bigregion, ncol=2)
```



注意上面的程序将 `x` 映射到了 `religion`，将 `group` 映射到了 `bigregion`。如果不映射 `group` 则图形还是错误的。直接对原始数据做并列条形图比较类内比例不够直观，更好的办法是先统计出交叉频数表，以交叉频数表作为图形函数的输入作图，见30.4.6。

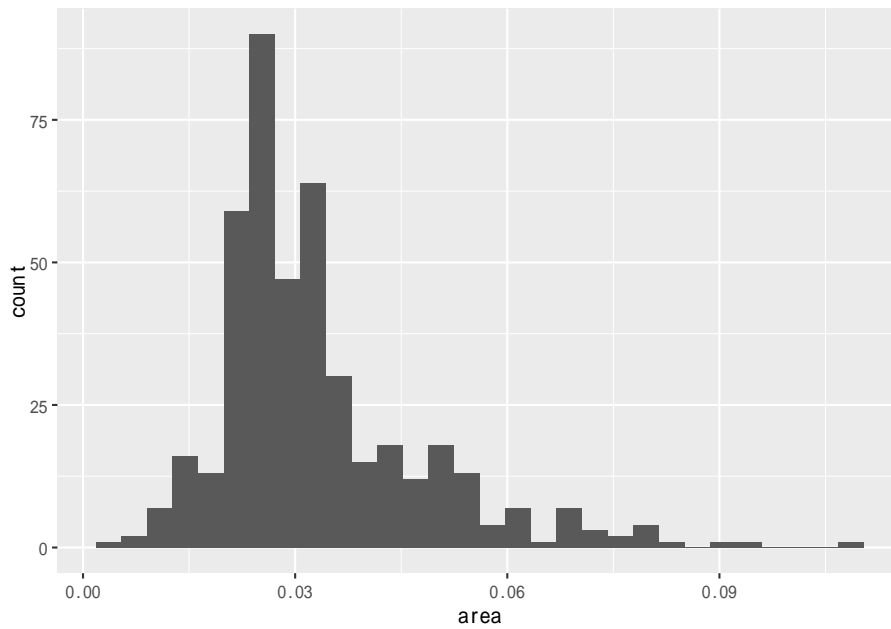
### 30.4.5 直方图与密度估计

条形图 (barplot) 反映分类变量的频数分布或者比例，直方图 (histogram) 反映连续取值的数值变量的分布。`geom_histogram()` 作直方图，可以自动选取合适的分组个数，也可以人为指定分组个数。

`ggplot2` 包中的 `midwest` 数据集包含了美国中西部的一些县的统计数据，如面积（单位：平方英里）。下面的程序对连续取值的数值型变量 `area` 作频数直方图，自动确定分组个数：

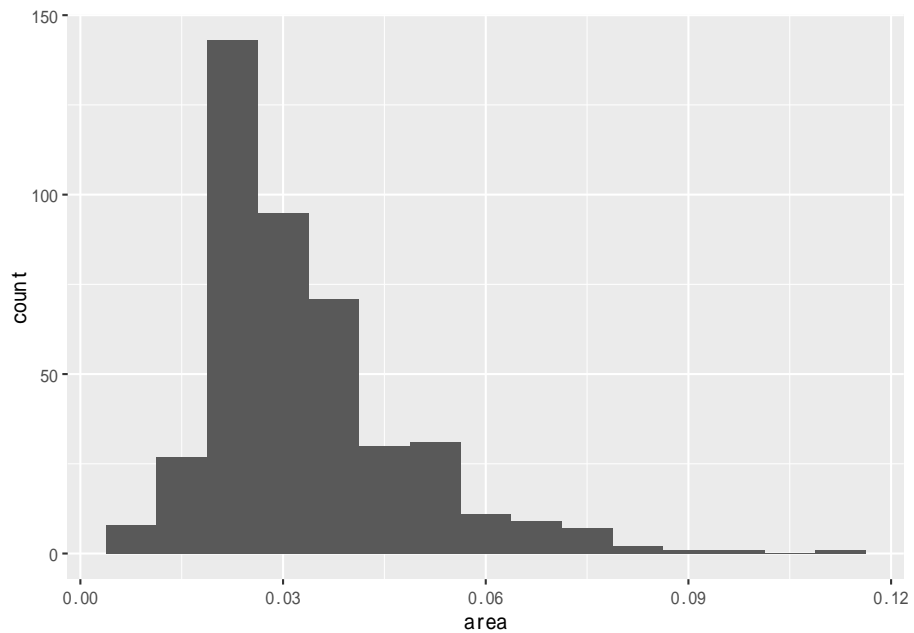
```
p <- ggplot(data = midwest,
 mapping = aes(x = area))
p + geom_histogram()

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



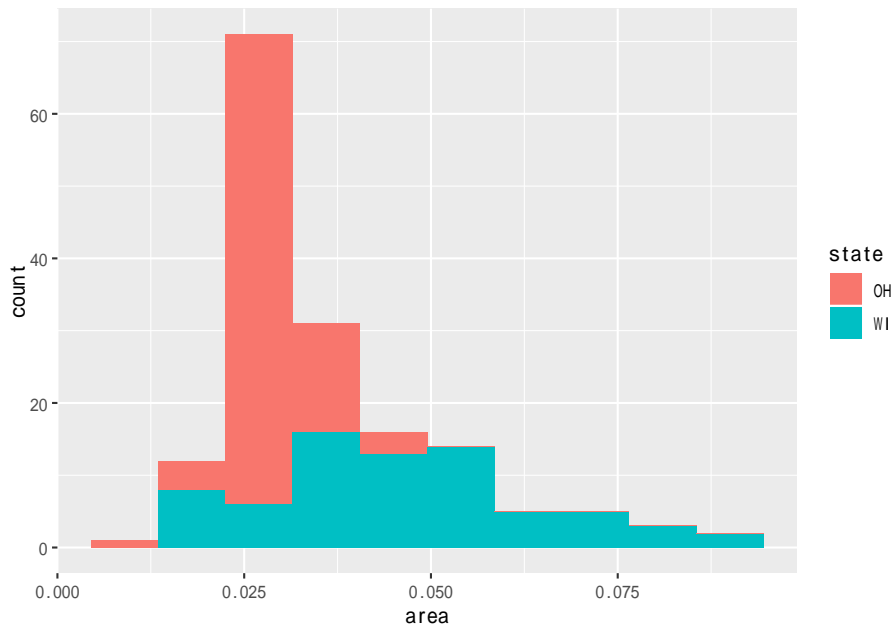
上面图形的纵坐标是频数 (count), 是每个组的频数。`geom_histogram()` 默认调用 `stat_bin()` 进行分组及频数统计。直方图的形状比较依赖于分组数与分组起始点位置, 可以用 `bins` 参数控制分组数, 用 `binwidth` 参数控制分组宽度, 用 `center` 或者 `boundary` 参数控制组中心或者组边界对齐位置, 如:

```
p + geom_histogram(bins = 15)
```



可以利用 `fill` 映射将构成直方图的观测按照某个分类变量分组，然后每个条形内部按照该分类变量的值分段染色，段内各颜色的长度代表该条形所在组某一类的频数，如：

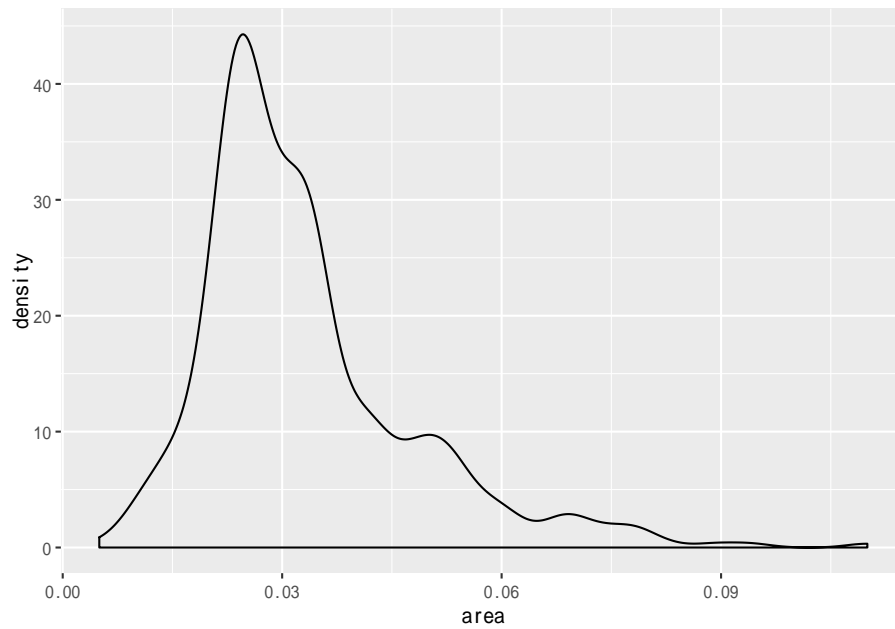
```
midwest_sub <- midwest %>%
 filter(state %in% c("OH", "WI"))
p <- ggplot(data=midwest_sub,
 mapping = aes(x = area, fill = state))
p + geom_histogram(bins = 10)
```



可见面积较小的县主要来自 OH 州，面积较大的县主要来自 WI 州。

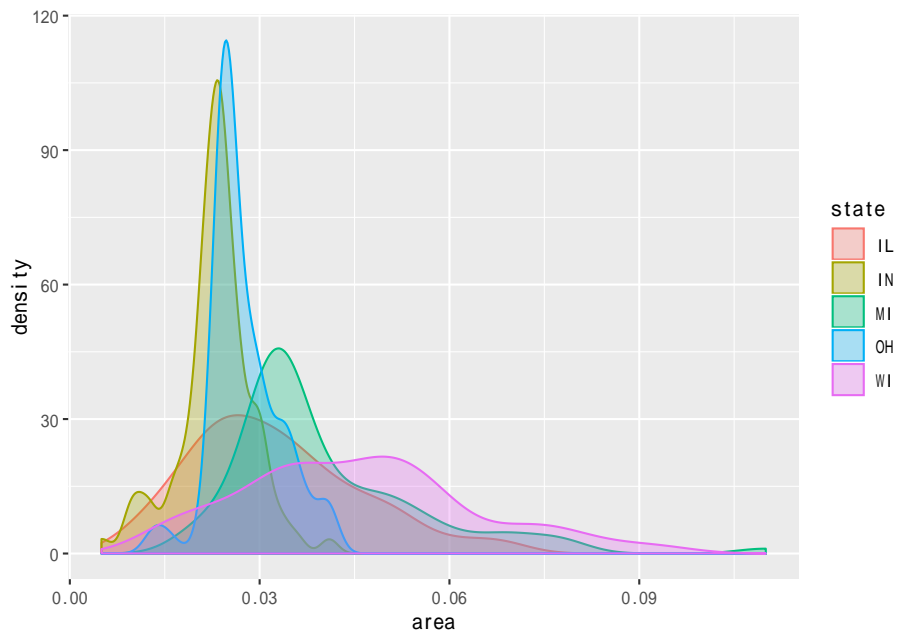
`geom_density()` 可以对连续变量绘制密度估计曲线，如：

```
p <- ggplot(data = midwest,
 mapping = aes(x = area))
p + geom_density()
```



下面的程序写法制作每个州的各县的面积密度估计，画在同一坐标系中：

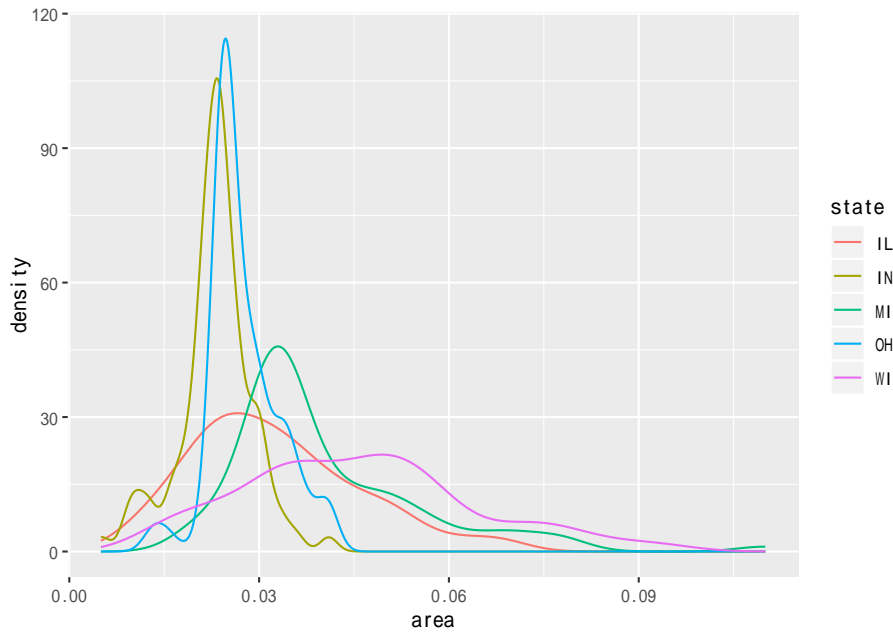
```
p <- ggplot(data = midwest,
 mapping = aes(
 x = area,
 color = state,
 fill = state))
p + geom_density(alpha = 0.3)
```



可以看出，IN 与 MI 州各县的面积偏小。WI 州各县的面积较大。

上面的图形可以借助于 `geom_line(stat = "density")` 改成仅有多条曲线：

```
p <- ggplot(data = midwest,
 mapping = aes(
 x = area,
 color = state))
p + geom_line(stat = "density")
```

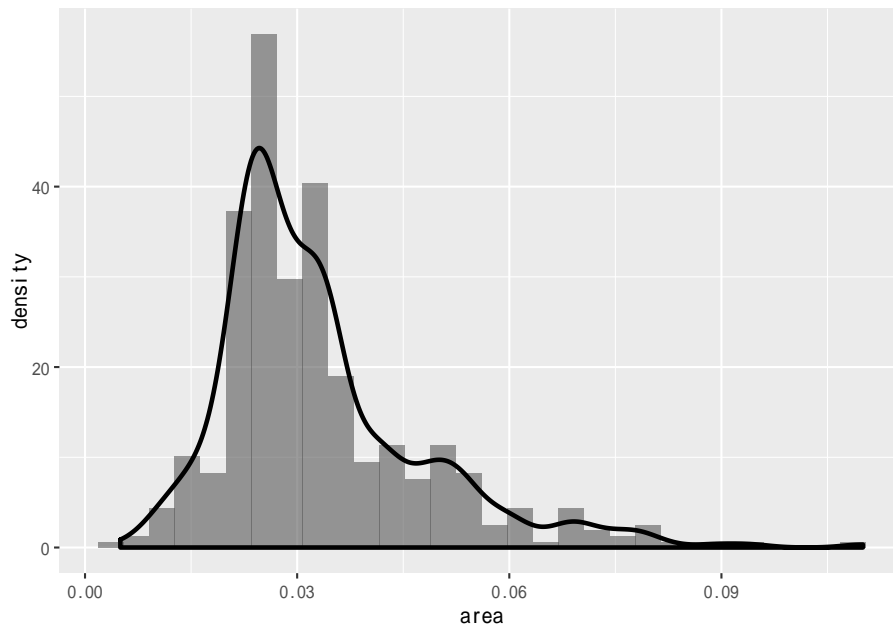


`geom_density()` 的纵轴是密度估计。为了能够将直方图与密度估计画在同一坐标系中，需要将直方图的纵轴也改为密度估计，如：

```
p <- ggplot(data = midwest,
 mapping = aes(x = area))
p + geom_histogram(mapping = aes(y = ..density..), alpha = 0.6) +
 geom_density(size = 1.1)
```

```
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```





进一步地, `geom_freqpoly()` 将直方图做成折线格式; `geom_bin2d()` 作二维的直方图, 用不同颜色代表密度; `geom_density_2d()` 作二维密度估计等值线图。

### 30.4.6 频数表的条形图

有时用来绘图的数据已经是一个频数表, 比如泰坦尼克号乘客生存与性别的频数表:

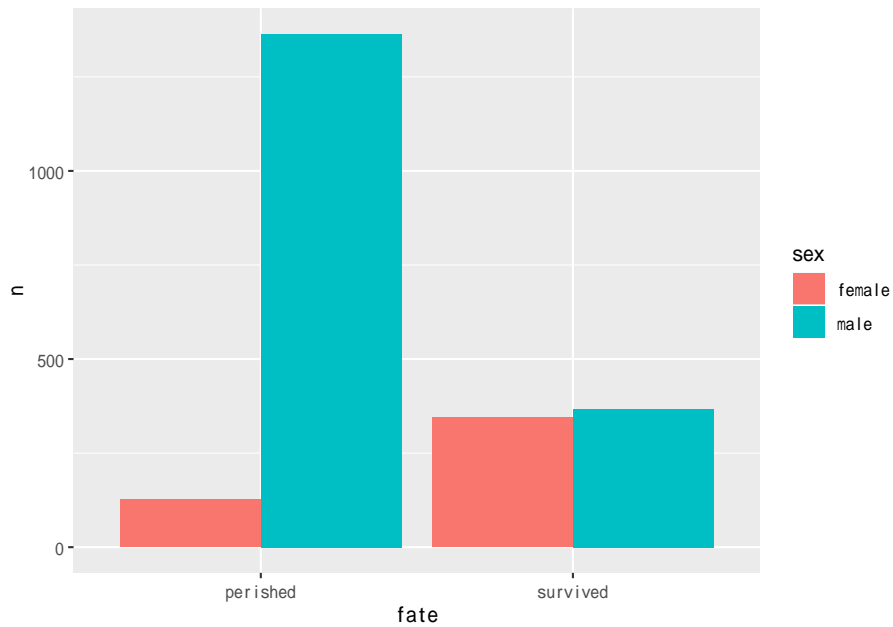
```
titanic
```

```
fate sex n percent
1 perished male 1364 62.0
2 perished female 126 5.7
3 survived male 367 16.7
4 survived female 344 15.6
```

这是一个长表格的列联表, 对于 `table()` 生成的列联表可以用 `as.data.frame` 将其转换为长表格式。

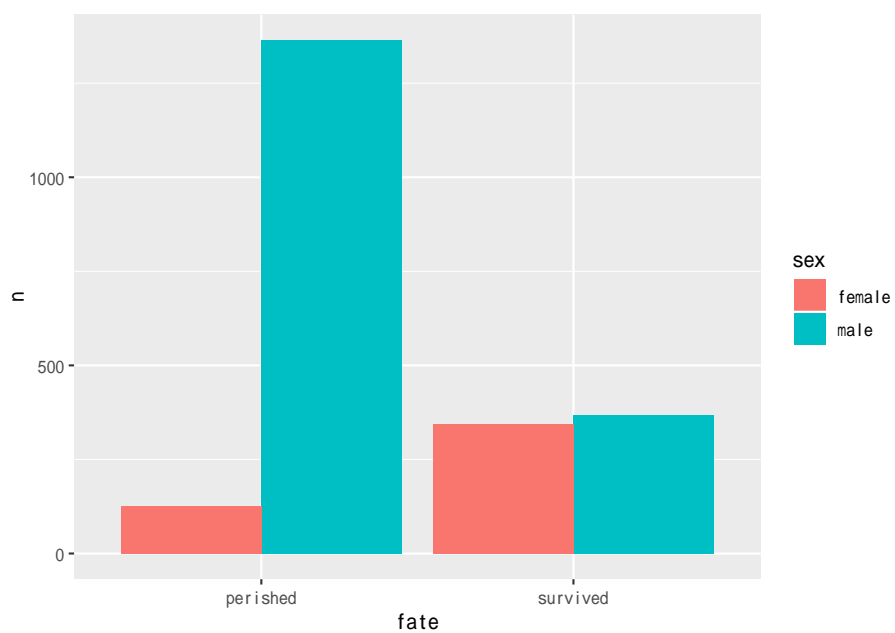
在 `geom_bar()` 中指定 `stat = "identity"` 就表示频数(或者比例)已有,用 `y` 维进行映射,或者使用 `geom_col()` 函数,这时不需要 `stat = "identity"`。如:

```
p <- ggplot(data=titanic,
 mapping = aes(
 x = fate,
 y = n,
 fill = sex))
p + geom_bar(position = "dodge", stat = "identity")
```



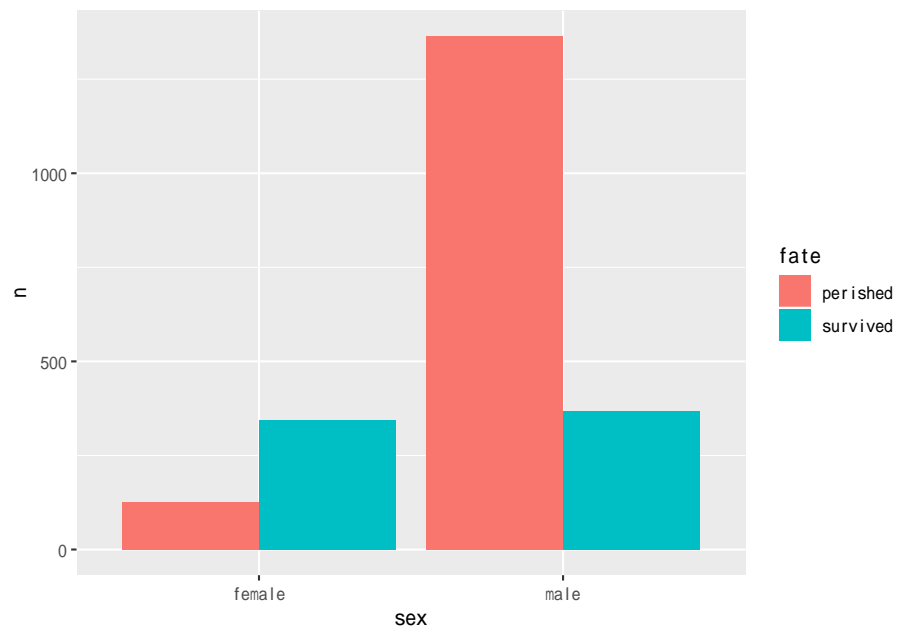
或:

```
p + geom_col(position = "dodge")
```



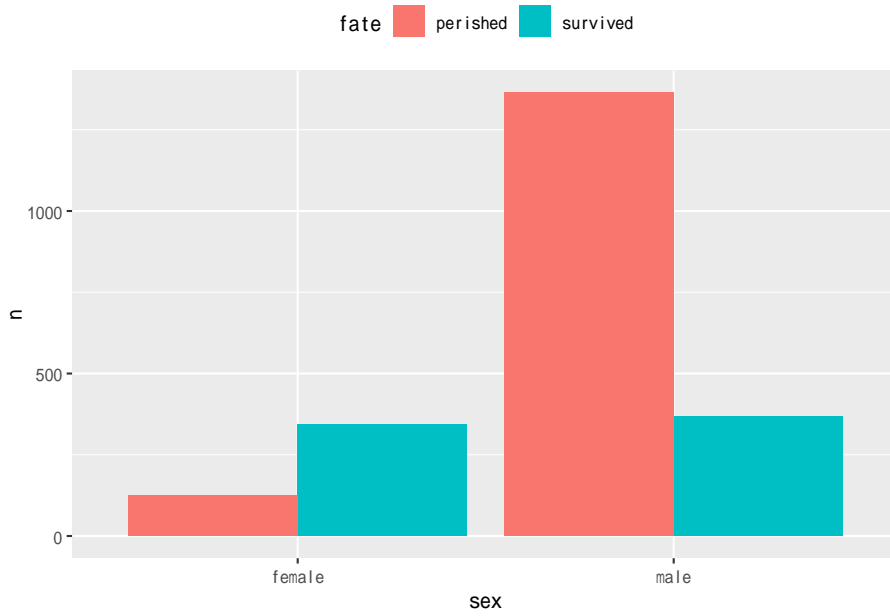
也可以按照性别分成大组:

```
p <- ggplot(data=titanic,
 mapping = aes(
 x = sex,
 y = n,
 fill = fate))
p + geom_col(position = "dodge")
```



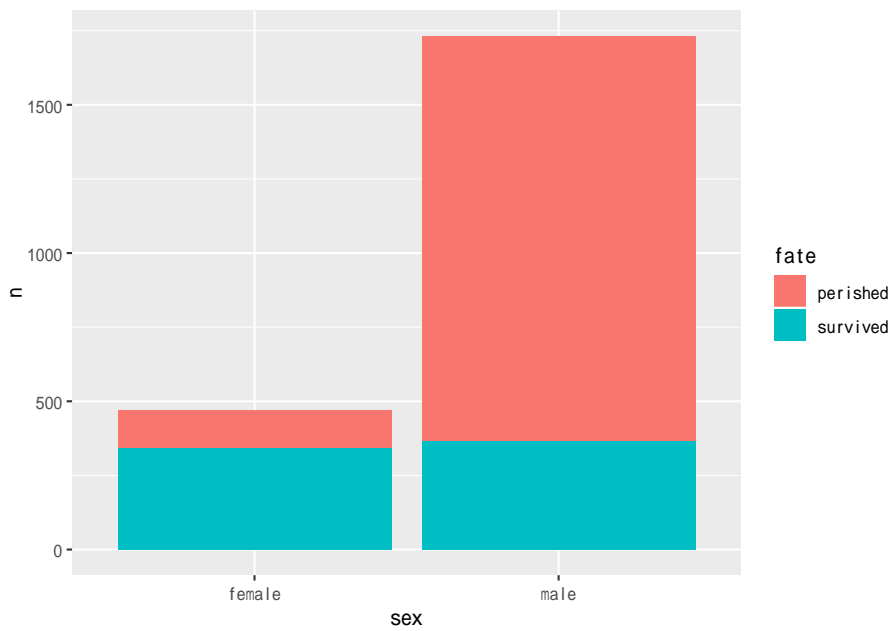
用 `theme()` 函数的 `legend.position` 参数可以指定图例的位置，如：

```
p + geom_col(position = "dodge") +
 theme(legend.position = "top")
```



做成堆叠形式:

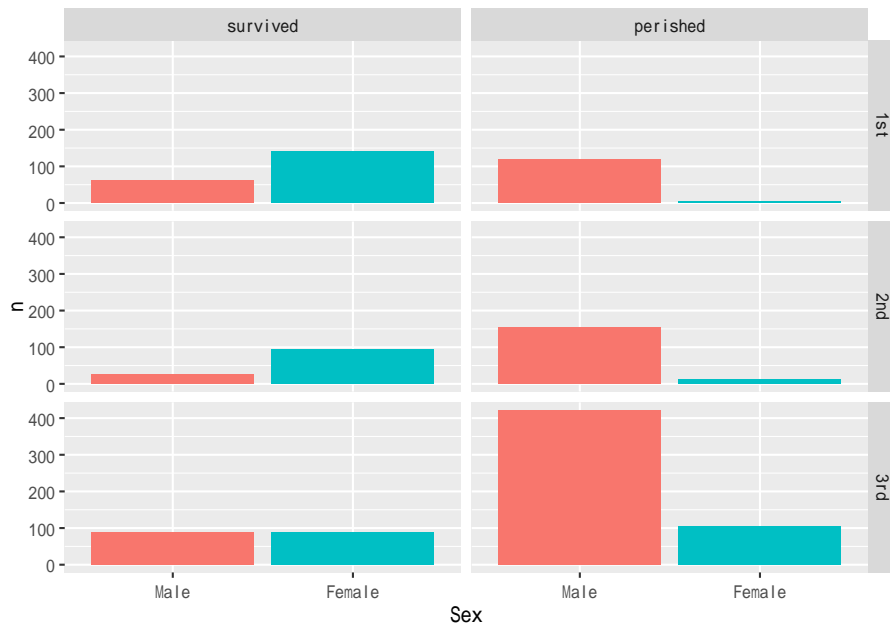
```
p + geom_col(position = "stack")
```



实际上，还可以用适当程序将存亡状态以及频数直接标在条形的色块内，`geom_text()` 函数可以在指定坐标位置标注指定的文字标签，见30.5.3。

`datasets` 包的 `Titanic` 数据集包含了泰坦尼克号乘客更详细的信息。我们按照存亡结果和舱位等级分小图作男女频数条形图：

```
titanic2 <- as.data.frame(Titanic) %>%
 group_by(Class, Sex, Survived) %>%
 summarise(n = sum(Freq)) %>%
 filter(Class != "Crew") %>%
 mutate(Survived = factor(
 Survived, levels = c("Yes", "No"),
 labels = c("survived", "perished")))
p <- ggplot(data = titanic2, mapping = aes(
 x = Sex, y = n, fill = Sex))
p + geom_col() +
 facet_grid(Class ~ Survived) +
 guides(fill = FALSE)
```



这里将 `fill` 映射到了 `Sex`，使得表示男女的条形填充了不同的颜色。如果不

满意上面的颜色，可以用 `scale_fill_manual()` 函数人为地指定颜色、对应离散值和图例标签。R 扩展包 `colourpicker` 提供了很好交互图形界面用来挑选颜色。

`geom_col()` 不仅限于画频数或者比例的条形图，此函数可以将一般用折线图表现的内容画成条形图，但一定要注意一点：`y` 坐标轴必须从 0 开始，这也是 `geom_col()` 和 `geom_bar()` 函数默认的设置。如果坐标轴不从零开始，则条形的长度就不能正确表示对应的 `y` 变量数值。

例如，`socviz` 包的 `oecd_sum` 数据集包含各年的美国以及 OECD 国家的期望寿命：

```
oecd_sum

A tibble: 57 x 5
Groups: year [57]
year other usa diff hi_lo
<int> <dbl> <dbl> <dbl> <chr>
1 1960 68.6 69.9 1.3 Below
2 1961 69.2 70.4 1.2 Below
3 1962 68.9 70.2 1.30 Below
4 1963 69.1 70 0.9 Below
5 1964 69.5 70.3 0.800 Below
6 1965 69.6 70.3 0.7 Below
7 1966 69.9 70.3 0.400 Below
8 1967 70.1 70.7 0.6 Below
9 1968 70.1 70.4 0.3 Below
10 1969 70.1 70.6 0.5 Below
... with 47 more rows
```

我们用 `geom_col()` 作 `diff` 变量的条形图，并按照 `hi_lo` 变量对正负差值分别使用不同颜色：

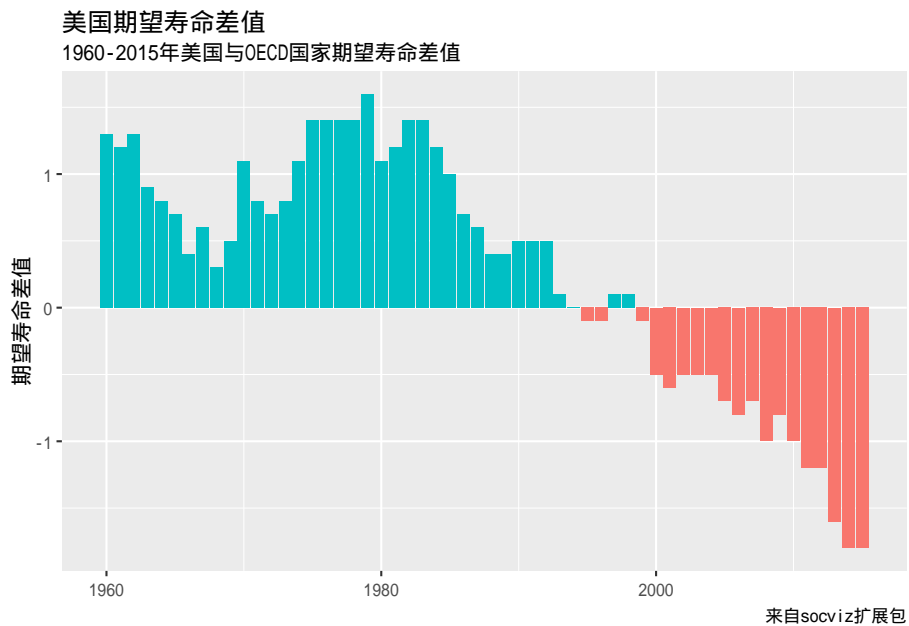
```
p <- ggplot(data = oecd_sum,
 mapping = aes(
 x = year,
 y = diff,
```

```

 fill = hi_lo))
p + geom_col() +
 guides(fill = FALSE) + # 正负号的不同颜色不使用图例标注
 labs(x = NULL,
 y = " 期望寿命差值",
 title = " 美国期望寿命差值",
 subtitle = "1960-2015 年美国与 OECD 国家期望寿命差值",
 caption=" 来自 socviz 扩展包")

```

```
Warning: Removed 1 rows containing missing values (position_stack).
```



## 30.5 更多图形种类

### 30.5.1 列联表生成与绘图

虽然 `ggplot2` 包能进行一些默认统计汇总，但是这些工作毕竟还是用专门的工具处理更好。`tidyverse` 系列的工具比较适用于进行数据整理和数据汇总，见第27章和第28章。



比如，将列联表用并列条形图表示，虽然 `ggplot2` 能直接从原始数据作图，但是用户比较容易混淆频数、比例、行比例、列比例的区别，最好是预先计算出作图所需的汇总表，用长表形式表示，即每个观测为某个分类或者交叉分类的频数。

考虑 `gss_sm` 数据集中分类变量 `bigregion` 与分类变量 `religion` 的列联表作图问题。频数表：

```
tab <- gss_sm %>%
 count(bigregion, religion)
```

```
Warning: Factor `religion` contains implicit NA, consider using
`forcats::fct_explicit_na`
```

```
tab
```

```
A tibble: 24 x 3
bigregion religion n
<fct> <fct> <int>
1 Northeast Protestant 158
2 Northeast Catholic 162
3 Northeast Jewish 27
4 Northeast None 112
5 Northeast Other 28
6 Northeast <NA> 1
7 Midwest Protestant 325
8 Midwest Catholic 172
9 Midwest Jewish 3
10 Midwest None 157
... with 14 more rows
```

生成了长表格式的列联表，保存成了 `tibble` 数据集，这也正是 `ggplot2` 包所需的格式。

将上述的长表显示为宽表：

```
tab2 <- xtabs(n ~ bigregion + religion, data=tab)
tab2
```

```
religion
```

```
bigregion Protestant Catholic Jewish None Other
Northeast 158 162 27 112 28
Midwest 325 172 3 157 33
South 650 160 11 170 50
West 238 155 10 180 48
```

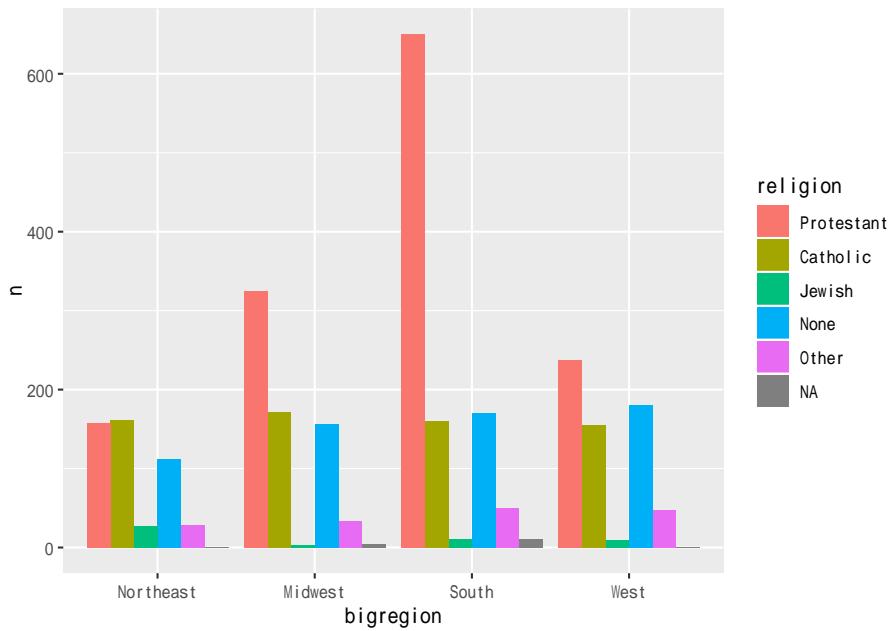
从原始数据生成这样的宽表，程序为：

```
tab2 <- gss_sm %>%
 table(bigregion, religion)
tab2
```

```
religion
bigregion Protestant Catholic Jewish None Other
Northeast 158 162 27 112 28
Midwest 325 172 3 157 33
South 650 160 11 170 50
West 238 155 10 180 48
```

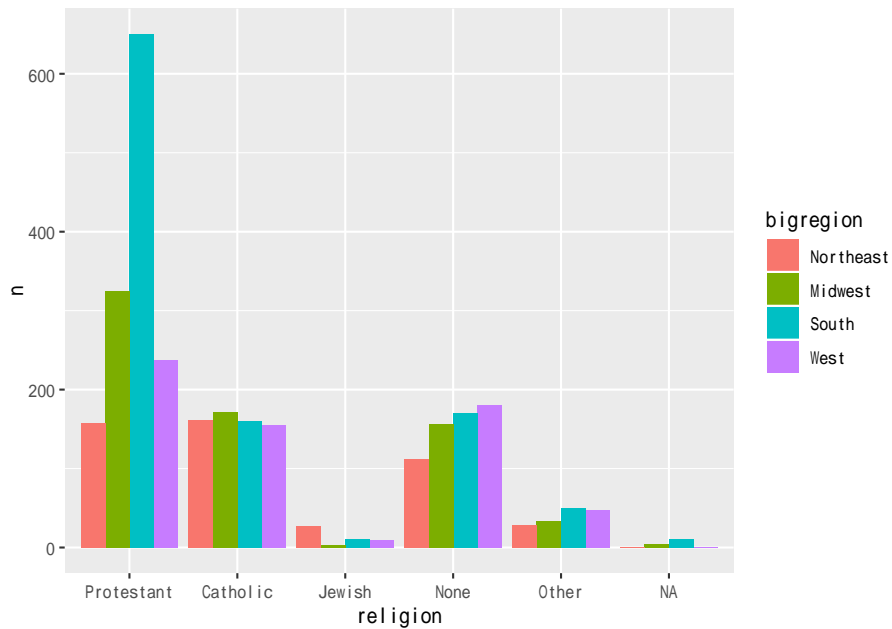
将 bigregion(大区) 作为大类，每个大区内做出各个宗教人数的频数条形图：

```
p <- ggplot(data=tab,
 mapping = aes(
 x = bigregion,
 fill = religion,
 y = n))
p + geom_bar(stat = "identity", position = "dodge")
```



使用同一个频数表，将 religion(宗教) 作为大类，在每种宗教中作不同大区的频数条形图：

```
p <- ggplot(data=tab,
 mapping = aes(
 x = religion,
 fill = bigregion,
 y = n))
p + geom_bar(stat = "identity", position = "dodge")
```



我们希望表现每个大区内的各种宗教的比例, 使得每个大区的比例之和等于 1, 从宽表格式可以用 `prop.table()` 函数直接得到列百分比表:

```
tab2a <- prop.table(tab2, 2)
round(tab2a, 2)
```

```
religion
bigregion Protestant Catholic Jewish None Other
Northeast 0.12 0.25 0.53 0.18 0.18
Midwest 0.24 0.27 0.06 0.25 0.21
South 0.47 0.25 0.22 0.27 0.31
West 0.17 0.24 0.20 0.29 0.30
```

为了得到 `ggplot2` 所需格式, 可以用 `as.data.frame(tab2a)` 转换。

如果从原始数据直接统计, 可以利用 `dplyr` 包的 `group_by()` 分组汇总功能:

```
tab3 <- gss_sm %>%
 group_by(bigregion, religion) %>%
 summarize(N = n()) %>%
 mutate(prop = N / sum(N),
```

```
pct = round(100 * prop, 2))

Warning: Factor `religion` contains implicit NA, consider using
`forcats::fct_explicit_na`

tab3

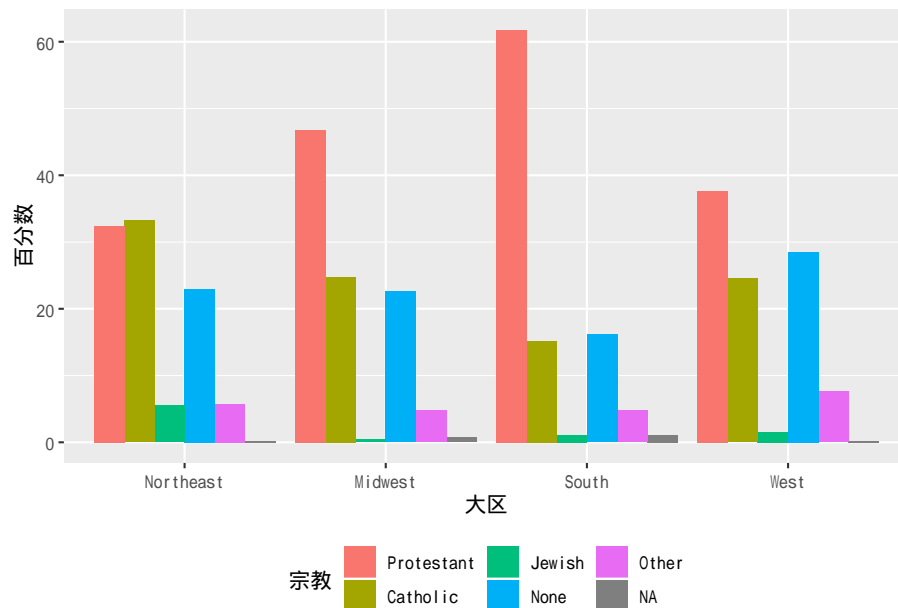
A tibble: 24 x 5
Groups: bigregion [4]
bigregion religion N prop pct
<fct> <fct> <int> <dbl> <dbl>
1 Northeast Protestant 158 0.324 32.4
2 Northeast Catholic 162 0.332 33.2
3 Northeast Jewish 27 0.0553 5.53
4 Northeast None 112 0.230 23.0
5 Northeast Other 28 0.0574 5.74
6 Northeast <NA> 1 0.00205 0.2
7 Midwest Protestant 325 0.468 46.8
8 Midwest Catholic 172 0.247 24.8
9 Midwest Jewish 3 0.00432 0.43
10 Midwest None 157 0.226 22.6
... with 14 more rows
```

程序中用 `group_by()` 做了交叉分组, `n()` 计算每个交叉组的频数, 而 `sum(N)` 并不是计算总频数, 而是每个 `bigregion` 内部的频数。在交叉分组时, `sum()`, `mean()` 这些汇总函数是针对最内层的数据进行汇总的。

利用上述的数据, 作每个大区内各种宗教人数百分数的条形图:

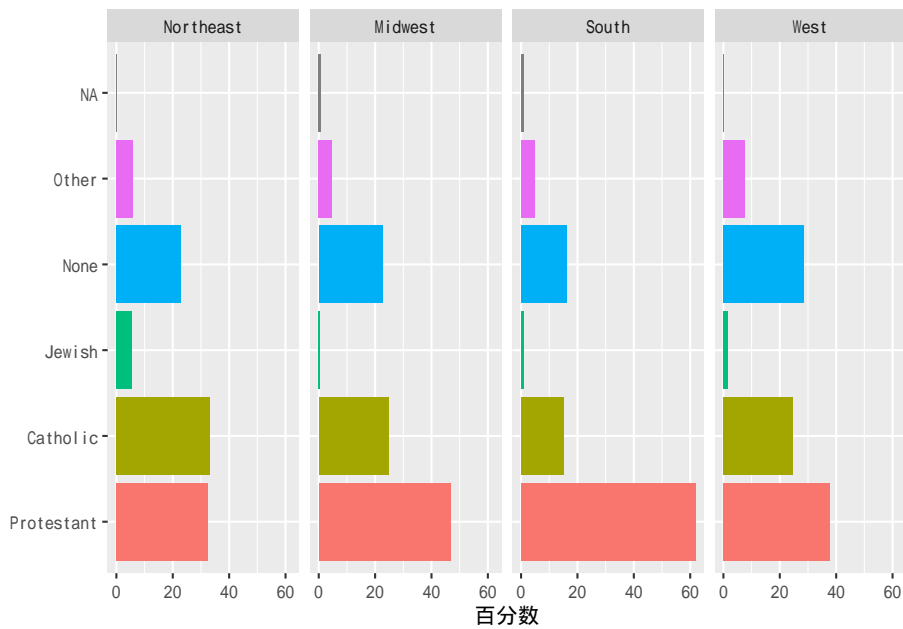
```
p <- ggplot(data=tab3,
 mapping = aes(
 x = bigregion,
 fill = religion,
 y = pct))
p + geom_bar(stat = "identity", position = "dodge") +
 labs(x = "大区", y = "百分数", fill="宗教") +
```

```
theme(legend.position = "bottom")
```



下面将每个大区的宗教百分数分布做成一幅小图，用横向条形图，作小图的程序中不涉及 bigregion 变量：

```
p <- ggplot(data=tab3,
 mapping = aes(
 x = religion,
 fill = religion,
 y = pct))
p + geom_bar(stat = "identity", position = "dodge") +
 labs(x = NULL, y = "百分数") +
 guides(fill = FALSE) +
 coord_flip() +
 facet_grid(~ bigregion)
```



### 30.5.2 连续变量的分组图形

选用 `socviz` 包的 `organdata` 数据集，这是 17 个 OECD 国家历年的器官捐献情况以及一些其他变量的记录。其中前 6 列的一些抽样数据：

```
organdata %>%
 select(1:6) %>%
 sample_n(size = 10)
```

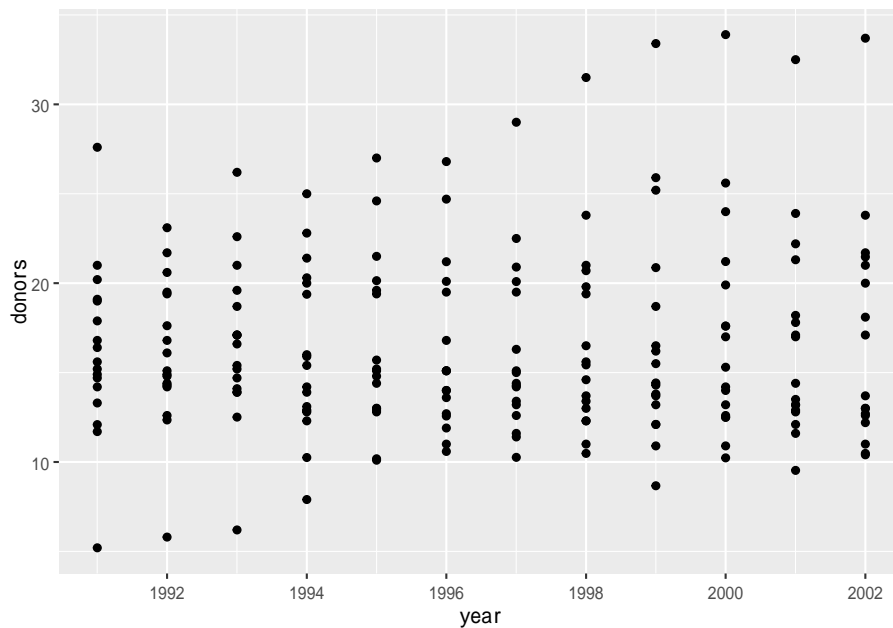
```
A tibble: 10 x 6
country year donors pop pop_dens gdp
<chr> <date> <dbl> <int> <dbl> <int>
1 United States 1993-01-01 18.7 259919 2.70 25327
2 Norway 2000-01-01 17.6 4491 1.39 35829
3 Spain 2000-01-01 33.9 40614 8.03 20017
4 Finland 1994-01-01 20 5088 1.50 17993
5 United States 1991-01-01 17.9 252981 2.63 23443
6 Spain 1996-01-01 26.8 39279 7.76 16416
```

```
7 Switzerland 2000-01-01 14 7184 17.4 29837
8 Sweden NA NA NA NA NA
9 Belgium 1995-01-01 19.6 10137 30.6 21679
10 Canada 2001-01-01 13.5 31111 0.312 29235
```

变量 `donors` 是每百万人中器官捐献数。作 `donors` 对 `year` 的散点图：

```
p <- ggplot(data = organdata,
 mapping = aes(
 x = year,
 y = donors))
p + geom_point()
```

```
Warning: Removed 34 rows containing missing values (geom_point).
```



每年有多个数值，是不同国家的捐献数。这个图形不能反映一种时间趋势，不太有用。

可以对每个国家画一条折线图：

```
p <- ggplot(data = organdata,
 mapping = aes(
```

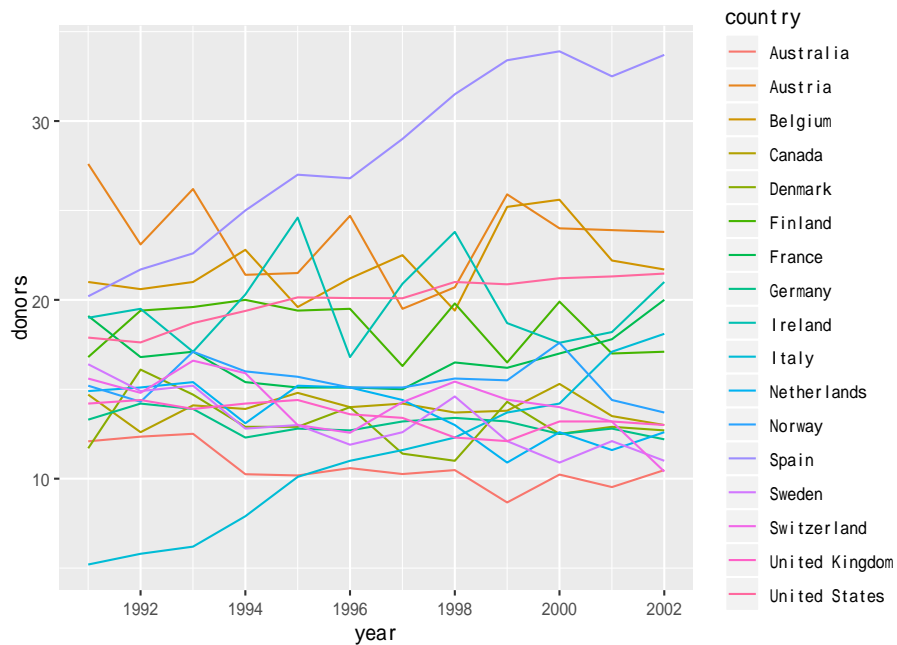


```

 x = year,
 y = donors,
 color = country))
p + geom_line()

```

```
Warning: Removed 34 rows containing missing values (geom_path).
```



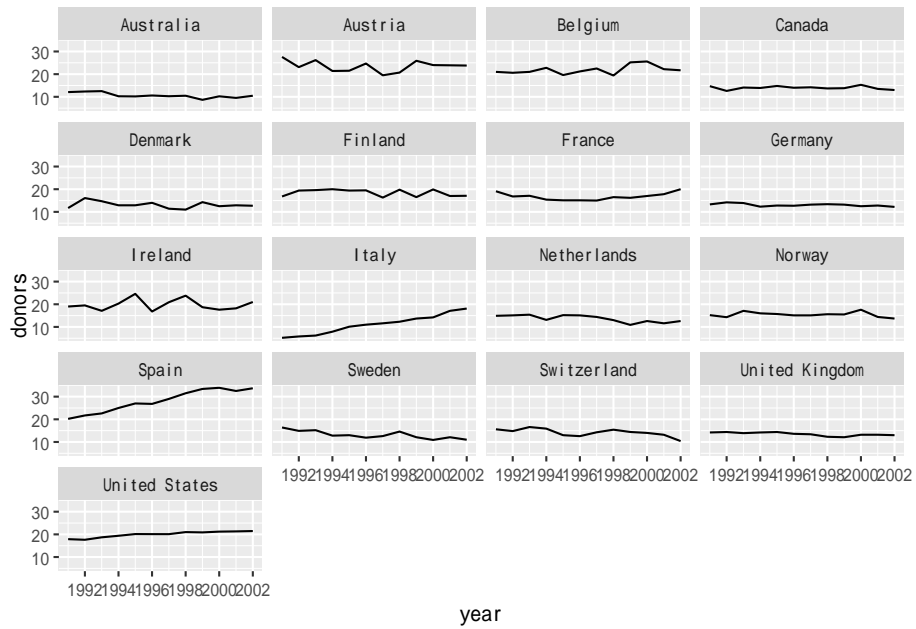
共有 17 个国家，每个国家做了器官捐赠率随时间变化的折线图。用小图的方法将其分配到不同的小图：

```

p <- ggplot(data = organdata,
 mapping = aes(
 x = year,
 y = donors))
p + geom_line() +
 facet_wrap(~ country, ncol=4)

```

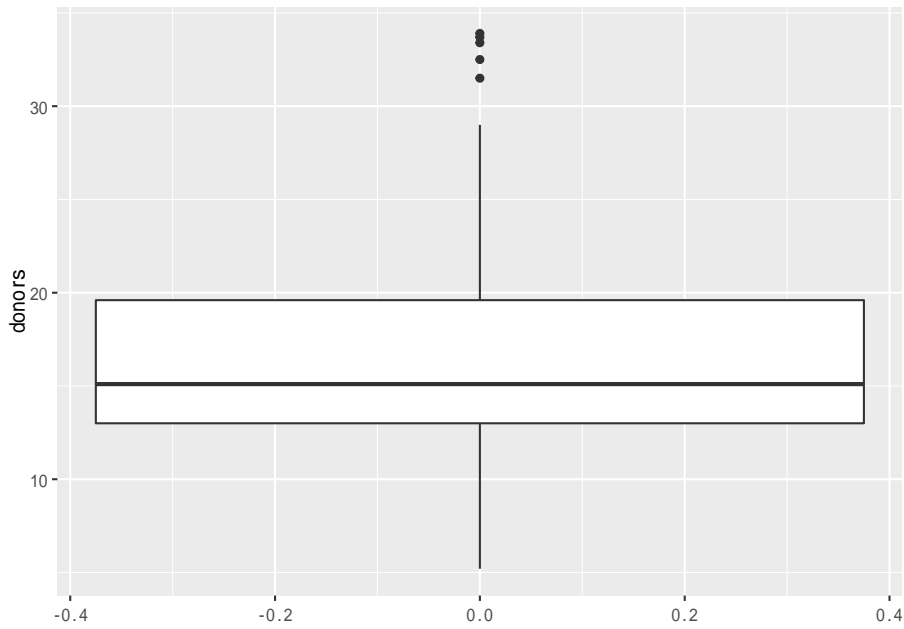
```
Warning: Removed 2 rows containing missing values (geom_path).
```



用 `geom_boxplot()` 可以做盒形图，能够画出连续型变量的主要分位数，表现变量分布，如：

```
p <- ggplot(data = organdata,
 mapping = aes(y = donors))
p + geom_boxplot()
```

```
Warning: Removed 34 rows containing non-finite values (stat_boxplot).
```

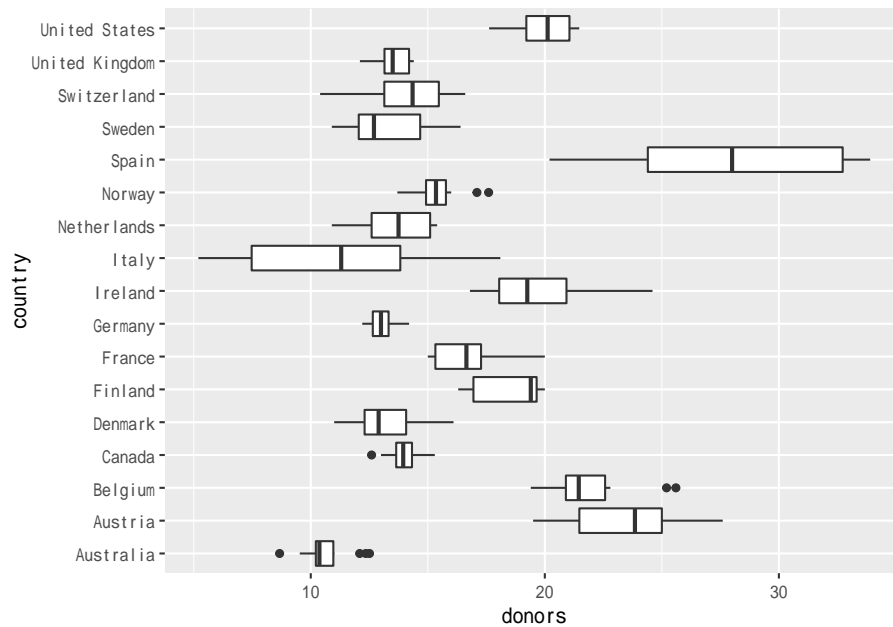


这是所有国家所有年的捐献率分布情况。类似函数还有 `geom_violin()`。

每个国家的捐献率单独做盒形图并且放在同一坐标系中：

```
p <- ggplot(data = organdata,
 mapping = aes(y = donors, x = country))
p + geom_boxplot() + coord_flip()
```

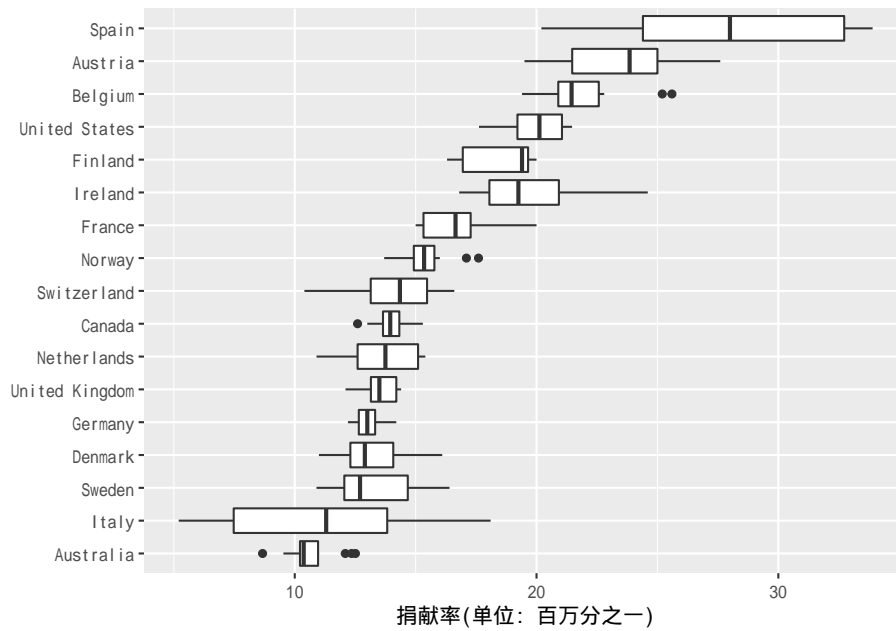
```
Warning: Removed 34 rows containing non-finite values (stat_boxplot).
```



这个图形很好地比较了不同国家的历年捐献率的分布，比如，Spain 的捐献率最高。为了将图形中的各个国家按照捐献率的某个统计量排序，可以使用 stats 包的 `reorder()` 函数，调整因子的水平次序。为了能够比较容易地标出国家名称，交换 x 轴与 y 轴的作用，如：

```
p <- ggplot(data = organdata,
 mapping = aes(
 y = donors,
 x = reorder(country, donors, median, na.rm=TRUE)))
p + geom_boxplot() +
 coord_flip() +
 labs(y = " 捐献率 (单位: 百万分之一)",
 x = NULL)
```

```
Warning: Removed 34 rows containing non-finite values (stat_boxplot).
```

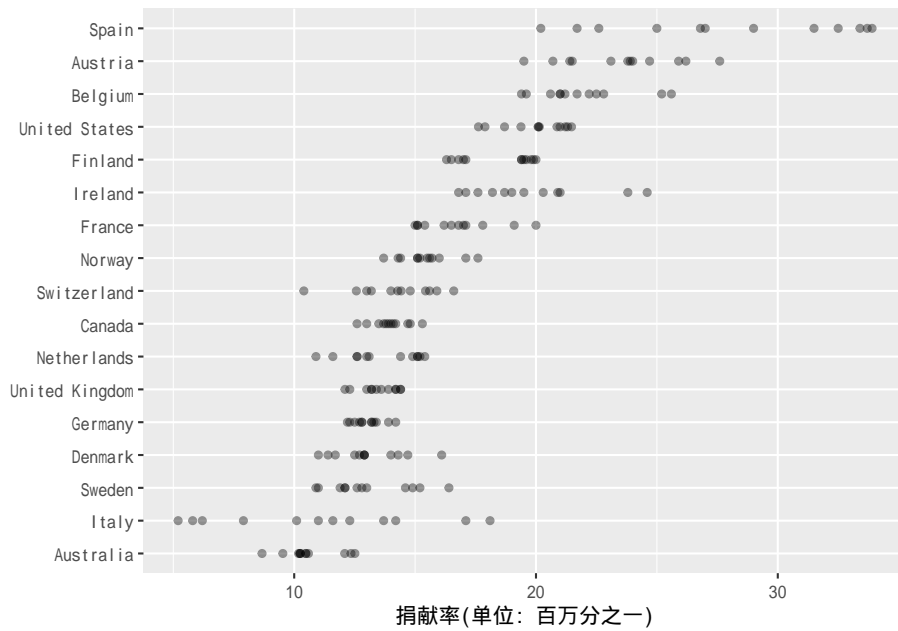


这里盒形图是横向的，如果仍然纵向作图，国家名称在横轴，许多个国家名称就会重叠在一起，只好仅显示其中一部分名称。

当每组（这里是每个国家）的观测个数很少时，也可以做成散点图，如：

```
p <- ggplot(data = organdata,
 mapping = aes(
 y = donors,
 x = reorder(country, donors, median, na.rm=TRUE)))
p + geom_point(alpha = 0.4) +
 coord_flip() +
 labs(y = " 捐献率 (单位: 百万分之一)",
 x = NULL)
```

```
Warning: Removed 34 rows containing missing values (geom_point).
```

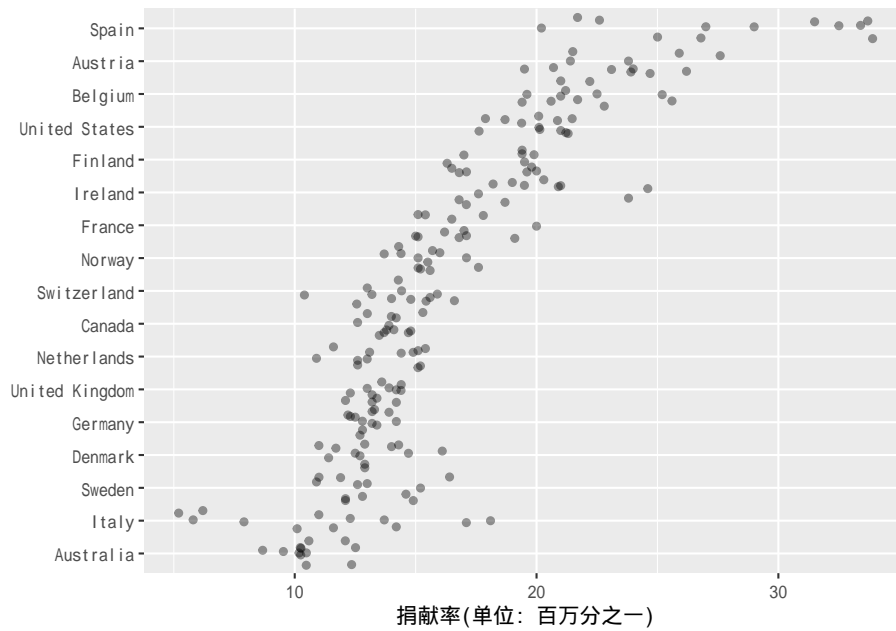


因为有的观测点完全重叠，所以用了 `alpha` 参数指定一定的透明度，重叠越多的点显示的颜色越深。但是，如果两个不同颜色的点完全重叠，半透明不能显示两个不同颜色的效果。

在作这样的散点图时，为了避免重叠的点，可以将 `geom_point()` 改为 `geom_jitter()`，如：

```
p + geom_jitter(alpha = 0.4) +
 coord_flip() +
 labs(y = " 捐献率 (单位: 百万分之一)",
 x = NULL)
```

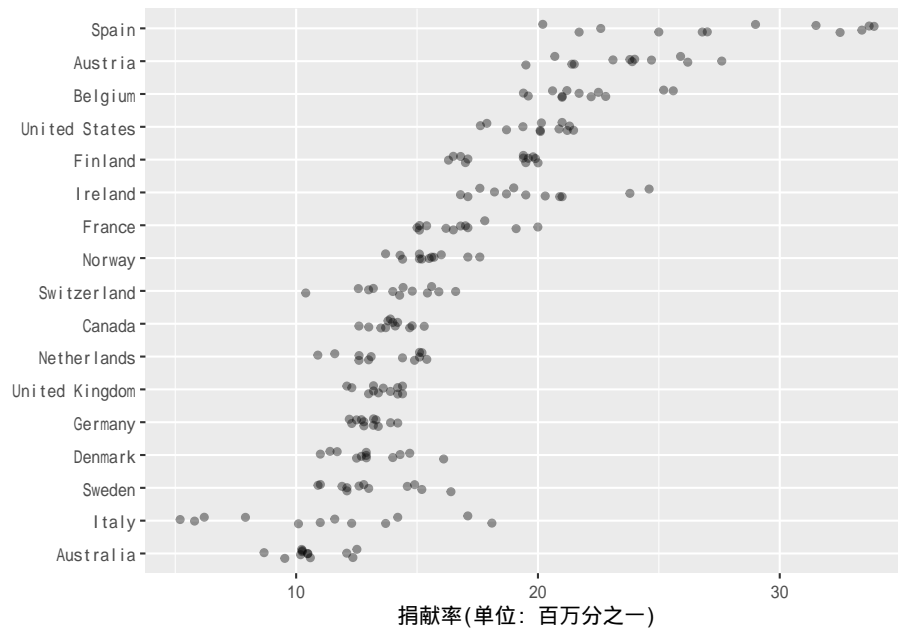
```
Warning: Removed 34 rows containing missing values (geom_point).
```



上图的点的扰动过大了，使得不同国家的区分不明显了。作扰动的散点图时，可以用 `width` 指定左右扰动范围，用 `height` 指定上下扰动范围，这里只需要指定左右扰动范围，因为坐标轴对调所以就变成了上下扰动：

```
p + geom_jitter(alpha = 0.4, width = 0.15) +
 coord_flip() +
 labs(y = " 捐献率 (单位: 百万分之一)",
 x = NULL)
```

```
Warning: Removed 34 rows containing missing values (geom_point).
```

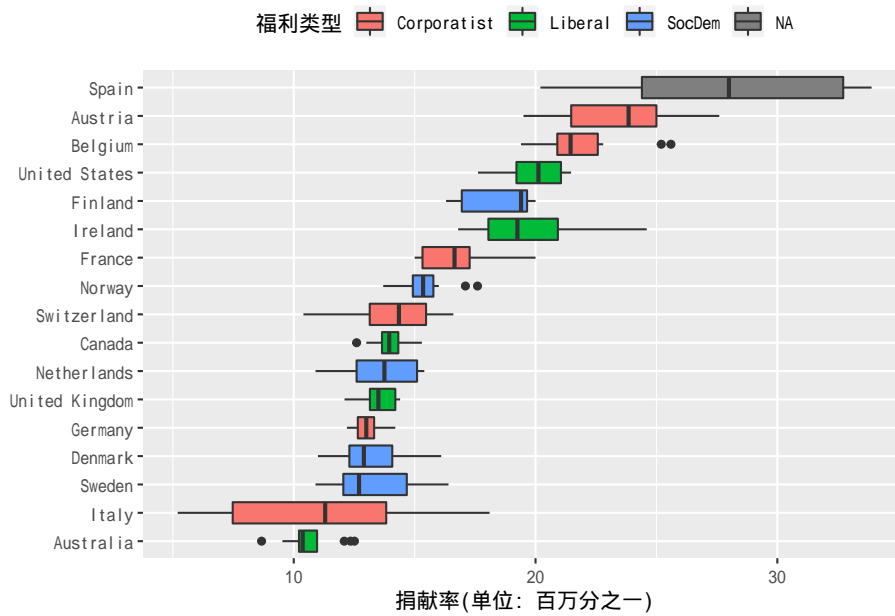


`geom_boxplot()` 也支持 `color`, `fill` 维度。`organdata` 中的变量 `world` 是一个国家的福利类型，用不同填充色表示 `world` 变量：

```
p <- ggplot(data = organdata,
 mapping = aes(
 y = donors,
 x = reorder(country, donors, median, na.rm=TRUE),
 fill = world))
p + geom_boxplot() +
 coord_flip() +
 labs(y = " 捐献率 (单位: 百万分之一)",
 x = NULL,
 fill = " 福利类型") +
 theme(legend.position = "top")
```

```
Warning: Removed 34 rows containing non-finite values (stat_boxplot).
```





下面做不同国家的平均捐赠率的图形。首先得到统计数据:

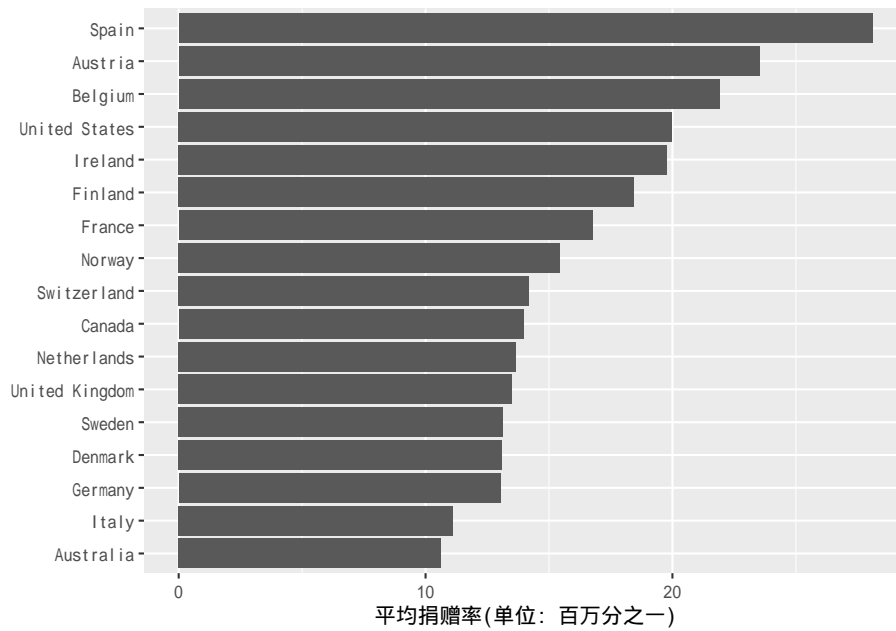
```
organdata2 <- organdata %>%
 group_by(country) %>%
 summarize(
 donors_mean = mean(donors, na.rm=TRUE),
 donors_sd = sd(donors, na.rm=TRUE))
organdata2
```

```
A tibble: 17 x 3
country donors_mean donors_sd
<chr> <dbl> <dbl>
1 Australia 10.6 1.14
2 Austria 23.5 2.42
3 Belgium 21.9 1.94
4 Canada 14.0 0.751
5 Denmark 13.1 1.47
6 Finland 18.4 1.53
7 France 16.8 1.60
```

##	8	Germany	13.0	0.611
##	9	Ireland	19.8	2.48
##	10	Italy	11.1	4.28
##	11	Netherlands	13.7	1.55
##	12	Norway	15.4	1.11
##	13	Spain	28.1	4.96
##	14	Sweden	13.1	1.75
##	15	Switzerland	14.2	1.71
##	16	United Kingdom	13.5	0.775
##	17	United States	20.0	1.33

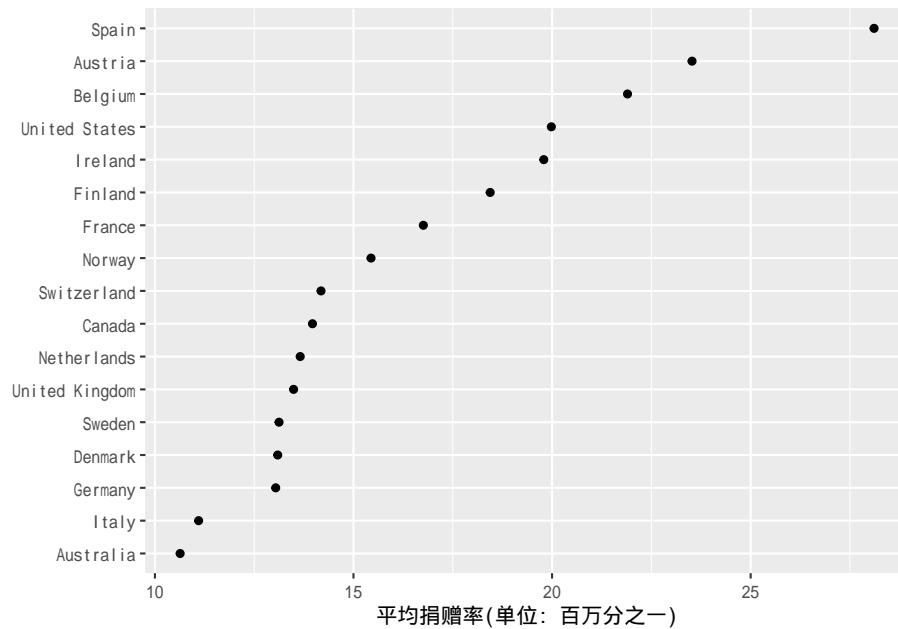
用条形图表现不同国家的平均捐献率:

```
p <- ggplot(data = organdata2,
 mapping = aes(
 x = reorder(country, donors_mean),
 y = donors_mean))
p + geom_col() +
 coord_flip() +
 labs(x = NULL, y = " 平均捐赠率 (单位: 百万分之一)")
```



这样的图形也可以做成点图，称为 Cleveland 点图。不需要再颠倒横纵坐标，直接规定 x 轴为平均捐赠率即可：

```
p <- ggplot(data = organdata2,
 mapping = aes(
 y = reorder(country, donors_mean),
 x = donors_mean))
p + geom_point() +
 labs(y = NULL, x = "平均捐赠率 (单位: 百万分之一)")
```



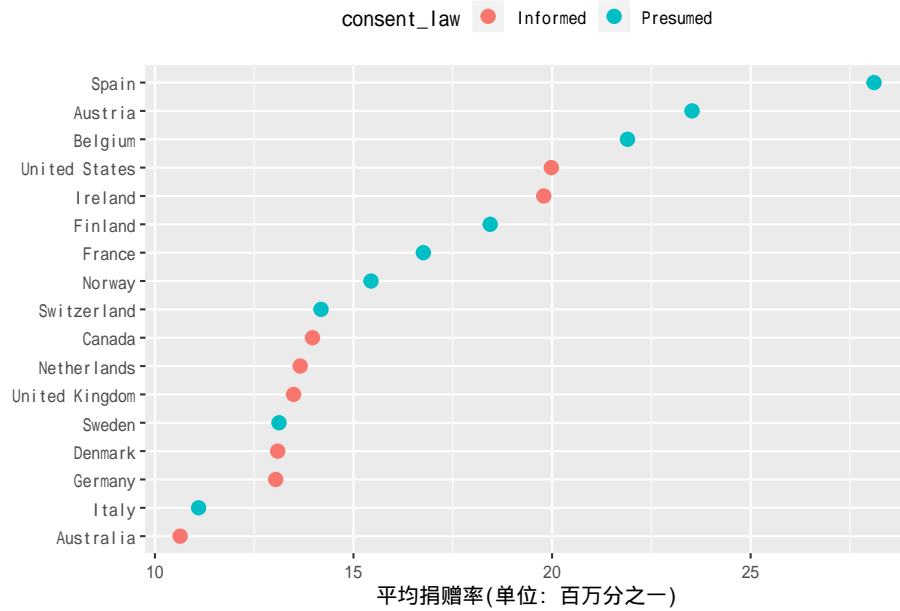
organdata 数据集中变量 `consent_law` 是关于一个国家中器官捐赠是必须告知还是默认捐赠的区别。为了在上面的点图中用不同颜色区分这两种做法，需要在分组汇总阶段就将 `consent_law` 也作为分组变量。这是因为 `summarize()` 函数会自动舍弃分组变量和统计结果之外的原有变量。

```
organdata3 <- organdata %>%
 group_by(consent_law, country) %>%
 summarize(
 donors_mean = mean(donors, na.rm=TRUE),
 donors_sd = sd(donors, na.rm=TRUE)) %>%
 ungroup()
organdata3
```

```
A tibble: 17 x 4
consent_law country donors_mean donors_sd
<chr> <chr> <dbl> <dbl>
1 Informed Australia 10.6 1.14
2 Informed Canada 14.0 0.751
3 Informed Denmark 13.1 1.47
```

##	4	Informed	Germany	13.0	0.611
##	5	Informed	Ireland	19.8	2.48
##	6	Informed	Netherlands	13.7	1.55
##	7	Informed	United Kingdom	13.5	0.775
##	8	Informed	United States	20.0	1.33
##	9	Presumed	Austria	23.5	2.42
##	10	Presumed	Belgium	21.9	1.94
##	11	Presumed	Finland	18.4	1.53
##	12	Presumed	France	16.8	1.60
##	13	Presumed	Italy	11.1	4.28
##	14	Presumed	Norway	15.4	1.11
##	15	Presumed	Spain	28.1	4.96
##	16	Presumed	Sweden	13.1	1.75
##	17	Presumed	Switzerland	14.2	1.71

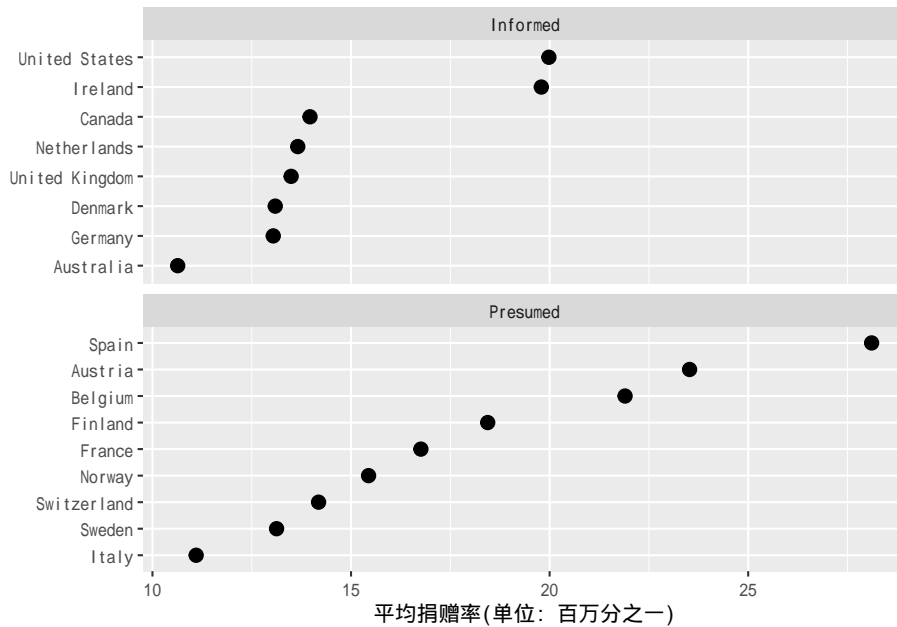
```
p <- ggplot(data = organdata3,
 mapping = aes(
 y = reorder(country, donors_mean),
 x = donors_mean,
 color = consent_law))
p + geom_point(size = 3) +
 labs(y = NULL, x = " 平均捐赠率 (单位: 百万分之一)") +
 theme(legend.position = "top")
```



平均捐赠率最高的三个国家都是不需告知预先假定同意的。

也可以将两种告知规定分成两个小图：

```
p <- ggplot(data = organdata3,
 mapping = aes(
 y = reorder(country, donors_mean),
 x = donors_mean))
p + geom_point(size = 3) +
 facet_wrap(~ consent_law, ncol=1, scales = "free_y") +
 labs(y = NULL, x = " 平均捐赠率 (单位: 百万分之一)")
```



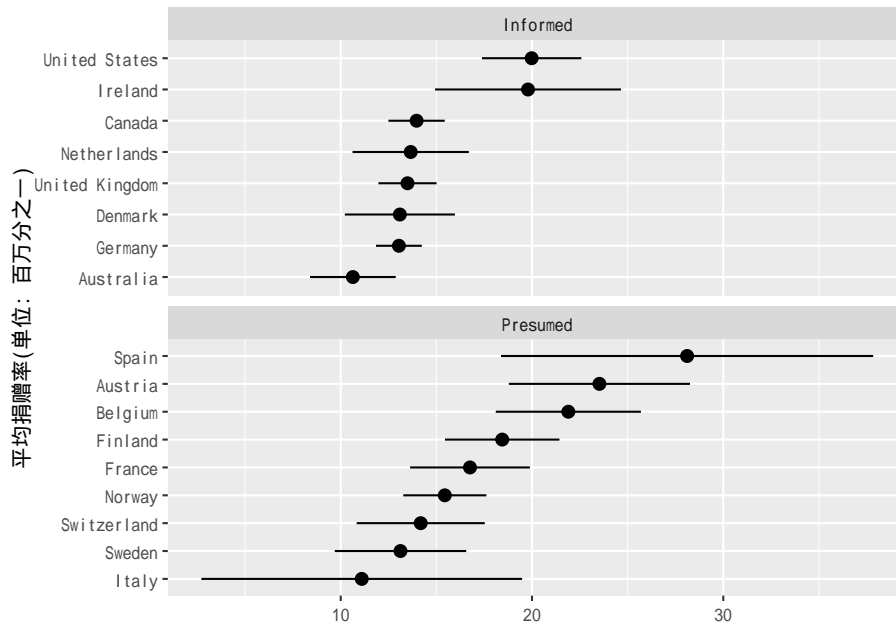
因为纵轴是分类变量，程序中的 `scales = "free_y"` 使得纵轴仅对存在的类留出空间。用了 `ncol = 1` 使得两种告知规定的小图上下排列，便于比较横坐标值。

当每个类别仅有一个数值时，一般推荐使用 Cleveland 点图，而不是条形图或者折线图。Cleveland 点图总是将类别值绘制在 y 轴，将要比较的数量值用 x 坐标表示，并将各类按照数量值大小次序排列。

可以表示平均值的点图上增加一条线，表示误差大小，所用函数为 `geom_pointrange()`。这个函数仅支持对 y 轴加误差线，所以需要交换坐标轴的办法将分类变量放在 y 轴。比如，画出近似 95% 置信区间范围：

```
p <- ggplot(data = organdata3,
 mapping = aes(
 x = reorder(country, donors_mean),
 y = donors_mean))
p + geom_pointrange(
 mapping = aes(ymin = donors_mean - 1.96*donors_sd,
 ymax = donors_mean + 1.96*donors_sd)) +
coord_flip() +
```

```
facet_wrap(~ consent_law, ncol=1, scales = "free_y") +
labs(y = NULL, x = " 平均捐赠率 (单位: 百万分之一)")
```



类似的函数还有 `geom_linerange()`、`geom_crossbar()`、`geom_errorbar()`。

### 30.5.3 坐标系中的文字

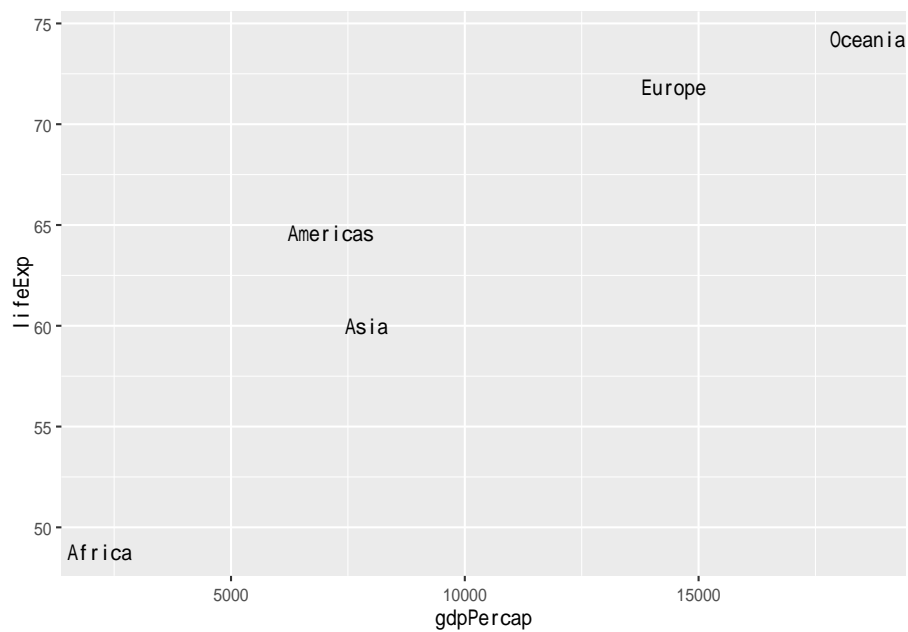
类似于散点图,可以将指定的文字绘制在指定的坐标位置,使用 `geom_text(mapping = aes(label = 字符型变量))`。

例如, `gapminder` 数据集中各大洲的平均寿命与平均 `gdp` 的文字散点图:

```
gapminder2 <- gapminder %>%
 group_by(continent) %>%
 summarize(lifeExp = mean(lifeExp, na.rm=TRUE),
 gdpPercap = mean(gdpPercap))
p <- ggplot(data=gapminder2,
 mapping = aes(
 x = gdpPercap,
```

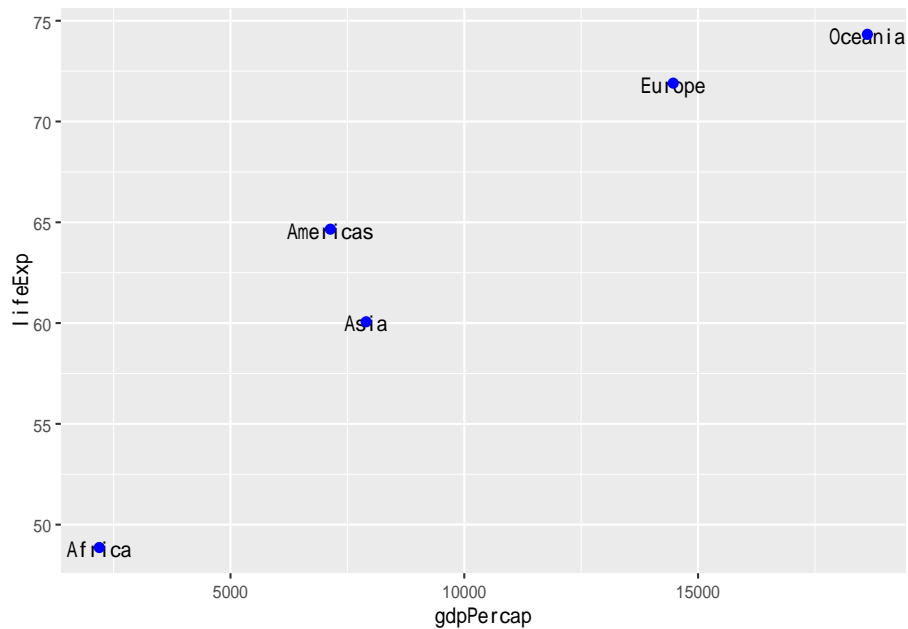


```
 y = lifeExp,
 label = continent))
p + geom_text()
```



可以同时绘制散点:

```
p + geom_text() +
 geom_point(size = 2, col="blue")
```



文字默认是居中对齐，加选项 `hjust = 0` 表示左对齐，`hjust = 1` 表示右对齐。

`ggrepel` 扩展包提供了增强的图形文本功能。`geom_text_repel()` 提供了与 `geom_text()` 类似的功能。

考虑 `socviz` 扩展包的 `elections_historic` 数据集，这是美国历次总统选举情况数据。部分数据显示：

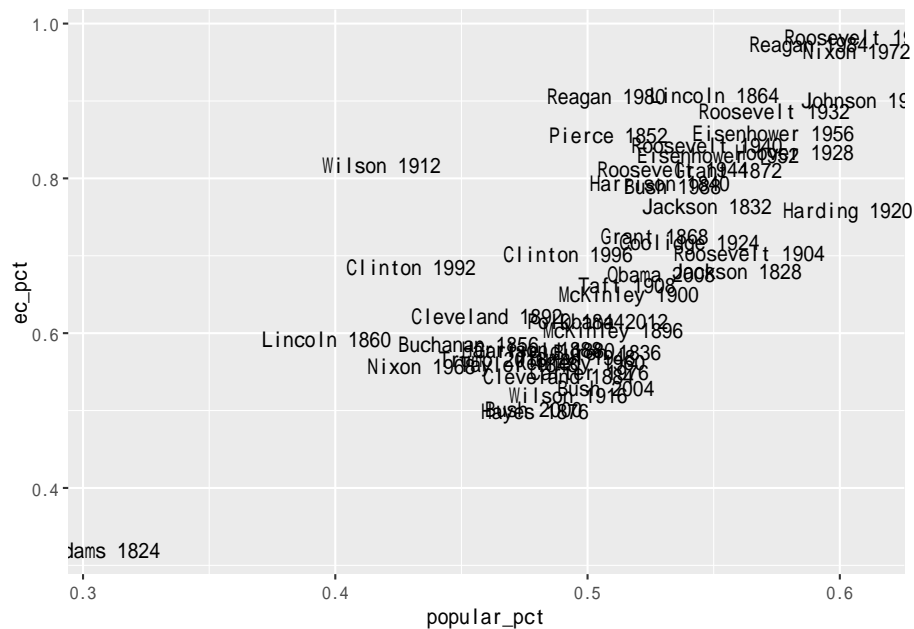
```
elections_historic %>%
 select(2:7) %>% head(10)
```

```
A tibble: 10 x 6
year winner win_party ec_pct popular_pct popular_margin
<int> <chr> <chr> <dbl> <dbl> <dbl>
1 1824 John Quincy Adams D.-R. 0.322 0.309 -0.104
2 1828 Andrew Jackson Dem. 0.682 0.559 0.122
3 1832 Andrew Jackson Dem. 0.766 0.547 0.178
4 1836 Martin Van Buren Dem. 0.578 0.508 0.142
5 1840 William Henry Harrison Whig 0.796 0.529 0.0605
```

##	6	1844	James Polk	Dem.	0.618	0.495	0.0145
##	7	1848	Zachary Taylor	Whig	0.562	0.473	0.0479
##	8	1852	Franklin Pierce	Dem.	0.858	0.508	0.0695
##	9	1856	James Buchanan	Dem.	0.588	0.453	0.122
##	10	1860	Abraham Lincoln	Rep.	0.594	0.396	0.101

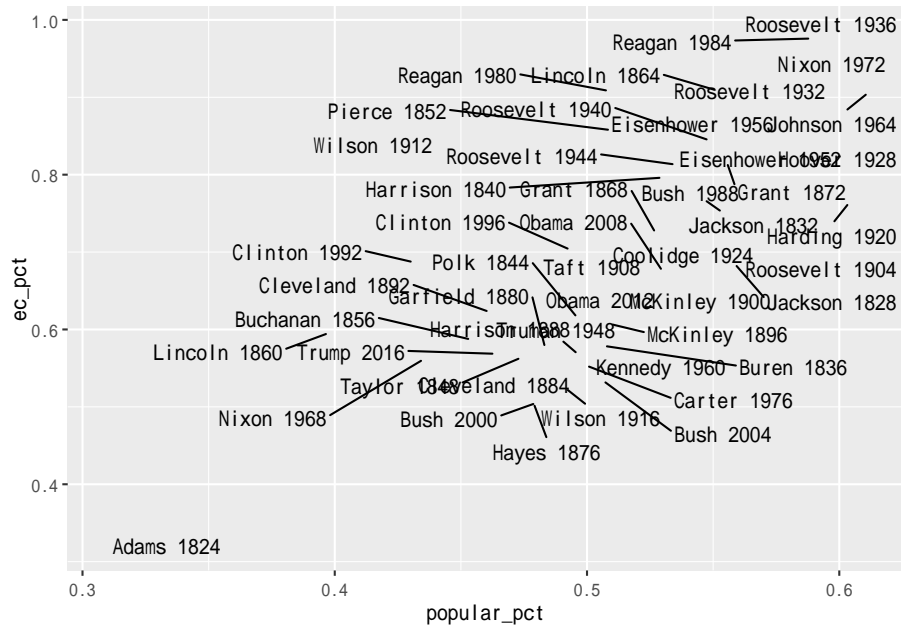
取 `popular_pct`(popular 投票支持率) 为横坐标, 取 `ec_pct`(election college 投票支持率) 为纵坐标, 将历次结果标在坐标系中:

```
p <- ggplot(data = elections_historic,
 mapping = aes(
 x = popular_pct,
 y = ec_pct,
 label = winner_label))
p + geom_text()
```



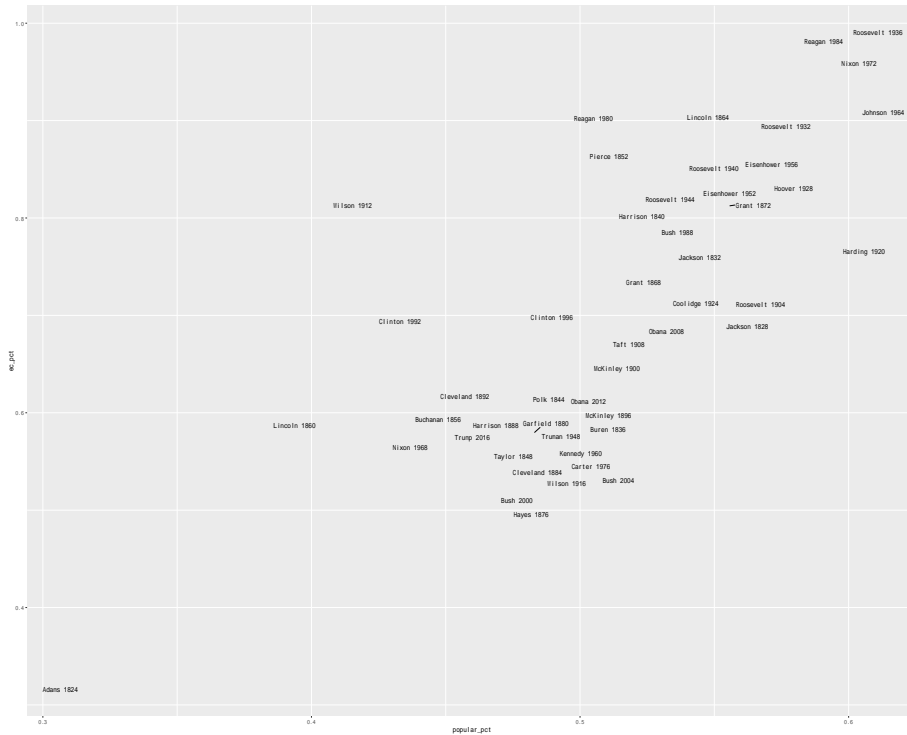
因为点比较多, 文字也比较长, 有很多重叠。ggrepel 包的 `geom_text_repel()` 则很好地处理了这个问题:

```
library(ggrepel)
p + geom_text_repel()
```



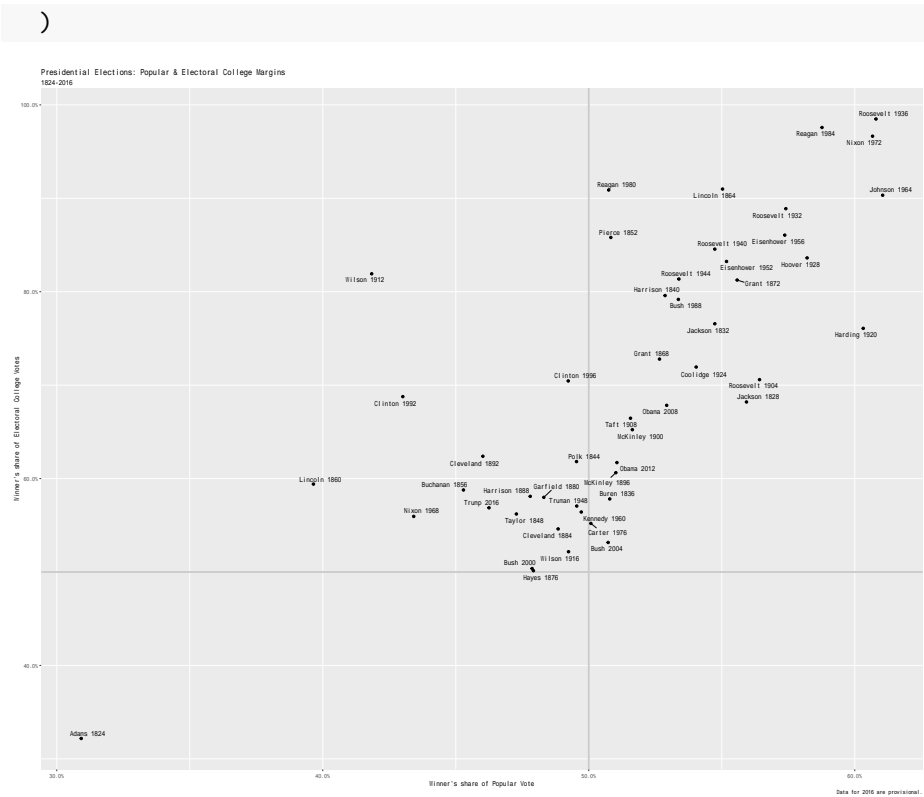
可以看出，其解决重叠问题的方式是用短线指向实际的坐标位置。在 Rmd 文件中，还是有少量的重叠，可以通过在 R 代码段选项中增大 `fig.width` 和 `fig.height` 参数实现，下面的代码段用了选项 `fig.width=20`, `fig.height=16`:

```
p + geom_text_repel()
```



可以用 `scale_x_continuous()` 和 `scale_y_continuous()` 将坐标轴的比例值转换成百分数，用 `geom_hline(yintercept)` 添加横线，用 `geom_vline(xintercept)` 添加竖线，适当地用标注改善图形：

```
p + geom_hline(yintercept = 0.5, size = 1.4, col = "gray80") +
 geom_vline(xintercept = 0.5, size = 1.4, col = "gray80") +
 geom_point() +
 geom_text_repel() +
 scale_x_continuous(labels = scales::percent) +
 scale_y_continuous(labels = scales::percent) +
 labs(
 x = "Winner's share of Popular Vote",
 y = "Winner's share of Electoral College Votes",
 title = "Presidential Elections: Popular & Electoral College Margins",
 subtitle = "1824-2016",
 caption = "Data for 2016 are provisional."
```



为了画斜线，可以用 `geom_abline()` 函数。

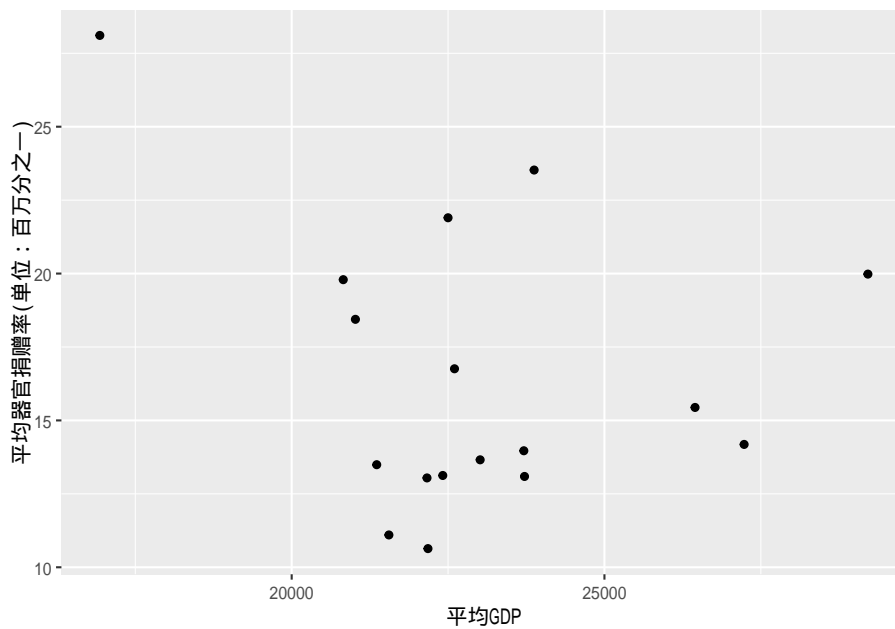
### 30.5.4 标出特殊点

在坐标系中标注文字的功能更经常用来标出图形中的特殊点。

考虑 `organdata` 中各国的平均捐赠率数据。作平均捐赠率对平均 gdp 的散点图：

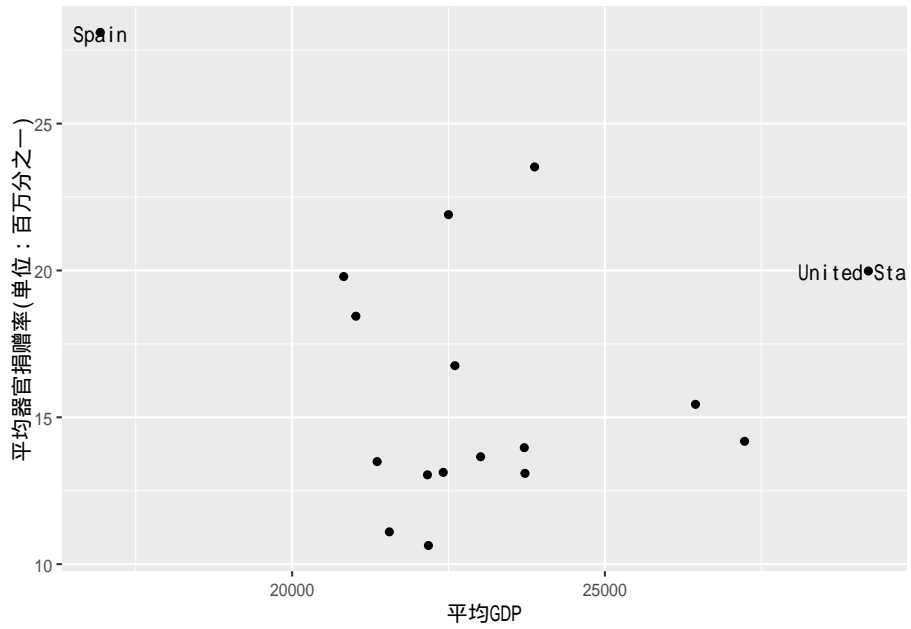
```
organdata4 <- organdata %>%
 group_by(country) %>%
 summarize(
 donors_mean = mean(donors, na.rm=TRUE),
 gdp_mean = mean(gdp, na.rm=TRUE)
)
```

```
p <- ggplot(data = organdata4,
 mapping = aes(x = gdp_mean, y = donors_mean))
p + geom_point() +
 labs(x = " 平均 GDP",
 y = " 平均器官捐赠率 (单位: 百万分之一)")
```



如果需要标出其中的特殊点，就要生成一个数据子集，并在 `geom_text()` 中指定输入数据为此子集：

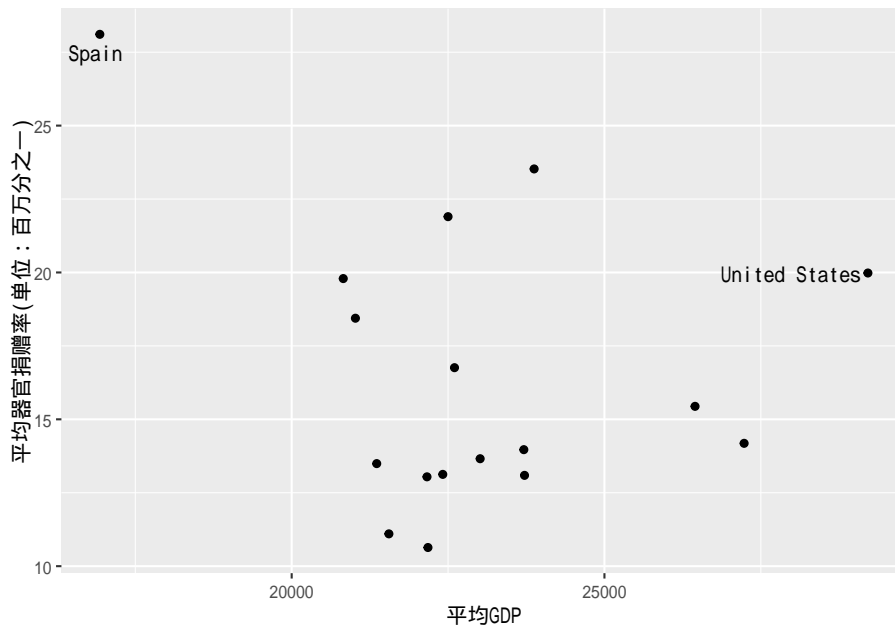
```
p + geom_point() +
 geom_text(data = subset(organdata4, gdp_mean > 27500 | donors_mean > 25),
 mapping = aes(label = country)) +
 labs(x = " 平均 GDP",
 y = " 平均器官捐赠率 (单位: 百万分之一)")
```



上面标的文字有超出边界的问题，ggrepel 包的 `geom_text_repel()` 可以改进：

```
library(ggrepel)
p + geom_point() +
 geom_text_repel(data = subset(organdata4, gdp_mean > 27500 | donors_mean > 25),
 mapping = aes(label = country)) +
 labs(x = "平均 GDP",
 y = "平均器官捐赠率 (单位：百万分之一)")
```





## 30.6 刻度 (scale)

在 `ggplot()` 的 `mapping` 参数中指定 `x` 维、`y` 维、`color` 维等，实际上每一维度都有一个对应的默认刻度 (scale)，即将数据值映射到图形中的映射方法。如果需要修改刻度对应的变换或者标度方法，可以调用相应的 `scale_xxx()` 函数。

可以映射的维度包括 `x`、`y`、`color`、`fill`、`shape`、`size` 等。其中 `x`、`y` 维度多用于表示连续变量，但是也可以用于分类变量，如 Cleveland 点图就是将 `y` 维度分配给一个分类变量。`color`、`fill` 可以用于连续变量，用于有序变量，也可以用于无序的分类变量。`shape` 只能用于无序的分类变量。

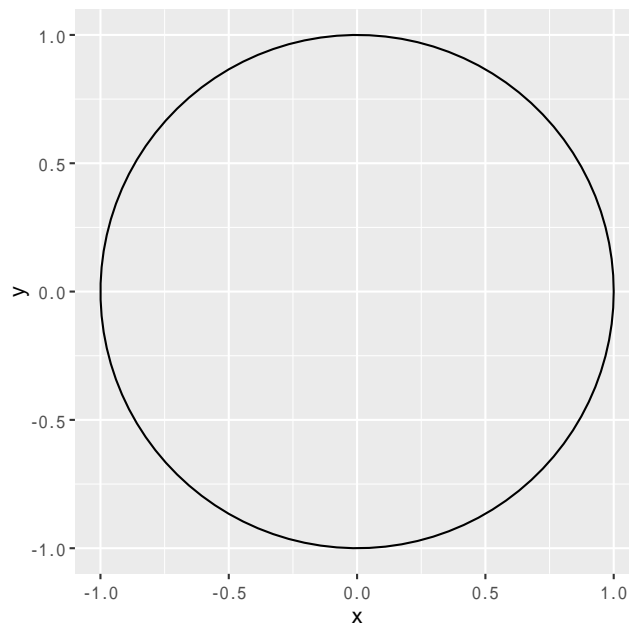
数值或者不同类别可以用平面上的不同位置表示，一般使用直角坐标系，有有一对相互垂直的 `x` 轴和 `y` 轴，但也可以有其他选择，比如 `y` 轴可以不与 `x` 轴垂直，某个轴的刻度可以不是线性的而是经过对数变换的，可以用极坐标系，等等。普通的直角坐标系不需要用 `scale_xxx()` 函数。

直角坐标系中的位置有 `x` 坐标与 `y` 坐标，可以分别代表两个变量，这两个变量经常是不同单位的，比如，身高与体重。这时，`x` 轴与 `y` 轴的数值之间没有可

比性，两个轴的数值范围与轴的实际长度只要不过于极端就没有什么关系。但是，应该尽可能使用比较合适的高宽比。

如果两个轴的变量含义相同，最好使用完全相同的坐标轴范围与坐标轴长度，使得水平与垂直方向上的等长距离在对应的坐标轴上也代表相等的距离。用 `coord_fixed()` 函数指定两个轴的单位为 1:1 长度的坐标系，其中参数 `xlim` 和 `ylim` 可以用来指定坐标范围。也可以用 `ratio` 参数指定一个其他的宽高比。如

```
d <- data.frame(
 t <- seq(0, 2*pi, length.out=100)
)
d$x <- cos(d$t)
d$y <- sin(d$t)
p <- ggplot(data = d, mapping = aes(
 x = x, y = y))
p + geom_path() +
 coord_fixed()
```



上面的程序中，`geom_path()` 按照输入数据集的次序将坐标点连接在一起。如

果改用 `geom_line()`，会将数据先按照 `x` 坐标值排序然后再顺序相连。可以看出，结果是一个正圆；如果不使用 `coord_fixed()`，结果不一定是正圆。

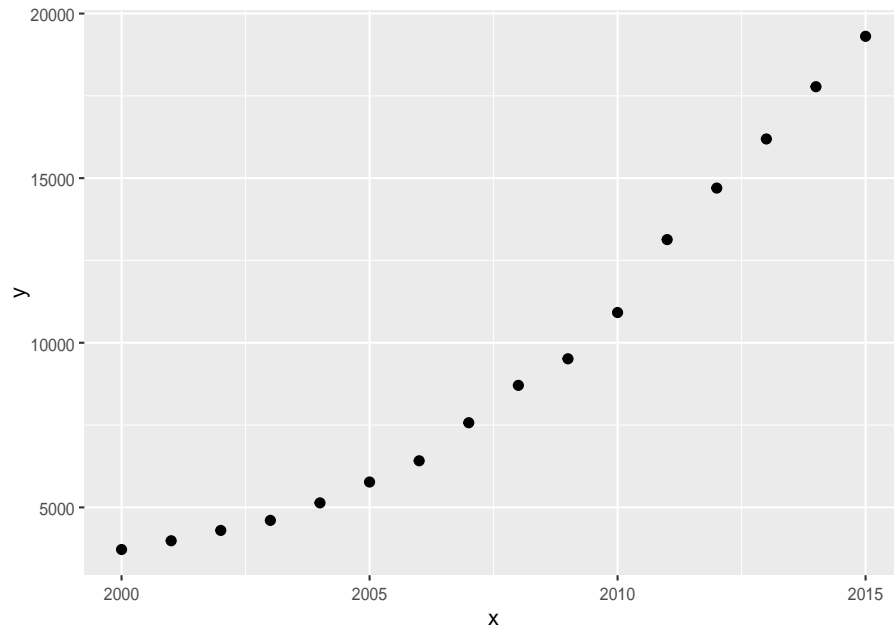
在 `coord_fixed()` 函数中，还可以用 `expand = FALSE` 要求坐标轴范围严格等于数据范围或 `xlim` 和 `ylim` 的规定，否则会比数据范围略宽一些。

`scale_x_log10()` 和 `scale_y_log10()` 可以将 `x` 轴或者 `y` 轴用对数刻度，这实际上是将数据做常用对数变换，但是相应的坐标轴刻度的数值还标成变换之前的原始值。对数轴方向相同的距离代表相差相同的倍数。经济、金融数据常常需要用对数刻度。

对数轴最好用来代表比例，尤其是使用条形图时，应该从 1 开始而不是从 0 开始。

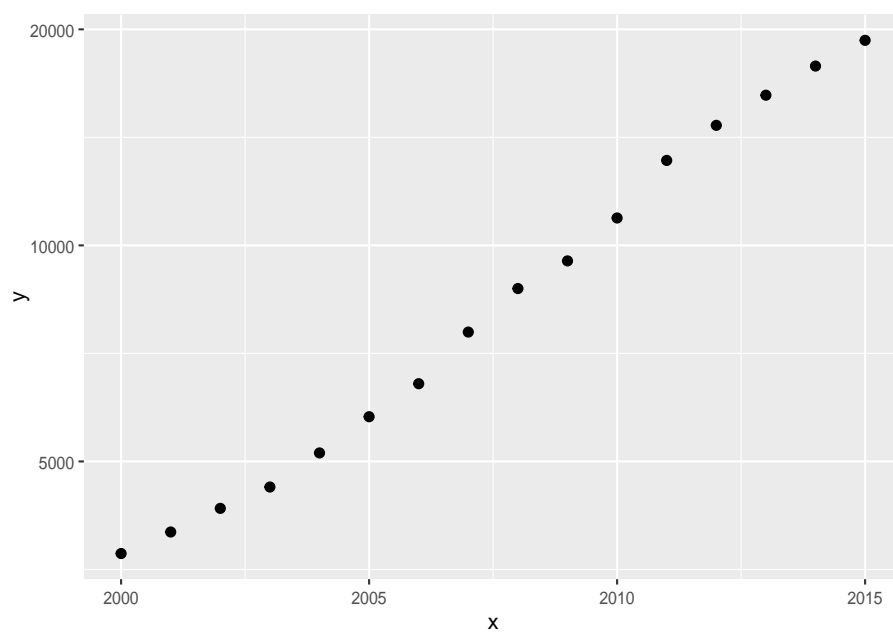
下图是我国 2000-2015 年的居民消费水平（元），可以看出，增长是指数型的：

```
d <- data.frame(
 x = 2000:2015,
 y = c(3721L, 3987L, 4301L, 4606L, 5138L, 5771L, 6416L, 7572L, 8707L,
 9514L, 10919L, 13134L, 14699L, 16190L, 17778L, 19308L))
p <- ggplot(data=d, mapping = aes(
 x = x, y = y))
p + geom_point(size = 2)
```



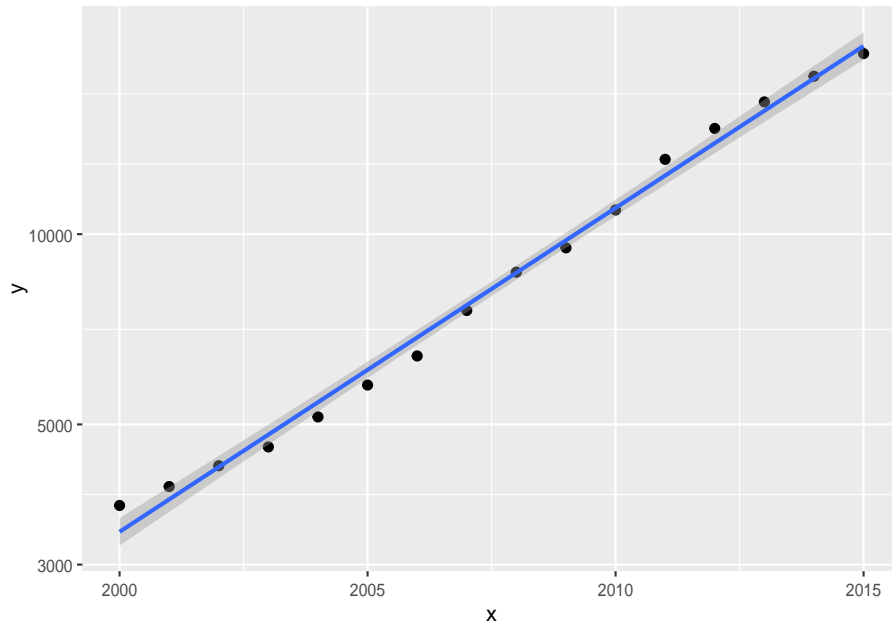
将 y 轴用对数刻度，则散点可以呈现为线性：

```
p + geom_point(size = 2) +
 scale_y_log10()
```



用 `geom_smooth()` 添加线性拟合线:

```
p + geom_point(size = 2) +
 geom_smooth(method="lm") +
 scale_y_log10()
```

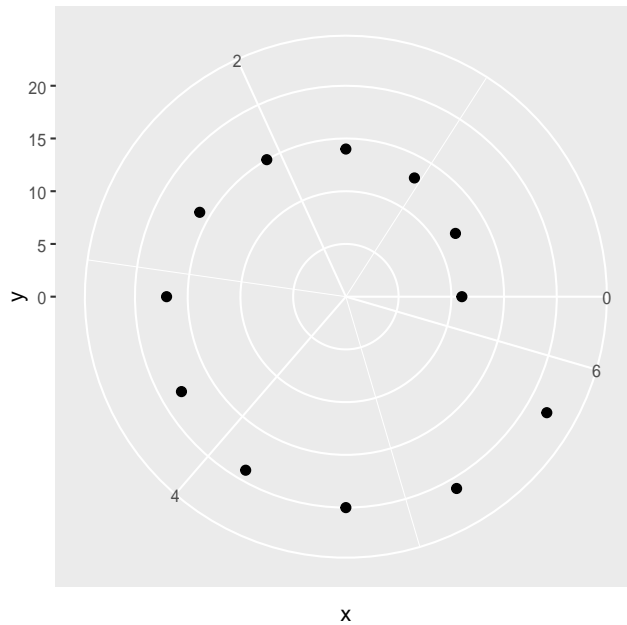


注意 `geom_smooth(method="lm")` 在计算拟合时会自动采用  $\log_{10}(y)$  的值。

`scale_x_sqrt()` 和 `scale_y_sqrt()` 与对数轴类似，做的是平方根变换。对于面积值，这种变换有一定意义。

用 `coord_polar()` 指定极坐标系，用 `theta="x"` 或者 `theta="y"` 指定那一维映射到极角。如：

```
d <- data.frame(
 x <- (0:11)/12*2*pi,
 y <- 11:22
)
p <- ggplot(data=d, mapping = aes(
 x = x, y = y))
p + geom_point(size=2) +
 scale_x_continuous(limits=c(0, 2*pi)) +
 scale_y_continuous(limits=c(0, 22)) +
 coord_polar(theta="x", start=-pi/2, direction=-1)
```



上图中  $y$  轴为极径的刻度，圆周上的数字表示极角刻度。程序中 `start` 指定极角为 0 的射线与 12 点方向的夹角弧度，`direction=1` 表示顺时针计算角度，`direction=-1` 表示逆时针。`scale_continuous_x()` 和 `scale_continuous_y()` 指定了两个维度的坐标范围。

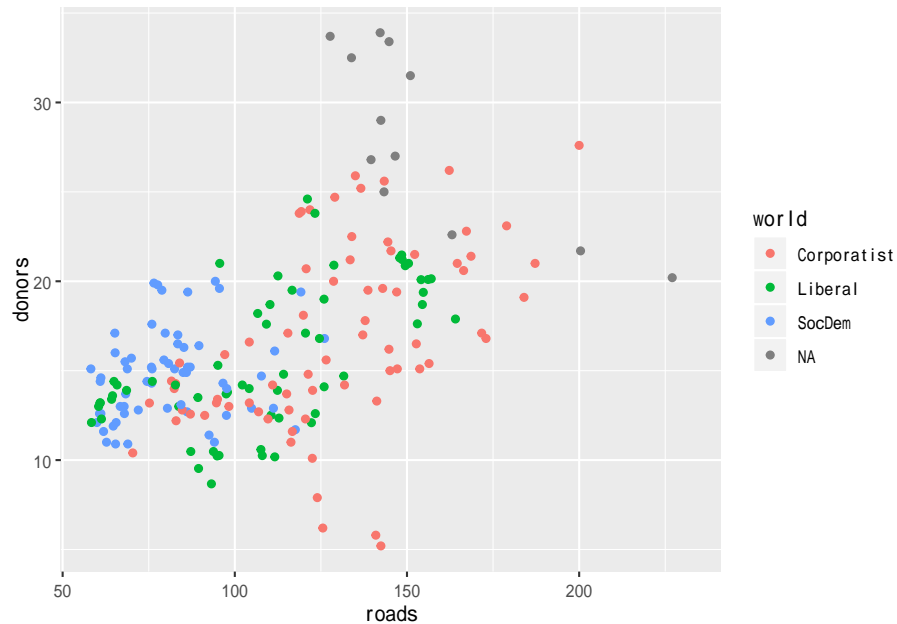
用 `scale_xxx()` 函数指定非默认的刻度，一般模式为 `scale_ 映射维度 _ 类型 ()`，“类型”包括 `continuous`、`discrete`、`log10` 等。可以用来人为指定坐标轴的刻度值，修改 `color` 或者 `fill` 维度所利用的颜色表，等等，用其中的参数指定这些内容。注意坐标轴的标签（标题）用 `labs()` 函数指定，而不是用 `scale_xxx()` 函数指定。

例如，在 `organdata` 中，作 `donors` 对 `roads`(每十万人的道路交通事故死亡率) 的散点图，并按 `world`(不同福利类型) 对散点染色，`x`、`y`、`color` 这三个维度都用了默认的刻度：

```
p <- ggplot(data = organdata,
 mapping = aes(
 x = roads,
 y = donors,
 color = world))
```

```
p + geom_point()
```

```
Warning: Removed 34 rows containing missing values (geom_point).
```



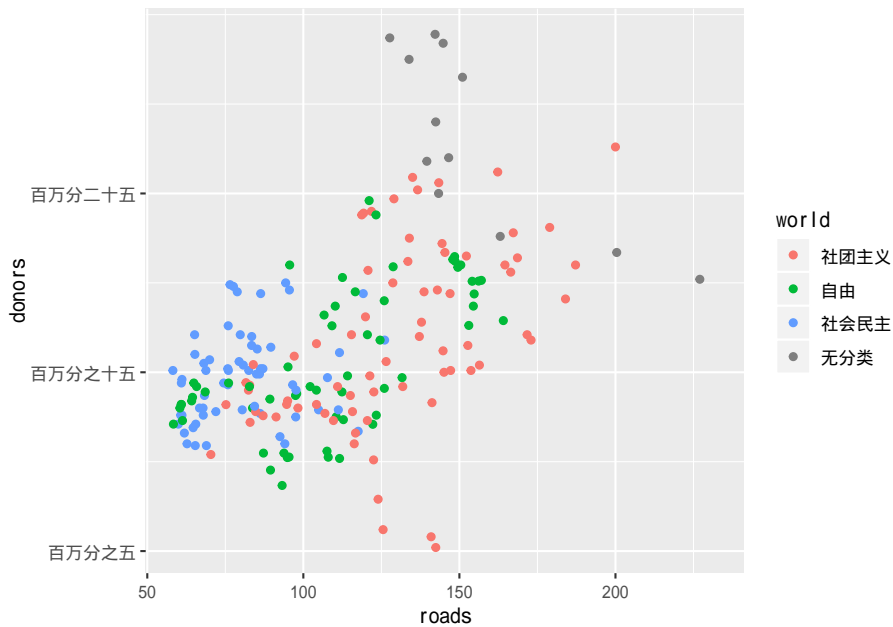
这里有三个维度：连续型的  $x$ 、 $y$  维度与无序分类值的  $color$  维度。 $color$  维度既可以取连续型，也可以取有序或者无序的分类型。 $x$ 、 $y$  维度有自动的坐标轴， $color$  维度有图例在图形右侧。如果没有特殊要求，不必调用 `scale_xxx()` 函数。

下面将  $y$  轴的刻度值进行人为的规定，将  $color$  维的标签人为指定：

```
p + geom_point() +
 scale_y_continuous(
 breaks = c(5, 15, 25),
 labels = c(" 百万分之五", " 百万分之十五", " 百万分二十五")) +
 scale_color_discrete(labels = c(
 " 社团主义", " 自由", " 社会民主", " 无分类"))
```

```
Warning: Removed 34 rows containing missing values (geom_point).
```





在 `scale_ 维度 _continuous()` 函数中, 可以用 `breaks` 指定坐标刻度标线的坐标数值表, 用 `minor_breaks` 指定细刻度的坐标数值表, 用 `labels` 指定在坐标刻度标线处标出的坐标值标签字符串表, 用 `limits` 指定坐标系的范围 (两个数的向量), 用 `expand` 指定在将坐标轴范围比数据的范围在两侧分别扩充多少, 取 `c(0,0)` 要求不扩充, 默认会左右各扩充 5%。

`ggplot2` 的散点图默认使用灰色背景并带有白色的网格线。`theme()` 函数的 `panel.background` 可以指定背景色, `panel.grid`、`panel.grid.xxx` 可以指定网格线做法。

## 30.7 如何使用颜色

可以用颜色来表示分组, 比如, 不同组的散点用不同颜色, 多条曲线用不同颜色; 可以用颜色表示数值, 用颜色深浅表示绝对值大小; 可以用颜色来突出某些要强调的图形元素。

将无序的分类值映射到颜色, 应该使用完全不相像、很容易区分的颜色, 各个颜色应该没有明显次序、没有哪一个与其他明显不同。可以用 `RColorBrewer` 扩展包提供的调色盘, 在 `ggplot2` 中用 `scale_color_brewer(palette)` 和

`scale_fill_brewer(palette)` 选择。图30.1为无序分类适用的调色板。称这样的调色板为名义型 (qualitative)。

色盲的人会分辨不出某些颜色，比如红绿色盲的人无法分辨红色和绿色，蓝绿色盲的人无法分辨蓝色和绿色。男性中有 8% 的人有色盲，所以绘图时应该考虑到这个问题。(Wilke, 2019) 提供了 8 种对色盲也可区分的颜色：

```
dcolorqua <- tribble(
 ~name, ~code,
 "black", "#000000",
 "orange", "#E69F00",
 "sky blue", "#56B4E9",
 "bluish green", "#009E73",
 "yellow", "#F0E442",
 "blue", "#0072B2",
 "vermilion", "#D55E00",
 "reddish purple", "#CC79A7",
 "black", "#000000"
)
knitr::kable(dcolorqua)
```

name	code
black	#000000
orange	#E69F00
sky blue	#56B4E9
bluish green	#009E73
yellow	#F0E442
blue	#0072B2
vermilion	#D55E00
reddish purple	#CC79A7
black	#000000

将连续的或者有序的值映射到颜色，应该使用渐变色，如果都是正值，可以使用从浅到深或者从深到浅的颜色，并且应该使用同一个颜色，只是深浅程度不同。图30.2为单向的有序分类适用的调色板。称这样的调色板为有序型 (sequential)。

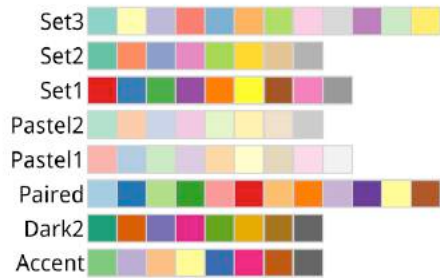


图 30.1: 无序分类适用的调色板

如果颜色代表有正有负的数值（比如反对、中立、支持），则应该以浅色为接近零值，正值与负值分别用两个不同的颜色。另外，如果要代表的数值有明显的中间值，为了强调较低的值与较高的值的对比，也可以使用这样的颜色刻度。图30.3为有正有负的有序分类适用的调色板。称这样的调色板为相异型 (diverging)。相异型的渐变色有可能对色盲的人不可辨识，Colorbrewer::PiYG 刻度对色盲人群也可以分辨。

例如，作 gapminder 数据集中各国在 2007 年期望寿命对人均 GDP 的散点图，不同大洲使用不同的颜色，指定 ColorBrewer 的 Set1 调色板：

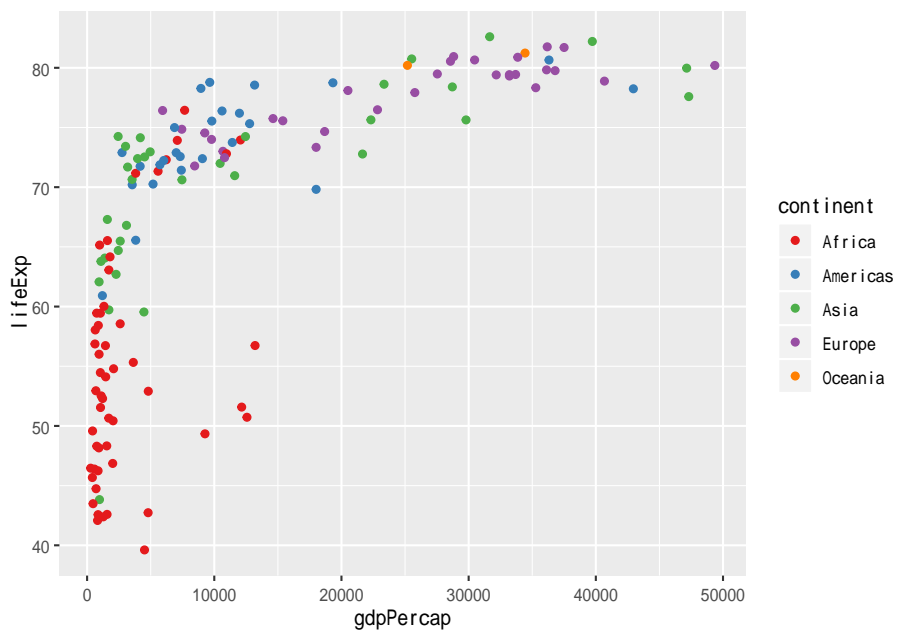
```
p <- ggplot(data = subset(gapminder, year == 2007),
 mapping = aes(
 x = gdpPerCap,
 y = lifeExp))
p + geom_point(mapping = aes(color = continent)) +
 scale_color_brewer(palette = "Set1")
```



图 30.2: 有序单向分类适用的调色板



图 30.3: 有序正负分类适用的调色板



使用颜色进行强调时，可以将要强调的点、线、条形用鲜明的颜色，而非强调的用褪色的颜色。最简单的做法是仅对要强调的内容指定一个鲜明的颜色。

## 30.8 标题、标注、指南、拼接

除了 `ggplot()` 指定数据与映射, `geom_xxx()` 作图, 还可以用许多辅助函数增强图形。

- `labs()` 可以设置适当的标题和标签。
- `annotate()` 函数可以直接在坐标系内进行文字、符号、线段、箭头、长方形的绘制。
- `guides()` 函数可以控制图例的取舍以及做法。
- `theme()` 函数可以控制一些整体的选项如背景色、字体类型、图例的摆放位置等。

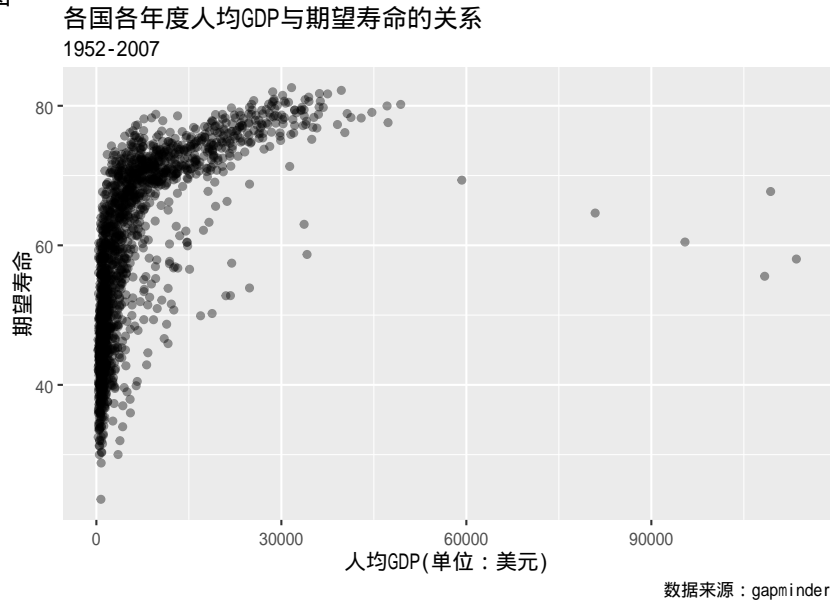
在需要修改图形时, 如果修改会影响到相应的 `geom_xxx()` 的主要结果, 一般需要在 `ggplot()` 或者该 `geom_xxx()` 函数中将适当的变量映射为某一维度, 或者用 `scale_xxx()` 函数进行变换。如果仅仅是一些显示效果的修改, 则一般作为 `geom_xxx()` 的选项, 或者调用 `labs()`、`theme()`、`guides()` 完成。

### 30.8.1 标题

函数 `labs()` 可以用来指定图形上方的标题 (title)、副标题 (subtitle)、右下方的标注 (caption)、左上方的标签以及坐标轴标题和其它维的名称。例如:

```
p <- ggplot(data = gapminder,
 mapping = aes(
 x = gdpPerCap,
 y = lifeExp))
p + geom_point(alpha = 0.4) +
 labs(
 title = " 各国各年度人均 GDP 与期望寿命的关系",
 subtitle = "1952-2007",
 tag = " 散点图",
 caption = " 数据来源: gapminder",
 x = " 人均 GDP(单位: 美元)",
 y = " 期望寿命"
)
```

## 散点图



在 `labs()` 中用 `x=`、`y=`、`color=` 之类的选项指定坐标轴的标签或者其它维的名称, 如果该维用了 `scale_xxx()` 函数, 则应该则 `scale_xxx()` 函数中用 `name=` 指定轴标签 (名字)。数值型的维度除非是显然的应在标签中包含单位。`labs()` 中或者 `scale_xxx(name=)` 中如果指定文字为 `NULL` 则取消该标题或标签, 但是取消坐标轴标签必须是取消标签后含义也显而易见, 没有任何疑问的情况。

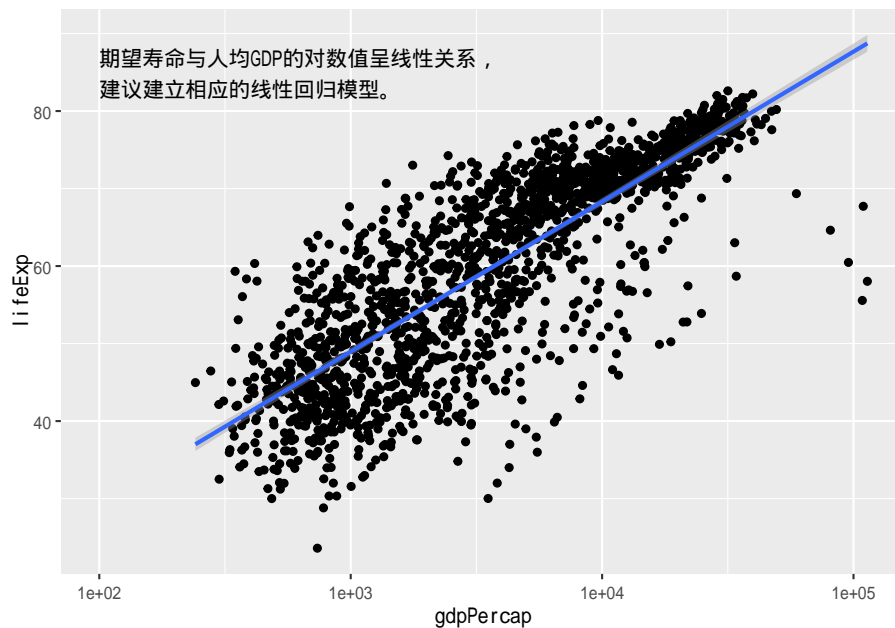
`labs()` 只是提供了这些标题功能, 一般并不会同时使用这些功能。在出版图书内, 图形下方一般伴随有图形说明, 这时一般就不再使用标题、副标题、标签、标注, 而只需写在图的伴随说明文字中, 当然, 坐标轴标签一般还是需要的。

有一点要注意, 默认的标题、坐标轴标签、图例标签中的文字往往偏小, 如果单独放大这些图形来看并没有问题, 但是作为插图放在书中或者网页中就偏小了。可以用 `theme()` 函数调整字体大小。

### 30.8.2 标注功能

通过 `annotate(geom = "text")` 调用 `geom_text()` 的功能, 可以在一个散点图中标注多行文字, 多行之间用 `"\n"` 分开:

```
p <- ggplot(data = gapminder,
 mapping = aes(
 x = gdpPercap,
 y = lifeExp))
p + geom_point() +
 geom_smooth(method="gam") +
 scale_x_log10() +
 annotate(
 geom = "text",
 x = 1E2, y = 85, hjust = 0,
 label=" 期望寿命与人均 GDP 的对数值呈线性关系, \n建议建立相应的线性回归模型。")
```



将上面图形中最右侧偏低的点用长方形填充标出:

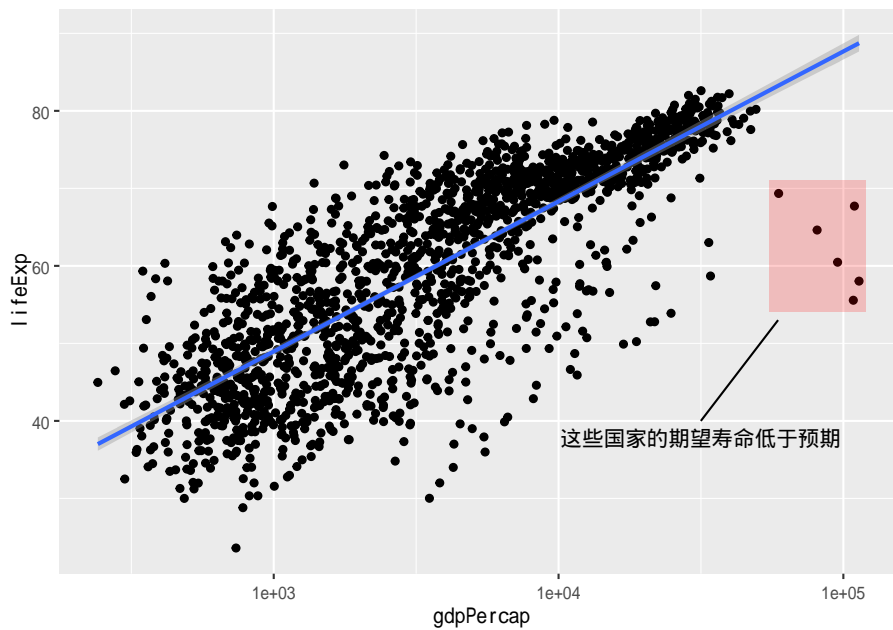
```
p + geom_point() +
 geom_smooth(method="gam") +
 scale_x_log10() +
 annotate(geom = "rect",
 xmin = 5.5E4, xmax = 1.2E5,
```



```

 ymin = 54, ymax = 71,
 fill = "red", alpha = 0.2) +
 annotate(geom = "line",
 x = c(5.9E4, 3.16E4),
 y = c(53, 40)) +
 annotate(geom = "text",
 x = 3.16E4, y = 38,
 label = " 这些国家的期望寿命低于预期")

```



这样标注的缺点是坐标都需要读图并试错摆放。

上述被标注的国家是：

```

gapminder %>%
 filter(gdpPerCap > 5.5E4 & gdpPerCap < 1.2E5,
 lifeExp > 54 & lifeExp < 71) %>%
 select(country, gdpPerCap, lifeExp)

```

```
A tibble: 6 x 3
```

```
country gdpPerCap lifeExp
```

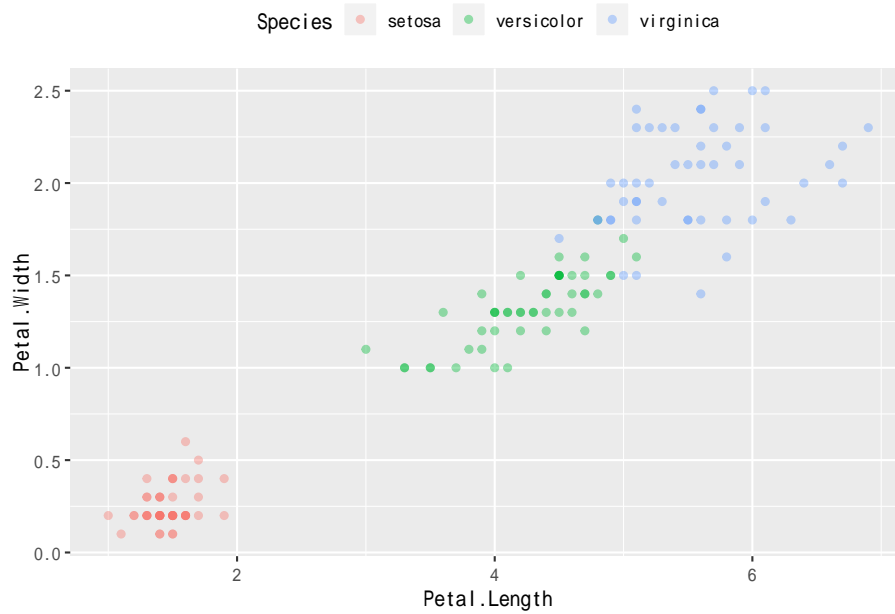
```
<fct> <dbl> <dbl>
1 Kuwait 108382. 55.6
2 Kuwait 113523. 58.0
3 Kuwait 95458. 60.5
4 Kuwait 80895. 64.6
5 Kuwait 109348. 67.7
6 Kuwait 59265. 69.3
```

### 30.8.3 指南 (guides)

对于颜色、填充色等维度，会自动生成图例。用 `guides(color = FALSE)` 这样的方法可以取消指定维度的图例。

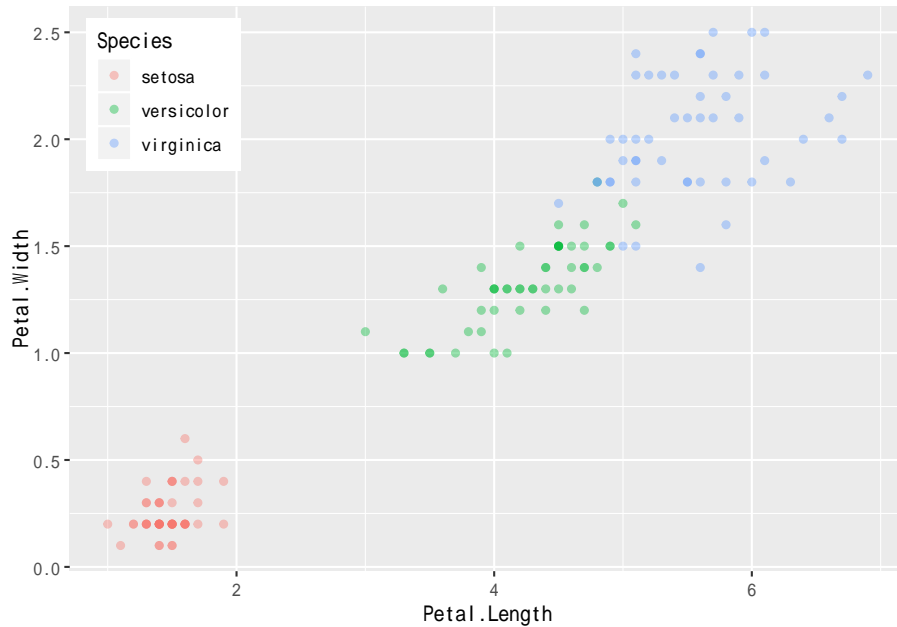
`theme()` 可以调整一些整体的设置，如背景色、字体、图例的摆放位置。用 `theme()` 的 `legend.position` 改变图例的位置，如 `theme(legend.position = "top")` 可以将图例放置在上方，默认是放置在右侧的。可取值有 "none"、"left"、"right"、"bottom"、"top"，如：

```
p <- ggplot(data = iris, mapping = aes(
 x = Petal.Length, y = Petal.Width, color = Species))
p + geom_point(alpha = 0.4) +
 theme(legend.position = "top")
```



图例位置 `legend.position` 还可以指定在作图区域内部用两个百分比数字给定，同时可以用 `legend.just` 指定位置对准图例框的哪一个角，也是用两个百分比数字，比如放在左上角内部：

```
p + geom_point(alpha = 0.4) +
 theme(
 legend.position = c(0.02, 0.98),
 legend.justification = c(0,1))
```



某一维的图例做法可以在相应的 `scale_xxx()` 函数中用 `guide=` 指定, 或者在 `guides()` 中用 `维名称 = 指南方法` 指定, 如 `guides(fill = "guides")` 或者 `guides(fill = "colorbar")`。指定的指南还可以是 `guide_legend()` 或者 `guide_colorbar()` 的结果, 在这两个函数中可以用多个选项指定指南或颜色条的具体做法。

### 30.8.4 拼接图形

`facet_wrap()` 和 `facet_grid()` 可以按照某一个或两个分类变量的值将输入数据集分为若干个子集, 将每个子集分别在一个小图上绘图。有时还需要将几幅不同图形拼在一起, 这些图形可以是同一数据的不同类型图形, 也可以是完全无关的图形。

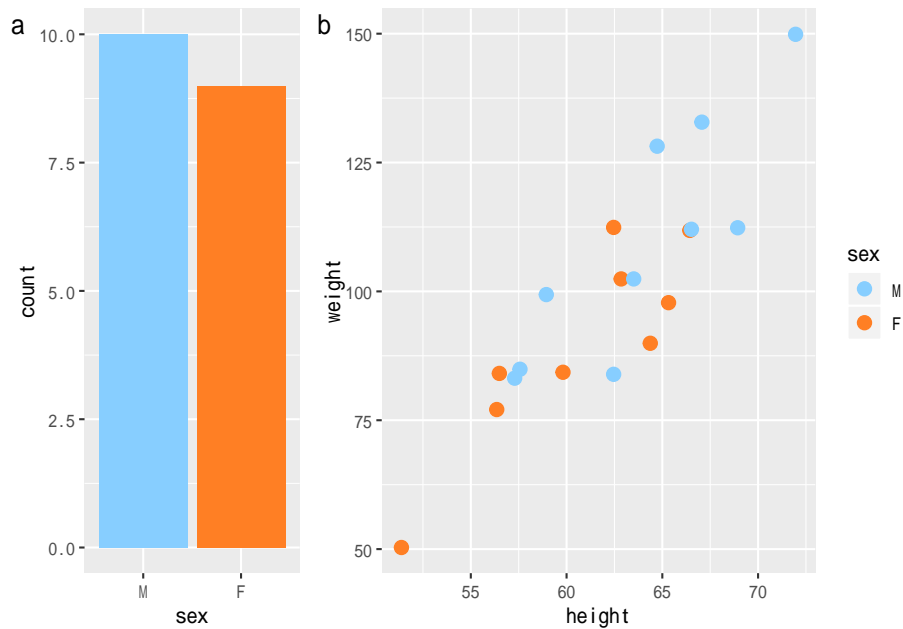
在拼接图形时, 各个小图应该具有类似的风格, 即背景、配色、字体等应该尽可能一致, 上下和左右的坐标轴应尽可能对齐, 小图的标签应该尽可能不显眼。

`cowplot` 包提供了拼接图形的办法, 使用时先将每个小图分别赋值给一个 R 变量, 然后用 `plot_grid()` 函数摆放在一幅图中。

例如, 19 个学生的性别、年龄、身高、体重数据:

```
library(cowplot)
dclass <- read_csv(
 "class.csv",
 col_types=cols(
 .default = col_double(),
 name=col_character(),
 sex=col_factor(levels=c("M", "F"))
))
p1 <- ggplot(data = dclass, mapping = aes(x = sex, fill = sex)) +
 geom_bar() +
 scale_fill_manual(
 guide = FALSE,
 values = c("F" = "chocolate1",
 "M" = "skyblue1"))
p2 <- ggplot(data = dclass, mapping = aes(
 x = height, y = weight, color = sex)) +
 geom_jitter(size = 3) +
 scale_color_manual(
 values = c("F" = "chocolate1",
 "M" = "skyblue1"))

plot_grid(
 p1, p2, labels = "auto",
 rel_widths = c(0.3, 0.6),
 align = "h")
```



`plot_grid()` 中用 `rel_widths` 指定了左右图的相对比例，默认是均分的。当图形上下排列时，可以用 `rel_height` 指定上下排的比例。程序中用了 `align = "h"` 使得左右图上下对齐，`align` 默认取 "none"，也可以取 "h"、"v" 或 "hv"。可以用 `ncol` 指定图的列数，用 `nrow` 指定图的行数。

还可以通过嵌套调用 `plot_grid()` 方式实现非网格的排列。

## 30.9 图形定制调整

`ggplot2` 的默认设置一般能够满足我们的要求，只有在有特殊的图形类型需求或者制作出版用的图形时，才需要对图形进行定制调整。图形调整可以包括：

- 配色、位置摆放等审美方面的调整；
- 面向目标出版物或者目标读者的调整；
- 增加有意义的标注；
- 调整整体的观感。

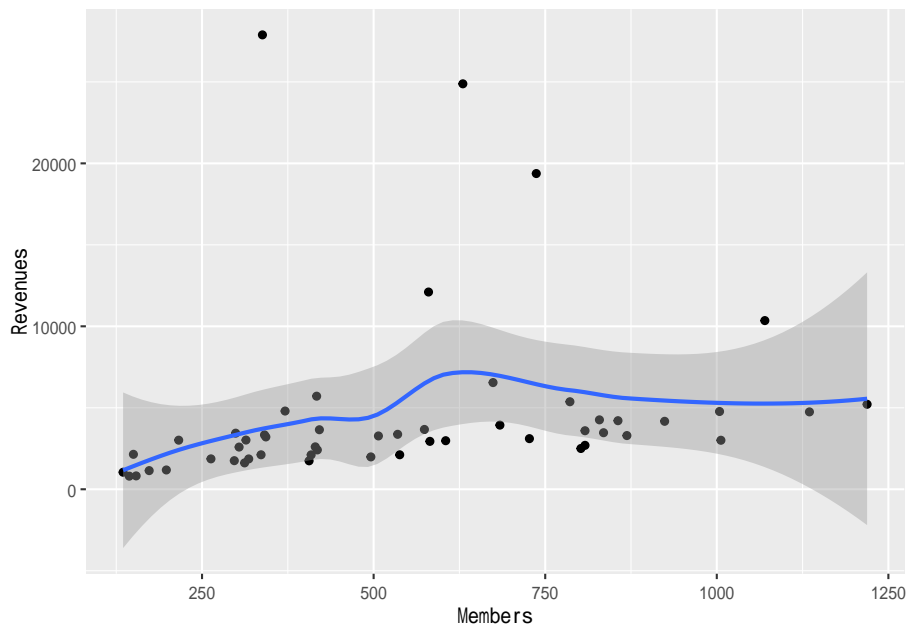
### 30.9.1 图形逐步调整例子

socviz 扩展包的 `asasec` 数据集是美国社会学学会 (ASA) 的各分会 2005 年到 2015 年的一些数据, 其中的财务数据 (`Beginning`, `Revenues`, `Expenses`, `Ending`) 虽然各年都有值, 但实际是用 2015 年的值填进去的。

作 2014 年收入对会员数的散点图与拟合曲线, 每个散点是一个分会:

```
p <- ggplot(data = subset(asasec, Year == 2014),
 mapping = aes(
 x = Members,
 y = Revenues,
 label = Sname))
p + geom_point() +
 geom_smooth()
```

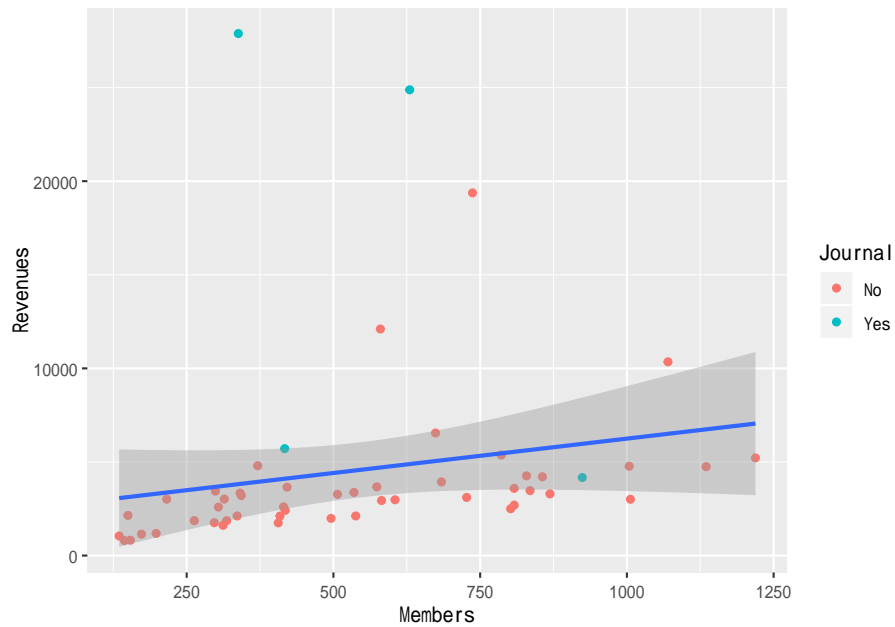
```
`geom_smooth()` using method = 'loess' and formula 'y ~ x'
```



程序中的映射 `label = Sname` 暂时不起作用, 在 `geom_text()` 中才需要这一维度。

下面按照有无期刊染色，将平滑方法改为线性回归：

```
p + geom_point(mapping = aes(color = Journal)) +
 geom_smooth(method = "lm")
```

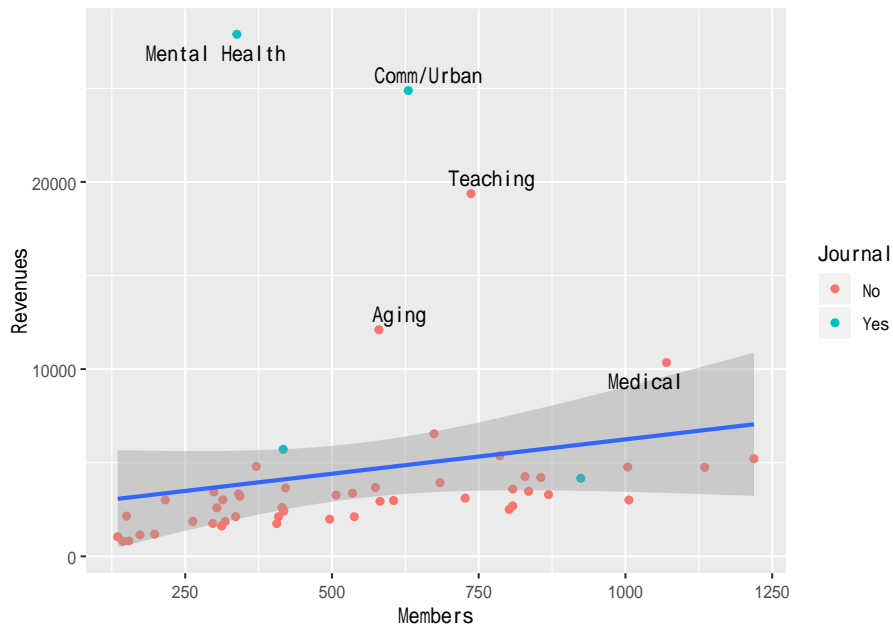


注意 `color = Journal` 的映射仅对 `geom_point()` 有效，如果写在 `ggplot()` 中，就对 `geom_smooth()` 也有效了。

下面标出异常值：

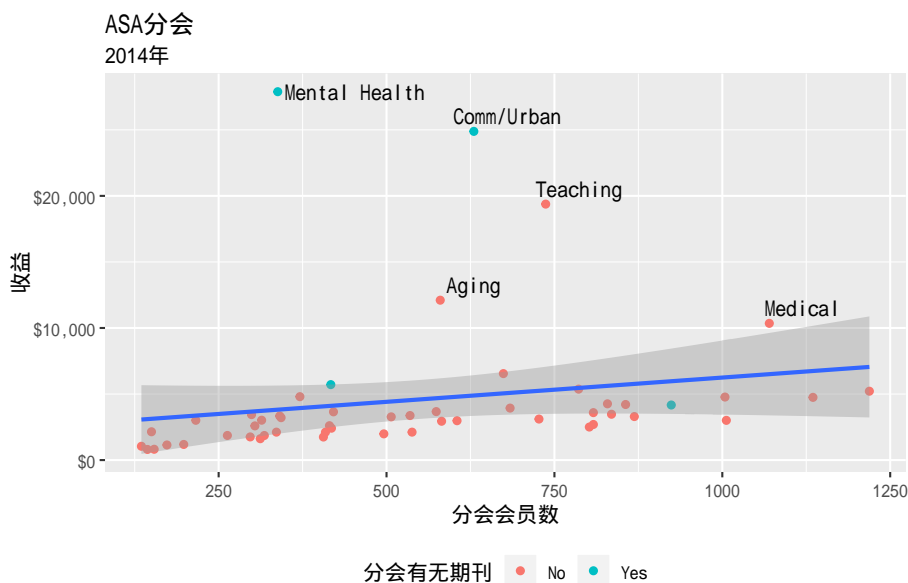
```
p + geom_point(mapping = aes(color = Journal)) +
 geom_smooth(method = "lm") +
 geom_text_repel(data = subset(asasec, Year == 2014 & Revenues > 7000))
```





下面用 `labs()` 调整标题、用 `scale_y_continuous()` 调整 y 轴刻度标法、用 `theme()` 调整图例位置:

```
p + geom_point(mapping = aes(color = Journal)) +
 geom_smooth(method = "lm") +
 geom_text_repel(data = subset(asasec, Year == 2014 & Revenues > 7000)) +
 labs(
 title = "ASA 分会",
 subtitle = "2014 年",
 x = "分会会员数",
 y = "收益",
 color = "分会是否有期刊",
 caption = "数据来源: ASA 年报") +
 scale_y_continuous(labels = scales::dollar) +
 theme(legend.position = "bottom")
```



数据来源：ASA年报

## 30.10 主题

ggplot2 包作图可以实现内容与设计的分离，这里内容就是指数据、映射、统计、图形类型等方面，而设计就是指背景色、颜色表、字体、坐标轴做法、图例位置等的安排。将作图任务分解为内容与设计两个方面，可以让数据科学家不必关心设计有关的元素，而设计可以让专门的艺术设计人才来处理。这种工作分配已经在图书出版、网站、游戏开发等行业发挥了重要作用。

`theme()` 函数用来指定设计元素，称为主题 (theme)，而且可以单独开发 R 扩展包来提供适当的主题。

`theme(legend.position)` 可以用来选择图例位置。`theme_set()` 可以改变后续 ggplot2 作图的主题 (配色、字体等)。如 `theme_set(theme_bw())`、`theme_set(theme_dark())` 等。对单次绘图，可以直接用加号连接 `theme_gray()` 等这些主题函数。主题包括 `theme_gray()` (默认主题)、`theme_minimal()`、`theme_classic()` 等。ggthemes 扩展包提供了更多的主题选择。

`theme()` 函数还可以直接指定颜色、字体、大小等设置。

## 30.11 参考文献

- Cleveland, W. S. (1993). *The elements of graphing data*. Hobart Press.
- Cleveland, W. S. (1994). *Visualizing data*. Hobart Press.
- Cleveland, W. S., & McGill, R. (1984). Graphical perception: Theory, experimentation, and application to the development of graphical methods. *Journal of the American Statistical Association*, 79, 531–534.
- Cleveland, W. S., & McGill, R. (1987). Graphical perception: The visual decoding of quantitative information on graphical displays of data. *Journal of the Royal Statistical Society Series A*, 150, 192–229.
- Edward R. Tufte((1983) *The Visual Display of Quantitative Information*
- Tufte, E. R. (1983). *The visual display of quantitative information*. Cheshire, CT: Graphics Press.
- Tufte, E. R. (1990). *Envisioning information*. Cheshire, CT: Graphics Press.
- Tufte, E. R. (1997). *Visual explanations: Images and quantities, evidence and narrative*. Cheshire, CT: Graphics Press.
- Wickham, H. (2016). *Ggplot2: Elegant graphics for data analysis*. New York: Springer.
- Wickham, H., & Chang, W. (2018). *Ggplot2: Create elegant data visualisations using the grammar of graphics*.
- Wilkinson, L. (2005). *The grammar of graphics (Second)*. New York: Springer.



# Chapter 31

## ggplot 的各种图形

### 31.1 介绍

ggplot2 包提供了许多种图形，其作用可以大致地分为：

- 表现数量；
- 表现一维或者二维分布；
- 表现两个变量之间的数量关系，等等。

下面按照其作用分别进行介绍。

主要参考：

- Claus O. Wilke(2019). Fundamentals of Data Visualization. O'Reilly Media. <https://serialmentor.com/dataviz/>

### 31.2 表现数量

#### 31.2.1 条形图

设有若干个类，每个类有一个数量属性值。经常用条形图表现数量。

### 31.2.1.1 简单的条形图

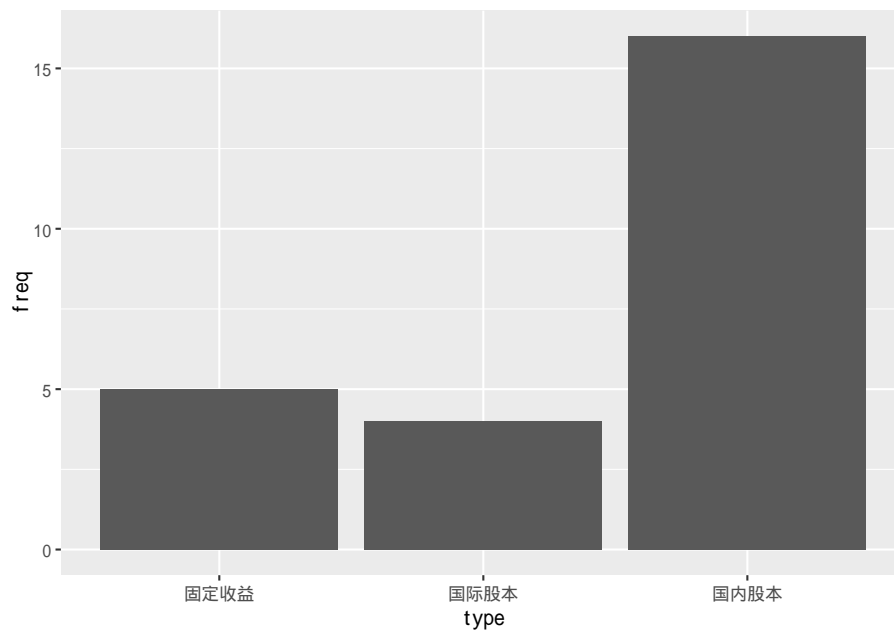
例如，有 25 个共同基金，分为三个类别，各类别的频数为：

```
d_funds_type <- tibble(
 type = c(" 国内股本", " 国际股本", " 固定收益"),
 freq = c(16, 4, 5)
)
knitr::kable(d_funds_type)
```

type	freq
国内股本	16
国际股本	4
固定收益	5

用 `ggplot()` 输入数据并指定映射，用 `geom_col()` 作条形图，x 轴为类别，y 轴为每个类别对应的数值，这里是输入数据中的频数：

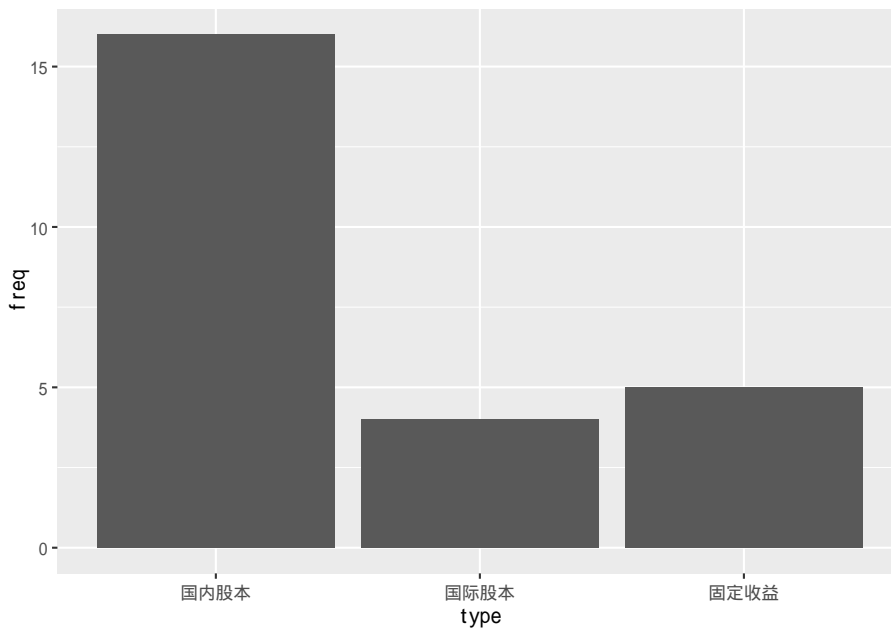
```
p <- ggplot(data = d_funds_type, mapping = aes(
 x = type, y = freq))
p + geom_col()
```



上面的最后一行程序也可以用 `geom_bar()` 函数写成 `p + geom_bar(stat = "identity")`。`geom_bar()` 可以对分类变量的原始值自动统计频数然后做频数条形图，如果输入数据已经是频数，就需要加 `stat = "identity"` 选项。

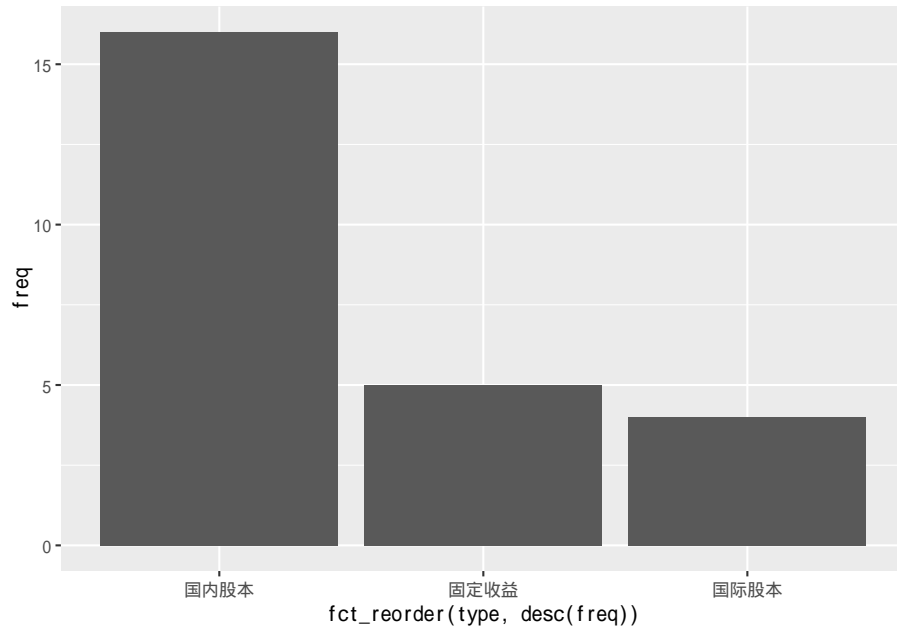
上面的三个类的次序并不是输入的次序，这是因为字符型变量会自动转换为因子型，因子水平一般按字典序排列。可以在输入时用 `factor(x, levels=<指定的因子水平表 >)` 定义因子并人为指定因子次序，或者用 `forcats` 包的 `fct_relevel()` 函数，如：

```
d <- d_funds_type %>%
 mutate(type = fct_relevel(type, " 国内股本", " 国际股本", " 固定收益"))
p <- ggplot(data = d, mapping = aes(
 x = type, y = freq))
p + geom_col()
```



可以用 `forcats` 包的 `fct_reorder()` 函数，将各个水平按照频数大小排序，如：

```
p <- ggplot(data = d, mapping = aes(
 x = fct_reorder(type, desc(freq)),
 y = freq))
p + geom_col()
```

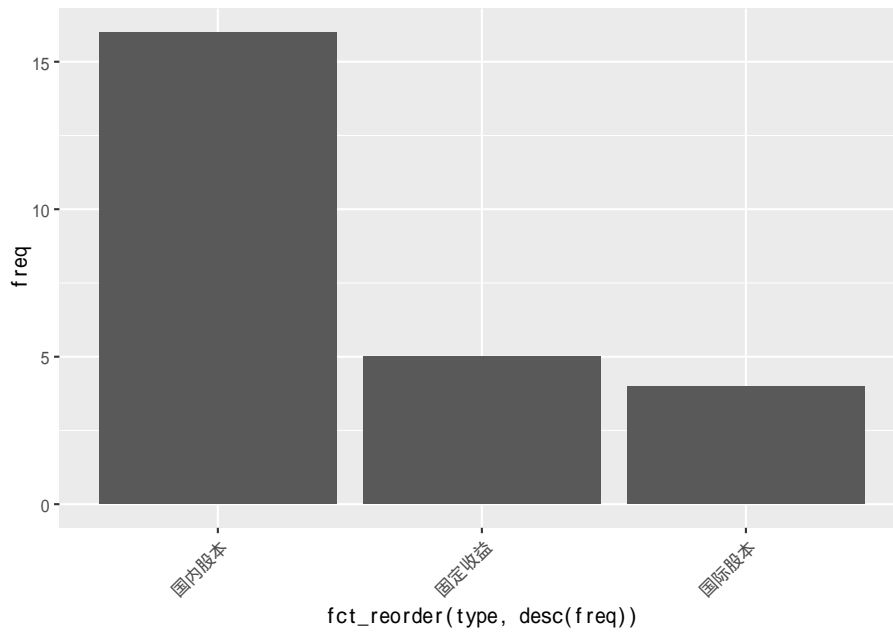


注意，条形图中这种对各条形按照高度重排序的做法仅适用于各类没有自然的次序的情形，如果各类有自然的次序，比如年份、月份、季度、年龄组、收入高低分组等，就应该按照各类的自然次序，可以在输入数据时用 `factor(x, levels=< 指定的因子水平表 >)` 定义因子并人为指定因子次序，或者在输入数据后用 `forcats` 包的 `fct_relevel()` 函数重排次序。

当分类的标签值（如“固定收益”）较长而且类个数较多时，横轴的宽度可能不够用，使得相邻两项的标签重叠，只好省略部分标签。这时，一种办法是将类标签方向变成 45 度角：

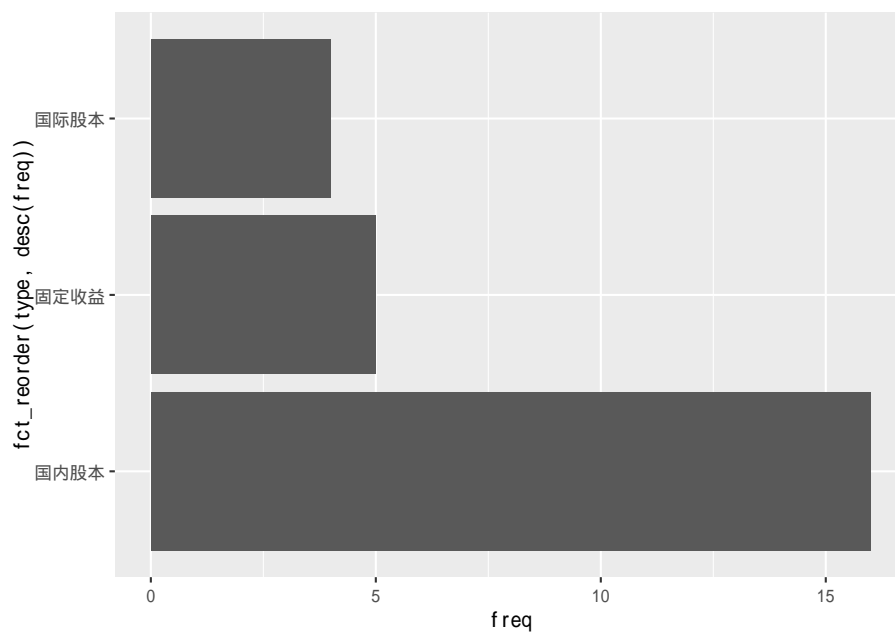
```
p + geom_col() +
 theme(
 axis.text.x = element_text(angle = 45, vjust = 1, hjust = 1)
)
```





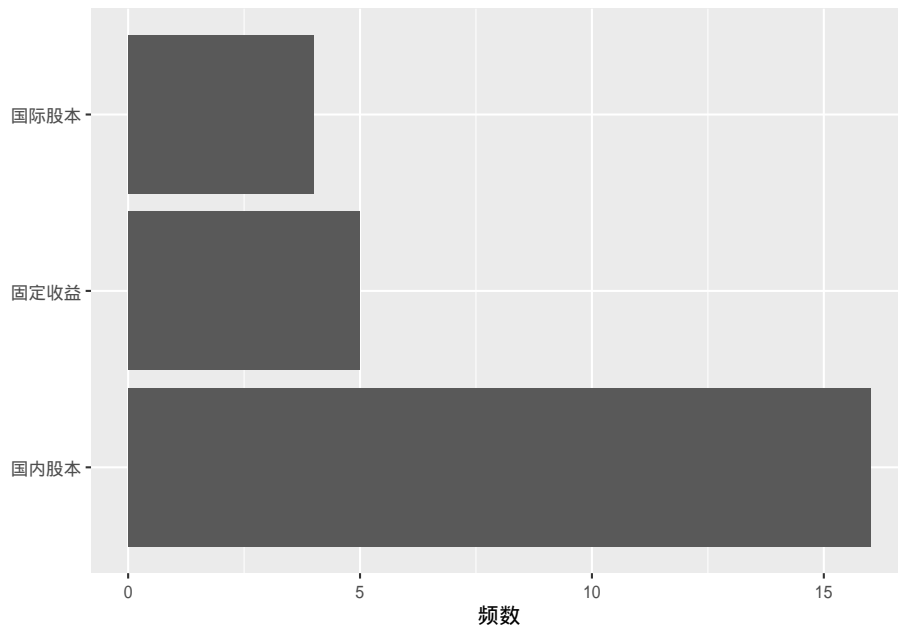
针对条形图的各条形的标签太长的的问题，更好的办法是将条形横向放置，这样类标签就到了 y 轴，可以横向显示。这可以在映射中直接修改，但最好还是用 `coord_flip()` 函数颠倒 x 轴与 y 轴的作用，因为有些图形是不允许修改映射的。程序如

```
p + geom_col() +
 coord_flip()
```



下面用 `labs()` 去掉分类标签轴的标题并修改频数轴的标题:

```
p + geom_col() +
 coord_flip() +
 labs(x = NULL,
 y = " 频数")
```



如果要将频数替换成比例或者百分数，只要修改输入数据增加相应的变量即可。

### 31.2.1.2 并列和堆叠的条形图

在某次对某物业公司服务的 4 个社区的居民的抽样调查中抽查了 80 位居民，其性别与所属社区如下：

```
mat_comm <- matrix(
 c(9, 10, 13, 4,
 18, 7, 8, 11),
 nrow=4, ncol=2,
 dimnames = list(c("A 社区", "B 社区", "C 社区", "D 社区"), c("男", "女"))
)
knitr::kable(mat_comm)
```

	男	女
A 社区	9	18
B 社区	10	7
C 社区	13	8
D 社区	4	11

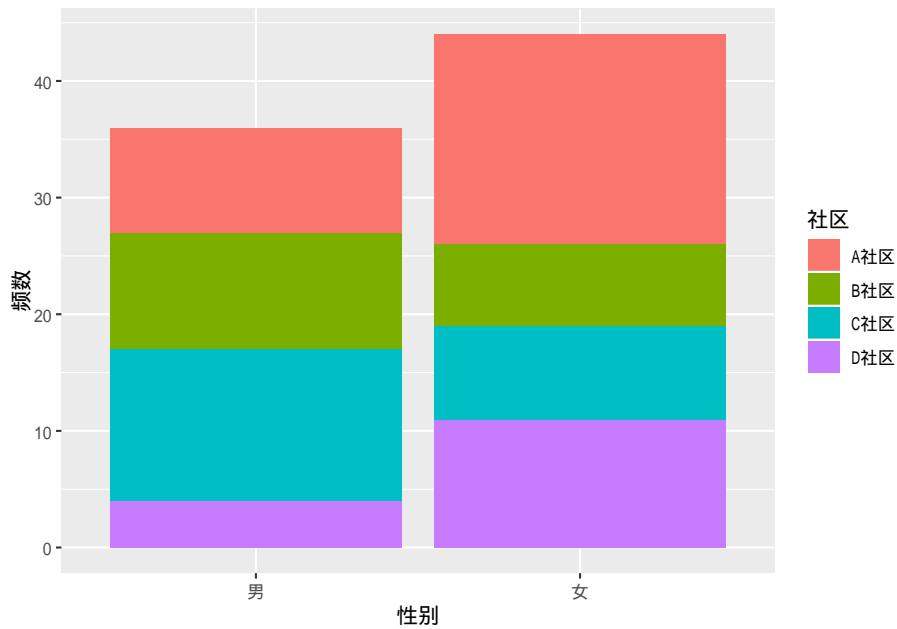
将表格数据转换为数据框:

```
d_comm <- as.data.frame(as.table(mat_comm))
names(d_comm) <- c("社区", "性别", "频数")
knitr::kable(d_comm)
```

社区	性别	频数
A 社区	男	9
B 社区	男	10
C 社区	男	13
D 社区	男	4
A 社区	女	18
B 社区	女	7
C 社区	女	8
D 社区	女	11

现在数据中有两个分类变量，频数是交叉分类的结果，这样的数据称为“列联表”。这时，将两个分类变量中一个作为大类，映射到 x 维度，另一个作为小类，映射到 fill 维度，频数映射到 y 维度。使用 `geom_col()` 或者 `geom_bar(stat = "identity")`。如:

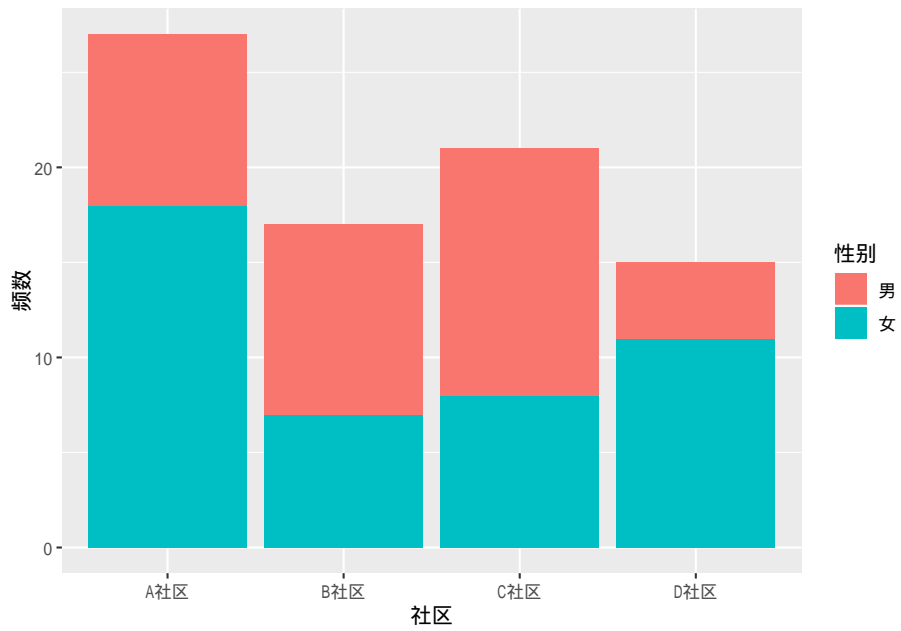
```
p <- ggplot(data = d_comm, mapping = aes(
 x = `性别`, fill = `社区`, y = `频数`
))
p + geom_col()
```



结果为堆叠的条形图，每个大类画一个分段的条形，条形的程度等于各段长度之和。这需要各段长度之和相当于每个大类的某种总计，对于频数，这当然是合理的，每段长度每个交叉小类的频数，而一个条形的总长度是一个大类的频数。当 y 轴映射到其它指标时堆叠条形图不一定合适，比如，如果 y 映射到每个交叉小类的某个指标的标准差，因为标准差的和并不是标准差，所以就不能用堆叠条形图表示。

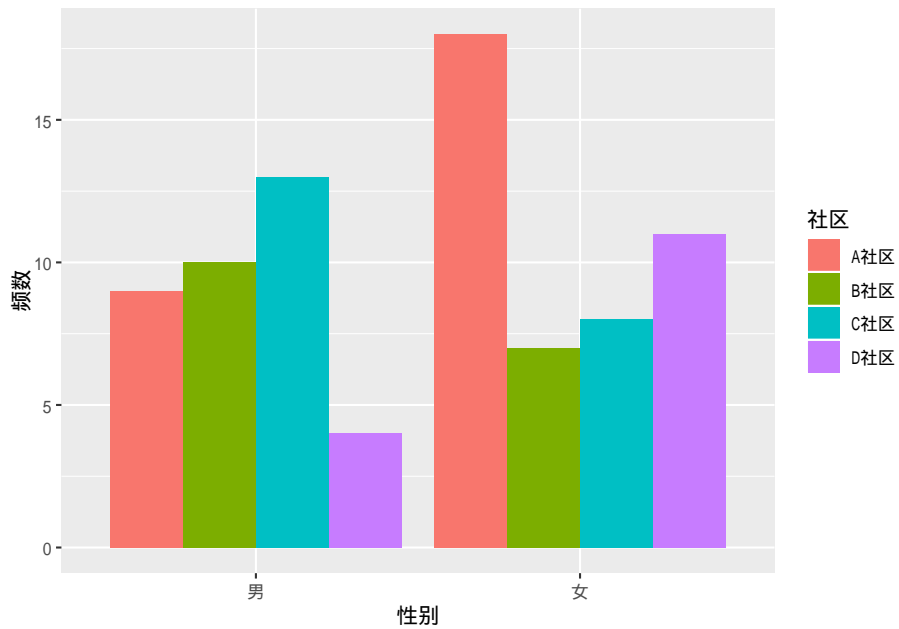
下面将社区作为大类，性别作为小类：

```
p <- ggplot(data = d_comm, mapping = aes(
 x = `社区`, fill = `性别`, y = `频数`
))
p + geom_col()
```



在 `geom_col()` 和 `geom_bar()` 中加选项 `position = "dodge"` 可以将小类并列放置:

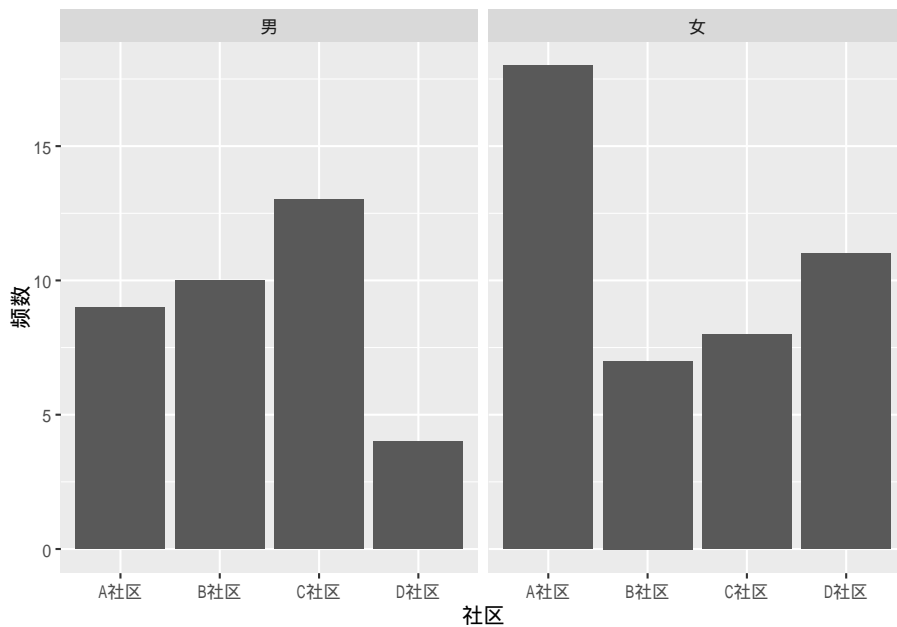
```
p <- ggplot(data = d_comm, mapping = aes(
 x = `性别`, fill = `社区`, y = `频数`
))
p + geom_col(position = "dodge")
```



这样得到并列的条形图。并列条形图就不需要像堆叠条形图那样要求各个小类的 y 轴变量和有意义。堆叠条形图相当于在 `geom_col()` 中加 `position = "stack"`。

也可以将不同的大类放置到不同的小图,用 `facet_wrap()` 或者 `facet_grid()` 函数:

```
p <- ggplot(data = d_comm, mapping = aes(
 x = `社区`, y = `频数`
))
p + geom_col() +
 facet_wrap(~ `性别`)
```



注意各小图默认使用相同的坐标范围。

### 31.2.1.3 百分数的条形图

对于只有一个分类变量的条形图，如果 y 轴代表频数，为了转换成比例或者百分比，修改输入数据增加像相应的变量，如：

```
d_funds_type %>%
 mutate(ratio = freq / sum(freq),
 pct = 100*ratio) -> d
```

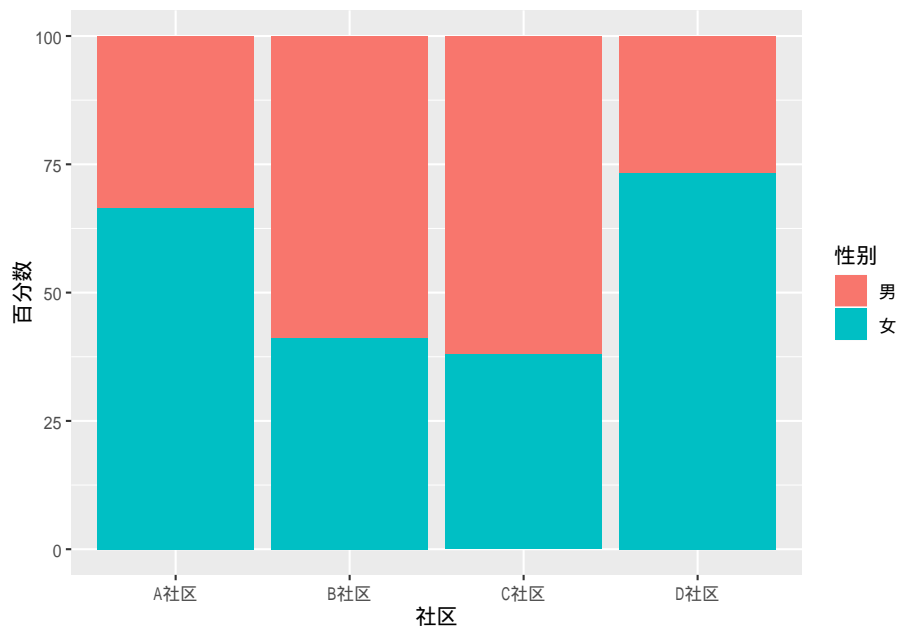
对于有两个分类变量的列联表，如果希望 y 轴表现总百分数，只要修改输入数据集增加相应的变量。如果 y 轴希望表现每个大类内的百分数，首先增加大类内的百分数，比如下面的程序将社区居民调查数据按社区分为大类，每个社区内计算不同性别的百分数：

```
d_comm %>%
 group_by(`社区`) %>%
 mutate(`百分数` = `频数` / sum(`频数`) * 100) %>%
 ungroup() -> d
```



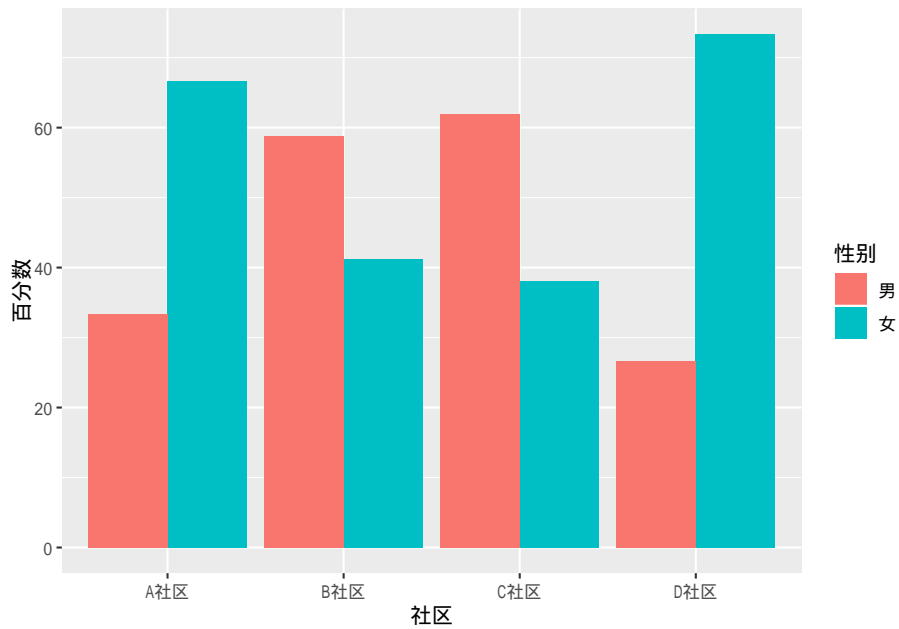
这样，与对频数作图类似就可以做堆叠条形图，每个社区为一个大类，每个大类总计 100%:

```
p <- ggplot(data = d, mapping = aes(
 x = `社区`, fill = `性别`, y = `百分数`
))
p + geom_col()
```



在 `geom_col()` 中加 `position = "dodge"` 可以变成并列条形图，每个社区为一个子类，每个子类中总计 100%:

```
p <- ggplot(data = d, mapping = aes(
 x = `社区`, fill = `性别`, y = `百分数`
))
p + geom_col(position = "dodge")
```

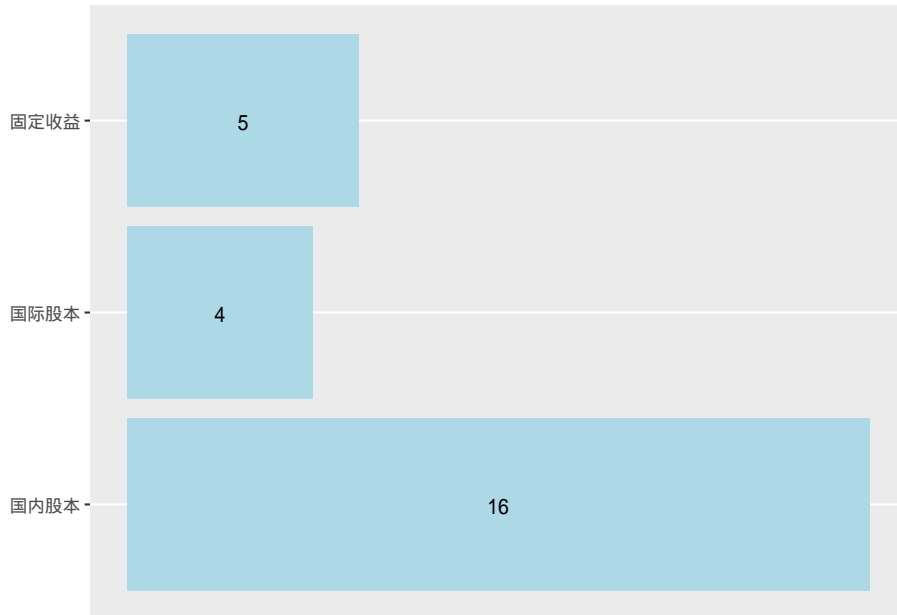


与前一章30.4.4中类似问题相比，对表格数据做比例或者百分数条形图，比基于原始数据直接用 `geom_bar()` 做图要简单而且不易出错。

#### 31.2.1.4 直接标记数字的条形图

比较简单的条形图，可以不使用  $y$  坐标轴而是对每个条形或者堆叠的每个段标出数字。用 `geom_text()` 函数标数字。如：

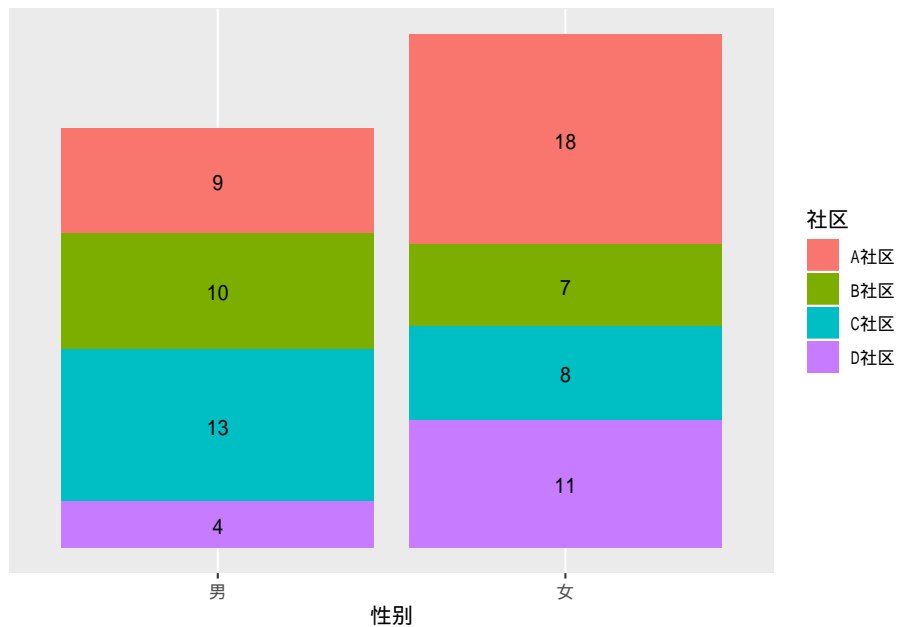
```
d <- d_funds_type %>%
 mutate(type = fct_relevel(type, " 国内股本", " 国际股本", " 固定收益"))
p <- ggplot(data = d, mapping = aes(
 x = type, y = freq))
p + geom_col(fill = "lightblue") +
 geom_text(mapping = aes(
 y = freq / 2, label = paste(freq))) +
 scale_y_continuous(breaks = NULL) +
 coord_flip() +
 labs(x = NULL,
 y = NULL)
```



上述程序中，用了 `scale_y_continuous()` 的 `breaks = NULL` 选项使得 y 轴（经过对换后画在了 x 轴）不画刻度和刻度值。

堆叠的条形图用数字而非坐标轴的例子如下：

```
d_comm %>%
 arrange(`性别`, desc(`社区`)) %>%
 group_by(`性别`) %>%
 mutate(ylabel = cumsum(`频数`) - 0.5*`频数`) -> d
p <- ggplot(data = d, mapping = aes(
 x = `性别`, fill = `社区`, y = `频数`
))
p + geom_col() +
 geom_text(mapping = aes(
 y = ylabel, label = paste(`频数`))) +
 scale_y_continuous(breaks = NULL) +
 labs(y = NULL)
```



上述程序中预先计算了每段的纵坐标值，保存在变量 `ylabel` 中。注意计算过程中的排序与分组操作。

### 31.2.2 点图

条形图的 y 轴一般都代表数量，最小值从 0 开始，一般不画 y 轴最小值大于 0 的条形图。

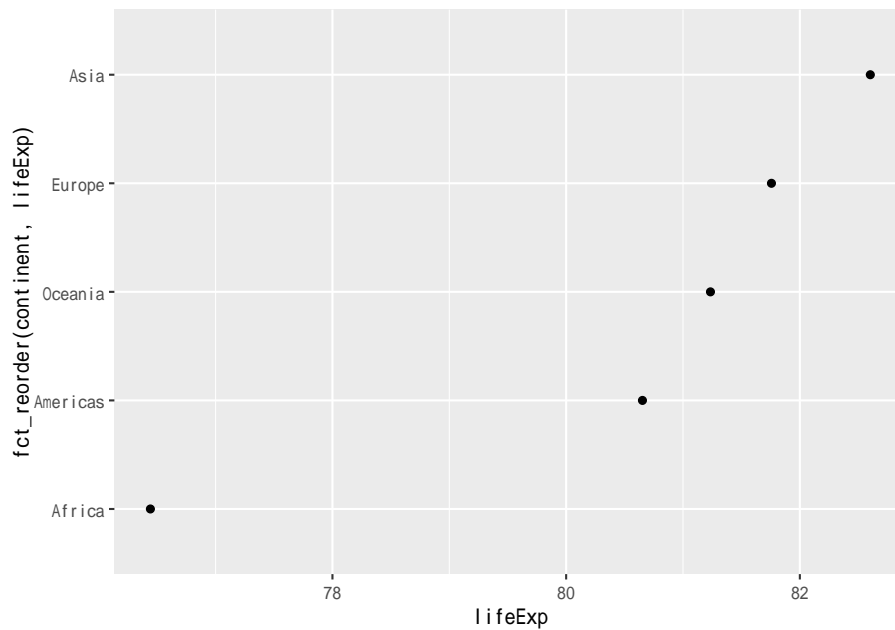
对于不会取 0 的量，可以取 y 轴的范围为某个区间而不必总是从 0 开始。这时，可以用散点位置表示 y 坐标。

考虑 `gapminder` 包的 `gapminder` 数据集中各国的平均期望寿命，先求出每个大洲的期望寿命最大值：

```
library(gapminder)
gapminder %>%
 select(continent, country, lifeExp) %>%
 group_by(continent) %>%
 summarise(lifeExp = max(lifeExp, na.rm=TRUE)) %>%
 ungroup() -> d_gap
```

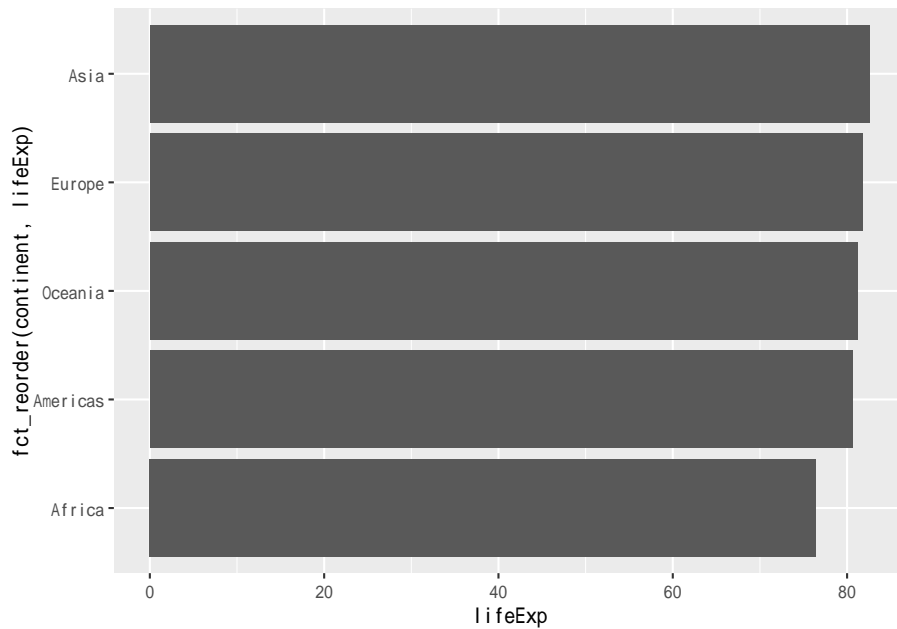
用点图表现各大洲的期望寿命:

```
p <- ggplot(data = d_gap, mapping = aes(
 x = fct_reorder(continent, lifeExp), y = lifeExp))
p + geom_point() +
 coord_flip()
```



程序中并没有人为指定 lifeExp 轴的坐标范围, 作图程序自动选择了合适的范围。对各大洲按照 lifeExp 的值做了排序。这个图形很好地体现了不同大洲的最大期望寿命的差距。如果使用条形图, 因为坐标轴从 0 开始, 各个条形的长度会很接近, 不利于体现差距:

```
p + geom_col() +
 coord_flip()
```



### 31.2.3 热力图

如果交叉分类得到许多个交叉组，每组有一个数量需要展示，用堆叠或者并排条形图可能过于复杂，结果很难判读；如果用点图，不同小类之间的区别也不容易认读。

可以用 x 轴和 y 轴分别表示两种分类，在坐标交叉处用色块颜色代表数量，称为热力图。用 `geom_tile()` 作这种图，将数量映射到 fill 维度。

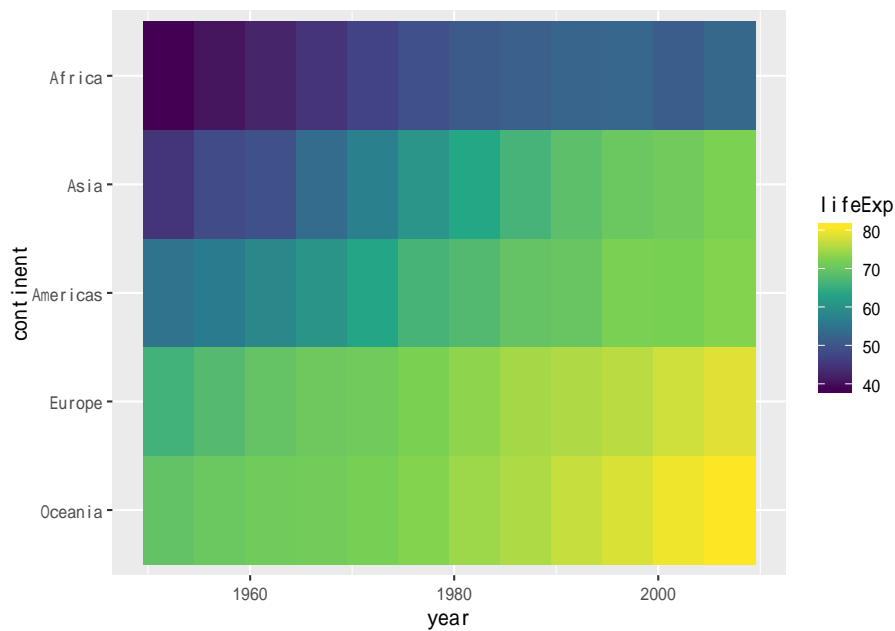
例如，考虑 `gapminder` 包的 `gapminder` 数据集中各国的平均期望寿命，先得到每个大洲每年内各国期望寿命的中位数：

```
gapminder %>%
 select(continent, year, lifeExp) %>%
 group_by(continent, year) %>%
 summarise(lifeExp = median(lifeExp, na.rm=TRUE)) %>%
 ungroup() -> d_gap2a
d_gap2a %>%
 filter(year == 2007) %>%
```

```
arrange(desc(lifeExp)) -> d_gap2b
d_gap2a %>%
 mutate(continent = factor(continent, levels = d_gap2b$continent)) -> d_gap2
```

上面的程序中将各大洲的次序调整为最后一年的中位期望寿命。作热力图：

```
p <- ggplot(data = d_gap2, mapping = aes(
 x = year, y = continent, fill = lifeExp))
p + geom_tile() +
 scale_fill_viridis_c()
```



函数 `scale_fill_viridis_c()` 提供了连续变量到颜色的一种映射，在黑白显示时仍适用，且使用色盲的读者。

热力图中不易读出每个交叉组的数量的具体数值，但是比较容易进行组间比较以及展示发展趋势。

## 31.3 表现分布

对于离散变量，可以用频数、比例、百分数的条形图表现单个离散变量分布，可以用热力图表现两个离散变量的分布。

对于连续型变量，可以用直方图、密度估计图表现单个变量分布，可以对多个变量同时做密度估计图。可以用正态 QQ 图、盒形图、经验分布函数图等表现一元分布。

### 31.3.1 单个一元分布的直方图与密度估计图

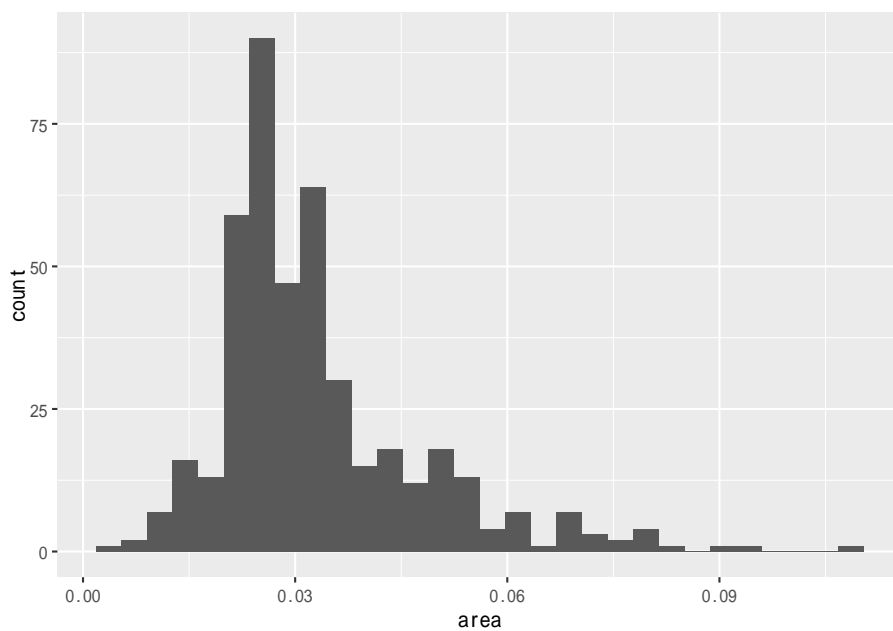
直方图是最常用的表现一元连续变量分布的方法。`geom_histogram()` 作直方图，可以自动选取合适的分组个数，也可以人为指定分组个数。直方图的结果受分组个数与分组的开始位置的影响较大，所以使用直方图时应该尝试选用不同的分组数，选择较适当的分组个数。

`ggplot2` 包中的 `midwest` 数据集包含了美国中西部的一些县的统计数据，如面积（单位：平方英里）。下面的程序对连续取值的数值型变量 `area` 作频数直方图，自动确定分组个数：

```
p <- ggplot(data = midwest, mapping = aes(
 x = area))
p + geom_histogram()

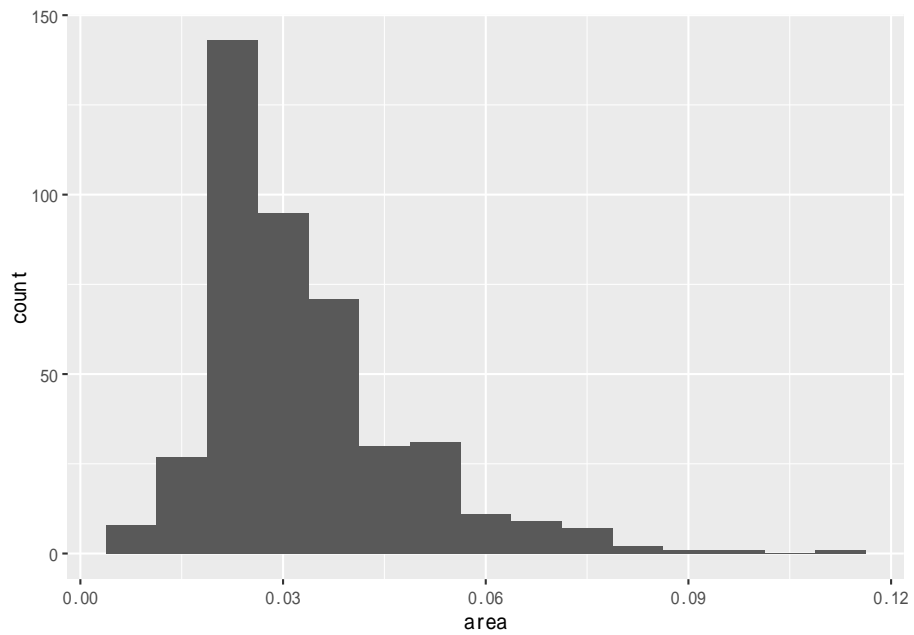
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```





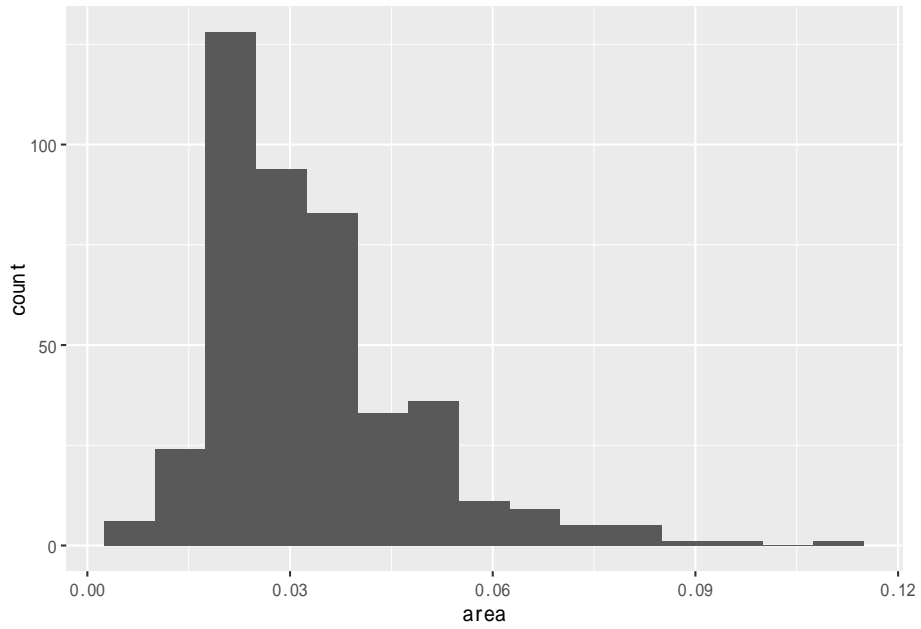
可以在 `geom_histogram()` 中用 `bins =` 指定分组个数, 或者用 `width` 指定分组宽度, 如:

```
p + geom_histogram(bins = 15)
```



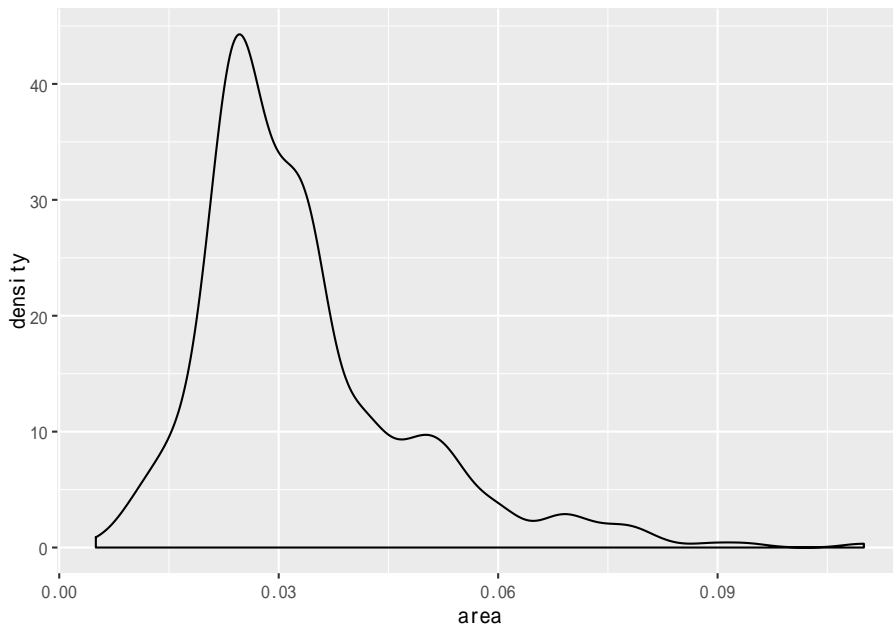
可以用 `boundary` 或者 `center` 指定一个条形的边界或者中心, 用 `breaks` 指定所有分点位置。如

```
p + geom_histogram(bins = 15, boundary = 0.0025)
```



`geom_density()` 可以对连续变量绘制密度的核估计曲线，如：

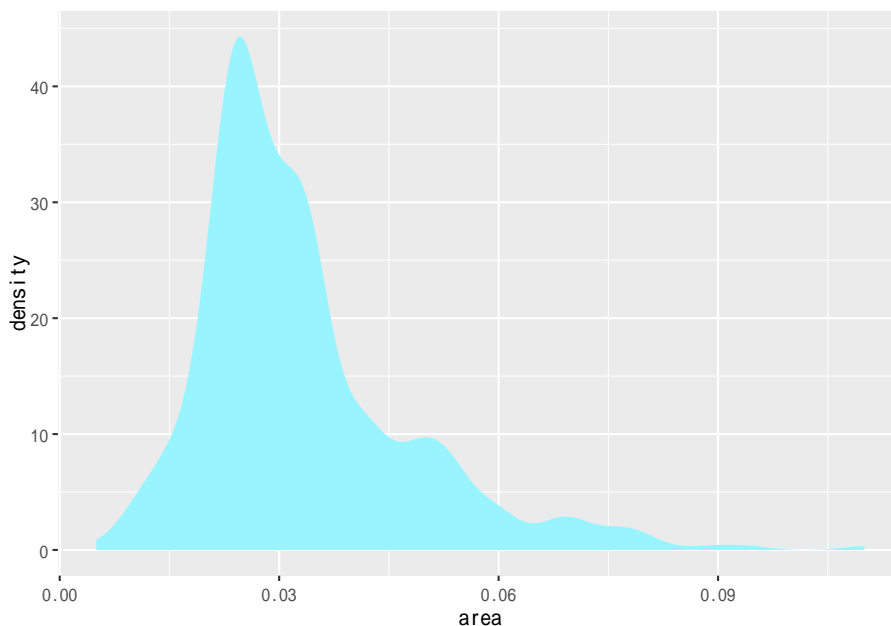
```
p + geom_density()
```



密度曲线估计受到窗宽 `bw` 与核函数 `kernel` 的影响。默认核函数为高斯核（标准正态分布密度作为核函数）。

除了估计曲线，还可以用涂色区域表示密度估计，用 `geom_area(stat = "density")` 函数：

```
p + geom_area(stat = "density", fill = "cadetblue1")
```



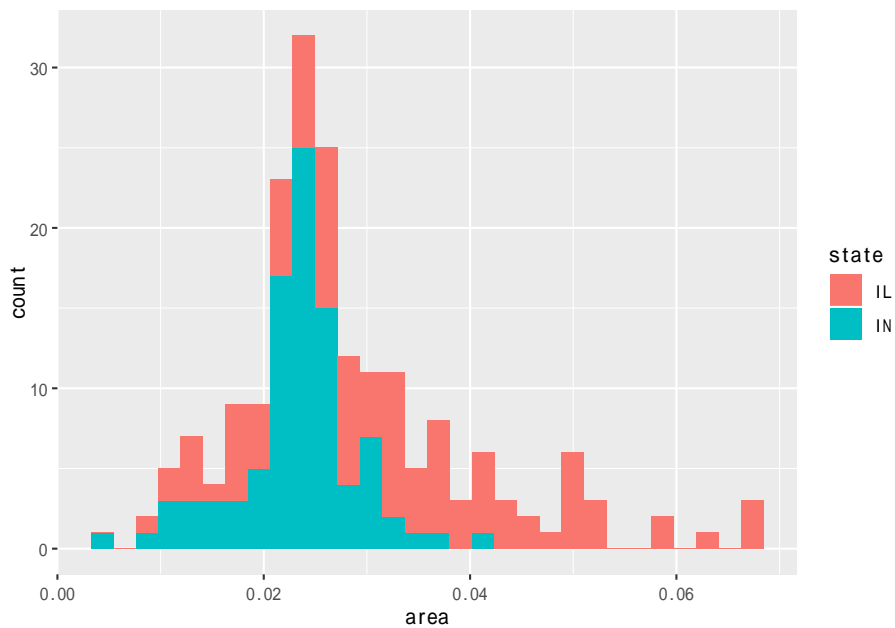
因为核密度估计一般是连续曲线，所以估计的曲线很可能在密度为零的地方估计为正数。ggplot2 在设计时已经考虑了避免这个问题，但是用户自己用 R 的 `density()` 函数作核密度估计时还有这个问题。

### 31.3.2 多个一元分布的直方图和密度图形

考虑 ggplot2 包的 `midwest` 数据集，这是美国中西部若干个县的数据。选取其中的 IL 和 IN 两个州的数据，只要将州名映射为 `fill` 维，就可以在 `geom_histogram()` 或者 `geom_line(stat = "density")` 对各县面积作直方图时按不同州分段显示：

```
midwest %>%
 filter(state %in% c("IL", "IN")) -> d
p <- ggplot(data = d, mapping = aes(
 x = area, fill = state))
p + geom_histogram()
```

```
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

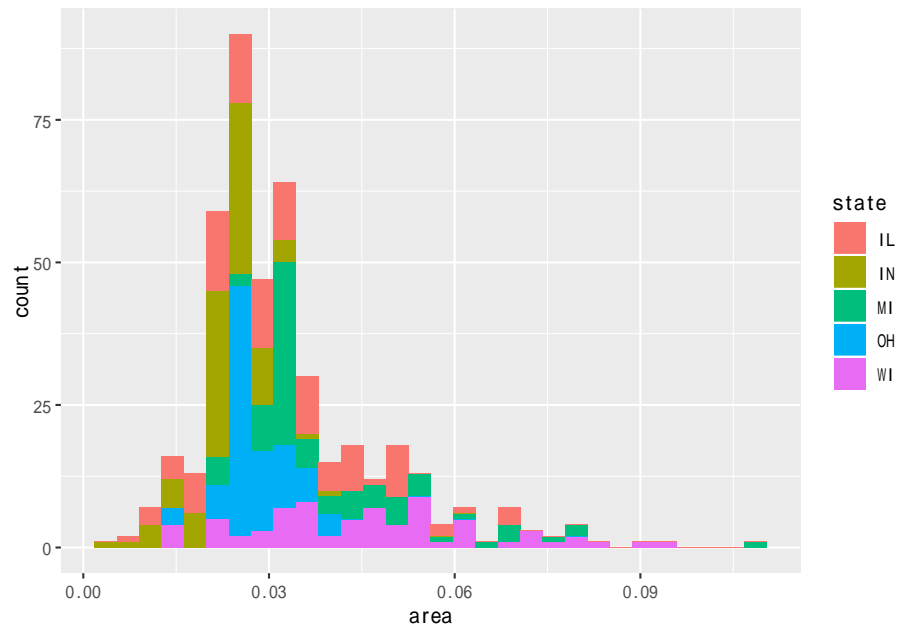


可以看出，IN 州的县面积偏小，IL 州的县面积较大。这样的条形分段的直方图的缺点是，下面的段的直方图形状是清楚的，但是上面一段（上图中为 IL 州）的直方图则看不出形状。

如果将数据中所有的 5 个州都包含并在直方图中分段，图形就变得更难认读：

```
p <- ggplot(data = midwest, mapping = aes(
 x = area, fill = state))
p + geom_histogram()
```

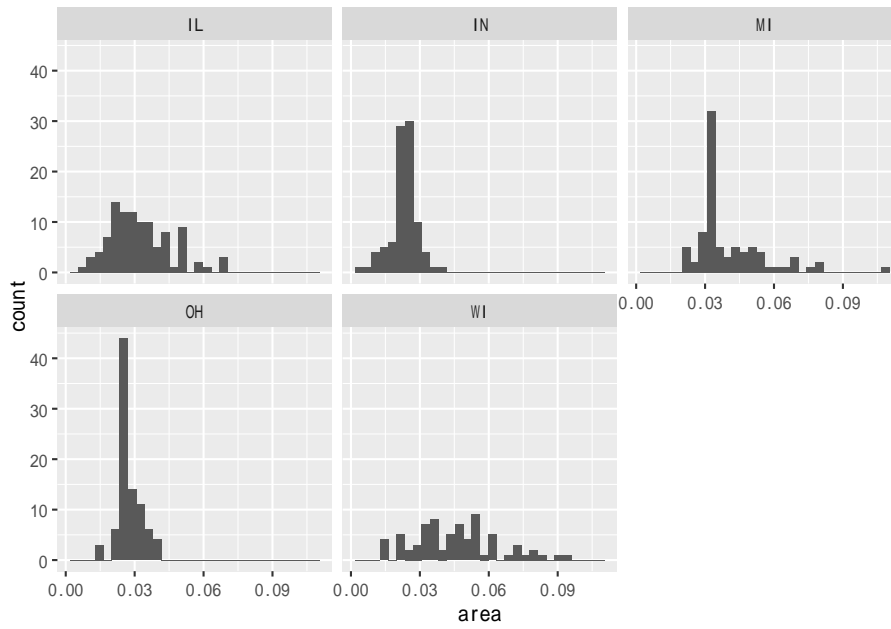
```
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



可以将不同的州分配到不同的小图，如：

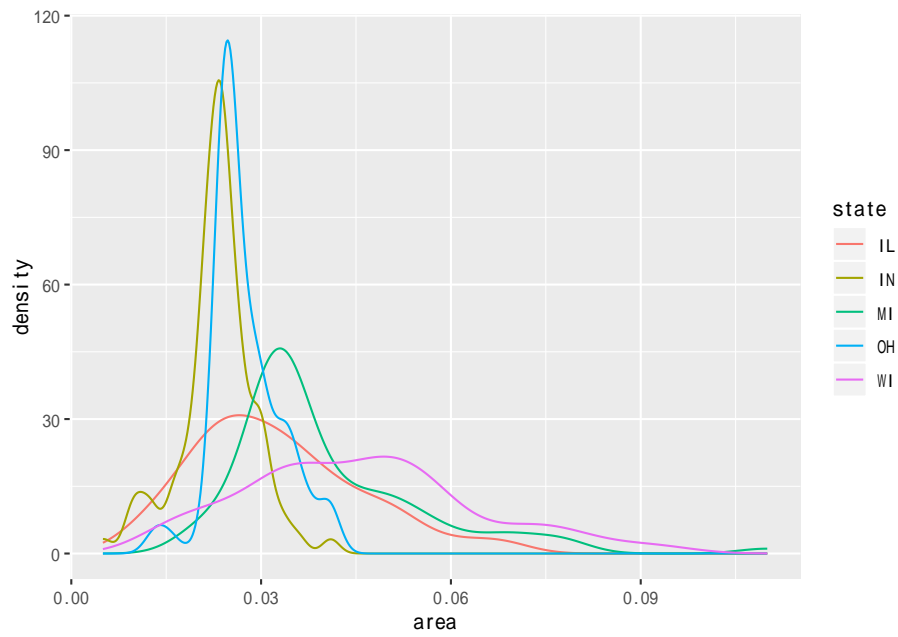
```
p <- ggplot(data = midwest, mapping = aes(
 x = area))
p + geom_histogram() +
 facet_wrap(~ state)
```

```
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



对每个州作密度估计并将多条密度曲线画在同一坐标系可以较好地反映不同州的县面积分布：

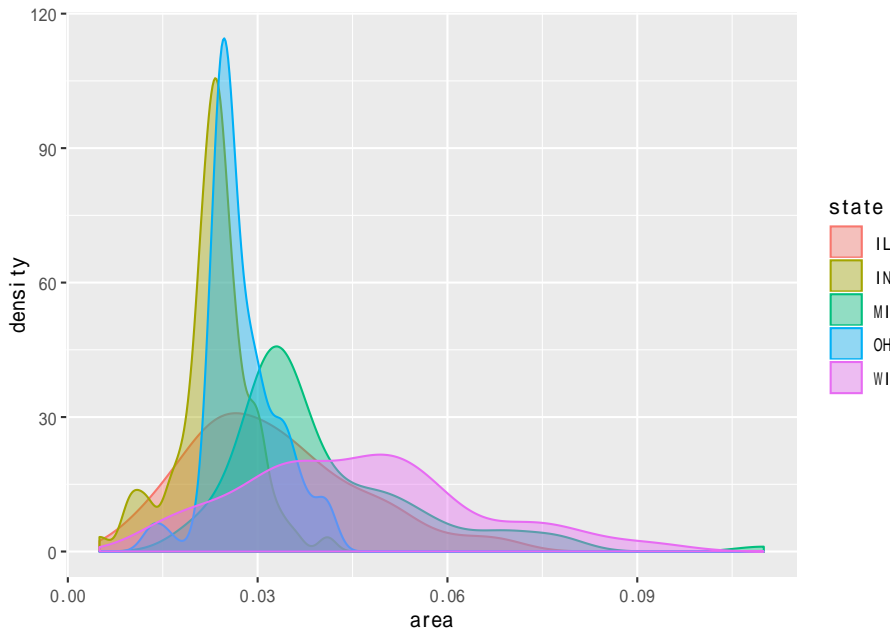
```
p <- ggplot(data = midwest, mapping = aes(
 x = area, color = state))
p + geom_line(stat = "density")
```



增加涂色，用透明度 `alpha` 使其可以重叠显示：

```
p <- ggplot(data = midwest, mapping = aes(
 x = area, color = state, fill = state))
p + geom_density(alpha = 0.4)
```





可以利用 `geom_text()` 或者 `geom_text_repl()` 将各曲线的标签标在曲线附近，就可以取消图例显示。

### 31.3.3 经验分布函数图

直方图和密度估计图能够比较直观地了解数据的分布情况，但是会受到窗宽、分组位置的影响，所以可以看成是对数据的某种建模结果。可以将所有数据观测值画出来，用适当的透明度可以表现点的稠密程度，但是当数据量特别大的时候这种方法就不合适，而且对数据的概括性描述总是必要的。经验分布函数图很好地描述和概括了数据的分布，从中可以很容易得到样本分位数的估计，而且不会受到任何参数选择影响。一组观测值  $X_1, X_2, \dots, X_n$  的经验分布函数定义为

$$F_n(x) = \frac{\#\{X_i : X_i \leq x\}}{n} = \frac{1}{n} \sum_{i=1}^n I_{(-\infty, x]}(x)$$

其中  $\#\{\dots\}$  表示集合元素计数， $I_A(x)$  是集合  $A$  的示性函数，当  $x \in A$  是取 1，否则取 0。

设有如下的 10 位学生的成绩：

```
scores <- c(63, 72, 74, 79, 82, 82, 87, 89, 90, 90)
```

```
scores
```

```
[1] 63 72 74 79 82 82 87 89 90 90
```

可以制作如下的经验分布函数表格：

```
cdf_table <- function(x){
 x <- sort(x)
 n <- length(x)
 tab <- unname(c(table(x)))
 pct = tab / n
 d <- data.frame(
 x = sort(unique(x)),
 freq = tab,
 pct = pct,
 cumfreq = cumsum(tab),
 cumpct = cumsum(pct))
 d
}
knitr::kable(cdf_table(scores))
```

x	freq	pct	cumfreq	cumpct
63	1	0.1	1	0.1
72	1	0.1	2	0.2
74	1	0.1	3	0.3
79	1	0.1	4	0.4
82	2	0.2	6	0.6
87	1	0.1	7	0.7
89	1	0.1	8	0.8
90	2	0.2	10	1.0

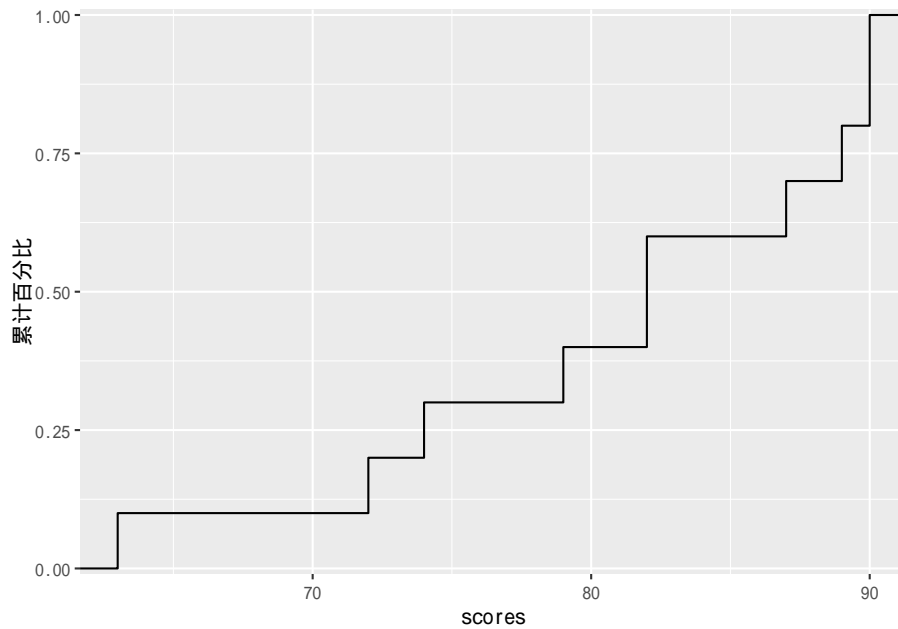
作经验分布函数 (ECDF) 图如下：

```
d <- data.frame(
 scores = c(63, 72, 74, 79, 82, 82, 87, 89, 90, 90))
p <- ggplot(data = d, mapping = aes(
```

```

x = scores, y = ..y..))
p + stat_ecdf(geom = "step") +
 scale_y_continuous(
 limits = c(-.01, 1.01),
 expand = c(0, 0),
 name = " 累计百分比")

```



从上图中，可以看出 80 分及以下的人数 40% 左右，最高分的 25% 学生分数至少为 88 分。当然，准确的数字还需要查看上面的频数分布表。

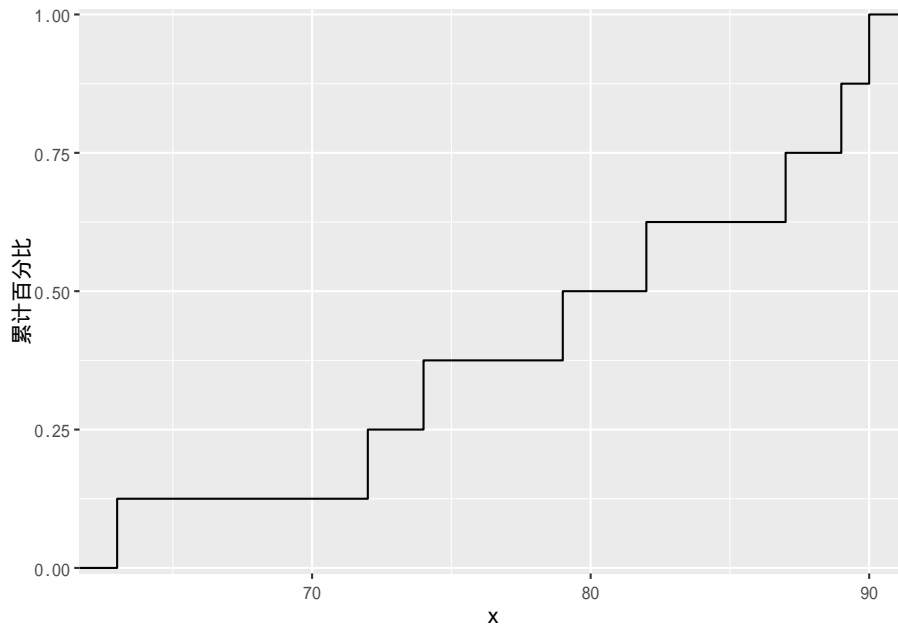
作经验分布函数图，用了 `stat_ecdf()` 统计函数，其中指定 `geom="step"` 选项，并将 y 轴映射为统计函数的结果 `..y..`。如果已经有经验分布函数，则可以直接做阶梯函数图：

```

d <- data.frame(
 scores = c(63, 72, 74, 79, 82, 82, 87, 89, 90, 90))
dfreq <- cdf_table(d$scores)
p <- ggplot(data = dfreq, mapping = aes(
 x = x, y = cumpct))
p + stat_ecdf(geom = "step") +

```

```
scale_y_continuous(
 limits = c(-.01, 1.01),
 expand = c(0, 0),
 name = " 累计百分比")
```



现实中有许多数据的分布呈现严重的偏态，即有比较多的值偏离数据的中心。比如，全国各省的人口数、居民收入，全世界各国的 GDP，等等。

下面是全国 31 个省（不含港澳台）2017 年的人口（单位：万人）、居民年均可支配收入（单位：千元）数据，来自统计年鉴：

```
dpop <- data.frame(
 province =c(
 " 北京", " 天津", " 河北", " 山西", " 内蒙古", " 辽宁", " 吉林", " 黑龙江",
 " 上海", " 江苏", " 浙江", " 安徽", " 福建", " 江西", " 山东", " 河南",
 " 湖北", " 湖南", " 广东", " 广西", " 海南", " 重庆", " 四川", " 贵州",
 " 云南", " 西藏", " 陕西", " 甘肃", " 青海", " 宁夏", " 新疆"),
 pop = c(
 2171L, 1557L, 7520L, 3702L, 2529L, 4369L, 2717L, 3789L, 2418L,
 8029L, 5657L, 6255L, 3911L, 4622L, 10006L, 9559L, 5902L, 6860L,
```

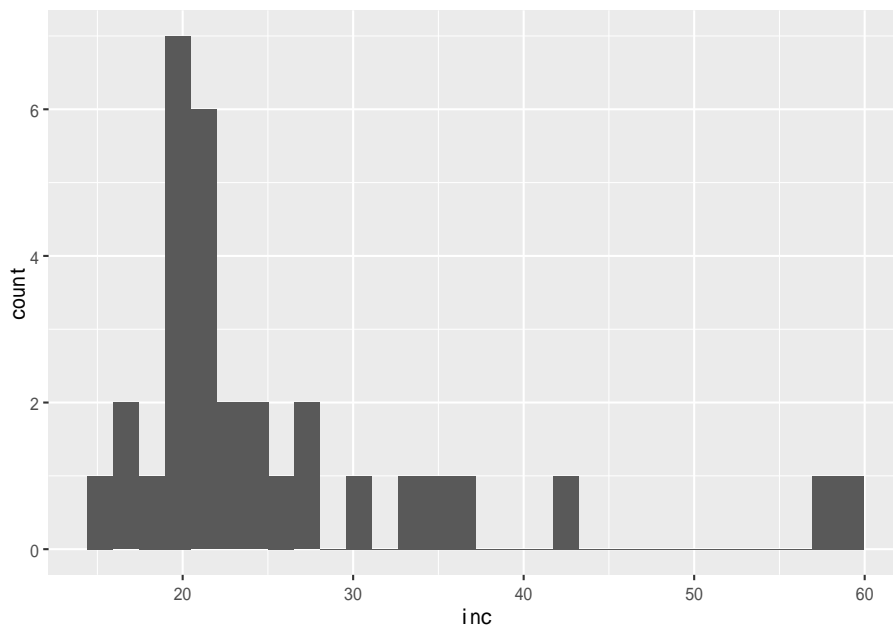
```
11169L, 4885L, 926L, 3075L, 8302L, 3580L, 4801L, 337L, 3835L,
2626L, 598L, 682L, 2445L),
 inc = c(57L, 37L, 21L, 20L, 26L, 28L, 21L, 21L, 59L, 35L, 42L, 22L,
30L, 22L, 27L, 20L, 24L, 23L, 33L, 20L, 23L, 24L, 20L, 17L, 18L,
15L, 20L, 16L, 19L, 21L, 20L))
knitr::kable(dpop)
```

province	pop	inc
北京	2171	57
天津	1557	37
河北	7520	21
山西	3702	20
内蒙古	2529	26
辽宁	4369	28
吉林	2717	21
黑龙江	3789	21
上海	2418	59
江苏	8029	35
浙江	5657	42
安徽	6255	22
福建	3911	30
江西	4622	22
山东	10006	27
河南	9559	20
湖北	5902	24
湖南	6860	23
广东	11169	33
广西	4885	20
海南	926	23
重庆	3075	24
四川	8302	20
贵州	3580	17
云南	4801	18
西藏	337	15
陕西	3835	20
甘肃	2626	16
青海	598	19
宁夏	682	21
新疆	2445	20

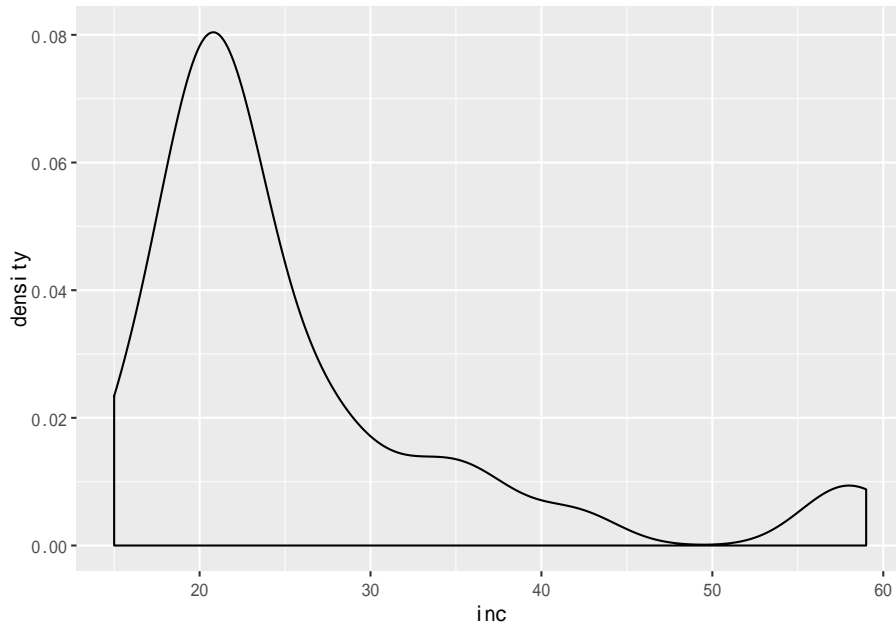
作各省人均可支配收入的直方图以及密度估计图:

```
p <- ggplot(data = dpop, mapping = aes(
 x = inc))
p + geom_histogram()
```

```
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



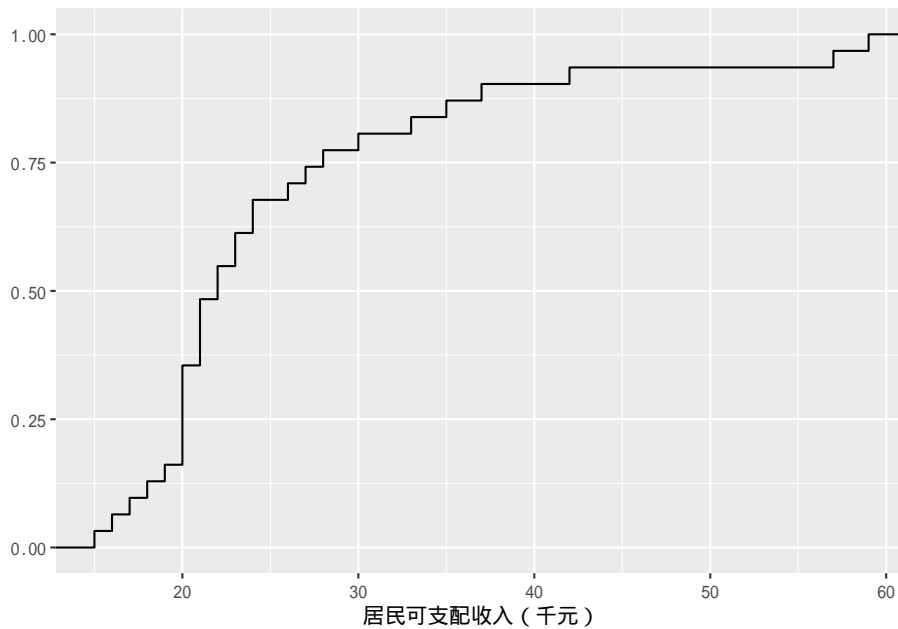
```
p + geom_density()
```



这样的分布称为右偏分布，特点是大部分数据取值集中在左侧，但是有不可忽略的多个数据点出现在远离数据取值中心的右侧。作收入的经验分布函数图：

```
p <- ggplot(data = dpop, mapping = aes(
 x = inc, y = ..y..))
p + stat_ecdf(geom = "step") +
 labs(x = "居民可支配收入(千元)", y = NULL)
```





如果将数据做适当的变换，比如对数变换，分布形状可能变得比较正态。

### 31.3.4 QQ 图

对连续型随机变量  $X$ ，设其分布函数为  $F(x)$ ，其反函数  $F^{-1}(p), p \in (0, 1)$  称为  $X$  的分位数函数 (quantile function)。为了估计  $F^{-1}(p)$ ，可以将观测值  $X_1, X_2, \dots, X_n$  从小到大排序为  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ ，令  $X_{(i)}$  作为  $F^{-1}(\frac{i}{n})$  的估计，或者做连续性修正，作为  $F^{-1}(\frac{i-0.375}{n+0.25})$  的估计。

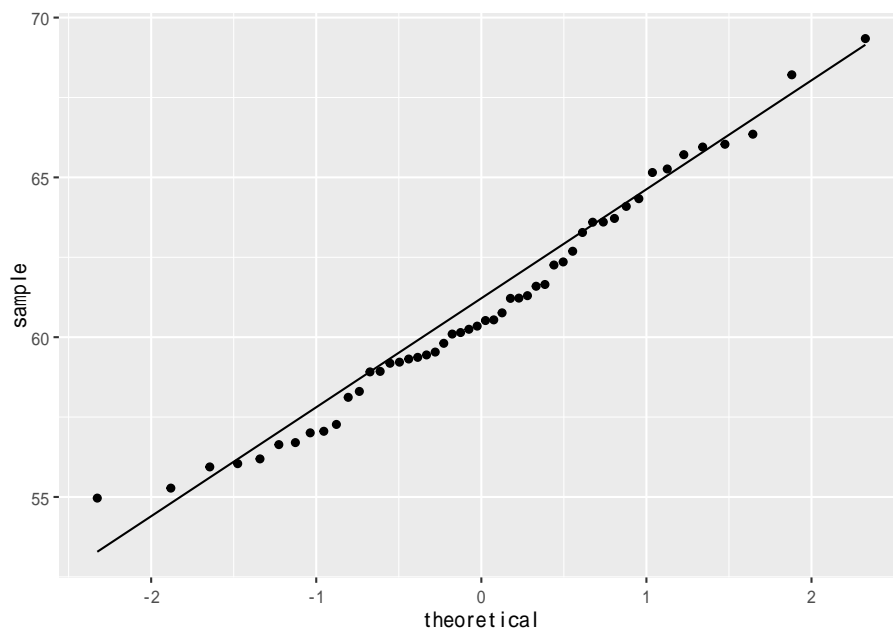
为了对比总体  $X$  的观测数据与某个理论分布比如正态分布之间的差距，可以用上述的样本分位数为散点纵坐标，以相应的理论分位数为横坐标作散点图，称为 QQ 图。如果总体服从所对比的分布，散点应该在一条直线附近，否则就会有明显的偏离。

最常见的 QQ 图是正态 QQ 图。下面对一些模拟数据做正态 QQ 图。

首先是模拟的正态分布数据：

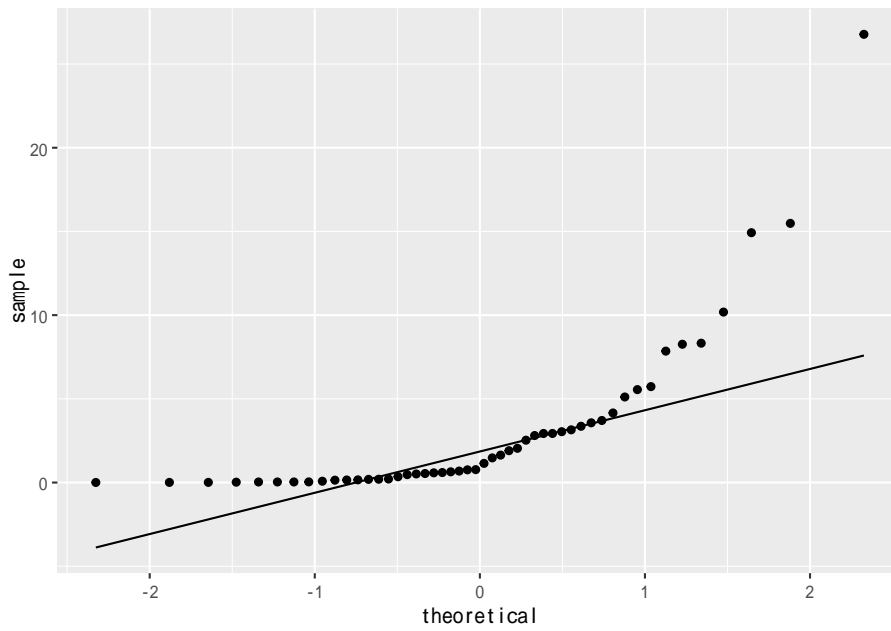
```
set.seed(102)
d <- data.frame(x = 60 + rnorm(50, sd = 3))
```

```
p <- ggplot(data = d, mapping = aes(
 sample = x))
p + geom_qq() +
 geom_qq_line()
```



模拟的对数正态数据:

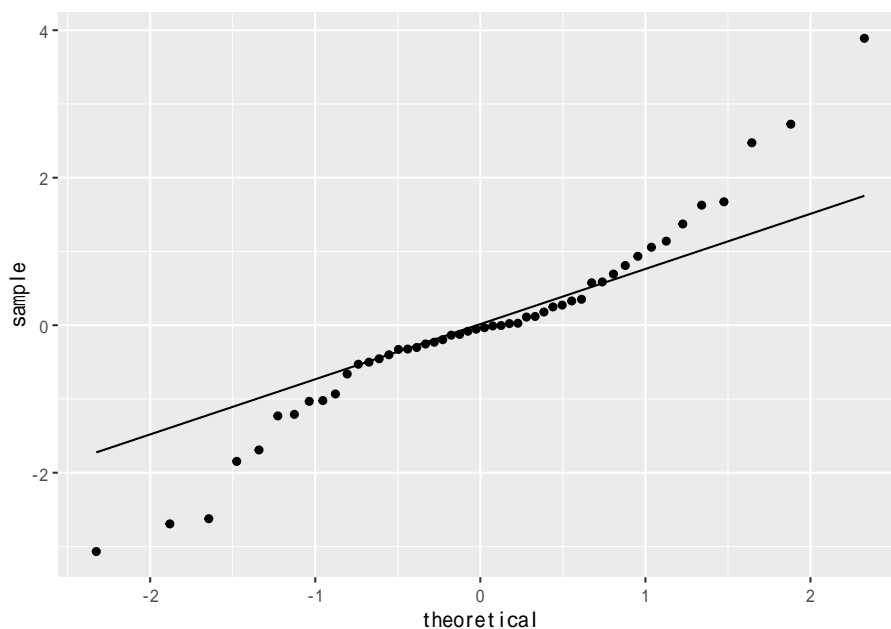
```
set.seed(101)
d <- data.frame(x = 10^(rnorm(50, sd = 1)))
p <- ggplot(data = d, mapping = aes(
 sample = x))
p + geom_qq() +
 geom_qq_line()
```



上图为典型的右偏分布正态 QQ 图形状。

作  $t(3)$  分布样本的正态 QQ 图:

```
set.seed(104)
d <- data.frame(x = rt(50, df=3))
p <- ggplot(data = d, mapping = aes(
 sample = x))
p + geom_qq() +
 geom_qq_line()
```



上图为典型的两侧都厚尾的分布的正态 QQ 图形状。

### 31.3.5 盒形图

如果需要同时表现多个分布，上面给出了同时绘制多条密度曲线的方法，但是这主要比较分布形状，而不能比较分布的主要数字特征。盒形图、小提琴图、脊线图 (ridgeline plot) 比较适合同时表现和比较多个分布。

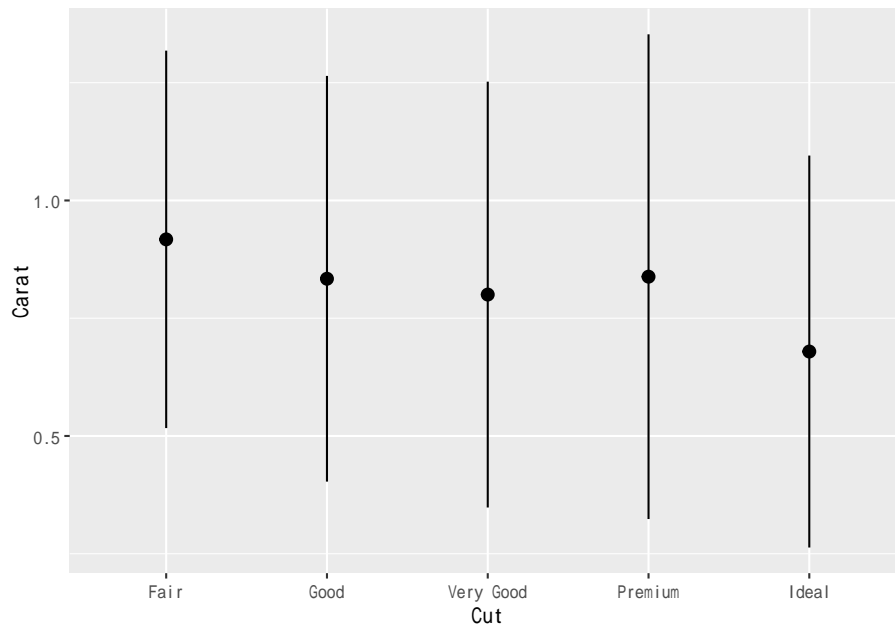
比较多个分布的需求，可能是比较多个类似的变量，比如，同一个学校的不同科目的成绩；也可以是比较不同组的同一属性，比如，比较不同省的城镇居民收入。

如果想仅表现每个分布的均值和标准差，可以对每个均值加减标准差画一个小的条形。比如，ggplot2 扩展包的 diamonds 数据集包含了 5 万多枚钻石的分级、大小、价格数据，我们从中抽取不含缺失数据的 1000 个观测，按 5 种不同品质 (Fair, Good, Very Good, Premium, Ideal) 分类，对克拉重量计算每类的均值和标准差：

```
set.seed(101)
diamonds %>%
 na.omit() %>%
 sample_n(1000) -> ddsmall
ddsmall %>%
 group_by(cut) %>%
 summarise(
 mean_carat = mean(carat),
 sd_carat = sd(carat)) -> ddsmall_summ
```

对每个类画以均值为中心、以均值加减标准差为两段的条形:

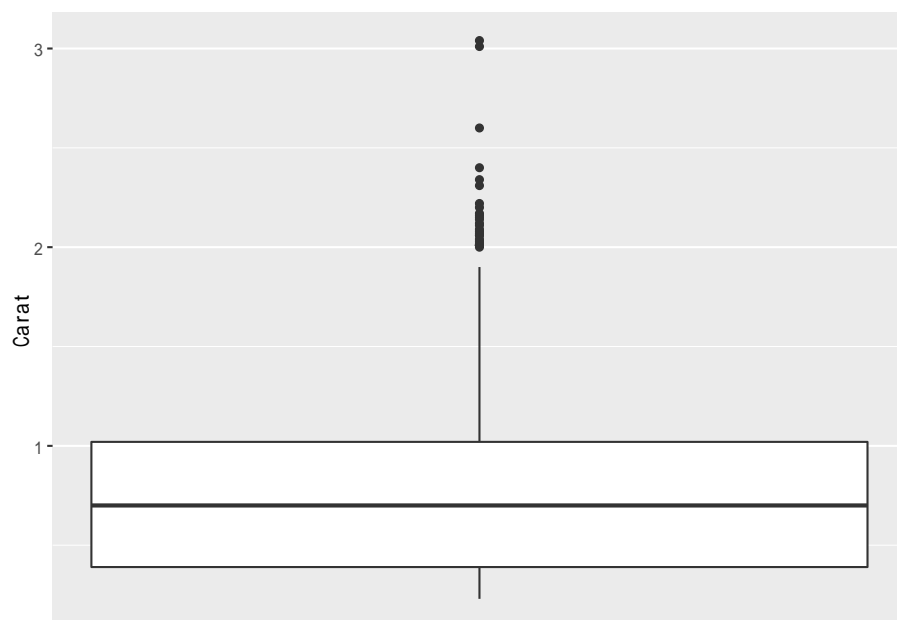
```
p <- ggplot(data = ddsmall_summ, mapping = aes(
 x = cut,
 y = mean_carat))
p + geom_point(size = 1.2) +
 geom_pointrange(mapping = aes(
 ymin = mean_carat - sd_carat,
 ymax = mean_carat + sd_carat)) +
 labs(x = "Cut",
 y = "Carat")
```



竖线的上下界也可以改为正态情形下的近似 95% 置信限，即  $\bar{x} \pm 1.96S$ 。可以用 `coord_flip()` 颠倒 x 轴和 y 轴。

这种图形对分布信息压缩过甚，对于非正态的分布则不能表现分布的形状。盒形图可以表现分布的中位数、四分之一和四分之三分位数，最小值、最大值和可能的离群值。`geom_boxplot()` 作盒形图。比如，`diamonds` 数据集抽样 1000 个观测的克拉重量的盒形图：

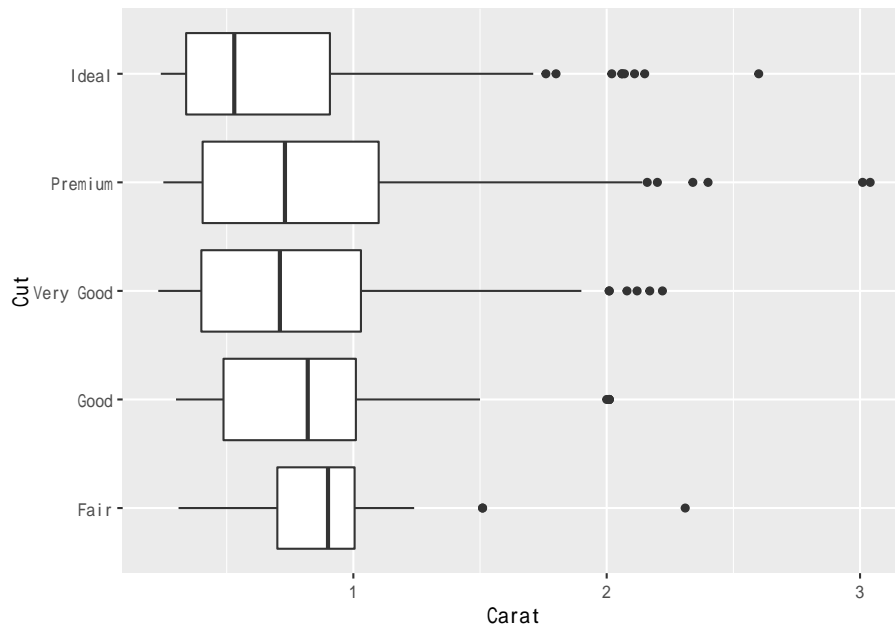
```
p <- ggplot(data = dds_small, mapping = aes(
 y = carat))
p + geom_boxplot() +
 scale_x_continuous(breaks = NULL) +
 labs(x = NULL,
 y = "Carat")
```



从这个盒形图可以看出克拉重量呈现右偏分布。

比较不同品质的钻石的克拉重量，并将盒形图横向画出：

```
p <- ggplot(data = ddsml, mapping = aes(
 x = cut,
 y = carat))
p + geom_boxplot() +
 coord_flip() +
 labs(x = "Cut",
 y = "Carat")
```



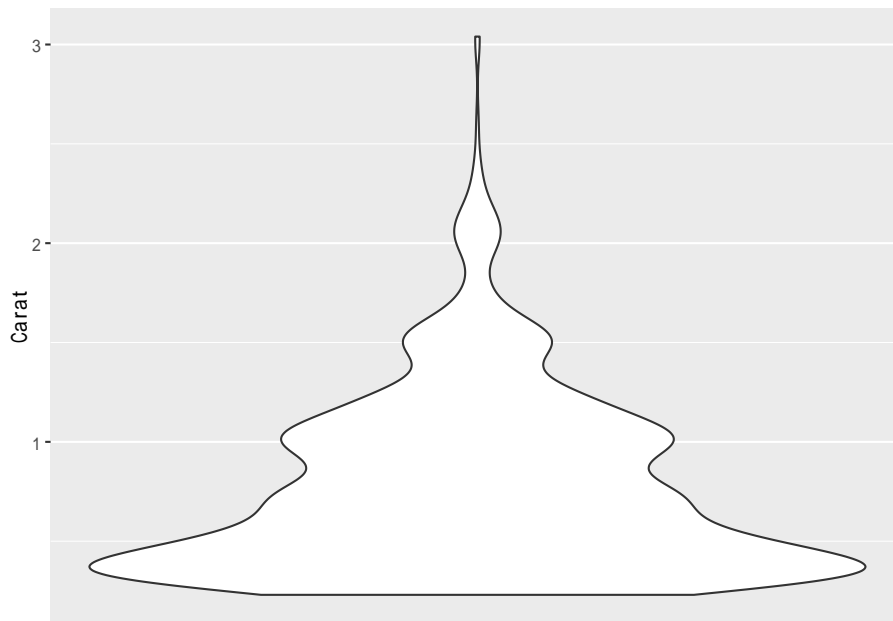
从克拉重量中位数看，品质从 Fair 到 Ideal 基本上是逐步下降的，只有 Very Good 级和 Premium 级的中位重量略有乱序。除了 Fair 和 Good 级，更好品质的三个等级的重量都呈右偏，而 Fair 和 Good 级的分布则相对来说不偏态。

### 31.3.6 小提琴图

小提琴图是盒形图的一个改进版本，比如，钻石数据抽样数据集的克拉重量的小提琴图：

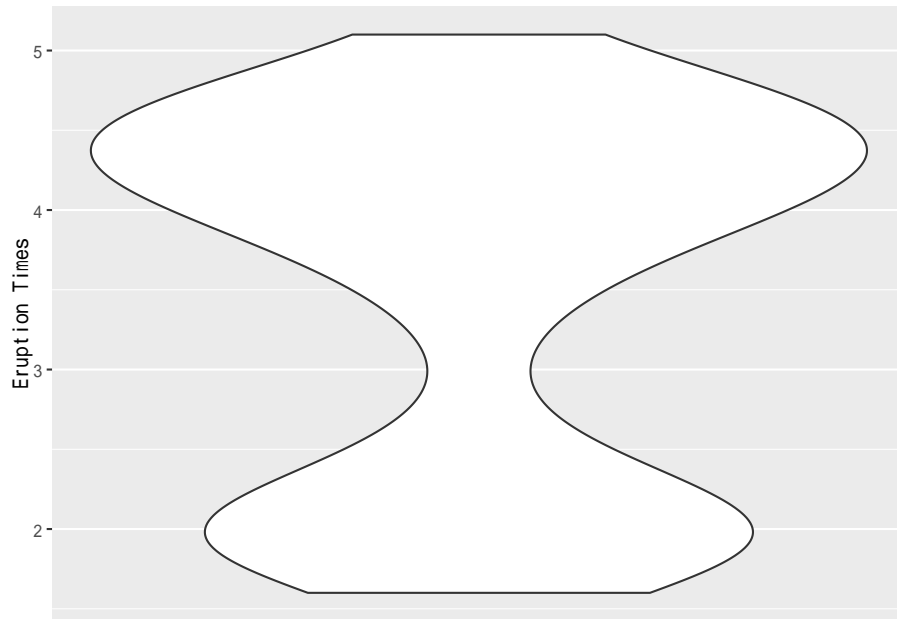
```
p <- ggplot(data = dds_small, mapping = aes(
 x = 0, y = carat))
p + geom_violin() +
 scale_x_continuous(breaks = NULL) +
 labs(x = NULL,
 y = "Carat")
```





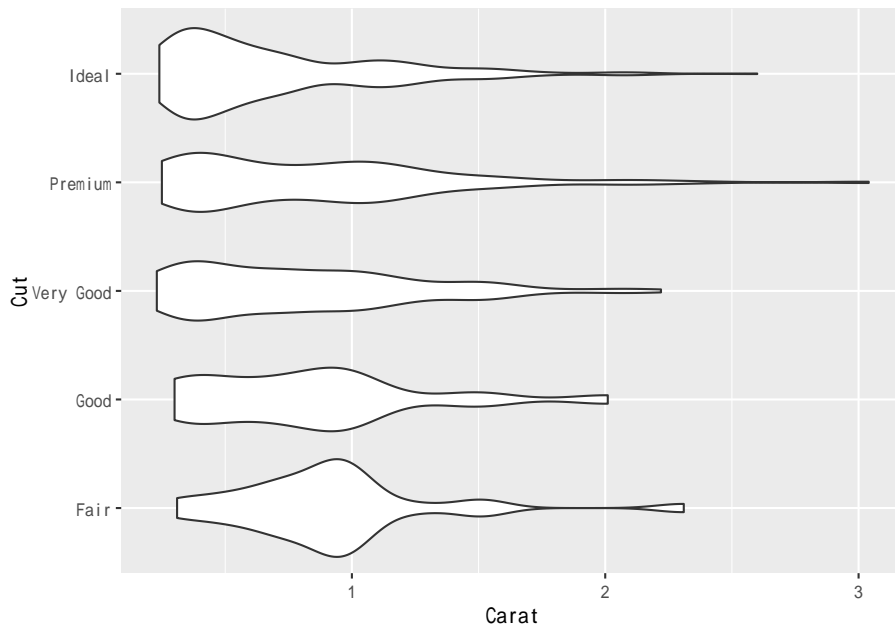
这实际是将密度估计曲线旋转 90 度画出，并左右镜像画两条。它提供了比盒形图更详细的分布信息，比如，可以表现出双峰分布。下面是 R 的 `faithful` 数据集中变量 `eruptions` 的小提琴图：

```
data(faithful)
p <- ggplot(data = faithful, mapping = aes(
 x = 0, y = eruptions))
p + geom_violin() +
 scale_x_continuous(breaks = NULL) +
 labs(x = NULL,
 y = "Eruption Times")
```



钻石数据小数据集中各品质级别的钻石克拉重量的小提琴图比较，横向放置：

```
p <- ggplot(data = ddsml, mapping = aes(
 x = cut,
 y = carat))
p + geom_violin() +
 coord_flip() +
 labs(x = "Cut",
 y = "Carat")
```

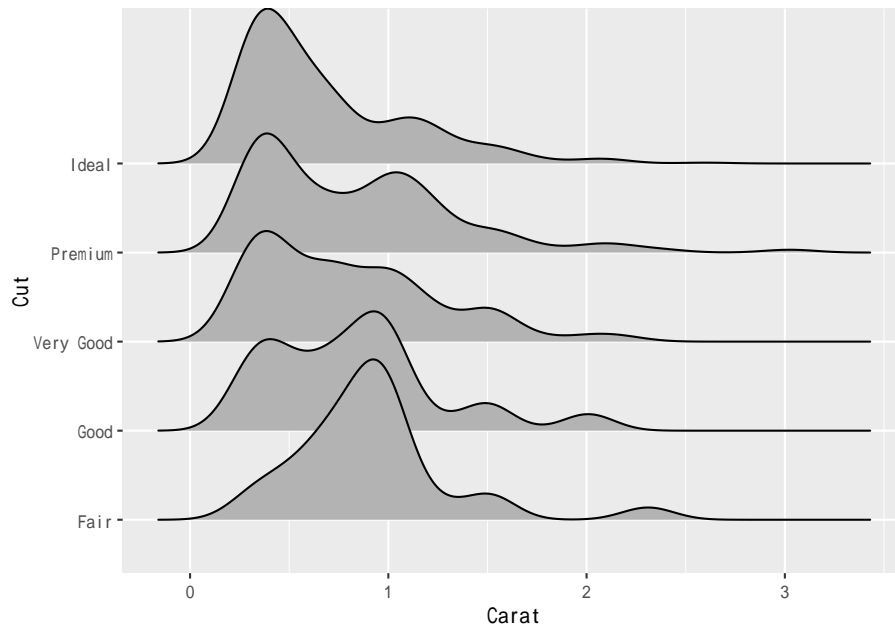


### 31.3.7 多个密度的其它画法

31.3.2给出了将多个密度曲线叠加画出的方法。在对多个组分别画某个连续变量的密度估计时，可以将组别放置在 y 轴上，仍用 x 轴表示该连续变量的取值范围。这种图形尤其适用于考察随着时间或者地点变化的变量分布。ggridges 包提供了 `geom_density_ridges()` 函数作脊线图。对钻石小数据集的各个级别的克拉重量作分布密度脊线图：

```
library(ggridges)
p <- ggplot(data = ddsml, mapping = aes(
 x = carat,
 y = cut))
p + geom_density_ridges() +
 labs(x = "Carat",
 y = "Cut")
```

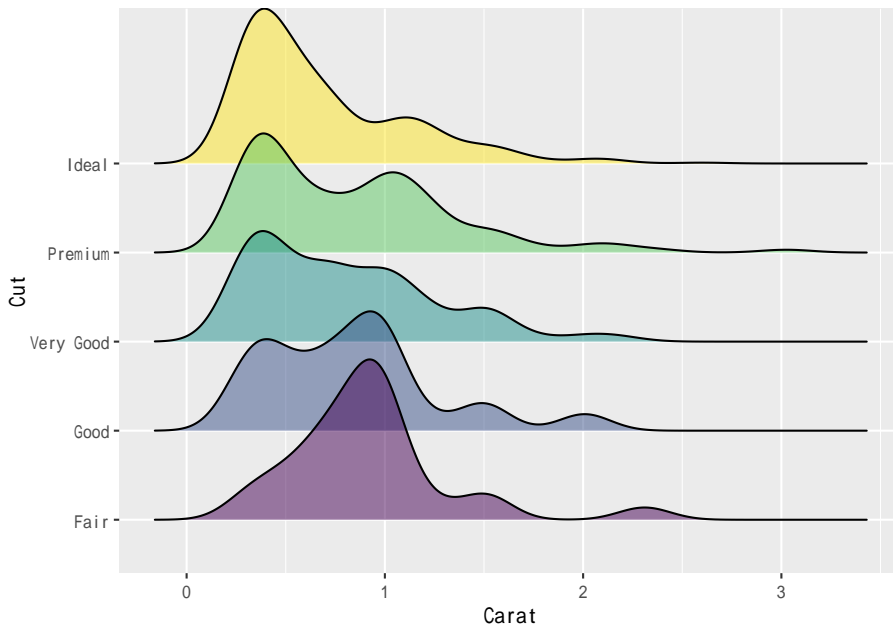
```
Picking joint bandwidth of 0.13
```



可以添加颜色和透明度:

```
p <- ggplot(data = dds_small, mapping = aes(
 x = carat,
 y = cut,
 fill = cut))
p + geom_density_ridges(alpha = 0.5) +
 guides(fill = FALSE) +
 labs(x = "Carat",
 y = "Cut")
```

```
Picking joint bandwidth of 0.13
```



可以看出不同级别的钻石的重量密度峰值的变化，以及一定的多峰分布特征。

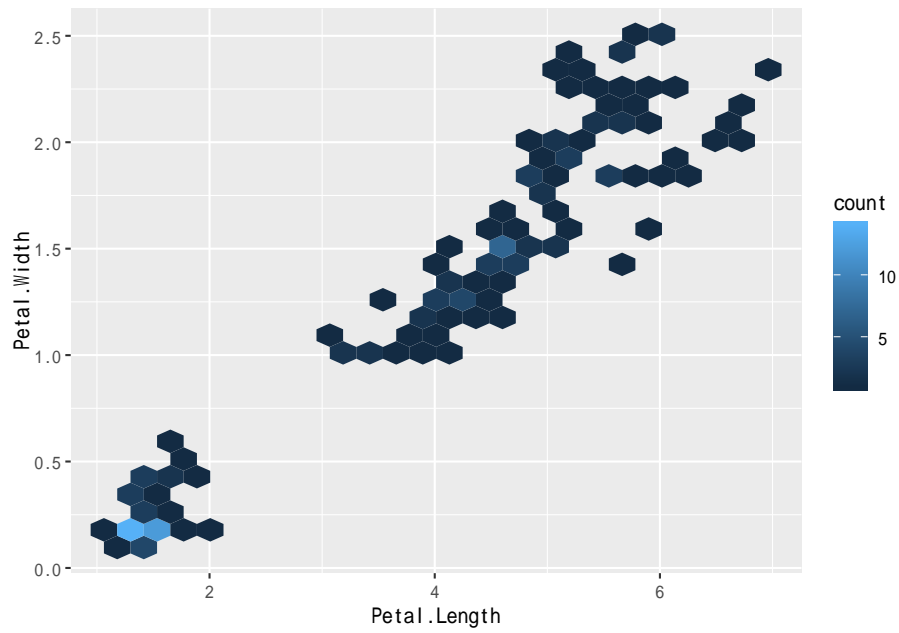
### 31.3.8 二元分布直方图和等值线图

对二元分布样本  $(x_i, y_i)$ ，可以用散点图表现其分布情况；当点比较多时，可以使用 `alpha` 参数设定一定的透明度，使得颜色深表示点比较密集，颜色浅表示点比较稀疏。但是，当点的数量很大时，设定透明度就效果不好了，这时稀疏的地方可能基本看不到颜色，稠密的地方就是一团黑，无法分别分布情况。

可以用类似直方图的想法将二元分布的取值区域分块，每块统计频数，然后用不同深浅的颜色为每个小块染色，颜色的深浅代表频数高低。这样的图形称为二维直方图。每个小块可以是正方形或者正六边形。

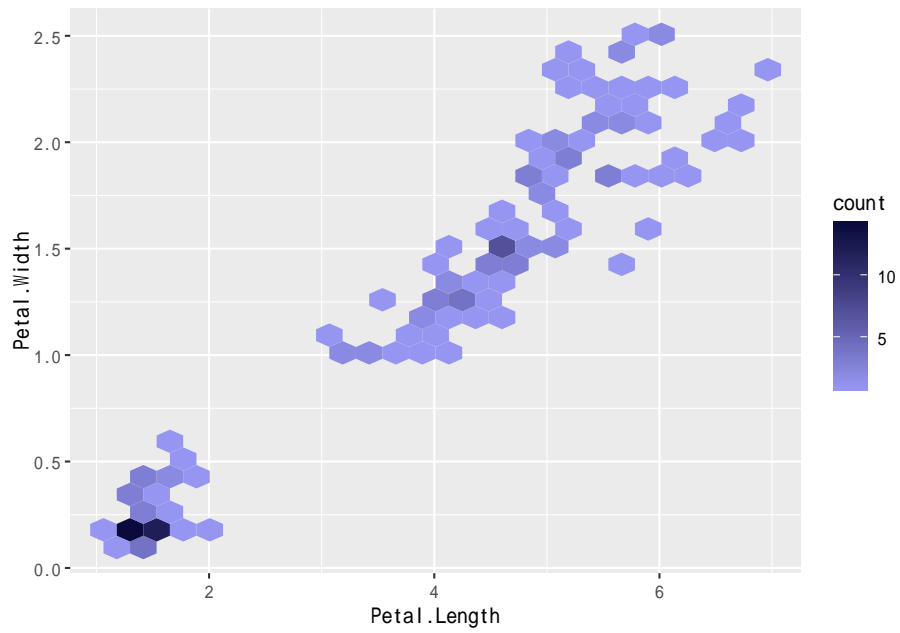
例如，对 `iris` 数据的花瓣长、宽变量画正六边形的二维直方图：

```
p <- ggplot(data = iris, mapping = aes(
 x = Petal.Length, y = Petal.Width))
p + geom_hex(bins = 25)
```



在 `geom_bins()` 中可以用 `bins=` 指定分块的个数。上图中颜色浅的值高，颜色浅的值低。可以用 `scale_color_gradient` 人为地指定代表频数的渐变色，可以改为用深色代表高值，如：

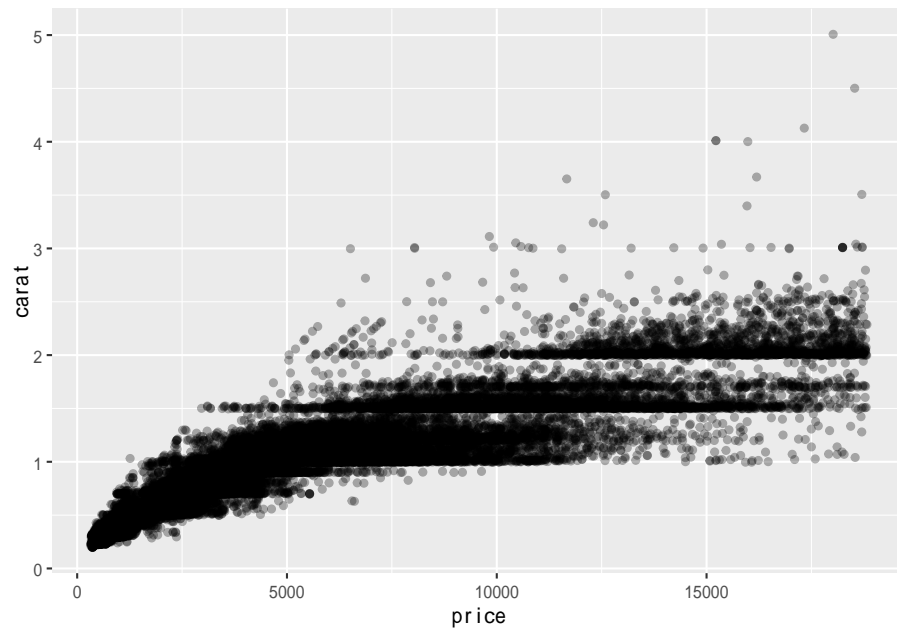
```
p + geom_hex(bins = 25) +
 scale_fill_gradient(
 low = "#9696F2",
 high = "#0A0A3D")
```



可以利用 `colourpicker` 包的 `colourPicker` 函数或者该包的 RStudio Addin – `colourPicker` 挑选颜色。

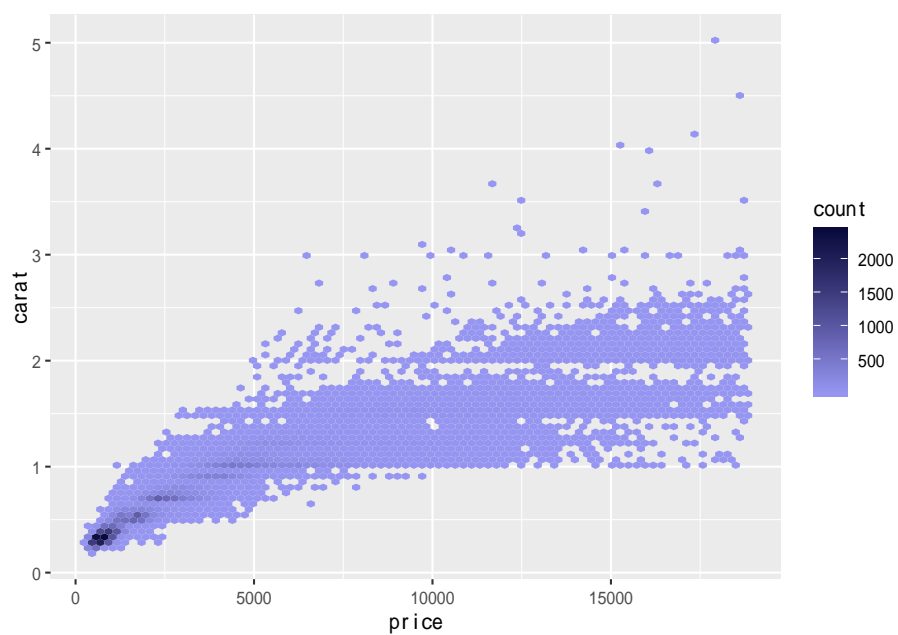
下面做 `ggplot2` 包的 `diamonds` 数据集中的 `price` 和 `carat` 的散点图和二维直方图：

```
p <- ggplot(data = diamonds, mapping = aes(
 x = price, y = carat))
p + geom_jitter(alpha = 0.3)
```



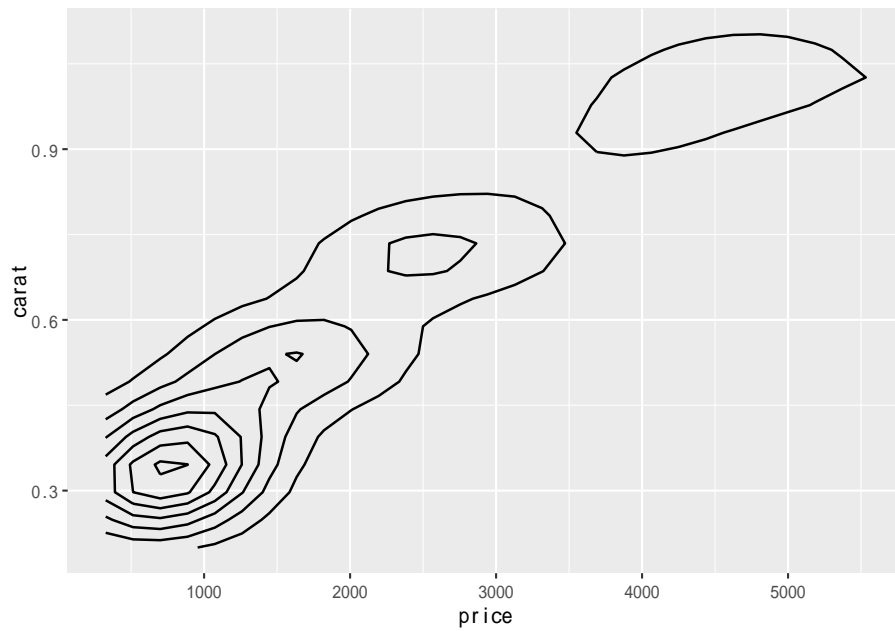
```
p + geom_hex(bins = 80) +
 scale_fill_gradient(
 low = "#9696F2",
 high = "#0A0A3D")
```





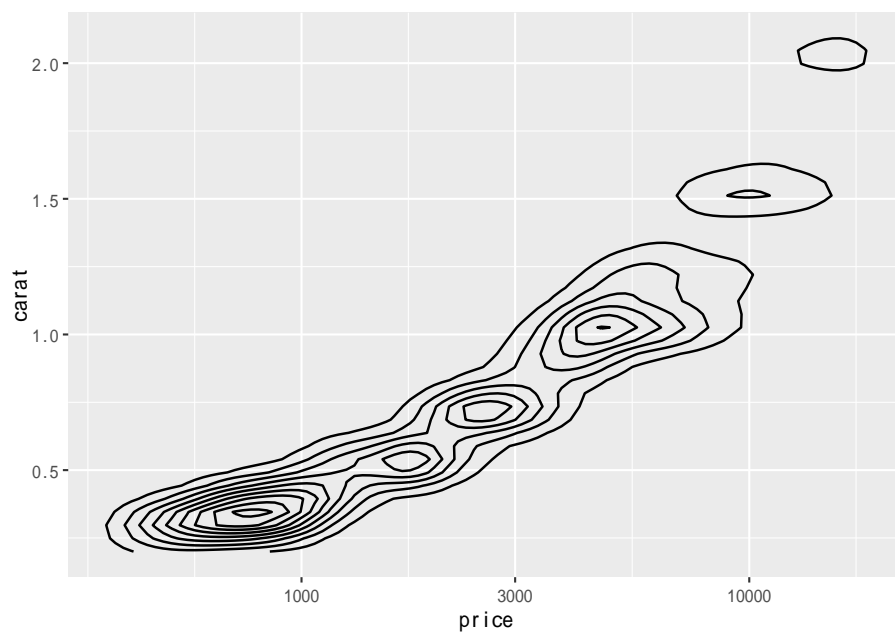
二维直方图估计的是二元阶梯函数，也可以像一元情况估计分布密度那样估计连续的二元分布密度，并用等值线图表示（类似于地图中的等高线）。例如 diamonds 数据集中价格和重量的联合密度估计：

```
p + stat_density_2d(color = "black", size = 0.6)
```



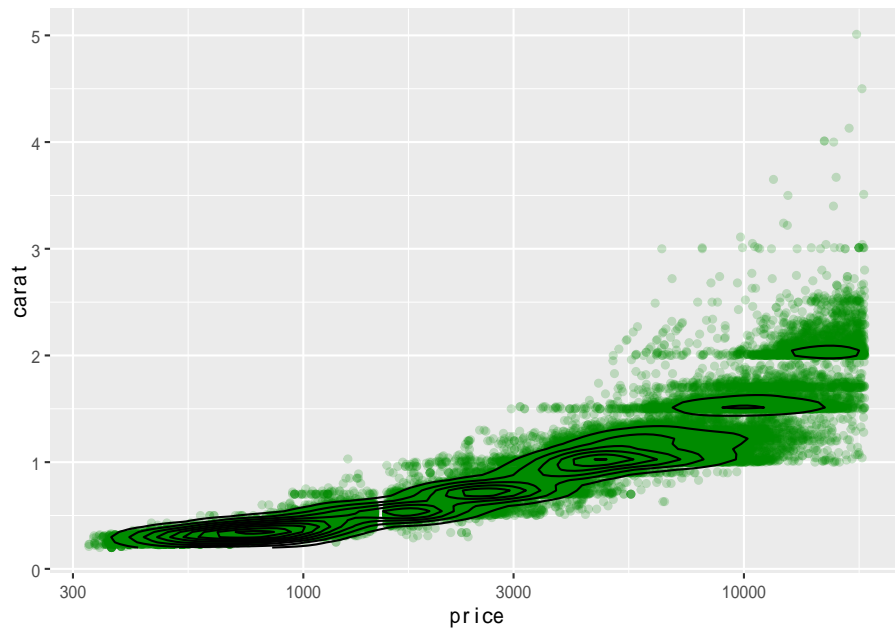
其中 `color` 是等值线颜色，`size` 是等值线粗细，单位是毫米。密度很低的区域被忽略了，为此，将 `price` 用对数轴：

```
p + stat_density_2d(color = "black", size = 0.6) +
 scale_x_log10()
```



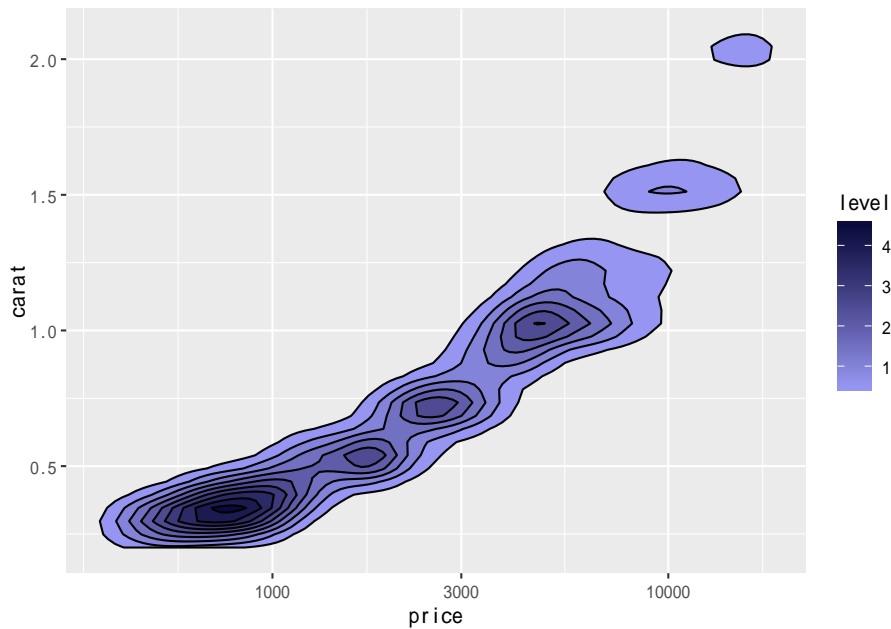
等值线图还可以与散点图叠加在一起:

```
p + geom_point(alpha = 0.2, color = "green4") +
 stat_density_2d(color = "black") +
 scale_x_log10()
```



等值线图还可以配合不同密度区域的不同颜色填充，如：

```
p + stat_density_2d(
 mapping = aes(fill = ..level..),
 color = "black", size = 0.5,
 geom = "polygon") +
 scale_fill_gradient(
 low = "#9696F2",
 high = "#0A0A3D") +
 scale_x_log10()
```



## 31.4 表现比例

条形图、堆叠条形图、并排条形图可以用来表现比例以及分组内的比例。饼图经常用来表现单个分组方式的各组比例。

### 31.4.1 单个比例分布

考虑 diamonds 数据集中各级别，首先计算其频数分布：

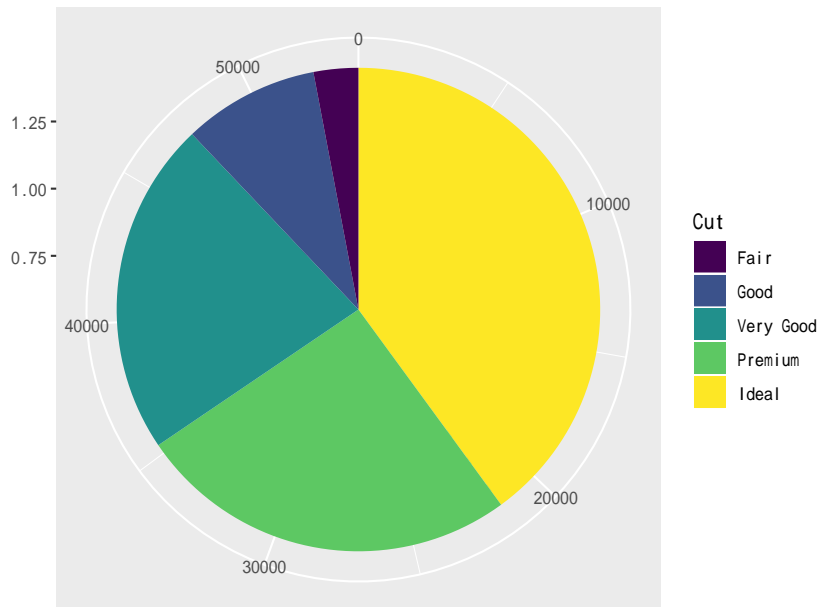
```
ddstab <- cdf_table(diamonds$cut)
knitr::kable(ddstab)
```

x	freq	pct	cumfreq	cumpct
Fair	1610	0.0298480	1610	0.0298480
Good	4906	0.0909529	6516	0.1208009
Very Good	12082	0.2239896	18598	0.3447905
Premium	13791	0.2556730	32389	0.6004635
Ideal	21551	0.3995365	53940	1.0000000

ggplot2 没有专门的饼图函数，可以用 `geom_col()` 与 `coord_polar()` 配合，或者用 `ggforce` 包的 `geom_arc_bar()`。

下面作级别分布的饼图，用饼图中色块大小代表不同级别的比例大小：

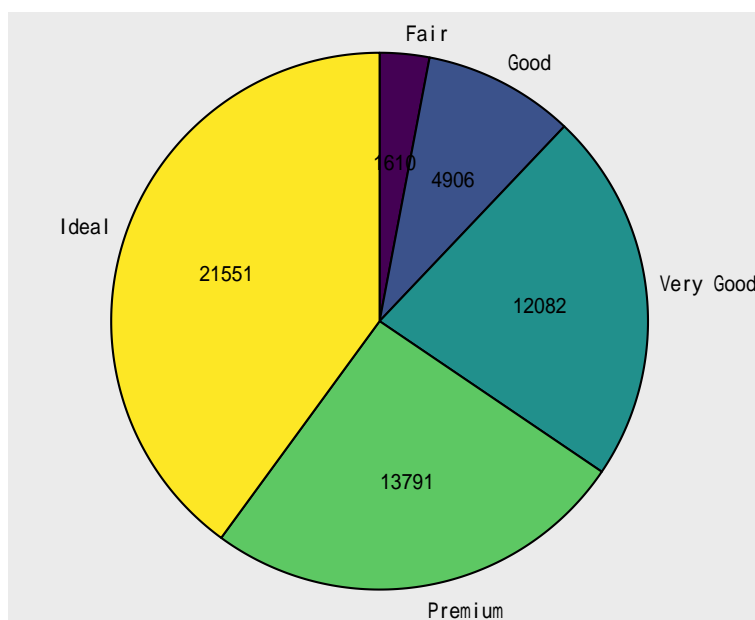
```
p <- ggplot(data = ddstab, mapping = aes(
 x = 1, y = freq, fill = x))
p + geom_col() +
 coord_polar(theta = "y") +
 labs(x = NULL, y = NULL, fill = "Cut")
```



用 `geom_arc_bar()` 函数需要直接指定各个分界线的角度。程序如：

```
library(ggforce)
d <- ddstab %>%
 mutate(
 end_angle = cumpct * 2*pi, # 每块的结束角度
 start_angle = lag(end_angle, default = 0), # 每块的开始角度
 mid_angle = 0.5*(start_angle + end_angle), # 每块的中间角度, 用于频数数值
 hjust = ifelse(mid_angle > pi, 1, 0),
 vjust = ifelse(mid_angle < pi/2 | mid_angle > 3*pi/2, 0, 1)
```

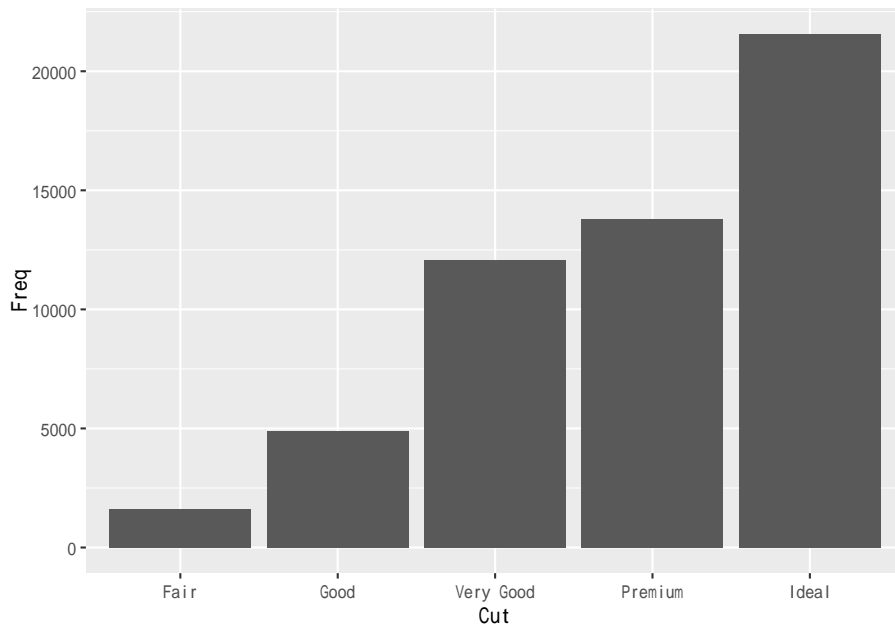
```
)
p <- ggplot(data = d)
p + geom_arc_bar(mapping = aes(
 x0 = 0, y0 = 0, r0 = 0, r = 1.0,
 start = start_angle,
 end = end_angle,
 fill = x)) +
 geom_text(aes(
 x = 1.05*sin(mid_angle),
 y = 1.05*cos(mid_angle),
 label = x,
 hjust = hjust, vjust = vjust)) +
 geom_text(
 aes(
 x = 0.6*sin(mid_angle),
 y = 0.6*cos(mid_angle),
 label = freq)) +
 coord_fixed() +
 scale_x_continuous(expand = c(0.1, 0.18), breaks = NULL, name=NULL) +
 scale_y_continuous(breaks = NULL, name=NULL) +
 guides(fill = FALSE)
```



这样的频数或者比例分布，当然可以用普通的条形图表示，如：

```
p <- ggplot(data = ddstab, mapping = aes(
 x = x, y = freq))
p + geom_col() +
 labs(x = "Cut", y = "Freq")
```





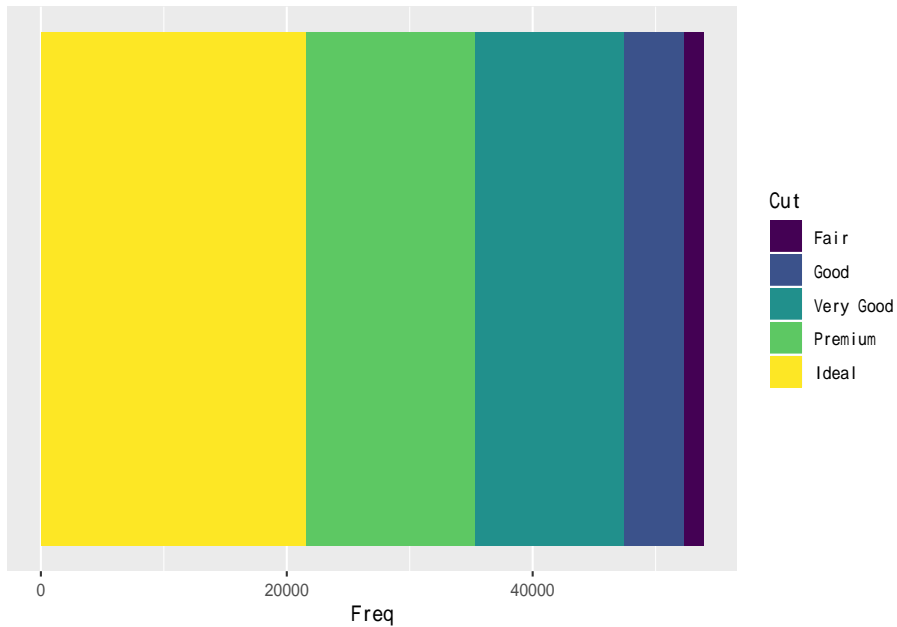
这样的图形比较容易在各类之间比较。也可以做成堆叠的条形图：

```
p <- ggplot(data = ddstab, mapping = aes(
 x = 1, fill = x, y = freq))
p + geom_col(position = "stack") +
 scale_x_continuous(breaks = NULL, name = NULL) +
 labs(fill = "Cut", y = "Freq")
```



也可以横向放置，用 `coord_flip()`:

```
p <- ggplot(data = ddstab, mapping = aes(
 x = 1, fill = x, y = freq))
p + geom_col(position = "stack") +
 scale_x_continuous(breaks = NULL, name = NULL) +
 coord_flip() +
 labs(fill = "Cut", y = "Freq")
```



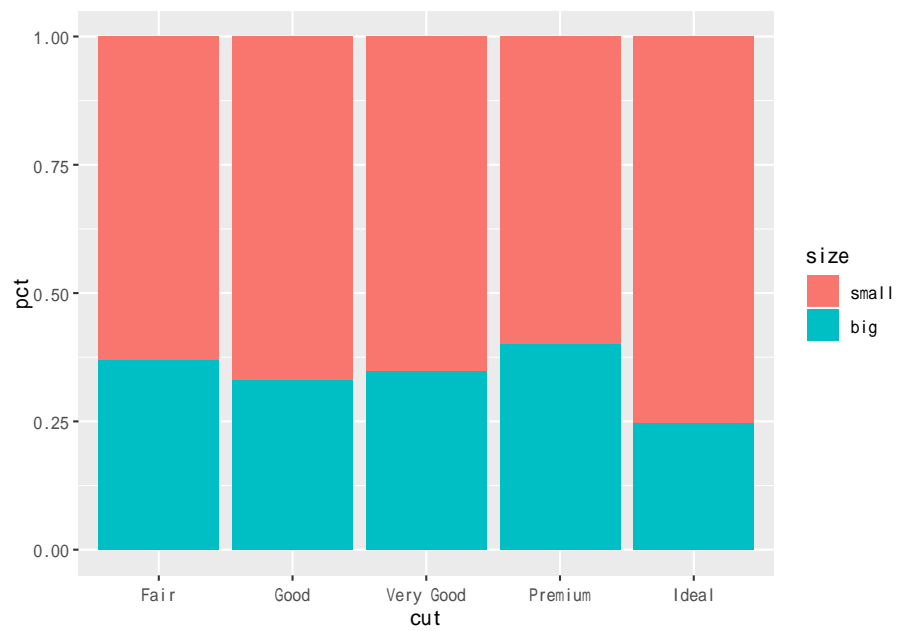
可以用 `geom_text()` 将各个类别标签和频数的数值直接标在色块上。堆叠的条形图比较容易看到每个类占总数的比例大小。

### 31.4.2 组间的比例分布的比较

有时需要比较不同组的某个变量的比例分布。比如，我们将钻石按重量分为大和小两种，在 5 种不同品质之间比较大小比例，可以用堆叠的条形图或者并列的条形图。程序：

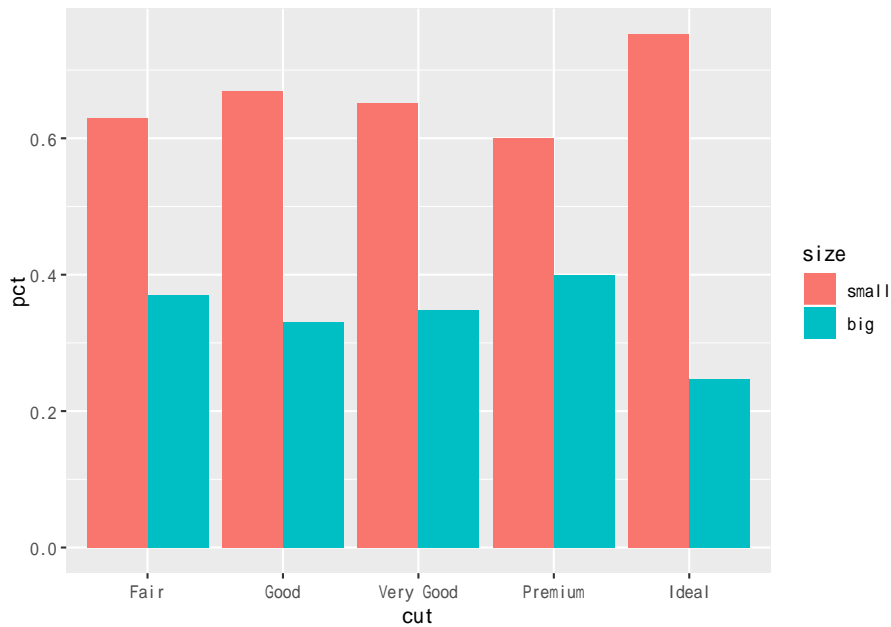
```
ddsmall %>%
 mutate(size = factor(
 ifelse(carat >= 1.0, "big", "small"),
 levels=c("small", "big"))) %>%
 count(cut, size) %>%
 group_by(cut) %>%
 mutate(pct = n / sum(n)) -> ddssize
p <- ggplot(data = ddssize, mapping = aes(
 x = cut, fill = size, y = pct))
```

```
p + geom_col(position = "stack")
```



可以很容易地在 5 个组直接比较两种大小各自的比例。但是，如果有多余两个的比例要比较，则堆叠条形中位于中间的比例不易比较。可以做并列的条形图，如：

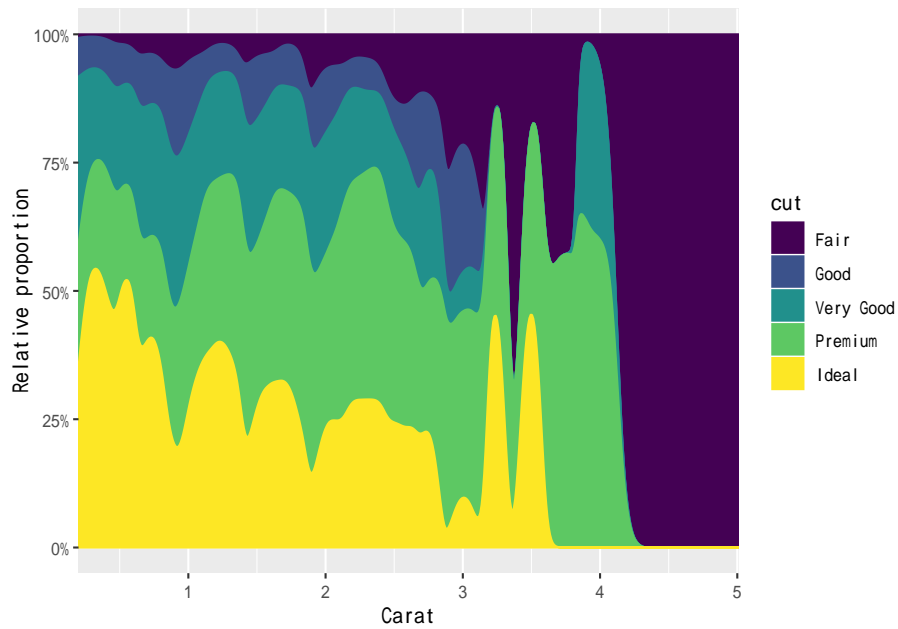
```
p + geom_col(position = "dodge")
```



这样可以比较容易地在不同大类之间比较比例，查看比例的变化趋势。

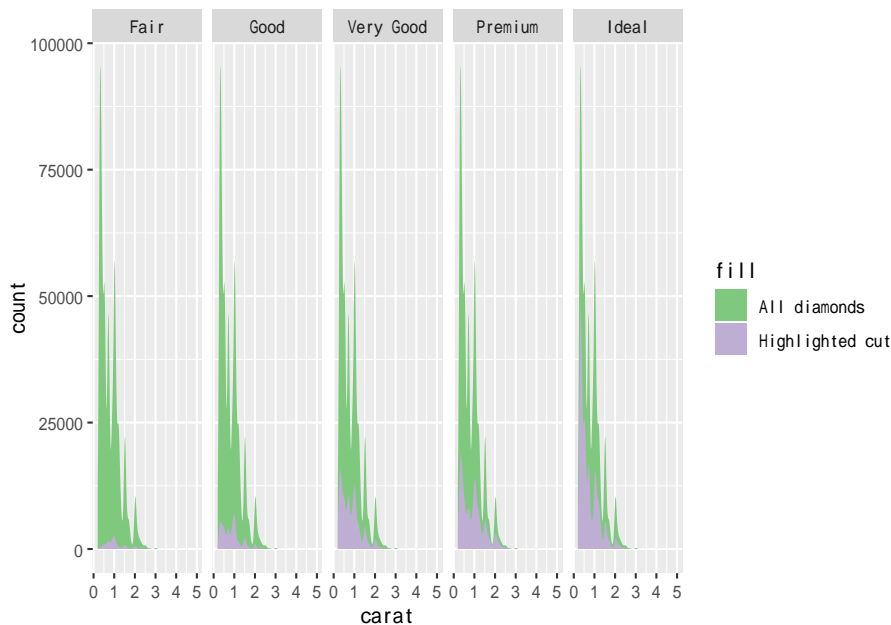
考虑 diamonds 数据集中随重量变化，品质分布的变化。因为重量是连续变化的，所以图形需要变成曲线或者阴影图，每条曲线表示某个品质的累积比例。

```
p <- ggplot(data = diamonds, mapping = aes(
 x = carat, y = ..count.., fill = cut, color = cut))
p + geom_density(position = "fill") +
 scale_x_continuous(name = "Carat", expand=c(0,0)) +
 scale_y_continuous(
 name = "Relative proportion",
 labels = scales::percent)
```



这样的图形仅表现了每个重量级别的品质比例，无法表现不同重量级别的数量。为此，可以每种品质单独画总的重量分布密度图，但在其中用阴影标出该品质所占比例：

```
p <- ggplot(data = diamonds, mapping = aes(
 x = carat, y = ..count..)
) + geom_density_line(data = select(diamonds, -cut), mapping = aes(
 fill = "All diamonds", color = "transparent") +
 geom_density_line(mapping = aes(
 fill = "Highlighted cut", color = "transparent") +
 scale_fill_brewer(type = "qual") +
 facet_wrap(~ cut, nrow = 1)
```



因为程序中用了 `facet_wrap()` 分组，所以密度图（纵轴为个数）本应只有某个品级的数据，但是用了删除 `cut` 变量的方法，将全集数据的密度图也画出来了。这个数据集中，高品质的钻石更多。

### 31.4.3 嵌套比例分布的比较

堆叠条形图、并列条形图都可以用来表现两重分类的列联表频数数据。堆叠条形图中每个条形的高度可以用来表现大类频数，每个条形内色块的大小可以用来表现大类内小类的频数和比例。不同大类之间的小类频数和比例相对来说不易比较。并列条形图中每个条形是一个交叉类的频数，所以在大类之间比较小类频数比较容易，但是每个大类的频数则表现得较为模糊。

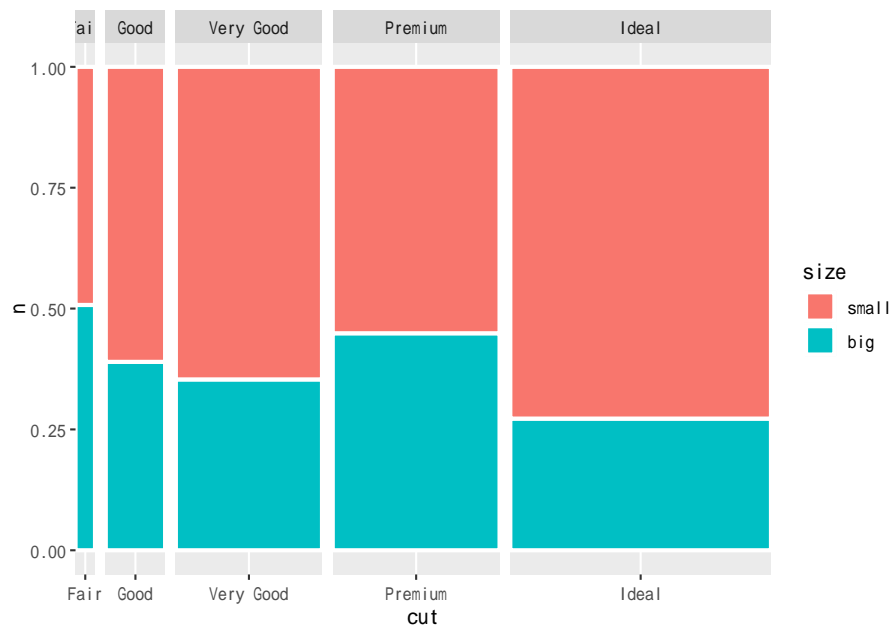
马赛格图是另外一种表现列联表频数的图形，用色块面积表示每个交叉类的频数。与堆叠条形图相比，这里用了条形宽度而不是条形高度表示大类频数。

```
diamonds %>%
 mutate(size = factor(
 ifelse(carat >= 1.0, "big", "small"),
 levels=c("small", "big"))) %>%
```

```

count(cut, size) %>%
group_by(cut) %>%
mutate(cutfreq = sum(n)) %>%
ungroup() -> ddssize0
p <- ggplot(data = ddssize0, mapping = aes(
 x = cut, fill = size, y = n, width = cutfreq))
p + geom_col(position = "fill",
 color="white", size=1) +
facet_grid(~ cut, scales = "free_x", space = "free_x")

```



ggplot2 中没有专门的马赛克图。上面的程序用了 `geom_col()`、`facet_grid()` 和 `width` 参数实现马赛克图。

与马赛克图类似的一种方法是在大类中再细分小类，而且细分小类的办法不是单向划分的。`treemapify` 包的 `geom_treemap()` 函数可以做树状分类图。如：

```

library(treemapify)
p <- ggplot(data = ddssize0, mapping = aes(
 subgroup = cut, fill = interaction(size, cut), area = n))
p + geom_treemap(color="white", size=0.5*.pt, alpha=NA) +

```



```
geom_treemap_subgroup_text(
 place = "center", alpha = 0.5, grow = TRUE) +
geom_treemap_text(mapping = aes(
 label = size),
 color = "white",
 place = "center", grow = FALSE) +
guides(fill = FALSE)
```



还可以用平行集 (parallel sets) 方法表现多个分类之间的关系，以一个分类为主分类染色，可以看出每个主分类与其他子类的关系。ggforce 包的 `geom_parallel_sets()` 函数作这种图。参见<https://serialmentor.com/dataviz/nested-proportions.html>。

### 31.5 表现多个变量间的关系

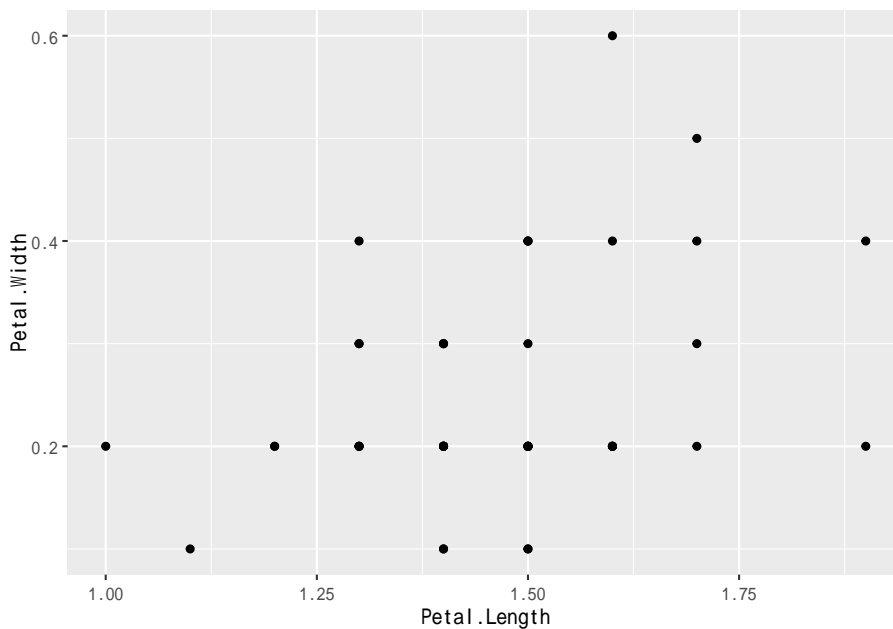
当数据集中有多个变量中，我们除了关心每一个变量的类型、取值集合、分布情况，还关心变量之间的关系，观测的分组情况等。

为表现两个变量之间的关系，最常用的是散点图。多个变量之间可以用散点图矩阵、相关图，可以在散点图中用符号大小、符号颜色、符号形状表示更多维数。对于高维数据，经常需要利用降维方法，如主成分分析 (PCA) 对数据降维，对降维数据作图。

### 31.5.1 散点图

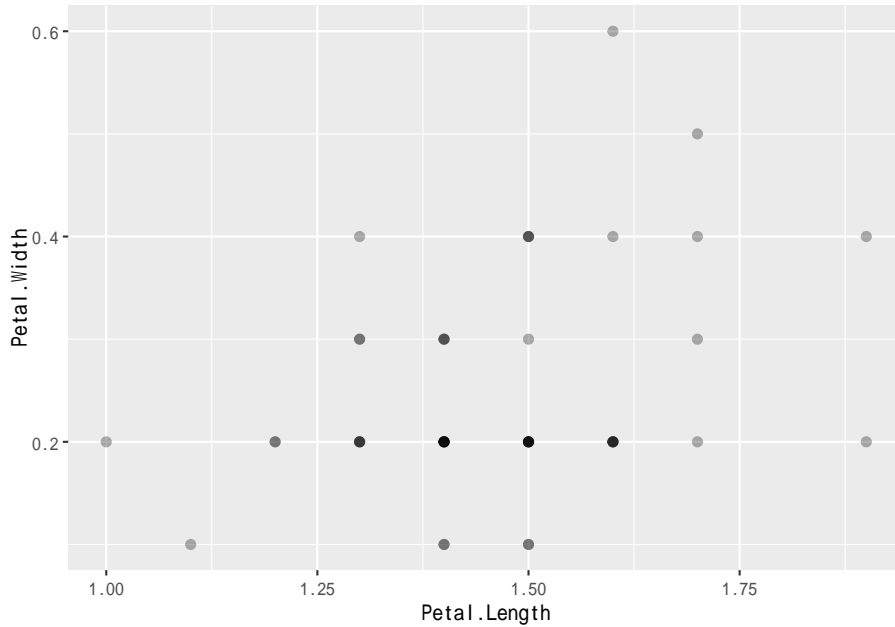
R 软件自带的 `iris` 数据集中包含了三种鸢尾花的 150 个样品的测量数据，每种各 50 个样品，每个样品测量了花瓣、花萼的长、宽。下面画 50 个 `setosa` 样品的花瓣长、宽的散点图，可以看出，两种有明显的线性相关关系：

```
p <- ggplot(
 data = filter(iris, Species == "setosa"),
 mapping = aes(x = Petal.Length, y = Petal.Width))
p + geom_point()
```



应该有 50 个点，但实际只看到 23 个点。这是因为许多点重叠在一起了。加上透明度参数可以使得重叠的点颜色更深：

```
p + geom_point(alpha = 0.3, size = 2.0)
```

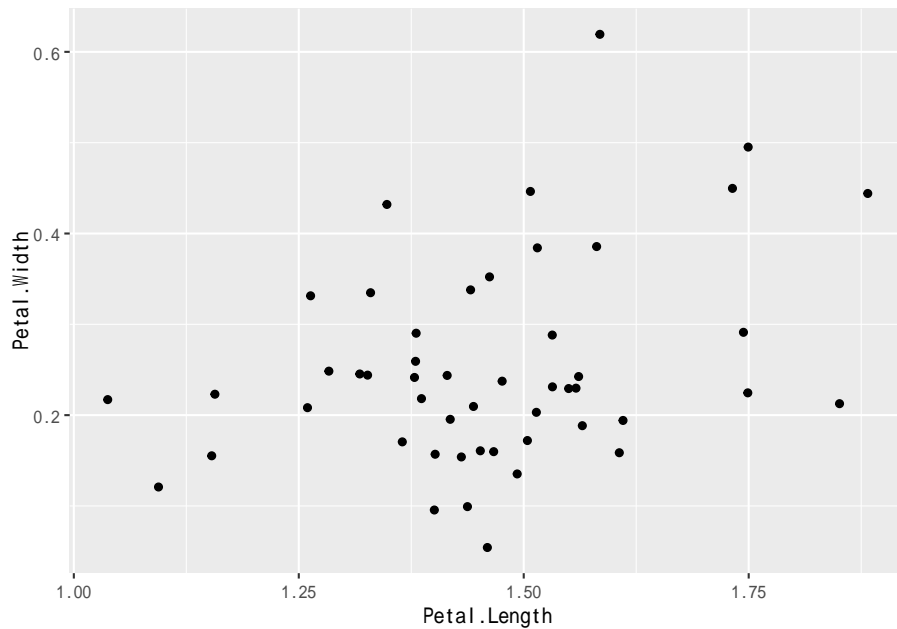


图的效果还是不够明显，不太好判断重叠了多少个点，另外如果不同颜色的点重叠，就无法判断有哪些颜色的点重叠在一起。如果不做说明，读者也不一定了解颜色深浅代表重叠点的多少。

当点数很多，如数百、数千个点时，可以用 `size` 参数画更小的点，这时使用透明度的效果会比较显著。

可以使用 `geom_jitter()`，使得每个点略有随机偏移，就可以使得各个点基本上不重叠。偏移可以是上下左右同时进行，也可以要求仅上下或者仅左右偏移。例如：

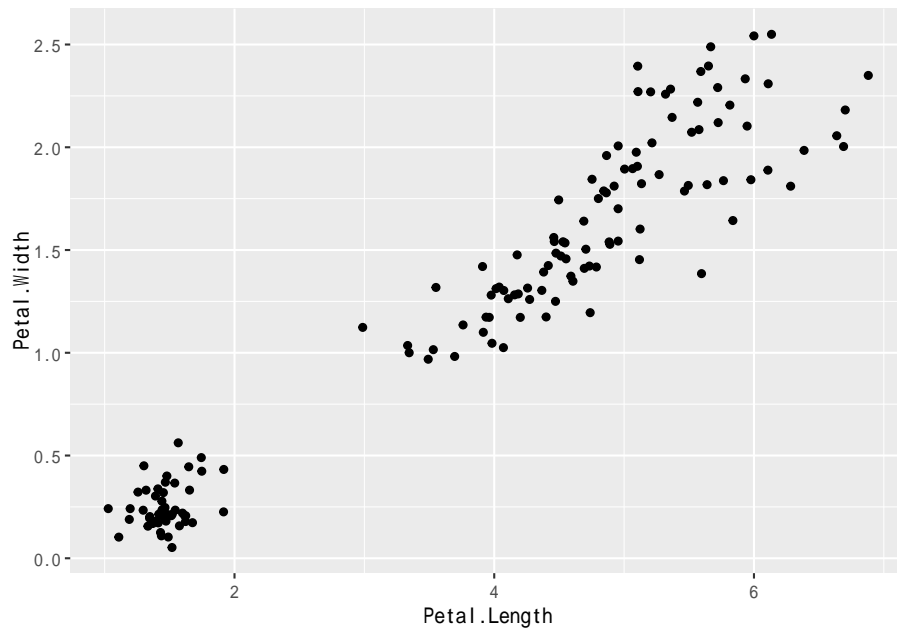
```
p + geom_jitter(width=0.05, height=0.05)
```



`geom_jitter()` 中的 `width` 和 `height` 是左右和上下抖动的幅度百分比，以数据间的最小差距为单位，默认的抖动比例是 40%，一般不超过 50%，否则原来是两个不同坐标值的点可能会看起来是同一坐标值。比如，假设 `x` 变量的值是精确到小数点后 1 位的，取值如 1.3, 1.4，抖动 40%，坐标可以变成 1.33, 1.46，如果允许抖动 60%，则 1.3 可能变成 1.36，1.4 可能变成 1.34，就无法区分原来不同的值了。

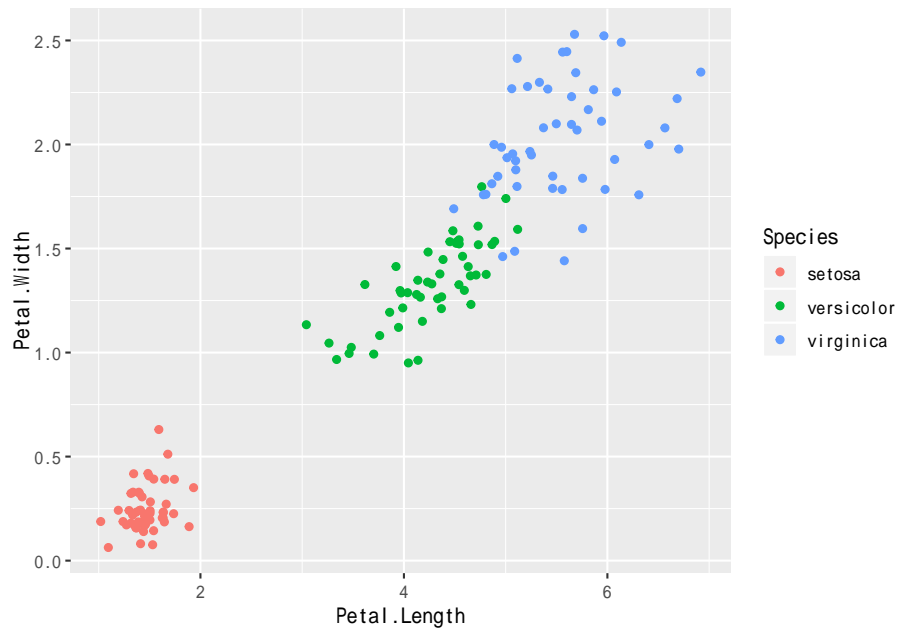
如果做所有的 150 个样品的散点图，则三点呈现出分组现象：

```
p <- ggplot(data = iris, mapping = aes(
 x = Petal.Length, y = Petal.Width))
p + geom_jitter(width=0.05, height=0.05)
```



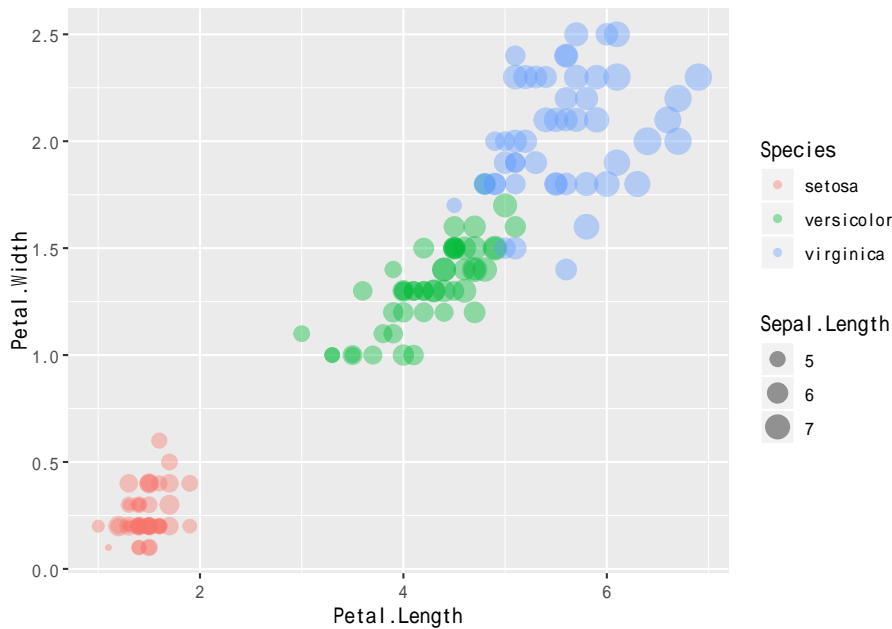
为此,应该设法在散点图中区分三种不同鸢尾花品种。将 `Species` 映射到 `color` 属性即可区分:

```
p <- ggplot(data = iris, mapping = aes(
 x = Petal.Length, y = Petal.Width,
 color = Species))
p + geom_jitter(width=0.05, height=0.05)
```



如果希望将花萼长度也添加到图形中，因为已经有 x 维、y 维和颜色，只能通过符号大小来添加。将花萼长度映射到 `size` 维：

```
p <- ggplot(data = iris, mapping = aes(
 x = Petal.Length, y = Petal.Width,
 size = Sepal.Length,
 color = Species))
p + geom_point(alpha = 0.4)
```



这样的用符号大小代表一个变量数据的图形称为气泡图 (bubble plot)。符号大小是比较难认读的图形刻度。

多个变量之间的关系经常用散点图矩阵表示。ggplot2 包没有提供专门的散点图矩阵，可以通过数据变换、散点图函数和 `facet_grid()` 函数编程实现。基础 R 图形中提供了 `pairs` 函数作散点图矩阵。

### 31.5.2 相关系数矩阵图

多个变量之间的相关系数矩阵可以用色块图表示，ggplot2 包没有提供专门的函数，可以用数据变换和 `geom_tile()` 函数实现。corrgram 包提供了 `corrgram` 函数作相关系数矩阵图，见29.6.1。

### 31.5.3 数据降维

当数据中有过多的变量时，比如变量达到数百、数千时，即使是散点图矩阵都信息量过大难以认读。这时常常用主成分分析等降维方法将数据降维到若干个新变量，对新变量作图。

主成分分析利用变量的协方差阵或者相关阵的特征值分解对原始变量进行线性组合，产生若干个新变量。当原始变量为相同单位且可比时，可以基于协方差阵，否则应该基于相关阵。例如，对 `iris` 数据集的 4 个测量值作主成分分析：

```
pca1 <- princomp(iris[,1:4], cor=FALSE)
summary(pca1)
```

```
Importance of components:
##
Comp.1 Comp.2 Comp.3 Comp.4
Standard deviation 2.0494032 0.49097143 0.27872586 0.153870700
Proportion of Variance 0.9246187 0.05306648 0.01710261 0.005212184
Cumulative Proportion 0.9246187 0.97768521 0.99478782 1.000000000
```

```
load1 <- loadings(pca1)
print(load1)
```

```
##
Loadings:
Comp.1 Comp.2 Comp.3 Comp.4
Sepal.Length 0.361 0.657 0.582 0.315
Sepal.Width 0.730 -0.598 -0.320
Petal.Length 0.857 -0.173 -0.480
Petal.Width 0.358 -0.546 0.754
##
Comp.1 Comp.2 Comp.3 Comp.4
SS loadings 1.00 1.00 1.00 1.00
Proportion Var 0.25 0.25 0.25 0.25
Cumulative Var 0.25 0.50 0.75 1.00
```

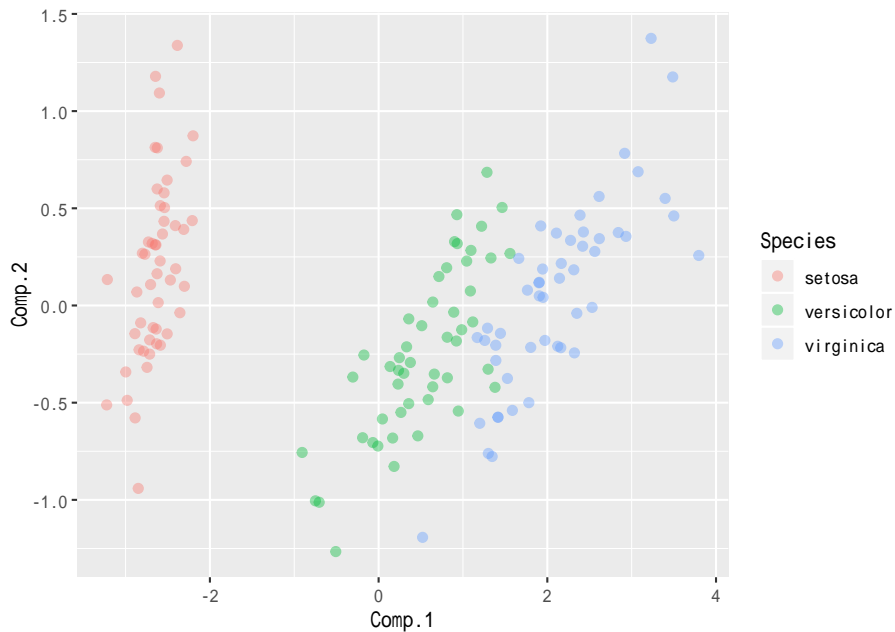
```
ps1 <- predict(pca1)
ps1 <- as.data.frame(ps1)
ps1 <- cbind(iris, ps1)
```

前两个主成了解释了原始变量中 98% 的方差。

作第一和第二主成分的散点图：

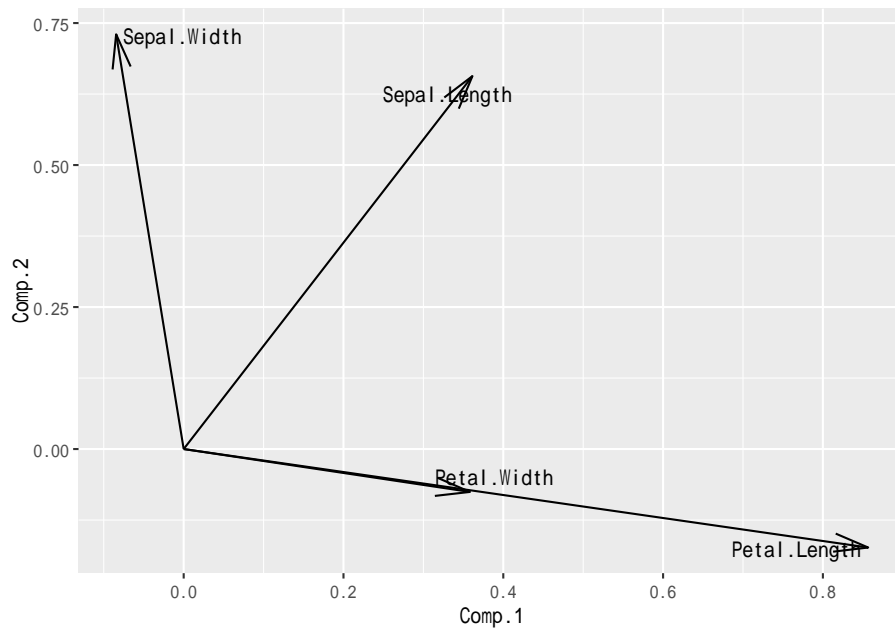


```
p <- ggplot(data = as.data.frame(ps1), mapping = aes(
 x = Comp.1, y = Comp.2, color = Species))
p + geom_point(size = 2.0, alpha = 0.4)
```



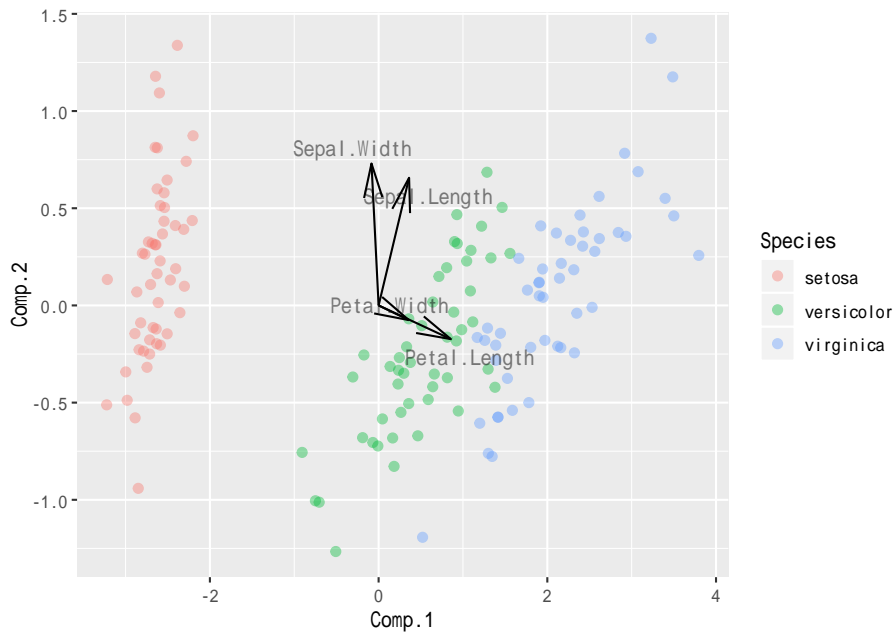
主成分分析取前两个主成分在各原始变量上的载荷可以做矢量图，用来表示各个原始变量与主成分的关系：

```
d <- as.data.frame(load1[,1:2])
d$vlabel <- rownames(d)
p <- ggplot(data = d, mapping = aes(
 x = Comp.1, y = Comp.2, label = vlabel))
p + geom_segment(mapping = aes(
 xend = Comp.1, yend = Comp.2),
 x = 0, y = 0,
 arrow = arrow(angle = 15)) +
 geom_text_repel()
```



可以将降维的散点图与变量的载荷图画在同一坐标系内：

```
d <- as.data.frame(load1[,1:2])
d$vlabel <- rownames(d)
ggplot() +
 geom_point(data = ps1, mapping = aes(
 x = Comp.1, y = Comp.2, color = Species),
 size = 2.0, alpha = 0.4) +
 geom_segment(data = d, mapping = aes(
 xend = Comp.1, yend = Comp.2),
 x = 0, y = 0,
 arrow = arrow(angle = 15)) +
 geom_text_repel(data = d, mapping = aes(
 x = Comp.1, y = Comp.2, label = vlabel),
 alpha = 0.5)
```

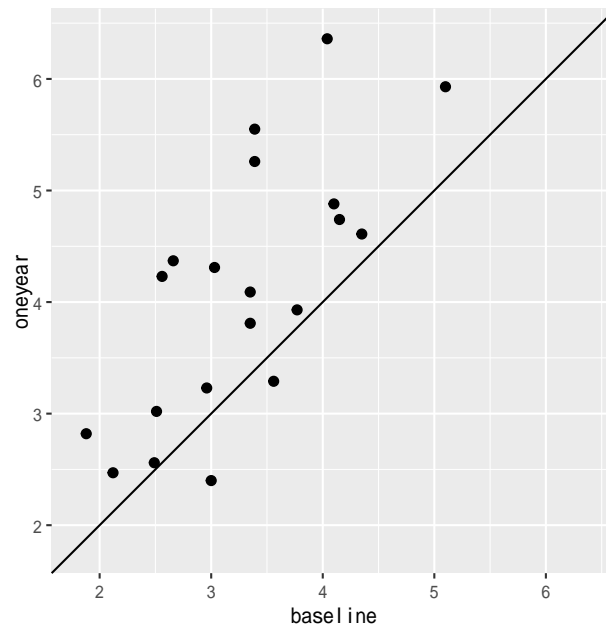


### 31.5.4 成对数据

成对数据指同一个变量在两个不同时间的测量值，或者同一个体两个取值基本相同的变量。作散点图时，增加  $y = x$  直线是显然的一个增强显示。另外，还应该强调两个变量之间的差别。

考虑 boot 扩展包的 cd4 数据集，这个数据集包含了 20 名临床试验受试者在开始时和一年以后的 cd4 指标测量值。作这两个指标的散点图：

```
data(cd4, package = "boot")
p <- ggplot(data = cd4, mapping = aes(
 x = baseline, y = oneyear))
p + geom_point(size = 2.0) +
 geom_abline(slope = 1, intercept = 0) +
 coord_fixed(xlim = c(1.80, 6.40), ylim = c(1.80, 6.40))
```



程序中用了 `geom_abline()` 作斜线，用了 `coord_fixed()` 指定 1:1 的宽高比并设定 x 轴和 y 轴相同的坐标范围。可以看出一年后的测量值普遍高于开始时的测量值。

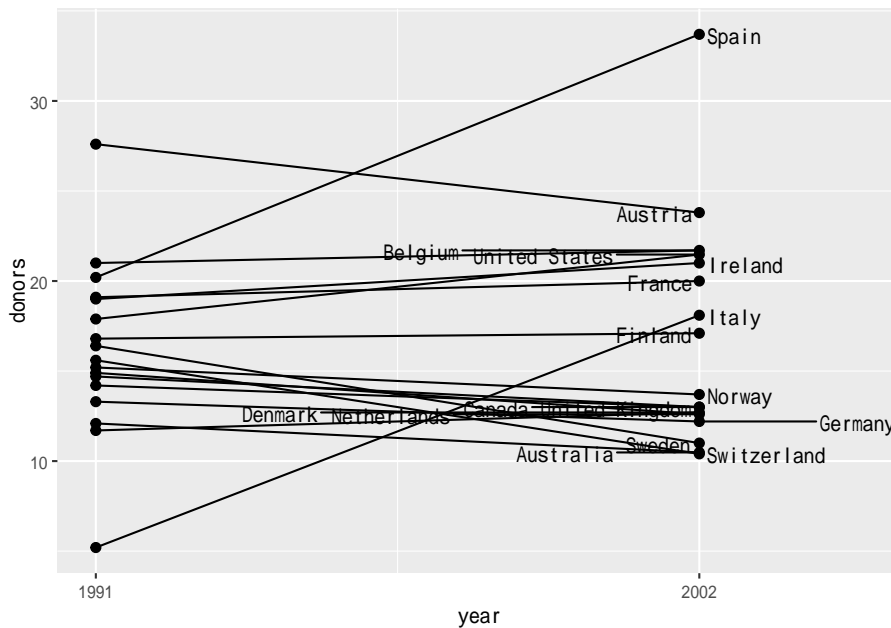
当观测个数较少时，可以用线段连接每个观测的两个变量值，并标出每个观测的标签。

```
data(organdata, package = "socviz")
organdata %>%
 select(country, year, donors) %>%
 mutate(year = lubridate::year(year)) %>%
 filter(year %in% c(1991, 2002)) -> dorgans
p <- ggplot(data = dorgans, mapping = aes(
 x = year, y = donors))
p + geom_line(aes(group = country)) +
 geom_point(size = 2.0) +
 geom_text_repel(data = filter(dorgans, year == 2002),
 mapping = aes(
 x = year, y = donors,
```

```

 label = country),
 direction = "x") +
scale_x_continuous(
 breaks = c(1991, 2002),
 limits = c(1991, 2005),
 labels = c("1991", "2002"))

```



可以看出，这种图可以表现超过两个变量。

## 31.6 时间序列图

某个指标随时间变化的数据称为时间序列，可以以时间为横坐标、该指标为纵坐标作折线图，称为时间序列图。其它有次序的变量也可以用作横轴变量。

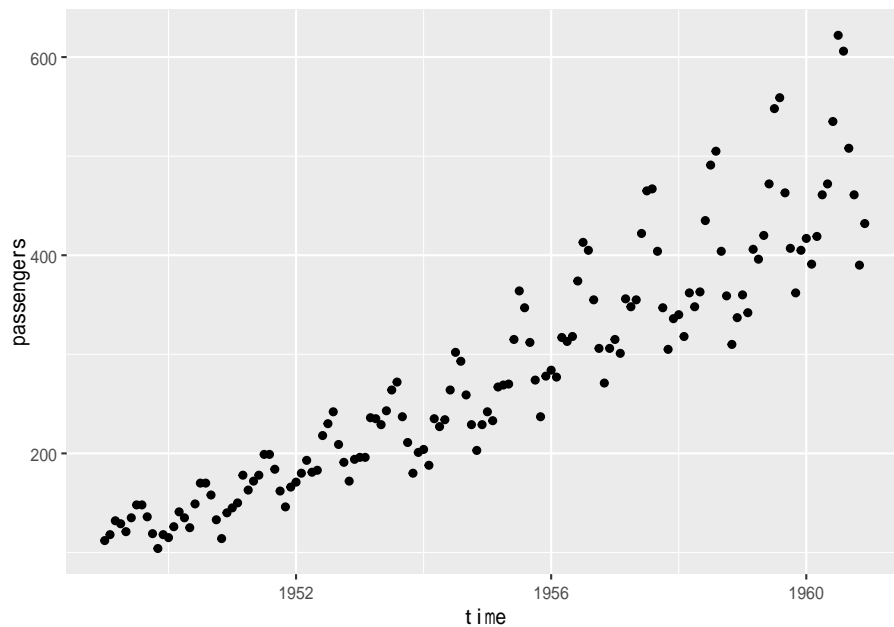
`geom_line()` 函数会自动按照 `x` 轴变量的次序连线，即使数据中的 `x` 变量没有从小到大排列都是如此。`geom_path()` 则按照数据集中原有次序连线。

### 31.6.1 一元时间序列

R 软件中提供了一个示例用的时间序列数据 `AirPassengers`，是美国泛美航空公司 1949-1960 的国际航班订票数的月度数据（单位：千人），12 年 144 个月的数据。

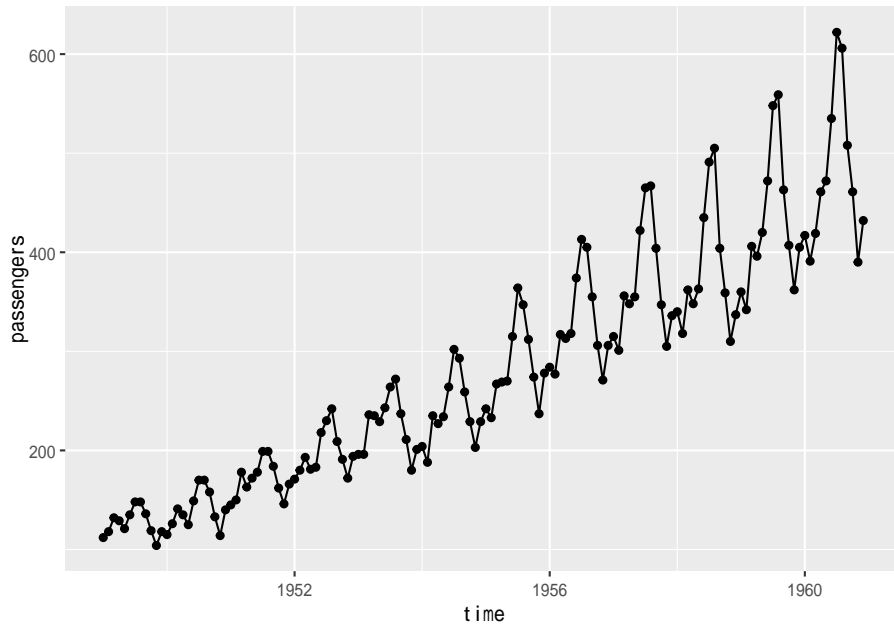
作其散点图：

```
dap <- data.frame(
 time = c(time(AirPassengers)),
 passengers = c(AirPassengers))
p <- ggplot(data = dap, mapping = aes(
 x = time, y = passengers))
p + geom_point()
```



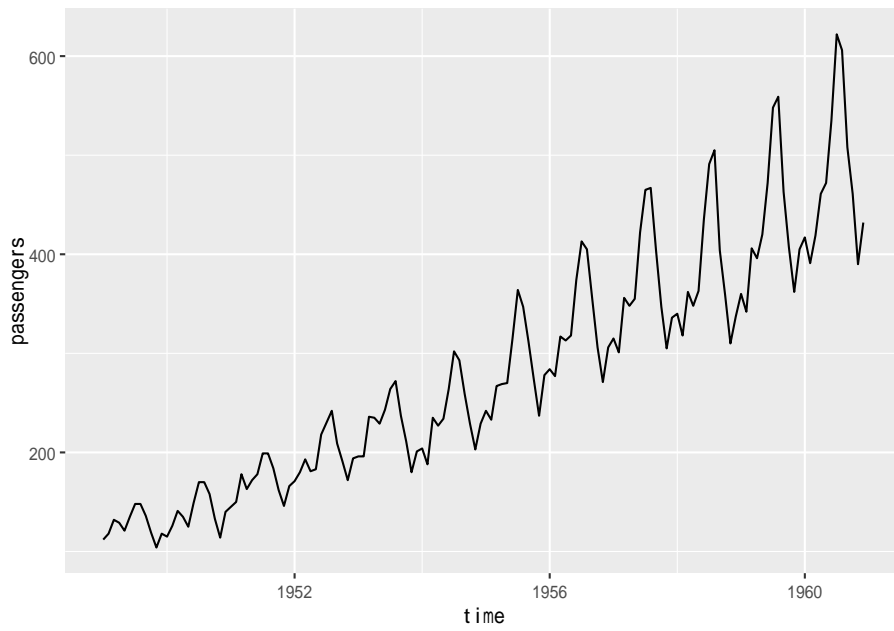
这样的图形在点数较少时还好，点数较多时不易分辨出局部随时间的变化情况。  
加上连线：

```
p + geom_point() +
 geom_line()
```



也可以仅连续，不画散点：

```
p + geom_line()
```

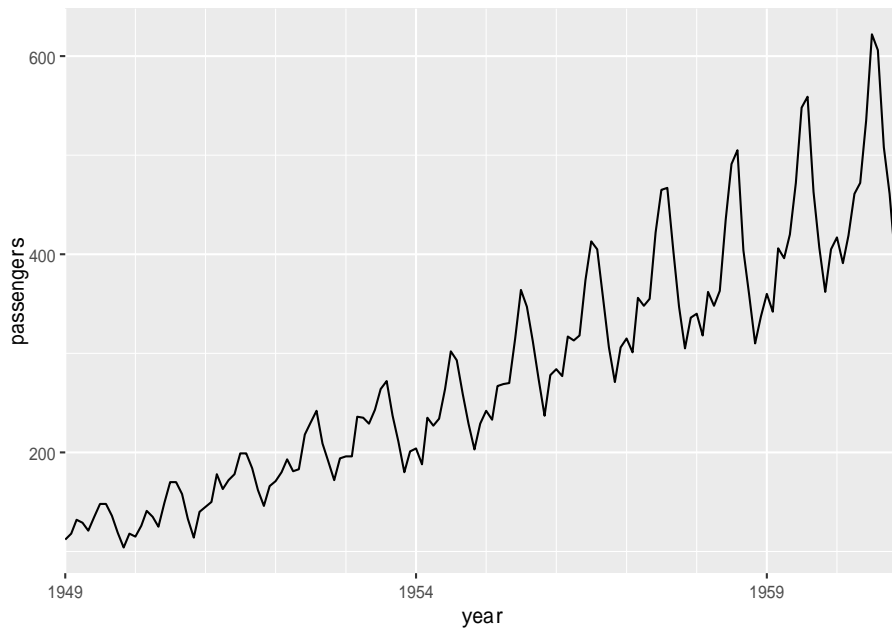


仅画线的时间序列图更强调变化趋势，不强调个别数据，当观测点很多时一般不画出散点。

这里的时间是 `time()` 函数返回的以年为单位的浮点数值。时间序列的时间往往是日期值，尤其是日数据，可以用 `scale_x_date()` 函数标日期，如：

```
library(xts)
xts_ap <- as.xts(AirPassengers)
dap2 <- data.frame(
 time = index(xts_ap),
 passengers = coredata(xts_ap))
p <- ggplot(data = dap2, mapping = aes(
 x = as.Date(time), y = passengers))
p + geom_line() +
 scale_x_date(
 name = "year",
 date_breaks = "5 years",
 date_labels = "%Y",
 date_minor_breaks = "1 year",
 expand = c(0, 0))
```

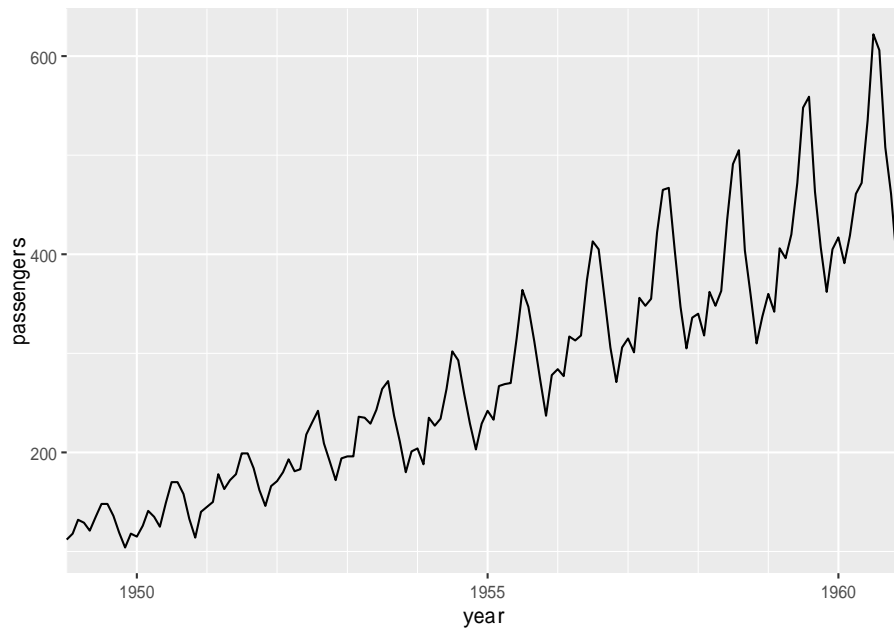




程序中的 `time` 变量是 `POSIXct` 类型的日期时间,用 `as.Date()` 将其转换成了 `Date` 类型以适应 `scale_x_time()` 的要求。用 `date_breaks` 选项指定日期轴刻度值的频率,“5 years”就是每 5 年一个刻度值。用 `date_minor_breaks` 选项指定细刻度的频率。`date_labels` 是日期刻度的一种格式规定,比如“%Y”是四位数的年份,“%Y-%m”是年月格式,等等。因为是月度数据,所以刻度线所在的数据点对应该年的 1 月份。

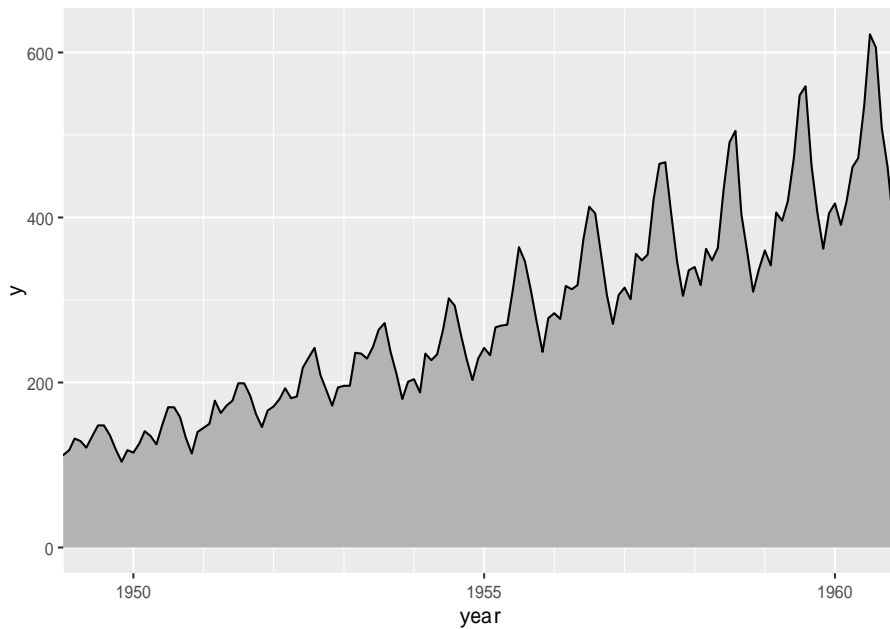
如果需要人为指定刻度位置,可以用 `breaks` 参数,如:

```
p + geom_line() +
 scale_x_date(
 name = "year",
 breaks = lubridate::make_date(c(1950, 1955, 1960)),
 date_labels = "%Y",
 date_minor_breaks = "1 year",
 expand = c(0, 0))
```



还可以在线下加阴影，但这时 y 轴应为正值变量且坐标范围从 0 开始，如：

```
library(ggribes)
p <- ggplot(data = dap2, mapping = aes(
 x = as.Date(time), height = passengers, y = 0))
p + geom_ridgeline() +
 scale_x_date(
 name = "year",
 breaks = lubridate::make_date(c(1950, 1955, 1960)),
 date_labels = "%Y",
 date_minor_breaks = "1 year",
 expand = c(0, 0))
```



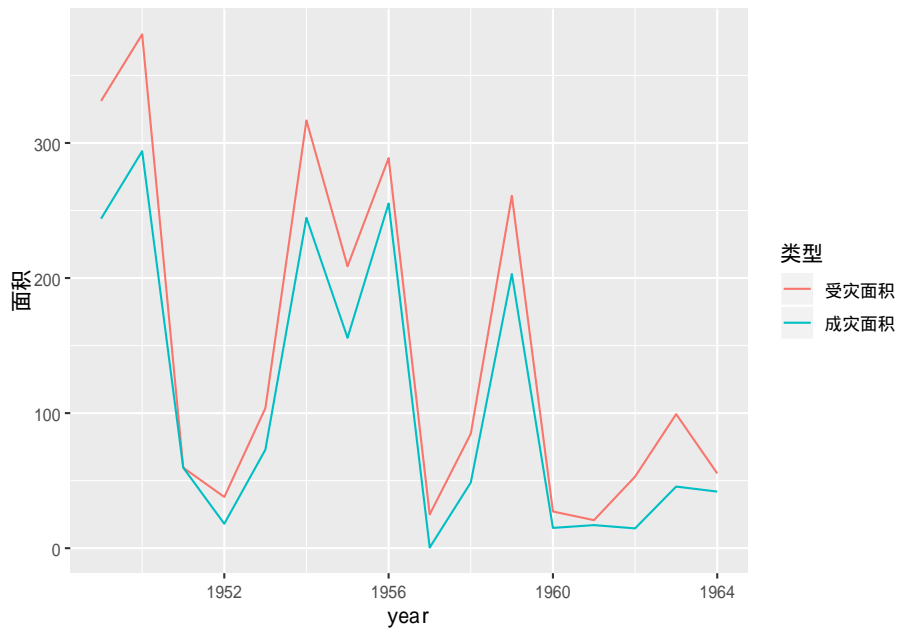
上面的程序用了 `ggridges` 包的 `geom_ridgeline()` 函数, `y` 轴对应于 `height` 维。

### 31.6.2 多元时间序列

随时间记录的指标可以有多个。当这些指标的取值范围相近时, 可以共用一个 `y` 轴, 每个指标单独做连线并画在同一坐标系中。要注意的是, `ggplot` 作多元时间序列图时需要长表数据而非宽表数据, 纵坐标变量必须保存在数据集的一个变量中, 用另外的变量区分不同的指标。例如, 北京地区从 1949 年到 1964 年的受灾面积与成灾面积 (单位: 万亩) 年数据。

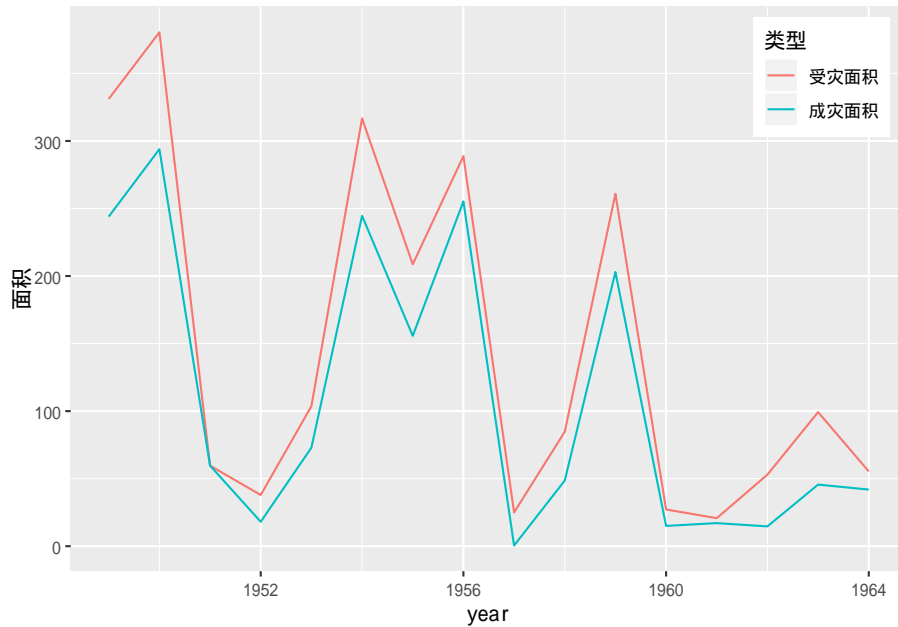
```
dflood <- as.data.frame(matrix(c(
 1949 , 1 , 331.12 , 243.96 ,
 1950 , 2 , 380.44 , 293.90 ,
 1951 , 3 , 59.63 , 59.63 ,
 1952 , 4 , 37.89 , 18.09 ,
 1953 , 5 , 103.66 , 72.92 ,
 1954 , 6 , 316.67 , 244.57 ,
```

```
1955 , 7 , 208.72 , 155.77 ,
1956 , 8 , 288.79 , 255.22 ,
1957 , 9 , 25.00 , 0.50 ,
1958 , 10 , 84.72 , 48.59 ,
1959 , 11 , 260.89 , 202.96 ,
1960 , 12 , 27.18 , 15.02 ,
1961 , 13 , 20.74 , 17.09 ,
1962 , 14 , 52.99 , 14.66 ,
1963 , 15 , 99.25 , 45.61 ,
1964 , 16 , 55.36 , 41.90 ,
byrow=TRUE, ncol=4,
dimnames=list(1949:1964,
c("year", "t", "area1", "area2"))))
dflood2 <- dflood %>%
 gather(area1, area2, key = " 类型", value = " 面积") %>%
 mutate(`类型` = factor(
 `类型`, levels = c("area1", "area2"),
 labels = c(" 受灾面积", " 成灾面积")))
p <- ggplot(data = dflood2, mapping = aes(
 x = year, y = `面积`, color = `类型`))
p + geom_line()
```



可以用 `guides(label.position)` 指定 "top", "bottom", "left", "right" 这样的图例位置，还可以用 `theme()` 函数将图例指定在坐标系内部：

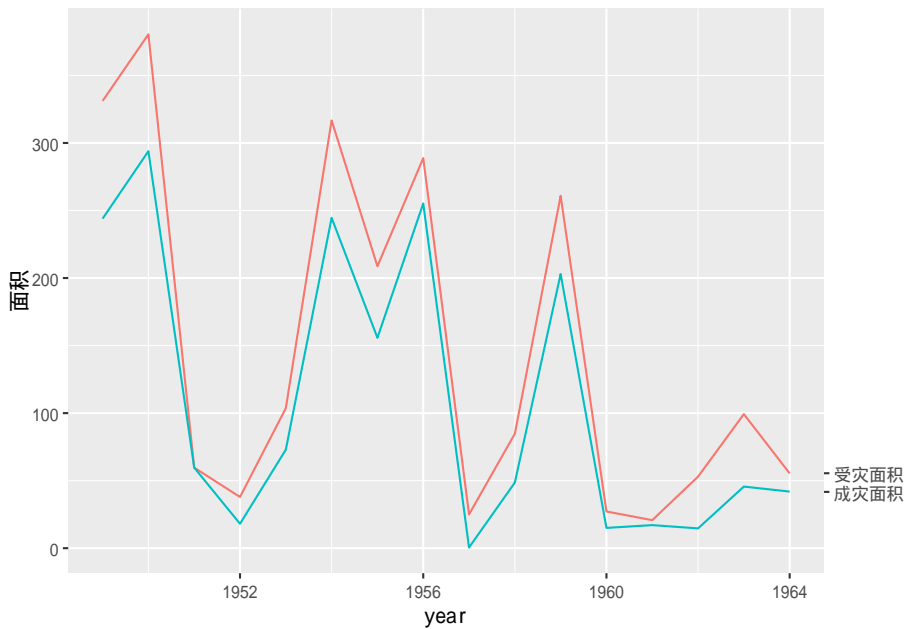
```
p + geom_line() +
 theme(
 legend.position = c(0.98, 0.98),
 legend.justification = c(1,1))
```



`theme()` 函数中 `legend.position` 给出的是图例框在坐标系内的百分比位置, `legend.justification` 给出的是图例框坐标的百分比对齐位置, `c(1,1)` 是将图例框的右上角对齐到给定百分比坐标位置。

也可以去掉图例, 而是用 `y` 轴的第二坐标轴 (画在右侧) 的办法直接标出两个序列的标签。如:

```
p + geom_line() +
 guides(color = FALSE) +
 scale_y_continuous(
 sec.axis = dup_axis(
 name = NULL,
 breaks = subset(dflood2, year == 1964)[["面积"]],
 labels = c("受灾面积", "成灾面积")))
```

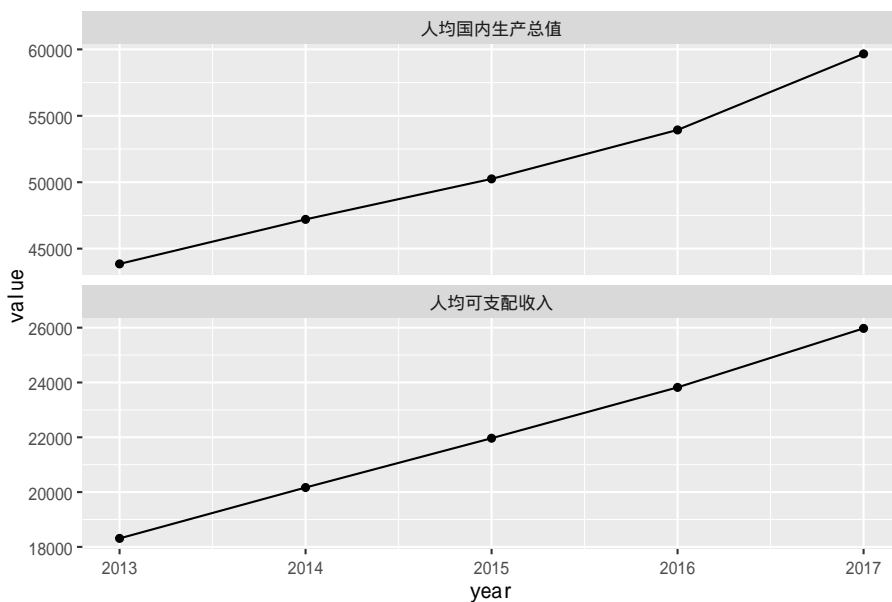


### 31.6.3 多个不同类型指标的时间序列

有时沿时间变化的多个指标是不可比的，这时就不应将数据画在同一坐标系内，一般可以画多个小图并上下按时间轴对齐，`facet_wrap()` 可以做小图。

例如，考虑我国 2013-2017 年人均国内生产总值和人均可支配收入数据。

```
dinc <- data.frame(
 year = 2013:2017,
 `人均可支配收入` = c(18311, 20167, 21966, 23821, 25974),
 `人均国内生产总值` = c(43852, 47203, 50251, 53935, 59660))
dinc2 <- dinc %>%
 gather(`人均可支配收入`, `人均国内生产总值`, key = "指标", value = "value")
p <- ggplot(data = dinc2, mapping = aes(
 x = year, y = value))
p + geom_line() + geom_point() +
 facet_wrap(~ `指标`, ncol = 1, scales = "free_y") +
 labs(caption = "数据来源: 国家统计局")
```



数据来源：国家统计局

上述程序在 `facet_wrap()` 中用 `scales = "free_y"` 使得两幅小图的纵轴不需要使用相同坐标范围。横轴默认还使用相同坐标范围并上下对齐。

对二元时间序列  $(x_t, y_t)$ ，还可以画以  $(x_t, y_t)$  为节点的轨迹图，但需要标出关键的时间点的时间。用 `geom_path()` 画轨迹，用 `ggrepel` 包的 `geom_text_repel()` 标时间。

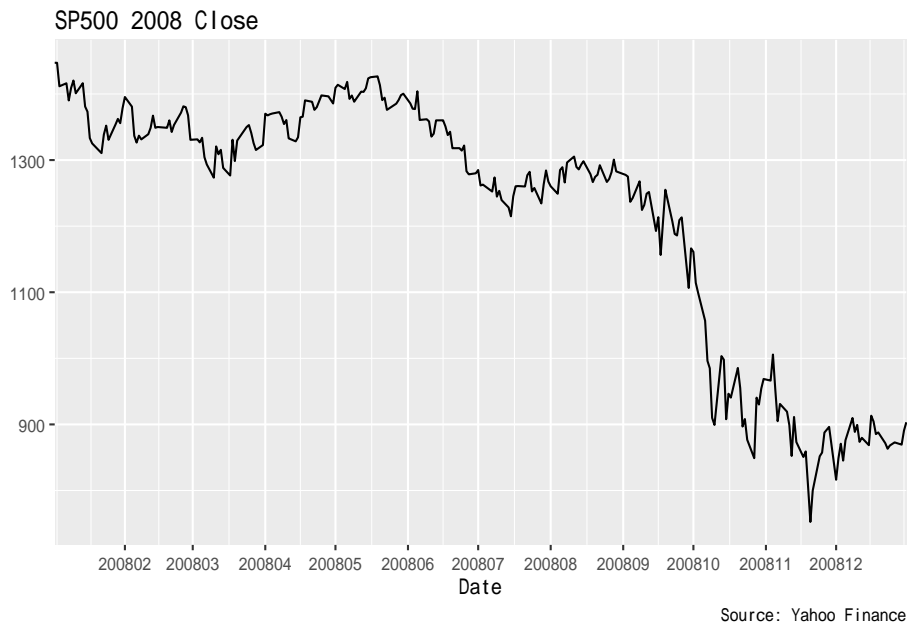
## 31.7 拟合曲线图

考虑两个变量之间的关系，将其中的一个作为因变量，另一个作为自变量，散点图可以反映其互相影响，可以用拟合曲线进一步概括两者的关系。拟合曲线的方法可以是曲线平滑方法如核回归、局部多项式回归、样条平滑，也可以用一些参数模型如线性回归、二次回归。对于时间序列，还经常分解成趋势项、季节项和不规则项。

### 31.7.1 平滑方法

考虑标普 500 指数 2008 年的收盘价。





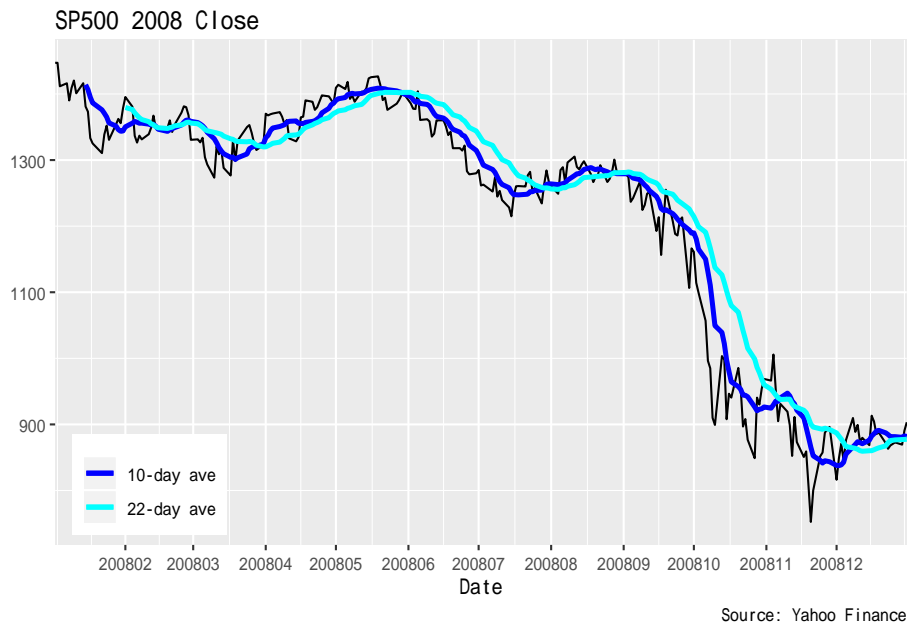
为了平滑上面的曲线，金融界常用滑动平均方法，比如，每 10 天画计算一个平均值并将该平均值的时间对准在 10 天时间窗的最后一天。这样得到的滑动平均线会滞后于原始数据的变化。

先写一个计算滑动平均的函数：

计算 10 日、22 日滑动平均并绘图：

```
Warning: Removed 9 rows containing missing values (geom_path).
```

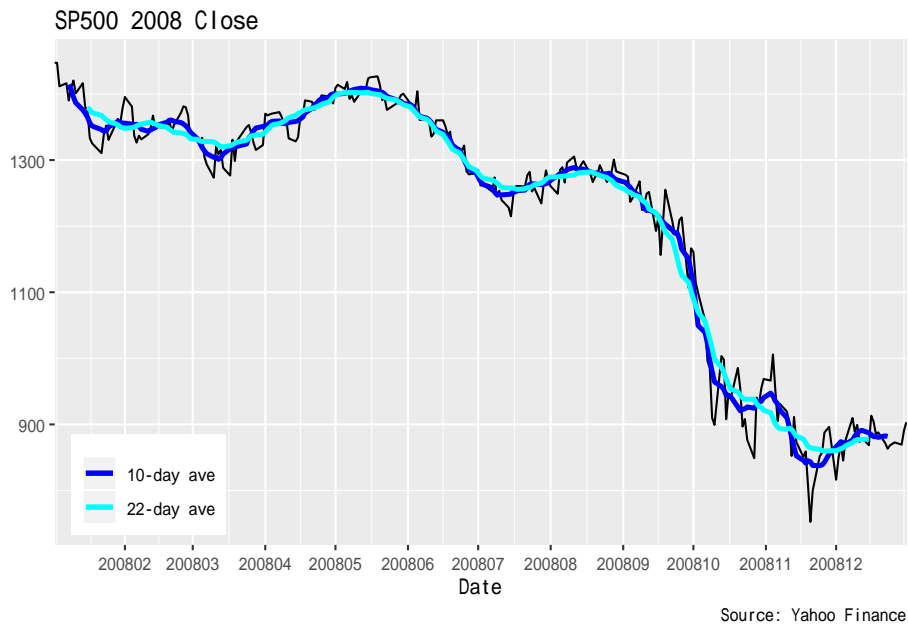
```
Warning: Removed 21 rows containing missing values (geom_path).
```



将滑动平均结果对齐在时间窗口末尾的优点是可以实时地更新。对历史数据，可以将滑动平均结果对齐在时间窗口中间，如果时间窗口中有偶数个点，可以选择中间两个点的任何一个，上面的滑动平均计算程序在偶数个点时选择与中间两个点当中第一个点对齐。图形如下：

```
Warning: Removed 9 rows containing missing values (geom_path).
```

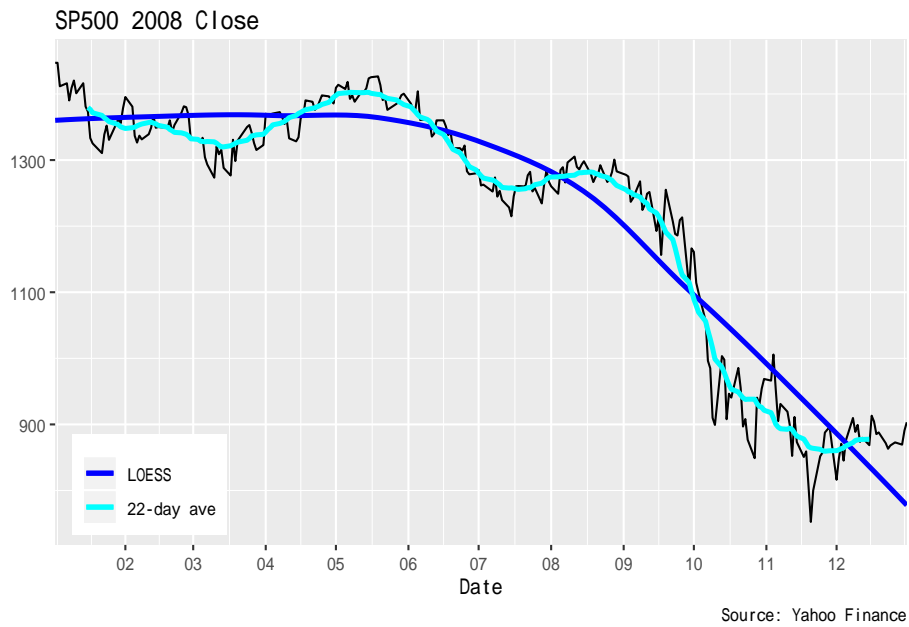
```
Warning: Removed 21 rows containing missing values (geom_path).
```



滑动平均平滑是金融数据分析中常用的方法，但是也有明显的缺点。因为需要多个点生成一个平均值，所以在序列的开头和末尾有一部分点是没有平均值的。当时间窗变宽时结果会更光滑，但有时还是不够光滑，这与滑动平均作为加窗平滑是等权而且在边缘处为间断有关。

统计学中提出了各种更复杂但效果更好而平滑方法，比如 LOESS(locally estimated scatterplot smoothing, 参见 (Cleveland, 1979))。这是许多局部多项式回归方法中的一个典型代表，每个  $x$  处用附近的点拟合一个局部模型，距离  $x$  越远的点贡献越小。这样的做法可以得到比滑动平均更平滑的结果。局部多项式方法得到的结果类似于中心对准而且滑动窗口较宽的滑动平均。

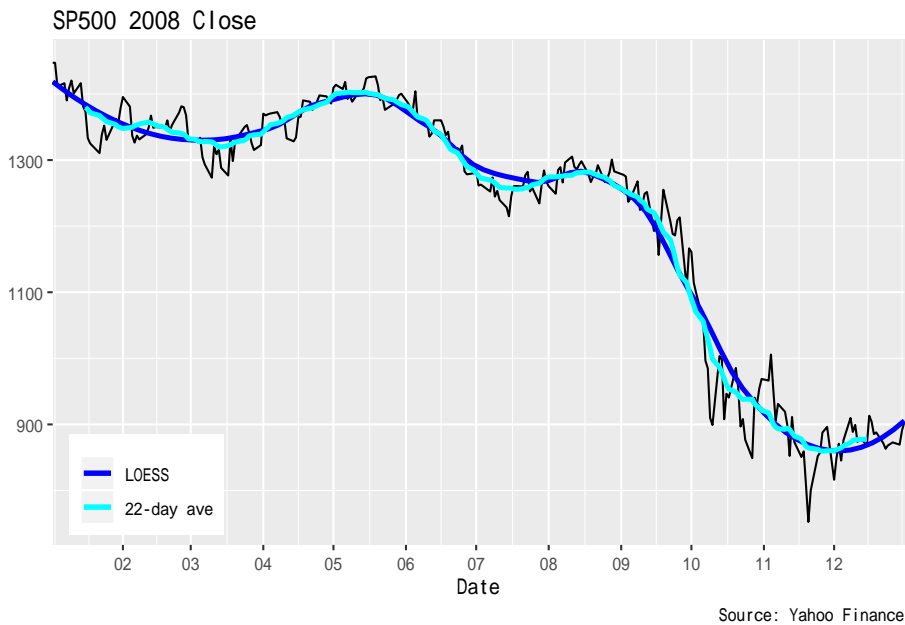
```
Warning: Removed 21 rows containing missing values (geom_path).
```



`geom_smooth()` 可以对输入的  $(x, y)$  坐标数据进行平滑, 选项 `se = FALSE` 要求不做置信区间。实际上, 不仅仅是时间序列数据可以添加如此的平滑曲线, 其它的散点图也都可以这样添加平滑曲线。

上图的结果是程序自动选择平滑程度的结果, 有些过于平滑了。可以人为地用参数控制结果的平滑程度。用参数 `span` 控制平滑程度, 越大的 `span` 值越光滑:

```
Warning: Removed 21 rows containing missing values (geom_path).
```



要注意的是，平滑结果不仅依赖于平滑参数的选择，也依赖于平滑方法和模型的选择，不同模型可能给出差别很大的结果，所以使用平滑结果做预测时一定要特别谨慎。

### 31.7.2 参数模型方法

LOESS 等平滑方法属于非参数方法，这些方法本质上没有一个不依赖于数据值的参数模型。另一类建模和平滑方法是假定  $y_i$  与  $x_i$  之间服从一定的参数模型，如线性回归模型

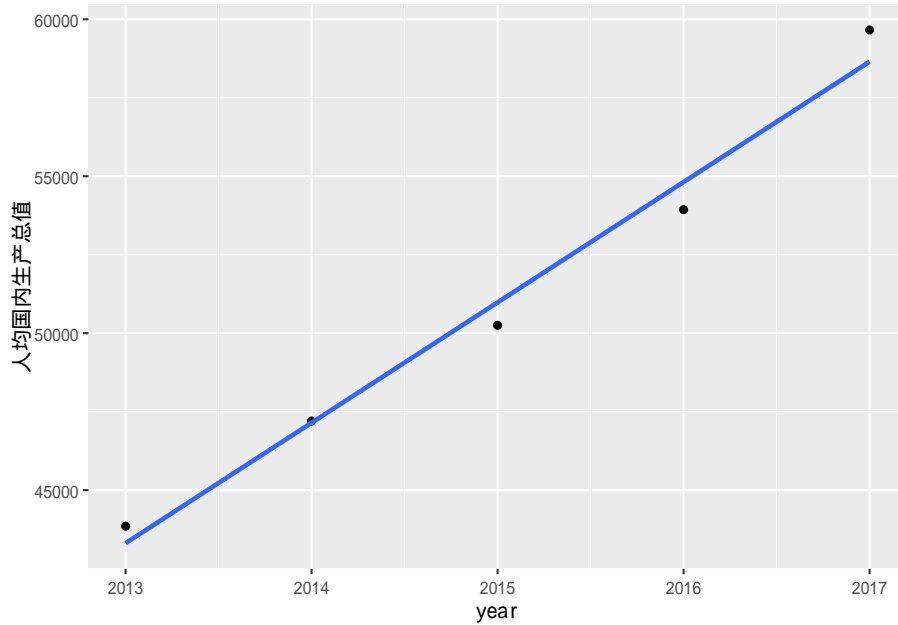
$$y_i = a + bx_i + \varepsilon_i$$

从数据估计模型参数，然后可以得到模型的拟合值（或称预测值）。

`geom_smooth()` 函数已经支持线性回归，比如对 2013-2017 人均国内生产总值数据进行平滑：

```
p <- ggplot(data = dinc, mapping = aes(
 x = year, y = `人均国内生产总值`))
p + geom_point() +
```

```
geom_smooth(
 method = "lm", size = 1.1, se=FALSE)
```



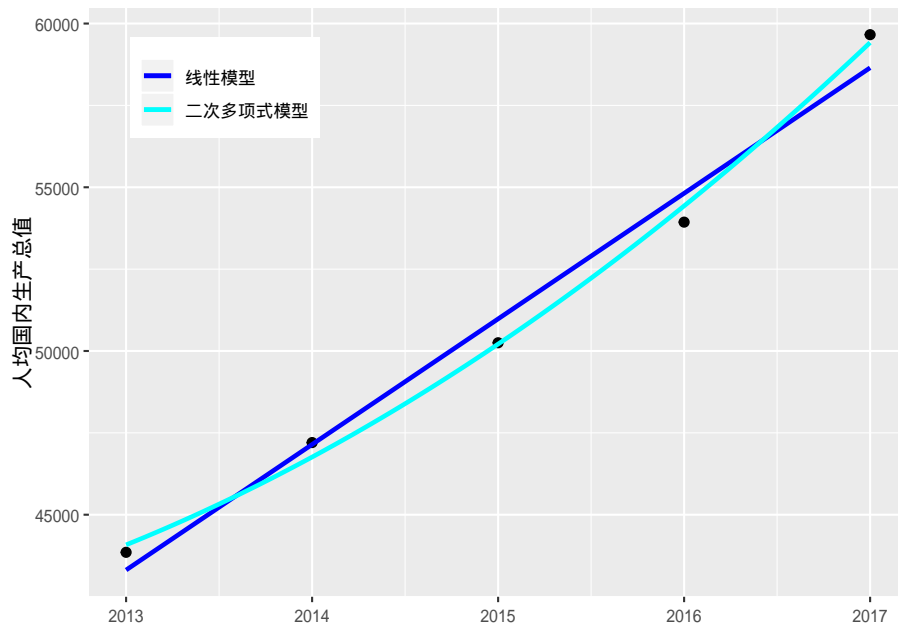
可以用 `formula = poly(x, 2)` 将模型变成二次多项式曲线，如：

```
p <- ggplot(data = dinc, mapping = aes(
 x = year, y = `人均国内生产总值`))
p + geom_point(size = 2) +
 geom_smooth(mapping = aes(col = "lm"),
 method = "lm", size = 1.1, se=FALSE) +
 geom_smooth(mapping = aes(col = "quadreg"),
 method = "lm", formula = y ~ poly(x, 2),
 size = 1.1, se=FALSE) +
 scale_color_manual(
 name = NULL,
 values = c(`lm` = "blue", `quadreg` = "cyan"),
 breaks = c("lm", "quadreg"),
 labels = c(" 线性模型", " 二次多项式模型")) +
 theme(
```

```

legend.position = c(0.05, 0.95),
legend.justification = c(0, 1)) +
labs(x = NULL)

```



`geom_smooth()` 还支持 `gam()` 函数能支持的模型。更一般的模型可以先用统计建模函数建模，然后在自变量取值范围内取一个网格计算拟合值，最后将原始数据与拟合曲线画在一起。这里拟合曲线实际是将网格点对应的拟合值连接成一条折线。

比如，对上面的数据拟合如下的指数增长模型：

$$y_t = Ae^{bt} + \varepsilon_t$$

```

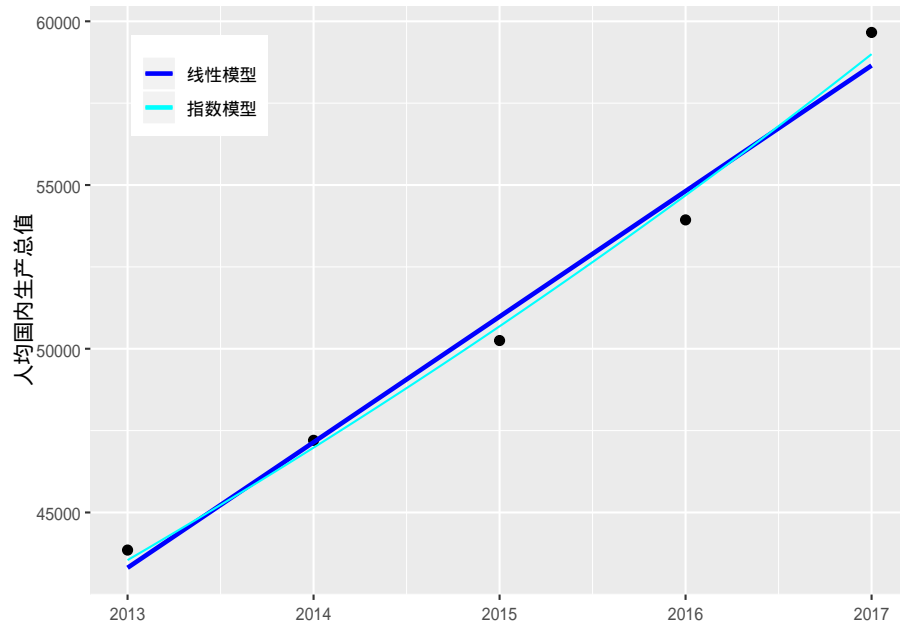
nls.inc <- nls(
 `人均国内生产总值` ~ A * exp(b * (year-2000)),
 data = dinc,
 start = c(A = 4E4, b = 1E-6))
dinc.pred1 <- data.frame(
 year = seq(2013, 2017, length.out = 100))
dinc.pred1 <- dinc.pred1 %>%

```

```

mutate(pred = predict(
 nls.inc, newdata = data.frame(year = year)))
p <- ggplot(data = dinc, mapping = aes(
 x = year, y = `人均国内生产总值`)
) + geom_point(size = 2) +
 geom_smooth(mapping = aes(col = "lm"),
 method = "lm", size = 1.1, se=FALSE) +
 geom_line(data = dinc.pred1, mapping = aes(
 x = year, y = pred, col = "Expmod")) +
 scale_color_manual(
 name = NULL,
 values = c(`lm` = "blue", `Expmod` = "cyan"),
 breaks = c("lm", "Expmod"),
 labels = c(" 线性模型", " 指数模型")) +
 theme(
 legend.position = c(0.05, 0.95),
 legend.justification = c(0, 1)) +
 labs(x = NULL)

```





$y = Ae^{bx}$  这样的模型可以将 y 轴改为对数尺度，这时指数曲线可以表现为直线。 $y = Ax^b$  这样的模型可以将 x 轴和 y 轴都改为对数尺度，曲线关系可以表现为直线关系。

## 31.8 表现不确定性

统计数据数据分析的建模结果都有一定的不确定性，参数估计和模型预测的不确定性经常用置信区间和置信带表现，绘图时可以画误差条 (error bar) 或置信带。这些方法可以比较精确地描绘不确定性大小，但需要一些概率统计知识才能正确理解。也可以开发一些更直观的图形用来表现不确定性。

### 31.8.1 表现概率

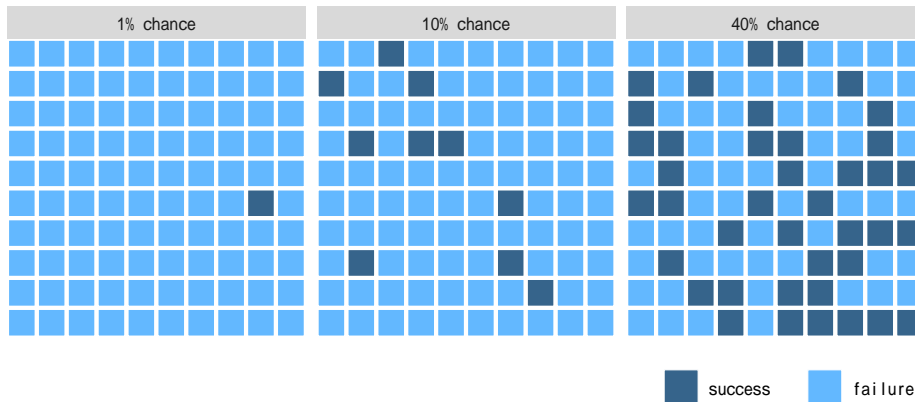
不确定性用概率描述，而概率可以通过多次独立重复试验的频率来理解。可以用一个饼图来表现概率。

一种容易理解的做法是画 100 个色块，用加亮显示部分表示成功部分，加亮显示部分随机摆放以显示结果的随机性。

```
#dg <- expand.grid(x = 1:10, y = 1:10)
set.seed(101)
dg <- expand.grid(ratio = c(0.01, 0.1, 0.4), x = 1:10, y = 1:10)
dg2 <- dg %>%
 group_by(ratio) %>%
 mutate(value = {
 n <- n()
 i <- round(n * ratio[1])
 sample(c(rep("S", i), rep("F", n-i)))
 }) %>%
 ungroup() %>%
 mutate(label = paste0(round(100*ratio), "% chance")
#set.seed(84524)

ggplot(data = dg2, mapping = aes(
```

```
x = x, y = y, fill = value)) +
geom_tile(color = "white", size = 1) +
coord_fixed(expand = FALSE, clip = "off") +
scale_x_continuous(name = NULL, breaks = NULL) +
scale_y_continuous(name = NULL, breaks = NULL) +
scale_fill_manual(
 name = NULL,
 breaks = c("S", "F"),
 labels = c("success", "failure"),
 values = c(
 "S" = "steelblue4",
 "F" = "steelblue1"
),
 guide = guide_legend(override.aes = list(size = 0))
) +
facet_wrap(~label) +
theme(
 legend.position = "bottom",
 legend.direction = "horizontal",
 legend.justification = "right",
 plot.margin = margin(0, 0, 3.5, 0) # crop plot a little more tightly
)
```



在只有两种随机结果时上图可以比较直观地表现两种结果的比例和随机性。在有許多可能结果，甚至于随机结果可以在一个区间内随机取值时，上面的方法不再有效，概率分布和概率密度能更准确地描述这种随机取值，但不一定容易理解。分布密度下方的阴影面积比例可以代表某个事件的概率，还可以用密度曲线下面画堆叠的圆的办法更直观地表现不同取值区域的概率大小。

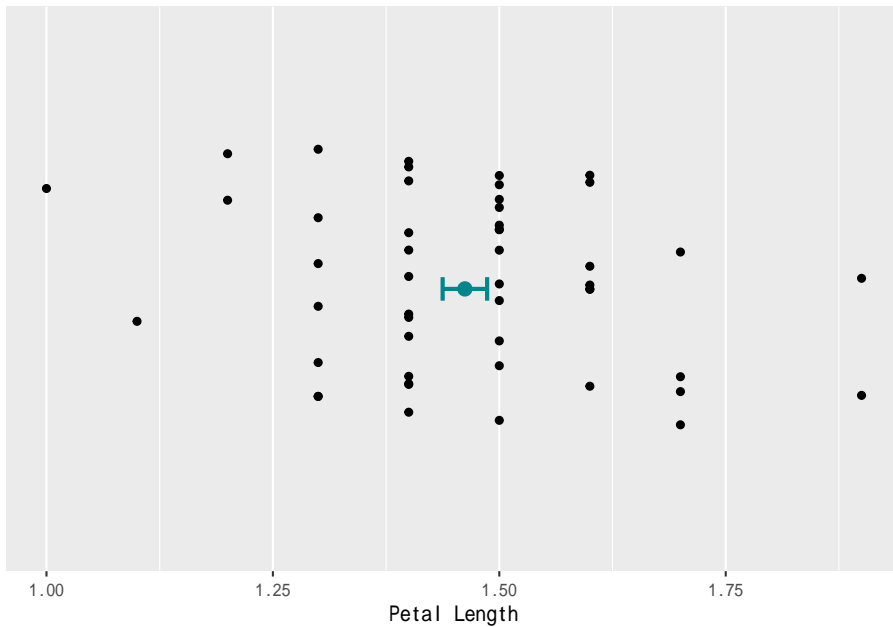
### 31.8.2 表现点估计精度

统计学中将要了解的分布称为总体，用参数概括总体的分布，从总体中抽取的一部分称为样本，从样本中计算统计量作为总体参数的估计，标准误差是估计量的精度的度量。可以用误差条同时画出参数估计值及对应的标准误差，误差条的两端是点估计加减标准误差。

例如，iris 数据集中 setosa 品种的花瓣长度均值及标准误差，同时画了对应的数据值：

```
d1 <- iris %>%
 filter(Species == "setosa")
d1summ <- d1 %>%
```

```
summarise(mpl = mean(Petal.Length),
 sd = sd(Petal.Length),
 se = sd / sqrt(n()))
p <- ggplot(data = d1summ)
p + geom_errorbar(
 mapping = aes(
 ymin = mpl - se, ymax = mpl + se,
 x = "setosa"),
 color = "turquoise4", width = 0.05, size = 1) +
geom_point(
 mapping = aes(
 y = mpl,
 x = "setosa"),
 color = "turquoise4", size = 3) +
geom_jitter(data = d1, mapping = aes(
 y = Petal.Length, x = "setosa"),
 width = 0.3, height = 0) +
scale_x_discrete(name = NULL, breaks = NULL) +
labs(y = "Petal Length") +
coord_flip()
```

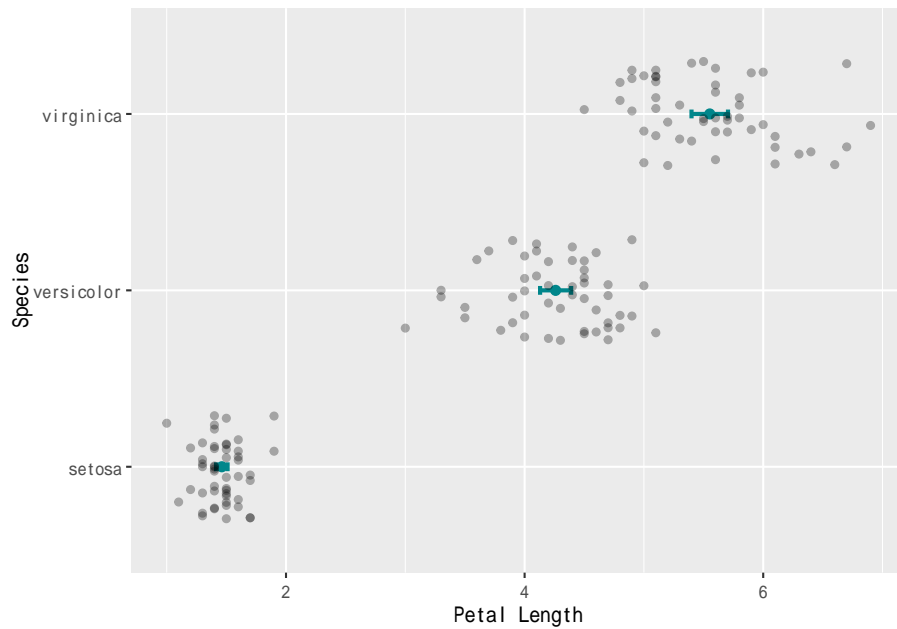


这里用了 `coord_flip()` 将图形 x、y 坐标对调，使得原来纵向的误差条变成横向。在 `geom_errorbar()` 中需要映射 `ymin`, `ymax` 和 `x` 维，分别表示误差条的上下限与 `x` 位置。

置信区间也可以类似地画出。下面对 `iris` 数据集中三个品种分别估计花瓣长度均值并画 95% 置信区间。

```
d2summ <- iris %>%
 group_by(Species) %>%
 summarise(mpl = mean(Petal.Length),
 sd = sd(Petal.Length),
 se = sd / sqrt(n()))
p <- ggplot(data = d2summ, mapping = aes(x = Species))
p + geom_errorbar(
 mapping = aes(
 ymin = mpl - 1.96*se, ymax = mpl + 1.96*se),
 color = "turquoise4", width = 0.05, size = 1) +
 geom_point(
 mapping = aes(y = mpl),
```

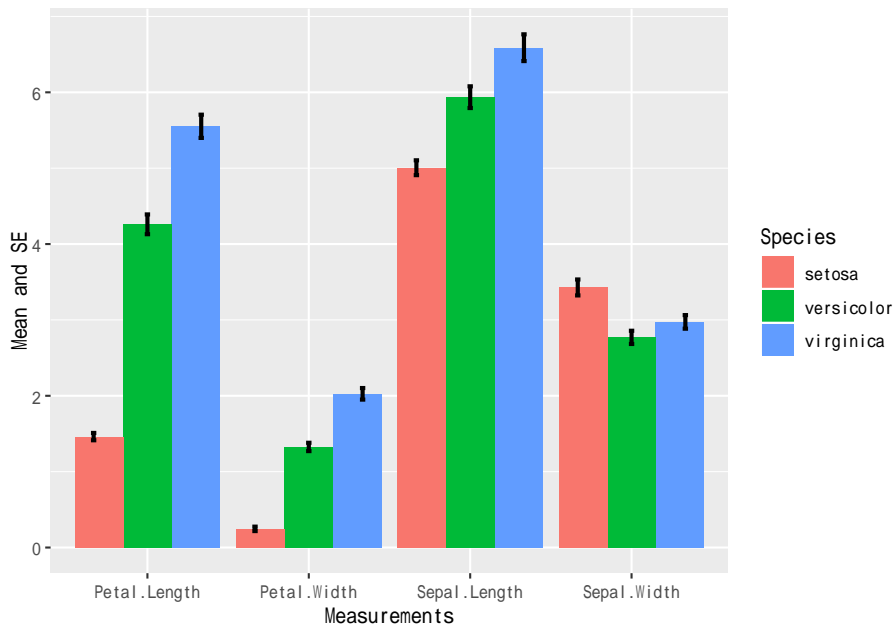
```
color = "turquoise4", size = 2) +
geom_jitter(data = iris, mapping = aes(
 y = Petal.Length, x = Species),
 alpha = 0.3, width = 0.3, height = 0) +
labs(y = "Petal Length") +
coord_flip()
```



置信区间是频率学派统计学的概念，可以认为在多次独立重复同一置信区间问题时真实参数落入置信区间的比例为给定的置信度，或者零假设下的参数值落入置信区间是对应的假设检验的接受域。而贝叶斯学派统计学也有类似概念，称为 *credible interval*，这是认为参数作为服从后验分布的随机变量的取值概率为给定置信度最可能取值区间。

误差条比较容易被不了解相关概念的人误认为数据范围或者参数估计取值范围，所以使用时需要解释其含义。误差条的一个优势是很容易与其他图形一起使用，上面的误差条是与散点图一起使用的，也可以与条形图配合给出条形图表示的比例的误差大小。例如，iris 数据中三个类别的 4 个测量值的均值用条形高度表示，标准误差用误差条表示：

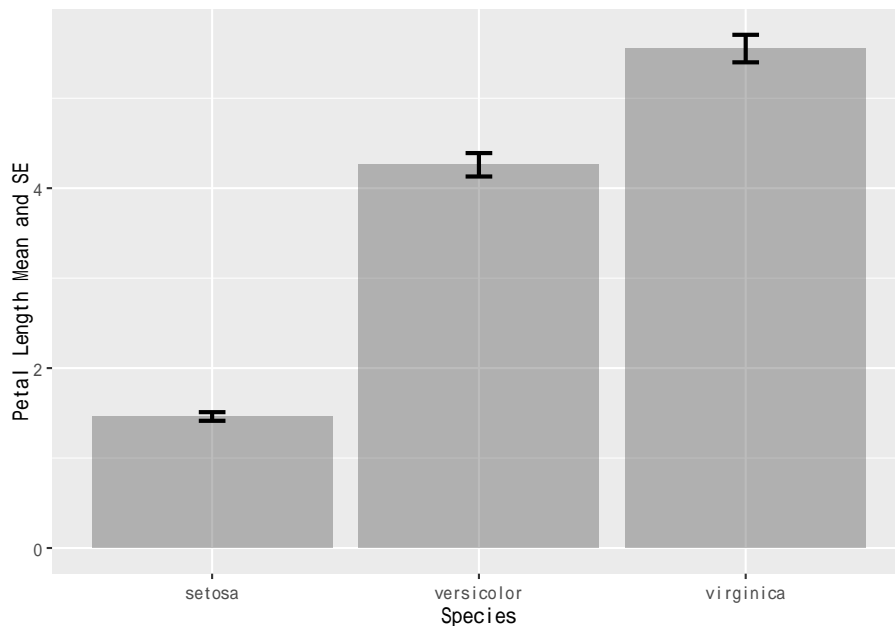
```
d3summ <- iris %>%
 gather(Sepal.Length, Sepal.Width, Petal.Length, Petal.Width,
 key = "variable", value = "value") %>%
 group_by(Species, variable) %>%
 summarise(
 mean = mean(value),
 se = sd(value) / sqrt(n()),
 lbar = mean - 1.96*se,
 hbar = mean + 1.96*se)
p <- ggplot(data = d3summ, mapping = aes(
 x = variable, fill = Species, y = mean))
p + geom_col(position = "dodge") +
 geom_errorbar(mapping = aes(
 ymin = lbar, ymax = hbar),
 position = position_dodge(width = 1),
 width = 0.1, size = 1) +
 labs(x = "Measurements", y = "Mean and SE")
```



这里 `geom_errorbar()` 的 `position` 参数用了 `position_dodge()` 函数说明两个条形之间的距离。`geom_errorbar()` 的 `width` 参数是条形端线的宽度, 如果 `width=0` 就没有了端线。

如果仅有一个变量则会比较简单, 比如, 花瓣长度的均值和加减标准误差:

```
d3b <- d3summ %>%
 filter(variable == "Petal.Length")
p <- ggplot(data = d3b, mapping = aes(
 x = Species, y = mean))
p + geom_col(alpha = 0.4) +
 geom_errorbar(mapping = aes(
 ymin = lbar, ymax = hbar),
 width = 0.1, size = 1) +
 labs(x = "Species", y = "Petal Length Mean and SE")
```



还可以同时绘制纵向与横向的误差条, 分别使用 `geom_errorbar()` 和 `geom_errorbarh()`。

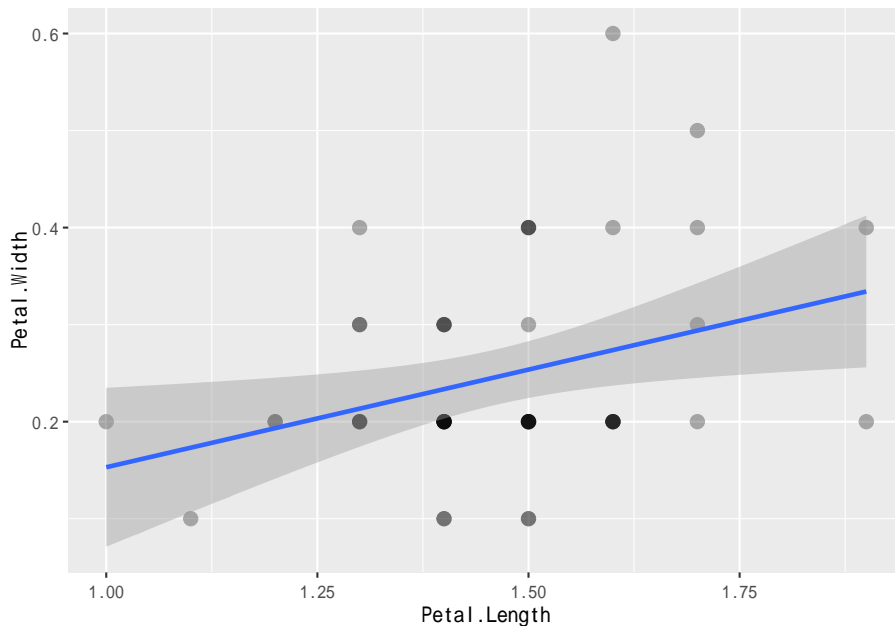


### 31.8.3 拟合曲线的置信带

对成对变量数据  $(x_i, y_i)$  可以拟合平滑曲线，同时往往希望给出曲线类似于置信区间或者 credible interval。作为置信区间解释时，虽然是置信带，但是置信度是针对每个单独的  $x$  值成立的；如果是贝叶斯置信区间 (credible interval)，则置信带可以理解为在模型参数的后验分布下取值概率为对应置信度的参数所对应的曲线集合。

`geom_smooth()` 函数中加 `se = TRUE` 就可以做置信带。例如，对 iris 数据集中 setosa 品种的花瓣长、宽作线性回归并加置信带：

```
diris1 <- iris %>%
 filter(Species == "setosa")
p <- ggplot(data = diris1, mapping = aes(
 x = Petal.Length, y = Petal.Width))
p + geom_point(size = 3, alpha = 0.3) +
 geom_smooth(
 method = "lm", size = 1.1, se=TRUE)
```

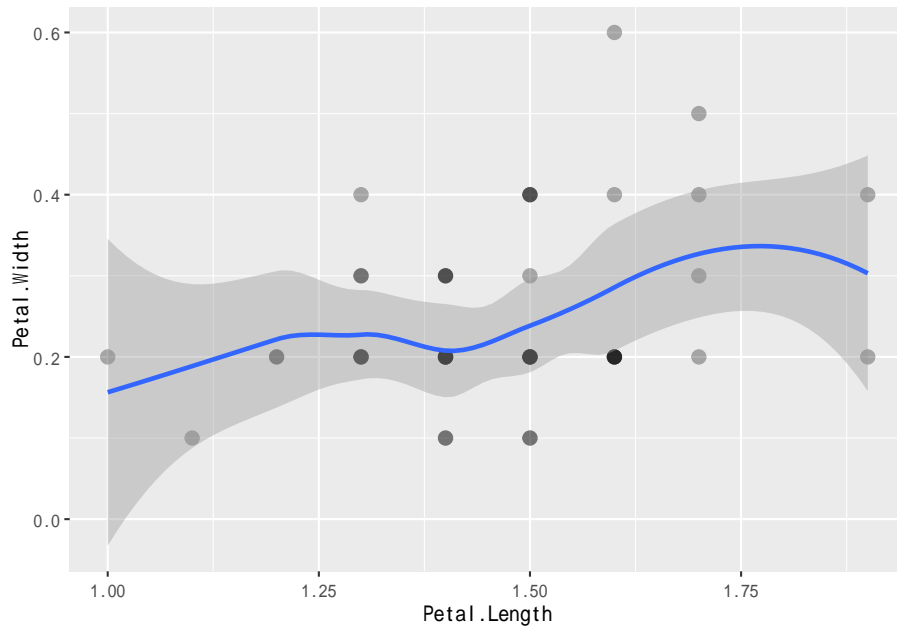


默认的置信度是 0.95，用参数 `level` 控制。`geom_smooth()` 的置信带是逐点

意义的。

LOESS 等平滑曲线也可以加置信带，如：

```
p + geom_point(size = 3, alpha = 0.3) +
 geom_smooth(
 method = "loess", size = 1.1, se=TRUE)
```



#### 31.8.4 随机结果动态演示

随机结果是概率分布确定但每个结果预先未知的。在讲课时，可以用动图（gif 或者 mp4 等格式）动态地演示不同的随机结果。gganimate 扩展包提供了制作动图的功能。

## Part VII

## 统计分析



## Chapter 32

# R 初等统计分析

这一部分讲授如何用 R 进行统计分析，包括基本概括统计和探索性数据分析，置信区间和假设检验，回归分析与各种回归方法，广义线性模型，非线性回归与平滑，判别树和回归树，等等。

主要参考书：

- (Venables and Ripley, 2002)
- (Kabacoff, 2012)

### 32.1 概率分布

R 中与 xxx 分布有关的函数包括：

- `dxxx(x)`，即 xxx 分布的分布密度函数 (PDF) 或概率函数 (PMF) $p(x)$ 。
- `pxxx(q)`，即 xxx 分布的分布函数 (CDF) $F(q) = P(X \leq q)$ 。
- `qxxx(p)`，即 xxx 分布的分位数函数  $q(p)$ ， $p \in (0, 1)$ ，对连续型分布， $q(p) = F^{-1}(p)$ ，即  $F(x) = p$  的解  $x$ 。
- `rxxx(n)`，即 xxx 的随机数函数，可以生成  $n$  个 xxx 的随机数。

`dxxx(x)` 函数可以加选项 `log=TRUE`，用来计算  $\ln p(x)$ ，这在计算对数似然函数时有用，比 `log(dxxx(x))` 更精确。

`pxxx(q)` 可以加选项 `lower.tail=FALSE`, 变成计算  $P(X > q) = 1 - F(q)$ 。

`qxxx(p)` 可以加选项 `lower.tail=TRUE`, 表示求  $P(X > x) = p$  的解  $x$ ; 可以加选项 `log.p=TRUE`, 表示输入的  $p$  实际是  $\ln p$ 。

这些函数都可以带有自己的分布参数, 有些分布参数有缺省值, 比如正态分布的缺省参数值为零均值单位标准差。

具体的分布类型可以在 R 命令行用 `?Distributions` 查看列表。常用的分布密度有:

- 离散分布有 `dbinom` 二项分布, `dpois` 泊松分布, `dgeom` 几何分布, `dnbinom` 负二项分布, `dmultinom` 多项分布, `dhyper` 超几何分布。
- 连续分布有 `dunif` 均匀分布, `dnorm` 正态分布, `dchisq` 卡方分布, `dt`  $t$  分布 (包括非中心  $t$ ), `df`  $F$  分布, `dexp` 指数分布, `dweibull` 威布尔分布, `dgamma` 伽马分布, `dbeta` 贝塔分布, `dlnorm` 对数正态分布, `dcauchy` 柯西分布, `dlogis` 逻辑斯谛分布。

R 的扩展包提供了更多的分布, 参见 R 网站的如下链接:

- <https://cran.r-project.org/web/views/Distributions.html>

## 32.2 最大似然估计

R 函数 `optim()`、`nlm()`、`optimize()` 可以用来求函数极值, 因此可以用来计算最大似然估计。`optimize()` 只能求一元函数极值。

### 32.2.1 一元正态分布参数最大似然估计

正态分布最大似然估计有解析表达式。作为示例, 用 R 函数进行数值优化求解。

对数似然函数为:

$$\ln L(\mu, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum (X_i - \mu)^2$$

定义 R 的优化目标函数为上述对数似然函数去掉常数项以后乘以  $-2$ , 求其最小值点。目标函数为:

```
objf.norm1 <- function(theta, x){
 mu <- theta[1]
 s2 <- exp(theta[2])
 n <- length(x)
 res <- n*log(s2) + 1/s2*sum((x - mu)^2)
 res
}
```

其中  $\theta_1$  为均值参数  $\mu$ ,  $\theta_2$  为方差参数  $\sigma^2$  的对数值。x 是样本数值组成的 R 向量。

可以用 `optim` 函数来求极小值点。下面是一个模拟演示:

```
norm1d.mledemo1 <- function(n=30){
 mu0 <- 20
 sigma0 <- 2
 set.seed(1)
 x <- rnorm(n, mu0, sigma0)

 theta0 <- c(0,0)
 ores <- optim(theta0, objf.norm1, x=x)
 print(ores)
 theta <- ores$par
 mu <- theta[1]
 sigma <- exp(0.5*theta[2])
 cat(' 真实 mu=', mu0, ' 公式估计 mu=', mean(x),
 ' 数值优化估计 mu=', mu, '\n')
 cat(' 真实 sigma=', sigma0,
 ' 公式估计 sigma=', sqrt(var(x)*(n-1)/n),
 ' 数值优化估计 sigma=', sigma, '\n')
}
norm1d.mledemo1()
```

```
$par
[1] 20.166892 1.193593
```

```
##
$value
[1] 65.83709
##
$counts
function gradient
115 NA
##
$convergence
[1] 0
##
$message
NULL
##
真实mu= 20 公式估计mu= 20.16492 数值优化估计mu= 20.16689
真实sigma= 2 公式估计sigma= 1.817177 数值优化估计sigma= 1.816291
```

也可以用 `nlm()` 函数求最小值点，下面是正态分布最大似然估计的另一个演示：

```
norm1d.mledemo2 <- function(){
 set.seed(1)
 n <- 30
 mu <- 20
 sig <- 2
 z <- rnorm(n, mean=mu, sd=sig)
 neglogL <- function(parm) {
 -sum(dnorm(z, mean=parm[1],
 sd=exp(parm[2]), log=TRUE))
 }

 res <- nlm(neglogL, c(10, log(10)))
 print(res)
 sig2 <- exp(res$estimate[2])
 cat(' 真实 mu=', mu, ' 公式估计 mu=', mean(z),
```



```

 ' 数值优化估计 mu=', res$estimate[1], '\n')
cat(' 真实 sigma=', sig,
 ' 公式估计 sigma=', sqrt(var(z)*(n-1)/n),
 ' 数值优化估计 sigma=', sig2, '\n')
invisible()
}
norm1d.mledemo2()

$minimum
[1] 60.48667
##
$estimate
[1] 20.1649179 0.5972841
##
$gradient
[1] 1.444576e-05 9.094805e-06
##
$code
[1] 2
##
$iterations
[1] 28
##
真实mu= 20 公式估计mu= 20.16492 数值优化估计mu= 20.16492
真实sigma= 2 公式估计sigma= 1.817177 数值优化估计sigma= 1.817177

```

函数 `optim()` 缺省使用 Nelder-Mead 单纯型搜索算法，此算法不要求计算梯度和海色阵，算法稳定性好，但是收敛速度比较慢。可以用选项 `method="BFGS"` 指定 BFGS 拟牛顿法。这时可以用 `gr=` 选项输入梯度函数，缺省使用数值微分计算梯度。如：

```

norm1d.mledemo1b <- function(n=30){
 mu0 <- 20
 sigma0 <- 2
 set.seed(1)

```

```

x <- rnorm(n, mu0, sigma0)

theta0 <- c(1,1)
ores <- optim(theta0, objf.norm1, x=x, method="BFGS")
print(ores)
theta <- ores$par
mu <- theta[1]
sigma <- exp(0.5*theta[2])
cat(' 真实 mu=', mu0, ' 公式估计 mu=', mean(x),
 ' 数值优化估计 mu=', mu, '\n')
cat(' 真实 sigma=', sigma0,
 ' 公式估计 sigma=', sqrt(var(x)*(n-1)/n),
 ' 数值优化估计 sigma=', sigma, '\n')
}
norm1d.mledemo1b()

$par
[1] 20.164916 1.194568
##
$value
[1] 65.83704
##
$counts
function gradient
41 17
##
$convergence
[1] 0
##
$message
NULL
##
真实mu= 20 公式估计mu= 20.16492 数值优化估计mu= 20.16492

```

```
真实sigma= 2 公式估计sigma= 1.817177 数值优化估计sigma= 1.817177
```

注意使用数值微分时，参数的初值的数量级应该与最终结果相差不大，否则可能不收敛，比如上面程序取初值 (0,0) 就会不收敛。

`optim()` 函数支持的方法还有"CG" 是一种共轭梯度法，此方法不如 BFGS 稳健，但是可以减少存储量使用，对大规模问题可能更适用。"L-BFGS-B" 是一种区间约束的 BFGS 优化方法。"SANN" 是模拟退火算法，有利于求得全局最小值点，速度比其它方法慢得多。

函数 `optimize()` 进行一元函数数值优化计算。例如，在一元正态分布最大似然估计中，在对数似然函数中代入  $\mu = \bar{x}$ ，目标函数（负 2 倍对数似然函数）为

$$h(\sigma^2) = \ln(\sigma^2) + \frac{1}{\sigma^2} \sum (X_i - \bar{X})^2$$

下面的例子模拟计算  $\sigma^2$  的最大似然估计：

```
norm1d.mledemo3 <- function(n=30){
 mu0 <- 20
 sigma0 <- 2
 set.seed(1)
 x <- rnorm(n, mu0, sigma0)

 mu <- mean(x)
 ss <- sum((x - mu)^2)/length(x)
 objf <- function(delta, ss) log(delta) + 1/delta*ss
 ores <- optimize(objf, lower=0.0001,
 upper=1000, ss=ss)
 delta <- ores$minimum
 sigma <- sqrt(delta)

 print(ores)
 cat(' 真实 sigma=', sigma0,
 ' 公式估计 sigma=', sqrt(var(x)*(n-1)/n),
 ' 数值优化估计 sigma=', sigma, '\n')
}
norm1d.mledemo3()
```

```
$minimum
[1] 3.30214
##
$objective
[1] 2.194568
##
真实sigma= 2 公式估计sigma= 1.817177 数值优化估计sigma= 1.817179
```

## 32.3 假设检验和置信区间

### 32.3.1 均值的假设检验和置信区间

#### 32.3.1.1 单样本均值

##### 32.3.1.1.1 大样本情形

设总体  $X$  期望为  $\mu$ , 方差为  $\sigma^2$ 。  $X_1, X_2, \dots, X_n$  是一组样本。

在大样本时, 令

$$Z = \frac{\bar{x} - \mu}{S/\sqrt{n}} \overset{\bullet}{\sim} N(0, 1)$$

其中  $\overset{\bullet}{\sim}$  表示近似服从某分布。 $\bar{x}$  为样本均值,  $S$  为样本标准差。 $\sigma$  未知且大样本时,  $\mu$  近似  $1 - \alpha$  置信区间为

$$\bar{x} \pm \text{qnorm}(1 - \alpha/2)S/\sqrt{n}$$

用 `BSDA::z.test(x, sigma.x=sd(x), conf.level=1-alpha)` 可以计算如上的近似置信区间。

对于假设检验问题

$$H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$$

$$H_0 : \mu \leq \mu_0 \longleftrightarrow H_a : \mu > \mu_0$$

$$H_0 : \mu \geq \mu_0 \longleftrightarrow H_a : \mu < \mu_0$$

定义检验统计量

$$Z = \frac{\bar{x} - \mu_0}{S/\sqrt{n}}$$

当  $\mu = \mu_0$  时  $Z$  近似服从  $N(0,1)$  分布。可以用 `BSDA::z.test(x, mu=mu0, sigma.x=sd(x))` 进行双侧  $Z$  检验。加选项 `alternative="greater"` 作右侧检验，加选项 `alternative="less"` 作左侧检验。

如果上述问题中标准差  $\sigma$  为已知，应该在  $Z$  的公式用  $\sigma$  替换  $S$ ，即

$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

计算  $\mu$  的置信区间程序为 `BSDA::z.test(x, sigma.x=sigma0)`，计算  $Z$  检验的程序为 `BSDA::z.test(x, mu=mu0, sigma.x=sigma0)`，仍可用选项 `alternative=` 指定双侧、右侧、左侧。

### 32.3.1.1.2 小样本正态情形

如果样本量较小（小于 30），则需要假定总体服从正态分布  $N(\mu, \sigma^2)$ 。

当  $\sigma$  已知时，有

$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \sim N(0,1)$$

计算  $\mu$  的置信区间程序为 `BSDA::z.test(x, sigma.x=sigma0)`，计算  $Z$  检验的程序为 `BSDA::z.test(x, mu=mu0, sigma.x=sigma0)`。

可以自己写一个计算  $Z$  检验的较通用的 R 函数：

```
z.test.ls <- function(
 x, n=length(x), mu=0, sigma=sd(x), alternative="two.sided"){
 z <- (mean(x) - mu) / (sigma / sqrt(n))
 if(alternative=="two.sided"){ # 双侧检验
 pvalue <- 2*(1 - pnorm(abs(z)))
 } else if(alternative=="less"){ # 左侧检验
 pvalue <- pnorm(z)
 } else if(alternative=="greater"){ # 右侧检验
 pvalue <- 1 - pnorm(z)
 } else {
 stop("alternative unknown!")
 }
}
```

```
c(stat=z, pvalue=pvalue)
}
```

这个函数输入样本  $x$  和  $\sigma = \sigma_0$ , 或者输入  $\bar{x}$  为  $x$  和  $\sigma = \sigma_0$ , 输出检验统计量值和  $p$  值。

如果  $\sigma$  未知, 应使用  $t$  统计量检验

$$T = \frac{\bar{x} - \mu}{S/\sqrt{n}}$$

当  $\mu = \mu_0$  时  $T \sim t(n-1)$ 。如果  $x$  保存了样本, 用 `t.test(x, conf.level=1-alpha)` 计算  $\mu$  的置信区间; 用 `t.test(x, mu=mu0, alternative="two.sided")` 计算  $H_0: \mu = \mu_0$  的双侧检验, 称为单样本  $t$  检验。

如果已知的  $\bar{x}$  和  $S$  而非具体样本数据, 可以自己写一个 R 函数计算  $t$  检验:

```
t.test.1s <- function(
 x, n=length(x), sigma=sd(x), mu=0,
 alternative="two.sided"){
 tstat <- (mean(x) - mu) / (sigma / sqrt(n))
 if(alternative=="two.sided"){ # 双侧检验
 pvalue <- 2*(1 - pt(abs(tstat), n-1))
 } else if(alternative=="less"){ # 左侧检验
 pvalue <- pt(tstat, n-1)
 } else if(alternative=="greater"){ # 右侧检验
 pvalue <- 1 - pt(tstat, n-1)
 } else {
 stop("alternative unknown!")
 }

 c(stat=tstat, pvalue=pvalue)
}
```

这个函数允许输入  $\bar{x}$  到  $x$  中, 输入  $n$ , 输入  $S$  到 `sigma` 中, 然后计算  $t$  检验。

### 32.3.1.1.3 统计假设检验显著性的解释

统计假设检验如果拒绝了  $H_0$ ，也称“结果显著”、“有显著差异”等。这只代表样本与零假设的差距足够大，可以以比较小的代价（第一类错误概率）排除这样的差距是由于随机样本的随机性引起的，但是并不表明这样的差距是有科学意义或者有现实意义的。

尤其是超大样本时，完全没有实际意义的差距也会检验出来是统计显著的。

这时，可以设定一个“有实际意义的差距” $\delta$ ，当差距超过  $\delta$  时才认为有实际差距，如：

$$H_0 : |\mu - \mu_0| \leq \delta \longleftrightarrow |\mu - \mu_0| > \delta$$

统计假设检验如果不拒绝  $H_0$ ，也称为“结果不显著”、“没有显著差异”，通常不认为  $H_0$  就是对的，看法是没有充分理由推翻  $H_0$ 。

进一步收集更多的数据，有可能就会发现显著差异。

例如， $H_0$  是嫌疑人无罪，最后不拒绝  $H_0$ ，只能是没有充分证据证明有罪，如果后来收集到了关键的罪证，就还可以再判决嫌疑人有罪。

#### 32.3.1.1.4 置信区间与双侧检验的关系

设总体为正态分布  $N(\mu, \sigma^2)$ ， $\sigma$  已知。随机抽取了  $n$  个样品，计算得  $\bar{x}$ 。 $\mu$  的  $1 - \alpha$  置信区间为

$$\bar{x} \pm z_{\alpha/2} \sigma / \sqrt{n}$$

置信区间理论依据：

$$P\left(\frac{|\bar{x} - \mu|}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = 1 - \alpha \quad (*)$$

对双侧检验

$$H_0 : \mu = \mu_0 \longleftrightarrow H_a : \mu \neq \mu_0$$

当

$$\frac{|\bar{x} - \mu_0|}{\sigma/\sqrt{n}} \leq z_{\alpha/2}$$

时接受  $H_0$ 。

接受  $H_0$  的条件，就是  $\mu_0$  满足 (\*) 式的大括号中的条件，即当且仅当  $\mu_0$  落入  $\mu$  的  $1 - \alpha$  置信区间时接受  $H_0$ 。

其它检验也与置信区间有类似关系，但是单侧检验对应于单侧置信区间。

### 32.3.1.1.5 单样本均值比较例子

载入假设检验部分示例程序所用数据：

```
load("hyptest-data.RData")
```

#### 32.3.1.1.5.1 咖啡标重的单侧检验

美国联邦贸易委员会定期检查商品标签是否与实际相符。Hiltop Coffee 大罐咖啡标签注明含量为 3 磅。每罐具体多一些或少一些并不要紧，只要总体平均值不少于 3 磅，消费者利益就未受侵害。

只有在有充分证据时才应该做出分量不够的判决进而对厂家进行处罚。所以零假设是  $H_0: \mu \geq 3$ ，对立假设是  $H_a: \mu < 3$ 。

收集数据进行统计判决后，如果不拒绝  $H_0$ ，就不需要采取行动；如果拒绝了  $H_0$ ，就说明有足够证据认为厂商灌装分量不足，应进行处罚。

随机抽取了 36 件样品，计算平均净重  $\bar{x}$  作为总体均值  $\mu$  的估计。

如果  $\bar{x} < 3$ ，则  $H_0$  有理由被怀疑。

$\bar{x}$  比  $\mu_0 = 3$  少多少，才让我们愿意拒绝  $H_0$ ？做出拒绝  $H_0$  的决定可能会犯第一类错误，如果差距很小，这样的差距很可能是随机抽样的随机性引起的，拒绝  $H_0$  有很大的机会犯错。

必须先确定一个能够容忍的第一类错误概率。调查员认为，可以容忍 1% 的错误拒绝  $H_0$  的机会（不该处罚时给出了处罚）以换得能够在厂商分量不足时正确地做出处罚决定的机会。即选检验水平  $\alpha = 0.01$ 。

当  $H_0$  实际成立时，边界  $\mu = \mu_0 = 3$  处犯第一类错误的概率最大（这里最不容易区分），所以计算第一类错误概率只要在边界处计算。

统计假设检验第一步是确定零假设与对立假设，第二步是确定检验水平  $\alpha$ 。第三步是抽取随机样本，计算检验统计量。

以往经验证明净重的标准差是  $\sigma = 0.18$ ，且净重的分布为正态分布。

则  $\bar{x}$  的抽样分布为  $N(\mu, \sigma^2/n)$ ， $\bar{x}$  的标准误差为  $SE(\bar{x}) = 0.18/\sqrt{36} = 0.03$ 。

令

$$z = \frac{\bar{x} - \mu_0}{SE(\bar{x})}$$



则  $H_0$  成立且  $\mu$  取边界处的  $\mu_0$  时  $z \sim N(0, 1)$ 。

当  $\bar{x}$  比  $\mu_0 = 3$  小很多时应拒绝  $H_0$ 。即检验统计量  $z$  很小时拒绝  $H_0$ 。

如果  $z = -1$ ，标准正态分布取-1 或比-1 更小的值的概率是 0.1586 (用 R 程序 `pnorm(-1)` 计算)，这个概率不小，不小于  $\alpha = 0.01$ ，这时不应该拒绝  $H_0$ 。

如果  $z = -2$ ，标准正态分布取-2 或比-2 更小的值的概率是 0.0228，(用 R 程序 `pnorm(-2)` 计算)，这个概率已经比较小了，但还不小于预先确定的第一类错误概率  $\alpha = 0.01$ ，这时不应该拒绝  $H_0$ 。

如果  $z = -3$ ，标准正态分布取-3 或比-3 更小的值的概率是 0.0013，(用 R 程序 `pnorm(-3)` 计算)，这个概率已经很小了，也小于预先确定的第一类错误概率  $\alpha = 0.01$ ，这时应该拒绝  $H_0$ 。

对左侧检验问题  $H_0 : \mu \geq \mu_0 \longleftrightarrow H_a : \mu < \mu_0$ ，若  $Z$  表示标准正态分布随机变量， $z$  是检验统计量的值，在  $\mu = \mu_0$  时  $z$  的抽样分布是标准正态分布，可以计算抽样分布中取到等于  $z$  以及比  $z$  更小的值的概率  $P(Z \leq z)$ ，称为检验的 p 值。

p 值越小，说明  $\bar{x}$  比  $\mu_0$  小得越多，拒绝  $H_0$  也越有充分证据。当且仅当 p 值小于检验水平  $\alpha$  时拒绝  $H_0$ 。

p 值小于  $\alpha$ ，说明如果  $H_0$  真的成立，出现  $\bar{x}$  那么小的检验统计量值的概率也很小，小于  $\alpha$ ，所以第一类错误概率小于等于  $\alpha$ 。

一般地，p 值是从样本计算的一个在假定零假设成立情况下，检验统计量取值更倾向于对立假设成立的概率。p 值越小，拒绝  $H_0$  越有依据。

设咖啡标称重量问题中抽取了 36 件样品，测得  $\bar{x} = 2.92$  磅，计算得

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{2.92 - 3}{0.18/\sqrt{36}} = -2.67$$

p 值为  $P(Z \leq -2.67)$  ( $Z$  是标准正态分布随机变量)，用 R 计算得 `pnorm(-2.67) = 0.0038`  $\leq \alpha = 0.01$ ，拒绝  $H_0$ ，应对该厂商进行处罚。

已经预先确定了检验水平  $\alpha = 0.01$  就不能在看到检验结果后更改检验水平。但是，报告出 p 值后，其他的人可以利用这样的信息，不同的人或利益方可能有不同的检验水平的选择，这通过同一个 p 值都可以给出需要的结果。

p 值也称为“观测显著性水平”，因为它是能够拒绝  $H_0$  的最小的可选检验水平，选取比 p 值更小的检验水平就不能拒绝  $H_0$  了。

用自定义的 `z.test.ls()` 计算：

```
z.test.ls(x=2.92, mu=3, n=36, sigma=0.18, alternative="less")
```

```
stat pvalue
-2.666666667 0.003830381
```

p 值为 0.004，在 0.01 水平下拒绝原假设，认为咖啡平均重量显著地低于标称的 3 磅。

设 `Coffee` 数据框的 `Weight` 是这 36 个样本值，也可以将程序写成：

```
z.test.ls(Coffee[["Weight"]], mu=3, sigma=0.18, alternative="less")
```

```
stat pvalue
-2.666666667 0.003830381
```

或

```
BSDA::z.test(Coffee[["Weight"]], mu=3, sigma.x=0.18, alternative="less")
```

```
##
One-sample z-Test
##
data: Coffee[["Weight"]]
z = -2.6667, p-value = 0.00383
alternative hypothesis: true mean is less than 3
95 percent confidence interval:
NA 2.969346
sample estimates:
mean of x
2.92
```

### 32.3.1.1.5.2 例：高尔夫球性能的双侧检验

美国高尔夫协会对高尔夫运动器械有一系列的规定。MaxFlight 生产高品质的

高尔夫球，平均的飞行距离是 295 码。

生产线有时会有波动，当生产的球平均飞行距离不足 295 码时，用户会感觉不好用；当生产的球平均飞行距离超过 295 码时，美国高尔夫协会禁止使用这样的球。

该厂商定期抽检生产线上下来的产品，每次抽取 50 只。

当平均飞行距离与 295 码没有显著差距时，不需要采取措施；如果有显著差距，就需要调整生产线。

所以取  $H_0 : \mu = 295$ ,  $H_a : \mu \neq 295$ 。选检验水平  $\alpha = 0.05$ 。

随机抽取  $n$  个样品，在可控条件下测试出每个的飞行距离，计算出平均值  $\bar{x}$ 。

当  $\bar{x}$  比  $\mu_0$  (这里是 295) 小很多或者  $\bar{x}$  比  $\mu_0$  大很多时拒绝  $H_0$ ，否则不拒绝  $H_0$ 。

设已知总体标准差  $\sigma = 12$ ，且飞行距离服从正态分布。

取统计量

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

在  $H_0$  成立时， $z$  的抽样分布为标准正态分布。

- 设测得  $\bar{x} = 297.6$ ，超过了  $\mu_0 = 295$ 。

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{297.6 - 295}{12/\sqrt{50}} = 1.53$$

$H_0$  成立时  $z$  抽样分布为标准正态分布，设  $Z$  是标准正态分布随机变量， $Z$  取到 1.53 或者比 1.53 更倾向于反对  $H_0$  的值概率？这包括  $P(Z \geq 1.53)$ ，还应包括  $P(Z \leq -1.53)$ 。

所以 p 值是  $P(|Z| \geq |z|)$ ，R 中计算为  $2*(1 - \text{pnorm}(\text{abs}(z)))$ 。p 值等于 0.1260，超过检验水平，不拒绝  $H_0$ ，不需要调整生产线。

用自定义的 `z.test.1s()` 检验：

```
z.test.1s(GolfTest[['Yards']], mu=295, sigma=12, alternative="two.sided")
```

```
stat pvalue
1.5320647 0.1255065
```

或用 `BSDA::z.test()`:

```
BSDA::z.test(GolfTest[['Yards']], mu=295, sigma.x=12, alternative="two.sided")

##
One-sample z-Test
##
data: GolfTest[["Yards"]]
z = 1.5321, p-value = 0.1255
alternative hypothesis: true mean is not equal to 295
95 percent confidence interval:
294.2738 300.9262
sample estimates:
mean of x
297.6
```

p 值为 0.13, 在 0.05 水平下不拒绝  $H_0$ , 生产的球的平均飞行距离与标准的 295 码之间没有显著差异。

### 32.3.1.1.5.3 Heathrow 机场打分的检验

某旅行杂志希望根据乘机进行商务旅行的意见给跨大西洋门户机场评分。乘客评分取 0 到 10 分, 超过 7 分算作是服务一流的机场。在每个机场随机选取 60 位商务旅行乘客打分评价。英国 Heathrow 机场的打分样本  $\bar{x} = 7.25$ ,  $S = 1.052$ 。英国 Heathrow 机场可以算是服务一流吗?

由于随机样本的随机性, 要判断的是  $\mu > 7$  是否成立, 这里  $\bar{x} = 7.25$  并不能很确定地说  $\mu > 7$  成立。

因为将机场评为一流需要有充分证据, 所以将对立假设  $H_a$  设为  $\mu > 7$ , 然后取  $H_0: \mu \leq 7$ 。

取检验水平  $\alpha = 0.05$ 。

计算检验统计量

$$t = \frac{\bar{x} - \mu_0}{S/\sqrt{n}} = \frac{7.25 - 7}{1.052/\sqrt{60}} = 1.84$$

设  $T$  是服从  $t(n-1)$  分布的随机变量, p 值为  $P(T \geq 1.84)$ 。用 R 软件计算

为  $1 - pt(1.84, 60-1) = 0.0354$ 。p 值小于  $\alpha$ ，拒绝  $H_0$ ，认为平均评分显著高于 7，可以认为 Heathrow 机场是服务一流。

从原始样本数据用 `t.test()` 检验：

```
t.test(AirRating[["Rating"]], mu=7, alternative="greater")
```

```
##
One Sample t-test
##
data: AirRating[["Rating"]]
t = 1.8414, df = 59, p-value = 0.03529
alternative hypothesis: true mean is greater than 7
95 percent confidence interval:
7.023124 Inf
sample estimates:
mean of x
7.25
```

用自定义的 `t.test.1s()` 和样本统计量计算：

```
t.test.1s(x=7.25, sigma=1.052, n=60, mu=7, alternative="greater")
```

```
stat pvalue
1.84077155 0.03534255
```

#### 32.3.1.1.5.4 玩具厂商订货量的假设检验

玩具厂商 Holiday Toys 通过超过 1000 家玩具商店售货。为了准备即将到来的冬季销售季节，销售主管需要确定今年某新款玩具的生产量。凭经验判断平均每家商店需要 40 件。

为了确定，随机抽取了 25 家商店，提供了新款玩具的特点、进货价和建议销售价，让他们给出订货量估计。目的是判断按每家 40 件来生产是否合适。

因为多生产或少生产都不合适，所以取零假设为  $H_0 : \mu = 40$ 。对立假设为  $H_a : \mu \neq 40$ 。只有拒绝  $H_0$  时才需要更改生产计划。

取检验水平  $\alpha = 0.05$ 。样本计算得  $n = 25$ ， $\bar{x} = 37.4$ ， $S = 11.79$ 。

检验统计量

$$t = \frac{\bar{x} - \mu_0}{S/\sqrt{n}} = \frac{37.4 - 40}{11.79/\sqrt{25}} = -1.10$$

设  $T$  为服从  $t(n-1)$  分布的随机变量,  $p$  值为  $2P(|T| > |-1.10|)$ , 用 R 计算为  $2*(1 - pt(abs(-1.10), 25-1)) = 0.2822$ 。

不拒绝  $H_0$ , 所以不需要更改生产计划。

从原始样本数据用 `t.test()` 检验:

```
t.test(Orders[["Units"]], mu=40, alternative="two.sided")
```

```
##
One Sample t-test
##
data: Orders[["Units"]]
t = -1.1026, df = 24, p-value = 0.2811
alternative hypothesis: true mean is not equal to 40
95 percent confidence interval:
32.5334 42.2666
sample estimates:
mean of x
37.4
```

用自定义的 `t.test.1s()` 和样本统计量计算:

```
t.test.1s(x=37.4, sigma=11.79, n=25, mu=40, alternative="two.sided")
```

```
stat pvalue
-1.1026293 0.2811226
```

### 32.3.1.2 均值比较

#### 32.3.1.2.1 独立两样本 Z 检验

假设有两个独立的总体  $X, Y$ , 例如, 男性与女性的寿命。要比较两个总体的期望  $\mu_1 = EX$  和  $\mu_2 = EY$ 。设两个总体的方差分别为  $\sigma_1^2$  和  $\sigma_2^2$ 。可以对均值作

双侧、左侧、右侧检验。比如，双侧检验为

$$H_0: \mu_1 = \mu_2 \longleftrightarrow H_a: \mu_1 \neq \mu_2$$

如果两个总体都服从正态分布，而且分别的方差  $\sigma_1^2$  和  $\sigma_2^2$  已知，设  $X$  的样本量为  $n_1$  的样本得到样本均值  $\bar{x}$ ，设  $Y$  的样本量为  $n_2$  的样本得到样本均值  $\bar{y}$ ，取统计量

$$Z = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

其中分母是分子的标准误差。当  $\mu_1 = \mu_2$  时  $Z \sim N(0, 1)$ 。

如果分布不一定是正态分布但是样本量足够大，可以令

$$Z = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{S_x^2}{n_1} + \frac{S_y^2}{n_2}}}$$

当  $\mu_1 = \mu_2$  时  $Z$  近似服从  $N(0, 1)$ 。

检验问题还可以变成  $\mu_1 - \mu_2$  与某个  $\delta$  的比较，如

$$H_0: \mu_1 \leq \mu_2 - \delta \longleftrightarrow H_a: \mu_1 > \mu_2 - \delta$$

这个对立假设是“不次于”的结论。

可以写一个大样本方差未知或已知，小样本独立两正态总体方差已知情况做  $Z$  检验的 R 函数：

```
z.test.2s <- function(
 x, y, n1=length(x), n2=length(y), delta=0,
 sigma1=sd(x), sigma2=sd(y), alternative="two.sided"){
 z <- (mean(x) - mean(y) - delta) / sqrt(sigma1^2 / n1 + sigma2^2 / n2)
 if(alternative=="two.sided"){ # 双侧检验
 pvalue <- 2*(1 - pnorm(abs(z)))
 } else if(alternative=="less"){ # 左侧检验
 pvalue <- pnorm(z)
 } else if(alternative=="greater"){ # 右侧检验
 pvalue <- 1 - pnorm(z)
 } else {
 stop("alternative unknown!")
 }
}
```

```

}

c(stat=z, pvalue=pvalue)
}

```

函数允许输入两组样本到  $x$  和  $y$  中, 输入  $\sigma_1$  和  $\sigma_2$  到 `sigma1` 和 `sigma2` 中, 这时不输入 `sigma1` 和 `sigma2` 则从样本中计算样本统计量代替。也允许输入  $\bar{x}$  和  $\bar{y}$  到  $x$  和  $y$  中, 输入  $\sigma_1$  和  $\sigma_2$  到 `sigma1` 和 `sigma2` 中, 或者输入 `S_x` 和 `S_y` 到 `sigma1` 和 `sigma2` 中, 输入 `n_1` 到  $n_1$  中, 输入 `n_2` 到  $n_2$  中, 计算独立两样本均值比较的  $Z$  检验。

### 32.3.1.2.2 独立两样本 $t$ 检验

如果样本是小样本, 但两个总体分别服从正态分布, 两个总体的方差未知, 但已知  $\sigma_1^2 = \sigma_2^2$ 。这时  $\sigma_1^2 = \sigma_2^2$  可以统一估计为

$$S_p^2 = \frac{1}{n_1 + n_2 - 2} ((n_1 - 1)S_x^2 + (n_2 - 1)S_y^2)$$

取检验统计量

$$T = \frac{\bar{x} - \bar{y}}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

当  $\mu_1 = \mu_2$  时  $T \sim t(n_1 + n_2 - 2)$ 。

当  $x$  和  $y$  存放了两组独立样本时, 可以用 `t.test(x, y, var.equal=TRUE)` 计算两样本  $t$  检验, 用 `alternative=` 选项指定双侧、右侧、左侧检验。

可以自己写一个这样的 R 函数作两样本  $t$  检验, 允许只输入统计量而非具体样本值:

```

t.test.2s <- function(
 x, y, n1=length(x), n2=length(y),
 sigma1=sd(x), sigma2=sd(y), delta=0,
 alternative="two.sided"){
 sp <- sqrt(1/(n1+n2-2) * ((n1-1)*sigma1^2 + (n2-1)*sigma2^2))
 t <- (mean(x) - mean(y) - delta) / (sp * sqrt(1 / n1 + 1 / n2))
 if(alternative=="two.sided"){ # 双侧检验
 pvalue <- 2*(1 - pt(abs(t), n1+n2-2))

```



```

} else if(alternative=="less"){ # 左侧检验
 pvalue <- pt(t, n1+n2-2)
} else if(alternative=="greater"){ # 右侧检验
 pvalue <- 1 - pt(t, n1+n2-2)
} else {
 stop("alternative unknown!")
}

c(stat=t, pvalue=pvalue)
}

```

如果两个独立正态总体的方差不相等，样本量不够大，可以用如下的在  $H_0$  下近似服从  $t$  分布的检验统计量：

$$T = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{S_x^2}{n_1} + \frac{S_y^2}{n_2}}}$$

自由度为

$$m = \frac{(S_x^2/n_1 + S_y^2/n_2)^2}{\frac{(S_x^2/n_1)^2}{n_1-1} + \frac{(S_y^2/n_2)^2}{n_2-1}}$$

这个检验称为 Welch 两样本  $t$  检验，或 Satterthwaite 两样本  $t$  检验当  $x$  和  $y$  存放了两组独立样本时，可以用 `t.test(x, y, var.equal=FALSE)` 计算 Welch 两样本  $t$  检验。

如果方差是否相等不易判断，可以直接选用 Welch 检验。

因为大样本时  $t$  分布与标准正态分布基本相同，所以非正态总体大样本  $Z$  检验也可以用 `t.test()` 计算。

### 32.3.1.2.3 成对均值比较问题

设总体包含两个  $X, Y$  分量，比如，一组病人的服药前血压与服药后血压。这两个变量是相关的，不能用独立两总体的均值比较方法比较  $\mu_1 = EX$  与  $\mu_2 = EY$ 。

若  $\xi = X - Y$  服从正态分布，则可以对  $\xi$  进行均值为零的单样本  $t$  检验。这样的检验称为成对均值的  $t$  检验。

若  $x$  和  $y$  分别为  $X$  和  $Y$  的样本值, 可以用 `t.test(x, y, paired=TRUE)` 计算成对均值  $t$  检验。

#### 32.3.1.2.4 均值比较的例子

##### 32.3.1.2.4.1 顾客平均年龄差别比较

HomeStyle 是一个连锁家具店, 两家分店一个在城区, 一个在郊区。经理发现同一种家具有可能在一个分店卖的很好, 另一个分店则不好。怀疑是顾客人群的人口学差异 (年龄、性别、种族等)。

独立地抽取两个分店的顾客样本, 比较平均年龄。设城内的分店顾客年龄  $\sigma_1 = 9$  已知, 郊区的分店顾客年龄  $\sigma_2 = 10$  已知。

城内调查了  $n_1 = 36$  位顾客, 年龄的样本平均值  $\bar{x} = 40$ ; 郊区调查了  $n_2 = 49$  位顾客, 年龄的样本平均值  $\bar{y} = 35$ 。  $\bar{x} - \bar{y} = 5$ 。

作双侧检验, 水平  $\alpha = 0.05$ 。

计算  $Z$  统计量:

$$Z = \frac{40 - 35}{\sqrt{\frac{9^2}{36} + \frac{10^2}{49}}} = 2.4138$$

$p$  值为  $P(|V| \geq |2.4138|)$ , 其中  $V$  为标准正态分布随机变量。用 R 计算为  $2*(1 - pnorm(abs(2.4138)))=0.016$ , 两个分店的顾客年龄有显著差异, 城里的顾客平均年龄显著地高。

基于样本统计量用自定义的 `z.test.2s()` 计算:

```
z.test.2s(n1=36, x=40, sigma1=9,
 n2=49, y=35, sigma2=10,
 alternative="two.sided")
```

```
stat pvalue
2.41379310 0.01578742
```

##### 32.3.1.2.4.2 两个银行营业所顾客平均存款比较

某银行要比较两个营业所的支票账户的平均存款额。作双侧检验，水平 0.05。没有方差信息时可以直接选用 Welch 检验。

设第一个营业所随机抽查了  $n_1 = 28$  个支票账户，均值为  $\bar{x} = 1025$ (美元)，样本标准差为  $S_1 = 150$ ，第二个营业所随机抽查了  $n_2 = 22$  个支票账户，均值为  $\bar{y} = 910$ ，样本标准差为  $S_2 = 125$ 。

从原始数据出发计算检验：

```
t.test(CheckAcct[[1]], CheckAcct[[2]], var.equal = FALSE, alternative = "two.sided")

##
Welch Two Sample t-test
##
data: CheckAcct[[1]] and CheckAcct[[2]]
t = 2.956, df = 47.805, p-value = 0.004828
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
36.77503 193.25224
sample estimates:
mean of x mean of y
1025.0000 909.9864
```

p 值为 0.005，在 0.05 水平下显著，两个营业所的账户平均余额有显著差异。

#### 32.3.1.2.4.3 两种工具软件的比较

考虑开发信息系统的两种工具软件的比较。新的软件声称比旧的软件能加快进度。设使用旧软件时平均项目完成时间为  $\mu_1$ ，使用新软件时为  $\mu_2$ 。

检验：

$$H_0 : \mu_1 \leq \mu_2 \longleftrightarrow H_a : \mu_1 > \mu_2$$

(对立假设是新软件的工期更短) 取 0.05 水平。

```
t.test(SoftwareTest[['Current']], SoftwareTest[['New']],
 var.equal = FALSE, alternative = 'greater')
```

```
##
```

```
Welch Two Sample t-test
##
data: SoftwareTest[["Current"]] and SoftwareTest[["New"]]
t = 2.2721, df = 21.803, p-value = 0.01665
alternative hypothesis: true difference in means is greater than 0
95 percent confidence interval:
9.514309 Inf
sample estimates:
mean of x mean of y
325 286
```

p 值为 0.017, 使用新的工具软件的平均项目完成时间显著地少于使用旧工具软件。

#### 32.3.1.2.4.4 两种工艺所需时间的比较

工厂希望比较同一生产问题两种工艺各自需要的时间, 将选用时间较短的工艺。可以一组工人用工艺 1, 一组工人用工艺 2, 随机分组。两组之间工人的个体差异会使得差异的比较误差较大。

所以, 让同一个工人分别采用两种工艺, 次序先后随机。这样不会引入个体差异的影响, 精度较高。使用成对检验。用双侧检验, 水平 0.05。

```
t.test(Matched[["Method 1"]], Matched[["Method 2"]],
 paired=TRUE, alternative = "two.sided")
```

```
##
Paired t-test
##
data: Matched[["Method 1"]] and Matched[["Method 2"]]
t = 2.1958, df = 5, p-value = 0.07952
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.05120834 0.65120834
sample estimates:
mean of the differences
```

## 0.3

### 32.3.1.3 多元均值置信域

待补充。

## 32.3.2 比例的假设检验和置信区间

### 32.3.2.1 单个比例的问题

设要比较的是总体的某个比例，如选民中对候选人甲的支持率。这时总体是两点分布总体，支持为 1，不支持为 0，参数为  $p = P(X = 1)$ 。其它的比例问题类似， $p$  是“成功概率”。

关于单总体的比例，也有双侧、右侧、左侧检验问题：

$$H_0 : p = p_0 \longleftrightarrow H_a : p \neq p_0$$

$$H_0 : p \leq p_0 \longleftrightarrow H_a : p > p_0$$

$$H_0 : p \geq p_0 \longleftrightarrow H_a : p < p_0$$

在 R 中，大样本情况下可以用 `prop.test(x, n, p=p0)` 作单个比例的检验，`n` 是样本量，`x` 是样本中符合条件（如支持）的个数，检验使用近似卡方统计量，用 `alternative` 选项选择双侧、右侧、左侧检验。这个函数也给出比例  $p$  的近似  $1 - \alpha$  置信区间，用 `conf.level` 选项指定置信度。

大样本时还有

$$\frac{\hat{p} - p}{\sqrt{p(1-p)/n}}$$

近似服从标准正态分布，可以据此计算 Z 检验。自定义的 R 函数如下：

```
prop.test.1s <- function(x, n, p=0.5, alternative="two.sided"){
 phat <- x/n
 zstat <- (phat - p)/sqrt(p*(1-p)/n)
 if(alternative=="two.sided"){ # 双侧检验
 pvalue <- 2*(1 - pnorm(abs(zstat)))
 } else if(alternative=="less"){ # 左侧检验
```

```

 pvalue <- pnorm(zstat)
 } else if(alternative=="greater"){ # 右侧检验
 pvalue <- 1 - pnorm(zstat)
 } else {
 stop("alternative unknown!")
 }

 c(stat=zstat, pvalue=pvalue)
}

```

因为比例问题本质上是一个二项分布推断问题，所以 R 还提供了 `binom.test(x, n, p=p0)` 函数进行精确检验并计算精确置信区间，但其结果略保守。用法与 `prop.test()` 类似。

### 32.3.2.2 两个比例的比较

设有  $X$  和  $Y$  两个独立的总体，如男性和女性，比较  $p_1 = P(X = 1)$  和  $p_2 = P(Y = 1)$ 。

检验问题包括：

$$H_0 : p_1 = p_2 \longleftrightarrow H_a : p_1 \neq p_2$$

$$H_0 : p_1 \leq p_2 \longleftrightarrow H_a : p_1 > p_2$$

$$H_0 : p_1 \geq p_2 \longleftrightarrow H_a : p_1 < p_2$$

在大样本情况下，可以用 R 函数 `prop.test(c(nsucc1, nsucc2), c(n1,n2), alternative=...)` 进行独立两样本比例的比较检验，其中  $n_1$  和  $n_2$  是两个样本量， $nsucc1$  和  $nsucc2$  是两个样本中符合特征的个数（个数除以样本量即样本中的比例），`alternative` 的选择也是 `"two.sided"`、`"less"`、`"greater"`。

还可以比较  $p_1 - p_2$  与某个  $\delta$ 。设两个总体中分别抽取了  $n_1$  和  $n_2$  个样本点，样本比例分别为  $\hat{p}_1$  和  $\hat{p}_2$ 。估计共同的样本比例

$$\hat{p} = \frac{n_1\hat{p}_1 + n_2\hat{p}_2}{n_1 + n_2}$$

当  $\delta = 0$  时取统计量

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

当  $p_1 = p_2$  时  $Z$  在大样本情况下近似服从标准正态分布。当  $\delta \neq 0$  时取统计量

$$Z = \frac{\hat{p}_1 - \hat{p}_2 - \delta}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}$$

当  $p_1 - p_2 = \delta$  时  $Z$  在大样本情况下近似服从  $N(0,1)$  分布。

可以定义  $Z$  检验 R 函数为:

```
prop.test.2s <- function(x, n, delta=0.0, alternative="two.sided"){
 phat <- sum(x)/sum(n)
 p <- x / n
 if(delta==0.0){
 zstat <- (p[1] - p[2])/sqrt(phat*(1-phat)*(1/n[1] + 1/n[2]))
 } else {
 zstat <- (p[1] - p[2] - delta)/sqrt(p[1]*(1-p[1])/n[1] + p[2]*(1-p[2])/n[2])
 }
 if(alternative=="two.sided"){ # 双侧检验
 pvalue <- 2*(1 - pnorm(abs(zstat)))
 } else if(alternative=="less"){ # 左侧检验
 pvalue <- pnorm(zstat)
 } else if(alternative=="greater"){ # 右侧检验
 pvalue <- 1 - pnorm(zstat)
 } else {
 stop("alternative unknown!")
 }

 c(stat=zstat, pvalue=pvalue)
}
```

独立两总体比例比较的小样本情形意义不大, 可考虑使用 Fisher 精确检验, R 函数为 `fisher.test()`。

### 32.3.2.3 比例检验的例子

#### 32.3.2.3.1 高尔夫培训女生比例检验例子

橡树溪高尔夫培训机构的女性学员比例较少，只有 20%。为了增加女性学员，设计了促销措施。一个月后调查，希望证明女性学员比例增加了。

设女性学员比例为  $p$ ,

$$H_0 : p \leq 0.20 \longleftrightarrow H_a : p > 0.20$$

取检验水平 0.05。抽查了 400 名学员，其中 100 名是女性学员， $\hat{p} = 0.25$ 。

用 `prop.test()` 检验：

```
prop.test(100, 400, p=0.20, alternative = "greater")
```

```
##
1-sample proportions test with continuity correction
##
data: 100 out of 400, null probability 0.2
X-squared = 5.9414, df = 1, p-value = 0.007395
alternative hypothesis: true p is greater than 0.2
95 percent confidence interval:
0.2149649 1.0000000
sample estimates:
p
0.25
```

检验  $p$  值为 0.0074，小于检验水平 0.05，所以女性学员比例已经显著地高于过去的 20% 了。

用自定义的大样本  $Z$  检验函数 `prop.test.1s()`：

```
prop.test.1s(100, 400, p=0.20, alternative = "greater")
```

```
stat pvalue
2.500000000 0.006209665
```

检验  $p$  值为 0.0062，与基于卡方检验的 `prop.test()` 结果略有差别。



用基于二项分布的 `binom.test()` 检验:

```
binom.test(100, 400, p=0.20, alternative = "greater")

##
Exact binomial test
##
data: 100 and 400
number of successes = 100, number of trials = 400, p-value =
0.008595
alternative hypothesis: true probability of success is greater than 0.2
95 percent confidence interval:
0.2146019 1.0000000
sample estimates:
probability of success
0.25
```

p 值为 0.0086。

### 32.3.2.3.2 报税代理分理处的错误率比较

某个代理报税的机构希望比较下属的两个分理处申报纳税返还的出错率。各自随机选取取了  $n_1 = 250$  和  $n_2 = 300$  件申报, 第一个分理处错了 35 件, 错误率  $\hat{p}_1 = 0.14$ ; 第二个分理处错了 27 件, 错误率  $\hat{p}_2 = 0.09$ 。

作双侧检验, 水平 0.10。

用 `prop.test()`:

```
prop.test(c(35,27), c(250,300), alternative = "two.sided")

##
2-sample test for equality of proportions with continuity
correction
##
data: c(35, 27) out of c(250, 300)
X-squared = 2.9268, df = 1, p-value = 0.08712
alternative hypothesis: two.sided
```

```
95 percent confidence interval:
-0.007506845 0.107506845
sample estimates:
prop 1 prop 2
0.14 0.09
```

p 值 0.087, 有显著差异, 第一处的错误率高。因为比例需要较大样本量, 所以样本量不太大时, 有人用 0.10 这样的检验水平。

用自定义的 `prop.test.2s()`:

```
prop.test.2s(c(35,27), c(250,300), alternative = "two.sided")
```

```
stat pvalue
1.84618928 0.06486473
```

### 32.3.3 方差的假设检验和置信区间

#### 32.3.3.1 单总体方差的假设检验和置信区间

方差和标准差是总体的重要数字特征, 在金融应用中标准差代表波动率, 在工业生产中标准差也是重要的质量指标, 如机床加工精度可以用标准差表示。

对总体方差  $\sigma^2$  可以做双侧、左侧、右侧检验, 与某个  $\sigma_0^2$  比较。

设总体为正态总体, 检验问题:

$$H_0 : \sigma^2 = \sigma_0^2 \longleftrightarrow H_a : \sigma^2 \neq \sigma_0^2$$

$$H_0 : \sigma^2 \geq \sigma_0^2 \longleftrightarrow H_a : \sigma^2 < \sigma_0^2$$

$$H_0 : \sigma^2 \leq \sigma_0^2 \longleftrightarrow H_a : \sigma^2 > \sigma_0^2$$

可以用如下的自定义函数计算单正态总体方差的检验:

```
var.test.1s <- function(x, n=length(x), var0, alternative="two.sided"){
 if(length(x)==1){ # 输入的是方差
 varx <- x
 } else {
 varx <- var(x)
 }
}
```

```

}
xi <- (n-1)*varx/var0

if(alternative=="less"){
 pvalue <- pchisq(xi, n-1)
} else if (alternative=="right"){
 pvalue <- pchisq(xi, n-1, lower.tail=FALSE)
} else if(alternative=="two.sided"){
 pvalue <- 2*min(pchisq(xi, n-1), pchisq(xi, n-1, lower.tail=FALSE))
}

c(statistic=xi, pvalue=pvalue)
}

```

可以输入样本数据到 `x` 中，输入  $\sigma_0^2$  到 `var0` 中，进行检验；也可以输入  $S^2$  到 `x` 中，输入样本量  $n$  到 `n` 中，输入  $\sigma_0^2$  到 `var0` 中，进行检验。

### 32.3.3.2 独立两总体方差的比较

设两个独立的正态总体的方差分别为  $\sigma_1^2, \sigma_2^2$ ，有双侧、左侧、右侧检验：

$$H_0 : \sigma_1^2 = \sigma_2^2 \longleftrightarrow H_a : \sigma_1^2 \neq \sigma_2^2$$

$$H_0 : \sigma_1^2 \geq \sigma_2^2 \longleftrightarrow H_a : \sigma_1^2 < \sigma_2^2$$

$$H_0 : \sigma_1^2 \leq \sigma_2^2 \longleftrightarrow H_a : \sigma_1^2 > \sigma_2^2$$

设两个总体分别抽取  $n_1$  和  $n_2$  个样本点，样本方差分别为  $S_x^2$  和  $S_y^2$ 。

检验统计量

$$F = \frac{S_x^2}{S_y^2}$$

在  $\sigma_1^2 = \sigma_2^2$  时  $F$  服从  $F(n_1 - 1, n_2 - 1)$  分布。

如果输入 `x, y` 为两个样本的具体值，可以用 `var.test(x, y, alternative=...)` 检验，其中 `alternative` 的选择是 `"two.sided"`、`"less"`、`"greater"`。

`var.test()` 假定总体都服从正态分布。R 中 `mood.test()` 是一种不要求正态分布假定的检验方法。`bartlett.test()` 检验多个独立正态总体方差是否全相等。

### 32.3.3.3 方差检验例子

(待补充)

## 32.3.4 拟合优度检验

### 32.3.4.1 各类比例相等的检验

设类别变量  $X$  有  $m$  个不同类别，抽取  $n$  个样本点后，各类的频数分别为  $f_1, f_2, \dots, f_m$ 。

零假设为所有类的比例相同；对立假设为各类比例不完全相同。

在零假设下，每个类的期望频数为  $n/m$ 。

检验统计量为

$$\chi^2 = \sum_{i=1}^m \frac{(f_i - \frac{n}{m})^2}{\frac{n}{m}}$$

当  $\chi^2$  超过某一临界值时拒绝零假设。在  $H_0$  成立且大样本情况下  $\chi^2$  统计量近似服从  $\chi^2(m-1)$  分布（自由度为组数减 1）。

设  $\mathbf{x}$  中包含各类的频数，R 函数 `chisq.test(x)` 作各类的总体比例相等的拟合优度卡方检验。

### 32.3.4.2 各类比例为指定值的检验

设分类变量  $X$  的各类每类有一个总体比例的假设，希望作为零假设检验。

设共有  $m$  个类，在零假设中，各类的比例分别为  $p_1, p_2, \dots, p_m$ 。

取了  $n$  个样本点组成的样本，各类的频数分别为  $f_1, f_2, \dots, f_m$ ，期望频数分别为  $np_1, np_2, \dots, np_m$ 。

检验统计量为

$$\chi^2 = \sum_{i=1}^m \frac{(f_i - np_i)^2}{np_i}$$

在  $H_0$  下成立且大样本情况下  $\chi^2$  近似服从  $\chi^2(m-1)$ ，自由度为组数减 1。因为是大样本检验法，每个组的期望频数不能低于 5，否则可以合并较小的类。

设  $\mathbf{x}$  中包含各类的频数， $\mathbf{p}$  中包含  $H_0$  假定的各类的概率，R 函数 `chisq.test(x, p)` 作各类的总体比例为指定概率的拟合优度卡方检验。

#### 32.3.4.3 带有未知参数的单分类变量的拟合优度假设检验

如果  $H_0$  下的概率有未知参数，先用最大似然估计法估计未知参数，然后将未知参数代入计算各类的概率，从而计算期望频数，计算卡方统计量。求  $p$  值时自由度要改为“组数减 1 减去未知参数个数”。

**例：**设某次选举有 5 位候选人，其中呼声最高的是甲和乙。问：甲、乙的支持率相等吗？

这不能用独立的两个比例的比较来解决，因为这两个比例是互斥的，不属于独立情况。

将类别简化为：甲，乙，其他。零假设为“甲、乙的支持率都等于  $p$ ，其他三位候选人的支持率为  $1 - 2p$ ， $p$  未知”；对立假设是甲、乙的支持率不相等。

假设调查了 1000 位选民，其中 300 名支持甲，200 名支持乙，500 名支持其他三位候选人。

在假定零假设成立的情况下先用最大似然估计法估计未知参数  $p$ ： $\hat{p} = (300 + 200)/1000/2 = 0.25$ 。这样，三个类的概率估计为  $(0.25, 0.25, 0.50)$ 。

然后，计算检验统计量：

```
chisq.test(c(300, 200, 500), c(0.25, 0.25, 0.50))
```

```
Warning in chisq.test(c(300, 200, 500), c(0.25, 0.25, 0.5)): Chi-squared
approximation may be incorrect
##
Pearson's Chi-squared test
```

```
##
data: c(300, 200, 500) and c(0.25, 0.25, 0.5)
X-squared = 3, df = 2, p-value = 0.2231
```

其中的  $p$  值所用的自由度错误, 我们需要自己计算  $p$  值:

```
res <- chisq.test(c(300, 200, 500), c(0.25, 0.25, 0.5))
```

```
Warning in chisq.test(c(300, 200, 500), c(0.25, 0.25, 0.5)): Chi-squared
approximation may be incorrect
```

```
c(statistic=res$statistic, pvalue=pchisq(res$statistic, res$parameter - 1, lower.t
```

```
statistic.X-squared pvalue.X-squared
3.00000000 0.08326452
```

检验  $p$  值为 0.08 而非原来的 0.22。

### 32.3.4.4 拟合优度检验例子

#### 32.3.4.4.1 市场占有率的检验例子

某类产品中公司 A 占 30%, 公司 B 占 50%, 公司 C 占 20%。近期 C 公司推出了新产品代替本公司的老产品。Scott Marketing Research 是一个调查公司, 公司 C 委托 Scott Marketing Research 调查分析新产品是否导致了市场占有率的变动。

设  $p_1, p_2, p_3$  分别为三个公司现在的市场占有率, 零假设为:

$$H_0 : p_1 = 0.30, p_2 = 0.50, p_3 = 0.20$$

取检验水平 0.05。抽查了 200 位消费者, 购买三个公司的产品的人数分配为:

$$(48, 98, 54)$$

卡方统计量

$$\chi^2 = \frac{(48 - 200 \times 0.3)^2}{200 \times 0.3} + \frac{(98 - 200 \times 0.5)^2}{200 \times 0.5} + \frac{(54 - 200 \times 0.2)^2}{200 \times 0.2} = 7.34$$

p 值 =  $P(\chi^2(3-1) > 7.34) = 0.02548$ , 拒绝零假设, 认为市场占有率有变化。

用 `chisq.test()` 计算:

```
chisq.test(c(48, 98, 54), p=c(0.3, 0.5, 0.2))
```

```
##
Chi-squared test for given probabilities
##
data: c(48, 98, 54)
X-squared = 7.34, df = 2, p-value = 0.02548
```

### 32.3.5 列联表独立性卡方检验

设分类变量  $X$  有  $m$  个类, 分类变量  $Y$  有  $k$  个类, 抽取了样本量为  $n$  的样本, 设  $X$  的第  $i$  类与  $Y$  的第  $j$  类交叉的频数为  $f_{ij}$ 。

零假设为  $X$  与  $Y$  相互独立, 即行变量与列变量相互独立。

交叉频数列成列联表形式, 第  $i$  行第  $j$  列为  $f_{ij}$ 。

第  $i$  行的行和为  $f_{i.}$ ,  $X$  的第  $i$  类的百分比为  $r_i = f_{i.}/n$ ;

第  $j$  列的列和为  $f_{.j}$ ,  $Y$  的第  $j$  类的百分比为  $c_j = f_{.j}/n$ 。

当  $X, Y$  独立时,  $(i, j)$  格子的期望频数为  $E_{ij} = n \times r_i \times c_j$ 。

列联表独立性检验统计量

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^k \frac{(f_{ij} - E_{ij})^2}{E_{ij}}$$

其中  $E_{ij} = nr_i c_j$ ,  $r_i$  为  $X$  第  $i$  类的百分比,  $c_j$  为  $Y$  第  $j$  类的百分比。

在  $H_0$  下  $\chi^2$  近似服从  $(m-1)(k-1)$  自由度的卡方分布。设统计量值为  $c$ , p 值为  $P(\chi^2((m-1)(k-1)) > c)$ 。

如果  $x, y$  分别是两个变量的原始观测值, `chisq.test(x, y)` 可以做列联表独立性检验。如果  $x$  保存了矩阵格式的列联表, 矩阵行名是  $X$  各个类的名称, 矩阵列名是  $Y$  各个类的名称, 则 `chisq.test(x)` 可以做列联表独立性检验。

列联表卡方检验法的检验统计量在零假设下的卡方分布是大样本情况的近似分布。如果每个变量仅有两个类，每个类的期望频数不能少于 5；如果有多个单元格，期望频数少于 5 的单元格的个数不能超过 20%，否则应该合并较小的类。

### 32.3.5.1 列联表独立性卡方检验例子

#### 32.3.5.1.1 性别与啤酒种类的独立性检验

Alber's Brewery of Tucson, Arizona 是啤酒制造与销售商。有三类啤酒产品：淡啤酒，普通啤酒，黑啤酒。

了解不同顾客喜好有利于制定更精准的销售策略。希望了解男女顾客对不同类型的偏好有没有显著差异，实际就是检验性别与啤酒类型偏好的独立性。

抽查了 150 位顾客，得到如下的列联表：

```
ctab.beer <- rbind(c(
 20, 40, 20),
 c(30,30,10))
colnames(ctab.beer) <- c("Light", "Regular", "Dark")
rownames(ctab.beer) <- c("Male", "Female")
addmargins(ctab.beer)
```

```
Light Regular Dark Sum
Male 20 40 20 80
Female 30 30 10 70
Sum 50 70 30 150
```

列联表独立性检验：

```
chisq.test(ctab.beer)

##
Pearson's Chi-squared test
##
data: ctab.beer
X-squared = 6.1224, df = 2, p-value = 0.04683
```



在 0.05 水平下认为啤酒类型偏好与性别有关。男性组的偏好分布、女性组的偏好分布、所有人的偏好分布：

```
tab2 <- round(prop.table(addmargins(ctab.beer, 1), 1), 3)
rownames(tab2)[3] <- "All"
tab2
```

```
Light Regular Dark
Male 0.250 0.500 0.250
Female 0.429 0.429 0.143
All 0.333 0.467 0.200
```

可以看出男性中对淡啤酒偏好偏少，女性中对淡啤酒偏好偏多。

### 32.3.6 非参数检验

待完成。



## Chapter 33

# R 相关与回归

### 33.1 相关分析

考虑连续型随机变量之间的关系。相关系数定义为

$$\rho(X, Y) = \frac{E[(X - EX)(Y - EY)]}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

又称 Pearson 相关系数。

$-1 \leq \rho \leq 1$ 。  $\rho$  接近于 +1 表示  $X$  和  $Y$  有正向的相关；  $\rho$  接近于 -1 表示  $X$  和  $Y$  有负向的相关。

相关系数代表的是线性相关性，对于  $X$  和  $Y$  的其它相关可能反映不出来，比如  $X \sim N(0, 1)$ ，  $Y = X^2$ ， 有  $\rho(X, Y) = 0$ 。

给定样本  $(X_i, Y_i), i = 1, 2, \dots, n$ ， 样本相关系数为

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

用散点图和散点图矩阵直观地查看变量间的相关。

例如，线性相关的模拟数据的散点图：

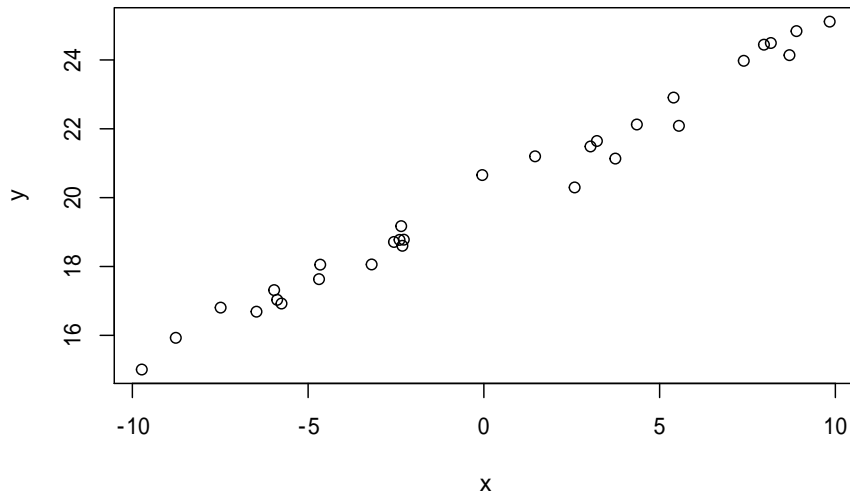


图 33.1: 线性相关数据

```
set.seed(1)
nsamp <- 30
x <- runif(nsamp, -10, 10)
y <- 20 + 0.5*x + rnorm(nsamp,0,0.5)
plot(x, y)
```

二次曲线相关的模拟数据散点图:

```
set.seed(1)
y2 <- 0.5*x^2 + rnorm(nsamp,0,2)
plot(x, y2)
```

指数关系的例子:

```
set.seed(1)
y3 <- exp(0.2*(x+10)) + rnorm(nsamp,0,2)
plot(x, y3)
```

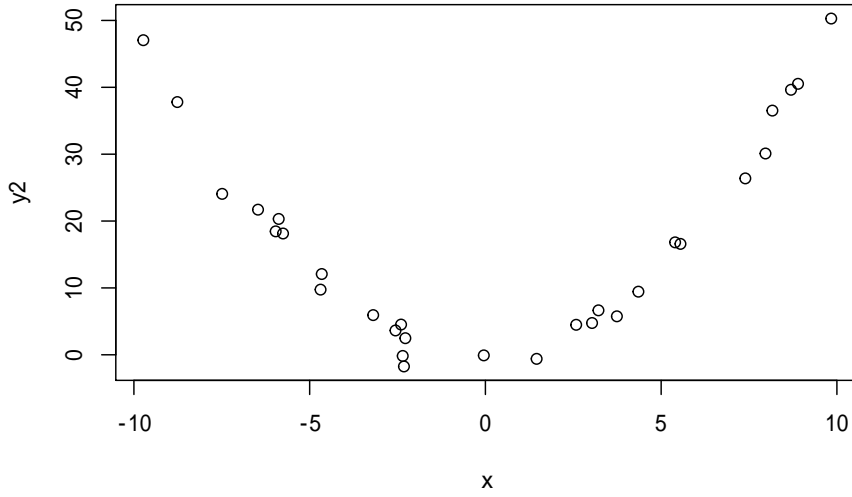


图 33.2: 二次曲线相关数据

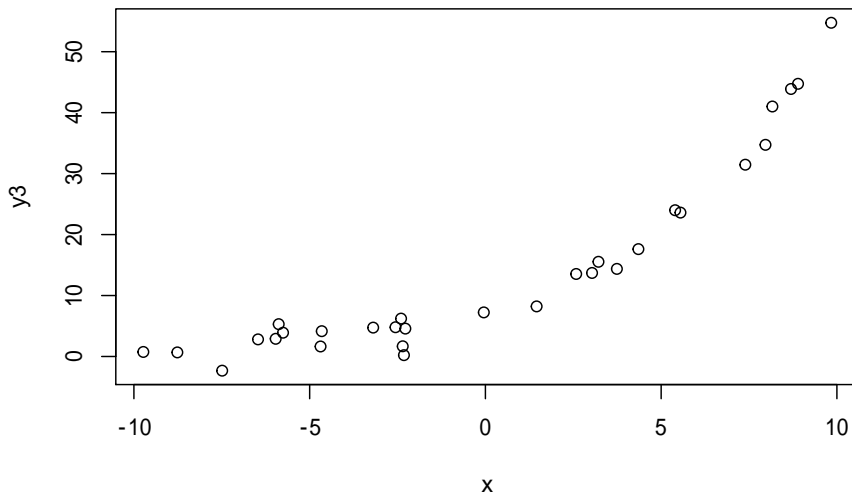


图 33.3: 指数关系相关数据

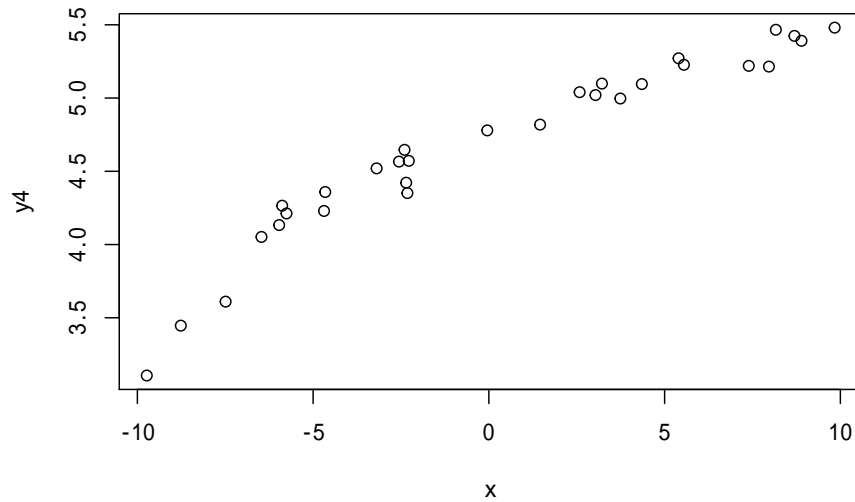


图 33.4: 对数关系相关数据

对数关系的例子:

```
set.seed(1)
y4 <- log(10*(x+12)) + rnorm(nsamp,0,0.1)
plot(x, y4)
```

### 33.1.1 相关系数的性质

- 取值  $-1 \leq r \leq 1$ 。
- $r > 0$  表示正线性相关;  $r < 0$  表示负线性相关。
- 当  $r = 1$  时  $(x_i, y_i), i = 1, \dots, n$  这  $n$  个点完全在一条正斜率的直线上, 属于确定性关系。
- 当  $r = -1$  时  $(x_i, y_i), i = 1, \dots, n$  这  $n$  个点完全在一条负斜率的直线上, 属于确定性关系。
- 计算相关系数时与  $x, y$  的次序或记号无关。
- 如果对变量  $x$  或者  $y$  分别乘以倍数, 再分别加上不同的平移量, 相关系

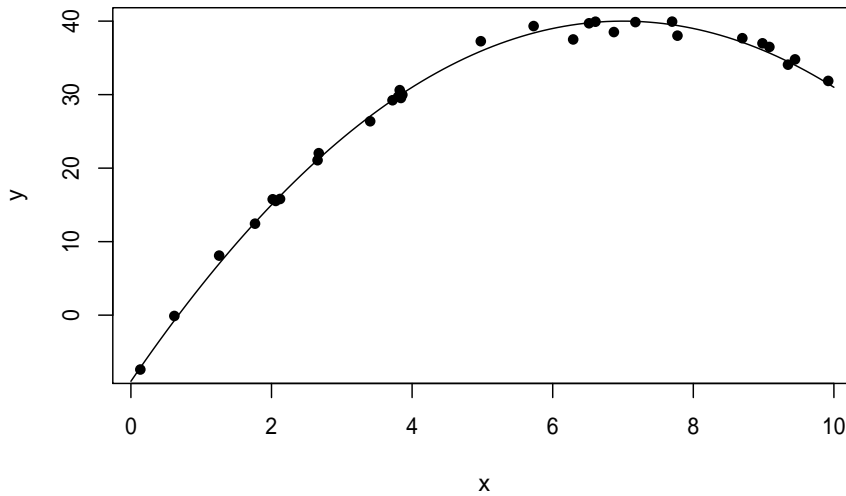


图 33.5: 曲线相关数据的相关系数例 1

数不变。即： $a + bx$  与  $c + dy$  的相关系数，等于  $x, y$  的相关系数。这样，相关系数不受单位、量纲的影响。相关系数是无量纲的。

- 如果  $x$  和  $y$  之间是非线性相关，相关系数不一定能准确反映其关系，只能作为一定程度的近似。

设变量  $X$  和  $Y$  的样本存放于 R 向量  $x$  和  $y$  中，用 `cor(x,y)` 计算样本相关系数。

```
set.seed(1)
x <- runif(30, 0, 10)
xx <- seq(0, 10, length.out = 100)
y <- 40 - (x-7)^2 + rnorm(30)
yy <- 40 - (xx-7)^2
plot(x, y, pch=16)
lines(xx, yy)
```

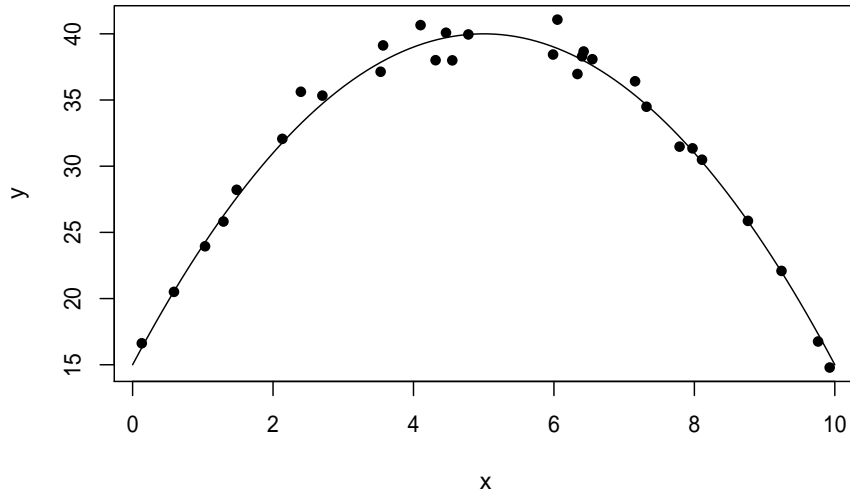


图 33.6: 曲线相关数据的相关系数例 2

```
cor(x, y)
```

```
[1] 0.8244374
```

```
x <- runif(30, 0, 10)
xx <- seq(0, 10, length.out = 100)
y <- 40 - (x-5)^2 + rnorm(30)
yy <- 40 - (xx-5)^2
plot(x, y, pch=16)
lines(xx, yy)
```

```
cor(x, y)
```

```
[1] -0.042684
```

```
x <- runif(30, 0, 10)
xx <- seq(0, 10, length.out = 100)
y <- 40*exp(-x/2) + rnorm(30)
```



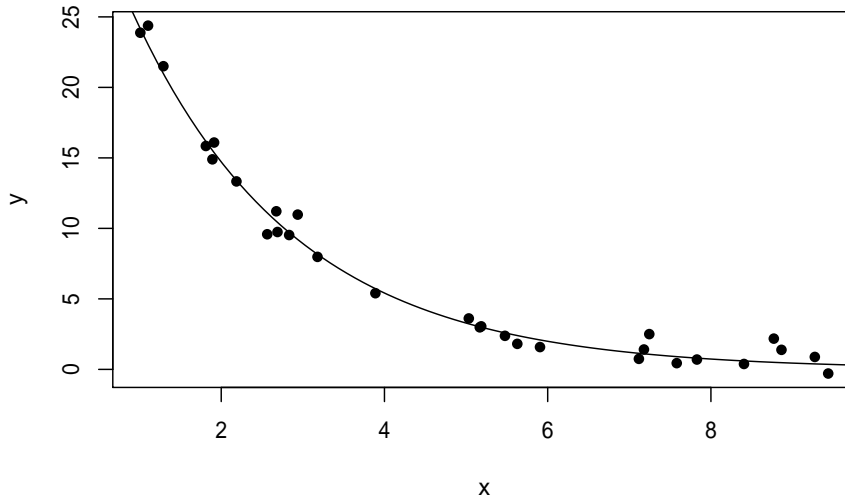


图 33.7: 曲线相关数据的相关系数例 3

```
yy <- 40*exp(-xx/2)
plot(x, y, pch=16)
lines(xx, yy)
```

```
cor(x, y)
```

```
[1] -0.8841117
```

### 33.1.2 相关与因果

相关系数是从数据中得到的  $x$  和  $y$  两个变量关系的一种描述，不能直接引申为因果关系。

例如，增加身高，不一定增长体重。

实际中有许多错误理解相关与因果性的例子。例如：研究发现喝咖啡的人群比不喝咖啡的人群心脏病发病率高。于是断言喝咖啡导致心脏病危险增加。进一步研究更多的因素发现，喝咖啡放糖很多的人心脏病发病率才增高了。

### 33.1.3 相关系数大小

- 相关系数绝对值在 0.8 以上认为高度相关。
- 在 0.5 到 0.8 之间认为中度相关。
- 在 0.3 到 0.5 之间认为低度相关。
- 在 0.3 以下认为不相关或相关性很弱以至于没有实际价值。
- 当然，在特别重要的问题中，只要经过检验显著不等于零的相关都认为是有意义的。

### 33.1.4 相关系数的检验

相关系数  $r$  是总体相关系数  $\rho$  的估计。

$$H_0 : \rho = 0 \longleftrightarrow H_a : \rho \neq 0$$

检验统计量

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

当总体  $(X, Y)$  服从二元正态分布且  $H_0 : \rho = 0$  成立时， $t$  服从  $t(n-2)$  分布。设  $t$  统计量值为  $t_0$ ， $p$  值为

$$P(|t(n-2)| > |t_0|)$$

在 R 中用 `cor.test(x,y)` 计算检验。

例如：

```
cor.test(d.class$height, d.class$weight)

##
Pearson's product-moment correlation
##
data: d.class$height and d.class$weight
t = 7.5549, df = 17, p-value = 7.887e-07
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
0.7044314 0.9523101
```

```
sample estimates:
cor
0.8777852
```

得身高与体重的相关系数为 0.88，检验的 p 值为  $7.887e-07$ ，95% 置信区间为 [0.70, 0.95]。

再例如：

```
set.seed(1)
nsamp <- 30
x <- runif(nsamp, -10, 10)
y <- 0.5*x^2 + rnorm(nsamp, 0, 2)
cor.test(x, y)

Pearson's product-moment correlation

data: x and y
t = 0.93076, df = 28, p-value = 0.3599
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
-0.1994815 0.5021658
sample estimates:
cor
0.1732378
```

得 x 与 y 的样本相关系数为 0.17，p 值为 0.36，不显著。

### 33.1.5 相关阵

如果  $x$  是一个仅包含数值型列的数据框，则 `cor(x)` 计算样本相关系数矩阵；`var(x)` 计算样本协方差阵。

`corrgram::corrgram()` 可以绘制相关系数矩阵的图形，用颜色和阴影浓度代表相关系数正负和绝对值大小，用饼图中阴影部分大小代表相关系数绝对值大小。如：

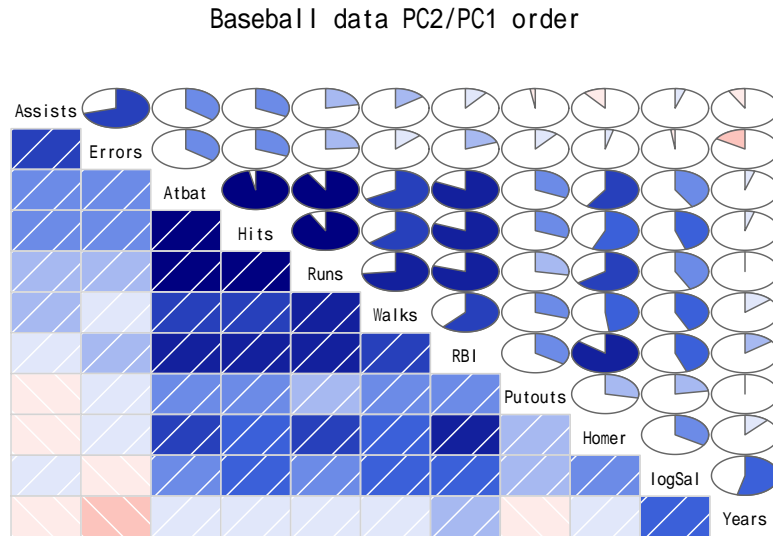


图 33.8: 相关阵图形

```
library(corrgram)

Registered S3 method overwritten by 'seriation':
method from
reorder.hclust gclus

corrgram(
 baseball[,c("Assists", "Atbat", "Errors", "Hits", "Homer", "logSal",
 "Putouts", "RBI", "Runs", "Walks", "Years")],
 order=TRUE, main="Baseball data PC2/PC1 order",
 lower.panel=panel.shade, upper.panel=panel.pie)
```

其中选项 `order=TRUE` 重排各列的次序使得较接近的列排列在相邻位置。

## 33.2 一元回归分析

### 33.2.1 模型

考虑两个变量  $Y$  与  $X$  的关系, 希望用  $X$  值的变化解释  $Y$  值的变化。  $X$  称为自变量 (independent variable),  $Y$  称为因变量 (response variable)。

假设模型

$$Y = a + bX + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

设观测值为  $(X_i, Y_i), i = 1, 2, \dots, n$ , 假设观测值满足上模型。即

$$Y_i = a + bX_i + \varepsilon_i, \quad \varepsilon_i \text{ iid } \sim N(0, \sigma^2)$$

### 33.2.2 最小二乘法

直观上看, 要找最优的直线  $y = a + bx$  使得直线与观测到的点最接近。例如:

```
plot(x, y)
abline(lm(y ~ x), col="red", lwd=2)
```

$a, b$  的解用最小二乘法得到:

$$\min_{a,b} \sum_{i=1}^n (y_i - a - bx_i)^2$$

最小二乘解表达式为

$$\hat{b} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = r_{xy} \frac{S_y}{S_x}$$

$$\hat{a} = \bar{y} - \hat{b}\bar{x}$$

其中  $r_{xy}$  为  $X$  与  $Y$  的样本相关系数,  $S_x$  与  $S_y$  分别为  $X$  和  $Y$  的样本标准差。

随机误差方差  $\sigma^2$  的估计取为

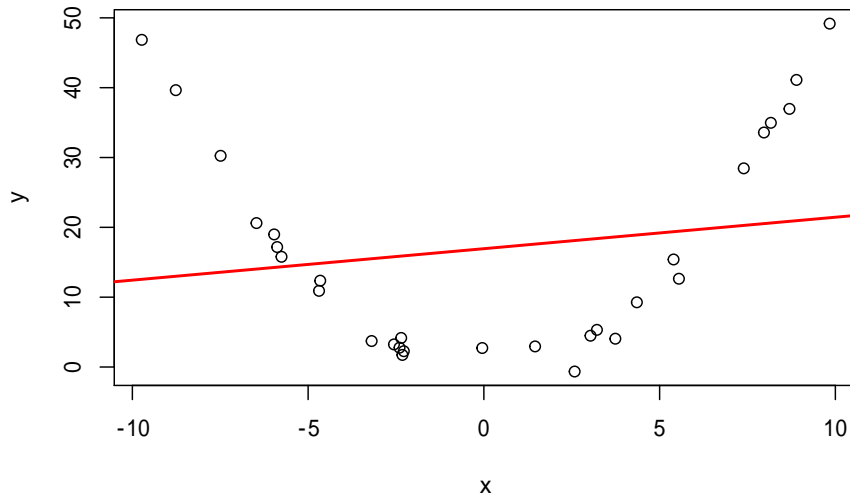


图 33.9: 一元线性回归最小二乘估计

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_i (y_i - \hat{a} - \hat{b}x_i)^2$$

标准误差: 为衡量  $\hat{a}$  和  $\hat{b}$  的估计精度, 计算  $SE(\hat{a})$  和  $SE(\hat{b})$ 。

估计结果:

- 拟合值 (预测值)

$$\hat{y}_i = \hat{a} + \hat{b}x_i, \quad i = 1, 2, \dots, n$$

- 残差

$$e_i = y_i - \hat{y}_i, \quad i = 1, 2, \dots, n$$

- $SSE = \sum_i e_i^2$  称为残差平方和, 残差平方和小则拟合优。

### 33.2.3 回归有效性

#### 33.2.3.1 拟合优度指标

按照最小二乘估计公式，只要  $x_1, x_2, \dots, x_n$  不全相等，不论  $x, y$  之间有没有相关关系，都能计算出参数估计值。得到的回归直线与观测数据的散点之间的接近程度代表了回归结果的优劣。

残差平方和

$$SSE = \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i)^2$$

残差平方和越小，说明回归直线与观测数据点吻合得越好。

$y$  是因变量，因变量的数据的离差平方和

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

代表了因变量的未知变动大小，需要对这样的变动用自变量进行解释，称 SST 为总平方和。

可以证明如下平方和分解公式：

$$SST = SSR + SSE$$

其中 SSR 称为回归平方和。

令  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  称为第  $i$  个拟合值，回归平方和为

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \hat{b}^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

在总平方和中，回归平方和是能够用回归斜率与自变量的变动解释的部分，残差平方和是自变量不能解释的变动。分解中，回归平方和越大，误差平方和越小，拟合越好。令

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

称为回归的复相关系数，或判定系数。

$0 \leq R^2 \leq 1$ 。判定系数越大，回归拟合越好。 $R^2 = 1$  时，观测点完全落在一条直线上。一元回归中  $R^2$  是  $x$  和  $y$  的样本相关系数的平方。在多元回归时只能用复相关系数。

估计误差项方差  $\sigma^2$  为

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i)^2 = \frac{1}{n-2} \text{SSE}$$

$\hat{\sigma}$  是  $\sigma$  的估计, R 的回归结果中显示为 “residual standard error” (残差标准误差)。 $\hat{\sigma}$  是残差  $e_i = y_i - \hat{y}_i$  的标准差的一个粗略估计。

### 33.2.3.2 线性关系显著性检验

当  $b = 0$  时, 模型退化为  $Y = a + \varepsilon$ ,  $X$  不出现在模型中, 说明  $Y$  与  $X$  不相关。检验

$$H_0 : b = 0 \longleftrightarrow H_a : b \neq 0$$

取统计量

$$F = \frac{\text{SSR}}{\text{SSE}/(n-2)}$$

检验 p 值为  $P(F(1, n-2) > c)$  (设  $c$  为  $F$  统计量的值)。取检验水平  $\alpha$ , p 值小于等于  $\alpha$  时拒绝  $H_0$ , 认为  $y$  与  $x$  有显著的线性相关关系, 否则认为  $y$  与  $x$  没有显著的线性相关关系。

也可以使用 t 统计量

$$t = \frac{\hat{b}}{\text{SE}_{\hat{b}}}$$

其中

$$\text{SE}_{\hat{b}} = \hat{\sigma} / \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}$$

设 t 统计量的值为  $c$ , 检验 p 值为  $P(|t(n-2)| > |c|)$ 。

### 33.2.4 R 程序

设数据保存在数据框  $d$  中, 变量名为  $y$  和  $x$ , 用 R 的 `lm()` 函数计算回归, 如:

```
set.seed(1)
nsamp <- 30
x <- runif(nsamp, -10, 10)
```



```
y <- 20 + 0.5*x + rnorm(nsamp,0,0.5)
d <- data.frame(x=x, y=y)
lm1 <- lm(y ~ x, data=d); lm1
```

```
##
Call:
lm(formula = y ~ x, data = d)
##
Coefficients:
(Intercept) x
20.0388 0.4988
```

结果只有回归系数。需要用 `summary()` 显示较详细的结果。如

```
summary(lm1)
```

```
##
Call:
lm(formula = y ~ x, data = d)
##
Residuals:
Min 1Q Median 3Q Max
-1.03030 -0.16436 -0.05741 0.29511 0.64046
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 20.03883 0.07226 277.3 <2e-16 ***
x 0.49876 0.01244 40.1 <2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 0.3956 on 28 degrees of freedom
Multiple R-squared: 0.9829, Adjusted R-squared: 0.9823
F-statistic: 1608 on 1 and 28 DF, p-value: < 2.2e-16
```

结果中  $H_0: b = 0$  的检验结果  $p$  值为  $< 2e - 16$ 。

又如对 d.class 数据集，建立体重对身高的回归方程：

```
lm2 <- lm(weight ~ height, data=d.class)
summary(lm2)

##
Call:
lm(formula = weight ~ height, data = d.class)
##
Residuals:
Min 1Q Median 3Q Max
-17.6807 -6.0642 0.5115 9.2846 18.3698
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -143.0269 32.2746 -4.432 0.000366 ***
height 3.8990 0.5161 7.555 7.89e-07 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 11.23 on 17 degrees of freedom
Multiple R-squared: 0.7705, Adjusted R-squared: 0.757
F-statistic: 57.08 on 1 and 17 DF, p-value: 7.887e-07
```

身高对体重的影响是显著的 (p 值  $7.89e-7$ )。

### 33.2.5 回归诊断

除了系数的显著性检验以外，对误差项是否符合回归模型假定中的独立性、方差齐性、正态性等，可以利用回归残差进行一系列的回归诊断。还可以计算一些异常值、强影响点的诊断。

对  $n$  组观测值，回归模型为

$$y_i = a + bx_i + \varepsilon_i$$

拟合值为

$$\hat{y}_i = \hat{a} + \hat{b}x_i$$

残差 (residual) 为

$$e_i = y_i - \hat{y}_i$$

残差相当于模型中的随机误差，但仅在模型假定成立时才是随机误差的合理估计。所以残差能够反映违反模型假定的情况。

$e_i$  有量纲，不好比较大小。定义标准化残差（内部学生化残差）：

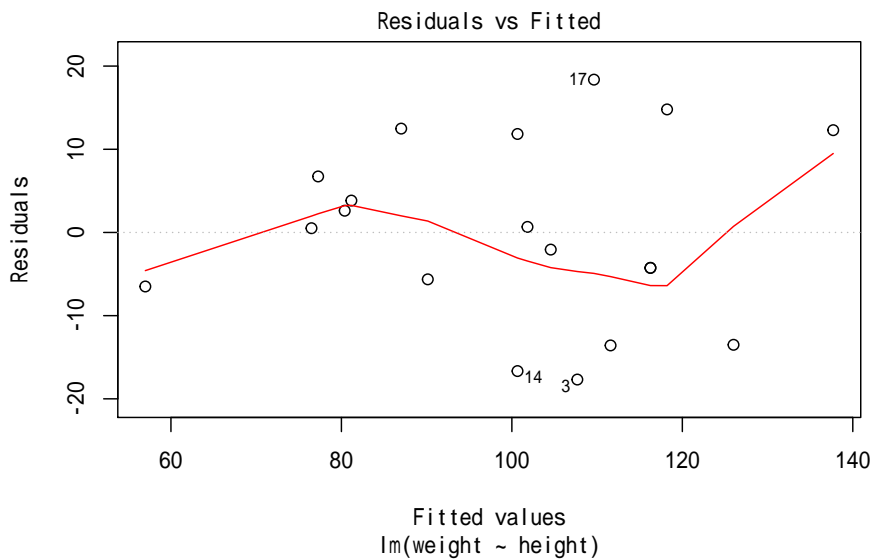
$$r_i = \frac{e_i}{S_e \sqrt{1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{\sum_{k=1}^n (x_k - \bar{x})^2}}}$$

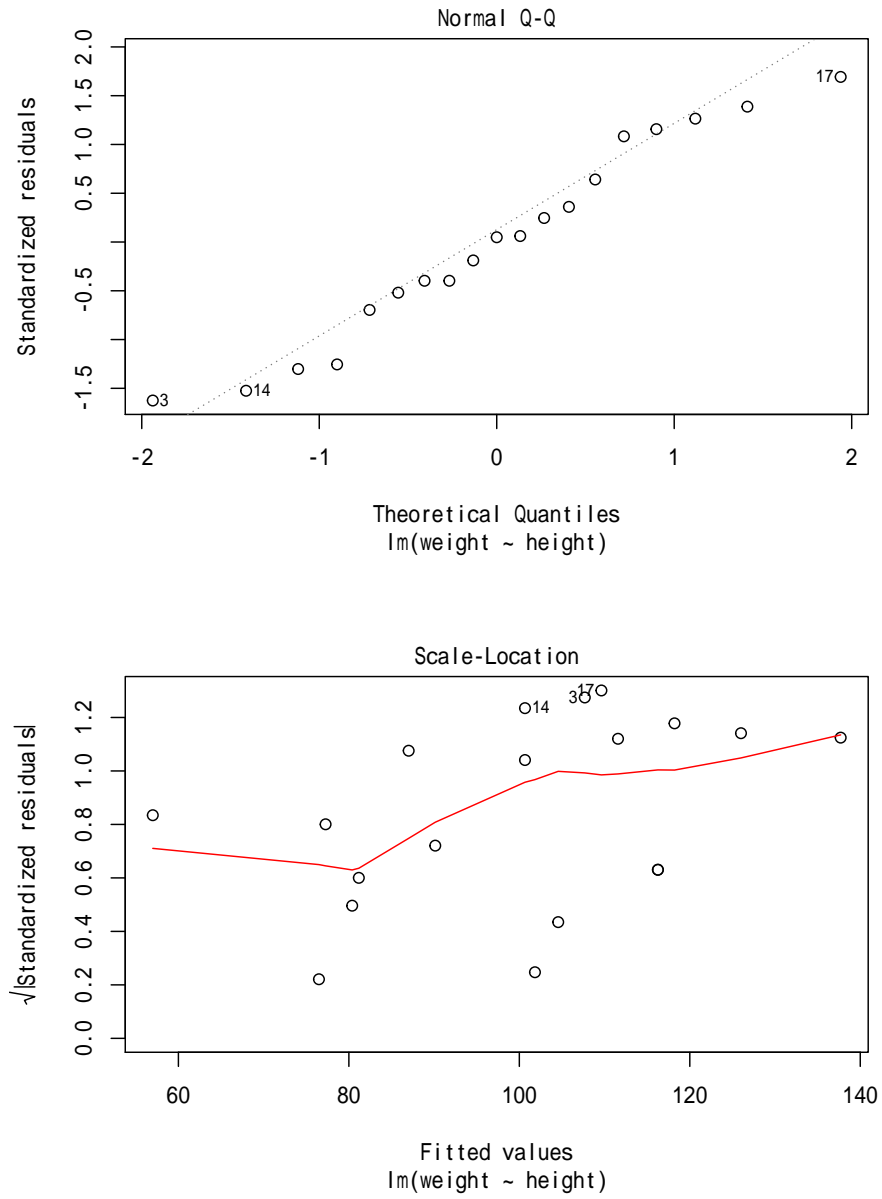
样本量较大时标准化残差近似服从标准正态分布，于是取值于  $[-2, 2]$  范围外的点是可疑的离群点。

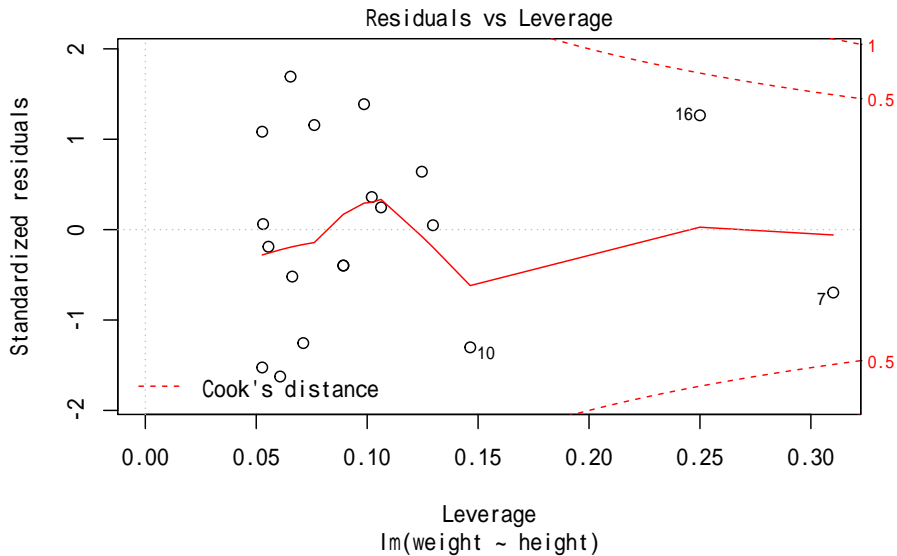
R 中用 `residuals()` 从回归结果计算残差，用 `rstandard()` 从回归结果计算标准化残差。

R 中 `plot(lmres)` (`lmres` 为回归结果变量) 可以做简单回归诊断。

```
plot(lm2)
```







第一个图是残差对预测值散点图，散点应该随机在 0 线上下波动，不应该有曲线模式、分散程度增大模式、特别突出的离群点等情况。

第二个图是残差的正态 QQ 图，散点接近于直线时可以认为模型误差项的正态分布假定是合理的。

第三个图是误差大小 (标准化残差绝对值的平方根) 对拟合值的图形，可以判断方差齐性假设 (方差  $\sigma^2$  不变)。

第四个图是残差对杠杆量图，并叠加了 Cook 距离等值线。杠杆量代表了回归自变量对结果的影响大小，超过  $4/n$  的值是需要重视的。Cook 距离考察删去第  $i$  个观测对回归结果的影响。

```
cbind(residuals(lm2), rstandard(lm2))
```

```
[,1] [,2]
1 6.7317083 0.64090788
2 -13.5797581 -1.25514370
3 -17.6807278 -1.62510342
4 0.5115143 0.04884012
5 -5.6350916 -0.51945382
```

```
6 -4.2585944 -0.39749776
7 -6.4933343 -0.69635607
8 11.8375266 1.08337723
9 0.6678176 0.06113186
10 -13.5061701 -1.30222590
11 -2.0615036 -0.18894975
12 14.7918904 1.38780152
13 2.6124840 0.24615612
14 -16.6624734 -1.52495912
15 12.4841326 1.15698074
16 12.2967391 1.26478834
17 18.3697570 1.69265472
18 3.8326780 0.36028625
19 -4.2585944 -0.39749776
```

可以用  $e_i$  或  $r_i$  作为纵坐标,  $x_i$  或者  $\hat{y}_i$  作为横坐标, 作残差图。

下面给出残差图的几种常见的缺陷:

非线性:

```
set.seed(1)
x <- runif(30, 0, 10)
xx <- seq(0, 10, length.out = 100)
y <- 40 - (x-7)^2 + rnorm(30)
yy <- 40 - (xx-7)^2
lms1 <- lm(y ~ x)
opar <- par(mfrow=c(1,2))
plot(x, y, pch=16, main=" 数据和真实模型")
lines(xx, yy)
plot(x, rstandard(lms1), main=" 使用线性回归的残差")
```

```
par(opar)
```

异方差:

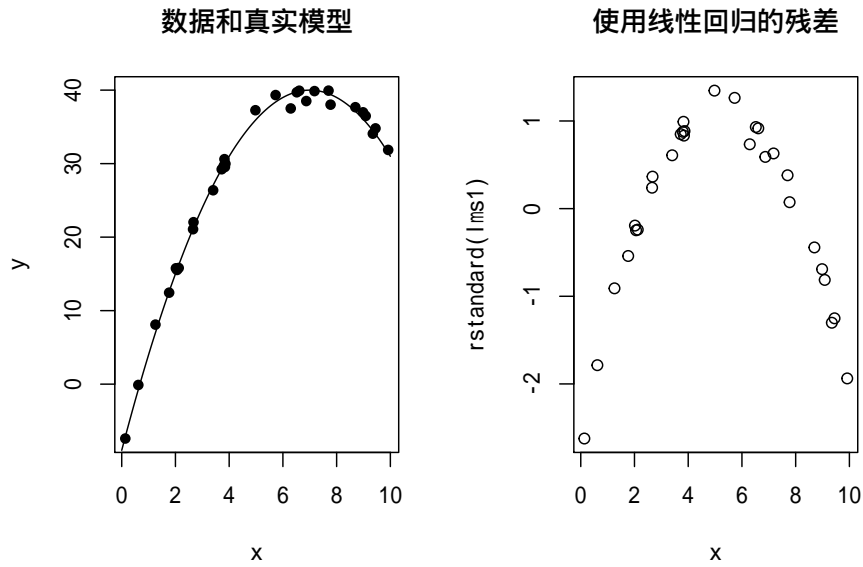
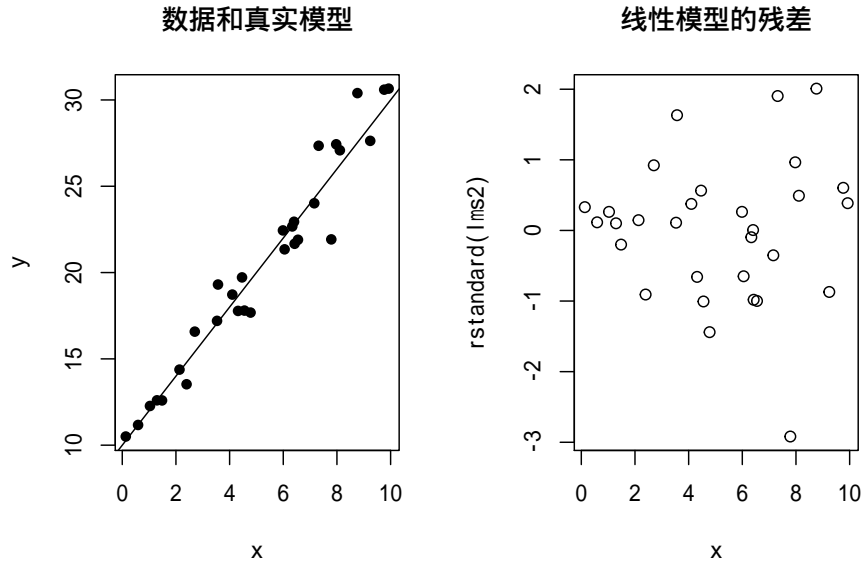


图 33.10: 残差中非线性模式

```
x <- sort(runif(30, 0, 10))
y <- 10 + 2*x + rnorm(30)*(seq(30)/10)
lms2 <- lm(y ~ x)
opar <- par(mfrow=c(1,2))
plot(x, y, pch=16, main=" 数据和真实模型")
abline(a=10, b=2)
plot(x, rstandard(lms2), main=" 线性模型的残差")
```



```
par(opar)
```

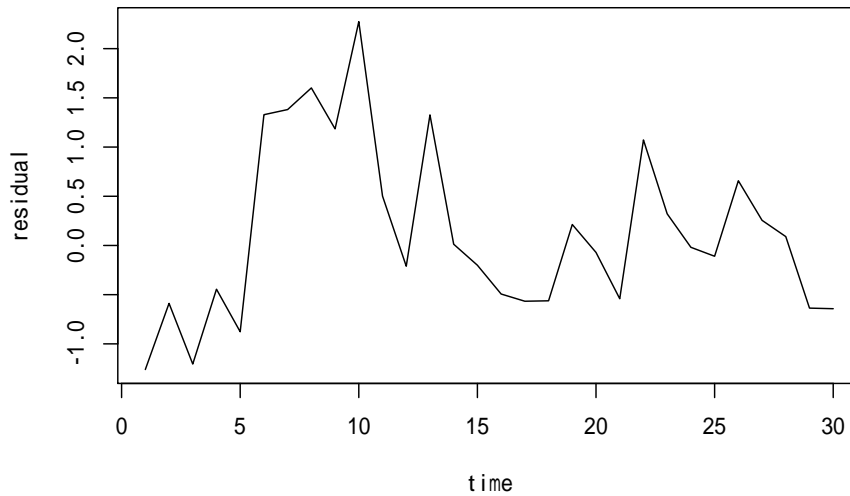
序列自相关:

当数据随时间变化时, 还需要考虑前后是否有序列自相关。自相关残差图形例子:

```
ar1 <- arima.sim(list(ar=c(0.5)), n=30)
plot(1:30, ar1[], type="l", xlab="time", ylab="residual",
 main="有序列自相关的残差图形")
```



有序列自相关的残差图形



图中横坐标是观测序号。

`car::ncvTest()` 检验方差齐性。零假设是方差齐性成立。如

```
car::ncvTest(lm2)
```

```
Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 1.664307, Df = 1, p = 0.19702
```

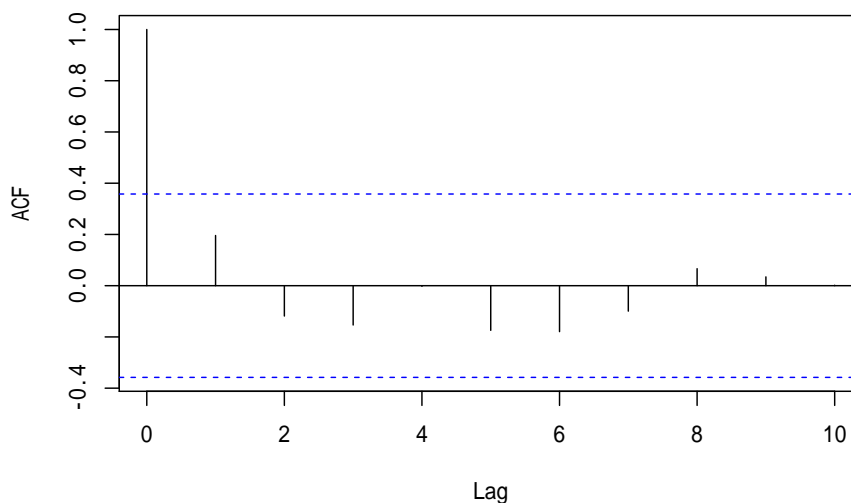
用 Durbin-Watson 检验 (DW 检验) 可以检验残差中是否有序列自相关, 零假设是没有序列自相关。R 中用 `car::durbinWatsonTest()` 检验。序列自相关检验仅当数据中数据是延时间等间隔记录时有意义。如

```
car::durbinWatsonTest(lm1)
```

```
lag Autocorrelation D-W Statistic p-value
1 0.1959139 1.57293 0.21
Alternative hypothesis: rho != 0
```

也可以对回归残差绘制 ACF 图, 如果除了横坐标 0 之外都落在两条水平界限内则认为没有序列自相关, 如果有明显超出界限的就认为有序列自相关。如

```
acf(residuals(lm1), lag.max = 10, main="")
```



### 33.2.6 预测区间

设  $X$  取  $x_0$ ,  $Y$  的预测值为  $\hat{y}_0 = \hat{a} + \hat{b}x_0$ 。置信水平为  $1 - \alpha$  的预测区间为

$$\hat{y}_0 \pm \lambda \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum_i (x_i - \bar{x})^2}}$$

其含义是  $P(y_0 \in \text{预测区间}) = 1 - \alpha$ 。 $\lambda$  是标准正态分布双侧  $\alpha$  分位数  $qnorm(1 - \alpha/2)$ 。

用 `predict(lmres)` 得到  $\hat{y}_i, i = 1, \dots, n$  的值, `lmres` 表示回归结果变量。

用 `predict(lmres, interval="prediction")` 同时得到预测的置信区间, 需要的话加入 `level=` 选项设置置信度。

回归的置信区间有两种, 上述的区间是“预测区间”(CLI), 使得  $P(y_0 \in \text{预测区间}) = 1 - \alpha$ ; 另一种区间称为均值的置信区间 (CLM), 使得  $P(Ey_0 \in \text{均值置信区间}) = 1 - \alpha$ 。

### 33.2.7 控制

如果需要把  $Y$  的值控制在  $[y_l, y_u]$  范围内，问如何控制  $X$  的范围，可以求解  $x_0$  的范围使上面的置信区间包含在  $[y_l, y_u]$  内。

近似地可以解不等式

$$\hat{a} + \hat{b}x_0 - z_{1-\frac{\alpha}{2}}\hat{\sigma} \geq y_l$$

$$\hat{a} + \hat{b}x_0 + z_{1-\frac{\alpha}{2}}\hat{\sigma} \leq y_u$$

其中  $z_{1-\frac{\alpha}{2}}$  为标准正态分布双侧  $\alpha$  分位数（用 `qnorm(1-alpha/2)` 计算）。

## 33.3 多元线性回归

建模步骤：

- 确定要研究的因变量，以及可能对因变量有影响并且数据可以获得的自变量集合
- 假定因变量  $y$  与  $p$  个自变量之间为线性相关关系，建立模型
- 对模型进行评估和检验
- 判断模型中是否存在多重共线性，如果存在，应进行处理
- 预测
- 残差分析

### 33.3.1 模型

模型

$$y = \beta_0 + \beta_1x_1 + \cdots + \beta_px_p + \varepsilon$$

其中

- $y$ : 因变量，随机变量
- $x_1, \dots, x_p$ : 自变量，非随机
- $\varepsilon$ : 随机误差项，随机变量
- $\beta_0$ : 截距项
- $\beta_j$ : 对应于  $x_j$  的斜率项
- 模型表明因变量  $y$  近似等于自变量的线性组合

- $E\varepsilon = 0, \text{Var}(\varepsilon) = \sigma^2$

对  $n$  组观测数据, 有

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i, \quad i = 1, 2, \dots, n$$

对其中的随机误差项  $\varepsilon_i, i = 1, 2, \dots, n$ , 假定:

- 期望为零, 服从正态分布;
- 方差齐性: 方差为  $\sigma^2$  与自变量值无关;
- 相互独立。

总之,  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  相互独立, 服从  $N(0, \sigma^2)$  分布。

数据格式如:

$$\left( \begin{array}{cccc|c} x_1 & x_2 & \cdots & x_p & y \\ \hline x_{11} & x_{12} & \cdots & x_{1p} & y_1 \\ x_{21} & x_{22} & \cdots & x_{2p} & y_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} & y_n \end{array} \right)$$

R 中回归数据放在一个数据框中, 有一列  $y$  和  $p$  列自变量。

根据模型有

$$Ey = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$$

### 33.3.2 参数估计

未知的参数有回归系数  $\beta_0, \beta_1, \dots, \beta_p, \sigma^2$ , 另外误差项也是未知的。

模型估计仍使用最小二乘法, 得到系数估计  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  及误差方差估计  $\hat{\sigma}^2$ 。得到估计的多元线性回归方程:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_p x_p$$

对第  $i$  组观测值, 将自变量值代入估计的回归方程中得拟合值

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_p x_{ip}$$

$e_i = y_i - \hat{y}_i$  称为残差。回归参数估计，用残差最小作为目标，令

$$Q = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip})]^2$$

取使得  $Q$  达到最小的  $(\beta_0, \beta_1, \dots, \beta_p)$  作为参数估计，称为最小二乘估计，记为  $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ 。称得到的最小值为 SSE， $\sigma^2$  的估计为

$$\hat{\sigma}^2 = \frac{1}{n-p-1} \text{SSE} = \frac{1}{n-p-1} \sum_{i=1}^n e_i^2$$

### 33.3.3 R 的多元回归程序

在 R 中用 `lm(y ~ x1 + x2 + x3, data=d)` 这样的程序来做多元回归，数据集为 `d`，自变量为 `x1,x2,x3` 三列。

#### 33.3.3.1 例：体重对身高和年龄的回归

考虑 `d.class` 数据集中体重对身高和年龄的回归：

```
lm3 <- lm(weight ~ height + age, data=d.class)
summary(lm3)

##
Call:
lm(formula = weight ~ height + age, data = d.class)
##
Residuals:
Min 1Q Median 3Q Max
-17.962 -6.010 -0.067 7.553 20.796
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -141.2238 33.3831 -4.230 0.000637 ***
height 3.5970 0.9055 3.973 0.001093 **
age 1.2784 3.1101 0.411 0.686492

```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 11.51 on 16 degrees of freedom
Multiple R-squared: 0.7729, Adjusted R-squared: 0.7445
F-statistic: 27.23 on 2 and 16 DF, p-value: 7.074e-06
```

得到的回归模型可以写成

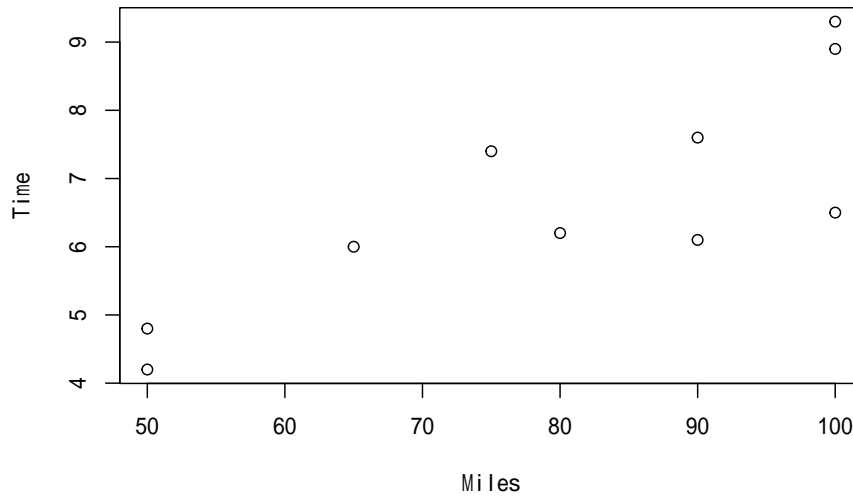
$$\widehat{\text{weight}} = -141.2 + 3.597 \times \text{height} + 1.278 \times \text{age}$$

其中身高的系数 3.597 表示, 如果两个学生年龄相同, 则身高增加 1 个单位 (这里是英寸), 体重平均增加 3.597 个单位 (这里是磅)。注意在仅有身高作为自变量时, 系数为 3.899。年龄的系数 1.278 也类似解释, 在两个学生身高相同时, 如果一个学生年龄大 1 岁, 则此学生的体重平均多 1.278 个单位。

### 33.3.3.2 例: 驾驶员行驶时间的模型

Butler Trucking Company 是一个短途货运公司。老板想研究每个驾驶员每天的行驶时间的影响因素。随机抽取了 10 名驾驶员一天的数据。考虑行驶里程 (Miles) 对行驶时间的影响。

```
with(Butler, plot(Miles, Time))
```



从散点图看，行驶里程与行驶时间有线性相关关系。作一元线性回归：

```
lmbut01 <- lm(Time ~ Miles, data=Butler)
summary(lmbut01)
```

```
##
Call:
lm(formula = Time ~ Miles, data = Butler)
##
Residuals:
Min 1Q Median 3Q Max
-1.5565 -0.4913 0.1783 0.7120 1.2435
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 1.27391 1.40074 0.909 0.38969
Miles 0.06783 0.01706 3.977 0.00408 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
Residual standard error: 1.002 on 8 degrees of freedom
Multiple R-squared: 0.6641, Adjusted R-squared: 0.6221
F-statistic: 15.81 on 1 and 8 DF, p-value: 0.00408
```

这里 Miles 的系数解释为：行驶里程增加 1 英里，时间平均增加 0.068 小时。老板想进一步改进对行驶时间的预测，增加“派送地点数”（Deliveries）作为第二个自变量。作多元回归：

```
lmbut02 <- lm(Time ~ Miles + Deliveries, data=Butler)
summary(lmbut02)
```

```
##
Call:
lm(formula = Time ~ Miles + Deliveries, data = Butler)
##
Residuals:
Min 1Q Median 3Q Max
-0.79875 -0.32477 0.06333 0.29739 0.91333
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.868701 0.951548 -0.913 0.391634
Miles 0.061135 0.009888 6.182 0.000453 ***
Deliveries 0.923425 0.221113 4.176 0.004157 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 0.5731 on 7 degrees of freedom
Multiple R-squared: 0.9038, Adjusted R-squared: 0.8763
F-statistic: 32.88 on 2 and 7 DF, p-value: 0.0002762
```

这里 Miles 的系数从一元回归时的 0.068 改成了 0.061，解释为：在派送地点数相同的条件下，行驶里程每增加 1 英里，行驶时间平均增加 0.061 小时。

Deliveries 的系数解释为：在行驶里程相同的条件下，派送地点数每增加一个地



点，行驶时间平均增加 0.92 小时。

### 33.3.4 模型的检验

#### 33.3.4.1 复相关系数平方

将总平方和分解为：

$$SST = SSR + SSE$$

其中 - 总平方和

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

- 回归平方和

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

是能够用回归系数和自变量的变量解释的因变量变化。- 残差平方和

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

是回归模型不能解释的因变量变化。

回归平方和越大，残差平方和越小，回归拟合越好。定义复相关系数平方（判定系数）

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

则  $0 \leq R^2 \leq 1$ 。  $R^2$  越大，说明数据中的自变量拟合因变量值拟合越好。

但是，“拟合好”不是唯一标准。仅考虑拟合好，可能产生很复杂的仅对建模用数据有效但是对其它数据无效的模型，这称为“过度拟合”。

定义调整的（修正的）复相关系数平方：

$$R^{*2} = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

克服了  $R^2$  的部分缺点。

在体重与身高、年龄的模型结果中，  $R^2 = 0.7729$ ，  $R^{*2} = 0.7445$ 。

### 33.3.4.2 残差标准误差

模型中  $\varepsilon$  的方差  $\sigma^2$  的估计为  $\hat{\sigma}_e^2$ ,

$$\hat{\sigma}^2 = \frac{1}{n-p-1} \text{SSE}$$

$\hat{\sigma}$  是随机误差的标准差  $\sigma$  的估计量, 称为“残差标准误差”(Residual standard error)。这是残差  $e_i = y_i - \hat{y}_i$  的标准差的一个较粗略的近似估计。

在体重与身高、年龄的模型结果中,  $\hat{\sigma} = 11.51$ 。

### 33.3.4.3 线性关系检验

为了检验整个回归模型是否都无效, 考虑假设检验:

$$H_0 : \beta_1 = \cdots = \beta_p = 0$$

当  $H_0$  成立时模型退化成  $y = \beta_0 + \varepsilon$ ,  $y$  与  $x_1, x_2, \dots, x_p$  之间不再具有线性相关关系。

取统计量为

$$F = \frac{\text{SSR}/p}{\text{SSE}/(n-p-1)}$$

在模型前提条件都满足且  $H_0$  成立时  $F$  服从  $F(p, n-p-1)$  分布。计算右侧  $p$  值, 给出检验结论。

当检验显著时, 各斜率项不全为零。结果不显著时, 当前回归模型不能使用。

在体重与身高、年龄的模型结果中,  $F = 27.23$ ,  $p$  值为  $7.074 \times 10^{-6}$ , 在 0.05 水平下显著, 模型有意义。

### 33.3.4.4 单个斜率项的显著性检验

为了检验某一个自变量  $X_j$  是否对因变量的解释有贡献 (在模型中已经包含了其它自变量的情况下), 检验

$$H_0 : \beta_j = 0$$

如果不拒绝  $H_0$ ，相当于  $x_j$  可以不出现在模型中。但是，在多元线性回归中这不代表  $x_j$  与  $y$  之间没有线性相关，可能是由于其它的自变量包含了  $x_j$  中的信息。

检验使用如下  $t$  统计量

$$t = \frac{\hat{\beta}_j}{\text{SE}_{\hat{\beta}_j}}$$

在模型前提条件都满足且  $H_0$  成立时  $t$  服从  $t(n - p - 1)$  分布。给定检验水平  $\alpha$ ，计算双侧  $p$  值，做出检验结论。

对单个回归系数的检验，检验水平可取得略高，如 0.10, 0.15。

在体重与身高、年龄的模型结果中，身高的斜率项显著性  $p$  值为 0.001，在 0.15 水平下显著；体重的斜率项显著性  $p$  值为 0.686，在 0.15 水平下不显著，可考虑从模型中去掉体重变量。

### 33.3.5 回归自变量筛选

在 19 个学生的体重与身高、年龄的线性回归模型结果中，发现关于年龄的系数为零的检验  $p$  值为 0.686，不显著，说明在模型中已经包含身高的情况下，年龄不提供对体重的额外信息。

但是如果体重对年龄单独建模的话，年龄的影响还是显著的：

```
lm4 <- lm(weight ~ age, data=d.class)
summary(lm4)

##
Call:
lm(formula = weight ~ age, data = d.class)
##
Residuals:
Min 1Q Median 3Q Max
-23.349 -7.609 -5.260 7.945 42.847
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
```

```
(Intercept) -50.493 33.290 -1.517 0.147706
age 11.304 2.485 4.548 0.000285 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 15.74 on 17 degrees of freedom
Multiple R-squared: 0.5489, Adjusted R-squared: 0.5224
F-statistic: 20.69 on 1 and 17 DF, p-value: 0.0002848
```

模型中不显著的自变量应该逐一剔除。可以用 `step` 函数进行逐步回归变量选择，如：

```
lm5 <- step(lm(weight ~ height + age + sex, data=d.class))
```

```
Start: AIC=94.93
weight ~ height + age + sex
##
Df Sum of Sq RSS AIC
- age 1 113.76 1957.8 94.067
<none> 1844.0 94.930
- sex 1 276.09 2120.1 95.581
- height 1 1020.61 2864.6 101.299
##
Step: AIC=94.07
weight ~ height + sex
##
Df Sum of Sq RSS AIC
- sex 1 184.7 2142.5 93.780
<none> 1957.8 94.067
- height 1 5696.8 7654.6 117.974
##
Step: AIC=93.78
weight ~ height
##
Df Sum of Sq RSS AIC
```

```
<none> 2142.5 93.78
- height 1 7193.2 9335.7 119.75

summary(lm5)

##
Call:
lm(formula = weight ~ height, data = d.class)
##
Residuals:
Min 1Q Median 3Q Max
-17.6807 -6.0642 0.5115 9.2846 18.3698
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -143.0269 32.2746 -4.432 0.000366 ***
height 3.8990 0.5161 7.555 7.89e-07 ***

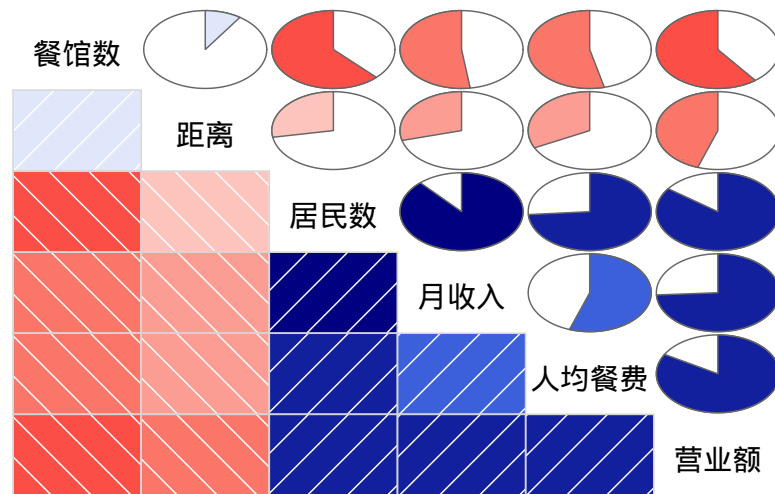
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 11.23 on 17 degrees of freedom
Multiple R-squared: 0.7705, Adjusted R-squared: 0.757
F-statistic: 57.08 on 1 and 17 DF, p-value: 7.887e-07
```

考虑另一个例子数据。在数据框 `Resturant` 中包含 25 个餐馆的一些信息，包括：

- $y$ , 营业额, 日均营业额 (万元)
- $x_1$ , 居民数, 周边居民人数 (万人)
- $x_2$ , 人均餐费, 用餐平均支出 (元/人)
- $x_3$ , 月收入, 周边居民月平均收入 (元)
- $x_4$ , 餐馆数, 周边餐馆数 (个)
- $x_5$ , 距离, 距市中心距离 (km)

对营业额进行回归建模研究。变量间的相关系数图：

```
corrgram::corrgram(
 Resturant, order=TRUE,
 lower.panel=corrgram::panel.shade,
 upper.panel = corrgram::panel.pie,
 text.panel = corrgram::panel.txt)
```



```
lmrst01 <- lm(`营业额` ~ `居民数` + `人均餐费` + `月收入` + `餐馆数` + `距离`,
 data=Resturant)
summary(lmrst01)
```

```
##
Call:
lm(formula = 营业额 ~ 居民数 + 人均餐费 + 月收入 + 餐馆数 + 距离,
data = Resturant)
##
Residuals:
Min 1Q Median 3Q Max
-16.7204 -6.0600 0.7152 3.2144 21.4805
##
```

```
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 4.2604768 10.4679833 0.407 0.68856
居民数 0.1273254 0.0959790 1.327 0.20037
人均餐费 0.1605660 0.0556834 2.884 0.00952 **
月收入 0.0007636 0.0013556 0.563 0.57982
餐馆数 -0.3331990 0.3986248 -0.836 0.41362
距离 -0.5746462 0.3087506 -1.861 0.07826 .

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 10.65 on 19 degrees of freedom
Multiple R-squared: 0.8518, Adjusted R-squared: 0.8128
F-statistic: 21.84 on 5 and 19 DF, p-value: 2.835e-07
```

在 0.15 水平上只有人均餐费和到市中心距离是显著的。进行逐步回归筛选：

```
lmrst02 <- step(lmrst01)
```

```
Start: AIC=123.39
营业额 ~ 居民数 + 人均餐费 + 月收入 + 餐馆数 + 距离
##
Df Sum of Sq RSS AIC
- 月收入 1 35.96 2189.0 121.81
- 餐馆数 1 79.17 2232.2 122.30
<none> 2153.0 123.39
- 居民数 1 199.42 2352.4 123.61
- 距离 1 392.54 2545.6 125.58
- 人均餐费 1 942.22 3095.2 130.47
##
Step: AIC=121.81
营业额 ~ 居民数 + 人均餐费 + 餐馆数 + 距离
##
Df Sum of Sq RSS AIC
- 餐馆数 1 78.22 2267.2 120.69
```

```

<none> 2189.0 121.81
- 距离 1 445.69 2634.7 124.44
- 人均餐费 1 925.88 3114.9 128.63
- 居民数 1 1133.27 3322.3 130.24
##
Step: AIC=120.69
营业额 ~ 居民数 + 人均餐费 + 距离
##
Df Sum of Sq RSS AIC
<none> 2267.2 120.69
- 距离 1 404.28 2671.5 122.79
- 人均餐费 1 1050.90 3318.1 128.21
- 居民数 1 1661.83 3929.0 132.43

```

`summary(lmrst02)`

```

##
Call:
lm(formula = 营业额 ~ 居民数 + 人均餐费 + 距离, data = Resturant)
##
Residuals:
Min 1Q Median 3Q Max
-14.027 -5.361 -1.560 2.304 23.001
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.68928 6.25242 -0.270 0.78966
居民数 0.19022 0.04848 3.923 0.00078 ***
人均餐费 0.15763 0.05052 3.120 0.00518 **
距离 -0.56979 0.29445 -1.935 0.06656 .

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 10.39 on 21 degrees of freedom

```



```
Multiple R-squared: 0.8439, Adjusted R-squared: 0.8216
F-statistic: 37.85 on 3 and 21 DF, p-value: 1.187e-08
```

最后进保留了周边居民人数、人均餐费、到市中心距离三个自变量，删去了周边居民月收入和周边餐馆数两个自变量。

### 33.3.6 哑变量与变截距项的模型

#### 33.3.6.1 哑变量

回归分析的因变量是连续型的，服从正态分布。回归分析的自变量是数值型的，可以连续取值也可以离散取值。但是，如果自变量是类别变量，用简单的 1,2,3 编码并不能正确地表现不同类别的作用。可以将类别变量转换成一个或者多个取 0、1 值的变量，称为哑变量 (dummy variables) 或虚拟变量。

##### 33.3.6.1.1 二值哑变量与平行线模型

如果  $f$  是一个二值分类变量，将其中一个类编码为 1，另一个类编码为零。例如， $f = 1$  表示处理组， $f = 0$  表示对照组。这样编码后二值分类变量可以直接用在回归模型中。

例如

$$Ey = \beta_0 + \beta_1 x + \beta_2 f$$

相当于

$$Ey = \begin{cases} \beta_0 + \beta_1 x, & \text{对照组} \\ (\beta_0 + \beta_2) + \beta_1 x, & \text{处理组} \end{cases}$$

这样的模型叫做平行线模型，处理组的数据与对照组的数据服从斜率相同、截距不同的一元线性回归模型。比每个组单独建模更有效。

如果检验：

$$H_0 : \beta_2 = 0$$

则不显著时，可以认为在处理组和对照组中  $y$  与  $x$  的关系没有显著差异。

进一步考虑

$$Ey = \beta_0 + \beta_1 x + \beta_2 f + \beta_3 f x$$

相当于

$$E y = \begin{cases} \beta_0 + \beta_1 x, & \text{对照组} \\ (\beta_0 + \beta_2) + (\beta_1 + \beta_3)x, & \text{处理组} \end{cases}$$

比每个组单独建立回归模型更有效。可以检验

$$H_0 : \beta_3 = 0$$

不显著时可以用平行线模型。

比如，为了表示男生和女生的体重有不同，可以在以体重为因变量的回归中加入自变量性别：

```
lm6 <- lm(weight ~ height + sex, data=d.class)
summary(lm6)

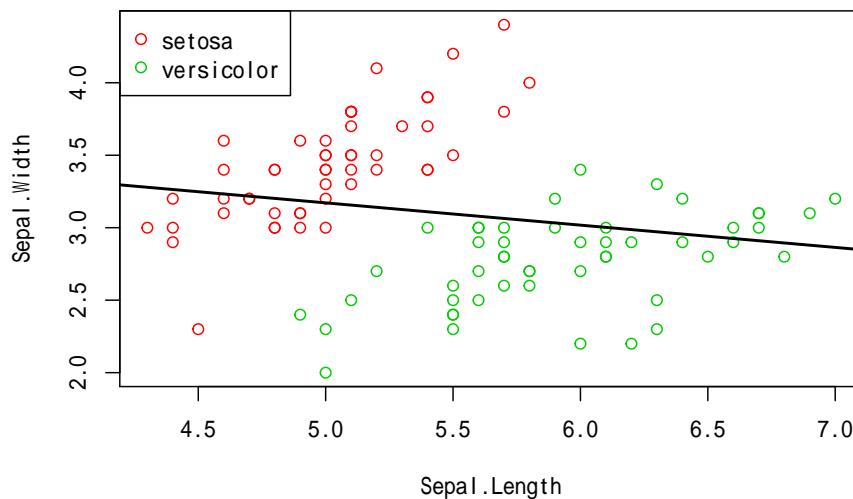
##
Call:
lm(formula = weight ~ height + sex, data = d.class)
##
Residuals:
Min 1Q Median 3Q Max
-19.7627 -5.9583 -0.3682 8.7725 15.7758
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -126.1687 34.6352 -3.643 0.00219 **
height 3.6789 0.5392 6.823 4.09e-06 ***
sexF -6.6208 5.3887 -1.229 0.23697

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 11.06 on 16 degrees of freedom
Multiple R-squared: 0.7903, Adjusted R-squared: 0.7641
F-statistic: 30.15 on 2 and 16 DF, p-value: 3.74e-06
```

结果中的 sexM 项表示以女生为基数，男生体重的平均增加量。这一项不显著。

有些例子中如果忽略了分类变量，结论可能是错误的。例如，考察 iris 数据中花瓣长宽之间的关系。数据中有三个品种的花，仅考虑其中的 setosa 和 versicolor 两个品种。

```
d <- iris[iris[["Species"]] %in% c("setosa", "versicolor"),
 c("Sepal.Length", "Sepal.Width", "Species")]
d$Species <- factor(as.character(d$Species))
lm7 <- lm(Sepal.Width ~ Sepal.Length, data=d)
with(d, plot(Sepal.Length, Sepal.Width, col=(2:3)[Species]))
with(d, legend("topleft", pch=1, col=2:3, legend=levels(Species)))
abline(lm7, lwd=2)
```



```
summary(lm7)
```

```
##
Call:
lm(formula = Sepal.Width ~ Sepal.Length, data = d)
##
Residuals:
Min 1Q Median 3Q Max
```

```
-1.17136 -0.26382 -0.04468 0.29966 1.33618
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 3.93951 0.40621 9.698 5.47e-16 ***
Sepal.Length -0.15363 0.07375 -2.083 0.0398 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 0.4709 on 98 degrees of freedom
Multiple R-squared: 0.04241, Adjusted R-squared: 0.03263
F-statistic: 4.34 on 1 and 98 DF, p-value: 0.03984
```

回归结果花萼长、宽是负相关的，这明显不合理，出现了“悖论”。实际是因为两个品种的样本混杂在一起了。

加入 `Species` 分类变量作为回归自变量：

```
lm8 <- lm(Sepal.Width ~ Species + Sepal.Length, data=d)
summary(lm8)

##
Call:
lm(formula = Sepal.Width ~ Species + Sepal.Length, data = d)
##
Residuals:
Min 1Q Median 3Q Max
-0.88917 -0.14677 -0.02517 0.16643 0.64444
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 1.06519 0.32272 3.301 0.00135 **
Speciesversicolor -1.09696 0.08170 -13.427 < 2e-16 ***
Sepal.Length 0.47200 0.06398 7.377 5.51e-11 ***

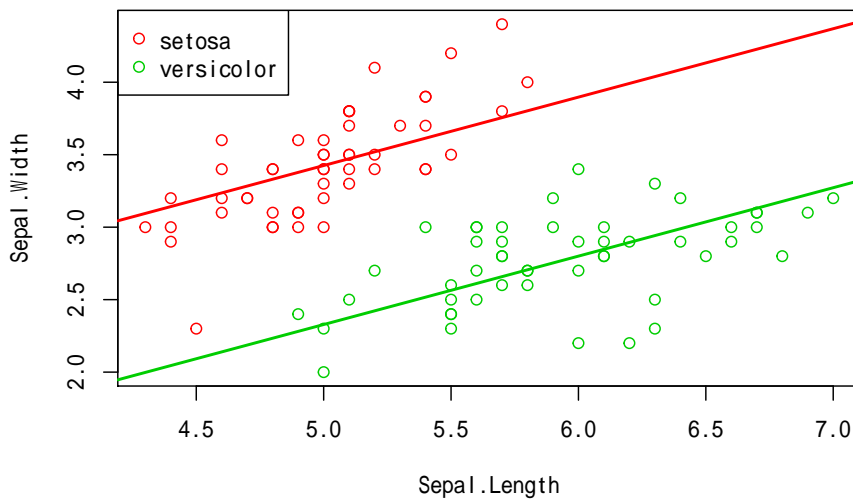
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
Residual standard error: 0.2799 on 97 degrees of freedom
Multiple R-squared: 0.665, Adjusted R-squared: 0.6581
F-statistic: 96.28 on 2 and 97 DF, p-value: < 2.2e-16
```

现在花萼长度变量的系数为正而且高度显著。两个品种区别的变量 `Speciesversicolor` 表示 `versicolor` 品种取 1, `setosa` 取 0 的哑变量。`versicolor` 品种的花萼宽度与 `setosa` 相比在长度相等的情况下平均较小。

作两条回归直线的图形:

```
with(d, plot(Sepal.Length, Sepal.Width, col=(2:3)[Species]))
with(d, legend("topleft", pch=1, col=2:3, legend=levels(Species)[1:2]))
abline(a=coef(lm8)[1], b=coef(lm8)[3], col=2, lwd=2)
abline(a=sum(coef(lm8)[1:2]), b=coef(lm8)[3], col=3, lwd=2)
```



### 33.3.6.1.2 多值哑变量

如果变量  $f$  是一个有  $k$  分类的变量, 可以将  $f$  变成  $k-1$  个取 0、1 值的哑变量  $z_1, z_2, \dots, z_{k-1}$ 。

例如,  $f$  取三种不同的类别, 可以转换为两个哑变量  $z_1, z_2$ , 常用的编码方式为:

$$f = 1 \Leftrightarrow (z_1, z_2) = (0, 0),$$

$$f = 2 \Leftrightarrow (z_1, z_2) = (1, 0),$$

$$f = 3 \Leftrightarrow (z_1, z_2) = (0, 1)$$

这是将  $f = 1$  看成对照组。

在 R 中, 用 `model.matrix(y ~ x1 + f, data=d)` 可以生成将因子  $f$  转换成哑变量后的包含因变量和自变量的矩阵。

### 33.3.7 残差诊断

先复习一些回归理论。将模型写成矩阵形式

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$\mathbf{Y}$  为  $n \times 1$  向量,  $\mathbf{X}$  为  $n \times p$  矩阵, 一般第一列元素全是 1, 代表截距项。  $\boldsymbol{\beta}$  为  $p \times 1$  未知参数向量;  $\boldsymbol{\varepsilon}$  为  $n \times 1$  随机误差向量,  $\boldsymbol{\varepsilon}$  的元素独立且方差为相等的  $\sigma^2$  (未知)。

假设矩阵  $\mathbf{X}$  满秩, 系数的估计为

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

拟合值 (或称预报值) 向量为

$$\hat{\mathbf{Y}} = \mathbf{X} (\mathbf{X} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}$$

其中  $\mathbf{H} = \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$  是  $R^n$  空间的向量向  $\mathbf{X}$  的列张成的线性空间  $\mu(\mathbf{X})$  投影的投影算子矩阵, 叫做帽子矩阵。设  $\mathbf{H} = (h_{ij})_{n \times n}$ 。

拟合残差向量为

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H}) \mathbf{Y}$$

残差平方和为

$$\text{ESS} = \mathbf{e}^T \mathbf{e} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

误差项方差的估计 (要求设计阵  $\mathbf{X}$  满秩) 为均方误差 (MSE)

$$\hat{\sigma}^2 = \text{MSE} = \frac{1}{n-p} \text{ESS}$$

(其中  $p$  在有截距项时是自变量个数加 1)

在线性模型的假设下, 若设计阵  $\mathbf{X}$  满秩,  $\hat{\beta}$  和  $\hat{\sigma}^2$  分别是  $\beta$  和  $\sigma^2$  的无偏估计。

系数估计的方差阵

$$\text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

回归残差及其方差为

$$e_i = y_i - \hat{y}_i, \quad i = 1, 2, \dots, n$$

$$\text{Var}(e_i) = \sigma^2(1 - h_{ii}) \quad (h_{ii} \text{ 是 } H \text{ 的主对角线元素})$$

若 `lmres` 是 R 中 `lm()` 的回归结果, 用 `residuals(lmres)` 可以求残差。

把  $e_i$  除以其标准差估计, 称为**标准化残差**, 或**内部学生化残差**:

$$r_i = \frac{e_i}{s\sqrt{1-h_{ii}}}, \quad i = 1, 2, \dots, n$$

$r_i$  渐近服从正态分布。

若 `lmres` 是 R 中 `lm()` 的回归结果, 用 `rstandard(lmres)` 可以求标准化残差。

如果计算  $y_i$  的预测值时, 删除第  $i$  个观测后建立回归模型得到  $\sigma^2$  的估计  $s_{(i)}^2$ , 则**外部学生化残差**为

$$t_i = \frac{e_i}{s_{(i)}\sqrt{1-h_{ii}}}$$

$t_i$  近似服从  $t(n-p-1)$  分布 (有截距项时  $p$  等于自变量个数加 1)。其中  $s_{(i)}$  有简单公式:

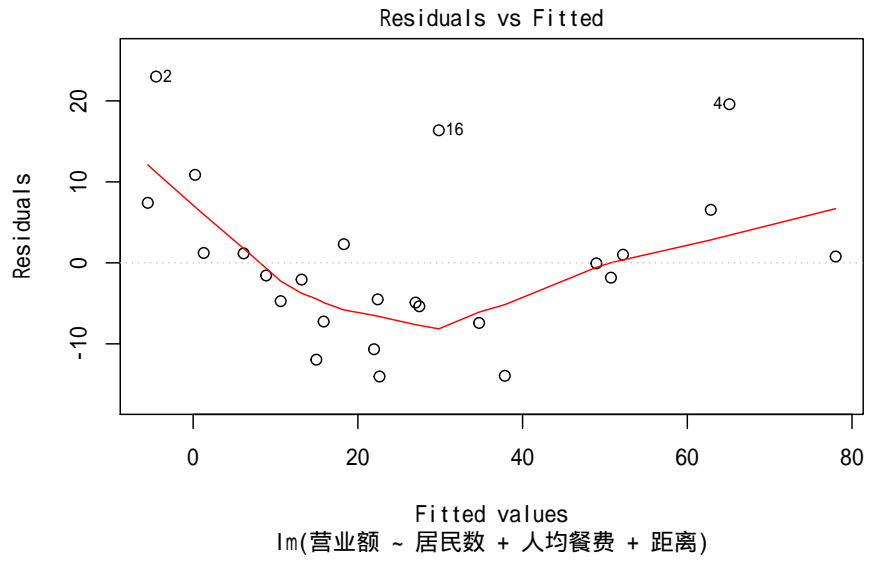
$$s_{(i)}^2 = \frac{n-p-r_i^2}{n-p-1} \hat{\sigma}^2$$

若 `lmres` 是 R 中 `lm()` 的回归结果, 用 `rstudent(lmres)` 可以求外部学生化残差。

在 R 中, 与一元回归的诊断类似。用 `plot()` 作 4 个残差诊断图。可以用 `which=1` 指定仅作第一幅图。

如餐馆营业额例子的残差诊断:

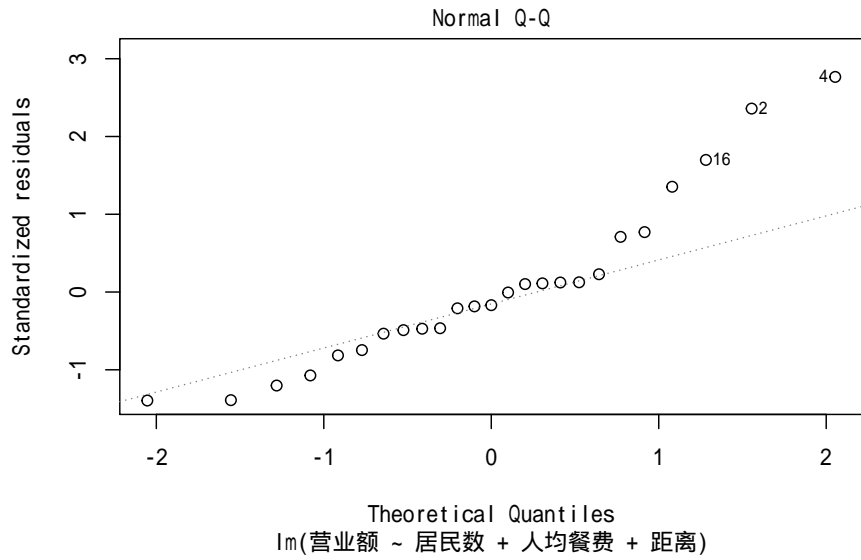
```
plot(lmrst02, which=1)
```



上图是残差对拟合值的散点图，可以查看有无非线性。有轻微的非线性关系。

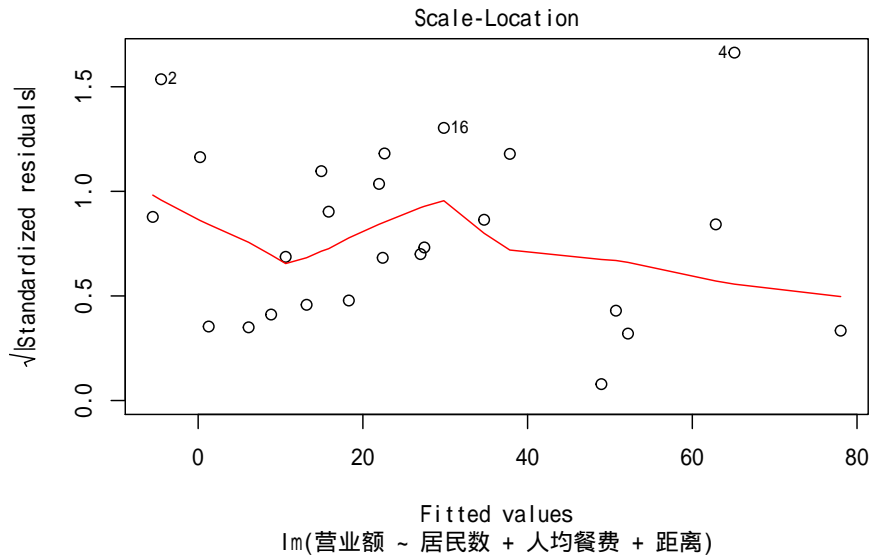
```
plot(lmrst02, which=2)
```





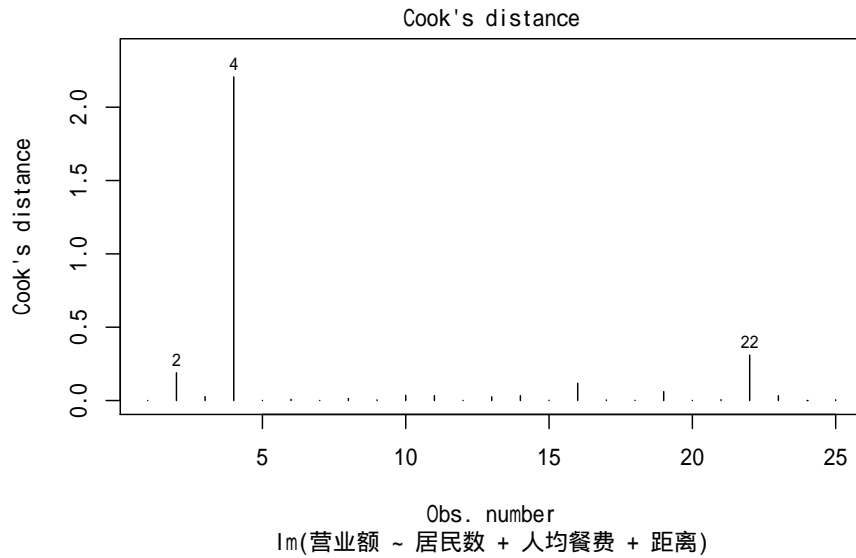
上图是残差的正态 QQ 图，查看残差是否正态分布。残差分布略有右偏，不算太严重。

```
plot(lmrst02, which=3)
```



上图是标准化残差绝对值平方根对拟合值的散点图，可以查看是否有异方差。没有明显的异方差。

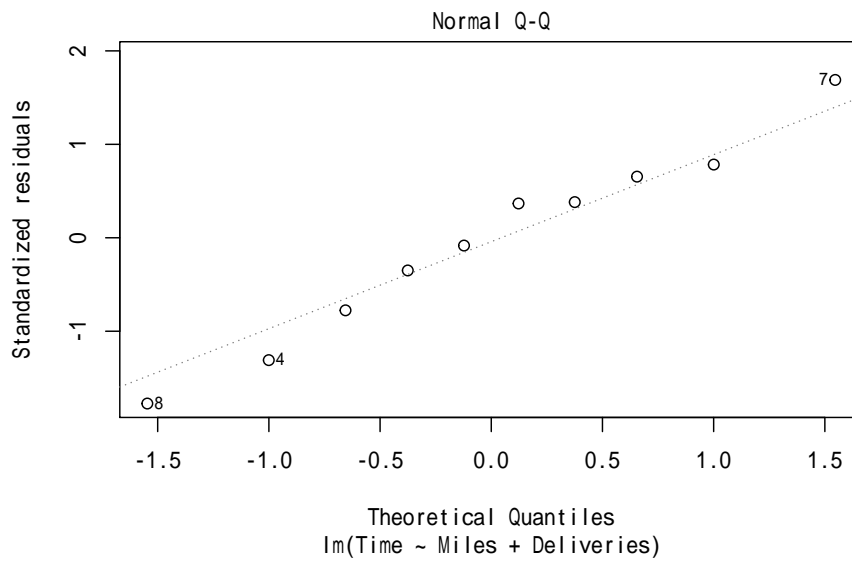
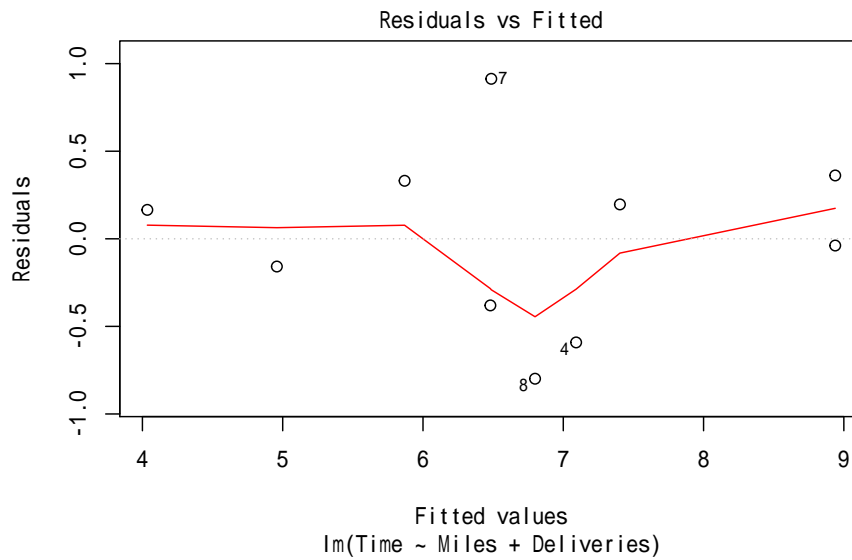
```
plot(lmrst02, which=4)
```

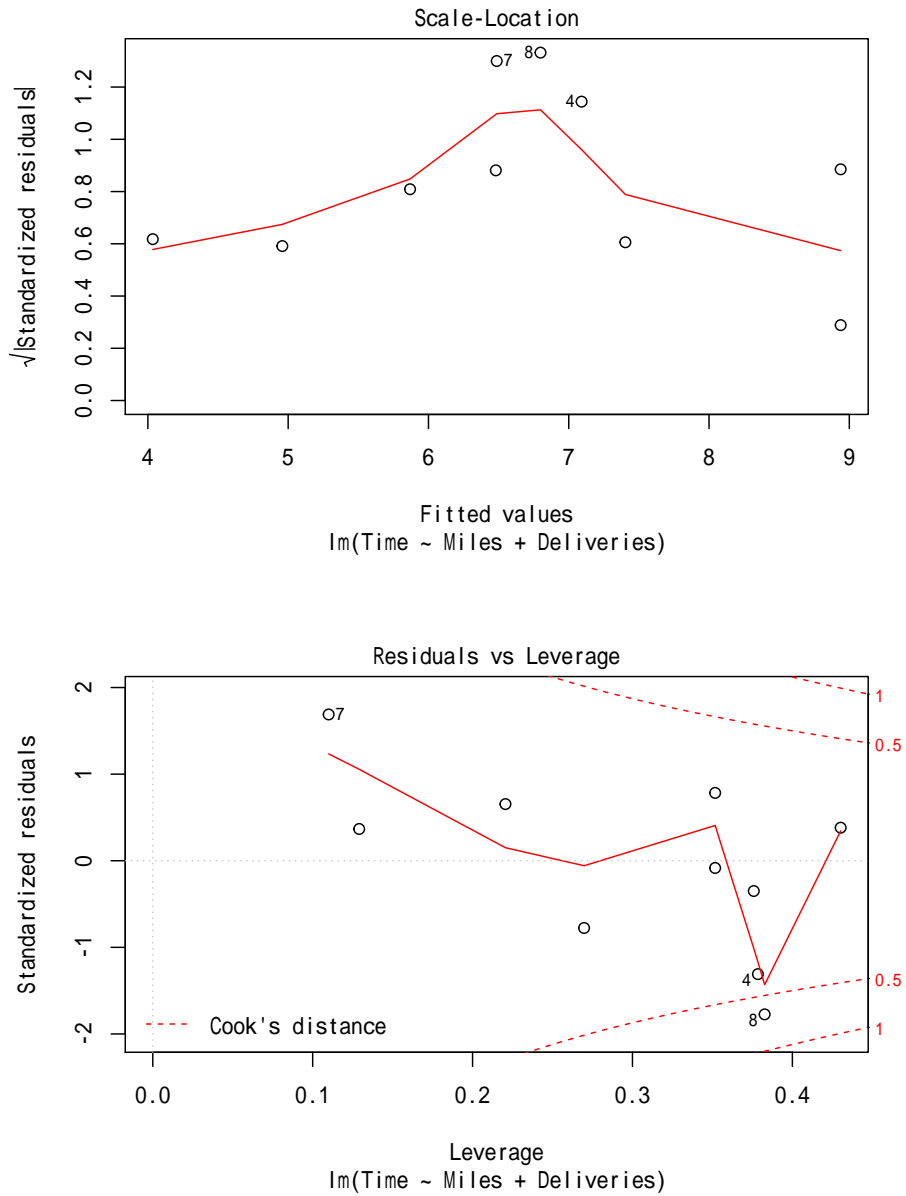


上图用来查看强影响点，4号观测是一个强影响点。

货运公司例子的多元回归的残差诊断：

```
plot(lmbut02)
```





有一定的异方差倾向，但是数据量不大就不做处理。

### 33.3.8 多重共线性

狭义的多重共线性 (multicollinearity): 自变量的数据存在线性组合近似地等于零, 使得矩阵求解回归系数时结果不稳定, 回归结果很差。

广义的多重共线性: 自变量之间存在较强的相关性, 这样自变量是联动的, 互相之间有替代作用。甚至于斜率项的正负号都因为这种替代作用而可能是错误的方向。

例如, 餐馆营业额例子中, F 检验显著, 5 个自变量如果用在一元回归中斜率项都显著, 但是在多元回归中, 在 0.15 水平下仅有人均餐费和到市中心的距离的系数是显著的, 月收入、餐馆数、居民数的系数不显著。

但是实际上, 单独使用这三个自变量作一元线性回归, 结果都是显著的:

```
summary(lm(`营业额` ~ `月收入`, data=Resturant))

##
Call:
lm(formula = 营业额 ~ 月收入, data = Resturant)
##
Residuals:
Min 1Q Median 3Q Max
-32.151 -10.725 -0.696 6.033 47.819
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.484580 5.955478 0.081 0.936
月收入 0.004995 0.000944 5.291 2.27e-05 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 16.88 on 23 degrees of freedom
Multiple R-squared: 0.549, Adjusted R-squared: 0.5294
F-statistic: 27.99 on 1 and 23 DF, p-value: 2.273e-05
```

如何识别多重共线性?

如果两个自变量之间的相关系数显著地不等于零，这两个自变量就有广义的共线性。

如果线性关系的 F 检验显著但是单个回归系数都不显著，可能是由于多重共线性。

如果有单个回归系数显著但是 F 检验不显著，可能是由于多重共线性。

如果某些回归系数的正负号与通常的认识相反，可能是由于多重共线性。

将第  $i$  个自变量  $x_i$  作为因变量，用其它的  $p-1$  个自变量作为自变量作多元线性回归，得到一个复相关系数平方  $R_i^2$ ，这个  $R_i^2$  接近于 1 时  $x_i$  与其他自变量之间存在多重共线性。令  $x_i$  的容忍度 (tolerance) 等于  $1 - R_i^2$ ，容忍度接近于 0 时存在多重共线性。

容忍度小于 0.1 时多重共线性为严重程度。

称容忍度的倒数为方差膨胀因子 (VIF)，VIF 大于 10 或者大于 5 作为严重的多重共线性。

为了计算 VIF，首先把矩阵  $X^T X$  看成一个协方差阵，把它转换为相关系数阵设为  $M$ ，则  $M^{-1}$  的各主对角线元素就是各个 VIF。

car 包的 `vif()` 函数计算方差膨胀因子。

如：

```
car::vif(lmrst01)

居民数 人均餐费 月收入 餐馆数 距离
8.233159 2.629940 5.184365 1.702361 1.174053
```

可以认为变量“居民数”和“月收入”有共线问题。

做共线诊断还可以用条件数 (Conditional Index): 这是一个正数，用来衡量  $(X^T X)^{-1}$  的稳定性，定义为  $X^T X$  的最大特征值与最小特征值之比。条件数在 0—100 之间时认为无共线性，在 100—1000 之间时认为自变量之间有中等或较强共线性，在 1000 以上认为自变量之间有强共线性。

解决多重共线性问题，最简单的方法是回归自变量选择，剔除掉有严重共线性的自变量，这些自变量的信息可以由其他变量代替。还可以对自变量作变换，如用主成分分析降维。可以用收缩方法如岭回归、lasso 回归等。

### 33.3.9 强影响点分析

强影响点是删去以后严重改变参数估计值的观测。包括自变量取值离群和因变量拟合离群的点。

杠杆 (leverage) 指帽子矩阵的对角线元素  $h_{ii}$ ,

$$\frac{1}{n} \leq h_{ii} \leq \frac{1}{d_i}$$

其中  $d_i$  是第  $i$  个观测的重复观测次数。某观测杠杆值高说明该观测自变量有异常值。杠杆值大于  $2p/n$  的观测需要仔细考察 (有截距项时  $p$  等于自变量个数加 1)。

若 `lmres` 是 R 中 `lm()` 的回归结果, 用 `hatvalues(lmres)` 可以求杠杆值。考察外学生化残差  $t_i$ , 绝对值超过 2 的观测拟合误差大, 在  $y$  方向离群, 需要关注。

若 `lmres` 是 R 中 `lm()` 的回归结果, 用 `rstudent(lmres)` 可以求外学生化。

Cook 距离统计量:

$$D_i = \frac{1}{p} r_i^2 \frac{h_{ii}}{1 - h_{ii}}$$

包含了  $y$  方向的离群  $r_i$  和  $x$  方向的离群  $h_{ii}$  的信息。超过  $\frac{4}{n}$  的值需要注意。

若 `lmres` 是 R 中 `lm()` 的回归结果, 用 `cooks.distance(lmres)` 可以求 Cook 距离。

R 中的强影响点诊断函数还有 `dfbetas()`, `dffits()`, `covratio()`。

**偏杠杆值** 衡量每个自变量 (包括截距项) 对杠杆的贡献。把第  $j$  个自变量关于其它自变量回归得到残差, 第  $i$  个残差的平方占总残差平方和的比例为第  $j$  自变量在第  $i$  观测处的偏杠杆值。偏杠杆值影响自变量选择时对该变量的选择。

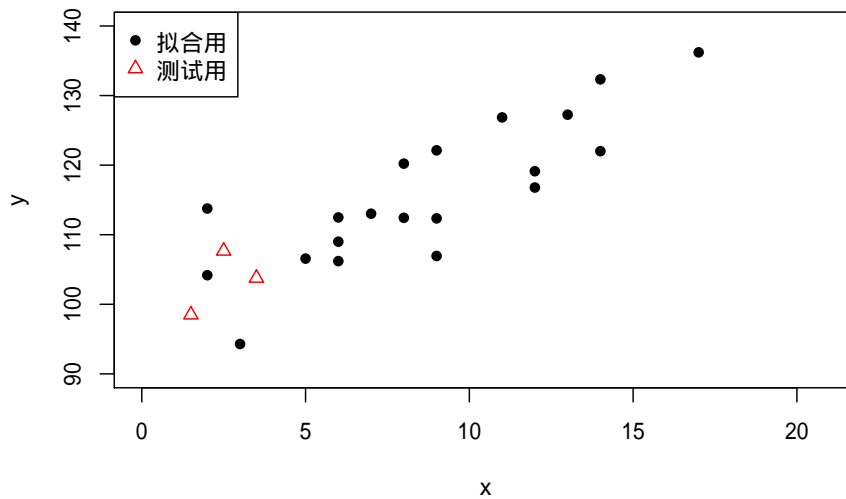
### 33.3.10 过度拟合示例

$R^2$  代表了模型对数据的拟合程度, 模型中加入的自变量越多,  $R^2$  越大。是不是模型中的自变量越多越好? 可能会发生“过度拟合”。用来建模的数据都拟合误差很小, 但是模型很难有合理解释, 对新的数据的预测效果很差甚至于完全错误。

```

set.seed(10)
n <- 20
x <- sample(1:n, size=n, replace=TRUE)
a <- 100
b <- 2
sigma <- 5
y <- a + b*x + rt(n, 4)*sigma
xnew <- c(1.5, 2.5, 3.5)
ynew <- a + b*xnew + rnorm(length(xnew), 0, sigma)
plot(x, y, pch=16, xlim=c(0, n+1), ylim=c(90,140))
points(xnew, ynew, pch=2, col="red")
legend("topleft", pch=c(16,2), col=c("black", "red"),
 legend=c("拟合用", "测试用"))

```



作线性回归:

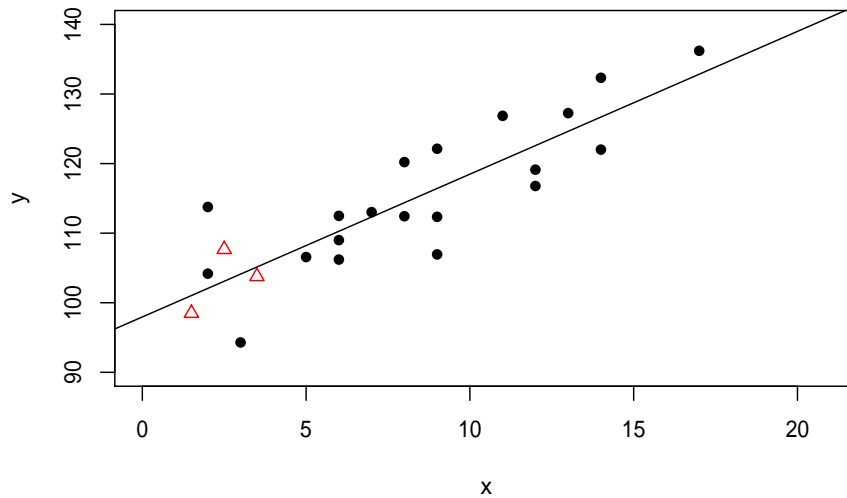
```

plot(x, y, pch=16, xlim=c(0, n+1), ylim=c(90,140))
points(xnew, ynew, pch=2, col="red")
lmof1 <- lm(y ~ x)

```



```
abline(lmof1)
```



回归系数:

```
summary(lmof1)
```

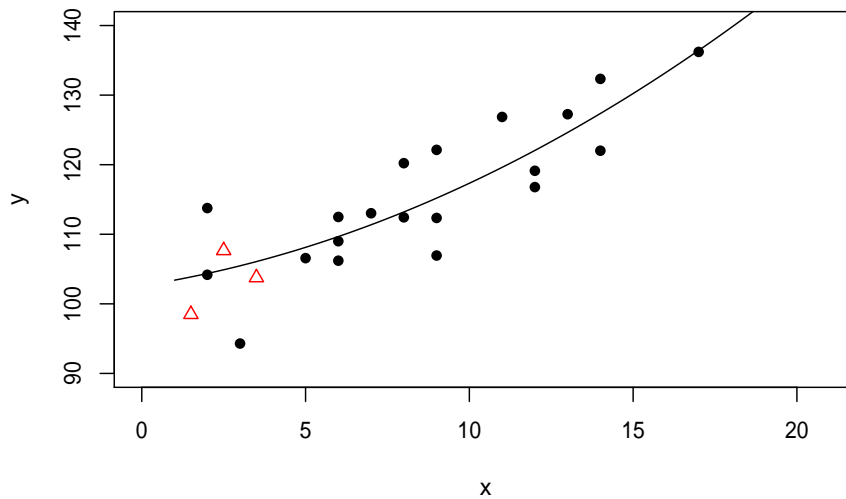
```
##
Call:
lm(formula = y ~ x)
##
Residuals:
Min 1Q Median 3Q Max
-9.8206 -4.0691 -0.2855 3.9371 11.7011
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 97.9570 3.0244 32.388 < 2e-16 ***
x 2.0521 0.3163 6.489 4.21e-06 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
Residual standard error: 5.767 on 18 degrees of freedom
Multiple R-squared: 0.7005, Adjusted R-squared: 0.6839
F-statistic: 42.1 on 1 and 18 DF, p-value: 4.209e-06
```

二次多项式回归:

```
plot(x, y, pch=16, xlim=c(0, n+1), ylim=c(90,140))
points(xnew, ynew, pch=2, col="red")
lmof2 <- lm(y ~ x + I(x^2))
xx <- seq(1, n, length=100)
yy <- predict(lmof2, newdata=data.frame(x=xx))
lines(xx, yy)
```



回归系数:

```
summary(lmof2)
```

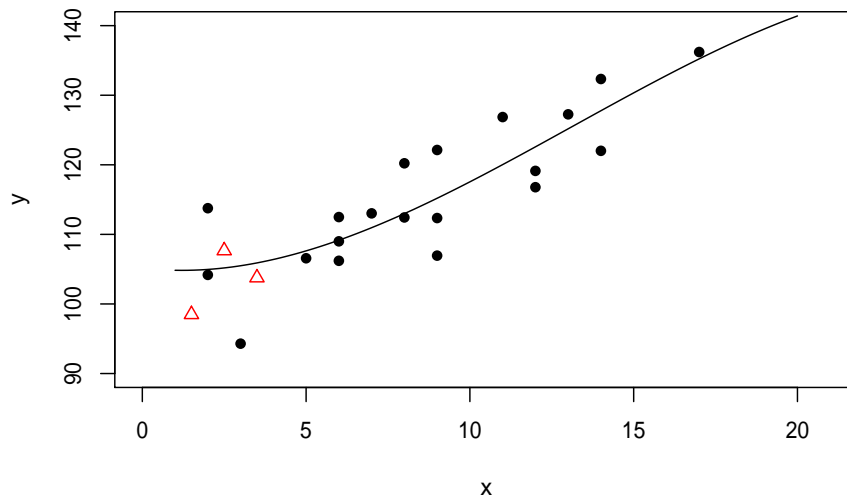
```
##
Call:
lm(formula = y ~ x + I(x^2))
```

```
##
Residuals:
Min 1Q Median 3Q Max
-11.175 -3.059 -0.439 3.358 9.400
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 102.59339 5.32854 19.254 5.57e-13 ***
x 0.73688 1.28561 0.573 0.574
I(x^2) 0.07371 0.06985 1.055 0.306

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 5.749 on 17 degrees of freedom
Multiple R-squared: 0.7189, Adjusted R-squared: 0.6859
F-statistic: 21.74 on 2 and 17 DF, p-value: 2.066e-05
```

这个回归结果出现了多重共线性问题。也已经过度拟合。

三次多项式回归:



回归系数:

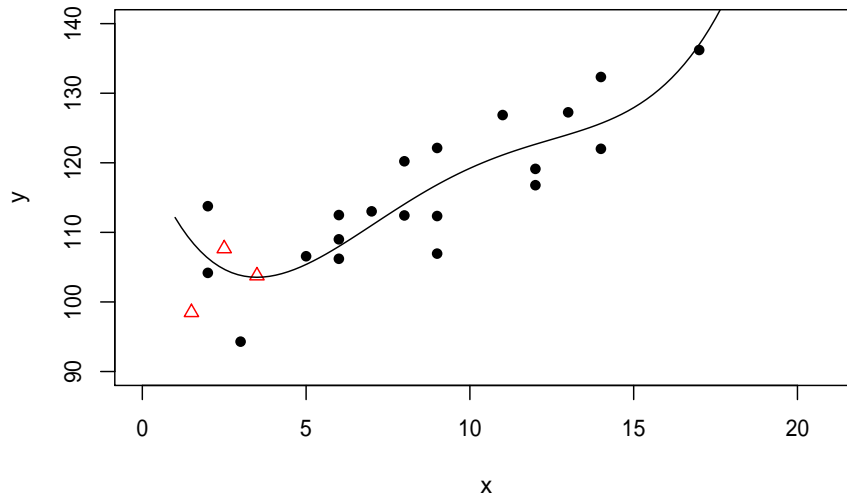
```
summary(lmof3)
```

```
##
Call:
lm(formula = y ~ x + I(x^2) + I(x^3))
##
Residuals:
Min 1Q Median 3Q Max
-11.1976 -3.0835 -0.3567 3.6324 8.8020
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 105.162401 8.827478 11.913 2.29e-09 ***
x -0.558534 3.734856 -0.150 0.883
I(x^2) 0.240987 0.456853 0.527 0.605
I(x^3) -0.006124 0.016518 -0.371 0.716

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 5.901 on 16 degrees of freedom
Multiple R-squared: 0.7213, Adjusted R-squared: 0.6691
F-statistic: 13.8 on 3 and 16 DF, p-value: 0.0001053
```

这个回归结果出现了多重共线性问题。也已经过度拟合。

四次多项式回归:



回归系数:

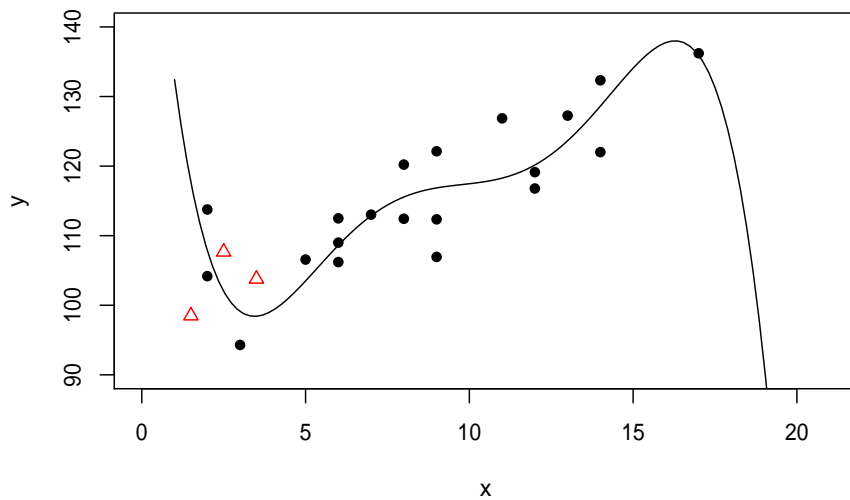
```
summary(lmof4)
```

```
##
Call:
lm(formula = y ~ x + I(x^2) + I(x^3) + I(x^4))
##
Residuals:
Min 1Q Median 3Q Max
-9.8975 -3.5647 0.1241 4.7063 7.4675
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 122.380891 18.117315 6.755 6.47e-06 ***
x -12.830942 11.890980 -1.079 0.298
I(x^2) 2.785077 2.385359 1.168 0.261
I(x^3) -0.206102 0.184800 -1.115 0.282
I(x^4) 0.005273 0.004854 1.086 0.294
```

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 5.868 on 15 degrees of freedom
Multiple R-squared: 0.7416, Adjusted R-squared: 0.6727
F-statistic: 10.76 on 4 and 15 DF, p-value: 0.0002563
```

五次多项式回归:



已经明显过度拟合。

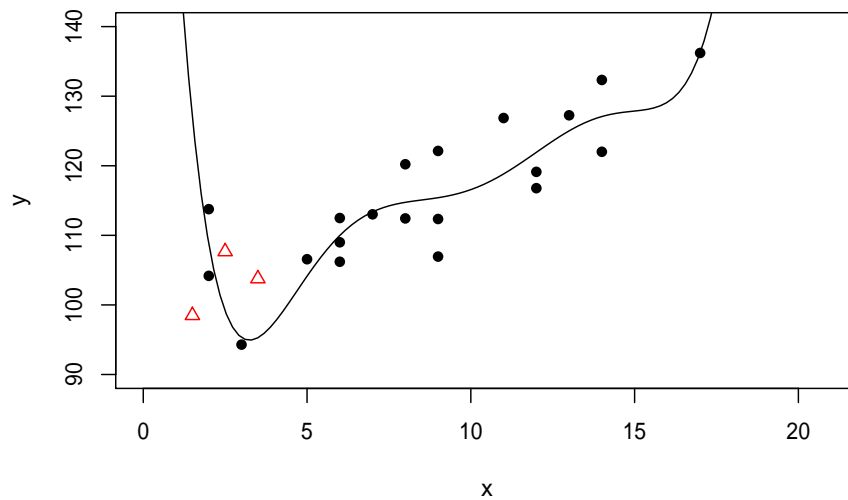
回归系数:

```
summary(lmof5)
```

```
##
Call:
lm(formula = y ~ x + I(x^2) + I(x^3) + I(x^4) + I(x^5))
##
Residuals:
Min 1Q Median 3Q Max
```

```
-9.9204 -3.4541 0.2201 3.7495 8.5922
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 181.846244 41.023546 4.433 0.000568 ***
x -65.530939 34.875410 -1.879 0.081224 .
I(x^2) 18.157919 9.886841 1.837 0.087594 .
I(x^3) -2.170388 1.242055 -1.747 0.102453
I(x^4) 0.118729 0.071167 1.668 0.117456
I(x^5) -0.002418 0.001513 -1.598 0.132453

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 5.586 on 14 degrees of freedom
Multiple R-squared: 0.7815, Adjusted R-squared: 0.7034
F-statistic: 10.01 on 5 and 14 DF, p-value: 0.0003082
```



六次多项式回归:

回归系数:

```
summary(lmof6)

##
Call:
lm(formula = y ~ x + I(x^2) + I(x^3) + I(x^4) + I(x^5) + I(x^6))
##
Residuals:
Min 1Q Median 3Q Max
-8.4627 -3.2409 -0.6654 3.1208 8.0410
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 2.699e+02 8.387e+01 3.218 0.00673 **
x -1.585e+02 8.483e+01 -1.868 0.08447 .
I(x^2) 5.346e+01 3.103e+01 1.723 0.10864
I(x^3) -8.572e+00 5.481e+00 -1.564 0.14188
I(x^4) 7.158e-01 5.033e-01 1.422 0.17849
I(x^5) -2.999e-02 2.307e-02 -1.300 0.21606
I(x^6) 4.982e-04 4.159e-04 1.198 0.25229

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 5.501 on 13 degrees of freedom
Multiple R-squared: 0.8032, Adjusted R-squared: 0.7124
F-statistic: 8.843 on 6 and 13 DF, p-value: 0.0005655
```

### 33.3.11 嵌套模型比较

如果两个多元线性回归模型，第一个模型中的自变量都在第二个模型中，且第二个模型具有更多的自变量，则称第一个模型嵌套在第二个模型中。

例如：第一个模型中自变量为  $x_1, x_2, x_5$ ，第二个模型自变量为  $x_1, x_2, x_3, x_4, x_5$ ，则第一个模型嵌套在第二个模型中。

又如：第一个模型自变量为  $x_1, x_2$ ，第二个模型自变量为  $x_1, x_2, x_1^2, x_2^2, x_1x_2$ ，则



第一个模型嵌套在第二个模型中。这时令  $x_3 = x_1^2$ ,  $x_4 = x_2^2$ ,  $x_5 = x_1x_2$ , 第二个模型变成有 5 个自变量的多元线性回归模型。

在嵌套的模型中, 相对而言自变量多的模型叫做完全模型 (full model), 自变量少的模型叫做简化模型 (reduced model)。

完全模型是否比简化模型更好? 在回归模型选择中贯彻一个“精简性”原则 (Ocam's razor): 在对建模数据拟合效果相近的情况下, 越简单的模型越好。

例如: 全模型是

$$Ey = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \beta_4x_4$$

精简模型是

$$Ey = \beta_0 + \beta_1x_1 + \beta_2x_2$$

判断两个模型是否没有显著差异, 只要检验:

$$H_0 : \beta_3 = \beta_4 = 0$$

这可以构造一个方差分析 F 检验,

设全模型有  $t$  个自变量, 模型得到的回归残差平方和为  $SSE_f$ ; 设精简模型有  $s$  个自变量 ( $s < t$ ), 模型得到的回归残差平方和为  $SSE_r$ ; 检验零假设为多出的自变量对应的系数都等于零。检验统计量为

$$F = \frac{(SSE_r - SSE_f)/(t - s)}{SSE_f/(n - t - 1)}$$

在全模型成立且  $H_0$  成立时,  $F$  服从  $F(t - s, n - t - 1)$  分布。计算右侧 p 值。

在 R 中用 `anova()` 函数比较两个嵌套的线性回归结果可以进行这样的方差分析 F 检验。

例如, 餐馆营业额回归模型的比较。`lmrst01` 是完全模型, 包含 5 个自变量; `lmrst02` 是嵌套的精简模型, 包含居民数、人均餐费、到市中心距离共 3 个自变量。用 `anova()` 函数可以检验多出的变量是否有意义:

```
anova(lmrst01, lmrst02)
```

```
Analysis of Variance Table
```

```
##
```

```
Model 1: 营业额 ~ 居民数 + 人均餐费 + 月收入 + 餐馆数 + 距离
```

```
Model 2: 营业额 ~ 居民数 + 人均餐费 + 距离
Res.Df RSS Df Sum of Sq F Pr(>F)
1 19 2153.0
2 21 2267.2 -2 -114.17 0.5038 0.6121
```

模型 p 值为 0.61, 在 0.05 水平下不拒绝多出的月收入 and 餐馆数的系数全为零的零假设, 两个模型的效果没有显著差异, 应选择更精简的模型。

方差分析检验仅能比较嵌套模型。对不同模型计算 AIC, 取 AIC 较小的模型, 这可以对非嵌套的模型进行比较选择。R 中用 `AIC()` 函数比较两个回归结果的 AIC 值。

如:

```
AIC(lmrst01, lmrst02)
```

```
df AIC
lmrst01 7 196.3408
lmrst02 5 193.6325
```

精简模型 `lmrst02` 的 AIC 更小, 是更好的模型。

### 33.3.12 拟合与预测

#### 33.3.12.1 拟合

有了参数最小二乘估计后, 对建模用的每个数据点计算

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_p x_{ip}$$

称为拟合值 (fitted value)。

得到回归模型结果 `lmres` 后, 要对原数据框中的观测值做预测, 只要使用 `predict(lmres)`。

#### 33.3.12.2 点预测

为了使用得到的模型结果 `lmres` 对新数据做预测, 建立包含了自变量的一组新的观测值的数据框 `dp`, 用 `predict(lmres, newdata=dp)` 做预测。

如餐馆营业额的选择自变量的回归模型 `lmrst02` 的拟合值:

```
predict(lmrst02)

1 2 3 4 5 6
52.1892541 -4.5010247 21.9626575 65.1136964 6.1284803 22.4083426
7 8 9 10 11 12
1.2783244 34.7189934 10.6275869 37.8433996 62.8524888 18.2959990
13 14 15 16 17 18
-5.5101343 14.9558131 8.8598588 29.8309866 78.0159456 13.1673581
19 20 21 22 23 24
15.8469287 50.7274217 27.4608977 0.2331759 22.6274030 48.9610796
25
27.0050675
```

如果是一元回归，一般还画数据的散点图并画回归直线。多元回归的图形无法在二维表现出来。

有了估计的回归方程后，对一组新的自变量值  $(x_{01}, \dots, x_{0p})$ ，可以计算对应的因变量的预测值：

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \dots + \hat{\beta}_p x_{0p}$$

在 R 中，设 `lmres` 保存了回归结果，`newd` 是一个保存了新的自变量值的数据框，此数据框结构与原建模用数据框类似但是自变量与原来不同，且不需要有因变量。这时用 `predict(lm1, data=newd)` 预测。

例如，利用包含居民数、人均餐费、到市中心距离的模型 `lmrst02`，求居民数 = 50(万居民)，人均餐费 = 100(元)，距市中心 10 千米的餐馆的日均营业额：

```
predict(
 lmrst02,
 newdata=data.frame(
 `居民数`=50, `人均餐费`=100, `距离`=10
))
```

```
1
17.88685
```

预测的日均营业额为 17.9 千元。

函数 `expand.grid()` 可以对若干个变量的指定值，生成包含所有组合的数据框，如：

```
newd <- expand.grid(
 `居民数`=c(60, 140),
 `人均餐费`=c(50, 130),
 `距离`=c(6, 16))
newd
```

```
居民数 人均餐费 距离
1 60 50 6
2 140 50 6
3 60 130 6
4 140 130 6
5 60 50 16
6 140 50 16
7 60 130 16
8 140 130 16
```

```
predict(lmrst02, newdata=newd)
```

```
1 2 3 4 5 6 7
14.186636 29.404103 26.797108 42.014574 8.488759 23.706225 21.099230
8
36.316696
```

### 33.3.12.3 均值的置信区间

对  $Ey = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$ ，可以计算置信度为  $1 - \alpha$  的置信区间，称为均值的置信区间。

在 `predict()` 中加选项 `interval="confidence"`，用 `level=` 指定置信度，可以计算均值的置信区间。

如

```
predict(
 lmrst02, interval="confidence", level=0.95,
 newdata=data.frame(
 `居民数`=50, `人均餐费`=100, `距离`=10
))
```

```
fit lwr upr
1 17.88685 10.98784 24.78585
```

其中 `fit` 是预测值, `lwr` 和 `upr` 分别是置信下限和置信上限。

#### 33.3.12.4 个别值的预测区间

对  $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$ , 可以计算置信度为  $1 - \alpha$  的预测区间, 称为预测区间, 预测区间比均值的置信区间要宽。

在 `predict()` 中加选项 `interval="prediction"`, 用 `level=` 指定置信度, 可以计算预测区间。

如

```
predict(
 lmrst02, interval="prediction", level=0.95,
 newdata=data.frame(
 `居民数`=50, `人均餐费`=100, `距离`=10
))
```

```
fit lwr upr
1 17.88685 -4.795935 40.56963
```

其中 `fit` 是预测值, `lwr` 和 `upr` 分别是预测下限和预测上限。

#### 33.3.13 利用线性回归模型做曲线拟合

某些非线性关系可以通过对因变量和自变量的简单变换变成线性回归模型。

例如，彩色显影中，染料光学密度  $Y$  与析出银的光学密度  $x$  有如下类型的关系

$$Y \approx Ae^{-B/x}, \quad B > 0$$

这不是线性关系。两边取对数得

$$\ln Y \approx \ln A - B \frac{1}{x}$$

令

$$Y^* = \ln Y, \quad x^* = \frac{1}{x}$$

则

$$Y^* \approx \ln A - Bx^*$$

为线性关系。

从  $n$  组数据  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  得到变换的数据  $(x_1^*, y_1^*), (x_2^*, y_2^*), \dots, (x_n^*, y_n^*)$ 。对变换后的数据建立线性回归方程

$$\hat{y}^* = \hat{a} + \hat{b}x^*$$

反变换得

$$\hat{A} = e^{\hat{a}}, \quad \hat{B} = -\hat{b}$$

则有

$$\hat{Y} = \hat{A}e^{-\hat{B}/x}$$

再考虑一个钢包容积的例子。炼钢钢包随使用次数增加而容积增大。测量了 13 组这样的数据：

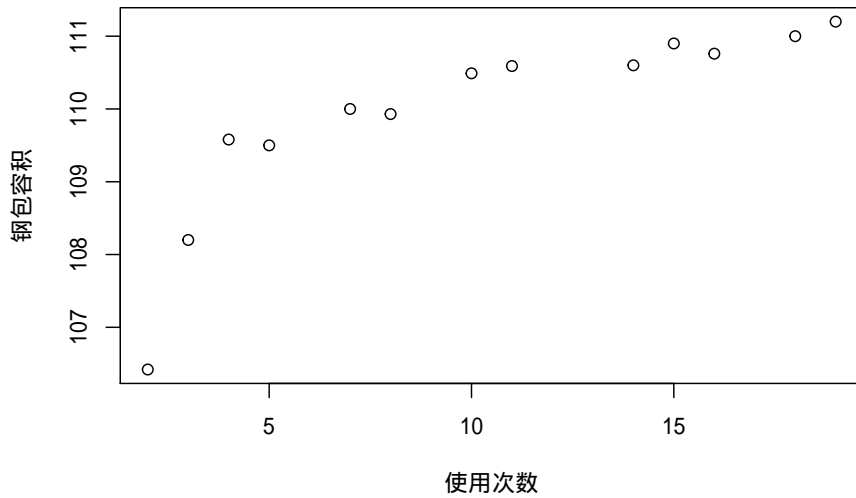
```
SteelBag <- data.frame(
 x = c(2, 3, 4, 5, 7, 8, 10,
 11, 14, 15, 16, 18, 19),
 y = c(106.42, 108.20, 109.58, 109.50, 110.0,
 109.93, 110.49, 110.59, 110.60, 110.90,
 110.76, 111.00, 111.20)
)
```

```
knitr::kable(SteelBag)
```

x	y
2	106.42
3	108.20
4	109.58
5	109.50
7	110.00
8	109.93
10	110.49
11	110.59
14	110.60
15	110.90
16	110.76
18	111.00
19	111.20

散点图:

```
with(SteelBag, plot(
 x, y, xlab=" 使用次数", ylab=" 钢包容积"
))
```



散点图呈现非线性。用线性回归近似：

```
lmsb1 <- lm(y ~ x, data=SteelBag)
summary(lmsb1)
```

```
##
Call:
lm(formula = y ~ x, data = SteelBag)
##
Residuals:
Min 1Q Median 3Q Max
-1.97615 -0.38502 0.04856 0.53724 0.80611
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 108.01842 0.44389 243.346 < 2e-16 ***
x 0.18887 0.03826 4.937 0.000445 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

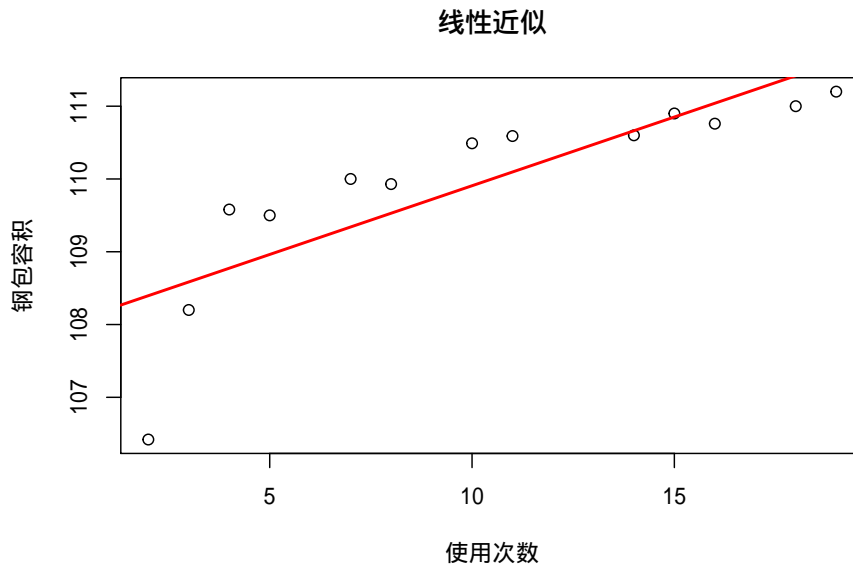


```

Residual standard error: 0.7744 on 11 degrees of freedom
Multiple R-squared: 0.689, Adjusted R-squared: 0.6607
F-statistic: 24.37 on 1 and 11 DF, p-value: 0.000445
```

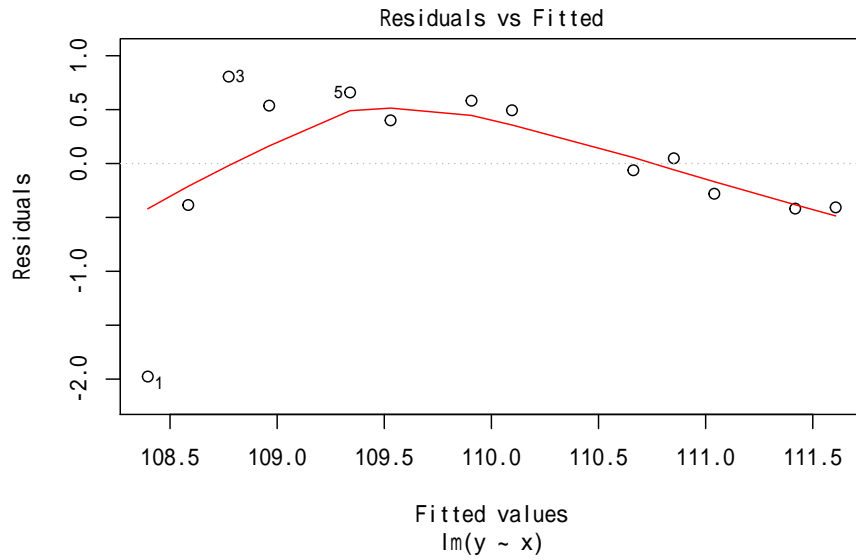
结果显著。 $R^2 = 0.69$ 。拟合图：

```
with(SteelBag, plot(
 x, y, xlab=" 使用次数", ylab=" 钢包容积",
 main=" 线性近似"
))
abline(lmsb1, col="red", lwd=2)
```



残差诊断：

```
plot(lmsb1, which=1)
```



残差图呈现非线性。

用双曲线模型：

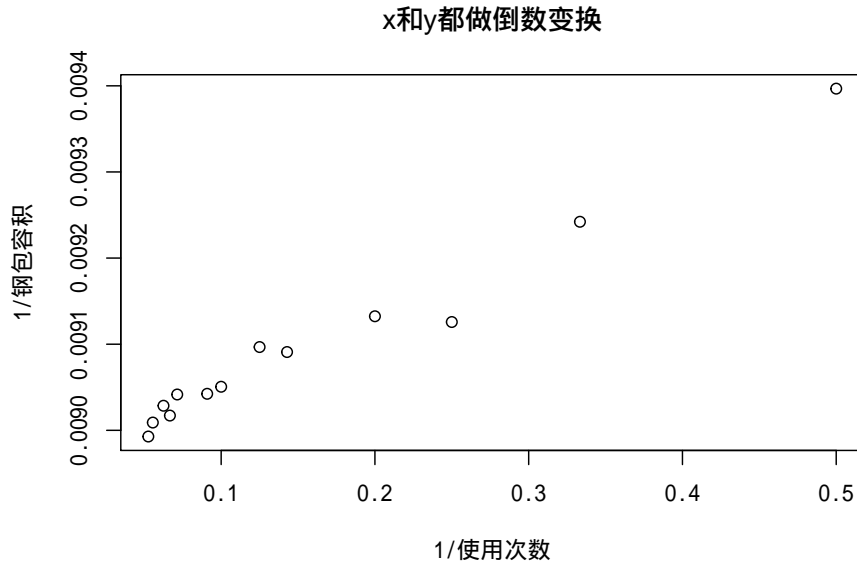
$$\frac{1}{y} \approx a + b\frac{1}{x}$$

令  $x^* = 1/x, y^* = 1/y$ ，化为线性模型

$$y^* \approx a + bx^*$$

$(x^*, y^*)$  的散点图：

```
with(SteelBag, plot(
 1/x, 1/y, xlab="1/使用次数", ylab="1/钢包容积",
 main="x 和 y 都做倒数变换"
))
```



$y^*$  对  $x^*$  的回归:

```
lmsb2 <- lm(I(1/y) ~ I(1/x), data=SteelBag)
summary(lmsb2)
```

```
##
Call:
lm(formula = I(1/y) ~ I(1/x), data = SteelBag)
##
Residuals:
Min 1Q Median 3Q Max
-4.817e-05 -3.686e-06 4.000e-07 1.008e-05 2.642e-05
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 8.967e-03 8.371e-06 1071.14 < 2e-16 ***
I(1/x) 8.292e-04 4.118e-05 20.14 4.97e-10 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

Residual standard error: 1.903e-05 on 11 degrees of freedom
Multiple R-squared: 0.9736, Adjusted R-squared: 0.9712
F-statistic: 405.4 on 1 and 11 DF, p-value: 4.97e-10
```

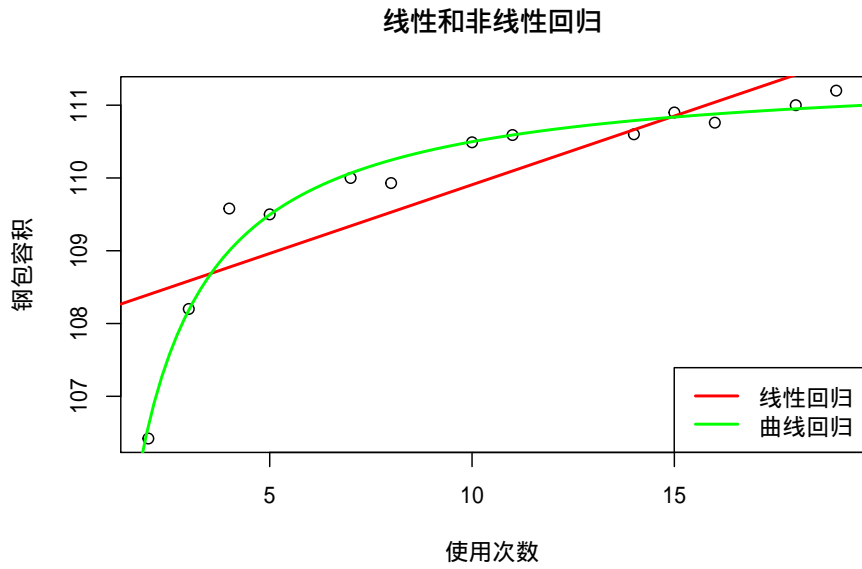
结果显著。解得  $\hat{a} = 0.008967$ ,  $\hat{b} = 0.0008292$ , 经验公式为

$$\frac{1}{\hat{y}} = 0.008967 + 0.0008292 \frac{1}{x}$$

$R^2$  从线性近似的 0.69 提高到了 0.97。

拟合图:

```
with(SteelBag, plot(
 x, y, xlab=" 使用次数", ylab=" 钢包容积",
 main=" 线性和非线性回归"
))
abline(lmsb1, col="red", lwd=2)
curve(1/(0.008967 + 0.0008292/x), 1, 20,
 col="green", lwd=2, add=TRUE)
legend("bottomright", lty=1, lwd=2,
 col=c("red", "green"),
 legend=c(" 线性回归", " 曲线回归"))
```



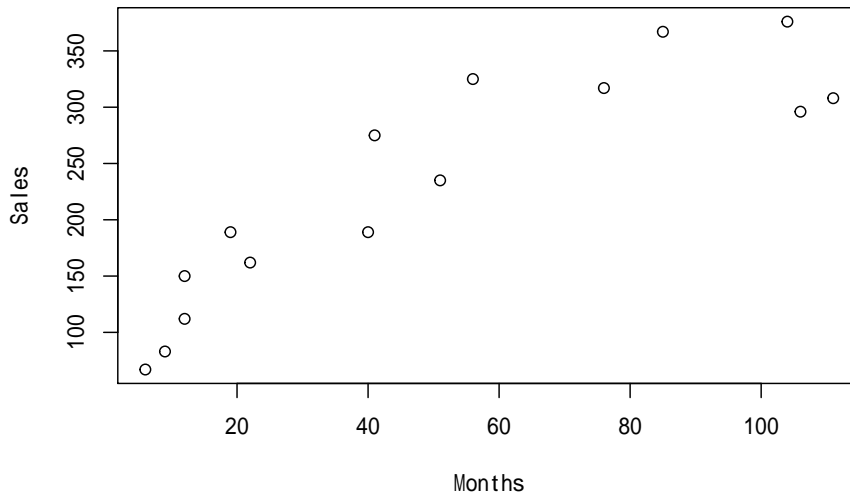
考虑 Reynolds, Inc. 销售业绩数据分析问题。Reynolds, Inc. 是一个工业和实验室量具厂商。为研究销售业务员的业绩，考察业务员从业年限 (Months, 单位: 月) 与其销售的电子量具数量 (Sales) 的关系。随机抽查了 15 名业务员。

```
knitr::kable(Reynolds)
```

Months	Sales
41	275
106	296
76	317
104	376
22	162
12	150
85	367
111	308
40	189
51	235
9	83
12	112
6	67
56	325
19	189

散点图:

```
with(Reynolds, plot(Months, Sales))
```



散点图呈现非线性。

用线性近似：

```
lmre1 <- lm(Sales ~ Months, data=Reynolds)
summary(lmre1)
```

```
##
Call:
lm(formula = Sales ~ Months, data = Reynolds)
##
Residuals:
Min 1Q Median 3Q Max
-67.166 -38.684 2.557 28.875 80.673
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 111.2279 21.6280 5.143 0.000189 ***
Months 2.3768 0.3489 6.812 1.24e-05 ***

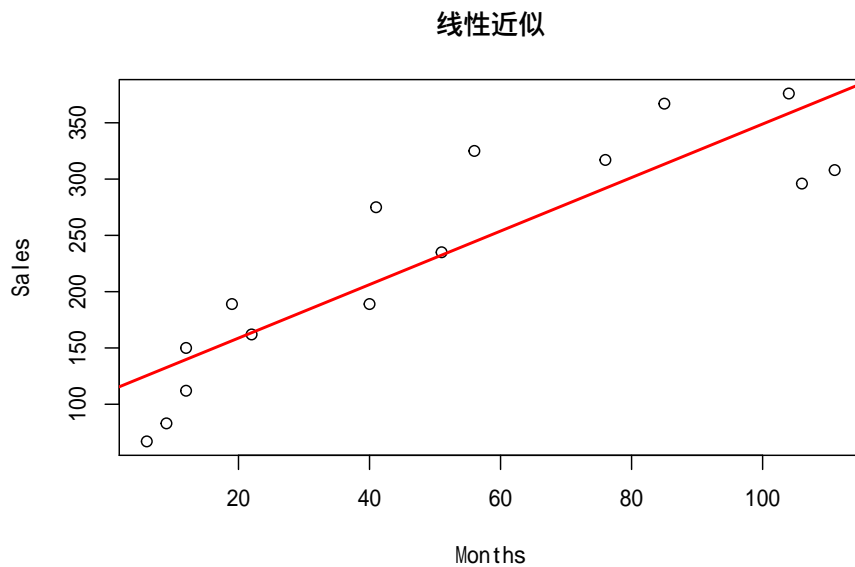
```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 49.52 on 13 degrees of freedom
Multiple R-squared: 0.7812, Adjusted R-squared: 0.7643
F-statistic: 46.41 on 1 and 13 DF, p-value: 1.239e-05
```

结果显著。  $R^2 = 0.78$ 。拟合图：

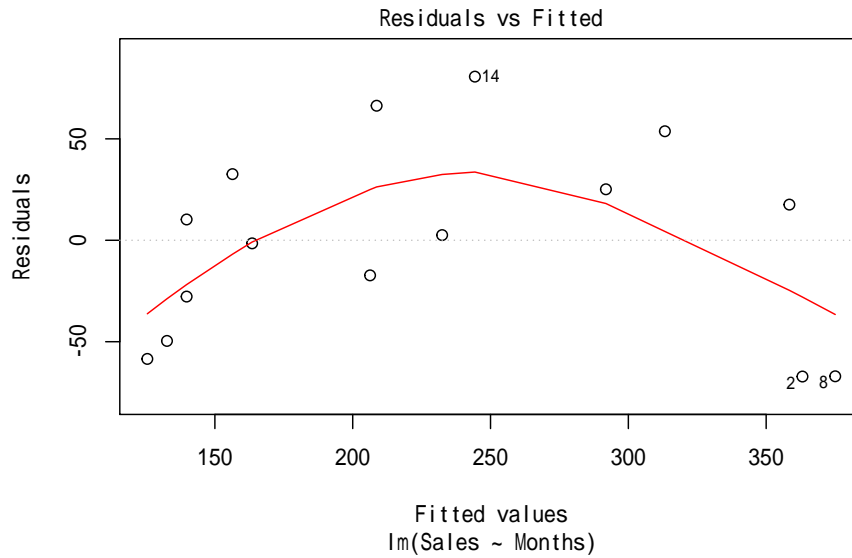
```
with(Reynolds, plot(Months, Sales, main=" 线性近似"))
abline(lmre1, col="red", lwd=2)
```



残差诊断：

```
plot(lmre1, which=1)
```





残差图有明显的非线性。

考虑最简单的非线性模型：

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$$

令  $x_1 = x$ ,  $x_2 = x^2$ , 有

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

是二元线性回归模型。作二次多项式回归：

```
lmre2 <- lm(Sales ~ Months + I(Months^2), data=Reynolds)
summary(lmre2)
```

```
##
Call:
lm(formula = Sales ~ Months + I(Months^2), data = Reynolds)
##
Residuals:
Min 1Q Median 3Q Max
-54.963 -16.691 -6.242 31.996 43.789
```

```
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 45.347579 22.774654 1.991 0.06973 .
Months 6.344807 1.057851 5.998 6.24e-05 ***
I(Months^2) -0.034486 0.008948 -3.854 0.00229 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 34.45 on 12 degrees of freedom
Multiple R-squared: 0.9022, Adjusted R-squared: 0.8859
F-statistic: 55.36 on 2 and 12 DF, p-value: 8.746e-07
```

模型显著。 $R^2$  从线性近似的 0.78 提高到 0.90。 $x^2$  项的系数的显著性检验  $p$  值为 0.002, 显著不等于零, 说明二次项是必要的。

这样添加二次项容易造成  $x$  与  $x^2$  之间的共线性, 所以添加中心化的二次项:

$$x_2 = (x - 60)^2$$

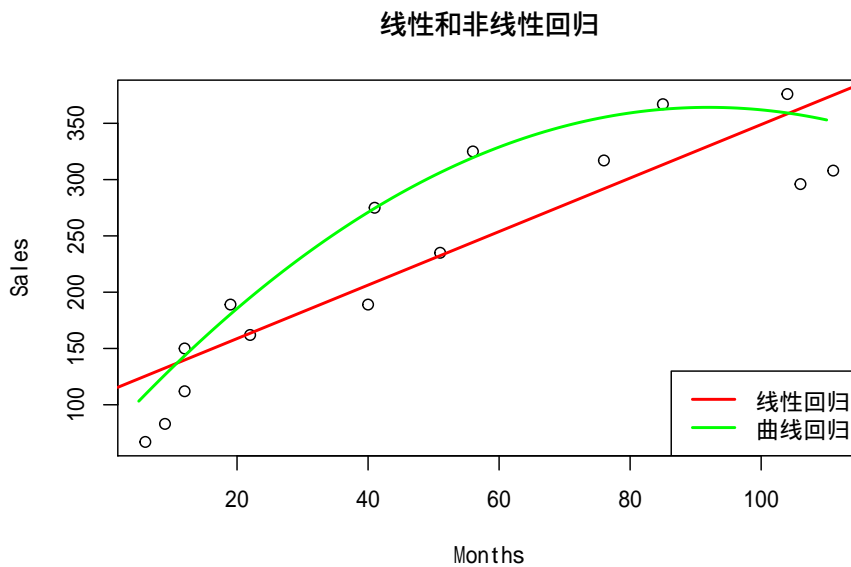
```
lmre3 <- lm(Sales ~ Months + I((Months-60)^2), data=Reynolds)
summary(lmre3)
```

```
##
Call:
lm(formula = Sales ~ Months + I((Months - 60)^2), data = Reynolds)
##
Residuals:
Min 1Q Median 3Q Max
-54.963 -16.691 -6.242 31.996 43.789
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 169.495742 21.331978 7.946 4.03e-06 ***
Months 2.206535 0.246744 8.943 1.18e-06 ***
I((Months - 60)^2) -0.034486 0.008948 -3.854 0.00229 **
```

```

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 34.45 on 12 degrees of freedom
Multiple R-squared: 0.9022, Adjusted R-squared: 0.8859
F-statistic: 55.36 on 2 and 12 DF, p-value: 8.746e-07

with(Reynolds, plot(Months, Sales, main=" 线性和非线性回归"))
abline(lmre1, col="red", lwd=2)
curve(196.50 + 2.2065*x - 0.03449*(x-60)^2, 5, 110,
 col="green", lwd=2, add=TRUE)
legend("bottomright", lty=1, lwd=2,
 col=c("red", "green"),
 legend=c(" 线性回归", " 曲线回归"))
```



### 33.3.14 分组建立多个模型

有时希望将数据按照一个或者几个分类变量分组，每一组分别建立模型，并将模型结果统一列表比较。`broom`包可以用来将模型结果转换成规范的数据框格

式。

以肺癌病人数据为例，建立  $v_1$  对  $v_0$  和  $age$  的多元线性回归模型：

```
print(d.cancer)
```

```
A tibble: 34 x 6
id age sex type v0 v1
<dbl> <dbl> <chr> <chr> <dbl> <dbl>
1 1 70 F 腺癌 26.5 2.91
2 2 70 F 腺癌 135. 35.1
3 3 69 F 腺癌 210. 74.4
4 4 68 M 腺癌 61 35.0
5 5 67 M 鳞癌 238. 128.
6 6 75 F 腺癌 330. 112.
7 7 52 M 鳞癌 105. 32.1
8 8 71 M 鳞癌 85.2 29.2
9 9 68 M 鳞癌 102. 22.2
10 10 79 M 鳞癌 65.5 21.9
... with 24 more rows
```

```
lmgr01 <- lm(v1 ~ v0 + age, data = d.cancer)
```

```
summary(lmgr01)
```

```
##
Call:
lm(formula = v1 ~ v0 + age, data = d.cancer)
##
Residuals:
Min 1Q Median 3Q Max
-42.757 -11.010 -2.475 11.907 52.941
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 7.90818 33.14895 0.239 0.814
v0 0.43288 0.05564 7.780 1.79e-07 ***
```

```
age -0.12846 0.53511 -0.240 0.813

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 21.76 on 20 degrees of freedom
(11 observations deleted due to missingness)
Multiple R-squared: 0.7683, Adjusted R-squared: 0.7451
F-statistic: 33.16 on 2 and 20 DF, p-value: 4.463e-07
```

用 broom 包的 tidy() 函数可以将系数估计结果转换成合适的 tibble 数据框格式:

```
library(broom)
tidy(lmgr01)
```

```
A tibble: 3 x 5
term estimate std.error statistic p.value
<chr> <dbl> <dbl> <dbl> <dbl>
1 (Intercept) 7.91 33.1 0.239 0.814
2 v0 0.433 0.0556 7.78 0.000000179
3 age -0.128 0.535 -0.240 0.813
```

用 broom 包的 augment() 函数可以获得拟合值、残差等每个观测的回归结果:

```
knitr::kable(augment(lmgr01), digits=2)
```

.rownames	v1	v0	age	.fitted	.se.fit	.resid	.hat	.sigma	.cooksd	.std.resid
1	2.91	26.51	70	10.39	7.92	-7.48	0.13	22.25	0.01	-0.37
2	35.08	135.48	70	57.56	5.43	-22.48	0.06	21.68	0.03	-1.07
3	74.44	209.74	69	89.84	6.92	-15.40	0.10	22.01	0.02	-0.75
4	34.97	61.00	68	25.58	6.06	9.39	0.08	22.21	0.01	0.45
5	128.34	237.75	67	102.22	8.06	26.12	0.14	21.37	0.09	1.29
6	112.34	330.24	75	141.23	12.50	-28.89	0.33	20.81	0.43	-1.62
7	32.10	104.90	52	46.64	7.83	-14.54	0.13	22.04	0.03	-0.72
8	29.15	85.15	71	35.65	6.31	-6.50	0.08	22.27	0.00	-0.31
9	22.15	101.65	68	43.18	5.10	-21.03	0.05	21.77	0.02	-0.99
10	21.94	65.54	79	26.13	10.19	-4.19	0.22	22.30	0.00	-0.22
11	12.33	125.31	55	55.09	6.88	-42.76	0.10	19.79	0.16	-2.07
12	99.44	224.36	54	98.09	10.54	1.35	0.23	22.32	0.00	0.07
13	2.30	12.93	55	6.44	7.58	-4.14	0.12	22.30	0.00	-0.20
14	23.96	40.21	75	15.68	9.24	8.28	0.18	22.23	0.01	0.42
15	7.39	12.58	61	5.52	6.93	1.87	0.10	22.32	0.00	0.09
16	112.58	231.04	76	98.16	8.81	14.42	0.16	22.03	0.03	0.72
17	91.62	172.13	65	74.07	5.57	17.55	0.07	21.93	0.02	0.83
18	13.95	39.26	66	16.42	6.37	-2.47	0.09	22.32	0.00	-0.12
20	122.45	161.00	63	69.51	5.43	52.94	0.06	18.47	0.14	2.51
21	68.35	105.26	67	44.87	4.84	23.48	0.05	21.63	0.02	1.11
22	5.25	13.25	51	7.09	8.67	-1.84	0.16	22.32	0.00	-0.09
23	3.34	18.70	49	9.71	9.27	-6.37	0.18	22.27	0.01	-0.32
24	50.36	60.23	49	27.69	8.91	22.67	0.17	21.59	0.09	1.14

用 broom 包的 `glance()` 函数可以将回归的复相关系数平方、F 检验 p 值等整体结果做成一行的数据框：

```
knitr::kable(glance(lmgr01), digits=2)
```

r.squared	adj.r.squared	sigma	statistic	p.value	df	logLik	AIC	BIC	deviance
0.77	0.75	21.76	33.16	0	3	-101.87	211.74	216.28	9470.34

下面将男病人与女病人分别建模，并以表格形式合并分组的建模结果。用 `dplyr` 包的 `group_by()` 函数分组，用 `tidyr` 包的 `nest()` 函数可以将分组后的数据

框转换成一个新的数据框，新数据框中有一列是列表类型，每个元素是一个子数据框：

```
d.cancer %>%
 group_by(sex) %>%
 nest()

A tibble: 2 x 2
sex data
<chr> <list>
1 F <tibble [13 x 5]>
2 M <tibble [21 x 5]>
```

这样，可以用 `purrr::map()` 函数对每一组分别建模，建模结果可以借助 `broom` 包制作成合适的数据框格式，然后用 `unnest()` 函数将不同组的结果合并成一个大数据框，如：

```
d.cancer %>%
 group_by(sex) %>%
 nest() %>%
 mutate(model = map(data, function(df) summary(lm(v1 ~ v0 + age, data=df))),
 tidied = map(model, tidy)) %>%
 unnest(tidied, .drop = TRUE)

A tibble: 6 x 6
sex term estimate std.error statistic p.value
<chr> <chr> <dbl> <dbl> <dbl> <dbl>
1 F (Intercept) 171. 147. 1.16 0.310
2 F v0 0.485 0.136 3.56 0.0235
3 F age -2.74 2.39 -1.15 0.316
4 M (Intercept) -17.2 30.5 -0.563 0.583
5 M v0 0.486 0.0657 7.39 0.00000528
6 M age 0.204 0.484 0.421 0.681
```

当需要分组拟合许多个模型时，这种方法比较方便。

## 33.4 非参数回归

### 33.4.1 模型

线性回归模型可以看成非线性回归模型的特例:

$$Y = f(X) + \varepsilon$$

其中  $f(x)$  为未知的回归函数。

参数方法: 假定  $f(x)$  具有某种形式, 如

- $f(x) = a + bx$ : 线性回归;
- $f(x) = a + bx + cx^2$ : 二次多项式回归;
- $f(x) = Ae^{bx}$ : 指数模型, 等等。

二次多项式回归可以令  $X_1 = x, X_2 = x^2$ , 变成二元回归模型来解决。指数模型可以令  $z = \ln Y$ , 模型化为  $z = a + bx$ 。有一些曲线模型可以通过变换化为线性回归。

在多元情形, 一般的非线性回归模型为

$$Y = f(x_1, x_2, \dots, x_p) + \varepsilon$$

但是这样可选的模型就过于复杂, 难以把握。比较简单的是不考虑变量之间交互作用的可加模型:

$$Y = \sum_{j=1}^p f_j(x_j) + \varepsilon$$

其中  $f_j(\cdot)$  是未知的回归函数, 需要满足一定的光滑性条件。

$f_j(\cdot)$  可以是参数形式的, 比如二次多项式、三次多项式、阶梯函数等。较好的一种选择是使用三次样条函数。

### 33.4.2 样条平滑

为了得到一般性的  $Y$  与  $X$  的曲线关系  $f(x)$  的估计, 可以使用样条函数。三次样条函数将实数轴用若干个节点 (knots)  $\{z_k\}$  分成几段, 每一段上  $\hat{f}(x)$  为



三次多项式，函数在节点处有连续的二阶导数。样条函数是光滑的分段三次多项式。

用样条函数估计  $f(x)$  的准则是曲线接近各观测值点  $(x_i, y_i)$ ，同时曲线足够光滑。

在 R 中用 `smooth.spline` 函数进行样条曲线拟合。取每个自变量  $x_i$  处为一个节点，对于给定的某个光滑度/模型复杂度系数值  $\lambda$ ，求函数  $f(x)$  使得

$$\min_{f(\cdot)} \left\{ \sum_i [y_i - f(x_i)]^2 + \lambda \int [f''(x)]^2 dx \right\}$$

$\lambda$  越大，所得的曲线越光滑。`smooth.spline()` 函数可以通过交叉验证方法自动取得一个对于预测最优的光滑参数  $\lambda$  值，也可以通过 `df=` 选项指定一个等效自由度，等效自由度越大，模型越复杂，曲线光滑程度越低。`df` 值相当于多元线性回归中的自变量个数。

如

```
set.seed(1)
nsamp <- 30
x <- runif(nsamp, -10, 10)
xx <- seq(-10, 10, length.out=100)
x <- sort(x)
y <- 10*sin(x/10*pi)^2 + rnorm(nsamp,0,0.3)
plot(x, y)
curve(10*sin(x/10*pi)^2, -10, 10, add=TRUE, lwd=2)

library(splines)
res <- smooth.spline(x, y)
lines(spline(x, res$y), col="red")
res2 <- loess(y ~ x, degree=2, span=0.3)
lines(xx, predict(res2, newdata=data.frame(x=xx)),
 col="blue")
legend("top", lwd=c(2,1,1),
 col=c("black", "red", "blue"),
 legend=c(" 真实函数关系", " 样条平滑结果", " 局部线性拟合"))
```

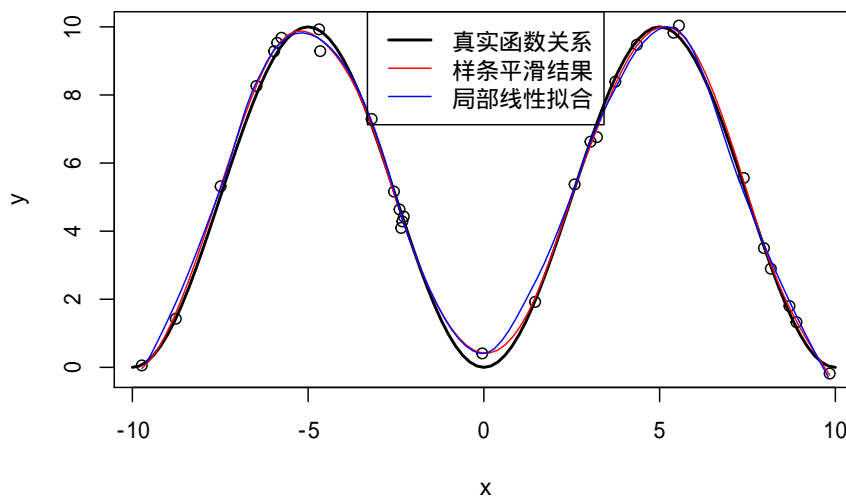


图 33.11: 曲线平滑示例

其中 `res` 的元素 `y` 为拟合值，用 `spline(x,y)` 从一组散点输出光滑曲线以便用 `lines()` 函数绘图。

R 函数 `splines(x,y)` 不是做样条平滑，而是做样条插值，其结果是在原始的自变量 `x` 范围内产生等间隔距离的格子点值，输出包含格子点上的样条插值 `x` 和 `y` 坐标的列表。样条平滑曲线不需要穿过输入的各个散点，但是插值则需要穿过输入的各个散点。

R 函数 `approx(x,y)` 用线性插值方法产生线性插值后的连续函数在等间隔的横坐标上的坐标值。

### 33.4.3 局部多项式曲线平滑

另一种非参数的拟合非线性回归曲线  $f(x)$  的方法是对每个自变量  $x$  处选一个局部的小邻域，用这个小邻域附近的  $(x_i, y_i)$  点拟合一个局部的零次、一次或二次多项式，用拟合的局部多项式在  $x$  处的值作为回归函数在  $x$  处的值的估计。

局部零次多项式的方法叫做核回归，公式为

$$\hat{f}(x) = \frac{\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) Y_i}{\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right)}$$

其中  $K(\cdot)$  称为核函数，是类似标准正态密度这样的中间高两边低的可以用来加权的函数，比如，双三次核函数：

$$K(x) = \begin{cases} (1 - |x|^3)^3 & |x| \leq 1 \\ 0 & \text{其它} \end{cases}$$

```
kernel.dcube <- function(u){
 y <- numeric(length(u))
 sele <- (abs(u) < 1)
 y[sele] <- (1 - abs(u[sele])^3)^3
 y
}
curve(kernel.dcube, -1.5, 1.5, main=" 双三次核函数")
```

参数  $h$  称为窗宽， $h$  越大，参与平均的点越多，曲线越光滑，回归复杂度越低。 $h$  选取可以用交叉验证方法。

R 扩展包 KernSmooth 的函数 `locpoly(x, y, degree=1, bandwidth=0.25)` 可以计算核回归，其中 `bandwidth` 是输入的窗宽，函数 `dpill(x,y)` 可以帮助确定窗宽。如 `locpoly(x, y, degree=1, bandwidth=dpill(x,y))`。`stats` 包的 `ksmooth()` 函数也可以计算核回归。

当局部拟合的是一次或者二次多项式时，这种曲线拟合方法叫做局部多项式回归。R 函数 `loess(y ~ x, degree=1, span=0.75)` 拟合局部线性函数，用 `span=` 控制结果的光滑度，`span` 是局部拟合所用的点的比例，越接近于 1 结果越光滑。取 `degree=2` 拟合局部二次多项式函数。见图33.11。

#### 33.4.4 样条函数变换

$m$  个节点的三次样条函数需要  $n + 4$  个参数，因为每段需要 4 个参数， $m + 1$  段需要  $4m + 4$  个参数，而在  $m$  个节点上连续、一阶导数连续、二阶导数连续构成三个约束条件，所以参数个数为  $m + 4$  个。

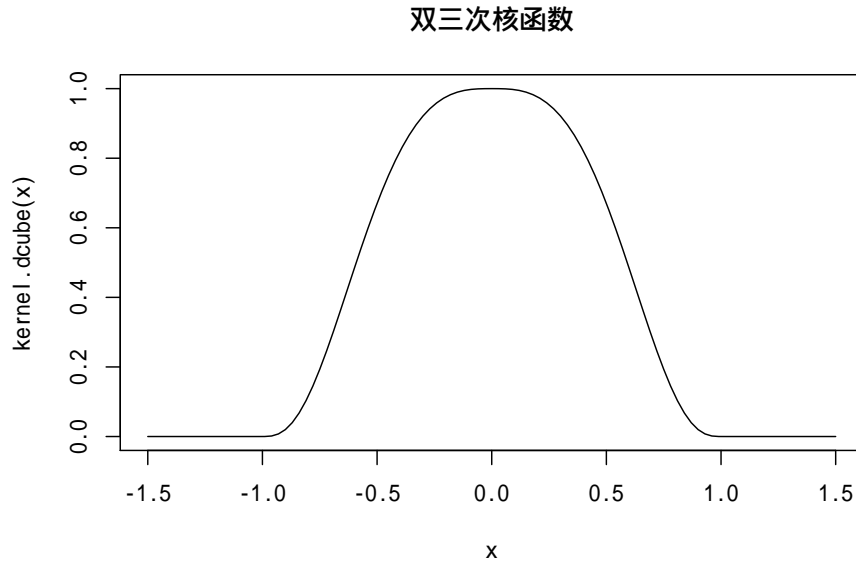


图 33.12: 双三次核函数

自然样条函数假定函数在最左边一段和最右边一段为线性函数，这样  $m$  个节点需要  $m+2$  个参数。R 的 `lm()` 函数中对自变量  $x$  指定 `ns(x)` 可以对输入的  $x$  指定作自然样条变换，`ns()` 可以用 `df=` 指定自由度作为曲线复杂度的度量。如

```
set.seed(1)
nsamp <- 30
x <- runif(nsamp, -10, 10)
xx <- seq(-10, 10, length.out=100)
x <- sort(x)
y <- 10*sin(x/10*pi)^2 + rnorm(nsamp,0,0.3)
plot(x, y)
curve(10*sin(x/10*pi)^2, -10, 10, add=TRUE, lwd=2)

res <- lm(y ~ ns(x, df=4))
lines(xx, predict(res, newdata=data.frame(x=xx)),
```

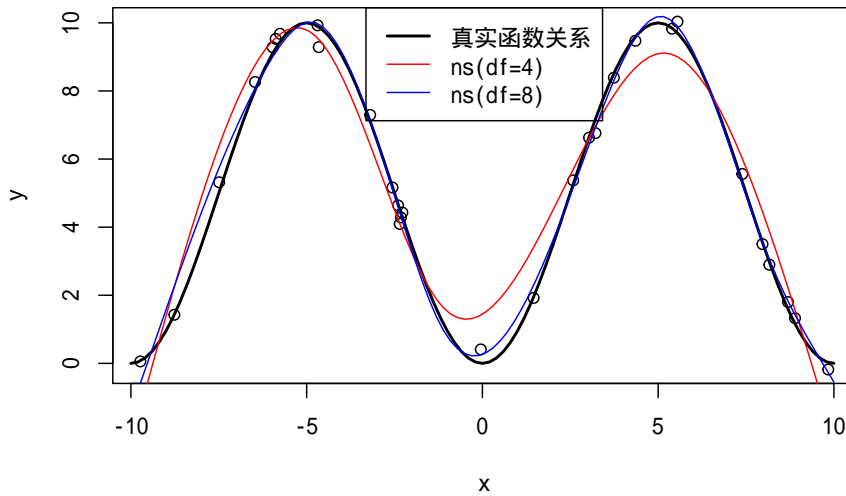


图 33.13: 自然样条回归示例

```

col="red")
res2 <- lm(y ~ ns(x, df=8))
lines(xx, predict(res2, newdata=data.frame(x=xx)),
 col="blue")
legend("top", lwd=c(2,1,1),
 col=c("black", "red", "blue"),
 legend=c(" 真实函数关系", "ns(df=4)", "ns(df=8)"))

```

在多元回归中也可以用 `ns(x)` 对单个自变量引入非线性。

### 33.4.5 线性可加模型

虽然在用 `lm()` 作多元回归时可以用 `ns()`、`poly()` 等函数引入非线性成分，但需要指定复杂度参数。对可加模型

$$Y = \sum_{j=1}^P f_j(x_j) + \varepsilon$$

最好能从数据中自动确定  $f_j(\cdot)$  的复杂度（光滑度）参数。

R 扩展包 `mgcv` 的 `gam()` 函数可以执行这样的可加模型的非参数回归拟合。模型中可以用 `s(x)` 指定  $x$  的样条平滑，用 `lo(x)` 指定  $x$  的局部多项式平滑。具体的计算是迭代地对每个自变量  $x_j$  进行，估计  $x_j$  的平滑函数  $f_j(\cdot)$  时，采用数据  $\tilde{y} = Y - \sum_{k \neq j} f_k(x_k)$ ，迭代估计到结果基本不变为止。

例如，`MASS` 包的 `rock` 数据框是石油勘探中 12 块岩石样本分别产生 4 个切片得到的测量数据，共 48 个观测，因变量是渗透率 (permeability)，自变量为 `area`, `peri`, `shape`。

先作线性回归：

```
lm.rock <- lm(log(perm) ~ area + peri + shape, data=rock)
summary(lm.rock)

##
Call:
lm(formula = log(perm) ~ area + peri + shape, data = rock)
##
Residuals:
Min 1Q Median 3Q Max
-1.8092 -0.5413 0.1734 0.6493 1.4788
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.333e+00 5.487e-01 9.720 1.59e-12 ***
area 4.850e-04 8.657e-05 5.602 1.29e-06 ***
peri -1.527e-03 1.770e-04 -8.623 5.24e-11 ***
shape 1.757e+00 1.756e+00 1.000 0.323

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 0.8521 on 44 degrees of freedom
Multiple R-squared: 0.7483, Adjusted R-squared: 0.7311
F-statistic: 43.6 on 3 and 44 DF, p-value: 3.094e-13
```

其中 shape 变量不显著。可能的原因有：

- shape 与因变量没有关系；
- 样本量不足；
- shape 与因变量之间有非线性关系，该非线性无法用线性近似。

用样条平滑的 `gam()` 建模：

```
library(mgcv)

载入需要的程辑包：nlme

Warning: 程辑包'nlme'是用R版本3.6.1 来建造的

##
载入程辑包：'nlme'

The following object is masked from 'package:dplyr':
##
collapse

This is mgcv 1.8-28. For overview type 'help("mgcv-package")'.

gam.rock1 <- gam(
 log(perm) ~ s(area) + s(peri) + s(shape), data=rock)
summary(gam.rock1)

##
Family: gaussian
Link function: identity
##
Formula:
log(perm) ~ s(area) + s(peri) + s(shape)
##
Parametric coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.1075 0.1222 41.81 <2e-16 ***

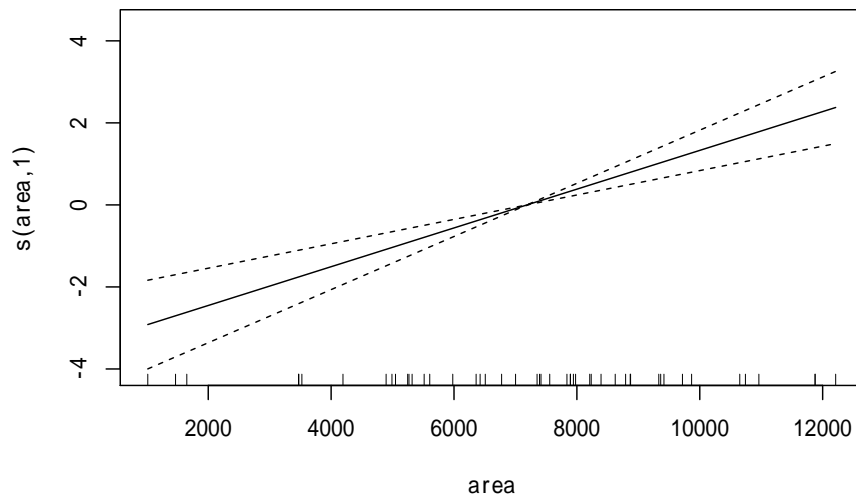
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
Approximate significance of smooth terms:
edf Ref.df F p-value
s(area) 1.000 1.000 29.13 1.96e-06 ***
s(peri) 1.000 1.000 71.30 4.12e-12 ***
s(shape) 1.402 1.705 0.58 0.425

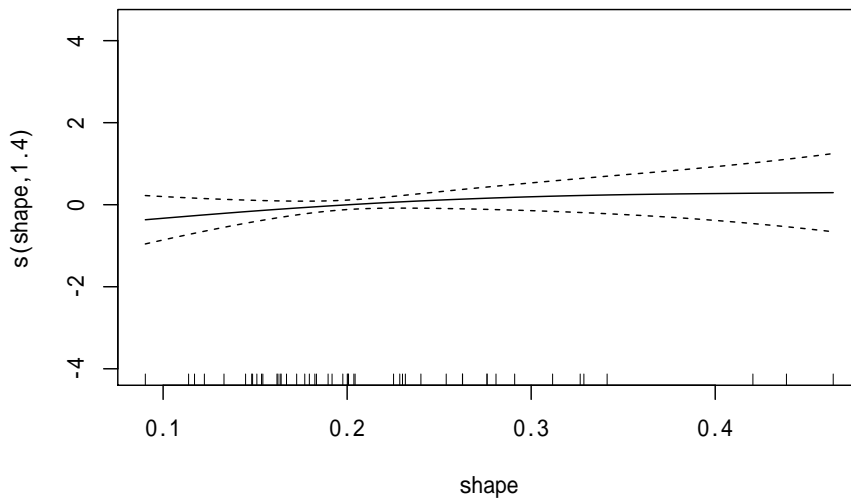
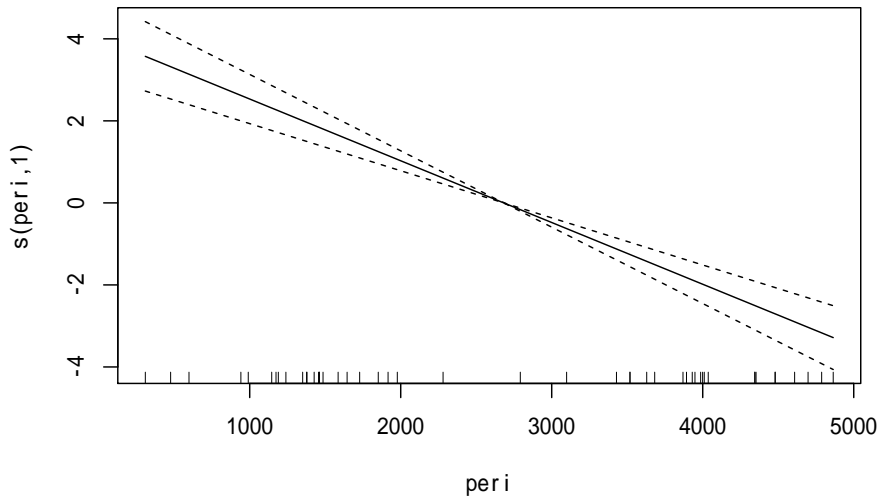
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
R-sq.(adj) = 0.735 Deviance explained = 75.4%
GCV = 0.78865 Scale est. = 0.71631 n = 48
```

对 `gam()` 回归结果做单个变量的曲线拟合图:

```
plot(gam.rock1)
```







可以看出三条曲线都基本不需要非线性。比较两个模型：

```
anova(lm.rock, gam.rock1)

Analysis of Variance Table
##
Model 1: log(perm) ~ area + peri + shape
Model 2: log(perm) ~ s(area) + s(peri) + s(shape)
Res.Df RSS Df Sum of Sq F Pr(>F)
1 44.000 31.949
2 43.598 31.230 0.40237 0.71914 2.4951 0.125
```

结果也不显著，说明线性模型是可取的。

## 33.5 Logistic 回归

### 33.5.1 模型

当因变量  $Y$  是零壹变量时，即  $Y$  表示分两类的类别，取值 1 和 0，我们关心的是  $P(Y = 1)$ 。这是一个  $[0, 1]$  区间内的值。如果把  $Y$  当作一般因变量做线性回归，会给出不合理的结果，比如负值，另外线性回归假定误差项为正态分布在这里也不适用。

为此考虑广义的回归模型 (广义线性模型):

$$Y \sim F(y; \theta)$$

$$g(\theta) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

特别地，定义 logit 函数

$$\text{logit}(p) = \ln \left( \frac{p}{1-p} \right)$$

Logit 模型:

$$Y \sim b(1, p)$$

$$\text{logit}(p) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

### 33.5.2 R 程序

进行 logistic 回归的 R 程序如 `glm(y ~ x1 + x2, family=binomial, data=d)`, 其中 `y` 为取 0 或 1 的向量, 输入数据集为 `d`。

`y` 也可以是一个两列的矩阵, 第一列为成功数, 第二列为失败数。

例如, “remiss.csv” 中保存了癌症病人康复的数据, 变量 `remiss` 为康复与否 (1 为康复, 0 为未康复), 另外的 6 个变量为可能影响康复概率的自变量。

程序:

```
d.remiss <- read.csv("remiss.csv", header=TRUE)
glm1 <- glm(
 remiss ~ cell+smear+infil+li+blast+temp,
 family=binomial, data=d.remiss)
summary(glm1)

##
Call:
glm(formula = remiss ~ cell + smear + infil + li + blast + temp,
family = binomial, data = d.remiss)
##
Deviance Residuals:
Min 1Q Median 3Q Max
-1.95165 -0.66491 -0.04372 0.74304 1.67069
##
Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) 58.0385 71.2364 0.815 0.4152
cell 24.6615 47.8377 0.516 0.6062
smear 19.2936 57.9500 0.333 0.7392
infil -19.6013 61.6815 -0.318 0.7507
li 3.8960 2.3371 1.667 0.0955
blast 0.1511 2.2786 0.066 0.9471
temp -87.4339 67.5735 -1.294 0.1957

```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
(Dispersion parameter for binomial family taken to be 1)
##
Null deviance: 34.372 on 26 degrees of freedom
Residual deviance: 21.751 on 20 degrees of freedom
AIC: 35.751
##
Number of Fisher Scoring iterations: 8
```

以  $p$  值 0.30 为界限, 逐步删去不显著的自变量:

```
glm1b <- glm(
 remiss ~ cell + smear + infil + li + temp,
 family=binomial, data=d.remiss)
summary(glm1b)
```

```
##
Call:
glm(formula = remiss ~ cell + smear + infil + li + temp, family = binomial,
data = d.remiss)
##
Deviance Residuals:
Min 1Q Median 3Q Max
-1.93753 -0.65855 -0.04328 0.75287 1.65475
##
Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) 57.128 69.977 0.816 0.414
cell 24.180 47.257 0.512 0.609
smear 18.370 56.218 0.327 0.744
infil -18.477 59.260 -0.312 0.755
li 3.987 1.902 2.097 0.036 *
temp -86.137 64.785 -1.330 0.184

```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
(Dispersion parameter for binomial family taken to be 1)
##
Null deviance: 34.372 on 26 degrees of freedom
Residual deviance: 21.755 on 21 degrees of freedom
AIC: 33.755
##
Number of Fisher Scoring iterations: 8

glm1c <- glm(
 remiss ~ cell + infil + li + temp,
 family=binomial, data=d.remiss)
summary(glm1c)

##
Call:
glm(formula = remiss ~ cell + infil + li + temp, family = binomial,
data = d.remiss)
##
Deviance Residuals:
Min 1Q Median 3Q Max
-1.88516 -0.66796 -0.07282 0.78697 1.72442
##
Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) 70.6171 59.1148 1.195 0.2323
cell 9.1434 7.9124 1.156 0.2479
infil 0.9088 3.1356 0.290 0.7719
li 3.9005 1.8108 2.154 0.0312 *
temp -85.2481 64.0897 -1.330 0.1835

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
(Dispersion parameter for binomial family taken to be 1)
##
Null deviance: 34.372 on 26 degrees of freedom
Residual deviance: 21.869 on 22 degrees of freedom
AIC: 31.869
##
Number of Fisher Scoring iterations: 7

glm1d <- glm(
 remiss ~ cell + li + temp,
 family=binomial, data=d.remiss)
summary(glm1d)

##
Call:
glm(formula = remiss ~ cell + li + temp, family = binomial, data = d.remiss)
##
Deviance Residuals:
Min 1Q Median 3Q Max
-2.02043 -0.66313 -0.08323 0.81282 1.65887
##
Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) 67.634 56.888 1.189 0.2345
cell 9.652 7.751 1.245 0.2130
li 3.867 1.778 2.175 0.0297 *
temp -82.074 61.712 -1.330 0.1835

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
(Dispersion parameter for binomial family taken to be 1)
##
Null deviance: 34.372 on 26 degrees of freedom
Residual deviance: 21.953 on 23 degrees of freedom
```

```
AIC: 29.953
##
Number of Fisher Scoring iterations: 7
```

或可用逐步回归:

```
glm2 <- step(glm(
 remiss ~ cell + smear + infil + li + blast + temp,
 family=binomial, data=d.remiss))
```

```
Start: AIC=35.75
remiss ~ cell + smear + infil + li + blast + temp
##
Df Deviance AIC
- blast 1 21.755 33.755
- infil 1 21.857 33.857
- smear 1 21.869 33.869
- cell 1 22.071 34.071
<none> 21.751 35.751
- temp 1 23.877 35.877
- li 1 26.095 38.095
##
Step: AIC=33.76
remiss ~ cell + smear + infil + li + temp
##
Df Deviance AIC
- infil 1 21.858 31.858
- smear 1 21.869 31.869
- cell 1 22.073 32.073
<none> 21.755 33.755
- temp 1 24.198 34.199
- li 1 30.216 40.216
##
Step: AIC=31.86
remiss ~ cell + smear + li + temp
```

```

##
Df Deviance AIC
- smear 1 21.953 29.953
<none> 21.858 31.858
- temp 1 24.292 32.292
- cell 1 24.477 32.477
- li 1 30.434 38.434
##
Step: AIC=29.95
remiss ~ cell + li + temp
##
Df Deviance AIC
<none> 21.953 29.953
- temp 1 24.341 30.341
- cell 1 24.648 30.648
- li 1 30.829 36.829

```

`summary(glm2)`

```

##
Call:
glm(formula = remiss ~ cell + li + temp, family = binomial, data = d.remiss)
##
Deviance Residuals:
Min 1Q Median 3Q Max
-2.02043 -0.66313 -0.08323 0.81282 1.65887
##
Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) 67.634 56.888 1.189 0.2345
cell 9.652 7.751 1.245 0.2130
li 3.867 1.778 2.175 0.0297 *
temp -82.074 61.712 -1.330 0.1835

```



```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
(Dispersion parameter for binomial family taken to be 1)
##
Null deviance: 34.372 on 26 degrees of freedom
Residual deviance: 21.953 on 23 degrees of freedom
AIC: 29.953
##
Number of Fisher Scoring iterations: 7
```



## Chapter 34

# 统计学习介绍

统计学习介绍的主要参考书: (James et al., 2013): Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani(2013) An Introduction to Statistical Learning: with Applications in R, Springer.

暂时只有运行例子, 缺少方法原理的讲解。参见 (James et al., 2013)。

调入需要的扩展包:

```
library(leaps) # 全子集回归
library(ISLR) # 参考书对应的包
library(glmnet) # 岭回归和 lasso
```

```
载入需要的程辑包: Matrix
```

```
载入需要的程辑包: foreach
```

```
Warning: 程辑包'foreach'是用R版本3.6.1 来建造的
```

```
Loaded glmnet 2.0-18
```

```
library(tree) # 树回归
library(randomForest) # 随机森林和装袋法
```

```
randomForest 4.6-14
```

```
Type rfNews() to see new features/changes/bug fixes.
```

```
library(MASS)
library(gbm) # boosting

Loaded gbm 2.1.5
```

## 34.1 Hitters 数据分析

考虑 ISLR 包的 Hitters 数据集。此数据集有 322 个运动员的 20 个变量的数据，其中的变量 Salary（工资）是我们关心的。变量包括：

```
names(Hitters)

[1] "AtBat" "Hits" "HmRun" "Runs" "RBI" "Walks" "Y
```

数据集的详细变量信息如下：

```
str(Hitters)

'data.frame': 322 obs. of 20 variables:
$ AtBat : int 293 315 479 496 321 594 185 298 323 401 ...
$ Hits : int 66 81 130 141 87 169 37 73 81 92 ...
$ HmRun : int 1 7 18 20 10 4 1 0 6 17 ...
$ Runs : int 30 24 66 65 39 74 23 24 26 49 ...
$ RBI : int 29 38 72 78 42 51 8 24 32 66 ...
$ Walks : int 14 39 76 37 30 35 21 7 8 65 ...
$ Years : int 1 14 3 11 2 11 2 3 2 13 ...
$ CAtBat : int 293 3449 1624 5628 396 4408 214 509 341 5206 ...
$ CHits : int 66 835 457 1575 101 1133 42 108 86 1332 ...
$ CHmRun : int 1 69 63 225 12 19 1 0 6 253 ...
$ CRuns : int 30 321 224 828 48 501 30 41 32 784 ...
$ CRBI : int 29 414 266 838 46 336 9 37 34 890 ...
$ CWalks : int 14 375 263 354 33 194 24 12 8 866 ...
$ League : Factor w/ 2 levels "A","N": 1 2 1 2 2 1 2 1 2 1 ...
$ Division : Factor w/ 2 levels "E","W": 1 2 2 1 1 2 1 2 2 1 ...
$ PutOuts : int 446 632 880 200 805 282 76 121 143 0 ...
$ Assists : int 33 43 82 11 40 421 127 283 290 0 ...
```

```
$ Errors : int 20 10 14 3 4 25 7 9 19 0 ...
$ Salary : num NA 475 480 500 91.5 750 70 100 75 1100 ...
$ NewLeague: Factor w/ 2 levels "A","N": 1 2 1 2 2 1 1 1 2 1 ...
```

希望以 Salary 为因变量，查看其缺失值个数：

```
sum(is.na(Hitters$Salary))
```

```
[1] 59
```

为简单起见，去掉有缺失值的观测：

```
d <- na.omit(Hitters); dim(d)
```

```
[1] 263 20
```

### 34.1.1 回归自变量选择

#### 34.1.1.1 最优子集选择

用 leaps 包的 `regsubsets()` 函数计算最优子集回归，办法是对某个试验性的子集自变量个数  $\hat{p}$  值，都找到  $\hat{p}$  固定情况下残差平方和最小的变量子集，这样只要在这些不同  $\hat{p}$  的最优子集中挑选就可以了。

可以先进行一个包含所有自变量的全集回归：

```
regfit.full <- regsubsets(Salary ~ ., data=d, nvmax=19)
reg.summary <- summary(regfit.full)
reg.summary
```

```
Subset selection object
Call: regsubsets.formula(Salary ~ ., data = d, nvmax = 19)
19 Variables (and intercept)
Forced in Forced out
AtBat FALSE FALSE
Hits FALSE FALSE
HmRun FALSE FALSE
Runs FALSE FALSE
RBI FALSE FALSE
```



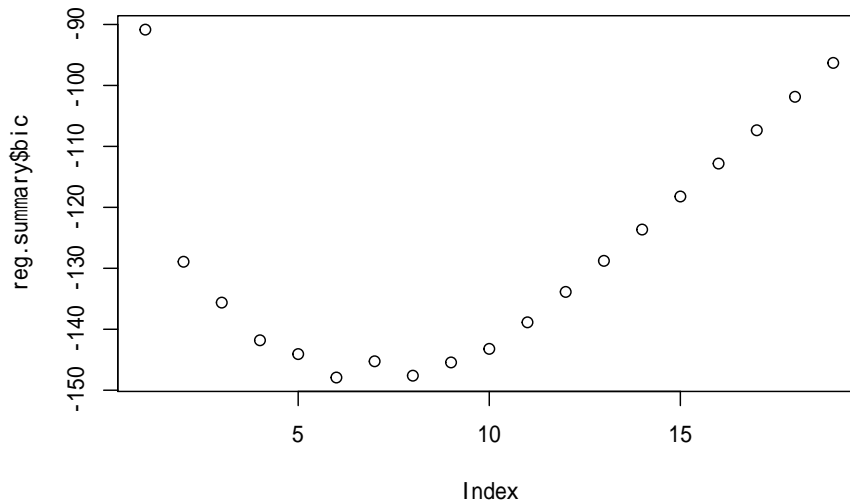


图 34.1: Hitters 数据最优子集回归 BIC

```
17 (1) "*" "*" "*" "*" "*" "*" " " "*" "*" " " "*" "*" "*" "*"
18 (1) "*" "*" "*" "*" "*" "*" "*" "*" "*" " " "*" "*" "*" "*"
19 (1) "*" "*" "*" "*" "*" "*" "*" "*" "*" "*" "*" "*" "*" "*" "
```

这里用 `nvmax=` 指定了允许所有的自变量都参加, 缺省行为是限制最多个数的。  
上述结果表格中每一行给出了固定  $\hat{p}$  条件下的最优子集。

试比较这些最优模型的 BIC 值:

```
reg.summary$bic
```

```
[1] -90.84637 -128.92622 -135.62693 -141.80892 -144.07143 -147.91690 -145.25594 -147.61
```

```
plot(reg.summary$bic)
```

其中  $\hat{p} = 6, 8$  的值相近, 都很低, 取  $\hat{p} = 6$ 。用 `coef()` 加 `id=6` 指定第六种子集:

```
coef(regfit.full, id=6)
```

```
(Intercept) AtBat Hits Walks CRBI DivisionW
91.5117981 -1.8685892 7.6043976 3.6976468 0.6430169 -122.9515338
```

这种方法实现了选取 BIC 最小的自变量子集。

### 34.1.1.2 逐步回归方法

在用 `lm()` 做了全集回归后，把全集回归结果输入到 `step()` 函数中可以执行逐步回归。如：

```
lm.full <- lm(Salary ~ ., data=d)
print(summary(lm.full))

##
Call:
lm(formula = Salary ~ ., data = d)
##
Residuals:
Min 1Q Median 3Q Max
-907.62 -178.35 -31.11 139.09 1877.04
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 163.10359 90.77854 1.797 0.073622 .
AtBat -1.97987 0.63398 -3.123 0.002008 **
Hits 7.50077 2.37753 3.155 0.001808 **
HmRun 4.33088 6.20145 0.698 0.485616
Runs -2.37621 2.98076 -0.797 0.426122
RBI -1.04496 2.60088 -0.402 0.688204
Walks 6.23129 1.82850 3.408 0.000766 ***
Years -3.48905 12.41219 -0.281 0.778874
CAtBat -0.17134 0.13524 -1.267 0.206380
CHits 0.13399 0.67455 0.199 0.842713
CHmRun -0.17286 1.61724 -0.107 0.914967
CRuns 1.45430 0.75046 1.938 0.053795 .
```



```
CRBI 0.80771 0.69262 1.166 0.244691
CWalks -0.81157 0.32808 -2.474 0.014057 *
LeagueN 62.59942 79.26140 0.790 0.430424
DivisionW -116.84925 40.36695 -2.895 0.004141 **
PutOuts 0.28189 0.07744 3.640 0.000333 ***
Assists 0.37107 0.22120 1.678 0.094723 .
Errors -3.36076 4.39163 -0.765 0.444857
NewLeagueN -24.76233 79.00263 -0.313 0.754218

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 315.6 on 243 degrees of freedom
Multiple R-squared: 0.5461, Adjusted R-squared: 0.5106
F-statistic: 15.39 on 19 and 243 DF, p-value: < 2.2e-16

lm.step <- step(lm.full)
```

```
Start: AIC=3046.02
Salary ~ AtBat + Hits + HmRun + Runs + RBI + Walks + Years +
CAtBat + CHits + CHmRun + CRuns + CRBI + CWalks + League +
Division + PutOuts + Assists + Errors + NewLeague
##
Df Sum of Sq RSS AIC
- CHmRun 1 1138 24201837 3044.0
- CHits 1 3930 24204629 3044.1
- Years 1 7869 24208569 3044.1
- NewLeague 1 9784 24210484 3044.1
- RBI 1 16076 24216776 3044.2
- HmRun 1 48572 24249272 3044.6
- Errors 1 58324 24259023 3044.7
- League 1 62121 24262821 3044.7
- Runs 1 63291 24263990 3044.7
- CRBI 1 135439 24336138 3045.5
- CAtBat 1 159864 24360564 3045.8
```

```

<none> 24200700 3046.0
- Assists 1 280263 24480963 3047.1
- CRuns 1 374007 24574707 3048.1
- CWalks 1 609408 24810108 3050.6
- Division 1 834491 25035190 3052.9
- AtBat 1 971288 25171987 3054.4
- Hits 1 991242 25191941 3054.6
- Walks 1 1156606 25357305 3056.3
- PutOuts 1 1319628 25520328 3058.0
##
Step: AIC=3044.03
Salary ~ AtBat + Hits + HmRun + Runs + RBI + Walks + Years +
CAtBat + CHits + CRuns + CRBI + CWalks + League + Division +
PutOuts + Assists + Errors + NewLeague
##
Df Sum of Sq RSS AIC
- Years 1 7609 24209447 3042.1
- NewLeague 1 10268 24212106 3042.2
- CHits 1 14003 24215840 3042.2
- RBI 1 14955 24216793 3042.2
- HmRun 1 52777 24254614 3042.6
- Errors 1 59530 24261367 3042.7
- League 1 63407 24265244 3042.7
- Runs 1 64860 24266698 3042.7
- CAtBat 1 174992 24376830 3043.9
<none> 24201837 3044.0
- Assists 1 285766 24487603 3045.1
- CRuns 1 611358 24813196 3048.6
- CWalks 1 645627 24847464 3049.0
- Division 1 834637 25036474 3050.9
- CRBI 1 864220 25066057 3051.3
- AtBat 1 970861 25172699 3052.4
- Hits 1 1025981 25227819 3052.9

```

```

- Walks 1 1167378 25369216 3054.4
- PutOuts 1 1325273 25527110 3056.1
##
Step: AIC=3042.12
Salary ~ AtBat + Hits + HmRun + Runs + RBI + Walks + CAtBat +
CHits + CRuns + CRBI + CWalks + League + Division + PutOuts +
Assists + Errors + NewLeague
##
Df Sum of Sq RSS AIC
- NewLeague 1 9931 24219377 3040.2
- RBI 1 15989 24225436 3040.3
- CHits 1 18291 24227738 3040.3
- HmRun 1 54144 24263591 3040.7
- Errors 1 57312 24266759 3040.7
- Runs 1 63172 24272619 3040.8
- League 1 65732 24275178 3040.8
<none> 24209447 3042.1
- CAtBat 1 266205 24475652 3043.0
- Assists 1 293479 24502926 3043.3
- CRuns 1 646350 24855797 3047.1
- CWalks 1 649269 24858716 3047.1
- Division 1 827511 25036958 3049.0
- CRBI 1 872121 25081568 3049.4
- AtBat 1 968713 25178160 3050.4
- Hits 1 1018379 25227825 3050.9
- Walks 1 1164536 25373983 3052.5
- PutOuts 1 1334525 25543972 3054.2
##
Step: AIC=3040.22
Salary ~ AtBat + Hits + HmRun + Runs + RBI + Walks + CAtBat +
CHits + CRuns + CRBI + CWalks + League + Division + PutOuts +
Assists + Errors
##

```

```

Df Sum of Sq RSS AIC
- RBI 1 15800 24235177 3038.4
- CHits 1 15859 24235237 3038.4
- Errors 1 54505 24273883 3038.8
- HmRun 1 54938 24274316 3038.8
- Runs 1 62294 24281671 3038.9
- League 1 107479 24326856 3039.4
<none> 24219377 3040.2
- CAtBat 1 261336 24480713 3041.1
- Assists 1 295536 24514914 3041.4
- CWalks 1 648860 24868237 3045.2
- CRuns 1 661449 24880826 3045.3
- Division 1 824672 25044049 3047.0
- CRBI 1 880429 25099806 3047.6
- AtBat 1 999057 25218434 3048.9
- Hits 1 1034463 25253840 3049.2
- Walks 1 1157205 25376583 3050.5
- PutOuts 1 1335173 25554550 3052.3
##
Step: AIC=3038.4
Salary ~ AtBat + Hits + HmRun + Runs + Walks + CAtBat + CHits +
CRuns + CRBI + CWalks + League + Division + PutOuts + Assists +
Errors
##
Df Sum of Sq RSS AIC
- CHits 1 13483 24248660 3036.5
- HmRun 1 44586 24279763 3036.9
- Runs 1 54057 24289234 3037.0
- Errors 1 57656 24292833 3037.0
- League 1 108644 24343821 3037.6
<none> 24235177 3038.4
- CAtBat 1 252756 24487934 3039.1
- Assists 1 294674 24529851 3039.6

```

```

- CWalks 1 639690 24874868 3043.2
- CRuns 1 693535 24928712 3043.8
- Division 1 808984 25044161 3045.0
- CRBI 1 893830 25129008 3045.9
- Hits 1 1034884 25270061 3047.4
- AtBat 1 1042798 25277975 3047.5
- Walks 1 1145013 25380191 3048.5
- PutOuts 1 1340713 25575890 3050.6
##
Step: AIC=3036.54
Salary ~ AtBat + Hits + HmRun + Runs + Walks + CAtBat + CRuns +
CRBI + CWalks + League + Division + PutOuts + Assists + Errors
##
Df Sum of Sq RSS AIC
- HmRun 1 40487 24289148 3035.0
- Errors 1 51930 24300590 3035.1
- Runs 1 79343 24328003 3035.4
- League 1 114742 24363402 3035.8
<none> 24248660 3036.5
- Assists 1 283442 24532103 3037.6
- CAtBat 1 613356 24862016 3041.1
- Division 1 801474 25050134 3043.1
- CRBI 1 903248 25151908 3044.2
- CWalks 1 1011953 25260613 3045.3
- Walks 1 1246164 25494824 3047.7
- AtBat 1 1339620 25588280 3048.7
- CRuns 1 1390808 25639469 3049.2
- PutOuts 1 1406023 25654684 3049.4
- Hits 1 1607990 25856650 3051.4
##
Step: AIC=3034.98
Salary ~ AtBat + Hits + Runs + Walks + CAtBat + CRuns + CRBI +
CWalks + League + Division + PutOuts + Assists + Errors

```

```

##
Df Sum of Sq RSS AIC
- Errors 1 44085 24333232 3033.5
- Runs 1 49068 24338215 3033.5
- League 1 103837 24392985 3034.1
<none> 24289148 3035.0
- Assists 1 247002 24536150 3035.6
- CAtBat 1 652746 24941894 3040.0
- Division 1 795643 25084791 3041.5
- CWalks 1 982896 25272044 3043.4
- Walks 1 1205823 25494971 3045.7
- AtBat 1 1300972 25590120 3046.7
- CRuns 1 1351200 25640348 3047.2
- CRBI 1 1353507 25642655 3047.2
- PutOuts 1 1429006 25718154 3048.0
- Hits 1 1574140 25863288 3049.5
##
Step: AIC=3033.46
Salary ~ AtBat + Hits + Runs + Walks + CAtBat + CRuns + CRBI +
CWalks + League + Division + PutOuts + Assists
##
Df Sum of Sq RSS AIC
- Runs 1 54113 24387345 3032.0
- League 1 91269 24424501 3032.4
<none> 24333232 3033.5
- Assists 1 220010 24553242 3033.8
- CAtBat 1 650513 24983746 3038.4
- Division 1 799455 25132687 3040.0
- CWalks 1 971260 25304493 3041.8
- Walks 1 1239533 25572765 3044.5
- CRBI 1 1331672 25664904 3045.5
- CRuns 1 1361070 25694302 3045.8
- AtBat 1 1378592 25711824 3045.9

```

```

- PutOuts 1 1391660 25724892 3046.1
- Hits 1 1649291 25982523 3048.7
##
Step: AIC=3032.04
Salary ~ AtBat + Hits + Walks + CAtBat + CRuns + CRBI + CWalks +
League + Division + PutOuts + Assists
##
Df Sum of Sq RSS AIC
- League 1 113056 24500402 3031.3
<none> 24387345 3032.0
- Assists 1 280689 24668034 3033.1
- CAtBat 1 596622 24983967 3036.4
- Division 1 780369 25167714 3038.3
- CWalks 1 946687 25334032 3040.1
- Walks 1 1212997 25600342 3042.8
- CRuns 1 1334397 25721742 3044.1
- CRBI 1 1361339 25748684 3044.3
- PutOuts 1 1455210 25842555 3045.3
- AtBat 1 1522760 25910105 3046.0
- Hits 1 1718870 26106215 3047.9
##
Step: AIC=3031.26
Salary ~ AtBat + Hits + Walks + CAtBat + CRuns + CRBI + CWalks +
Division + PutOuts + Assists
##
Df Sum of Sq RSS AIC
<none> 24500402 3031.3
- Assists 1 313650 24814051 3032.6
- CAtBat 1 534156 25034558 3034.9
- Division 1 798473 25298875 3037.7
- CWalks 1 965875 25466276 3039.4
- CRuns 1 1265082 25765484 3042.5
- Walks 1 1290168 25790569 3042.8

```

```
- CRBI 1 1326770 25827172 3043.1
- PutOuts 1 1551523 26051925 3045.4
- AtBat 1 1589780 26090181 3045.8
- Hits 1 1716068 26216469 3047.1
```

```
print(lm.step)
```

```
##
Call:
lm(formula = Salary ~ AtBat + Hits + Walks + CAtBat + CRuns +
CRBI + CWalks + Division + PutOuts + Assists, data = d)
##
Coefficients:
(Intercept) AtBat Hits Walks CAtBat CRuns
162.5354 -2.1687 6.9180 5.7732 -0.1301 1.4082
```

最后保留了 10 个自变量。

### 34.1.1.3 划分训练集与测试集

在整个数据集中随机选取一部分作为训练集，其余作为测试集。下面的程序把原始数据大约一分为二：

```
set.seed(1)
train <- sample(c(TRUE, FALSE), nrow(d), replace=TRUE)
test <- (!train)
```

仅用训练集估计模型。为了在测试集上用模型进行预报并估计预测均方误差，需要自己写一个预测函数：

```
predict.regsubsets <- function(object, newdata, id, ...){
 form <- as.formula(object$call[[2]])
 mat <- model.matrix(form, newdata)
 coefi <- coef(object, id=id)
 xvars <- names(coefi)
 mat[, xvars] %*% coefi
}
```



然后，对每个子集大小，用最优子集在测试集上进行预报，计算均方误差：

```
regfit.best <- regsubsets(Salary ~ ., data=d[train,], nvmax=19)
val.errors <- rep(as.numeric(NA), 19)
for(i in 1:19){
 #pred <- predict.regsubsets(regfit.best, newdata=d[test,], id=i)
 pred <- predict(regfit.best, newdata=d[test,], id=i)
 val.errors[i] <- mean((d[test, 'Salary'] - pred)^2)
}
print(val.errors)
```

```
[1] 164377.3 144405.5 152175.7 145198.4 137902.1 139175.7 126849.0 136191.4 132889.6 135
```

```
best.id <- which.min(val.errors); best.id
```

```
[1] 7
```

用测试集得到的最优子集大小为 7。模型子集和回归系数为：

```
coef(regfit.best, id=best.id)
```

```
(Intercept) AtBat Hits Walks CRuns CWalks Division
67.1085369 -2.1462987 7.0149547 8.0716640 1.2425113 -0.8337844 -118.436499
```

#### 34.1.1.4 用 10 折交叉验证方法选择最优子集

下列程序对数据中每一行分配一个折号：

```
k <- 10
set.seed(1)
folds <- sample(1:k, nrow(d), replace=TRUE)
```

下面，对 10 折中每一折都分别当作测试集一次，得到不同子集大小的均方误差：

```
cv.errors <- matrix(as.numeric(NA), k, 19, dimnames=list(NULL, paste(1:19)))
for(j in 1:k){ # 对
 best.fit <- regsubsets(Salary ~ ., data=d[folds != j,], nvmax=19)
 for(i in 1:19){
```

```

 pred <- predict(best.fit, d[folds==j,], id=i)
 cv.errors[j, i] <- mean((d[folds==j, 'Salary'] - pred)^2)
 }
}
head(cv.errors)

```

```

1 2 3 4 5 6 7
[1,] 98623.24 115600.61 120884.31 113831.63 120728.51 122922.93 155507.25 1377
[2,] 155320.11 100425.87 168838.35 159729.47 145895.71 123555.25 119983.35 966
[3,] 124151.77 68833.50 69392.29 77221.37 83802.82 70125.41 68997.77 641
[4,] 232191.41 279001.29 294568.10 288765.81 276972.83 260121.22 276413.09 2599
[5,] 115397.35 96807.44 108421.66 104933.55 99561.69 86103.05 89345.61 876
[6,] 103839.30 75652.50 69962.31 58291.91 65893.45 64215.56 65800.88 614

```

`cv.errors` 是一个  $10 \times 19$  矩阵，每行对应一折作为测试集的情形，每列是一个子集大小，元素值是测试均方误差。

对每列的 10 个元素求平均，可以得到每个子集大小的平均均方误差：

```

mean.cv.errors <- apply(cv.errors, 2, mean)
mean.cv.errors

```

```

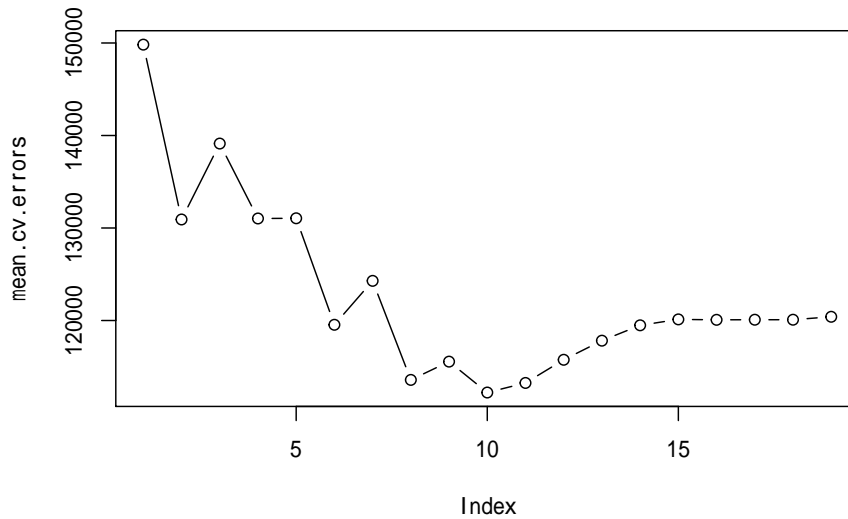
1 2 3 4 5 6 7 8
149821.1 130922.0 139127.0 131028.8 131050.2 119538.6 124286.1 113580.0 115556.

```

```

best.id <- which.min(mean.cv.errors)
plot(mean.cv.errors, type='b')

```



这样找到的最优子集大小是 10。用这种方法找到最优子集大小后，可以对全数据集重新建模但是选择最优子集大小为 10:

```
reg.best <- regsubsets(Salary ~ ., data=d, nvmax=19)
coef(reg.best, id=best.id)
```

```
(Intercept) AtBat Hits Walks CAtBat CRuns CRE
162.5354420 -2.1686501 6.9180175 5.7732246 -0.1300798 1.4082490 0.774312
```

事实上，划分训练集和验证集与交叉验证方法经常联合运用。取一个固定的较小规模的测试集，此测试集不用来作子集选择，对训练集用交叉验证方法选择最优子集，然后再测试集上验证。

### 34.1.2 岭回归

仍采用上面去掉了缺失值的 Hitters 数据集结果 d。需要使用 glmnet 包。

如下程序把回归的设计阵与因变量提取出来:

```
x <- model.matrix(Salary ~ ., d)[,-1]
y <- d$Salary
```

岭回归涉及到调节参数  $\lambda$  的选择，为了绘图，先选择  $\lambda$  的一个网格：

```
grid <- 10^seq(10, -2, length=100)
```

用所有数据针对这样的调节参数网格计算岭回归结果，注意 `glmnet()` 函数允许调节参数  $\lambda$  输入多个值：

```
ridge.mod <- glmnet(x, y, alpha=0, lambda=grid)
dim(coef(ridge.mod))
```

```
[1] 20 100
```

`glmnet()` 函数默认对数据进行标准化。

`coef()` 的结果是一个矩阵，每列对应一个调节参数值。

### 34.1.2.1 划分训练集与测试集

如下程序把数据分为一半训练、一半测试：

```
set.seed(1)
train <- sample(1:nrow(x), nrow(x)/2)
test <- (-train)
y.test <- y[test]
```

仅用测试集建立岭回归：

```
ridge.mod <- glmnet(x[train,], y[train], alpha=0, lambda=grid, thresh=1E-12)
```

用建立的模型对测试集进行预测，并计算调节参数等于 4 时的均方误差：

```
ridge.pred <- predict(ridge.mod, s=4, newx=x[test,])
mean((ridge.pred - y.test)^2)
```

```
[1] 142199.2
```

如果用因变量平均值作预测，这是最差的预测：

```
mean((mean(y[train]) - y.test)^2)
```

```
[1] 224669.9
```

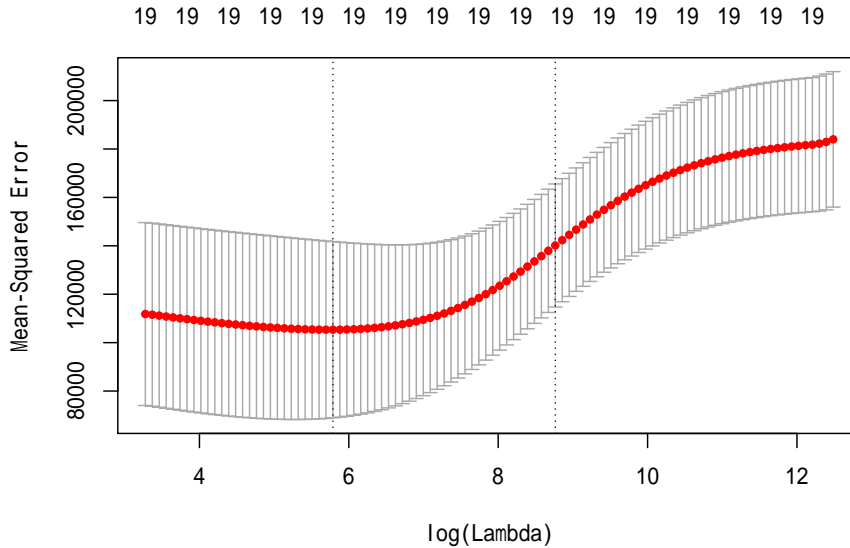


图 34.2: Hitters 数据岭回归参数选择

$\lambda = 4$  的结果要好得多。事实上，取  $\lambda$  接近正无穷时模型就相当于用因变量平均值预测。取  $\lambda = 0$  就相当于普通最小二乘回归（但是 `glmnet()` 是对输入数据要做标准化的）。

### 34.1.2.2 用 10 折交叉验证选取调节参数

仍使用训练集，但训练集再进行交叉验证。`cv.glmnet()` 函数可以执行交叉验证。

```
set.seed(1)
cv.out <- cv.glmnet(x[train,], y[train], alpha=0)
plot(cv.out)
```

```
bestlam <- cv.out$lambda.min
```

这样获得了最优调节参数  $\lambda = 326.0827885$ 。用最优调节参数对测试集作预测，得到预测均方误差：

```
ridge.pred <- predict(ridge.mod, s=bestlam, newx=x[test,])
mean((ridge.pred - y.test)^2)
```

```
[1] 139856.6
```

结果比  $\lambda = 4$  略有改进。

最后，用选取的最优调节系数对全数据集建模，得到相应的岭回归系数估计：

```
out <- glmnet(x, y, alpha=0)
predict(out, type='coefficients', s=bestlam)[1:20,]
```

```
(Intercept) AtBat Hits HmRun Runs RBI
15.44383135 0.07715547 0.85911581 0.60103107 1.06369007 0.87936105
```

### 34.1.3 Lasso 回归

仍使用 `glmnet` 包的 `glmnet()` 函数计算 Lasso 回归，指定一个调节参数网格（沿用前面的网格）：

```
lasso.mod <- glmnet(x[train,], y[train], alpha=1, lambda=grid)
plot(lasso.mod)
```

```
Warning in regularize.values(x, y, ties, missing(ties)): collapsing to unique '
```

对 lasso 结果使用 `plot()` 函数可以绘制调节参数网格变化的各回归系数估计，横坐标不是调节参数而是调节参数对应的系数绝对值和，可以看出随着系数绝对值和增大，实际是调节参数变小，更多地自变量进入模型。

#### 34.1.3.1 用交叉验证估计调节参数

按照前面划分的训练集与测试集，仅使用训练集数据做交叉验证估计最优调节参数：

```
set.seed(1)
cv.out <- cv.glmnet(x[train,], y[train], alpha=1)
plot(cv.out)
```

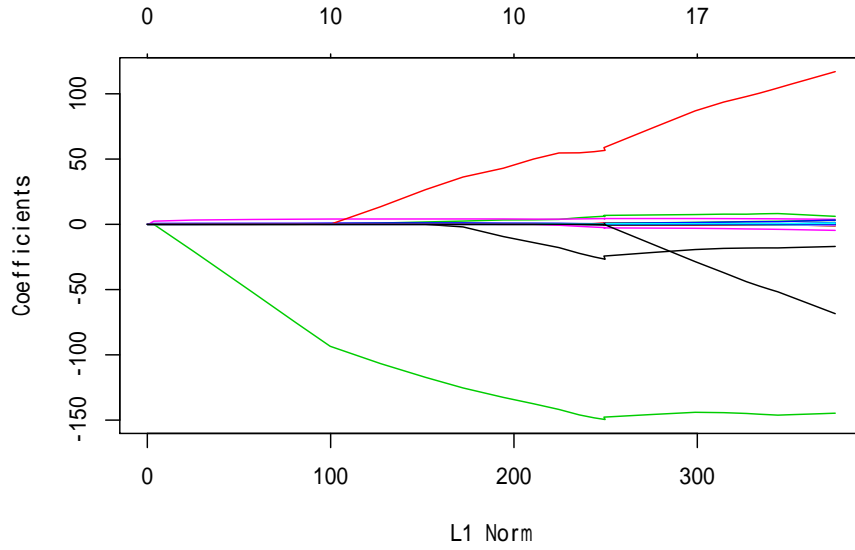
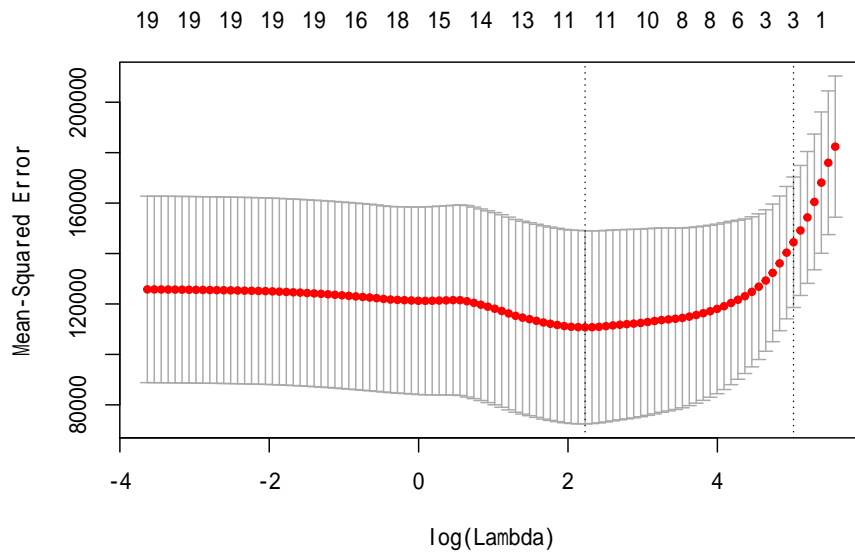


图 34.3: Hitters 数据 lasso 轨迹



```
bestlam <- cv.out$lambda.min; bestlam
```

```
[1] 9.286955
```

得到调节参数估计后，对测试集计算预测均方误差：

```
lasso.pred <- predict(lasso.mod, s=bestlam, newx=x[test,])
mean((lasso.pred - y.test)^2)
```

```
[1] 143673.6
```

这个效果比岭回归效果略差。

为了充分利用数据，使用前面获得的最优调节参数，对全数据集建模：

```
out <- glmnet(x, y, alpha=1, lambda=grid)
lasso.coef <- predict(out, type='coefficients', s=bestlam)[1:20,]; lasso.coef
```

```
(Intercept) AtBat Hits HmRun Runs
1.27479059 -0.05497143 2.18034583 0.00000000 0.00000000 0.0000
```

```
lasso.coef[lasso.coef != 0]
```

```
(Intercept) AtBat Hits Walks Years CH
1.27479059 -0.05497143 2.18034583 2.29192406 -0.33806109 0.0282
```

选择的自变量子集有 11 个自变量。

#### 34.1.4 树回归的简单演示

对 Hitters 数据，用 Years 和 Hits 作因变量预测  $\log(\text{Salary})$ 。

仅取 Hitters 数据集的 Salary, Years, Hits 三个变量，并仅保留完全观测：

```
d <- na.omit(Hitters[,c('Salary', 'Years', 'Hits')])
print(str(d))
```

```
'data.frame': 263 obs. of 3 variables:
$ Salary: num 475 480 500 91.5 750 ...
$ Years : int 14 3 11 2 11 2 3 2 13 10 ...
$ Hits : int 81 130 141 87 169 37 73 81 92 159 ...
```



```
- attr(*, "na.action")= 'omit' Named int 1 16 19 23 31 33 37 39 40 42 ...
..- attr(*, "names")= chr "-Andy Allanson" "-Billy Beane" "-Bruce Bochte" "-Bob Boone"
NULL
```

建立完整的树:

```
tr1 <- tree(log(Salary) ~ Years + Hits, data=d)
```

剪枝为只有 3 个叶结点:

```
tr1b <- prune.tree(tr1, best=3)
```

显示树:

```
print(tr1b)
```

```
node), split, n, deviance, yval
* denotes terminal node
##
1) root 263 207.20 5.927
2) Years < 4.5 90 42.35 5.107 *
3) Years > 4.5 173 72.71 6.354
6) Hits < 117.5 90 28.09 5.998 *
7) Hits > 117.5 83 20.88 6.740 *
```

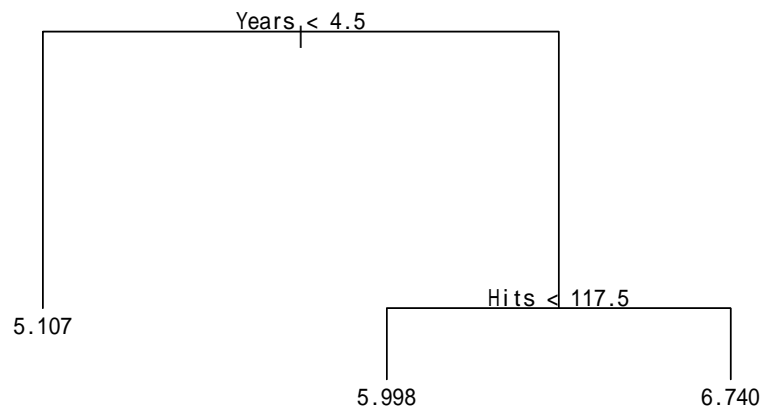
显示概括:

```
print(summary(tr1b))
```

```
##
Regression tree:
snip.tree(tree = tr1, nodes = c(6L, 2L))
Number of terminal nodes: 3
Residual mean deviance: 0.3513 = 91.33 / 260
Distribution of residuals:
Min. 1st Qu. Median Mean 3rd Qu. Max.
-2.24000 -0.39580 -0.03162 0.00000 0.33380 2.55600
```

做树图:

```
plot(tr1b); text(tr1b, pretty=0)
```



### 34.1.5 树回归

把数据随机地分成一半训练集，一半测试集：

```
d <- na.omit(Hitters)
set.seed(1)
train <- sample(c(TRUE, FALSE), nrow(d), replace=TRUE)
test <- (!train)
```

对训练集，建立未剪枝的树：

```
tr1 <- tree(log(Salary) ~ ., data=d, subset=train)
plot(tr1); text(tr1, pretty=0)
```

对训练集上的未剪枝树用交叉验证方法寻找最优大小：

```
cv1 <- cv.tree(tr1)
print(cv1)
```

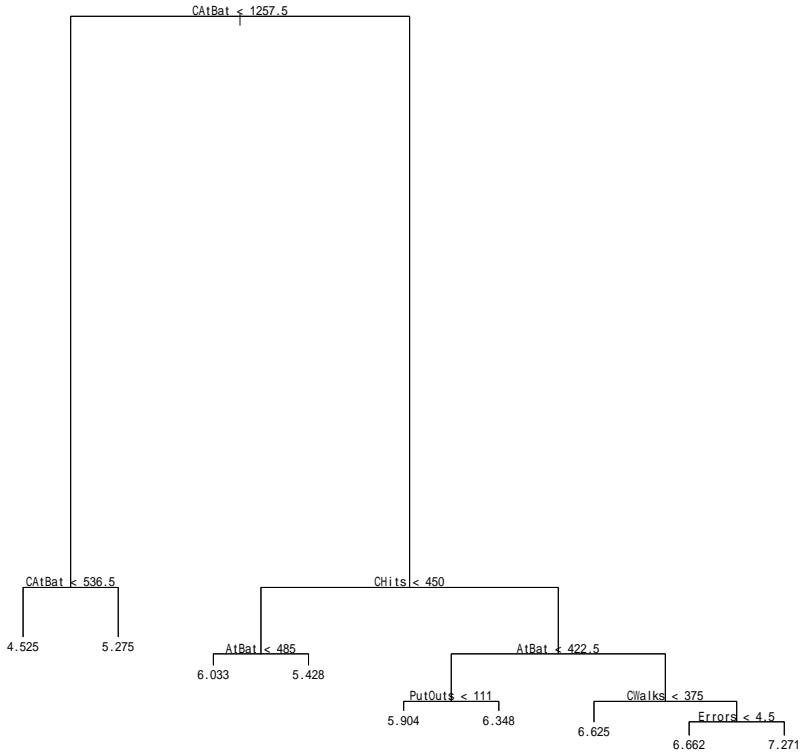
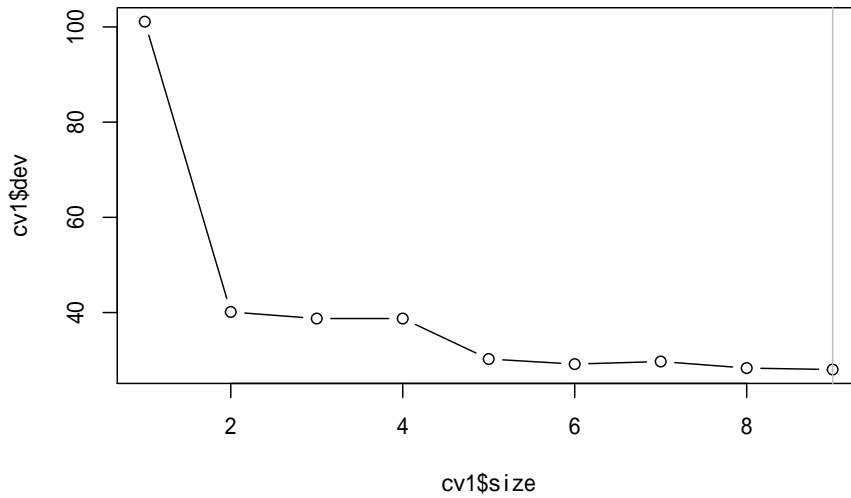


图 34.4: Hitters 数据训练集未剪枝树

```
$size
[1] 9 8 7 6 5 4 3 2 1
##
$dev
[1] 28.02658 28.33020 29.70951 29.16671 30.23694 38.74208 38.74208 40.1
##
$k
[1] -Inf 1.020847 1.256752 1.480873 2.283294 5.309848 5.505460 7.41
##
$method
[1] "deviance"
##
attr(,"class")
[1] "prune" "tree.sequence"
```

```
plot(cv1$size, cv1$dev, type='b')
best.size <- cv1$size[which.min(cv1$dev)[1]]
abline(v=best.size, col='gray')
```



最优大小为 9。获得训练集上构造的树剪枝后的结果：

```
best.size <- 4
tr1b <- prune.tree(tr1, best=best.size)
```

在测试集上计算预测均方误差：

```
pred.test <- predict(tr1b, newdata=d[test,])
test.mse <- mean((d[test, 'Salary'] - exp(pred.test))^2)
test.mse
```

```
[1] 142130.9
```

如果用训练集的因变量平均值估计测试集的因变量值，均方误差为：

```
worst.mse <- mean((d[test, 'Salary'] - mean(d[train, 'Salary']))^2)
worst.mse
```

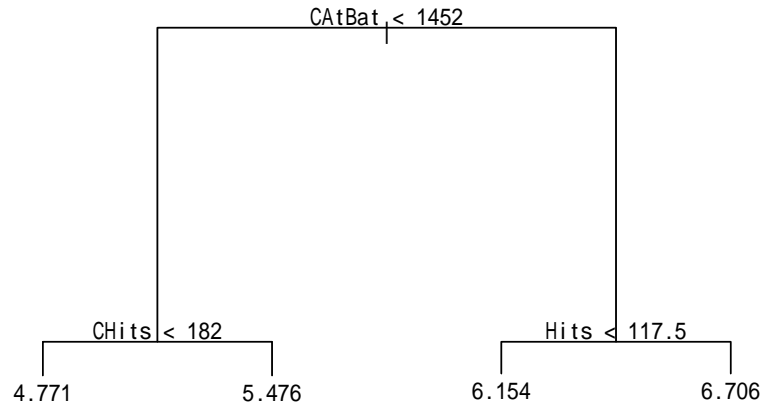
```
[1] 208338.9
```

用所有数据来构造未剪枝树：

```
tr2 <- tree(log(Salary) ~ ., data=d)
```

用训练集上得到的子树大小剪枝：

```
tr2b <- prune.tree(tr2, best=best.size)
plot(tr2b); text(tr2b, pretty=0)
```



### 34.1.6 装袋法

对训练集用装袋法:

```
bag1 <- randomForest(log(Salary) ~ ., data=d, subset=train, mtry=ncol(d)-1, import
bag1
```

```
##
Call:
randomForest(formula = log(Salary) ~ ., data = d, mtry = ncol(d) - 1, imp
Type of random forest: regression
Number of trees: 500
No. of variables tried at each split: 19
##
Mean of squared residuals: 0.1583502
% Var explained: 78.75
```

注意 `randomForest()` 函数实际是随机森林法, 但是当 `mtry` 的值取为所有自变量个数时就是装袋法。

对测试集进行预报:

```
pred2 <- predict(bag1, newdata=d[test,])
test.mse2 <- mean((d[test, 'Salary'] - exp(pred2))^2)
test.mse2
```

```
[1] 106686.6
```

结果与剪枝过的单课树相近。

在全集上使用装袋法:

```
bag2 <- randomForest(log(Salary) ~ ., data=d, mtry=ncol(d)-1, importance=TRUE)
bag2
```

```
##
```

```
Call:
```

```
randomForest(formula = log(Salary) ~ ., data = d, mtry = ncol(d) - 1, importance =
```

```
Type of random forest: regression
```

```
Number of trees: 500
```

```
No. of variables tried at each split: 19
```

```
##
```

```
Mean of squared residuals: 0.1907517
```

```
% Var explained: 75.78
```

变量的重要度数值和图形: 各变量的重要度数值及其图形:

```
importance(bag2)
```

```
%IncMSE IncNodePurity
AtBat 11.13812257 8.8307344
Hits 9.63201512 7.9785846
HmRun 3.22281779 2.0841188
Runs 7.20382600 3.7652715
RBI 7.79546620 5.6717192
Walks 8.96658890 7.1184930
Years 9.05717562 2.6408069
CAtBat 26.90155300 79.8147114
CHits 12.71002628 23.4661172
```

bag2

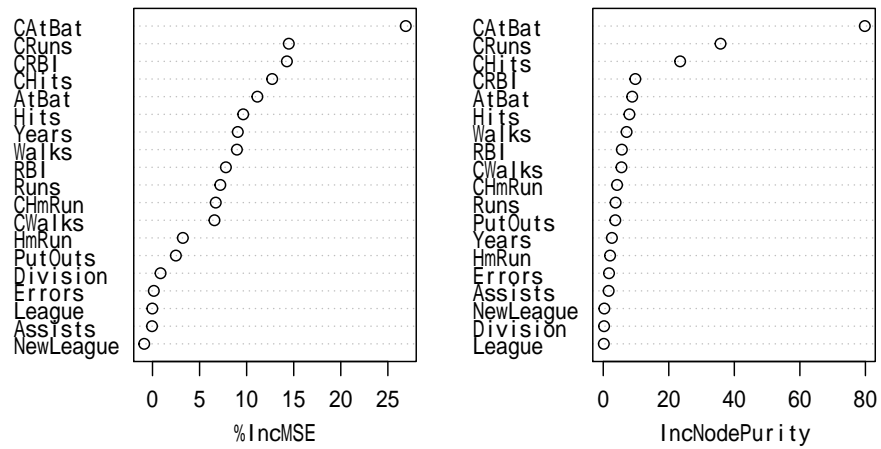


图 34.5: Hitters 数据装袋法的变量重要性结果

```
CHmRun 6.71648501 4.2139422
CRuns 14.47664565 35.7880779
CRBI 14.27594587 9.7903934
CWalks 6.58161583 5.5544426
League -0.02098083 0.1803449
Division 0.85552484 0.2508595
PutOuts 2.48153059 3.6745190
Assists -0.05144792 1.6443579
Errors 0.14020409 1.7566527
NewLeague -0.88103440 0.3243150
```

```
varImpPlot(bag2)
```

最重要的自变量是 CAtBats, 其次有 CRuns, CHits 等。



### 34.1.7 随机森林

对训练集用随机森林法:

```
rf1 <- randomForest(log(Salary) ~ ., data=d, subset=train, importance=TRUE)
rf1
```

```
##
Call:
randomForest(formula = log(Salary) ~ ., data = d, importance = TRUE, subset = train)
Type of random forest: regression
Number of trees: 500
No. of variables tried at each split: 6
##
Mean of squared residuals: 0.1518349
% Var explained: 79.63
```

当 `mtry` 的值取为缺省值时执行随机森林算法。

对测试集进行预报:

```
pred3 <- predict(rf1, newdata=d[test,])
test.mse3 <- mean((d[test, 'Salary'] - exp(pred3))^2)
test.mse3
```

```
[1] 104602.6
```

结果与剪枝过的单课树、装袋法相近。

在全集上使用随机森林:

```
rf2 <- randomForest(log(Salary) ~ ., data=d, importance=TRUE)
rf2

##
Call:
randomForest(formula = log(Salary) ~ ., data = d, importance = TRUE)
Type of random forest: regression
Number of trees: 500
No. of variables tried at each split: 6
```

```
##
Mean of squared residuals: 0.1787572
% Var explained: 77.31
```

各变量的重要度数值及其图形:

```
importance(rf2)
```

##	%IncMSE	IncNodePurity
## AtBat	9.3330183	7.4403391
## Hits	9.0323007	8.4407205
## HmRun	5.1945601	2.7092496
## Runs	7.9983833	4.3890115
## RBI	6.7481945	6.6695308
## Walks	8.7996281	6.0100606
## Years	12.0076993	7.9765603
## CAtBat	17.7313302	40.6807192
## CHits	15.9835937	32.7540381
## CHmRun	7.1611701	6.5471411
## CRuns	15.5172580	33.4208239
## CRBI	13.7590185	20.3374114
## CWalks	10.2698791	18.2179720
## League	-0.3625290	0.3213873
## Division	1.1810833	0.3230452
## PutOuts	2.0729531	3.3752484
## Assists	-0.7020171	1.8619767
## Errors	1.4288173	1.8079234
## NewLeague	0.8295664	0.3808202

```
varImpPlot(rf2)
```

最重要的自变量是 CAtBats, CRuns, CHits, CWalks, CRBI 等。

rf2

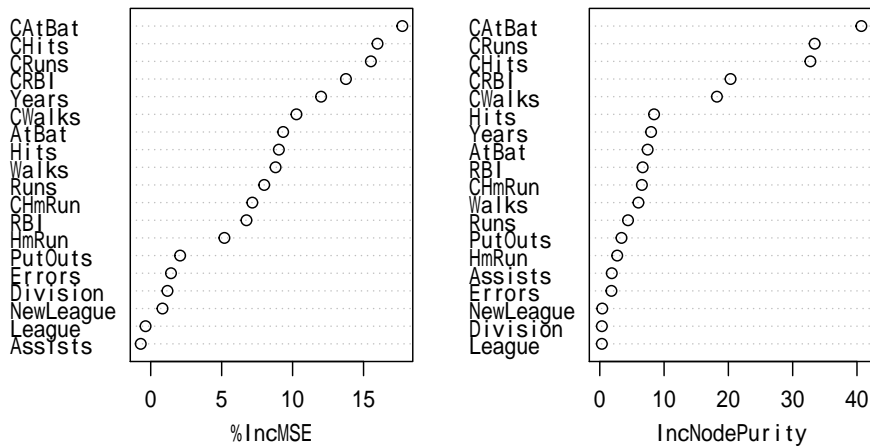


图 34.6: Hitters 数据随机森林法的变量重要度结果

## 34.2 Heart 数据分析

读入 Heart 数据集，并去掉有缺失值的观测：

```
Heart <- read.csv('Heart.csv', header=TRUE, row.names=1)
Heart <- na.omit(Heart)
str(Heart)
```

```
'data.frame': 297 obs. of 14 variables:
$ Age : int 63 67 67 37 41 56 62 57 63 53 ...
$ Sex : int 1 1 1 1 0 1 0 0 1 1 ...
$ ChestPain: Factor w/ 4 levels "asymptomatic",...: 4 1 1 2 3 3 1 1 1 1 ...
$ RestBP : int 145 160 120 130 130 120 140 120 130 140 ...
$ Chol : int 233 286 229 250 204 236 268 354 254 203 ...
$ Fbs : int 1 0 0 0 0 0 0 0 0 1 ...
$ RestECG : int 2 2 2 0 2 0 2 0 2 2 ...
$ MaxHR : int 150 108 129 187 172 178 160 163 147 155 ...
```

```
$ ExAng : int 0 1 1 0 0 0 0 1 0 1 ...
$ Oldpeak : num 2.3 1.5 2.6 3.5 1.4 0.8 3.6 0.6 1.4 3.1 ...
$ Slope : int 3 2 2 3 1 1 3 1 2 3 ...
$ Ca : int 0 3 2 0 0 0 2 0 1 0 ...
$ Thal : Factor w/ 3 levels "fixed","normal",...: 1 2 3 2 2 2 2 2 3 3 ...
$ AHD : Factor w/ 2 levels "No","Yes": 1 2 2 1 1 1 2 1 2 2 ...
- attr(*, "na.action")= 'omit' Named int 88 167 193 267 288 303
..- attr(*, "names")= chr "88" "167" "193" "267" ...
```

```
t(summary(Heart))
```

```
##
Age Min. :29.00 1st Qu.:48.00 Median :56.00 Mean
Sex Min. :0.0000 1st Qu.:0.0000 Median :1.0000 Mean
ChestPain asymptomatic:142 nonanginal : 83 nontypical : 49 typic
RestBP Min. : 94.0 1st Qu.:120.0 Median :130.0 Mean
Chol Min. :126.0 1st Qu.:211.0 Median :243.0 Mean
Fbs Min. :0.0000 1st Qu.:0.0000 Median :0.0000 Mean
RestECG Min. :0.0000 1st Qu.:0.0000 Median :1.0000 Mean
MaxHR Min. : 71.0 1st Qu.:133.0 Median :153.0 Mean
ExAng Min. :0.0000 1st Qu.:0.0000 Median :0.0000 Mean
Oldpeak Min. :0.000 1st Qu.:0.000 Median :0.800 Mean
Slope Min. :1.000 1st Qu.:1.000 Median :2.000 Mean
Ca Min. :0.0000 1st Qu.:0.0000 Median :0.0000 Mean
Thal fixed : 18 normal :164 reversable:115
AHD No :160 Yes:137
```

数据下载: Heart.csv

## 34.2.1 树回归

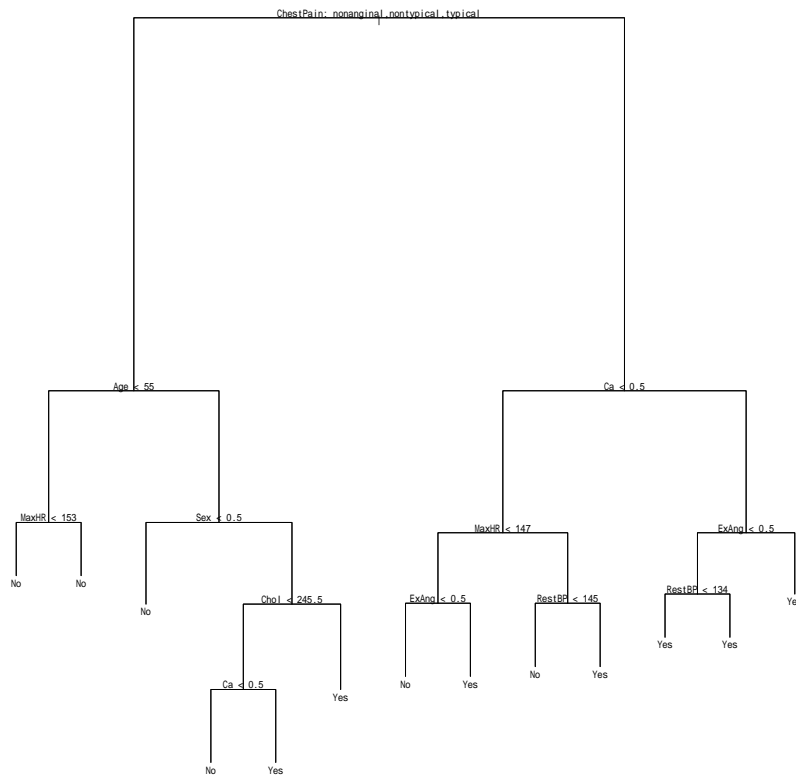
### 34.2.1.1 划分训练集与测试集

简单地把观测分为一半训练集、一半测试集:

```
set.seed(1)
train <- sample(c(TRUE, FALSE), size=nrow(Heart), replace=TRUE)
test <- !train
test.y <- Heart[test, 'AHD']
```

在训练集上建立未剪枝的判别树:

```
tr1 <- tree(AHD ~ ., data=Heart, subset=train)
plot(tr1); text(tr1, pretty=0)
```



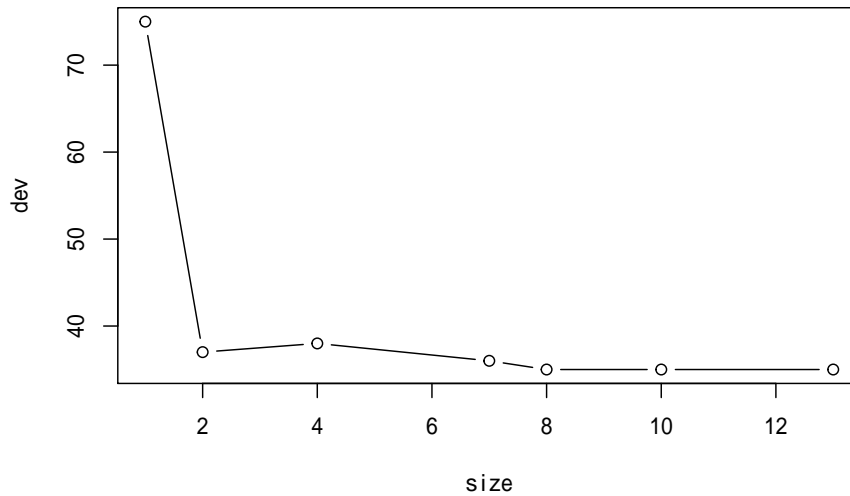
### 34.2.1.2 适当剪枝

用交叉验证方法确定剪枝保留的叶子个数，剪枝时按照错判率执行：

```
cv1 <- cv.tree(tr1, FUN=prune.misclass)
cv1
```

```
$size
[1] 13 10 8 7 4 2 1
##
$dev
[1] 35 35 35 36 38 37 75
##
$k
[1] -Inf 0.0 1.0 2.0 3.0 3.5 37.0
##
$method
[1] "misclass"
##
attr(,"class")
[1] "prune" "tree.sequence"
```

```
plot(cv1$size, cv1$dev, type='b', xlab='size', ylab='dev')
```



```
best.size <- cv1$size[which.min(cv1$dev)]
```

最优的大小是 13。但是从图上看，4 个叶结点已经足够好，所以取为 4。

对训练集生成剪枝结果：

```
best.size <- 4
tr1b <- prune.misclass(tr1, best=best.size)
plot(tr1b); text(tr1b, pretty=0)
```

注意剪枝后树的显示中，内部节点的自变量存在分类变量，这时按照这个自变量分叉时，取指定的某几个分类值时对应分支 Yes，取其它的分分类值时对应分支 No。

### 34.2.1.3 对测试集计算误判率

```
pred1 <- predict(tr1b, Heart[test,], type='class')
tab1 <- table(pred1, test.y); tab1

test.y
```

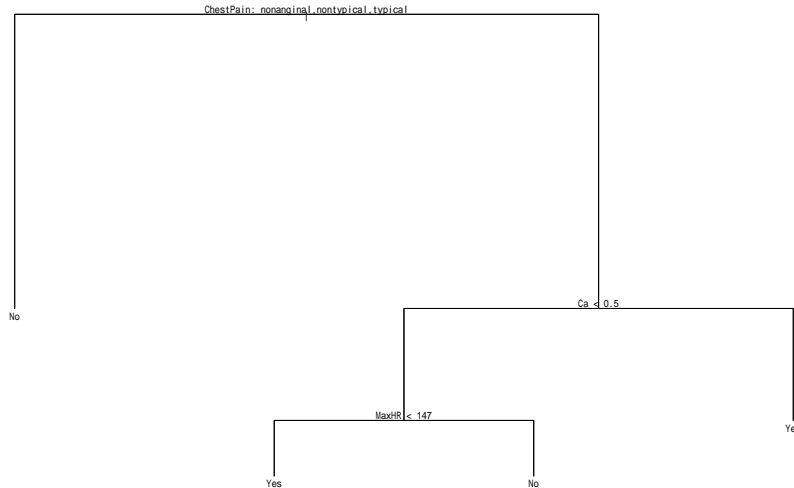


图 34.7: Heart 数据回归树

```
pred1 No Yes
No 78 23
Yes 8 43

test.err <- (tab1[1,2]+tab1[2,1])/sum(tab1[]); test.err

[1] 0.2039474
```

对测试集的错判率约 20%。

利用未剪枝的树对测试集进行预测, 一般比剪枝后的结果差:

```
pred1a <- predict(tr1, Heart[test,], type='class')
tab1a <- table(pred1a, test.y); tab1a
```

```
test.y
pred1a No Yes
No 79 18
Yes 7 48
```

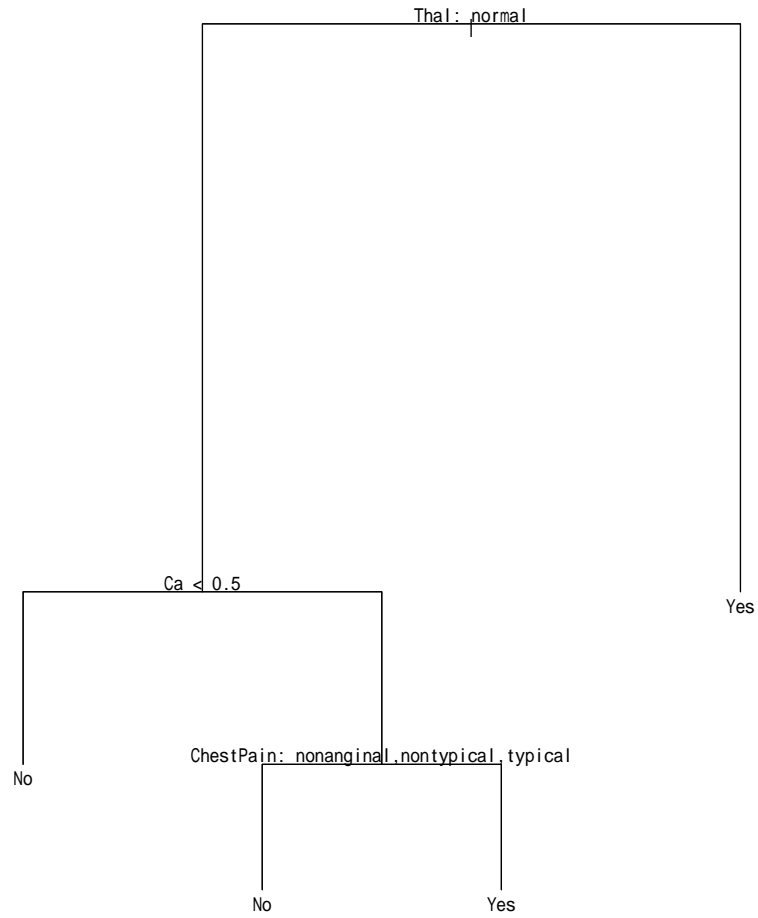


```
test.err1a <- (tab1a[1,2]+tab1a[2,1])/sum(tab1a[]); test.err1a
```

```
[1] 0.1644737
```

#### 34.2.1.4 利用全集数据建立剪枝判别树

```
tr2 <- tree(AHD ~ ., data=Heart)
tr2b <- prune.misclass(tr2, best=best.size)
plot(tr2b); text(tr2b, pretty=0)
```



### 34.2.2 用装袋法

对训练集用装袋法:

```
bag1 <- randomForest(AHD ~ ., data=Heart, subset=train, mtry=13, importance=TRUE)
bag1
```

```
##
Call:
randomForest(formula = AHD ~ ., data = Heart, mtry = 13, importance = TRUE, subset
Type of random forest: classification
Number of trees: 500
No. of variables tried at each split: 13
##
OOB estimate of error rate: 20%
Confusion matrix:
No Yes class.error
No 61 13 0.1756757
Yes 16 55 0.2253521
```

注意 `randomForest()` 函数实际是随机森林法，但是当 `mtry` 的值取为所有自变量个数时就是装袋法。袋外观测得到的错判率比较差。

对测试集进行预报:

```
pred2 <- predict(bag1, newdata=Heart[test,])
tab2 <- table(pred2, test.y); tab2
```

```
test.y
pred2 No Yes
No 76 20
Yes 10 46
```

```
test.err2 <- (tab2[1,2]+tab2[2,1])/sum(tab2[]); test.err2
```

```
[1] 0.1973684
```

测试集的错判率约为 20%。

对全集用装袋法:

```
bag1b <- randomForest(AHD ~ ., data=Heart, mtry=13, importance=TRUE)
bag1b
```

```
##
Call:
```

```
randomForest(formula = AHD ~ ., data = Heart, mtry = 13, importance = TRUE)
Type of random forest: classification
Number of trees: 500
No. of variables tried at each split: 13
##
OOB estimate of error rate: 19.87%
Confusion matrix:
No Yes class.error
No 134 26 0.1625000
Yes 33 104 0.2408759
```

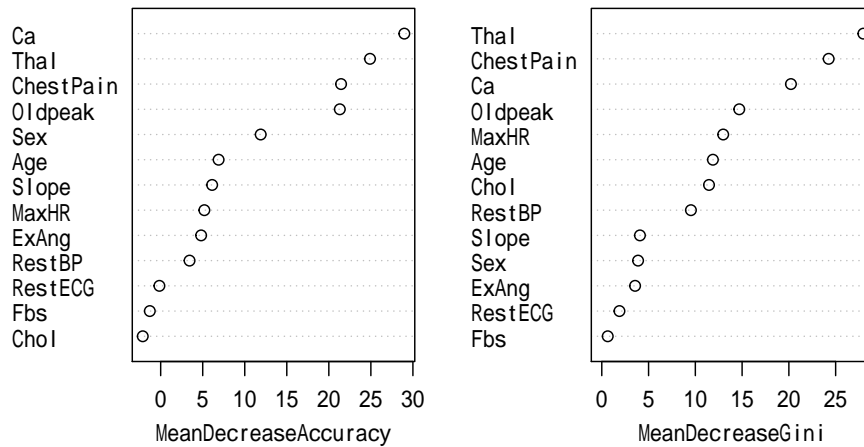
各变量的重要度数值及其图形:

```
importance(bag1b)
```

##	No	Yes	MeanDecreaseAccuracy	MeanDecreaseGini
## Age	5.01743888	4.54672209	6.8901260	11.8781808
## Sex	10.62273160	5.93775833	11.8960286	3.8780605
## ChestPain	12.12215333	18.01109966	21.4561091	24.2519995
## RestBP	2.14912897	2.57544213	3.4350095	9.5279960
## Chol	-0.06740331	-3.34866905	-2.1319263	11.4615191
## Fbs	-0.19159212	-1.68286268	-1.2936067	0.6356754
## RestECG	-0.83204215	0.61823420	-0.1523007	1.8873928
## MaxHR	6.95927941	-0.04284863	5.1956114	12.9777185
## ExAng	2.15397190	4.67398101	4.8112183	3.5598518
## Oldpeak	16.87683151	14.03391494	21.3056979	14.6968413
## Slope	2.97366941	5.17314908	6.1189587	4.0774248
## Ca	25.14524607	18.28904993	28.9768557	20.2200459
## Thal	18.69773426	17.80713339	24.8961726	27.9615428

```
varImpPlot(bag1b)
```

bag1b



最重要的变量是 Thal, ChestPain, Ca。

### 34.2.3 用随机森林

对训练集用随机森林法:

```
rf1 <- randomForest(AHD ~ ., data=Heart, subset=train, importance=TRUE)
rf1

##
Call:
randomForest(formula = AHD ~ ., data = Heart, importance = TRUE, subset = train)
Type of random forest: classification
Number of trees: 500
No. of variables tried at each split: 3
##
OOB estimate of error rate: 20%
Confusion matrix:
No Yes class.error
```

```
No 61 13 0.1756757
Yes 16 55 0.2253521
```

这里 `mtry` 取缺省值，对应于随机森林法。

对测试集进行预报：

```
pred3 <- predict(rf1, newdata=Heart[test,])
tab3 <- table(pred3, test.y); tab3
```

```
test.y
pred3 No Yes
No 78 20
Yes 8 46
```

```
test.err3 <- (tab3[1,2]+tab3[2,1])/sum(tab3[]); test.err3
```

```
[1] 0.1842105
```

测试集的错判率约为 18%。

对全集用随机森林：

```
rf1b <- randomForest(AHD ~ ., data=Heart, importance=TRUE)
rf1b
```

```
##
Call:
randomForest(formula = AHD ~ ., data = Heart, importance = TRUE)
Type of random forest: classification
Number of trees: 500
No. of variables tried at each split: 3
##
OOB estimate of error rate: 17.85%
Confusion matrix:
No Yes class.error
No 138 22 0.1375000
Yes 31 106 0.2262774
```

各变量的重要度数值及其图形：

```
importance(rf1b)
```

##		No	Yes	MeanDecreaseAccuracy	MeanDecreaseGini
## Age		6.647195	5.9142590	8.8592250	12.961695
## Sex		12.271439	6.6415765	14.3610223	4.523229
## ChestPain		11.068166	17.1308734	18.9981168	18.522923
## RestBP		1.965647	0.2372161	1.7415606	10.555002
## Chol		1.809381	-1.9169726	0.1992495	11.260342
## Fbs		1.554259	-2.3732219	-0.4023229	1.368904
## RestECG		1.160775	2.1910276	2.1228412	2.748856
## MaxHR		10.136446	6.7175368	12.0958262	17.994290
## ExAng		2.163296	8.7559583	8.0681465	7.536690
## Oldpeak		12.115693	11.8649191	16.9461505	15.307598
## Slope		3.082837	8.6440416	8.5035252	6.512301
## Ca		21.461418	18.1222940	25.5933528	17.283850
## Thal		19.599715	16.7413930	24.2539786	18.925007

```
varImpPlot(rf1b)
```

最重要的变量是 ChestPain, Thal, Ca。

### 34.3 汽车销量数据分析

Carseats 是 ISLR 包的一个数据集，基本情况如下：

```
str(Carseats)
summary(Carseats)
```

把 Sales 变量按照大于 8 与否分成两组，结果存入变量 High，以 High 为因变量作判别分析。

```
d <- na.omit(Carseats)
d$High <- factor(ifelse(d$Sales > 8, 'Yes', 'No'))
dim(d)
```

```
[1] 400 12
```

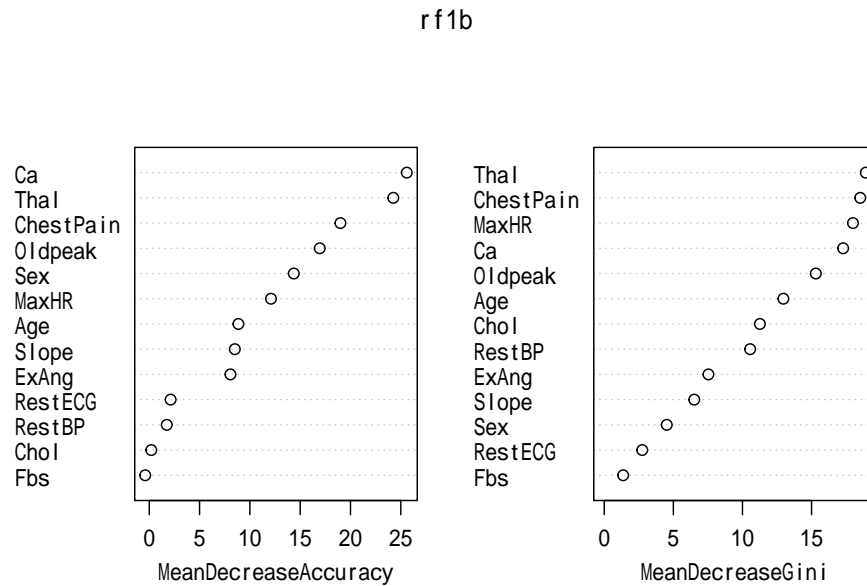


图 34.8: Heart 数据随机森林方法得到的变量重要度

### 34.3.1 判别树

#### 34.3.1.1 全体数据的判别树

对全体数据建立未剪枝的判别树:

```
tr1 <- tree(High ~ . - Sales, data=d)
summary(tr1)
```

```
##
```

```
Classification tree:
```

```
tree(formula = High ~ . - Sales, data = d)
```

```
Variables actually used in tree construction:
```

```
[1] "ShelveLoc" "Price" "Income" "CompPrice" "Population" "Adve
```

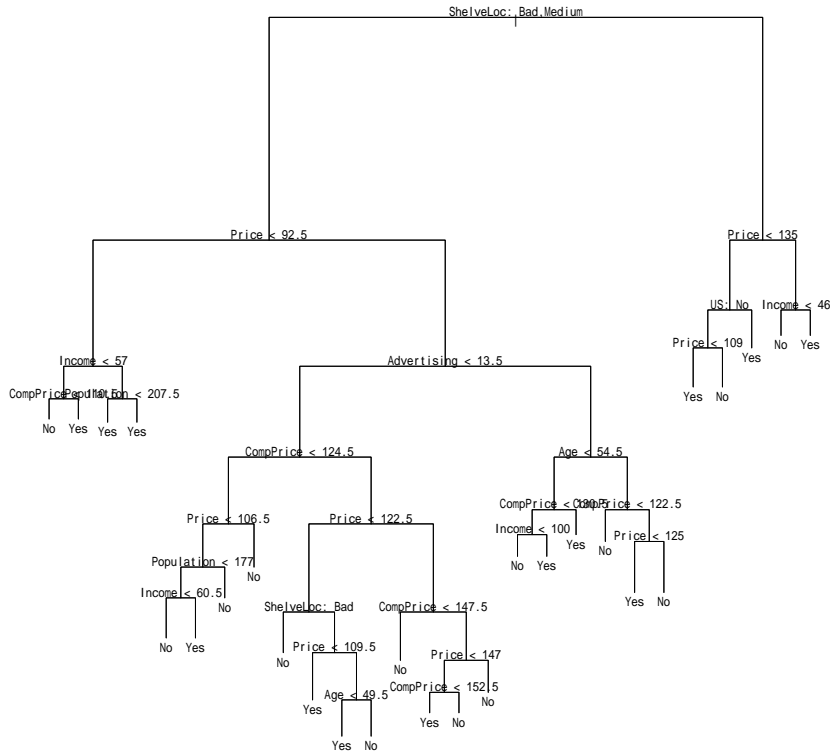
```
Number of terminal nodes: 27
```

```
Residual mean deviance: 0.4575 = 170.7 / 373
```

```
Misclassification error rate: 0.09 = 36 / 400
```



```
plot(tr1)
text(tr1, pretty=0)
```



### 34.3.1.2 划分训练集和测试集

把输入数据集随机地分一半当作训练集，另一半当作测试集：

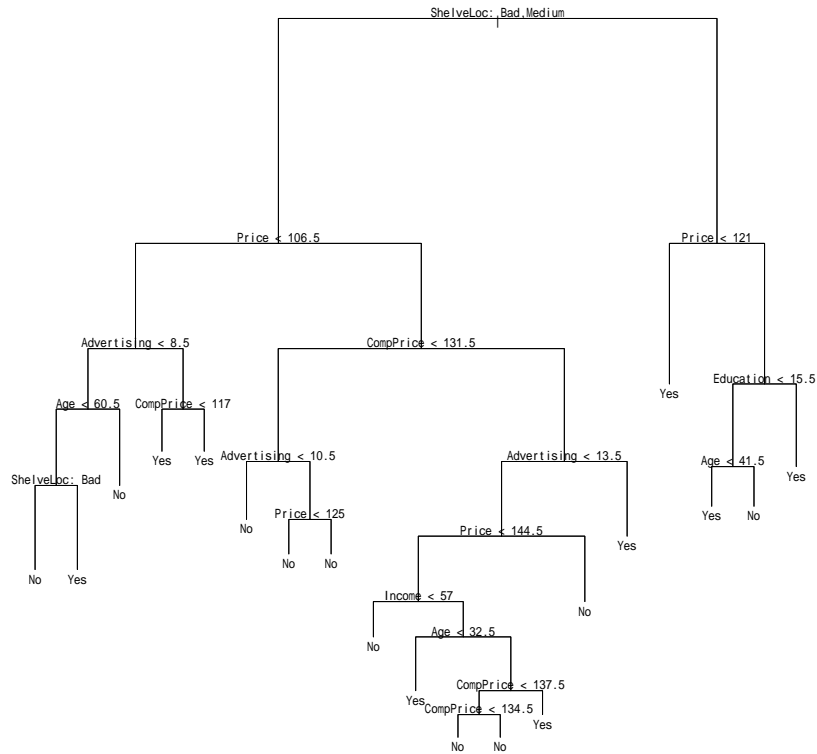
```
set.seed(2)
train <- sample(c(TRUE, FALSE), size=nrow(d), replace=TRUE)
test <- (!train)
test.high <- d[test, 'High']
```

用训练数据建立未剪枝的判别树:

```
tr2 <- tree(High ~ . - Sales, data=d, subset=train)
summary(tr2)
```

```
##
Classification tree:
tree(formula = High ~ . - Sales, data = d, subset = train)
Variables actually used in tree construction:
[1] "ShelveLoc" "Price" "Advertising" "Age" "CompPrice" "Inco
Number of terminal nodes: 19
Residual mean deviance: 0.4299 = 76.96 / 179
Misclassification error rate: 0.09596 = 19 / 198

plot(tr2)
text(tr2, pretty=0)
```



用未剪枝的树对测试集进行预测，并计算误判率：

```
pred2 <- predict(tr2, d[test,], type='class')
tab <- table(pred2, test.high); tab
```

```
test.high
pred2 No Yes
No 93 28
Yes 30 51
```

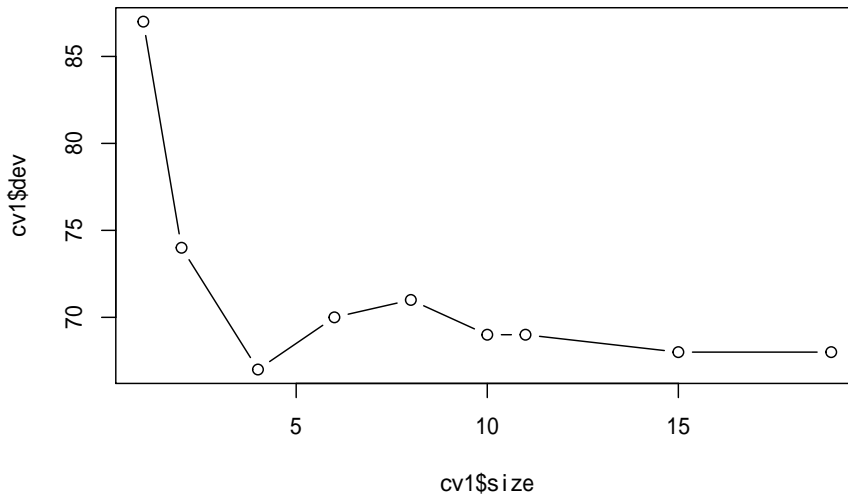
```
test.err2 <- (tab[1,2] + tab[2,1]) / sum(tab[]); test.err2
```

```
[1] 0.2871287
```

### 34.3.1.3 用交叉验证确定训练集的剪枝

```
set.seed(3)
cv1 <- cv.tree(tr2, FUN=prune.misclass)
cv1
```

```
$size
[1] 19 15 11 10 8 6 4 2 1
##
$dev
[1] 68 68 69 69 71 70 67 74 87
##
$k
[1] -Inf 0.00 1.75 2.00 2.50 3.00 4.00 5.50 27.00
##
$method
[1] "misclass"
##
attr(,"class")
[1] "prune" "tree.sequence"
plot(cv1$size, cv1$dev, type='b')
```



```
best.size <- cv1$size[which.min(cv1$dev)]
```

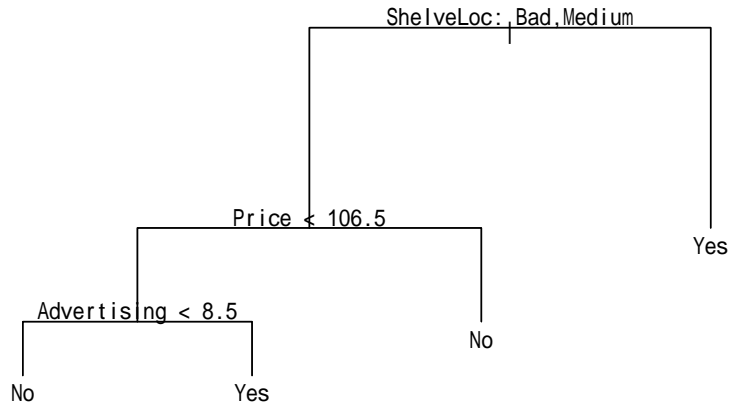
用交叉验证方法自动选择的最佳树大小为 4。

剪枝:

```
tr3 <- prune.misclass(tr2, best=best.size)
summary(tr3)
```

```
##
Classification tree:
snip.tree(tree = tr2, nodes = c(9L, 5L, 3L, 8L))
Variables actually used in tree construction:
[1] "ShelveLoc" "Price" "Advertising"
Number of terminal nodes: 4
Residual mean deviance: 1.103 = 214 / 194
Misclassification error rate: 0.2374 = 47 / 198
```

```
plot(tr3)
text(tr3, pretty=0)
```



用剪枝后的树对测试集进行预测，计算误判率：

```
pred3 <- predict(tr3, d[test,], type='class')
tab <- table(pred3, test.high); tab
```

```
test.high
pred3 No Yes
No 106 33
Yes 17 46
```

```
test.err3 <- (tab[1,2] + tab[2,1]) / sum(tab[]); test.err3
```

```
[1] 0.2475248
```

### 34.3.2 随机森林

对训练集用随机森林法：

```
rf4 <- randomForest(High ~ . - Sales, data=d, subset=train, importance=TRUE)
rf4
```

```
##
Call:
randomForest(formula = High ~ . - Sales, data = d, importance = TRUE, subset = trai
Type of random forest: classification
Number of trees: 500
No. of variables tried at each split: 3
##
OOB estimate of error rate: 22.22%
Confusion matrix:
No Yes class.error
No 94 19 0.1681416
Yes 25 60 0.2941176
```

这里 `mtry` 取缺省值，对应于随机森林法。

对测试集进行预报：

```
pred4 <- predict(rf4, newdata=d[test,])
tab <- table(pred4, test.high); tab
```

```
test.high
pred4 No Yes
No 113 26
Yes 10 53
```

```
test.err4 <- (tab[1,2]+tab[2,1])/sum(tab[,]); test.err4
```

```
[1] 0.1782178
```

注意错判率结果依赖于训练集和测试集的划分，另行选择训练集与测试集可能会得到很不一样的错判率结果。

对全集用随机森林：

```
rf5 <- randomForest(High ~ . - Sales, data=d, importance=TRUE)
rf5
```

```
##
Call:
```

```
randomForest(formula = High ~ . - Sales, data = d, importance = TRUE)
Type of random forest: classification
Number of trees: 500
No. of variables tried at each split: 3
##
OOB estimate of error rate: 18%
Confusion matrix:
No Yes class.error
No 214 22 0.09322034
Yes 50 114 0.30487805
```

各变量的重要度数值及其图形:

```
importance(rf5)
```

	No	Yes	MeanDecreaseAccuracy	MeanDecreaseGini
## CompPrice	9.853176	7.84106962	12.3403571	21.895753
## Income	3.901345	5.95859870	6.6820994	19.708699
## Advertising	10.859513	16.36240814	19.1737947	22.799306
## Population	-1.788738	-3.80236686	-4.1421409	15.366835
## Price	31.040118	27.34118336	37.7942812	43.588084
## ShelfLoc	31.493614	34.00428243	40.3330521	30.992966
## Age	9.033102	9.70947568	12.2385591	22.352854
## Education	1.829310	-0.01953518	1.4315544	10.190089
## Urban	0.193837	-1.09027048	-0.4997192	2.243931
## US	2.182323	6.29992327	6.0770295	3.542851

```
varImpPlot(rf5)
```

重要的自变量为 Price, ShelfLoc, 其次有 Age, Advertising, CompPrice, Income 等。



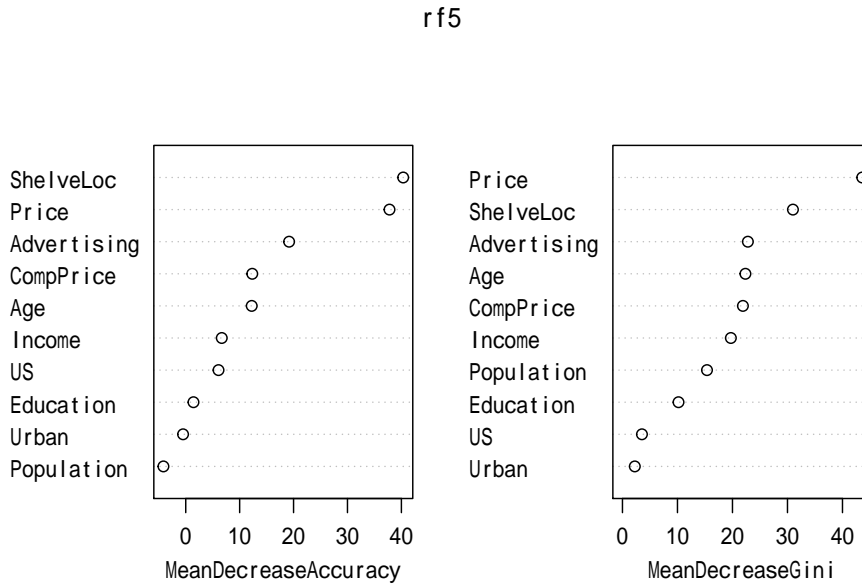


图 34.9: Carseats 数据随机森林法得到的变量重要度

## 34.4 波士顿郊区房价数据

MASS 包的 Boston 数据包含了波士顿地区郊区房价的若干数据。以中位房价 medv 为因变量建立回归模型。首先把缺失值去掉后存入数据集 d:

```
d <- na.omit(Boston)
```

数据集概况:

```
str(d)
```

```
'data.frame': 506 obs. of 14 variables:
$ crim : num 0.00632 0.02731 0.02729 0.03237 0.06905 ...
$ zn : num 18 0 0 0 0 0 12.5 12.5 12.5 12.5 ...
$ indus : num 2.31 7.07 7.07 2.18 2.18 2.18 7.87 7.87 7.87 7.87 ...
$ chas : int 0 0 0 0 0 0 0 0 0 0 ...
$ nox : num 0.538 0.469 0.469 0.458 0.458 0.458 0.524 0.524 0.524 0.524 ...
$ rm : num 6.58 6.42 7.18 7 7.15 ...
```

```
$ age : num 65.2 78.9 61.1 45.8 54.2 58.7 66.6 96.1 100 85.9 ...
$ dis : num 4.09 4.97 4.97 6.06 6.06 ...
$ rad : int 1 2 2 3 3 3 5 5 5 5 ...
$ tax : num 296 242 242 222 222 222 311 311 311 311 ...
$ ptratio: num 15.3 17.8 17.8 18.7 18.7 18.7 15.2 15.2 15.2 15.2 ...
$ black : num 397 397 393 395 397 ...
$ lstat : num 4.98 9.14 4.03 2.94 5.33 ...
$ medv : num 24 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 ...
```

```
summary(d)
```

```
crim zn indus chas nox
Min. : 0.00632 Min. : 0.00 Min. : 0.46 Min. :0.00000 Min. :
1st Qu.: 0.08204 1st Qu.: 0.00 1st Qu.: 5.19 1st Qu.:0.00000 1st Qu.:
Median : 0.25651 Median : 0.00 Median : 9.69 Median :0.00000 Median :
Mean : 3.61352 Mean : 11.36 Mean :11.14 Mean :0.06917 Mean :
3rd Qu.: 3.67708 3rd Qu.: 12.50 3rd Qu.:18.10 3rd Qu.:0.00000 3rd Qu.:
Max. :88.97620 Max. :100.00 Max. :27.74 Max. :1.00000 Max. :
```

### 34.4.1 回归树

#### 34.4.1.1 划分训练集和测试集

```
set.seed(1)
train <- sample(c(TRUE, FALSE), nrow(d), replace=TRUE)
test <- (!train)
```

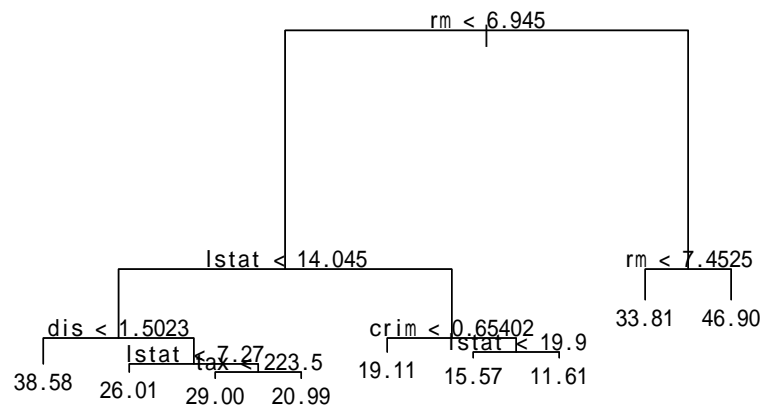
对训练集建立未剪枝的树:

```
tr1 <- tree(medv ~ ., d, subset=train)
summary(tr1)
```

```
##
Regression tree:
tree(formula = medv ~ ., data = d, subset = train)
Variables actually used in tree construction:
```

```
[1] "rm" "lstat" "dis" "tax" "crim"
Number of terminal nodes: 9
Residual mean deviance: 15.65 = 3960 / 253
Distribution of residuals:
Min. 1st Qu. Median Mean 3rd Qu. Max.
-23.5800 -2.0860 -0.1802 0.0000 2.0030 16.1900
```

```
plot(tr1)
text(tr1, pretty=0)
```



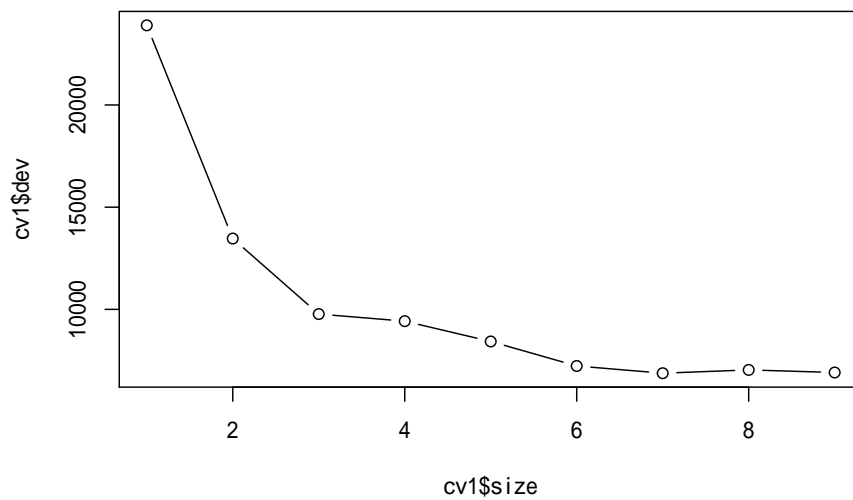
用未剪枝的树对测试集进行预测，计算均方误差：

```
yhat <- predict(tr1, newdata=d[test,])
mse1 <- mean((yhat - d[test, 'medv'])^2)
mse1
```

```
[1] 21.00586
```

## 34.4.1.2 用交叉验证方法确定剪枝复杂度

```
cv1 <- cv.tree(tr1)
plot(cv1$size, cv1$dev, type='b')
```

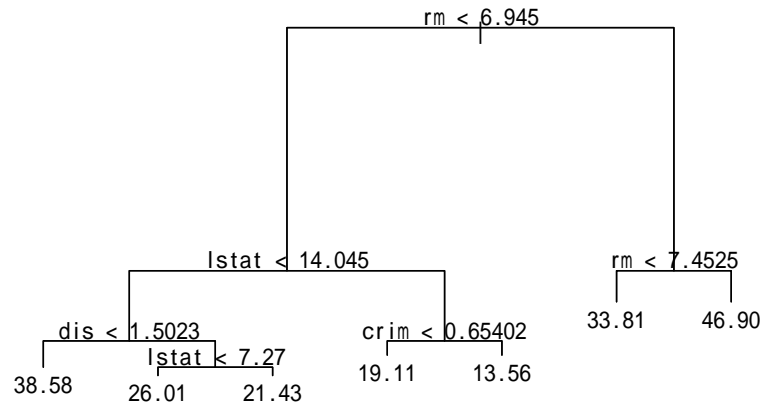


```
best.size <- cv1$size[which.min(cv1$dev)]; best.size
```

```
[1] 7
```

剪枝并对测试集进行预测:

```
tr2 <- prune.tree(tr1, best=best.size)
plot(tr2)
text(tr2, pretty=0)
```



```

yhat <- predict(tr2, newdata=d[test,])
mse2 <- mean((yhat - d[test, 'medv'])^2)
mse2

```

```
[1] 21.81789
```

剪枝后效果没有改善。

### 34.4.2 装袋法

用 `randomForest` 包计算。当参数 `mtry` 取为自变量个数时按照装袋法计算。对训练集计算。

```

set.seed(1)
bag1 <- randomForest(medv ~ ., data=d, subset=train, mtry=ncol(d)-1, importance=TRUE)
bag1

```

```
##
```

```
Call:
```

```
randomForest(formula = medv ~ ., data = d, mtry = ncol(d) - 1, importance = TRUE, s
```

```
Type of random forest: regression
Number of trees: 500
No. of variables tried at each split: 13
##
Mean of squared residuals: 14.62942
% Var explained: 83.56
```

在测试集上计算装袋法的均方误差:

```
yhat <- predict(bag1, newdata=d[test,])
mean((yhat - d[test, 'medv'])^2)
```

```
[1] 13.10016
```

比单棵树的结果有明显改善。

### 34.4.3 随机森林

用 randomForest 包计算。当参数 mtry 取为缺省值时按照随机森林方法计算。对训练集计算。

```
set.seed(1)
rf1 <- randomForest(medv ~ ., data=d, subset=train, importance=TRUE)
rf1
```

```
##
Call:
randomForest(formula = medv ~ ., data = d, importance = TRUE, subset = tr
Type of random forest: regression
Number of trees: 500
No. of variables tried at each split: 4
##
Mean of squared residuals: 13.78256
% Var explained: 84.52
```

在测试集上计算随机森林法的均方误差:

```
yhat <- predict(rf1, newdata=d[test,])
mean((yhat - d[test, 'medv'])^2)
```

```
[1] 10.75832
```

比单棵树的结果有明显改善, 但是不如装袋法的结果。

各变量的重要度数值及其图形:

```
importance(rf1)
```

```
%IncMSE IncNodePurity
crim 13.310976 1555.9199
zn 3.817851 196.5880
indus 11.156318 1583.5080
chas 2.710213 101.7295
nox 14.318933 1313.2498
rm 25.654938 6422.5521
age 9.119223 631.8902
dis 13.068722 1644.8226
rad 5.755673 206.6290
tax 10.895103 595.7777
ptratio 14.158137 1176.3553
black 9.203933 523.5443
lstat 29.054691 6839.5325
```

```
varImpPlot(rf1)
```

#### 34.4.4 提升法

使用 gbm 包。在训练集上拟合:

```
set.seed(1)
bst1 <- gbm(medv ~ ., data=d[train,], distribution='gaussian', n.trees=5000, interaction.
summary(bst1)
```

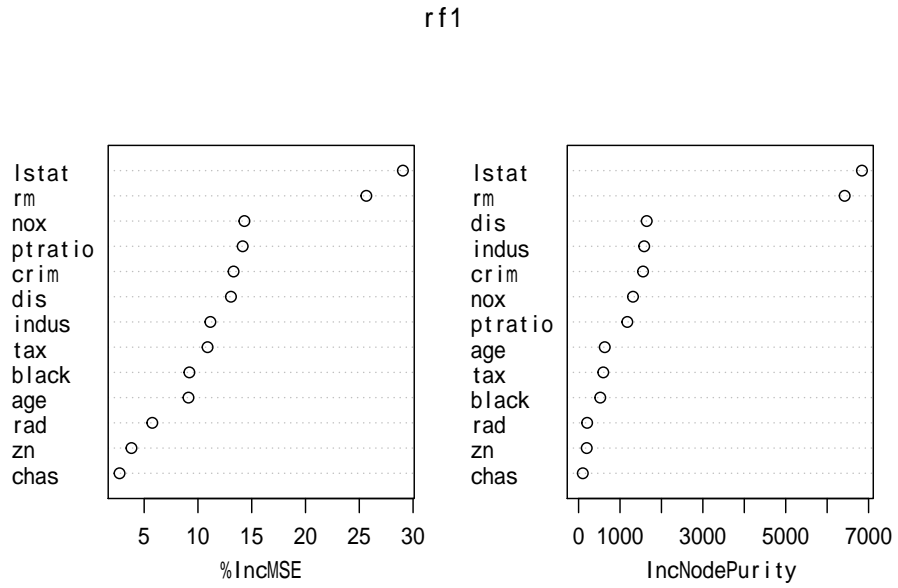
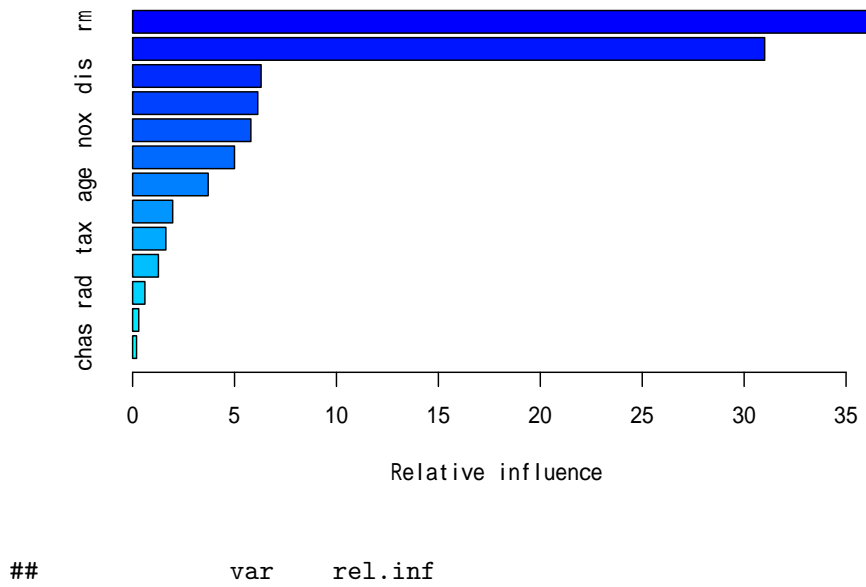


图 34.10: Boston 数据用随机森林法得到的变量重要度





```
rm rm 36.0855527
lstat lstat 31.0069766
dis dis 6.3021052
black black 6.1402767
nox nox 5.8010558
crim crim 4.9981102
age age 3.7062673
ptratio ptratio 1.9659316
tax tax 1.6314237
indus indus 1.2651374
rad rad 0.6021289
zn zn 0.3016741
chas chas 0.1933597
```

lstat 和 rm 是最重要的变量。

在测试集上预报，并计算均方误差：

```
yhat <- predict(bst1, newdata=d[test,], n.trees=5000)
mean((yhat - d[test, 'medv'])^2)
```

```
[1] 16.04224
```

与随机森林方法结果相近。

如果提高学习速度：

```
bst2 <- gbm(medv ~ ., data=d[train,], distribution='gaussian', n.trees=5000, interaction.
yhat <- predict(bst2, newdata=d[test,], n.trees=5000)
mean((yhat - d[test, 'medv'])^2)
```

```
[1] 14.87207
```

均方误差有改善。

## 34.5 附录

### 34.5.1 Hitters 数据

```
knitr::kable(Hitters)
```

	AtBat	Hits	HmRun	Runs	RBI	Walks	Years	CAtBat	CHits	CHmRun
-Andy Allanson	293	66	1	30	29	14	1	293	66	1
-Alan Ashby	315	81	7	24	38	39	14	3449	835	69
-Alvin Davis	479	130	18	66	72	76	3	1624	457	63
-Andre Dawson	496	141	20	65	78	37	11	5628	1575	225
-Andres Galarraga	321	87	10	39	42	30	2	396	101	12
-Alfredo Griffin	594	169	4	74	51	35	11	4408	1133	19
-Al Newman	185	37	1	23	8	21	2	214	42	1
-Argenis Salazar	298	73	0	24	24	7	3	509	108	0
-Andres Thomas	323	81	6	26	32	8	2	341	86	6
-Andre Thornton	401	92	17	49	66	65	13	5206	1332	253
-Alan Trammell	574	159	21	107	75	59	10	4631	1300	90
-Alex Trevino	202	53	4	31	26	27	9	1876	467	15
-Andy VanSlyke	418	113	13	48	61	47	4	1512	392	41
-Alan Wiggins	239	60	0	30	11	22	6	1941	510	4
-Bill Almon	196	43	7	29	27	30	13	3231	825	36
-Billy Beane	183	39	3	20	15	11	3	201	42	3
-Buddy Bell	568	158	20	89	75	73	15	8068	2273	177
-Buddy Biancalana	190	46	2	24	8	15	5	479	102	5
-Bruce Bochte	407	104	6	57	43	65	12	5233	1478	100
-Bruce Bochy	127	32	8	16	22	14	8	727	180	24
-Barry Bonds	413	92	16	72	48	65	1	413	92	16
-Bobby Bonilla	426	109	3	55	43	62	1	426	109	3
-Bob Boone	22	10	1	4	2	1	6	84	26	2
-Bob Brenly	472	116	16	60	62	74	6	1924	489	67
-Bill Buckner	629	168	18	73	102	40	18	8424	2464	164
-Brett Butler	587	163	4	92	51	70	6	2695	747	17
-Bob Dernier	324	73	4	32	18	22	7	1931	491	13
-Bo Diaz	474	129	10	50	56	40	10	2331	604	61
-Bill Doran	550	152	6	92	37	81	5	2308	633	32
-Brian Downing	513	137	20	90	95	90	14	5201	1382	166
-Bobby Grich	313	84	9	42	30	39	17	6890	1833	224
-Billy Hatcher	419	108	6	55	36	22	3	591	149	8
-Bob Horner	517	141	27	70	87	52	9	3571	994	215
-Brook Jacoby	583	168	17	83	80	56	5	1646	452	44
-Bob Kearney	204	49	6	23	25	12	7	1309	308	27
-Bill Madlock	379	106	10	38	60	30	14	6207	1906	146
-Bobby Meacham	161	36	0	19	10	17	4	1053	244	3
-Bob Melvin	268	60	5	24	25	15	2	350	78	5
-Ben Oglivie	346	98	5	31	53	30	16	5913	1615	235

### 34.5.2 Heart 数据

```
knitr::kable(Heart)
```

	Age	Sex	ChestPain	RestBP	Chol	Fbs	RestECG	MaxHR	ExAng	Oldpeak	Slope
1	63	1	typical	145	233	1	2	150	0	2.3	3
2	67	1	asymptomatic	160	286	0	2	108	1	1.5	2
3	67	1	asymptomatic	120	229	0	2	129	1	2.6	2
4	37	1	nonanginal	130	250	0	0	187	0	3.5	3
5	41	0	nontypical	130	204	0	2	172	0	1.4	1
6	56	1	nontypical	120	236	0	0	178	0	0.8	1
7	62	0	asymptomatic	140	268	0	2	160	0	3.6	3
8	57	0	asymptomatic	120	354	0	0	163	1	0.6	1
9	63	1	asymptomatic	130	254	0	2	147	0	1.4	2
10	53	1	asymptomatic	140	203	1	2	155	1	3.1	3
11	57	1	asymptomatic	140	192	0	0	148	0	0.4	2
12	56	0	nontypical	140	294	0	2	153	0	1.3	2
13	56	1	nonanginal	130	256	1	2	142	1	0.6	2
14	44	1	nontypical	120	263	0	0	173	0	0.0	1
15	52	1	nonanginal	172	199	1	0	162	0	0.5	1
16	57	1	nonanginal	150	168	0	0	174	0	1.6	1
17	48	1	nontypical	110	229	0	0	168	0	1.0	3
18	54	1	asymptomatic	140	239	0	0	160	0	1.2	1
19	48	0	nonanginal	130	275	0	0	139	0	0.2	1
20	49	1	nontypical	130	266	0	0	171	0	0.6	1
21	64	1	typical	110	211	0	2	144	1	1.8	2
22	58	0	typical	150	283	1	2	162	0	1.0	1
23	58	1	nontypical	120	284	0	2	160	0	1.8	2
24	58	1	nonanginal	132	224	0	2	173	0	3.2	1
25	60	1	asymptomatic	130	206	0	2	132	1	2.4	2
26	50	0	nonanginal	120	219	0	0	158	0	1.6	2
27	58	0	nonanginal	120	340	0	0	172	0	0.0	1
28	66	0	typical	150	226	0	0	114	0	2.6	3
29	43	1	asymptomatic	150	247	0	0	171	0	1.5	1
30	40	1	asymptomatic	110	167	0	2	114	1	2.0	2
31	69	0	typical	140	239	0	0	151	0	1.8	1
32	60	1	asymptomatic	117	230	1	0	160	1	1.4	1
33	64	1	nonanginal	140	335	0	0	158	0	0.0	1
34	59	1	asymptomatic	135	234	0	0	161	0	0.5	2
35	44	1	nonanginal	130	233	0	0	179	1	0.4	1
36	42	1	asymptomatic	140	226	0	0	178	0	0.0	1
37	43	1	asymptomatic	120	177	0	2	120	1	2.5	2
38	57	1	asymptomatic	150	276	0	2	112	1	0.6	2
39	55	1	asymptomatic	132	353	0	0	132	1	1.2	2

### 34.5.3 CarSeats 数据

```
knitr::kable(Carseats)
```

Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban	U
9.50	138	73	11	276	120	Bad	42	17	Yes	Y
11.22	111	48	16	260	83	Good	65	10	Yes	Y
10.06	113	35	10	269	80	Medium	59	12	Yes	Y
7.40	117	100	4	466	97	Medium	55	14	Yes	Y
4.15	141	64	3	340	128	Bad	38	13	Yes	N
10.81	124	113	13	501	72	Bad	78	16	No	Y
6.63	115	105	0	45	108	Medium	71	15	Yes	N
11.85	136	81	15	425	120	Good	67	10	Yes	Y
6.54	132	110	0	108	124	Medium	76	10	No	N
4.69	132	113	0	131	124	Medium	76	17	No	Y
9.01	121	78	9	150	100	Bad	26	10	No	Y
11.96	117	94	4	503	94	Good	50	13	Yes	Y
3.98	122	35	2	393	136	Medium	62	18	Yes	N
10.96	115	28	11	29	86	Good	53	18	Yes	Y
11.17	107	117	11	148	118	Good	52	18	Yes	Y
8.71	149	95	5	400	144	Medium	76	18	No	N
7.58	118	32	0	284	110	Good	63	13	Yes	N
12.29	147	74	13	251	131	Good	52	10	Yes	Y
13.91	110	110	0	408	68	Good	46	17	No	Y
8.73	129	76	16	58	121	Medium	69	12	Yes	Y
6.41	125	90	2	367	131	Medium	35	18	Yes	Y
12.13	134	29	12	239	109	Good	62	18	No	Y
5.08	128	46	6	497	138	Medium	42	13	Yes	N
5.87	121	31	0	292	109	Medium	79	10	Yes	N
10.14	145	119	16	294	113	Bad	42	12	Yes	Y
14.90	139	32	0	176	82	Good	54	11	No	N
8.33	107	115	11	496	131	Good	50	11	No	Y
5.27	98	118	0	19	107	Medium	64	17	Yes	N
2.99	103	74	0	359	97	Bad	55	11	Yes	Y
7.81	104	99	15	226	102	Bad	58	17	Yes	Y
13.55	125	94	0	447	89	Good	30	12	Yes	N
8.25	136	58	16	241	131	Medium	44	18	Yes	Y
6.20	107	32	12	236	137	Good	64	10	No	Y
8.77	114	38	13	317	128	Good	50	16	Yes	Y
2.67	115	54	0	406	128	Medium	42	17	Yes	Y
11.07	131	84	11	29	96	Medium	44	17	No	Y
8.89	122	76	0	270	100	Good	60	18	No	N
4.95	121	41	5	412	110	Medium	54	10	Yes	Y
6.59	109	73	0	454	102	Medium	65	15	Yes	N

### 34.5.4 Boston 数据

```
knitr::kable(Boston)
```



crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	black	lstat	m
0.00632	18.0	2.31	0	0.5380	6.575	65.2	4.0900	1	296	15.3	396.90	4.98	2
0.02731	0.0	7.07	0	0.4690	6.421	78.9	4.9671	2	242	17.8	396.90	9.14	2
0.02729	0.0	7.07	0	0.4690	7.185	61.1	4.9671	2	242	17.8	392.83	4.03	3
0.03237	0.0	2.18	0	0.4580	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	3
0.06905	0.0	2.18	0	0.4580	7.147	54.2	6.0622	3	222	18.7	396.90	5.33	3
0.02985	0.0	2.18	0	0.4580	6.430	58.7	6.0622	3	222	18.7	394.12	5.21	2
0.08829	12.5	7.87	0	0.5240	6.012	66.6	5.5605	5	311	15.2	395.60	12.43	2
0.14455	12.5	7.87	0	0.5240	6.172	96.1	5.9505	5	311	15.2	396.90	19.15	2
0.21124	12.5	7.87	0	0.5240	5.631	100.0	6.0821	5	311	15.2	386.63	29.93	1
0.17004	12.5	7.87	0	0.5240	6.004	85.9	6.5921	5	311	15.2	386.71	17.10	1
0.22489	12.5	7.87	0	0.5240	6.377	94.3	6.3467	5	311	15.2	392.52	20.45	1
0.11747	12.5	7.87	0	0.5240	6.009	82.9	6.2267	5	311	15.2	396.90	13.27	1
0.09378	12.5	7.87	0	0.5240	5.889	39.0	5.4509	5	311	15.2	390.50	15.71	2
0.62976	0.0	8.14	0	0.5380	5.949	61.8	4.7075	4	307	21.0	396.90	8.26	2
0.63796	0.0	8.14	0	0.5380	6.096	84.5	4.4619	4	307	21.0	380.02	10.26	1
0.62739	0.0	8.14	0	0.5380	5.834	56.5	4.4986	4	307	21.0	395.62	8.47	1
1.05393	0.0	8.14	0	0.5380	5.935	29.3	4.4986	4	307	21.0	386.85	6.58	2
0.78420	0.0	8.14	0	0.5380	5.990	81.7	4.2579	4	307	21.0	386.75	14.67	1
0.80271	0.0	8.14	0	0.5380	5.456	36.6	3.7965	4	307	21.0	288.99	11.69	2
0.72580	0.0	8.14	0	0.5380	5.727	69.5	3.7965	4	307	21.0	390.95	11.28	1
1.25179	0.0	8.14	0	0.5380	5.570	98.1	3.7979	4	307	21.0	376.57	21.02	1
0.85204	0.0	8.14	0	0.5380	5.965	89.2	4.0123	4	307	21.0	392.53	13.83	1
1.23247	0.0	8.14	0	0.5380	6.142	91.7	3.9769	4	307	21.0	396.90	18.72	1
0.98843	0.0	8.14	0	0.5380	5.813	100.0	4.0952	4	307	21.0	394.54	19.88	1
0.75026	0.0	8.14	0	0.5380	5.924	94.1	4.3996	4	307	21.0	394.33	16.30	1
0.84054	0.0	8.14	0	0.5380	5.599	85.7	4.4546	4	307	21.0	303.42	16.51	1
0.67191	0.0	8.14	0	0.5380	5.813	90.3	4.6820	4	307	21.0	376.88	14.81	1
0.95577	0.0	8.14	0	0.5380	6.047	88.8	4.4534	4	307	21.0	306.38	17.28	1
0.77299	0.0	8.14	0	0.5380	6.495	94.4	4.4547	4	307	21.0	387.94	12.80	1
1.00245	0.0	8.14	0	0.5380	6.674	87.3	4.2390	4	307	21.0	380.23	11.98	2
1.13081	0.0	8.14	0	0.5380	5.713	94.1	4.2330	4	307	21.0	360.17	22.60	1
1.35472	0.0	8.14	0	0.5380	6.072	100.0	4.1750	4	307	21.0	376.73	13.04	1
1.38799	0.0	8.14	0	0.5380	5.950	82.0	3.9900	4	307	21.0	232.60	27.71	1
1.15172	0.0	8.14	0	0.5380	5.701	95.0	3.7872	4	307	21.0	358.77	18.35	1
1.61282	0.0	8.14	0	0.5380	6.096	96.9	3.7598	4	307	21.0	248.31	20.34	1
0.06417	0.0	5.96	0	0.4990	5.933	68.2	3.3603	5	279	19.2	396.90	9.68	1
0.09744	0.0	5.96	0	0.4990	5.841	61.4	3.3779	5	279	19.2	377.56	11.41	2
0.08014	0.0	5.96	0	0.4990	5.850	41.5	3.9342	5	279	19.2	396.90	8.77	2
0.17505	0.0	5.96	0	0.4990	5.966	30.2	3.8473	5	279	19.2	393.43	10.13	2



## Part VIII

### 专题



## Chapter 35

# R 语言的文本处理

### 35.1 简单的文本处理

在信息爆炸性增长的今天，大量的信息是文本型的，如互联网上的大多数资源。R 具有基本的文本数据处理能力，而且因为 R 的向量语言特点和强大的统计计算和图形功能，用 R 处理文本数据是可行的。

#### 35.1.1 字符型常量与字符型向量

字符串常量写在两个双撇号或者两个单撇号中间，建议仅使用双撇号，因为这是大多数常见程序语言的做法。如果内容中有单撇号或者双撇号，可以在前面加反斜杠`\`。为了在字符串中写一个反斜杠，需要写成两个，比如路径 `C:\work` 写成 R 字符串，要写成 `"C:\\work"`。注意，这些规定都是针对程序中的字符串常量，数据中的文本类型数据是不需要遵照这些规定的。

在用 `print()` 显示字符串变量时，也会按照上述的办法显示，比如字符串内的双撇号会被自动加上前导反斜杠，但保存的实际内容中并没有反斜杠。

字符串中可以有一些特殊字符，如 `"\n"` 表示换行符，`"\t"` 表示制表符，`"\r"` 表示回车符，等等。

R 的字符型向量每个元素是一个字符串，如：

```
s <- c("123", "abc", " 张三李四", "@#$$%^&")
s

[1] "123" "abc" "张三李四" "@#$$%^&"
```

R 中处理文本型数据的函数有文件访问函数以及 `readLines`, `nchar`, `paste`, `sprintf`, `format`, `formatC`, `substring` 等函数。

R 支持正则表达式, 函数 `grep`, `grepl`, `sub`, `gsub`, `regexpr`, `gregexpr`, `strsplit` 与正则表达式有关。

字符型函数一般都是向量化的, 对输入的一个字符型向量的每个元素操作。

R 扩展包 `stringr` 和 `stringi` 提供了更方便、功能更强的字符串功能, 包括正则表达式功能。其中 `stringr` 是常用功能, `stringi` 是更基本、更灵活的功能, 一般使用 `stringr` 就足够了。`stringr` 包的函数名大多都以 `str_` 开头。

下面先介绍常用的较简单的字符串函数, 包括 `stringr` 包的函数与基本 R 函数。

```
library(stringr)
```

### 35.1.2 字符串连接、重复

`stringr::str_c()` 用来把多个输入自变量按照元素对应组合为一个字符型向量, 用 `sep` 指定分隔符, 默认为不分隔。类似于 R 中向量间运算的一般规则, 各自变量长度不同时短的自动循环使用。非字符串类型自动转换为字符型。如

```
str_c(c("x", "y"), c("a", "b"), sep="*")
```

```
[1] "x*a" "y*b"
```

```
str_c("data", 1:3, ".txt")
```

```
[1] "data1.txt" "data2.txt" "data3.txt"
```

字符型缺失值参与连接时, 结果变成缺失值; 可以用 `str_replace_na()` 函数将待连接的字符型向量中的缺失值转换成字符串 `"NA"` 再连接。

用 `collapse` 选项要求将连接后的字符型向量的所有元素连接在一起, `collapse` 的值为将多个元素合并时的分隔符。如

```
str_c(c("a", "bc", "def"), collapse="---")
```

```
[1] "a---bc---def"
```

在使用了 `collapse` 时如果有多个要连接的部分，`str_c()` 函数先将各部分连接成为一个字符型向量，然后再把结果的各个向量元素连接起来。如

```
str_c("data", 1:3, ".txt", sep="", collapse=";")
```

```
[1] "data1.txt;data2.txt;data3.txt"
```

`stringr::str_flatten()` 类似于 `stringr::str_c()` 仅有 `collapse` 参数作用一样，仅将一个字符型向量的各个元素按照 `collapse` 参数指定的分隔符连接成一个长字符串，`collapse` 默认值是空字符串，如：

```
str_flatten(c("a", "bc", "def"), collapse="---")
```

```
[1] "a---bc---def"
```

```
str_flatten(c("a", "bc", "def"))
```

```
[1] "abcdef"
```

基本 R 的 `paste()` 函数与 `stringr::str_c()` 函数有类似的用法，但是参数 `sep` 的默认值是空格。基本 R 的 `paste0()` 函数相当于 `stringr::str_c()` 函数固定 `sep` 参数为空字符串。如：

```
paste(c("x", "y"), c("a", "b"), sep="*")
```

```
[1] "x*a" "y*b"
```

```
paste("data", 1:3, ".txt", sep="")
```

```
[1] "data1.txt" "data2.txt" "data3.txt"
```

```
paste0("data", 1:3, ".txt")
```

```
[1] "data1.txt" "data2.txt" "data3.txt"
```

```
paste(c("a", "bc", "def"), collapse="---")
```

```
[1] "a---bc---def"
```

```
paste("data", 1:3, ".txt", sep="", collapse=";")
```

```
[1] "data1.txt;data2.txt;data3.txt"
```

`stringr::str_dup(string, times)` 类似于 `rep()` 函数，可以将字符型向量的元素按照 `times` 指定的次数在同一字符串内重复，如：

```
str_dup(c("abc", "长江"), 3)
```

```
[1] "abcabcabc" "长江长江长江"
```

也可以针对每个元素指定不同重复次数，如

```
str_dup(c("abc", "长江"), c(3, 2))
```

```
[1] "abcabcabc" "长江长江"
```

### 35.1.3 格式化输出

#### 35.1.3.1 `format()` 函数

`format()` 函数可以将一个数值型向量的各个元素按照统一格式转换为字符型，如：

```
as.character(1.000)
```

```
[1] "1"
```

```
as.character(1.2)
```

```
[1] "1.2"
```

```
as.character(1.23)
```

```
[1] "1.23"
```

```
format(c(1.000, 1.2, 1.23))
```

```
[1] "1.00" "1.20" "1.23"
```

选项 `digits` 与 `nsmall` 共同控制输出的精度，`nsmall` 控制非科学记数法显示时小数点后的至少要有的位数，`digits` 控制至少要有的有效位数。这使得输



出的宽度是不可控的，如：

```
format(c(pi, pi*10000), digits=8, nsmall=4)
```

```
[1] " 3.1415927" "31415.9265359"
```

width 参数指定至少要有的输出宽度，不足时默认在左侧用空格填充，如：

```
format(1.000, width=6, nsmall=2)
```

```
[1] " 1.00"
```

format() 还有许多选项，详见函数的帮助。

### 35.1.3.2 sprintf() 函数

format() 函数无法精确控制输出长度和格式。sprintf 是 C 语言中 sprintf 的向量化版本，可以把一个元素或一个向量的各个元素按照 C 语言输出格式转换为字符型向量。第一个自变量是 C 语言格式的输出格式字符串，其中 %d 表示输出整数，%f 表示输出实数，%02d 表示输出宽度为 2、不够左填 0 的整数，%6.2f 表示输出宽度为 6、宽度不足时左填空格、含两位小数的实数，等等。

比如，标量转换

```
sprintf("%6.2f", pi)
```

```
[1] " 3.14"
```

又如，向量转换：

```
sprintf("tour%03d.jpg", c(1, 5, 10, 15, 100))
```

```
[1] "tour001.jpg" "tour005.jpg" "tour010.jpg" "tour015.jpg" "tour100.jpg"
```

还可以支持多个向量同时转换，如：

```
sprintf("%1dx%1d=%2d", 1:5, 5:1, (1:5)*(5:1))
```

```
[1] "1x5= 5" "2x4= 8" "3x3= 9" "4x2= 8" "5x1= 5"
```

### 35.1.3.3 字符串插值函数

许多脚本型程序设计语言都有在字符串的内容中插入变量值的功能，R 本身不具有这样的功能，`sprintf()` 函数有类似作用但只是一个不方便使用的副作用。

`stringr::str_glue()` 和 `stringr::str_glue_data()` 提供了字符串插值的功能。只要在字符串内用大括号写变量名，则函数可以将字符串内容中的变量名替换成变量值，如：

```
name <- " 李明"
tele <- "13512345678"
str_glue(" 姓名: {name}\n电话号码: {tele}\n")
```

```
姓名: 李明
电话号码: 13512345678
```

上面的例子直接用了换行符"`\n`" 来分开不同内容。也可以输入多个字符串作为自变量，内容自动连接在一起，可以用参数`.sep` 指定分隔符：

```
name <- " 李明"
tele <- "13512345678"
str_glue(" 姓名: {name}, ", " 电话号码: {tele}")
```

```
姓名: 李明, 电话号码: 13512345678
```

```
str_glue(" 姓名: {name}", " 电话号码: {tele}", .sep="; ")
```

```
姓名: 李明; 电话号码: 13512345678
```

也可以直接在 `str_glue()` 中指定变量值，如：

```
str_glue(" 姓名: {name}", " 电话号码: {tele}", .sep="; ",
 name = " 张三", tele = "13588888888")
```

```
姓名: 张三; 电话号码: 13588888888
```

`stringr::str_glue_data()` 则以一个包含变量定义的对象`.x` 为第一自变量，类型可以是环境、列表、数据框等。如：

```
str_glue_data(list(name = "王五", tele = "13500000000"),
 "姓名: {name}", "电话号码: {tele}", .sep="; ")
```

```
姓名: 王五; 电话号码: 13500000000
```

### 35.1.4 字符串长度

`stringr::str_length(string)` 求字符型向量 `string` 每个元素的长度。一个汉字长度为 1。

```
str_length(c("a", "bc", "def", "北京"))
```

```
[1] 1 2 3 2
```

函数 `nchar(text)` 计算字符串长度，默认按照字符个数计算而不是按字节数计算，如

```
nchar(c("a", "bc", "def", "北京"))
```

```
[1] 1 2 3 2
```

注意函数对输入的字符型向量每个元素计算长度。

`nchar()` 加选项 `type="bytes"` 可用按字符串占用的字节数计算，这时一个汉字占用多个字节（具体占用多少与编码有关）。如

```
nchar(c("a", "bc", "def", "北京"), type="bytes")
```

```
[1] 1 2 3 4
```

### 35.1.5 取子串

`stringr::str_sub(string, start, end)` 字符串子串，用开始字符位置 `start` 和结束字符位置 `end` 设定子串位置。用负数表示倒数位置。默认开始位置为 1，默认结束位置为最后一个字符。

如：

```
str_sub("term2017", 5, 8)
```

```
[1] "2017"
```

```
str_sub(c("term2017", "term2018"), 5, 8)
```

```
[1] "2017" "2018"
```

```
str_sub("term2017", 5)
```

```
[1] "2017"
```

```
str_sub("term2017", -4, -1)
```

```
[1] "2017"
```

```
str_sub("term2017", end=4)
```

```
[1] "term"
```

取子串时，一般按照字符个数计算位置，如

```
str_sub(" 北京市海淀区颐和园路 5 号", 4, 6)
```

```
[1] "海淀区"
```

当起始位置超过总长度或结束位置超过第一个字符时返回空字符串；当起始位置超过结束位置是返回空字符串。如：

```
str_sub("term2017", 9)
```

```
[1] ""
```

```
str_sub("term2017", 1, -9)
```

```
[1] ""
```

```
str_sub("term2017", 8, 5)
```

```
[1] ""
```

可以对 `str_sub()` 结果赋值，表示修改子串内容，如：

```
s <- "term2017"
str_sub(s, 5, 8) <- "18"
s
```

```
[1] "term18"
```

字符串替换一般还是应该使用专用的替换函数如 `stringr::str_replace_all()`, `gsub()`。

基本 R 的 `substring(text, first, last)` 函数与 `stringr::str_sub()` 功能相同, 但 `first` 和 `last` 参数不允许用负数, `last` 的默认值是一个很大的数, 所以省略 `last` 时会取到字符串末尾。 `substring()` 对三个参数 `text`, `first`, `last` 都是向量化的, 长度不一致时按照一般的不等长向量间运算规则处理。如:

```
substring(c("term2017", "term2018"), first=c(1, 5), last=c(4, 8))
```

```
[1] "term" "2018"
```

```
substring("term2017", first=c(1, 5), last=c(4, 8))
```

```
[1] "term" "2017"
```

`substring()` 也允许修改某个字符串的指定子串的内容, 如

```
s <- "123456789"
substring(s, 3, 5) <- "abc"
s
```

```
[1] "12abc6789"
```

R 的 `substr(x, start, stop)` 作用类似, 但是仅支持 `x` 为字符型向量, `start` 和 `stop` 是标量。

## 35.1.6 字符串变换

### 35.1.6.1 大小写

`stringr::str_to_upper(string)` 将字符型向量 `string` 中的英文字母都转换为大写。类似函数有 `stringr::str_to_lower(string)` 转换

为小写, `stringr::str_to_title(string)` 转换为标题需要的大小写, `stringr::str_to_sentence(string)` 转换为句子需要的大小写。这都是针对英文的, 选项 `locale` 用来选语言, `locale="en"` 为默认值。

基本 R 的 `toupper()` 将字符型向量的每个元素中的小写字母转换为大写, `tolower()` 转小写。

### 35.1.6.2 字符变换表

基本 R 的 `chartr(old, new, x)` 函数指定一个字符对应关系, 旧字符在 `old` 中, 新字符在 `new` 中, `x` 是一个要进行替换的字符型向量。比如, 下面的例子把所有! 替换成., 把所有; 替换成,:

```
chartr("!", ".", c("Hi; boy!", "How do you do!"))
```

```
[1] "Hi, boy." "How do you do."
```

```
chartr("。; 县", ".; 区", " 昌平区, 大兴县; 固安县。")
```

```
[1] "昌平区, 大兴区; 固安区."
```

第二个例子中被替换的标点是中文标点, 替换成了相应的英文标点。

### 35.1.6.3 空白处理

`stringr::str_trim(string, side)` 返回删去字符型向量 `string` 每个元素的首尾空格的结果, 可以用 `side` 指定删除首尾空格 ("both")、开头空格 ("left")、末尾空格 ("right")。如:

```
str_trim(c(" 李明", " 李明 ", " 李明 ", " 李 明"))
```

```
[1] "李明" "李明" "李明" "李 明"
```

```
str_trim(c(" 李明", " 李明 ", " 李明 ", " 李 明"), side="left")
```

```
[1] "李明" "李明" "李明" "李 明"
```

```
str_trim(c(" 李明", " 李明 ", " 李明 ", " 李 明"), side="right")
```

```
[1] " 李明" "李明" " 李明" "李 明"
```

`stringr::str_squish(string)` 对字符型向量 `string` 每个元素，删去首尾空格，将重复空格变成单个，返回变换后的结果。如：

```
str_squish(c(" 李明", "李明 ", "李明 ", "李明"))
```

```
[1] "李明" "李明" "李明" "李明"
```

基本 R 函数 `trimws(x, which)` 与 `str_trim()` 作用类似，选项 `which="left"` 可以仅删去开头的空格，选项 `which="right"` 可以仅删去结尾的空格。

```
trimws(c(" 李明", "李明 ", "李明 ", "李明"))
```

```
[1] "李明" "李明" "李明" "李明"
```

```
trimws(c(" 李明", "李明 ", "李明 ", "李明"), which="left")
```

```
[1] "李明" "李明" "李明" "李明"
```

```
trimws(c(" 李明", "李明 ", "李明 ", "李明"), which="right")
```

```
[1] "李明" "李明" "李明" "李明"
```

为了去掉输入字符串中所有空格，可以用 `gsub()` 替换功能，如：

```
gsub(" ", "", c(" 李明", "李明 ", "李明 ", "李明"), fixed=TRUE)
```

```
[1] "李明" "李明" "李明" "李明"
```

`stringr::str_pad(string, width)` 可以将字符型向量 `string` 的每个元素加长到 `width` 个字符，不足时左补空格，已经达到或超过 `width` 的则不变，如：

```
str_pad(c("12", "1234"), 3)
```

```
[1] " 12" "1234"
```

可以用选项 `side` 选择在哪里填补空格，默认为 `"left"`，还可选 `"right"`，`"both"`。

`stringr::str_wrap()` 可以将作为字符型向量的长字符串拆分成近似等长的行，行之间用换行符分隔。

#### 35.1.6.4 排序

基本 R 函数 `sort()` 可以用来对字符型向量的各个元素按照字典序排序, 但是字符的先后顺序是按照操作系统的当前编码值次序, 见关于 `locales` 的帮助。

`str_sort(x)` 对字符型向量 `x` 排序。可以用 `locale` 选项指定所依据的 locale, 不同的 locale 下次序不同。默认为 "en" 即英语, 中国大陆的 GB 编码 (包括 GBK 和 GB18030) 对应的 locale 是 "zh"。

`str_order(x)` 返回将 `x` 的各个元素从小到大排序的下标序列。

#### 35.1.7 简单匹配与查找

##### 35.1.7.1 开头和结尾匹配

基本 R 的 `startsWith(x, prefix)` 可以判断字符型向量 `x` 的每个元素是否以 `prefix` 开头, 结果为一个与 `x` 长度相同的逻辑型向量。如

```
startsWith(c("xyz123", "tu004"), "tu")
```

```
[1] FALSE TRUE
```

`endsWith(x, suffix)` 可以判断字符型向量 `x` 的每个元素是否以 `suffix` 结尾, 如

```
endsWith(c("xyz123", "tu004"), "123")
```

```
[1] TRUE FALSE
```

`stringr` 包的 `str_starts(string, pattern)` 判断 `string` 的每个元素是否以模式 `pattern` 开头, 加选项 `negate=TRUE` 表示输出反面结果。`pattern` 是正则表达式, 如果需要用非正则表达式, 可以用 `fixed()` 或者 `coll()` 保护, 如:

```
str_starts(c("xyz123", "tu004"), fixed("tu"))
```

```
[1] FALSE TRUE
```

```
str_starts(c("xyz123", "tu004"), coll("tu"))
```



```
[1] FALSE TRUE
```

stringr 包的 `str_ends(string, pattern)` 判断是否以给定模式结尾。

### 35.1.7.2 中间匹配

函数 `grep()`, `grep1()` 等可以用于查找子字符串，位置不限于开头和结尾，详见“正则表达式”章节。

在 `grep1()` 函数中加 `fixed=TRUE` 选项表示查找一般文本内容（非正则表达式）。比如，查找字符串中是否含有 `our`:

```
grep1("our", c("flavor", "tournament"), fixed=TRUE)
```

```
[1] FALSE TRUE
```

### 35.1.8 字符串替换

用 `gsub(pattern, replacement, x, fixed=TRUE)` 把字符型向量 `x` 中每个元素中出现的子串 `pattern` 都替换为 `replacement`。如

```
gsub("the", "**",
 c("New theme", "Old times", "In the present theme"),
 fixed=TRUE)
```

```
[1] "New **me" "Old times" "In ** present **me"
```

设有些应用程序的输入要求使用逗号“,”分隔，但是用户可能输入了中文逗号“，”，就可以用 `gsub()` 来替换：

```
x <- c("15.34,14.11", "13.25, 16.92")
x <- gsub(",", ", ", x, fixed=TRUE); x
```

```
[1] "15.34,14.11" "13.25,16.92"
```

例子中 `x` 的第二个元素中的逗号是中文逗号。

函数 `sub()` 与 `gsub()` 类似，但是仅替换第一次出现的 `pattern`。

### 35.1.9 字符串拆分

`stringr::str_split(string, pattern)` 对字符型向量 `string` 的每一个元素按分隔符 `pattern` 进行拆分，每个元素拆分为一个字符型向量，结果是一个列表，列表元素为字符型向量。其中 `pattern` 是正则表达式，为了按照固定模式拆分，用 `fixed()` 进行保护。如

```
x <- c("11,12", "21,22,23", "31,32,33,34")
res1 <- str_split(x, fixed(","))
res1
```

```
[[1]]
[1] "11" "12"
##
[[2]]
[1] "21" "22" "23"
##
[[3]]
[1] "31" "32" "33" "34"
```

`str_split()` 可以用选项 `n` 指定仅拆分成成几项，最后一项合并不拆分，如：

```
x <- c("11,12", "21,22,23", "31,32,33,34")
res2 <- str_split(x, fixed(","), n=2)
res2
```

```
[[1]]
[1] "11" "12"
##
[[2]]
[1] "21" "22,23"
##
[[3]]
[1] "31" "32,33,34"
```

拆分的结果可以用 `lapply()`, `sapply()`, `vapply()` 等函数处理。例如，将每个元素的拆分结果转换成数值型：

```
lapply(res1, as.numeric)
```

```
[[1]]
[1] 11 12
##
[[2]]
[1] 21 22 23
##
[[3]]
[1] 31 32 33 34
```

可以用 `unlist()` 函数将列表中的各个向量连接成一个长向量，如：

```
unlist(res1)
```

```
[1] "11" "12" "21" "22" "23" "31" "32" "33" "34"
```

注意，即使输入只有一个字符串，`str_split()` 的结果也是列表，所以输入只有一个字符串时我们应该取出结果列表的第一个元素，如

```
strsplit("31,32,33,34", split=",", fixed=TRUE)[[1]]
```

```
[1] "31" "32" "33" "34"
```

如果确信每个字符串拆分出来的字符串个数都相同，可以用 `stringr::str_split_fixed()`，用参数 `n` 指定拆出来的项数，这时结果为一个字符型矩阵，原来的每个元素变成结果中的一行：

```
x <- c("11,12", "21,22", "31,32")
res3 <- str_split_fixed(x, fixed(","), n=2)
res3
```

```
[,1] [,2]
[1,] "11" "12"
[2,] "21" "22"
[3,] "31" "32"
```

基本 R 的 `strsplit(x,split,fixed=TRUE)` 可以把字符型向量 `x` 的每一个元素按分隔符 `split` 拆分为一个字符型向量，`strsplit` 的结果为一个列表，

每个列表元素对应于 `x` 的每个元素。

如

```
x <- c("11,12", "21,22,23", "31,32,33,34")
res4 <- strsplit(x, split=",", fixed=TRUE)
res4

[[1]]
[1] "11" "12"
##
[[2]]
[1] "21" "22" "23"
##
[[3]]
[1] "31" "32" "33" "34"
```

## 35.2 文本文件读写

文本文件是内容为普通文字、用换行分隔成多行的文件，与二进制文件有区别，二进制文件中换行符没有特殊含义，而且二进制文件的内容往往也不是文字内容。二进制文件的代表有图片、声音，以及各种专用软件的私有格式文件，如 Word 文件、Excel 文件。

对于文本文件，可以用 `readLines()` 函数将其各行的内容读入为一个字符型数组，字符型数组的每一个元素对应于文件中的一行，读入的字符型数组元素不包含分隔行用的换行符。

最简单的用法是读入一个本地的文本文件，一次性读入所有内容，用如

```
lines <- readLines("filename.ext")
```

其中 `filename.ext` 是文件名，也可以用全路径名或相对路径名。

当文本文件很大的时候，整体读入有时存不下，即使能存下处理速度也很慢，可以一次读入部分行，逐批读入并且逐批处理，这样程序效率更高。这样的程序要复杂一些，例如

```
infcon <- file("filename.ext", open="rt")
batch <- 1000
repeat{
 lines <- readLines(infcon, n=batch)
 if(length(lines)==0) break
 ## 处理读入的这些行
}
close(infcon)
```

以上程序先打开一个文件，`infcon` 是打开的文件的读写入口（称为一个“连接对象”）。每次读入指定的行并处理读入的行，直到读入了 0 行为止，最后关闭 `infcon` 连接。

对文本文件的典型处理是读入后作一些修改，另外保存。函数 `writeLines(lines, con="outfilename.txt")` 可以将字符型向量 `lines` 的各个元素变成输出文件的各行保存起来，自动添加分隔行的换行符。如果是分批读入分批处理的，则写入也需要分批写入，以上的分批处理程序变成：

```
infcon <- file("filename.ext", open="rt")
outfcon <- file("outfilename.txt", open="wt")
batch <- 1000
while(TRUE){
 lines <- readLines(infcon, n=batch)
 if(length(lines)==0) break
 ## 处理读入的这些行，变换成 outlines
 writeLines(outlines, con=outfcon)
}
close(outfcon)
close(infcon)
```

`readLines()` 也可以直接读取网站的网页文件，如

```
lines <- readLines(url("https://www.r-project.org/"))
length(lines)
[1] 116
head(lines)
```

```
[1] "<!DOCTYPE html>"
[2] "<html lang=\"en\""
[3] " <head>"
[4] " <meta charset=\"utf-8\""
[5] " <meta http-equiv=\"X-UA-Compatible\" content=\"IE=edge\""
[6] " <meta name=\"viewport\" content=\"width=device-width, initial-scale=1\"
```

`readr` 包的 `read_lines()` 和 `write_lines()` 函数起到与基本 R 中 `readLines()` 和 `writeLines()` 类似的作用, `read_file()` 和 `read_file_raw()` 可以将整个文件读入为一个字符串。

关于读写文件时的编码问题, 详见15.6。

### 35.3 正则表达式

在对字符串进行查找或替换时, 有时要查找替换的不是固定的子串而是某种模式。比如, 要查找或替换连续的三个数字, 正文中的电子邮件地址, 网址, 电话号码, 等等。正则表达式 (regular expressions) 用于表示各种复杂模式。基本 R 中的正则表达式规则可以用 POSIX 1003.2 标准或者 Perl 规则。建议使用 perl 语言的正则表达式, 在基本 R 的有关函数中规定参数 `perl=TRUE`。

`stringr` 包提供了更方便的正则表达式功能, 其正则表达式规则是 ICU 正则表达式规则, 针对 UTF-8 编码的文本数据, 基本与 perl 规则兼容。

在正则表达式的模式 (pattern) 中,

```
.*+?{\}\[\]^$()
```

等字符是特殊字符, 有特殊的解释。除了\之外的其它 12 个都称为“元字符” (meta characters)。

在 R 语言中使用正则表达式时, 需要注意 R 字符型常量中一个\要写成两个。

### 35.3.1 字面匹配与匹配显示

如果模式中不含特殊字符,匹配为原样的子串。也叫做字面 (literal) 匹配。`stringr` 包提供了定义正则表达式、匹配正则表达式、按正则表达式替换、抽取匹配结果、用富文本显示匹配结果等强大功能,其中 `str_view()` 函数可以在 HTML 输出中或者在 RStudio 软件中用富文本显示匹配结果,在源数据中加亮显示匹配。注意,如果你现在看的是 PDF 文件,结果可能无法显示。

下面的程序在字符型向量 `x` 的三个字符串元素中查找子字符串"the" 并加亮显示:

```
x <- c("New theme", "Old times", "In the present theme")
str_view(x, "the")
```

在 RStudio 中会在 Viewer 窗格显示匹配结果,匹配内容被高亮显示。

源数据中的第三项实际上有两处"the" 出现但结果只显示了第一处。用 `str_view_all()` 查看所有匹配,如:

```
x <- c("New theme", "Old times", "In the present theme")
str_view_all(x, "the")
```

### 35.3.2 不区分大小写匹配和 `regex` 函数

`str_view(string, pattern)` 中的 `pattern` 应该为正则表达式类型,如果输入了字符串,会自动被函数 `regex()` 转换成正则表达式类型。正则表达式的模式一般是区分大小写的,通过在 `regex()` 函数中加选项 `ignore_case=TRUE` 可以进行不区分大小写的匹配;在模式前面附加 `(?i)` 前缀式选项也可以实现不区分大小写匹配。如

```
str_view_all(c("Dr. Wang", "DR. WANG", "dR. W.R."), "Dr")
str_view_all(c("Dr. Wang", "DR. WANG", "dR. W.R."),
 regex("Dr", ignore_case=TRUE))
str_view_all(c("Dr. Wang", "DR. WANG", "dR. W.R."), "(?i)Dr")
```

### 35.3.3 用句点匹配单个字符

在模式中用 “.” 匹配任意一个字符（除了换行符“\n”，能否匹配此字符与选项有关）。如

```
s <- c("abc", "cabs", "lab")
str_view_all(s, "ab.")
```

像句点这样的字符称为元字符（meta characters），在正则表达式中有特殊作用。如果需要匹配句点本身，用 “[.]” 或者 “\.” 表示。比如，要匹配 a.txt 这个文件名，如下做法有错误：

```
str_view_all(c("a.txt", "a0txt"), "a.txt")
```

结果连 a0txt 也匹配了。用 “[.]” 表示句点则将句点不做特殊解释：

```
str_view_all(c("a.txt", "a0txt"), "a[.]txt")
str_view_all(c("a.txt", "a0txt"), "a\\.txt")
```

注意在 R 语言字符型常量中一个 \ 需要写成两个。

如果仅需按照原样进行查找，也可以将 pattern 的字符串用 fixed() 函数保护，如：如

```
str_view_all(c("a.txt", "a0txt"), fixed("a.txt"))
```

### 35.3.4 匹配一组字符中的某一个

模式中使用方括号给定一个字符类，单个字符与字符类中任何一个字符相同都算是匹配成功。比如，模式 “[ns]a.[.]xls” 表示匹配的第一个字符是 n 或 s，第二个字符是 a，第三个字符任意，第四个字符是句点，然后是 xls。测试：

```
str_view_all(c("sa1.xls", "dna2.xlss", "nat.xls"), "[ns]a.[.]xls")
```

注意匹配并不需要从开头匹配到结尾，中间匹配是允许的，类似于搜索符合某种规律的字串。在上例中第二个元素是从第二个字符开始匹配的，也没有匹配到末尾。

例：模式 [Rr]eg[Ee]x 可以匹配 RegEx 或 Regex 或 regex 或 regEx。



如果希望完全忽略大小写进行匹配，可以使用在 `regex()` 等函数中指定 `ignore_case=TRUE` 选项，或在模式字符串的最前面添加 `(?i)` 选项。

在“`[]`”中允许用`-`表示一个范围。如 `[a-z]` 匹配小写英文字母，`[A-Z]` 匹配大写英文字母，`[a-zA-Z]` 匹配大小写的英文字母，`[a-zA-Z0-9]` 匹配大小写的英文字母和数字。

为了匹配一个 16 进制数字，可以用 `[0-9A-Fa-f]`。

例：模式“`[ns]a[0-9][.]xls`”要求匹配的第三个字符为数字。

```
str_view_all(c("sa1.xls", "dna2.xlss", "nat.xls"), "[ns]a[0-9][.]xls")
```

在方括号内第一个位置的 `^` 表示对指定的范围取余集。例如，模式 `[ns]a[^0-9][.]xls` 要求匹配的第三个字符不能为数字：

```
str_view_all(c("sa1.xls", "dna2.xlss", "nat.xls"), "[ns]a[^0-9][.]xls")
```

### 35.3.5 原样匹配元字符

元字符 (meta characters) 是在正则表达式中有特殊含义的字符。比如句点可以匹配任意一个字符，左方括号代表字符集合的开始。所以元字符不能直接匹配自身，可以用“`[.]`”匹配一个句点。为匹配左方括号，在前面加上转义字符`\`变成`\[`，但是在 R 字符串中一个`\`必须用`\\`表示，所以模式“`\[`”在 R 中写成字符串常量，必须写成`"\\["`。其它的元字符如果要原样匹配也可以在前面加上转义字符`\`，比如匹配`\`本身可以用`\\`，但是在 R 字符型常量中需要写成`"\\\\"`。

例，匹配 `x[5]`，因为 `[`是元字符，需要写成：

```
str_view_all(c("int x;", "int x[5]"), "int x\\[5\\]")
```

Perl 中允许用“`[[]`”表示“`[`”，用“`[][]`”表示“`]`”，但 `stringr` 包不支持这种做法。

### 35.3.6 匹配空白

表示空白的元字符有：

`\f` 换页符

`\n` 换行符  
`\r` 回车符  
`\t` 制表符  
`\v` 垂直制表符

不同操作系统的文本文件的行分隔符不同，为了匹配 Windows 格式的文本文件中的空行，用 “`\r\n\r\n`”；为了匹配 Unix 格式的文本文件中的空行则用 “`\r\r`”。写成 R 的字符型常量时，这些表示本身也是 R 的相应字符的表示，所以在 R 字符型常量中这些字符不需要用两个 `\` 表示一个 `\`。

匹配任意一个空白字符用 “`\s`”，这等价于 “[ `\f\n\r\t\v`]”。大写的 “`\S`” 则匹配任意一个非空白的字符。

### 35.3.7 匹配数字

用 `\d` 匹配一个数字，相当于 [0-9]。用 `\D` 匹配一个非数字。如

```
str_view_all(c("n1.xls", "na.xls"), "n\\d[.]xls")
```

### 35.3.8 匹配开头和末尾

模式匹配相当于在字符串内部搜索某种模式，如果要从字符串开头匹配，在模式中取第一个模式规定为 `^` 或 `\A`。如果模式中最后一个字符是 `$` 或 `\Z`，则需要匹配到字符串末尾。用 `\Z` 匹配字符串末尾时如果末尾有一个换行符则匹配到换行符之前。

如

```
str_view_all(c("n1.xls", "na.xls", "cn1.xls", "n1.xlsx"), "^n\\d[.]xls$")
str_view_all(c("n1.xls", "na.xls", "cn1.xls", "n1.xlsx"), "\\An\\d[.]xls\\Z")
```

只匹配了第一个输入字符串。

有时候源文本的每个字符串保存了一个文本文件内容，各行用 `\n` 分隔，后面将给出匹配每行的行首与行尾的方法。

### 35.3.9 匹配字母、数字、下划线

匹配字母、数字、下划线字符用 `\w` (小写), 等价于 `[a-zA-Z0-9_]`。`\W` (大写) 匹配这些字符以外的字符。如

```
str_view_all(c("file-s1.xls", "s#.xls"), "s\\w[.]")
```

可以看出, 模式匹配了 `s1.` 而没有匹配 `s#.`。

### 35.3.10 十六进制和八进制数

在模式中可以用十六进制数和八进制数表示特殊的字符。十六进制数用 `\x` 引入, 比如 `\x0A` 对应 `\n` 字符。八进制数用 `\0` 引入, 比如 `\011` 表示 `\t` 字符。

例如, `str_view_all("abc\\nefg\\n", "\\x0A")` 可以匹配两个换行符。

### 35.3.11 POSIX 字符类

`\d`, `\w` 这样的字符类不方便用在方括号中组成字符集合, 而且也不容易记忆和认读。在模式中方括号内可以用 `[:alpha:]` 表示任意一个字母。比如, `[[:alpha:]]` 匹配任意一个字母 (外层的方括号表示字符集合, 内层的方括号是 POSIX 字符类的固有界定符)。

这样的 POSIX 字符类有:

- `[:alpha:]` 表示任意一个字母;
- `[:lower:]` 为小写字母;
- `[:upper:]` 为大写字母;
- `[:digit:]` 为数字;
- `[:xdigit:]` 为十六进制数字。
- `[:alnum:]` 为字母数字 (不包括下划线);
- `[:blank:]` 为空格或制表符;
- `[:space:]` 为任何一种空白字符, 包括空格、制表符、换页符、换行符、回车符;
- `[:print:]` 为可打印字符;
- `[:graph:]` 和 `[:print:]` 一样但不包括空格;

- `[:punct:]` 为 `[:print:]` 中除 `[:alnum:]` 和空白以外的所有字符;

例如:

```
str_view_all(c("x1", "_x", ".x", ".1"), "[[:alpha:]]_[:alnum:]]_")
```

模式匹配长度为 2 的字符串, 第一个字符是字母、下划线或者小数点, 第二个字符是字母、数字、下划线或者小数点。这个模式试图匹配由两个字符组成的合法 R 变量名, 但是最后一个非变量名 `.1` 也被匹配了。解决这样的问题可以采用后面讲到的 | 备择模式。

### 35.3.12 重复匹配

#### 35.3.12.1 加号重复匹配

模式中在一个字符或字符集合后加后缀 `+` 表示一个或多个前一字符。比如

```
str_view_all(c("sa1", "dsa123"), "sa[[:digit:]]+")
```

例如, 匹配电子邮件地址:

```
str_view_all("abc123@efg.com",
 "^[[:alnum:]]_+@[[:alnum:]]_+[.][:alnum:]]_+$")
```

匹配的电子邮件地址在 `@` 前面可以使用任意多个字母、数字、下划线, 在 `@` 后面由小数点分成两段, 每段可以使用任意多个字母、数字、下划线。这里用了 `^` 和 `$` 表示全字符串匹配。

#### 35.3.12.2 星号和问号重复匹配

在一个字符或字符集合后加后缀 `*` 表示零个或多个前一字符, 后缀 `?` 表示零个或一个前一字符。

比如, `^https?:/[[:alnum:]]_./]+$` 可以匹配 `http` 或 `https` 开始的网址。如

```
str_view_all(c("http://www.163.net", "https://123.456."),
 "^https?:/[[:alnum:]]_./]+$")
```

(注意第二个字符串不是合法网址但是按这个正则表达式也能匹配)

`x[[:digit:]]*` 能匹配 “x”, “x1”, “x123” 这样的变量名, 如:

```
str_view_all(c("x", "x1", "x123"), "x[[:digit:]]*")
str_view_all(c("x", "x1", "x123"), "x\\d*")
```

### 35.3.12.3 计数重复

问号可以表示零个或一个, 而加号、星号重复不能控制重复次数。在后缀大括号中写一个整数表示精确的重复次数。如

```
str_view_all(c("1", "12", "123", "1234"), "[[:digit:]]{3}")
```

模式匹配的是三位的数字。因为没有要求从开头一直匹配到末尾, 所以三位以上数字也能匹配其中开始的三位。

可以在后缀大括号中指定重复的最小和最大次数, 中间用逗号分隔。比如, 月日年的日期格式可以用

```
[[:digit:]]{1,2}[-/] [[:digit:]]{1,2}[-/] [[:digit:]]{2,4}
```

来匹配。如 (注意这个模式还会匹配非日期)

```
pat <- paste0(
 "[[:digit:]]{1,2}[-/]",
 "[[:digit:]]{1,2}[-/]",
 "[[:digit:]]{2,4}")
str_view_all(c("2/4/1998", "13/15/198"), pat)
```

重复数允许指定为 0。重复数的逗号后面空置表示重复数没有上限。例如, 后缀 `{3,}` 表示前一模式必须至少重复 3 次。

### 35.3.13 贪婪匹配和懒惰匹配

无上限的重复匹配如 `*`, `+`, `{3,}` 等缺省是贪婪型的, 重复直到文本中能匹配的最长范围。比如我们希望找出圆括号这样的结构, 很容易想到用 `\(.+\)` 这样的模式 (注意圆括号是元字符, 需要用反斜杠保护), 但是这不会恰好匹配一次, 模式会一直搜索到最后一个) 为止。

例如:

```
str_view_all("(1st) other (2nd)", "\\(.+\\)")
```

我们本来期望的是提取两个“(1st)”和“(2nd)”组合，不料整个地提取了“(1st) other (2nd)”。这就是因为.+的贪婪匹配。

如果要求尽可能短的匹配，使用\*?, +?, {3,}?等“懒惰型”重复模式。在无上限重复标志后面加问号表示懒惰性重复。

比如，上例中模式修改后得到了期望的结果:

```
str_view_all("(1st) other (2nd)", "\\(.+?\\)")
```

懒惰匹配会造成搜索效率降低，应仅在需要的时候使用。

### 35.3.14 单词边界

用\b匹配单词边界，这样可以查找作为单词而不是单词的一部分存在的内容。  
\B匹配非单词边界。如

```
str_view_all(c("a cat meaos", "the category"), "\\bcat\\b")
```

### 35.3.15 句点全匹配与多行模式

句点通配符一般不能匹配换行，如

```
str_view_all("(1,\n2)", "\\(.+?\\)")
```

跨行匹配失败。一种办法是预先用str\_replace\_all()或gsub()把所有换行符替换为空格。但是这只能解决部分问题。

解决方法是在将模式用regex()保护并加选项dotall=TRUE，或者在Perl正则表达式开头添加(?s)选项，这样使得句点通配符可以匹配合换行符，称为句点全匹配模式。如

```
str_view_all("(1,\n2)", regex("\\(.+?\\)", dotall=TRUE))
str_view_all("(1,\n2)", "(?s)\\(.+?\\)")
```

在 `regex()` 函数中加选项 `multiline=TRUE`, 或者在正则表达式开头用 `(?m)` 表示把整个输入字符串看成用换行符分开的多行。这时 `^` 和 `$` 匹配每行的开头和结尾, “每行”是指字符串中用换行符分开的各个字符子串。`(?s)` 与 `(?m)` 可以同时使用。

例:

```
str_view_all("(1,2)\n(3,4)\n", "^\\(\\.+?\\)$")
```

元数据中包含两行内容, 结果没有能够匹配, 这是因为模式要求从整个字符串开头一直匹配到末尾。增加 `multiline=TRUE` 或者 `(?m)` 选项则可以匹配两处:

```
str_view_all("(1,2)\n(3,4)\n", regex("^\\(\\.+?\\)$", multiline=TRUE))
str_view_all("(1,2)\n(3,4)\n", "(?m)^\\(\\.+?\\)$")
```

### 35.3.15.1 逐行处理

虽然正则表达式有多行和跨行选项, 但是当源数据很长时, 匹配效率会很低。

R 的 `readLines()` 函数可以把一整个文本文件读成一个字符型向量, 每个元素为一行, 元素中不包含换行符。R 的字符型函数可以对这样的字符型向量每个元素同时处理, 也就实现了逐行处理。

如果字符串 `x` 中包含了一整个文本文件内容, 其中以 `\n` 分隔各行, 为了实现逐行处理, 可以先用 `str_split()` 函数拆分成不同行:

```
cl <- str_split(x, "\r?\n")[[1]]
```

结果将是一个字符型向量, 每个元素是原来的一行, 最后一个元素是空字符串。如

```
x <- c("This is first line.\nThis is second line.\n")
cl <- str_split(x, "\r?\n")[[1]]
cl
```

```
[1] "This is first line." "This is second line." ""
```

### 35.3.16 备择模式

如果有两种模式都算正确匹配，则用 | 连接这两个模式表示两者都可以。例如，某人的名字用 James 和 Jim 都可以，表示为 James|Jim，如

```
str_view(c("James, Bond", "Jim boy"), "James|Jim")
```

两个字符的合法 R 变量名的匹配：

```
str_view_all(c("x1", "_x", ".x", ".1"),
 "[[:alpha:]]_|[[:alnum:]]_|[.][[:alpha:]]_")
```

### 35.3.17 分组与捕获

在正则表达式中用圆括号来分出组，作用是

- 确定优先规则
- 组成一个整体
- 拆分出模式中的部分内容（称为捕获）
- 定义一段供后续引用或者替换。

圆括号中的模式称为子模式，或者捕获。

在使用备择模式时，James|Jim 是在单词 James 和 Jim 之间选择。如果希望选择的是中间的 ms 和 Ji 怎么办？可以将备择模式保护起来，如 Jam(es|Ji)m，就可以确定备择模式的作用范围。

有时一个模式中部分内容仅用于定位，而实际有用的内容是其中的一部分，就可以将这部分有用的内容包在圆括号中作为一个捕获。

元字符问号、加号、星号、大括号等表示重复，前面的例子中都是重复一个字符或者字符类。如果需要重复由多个字符组成的模式，如 x[[:digit:]]{2} 怎么办？只要将该模式写在括号中，如：

```
str_view_all(c("x01x02", "_x11x9"),
 "(x[[:digit:]]{2})+")
```

上例的元数据中，第一个元素重复了两次括号中的模式，第二个元素仅有一次括号中的模式。



**注意：**用表示重复的元字符重复某个模式时，从第二次开始，并不是要重复前面的子字符串，而是重复前面的模式。比如上例中 `x01x02`。

如果想严格重复前面的某个子字符串怎么办？

分组是自动编号的，以左开括号的序号为准（除了作为选项、有名捕获等开扩号以外）。在替换或者向后引用时，可以用 `\1`，`\2` 等表示匹配中第一个开扩号对应的分组，第二个开扩号对应的分组，.....。

在模式中可以用 `\1`，`\2` 等表示严格重复前面捕获的子字符串。例如，`([a-z]{3})\1` 这样的模式可以匹配如 `abcabc`，`uxzuxz` 这样的三字母重复：

```
str_view_all(c("abcabc", "aabbcc"),
 "([a-z]{3})\\1")
```

又例如，下面的程序找出了年（后两位）、月、日数字相同的日期：

```
str_view_all(c("2008-08-08", "2017-01-18"),
 "\\d\\d(\\d\\d)-\\1-\\1")
```

`stringr` 的函数 `str_replace_all(string, pattern, replacement)` 可以指定一个查找模式 `pattern`，替换模式 `replacement`，对字符型向量 `string` 中的每个元素进行替换。`replacement` 中的元字符没有特殊解释，但是 `\1`，`\2` 等代表匹配的子模式（捕获）。

例：希望把带有前导零的数字的前导零删除，可以用如

```
x <- c("1204", "01204", "001204B")
pat <- "\\b0+([1-9][0-9]*)\\b"
repl <- "\\1"
str_view_all(x, pat)
str_replace_all(x, pat, repl)
```

上例的模式中的 `\b` 表示单词边界，所以中间的 `0` 不会被当作前导零，不是整个数字的也不会被修改。

上例中的 `str_view_all()` 仅用于调试目的，在进行替换时不必必要步骤。

例：为了交换纵横坐标，可以用如下替换

```
x <- "1st: (5,3.6), 2nd: (2.5, 1.1)"
pat <- paste0(
 "([[:digit:]]+)",
 "[[:space:]]*([[:digit:]]+)", sep=""
)
repl <- "(\\2, \\1)"
str_view_all(x, pat)
str_replace_all(x, pat, repl)
```

例: 要匹配 yyyy-mm-dd 这样的日期, 并将其改写为 mm/dd/yyyy, 就可以用这样的替换模式:

```
x <- c("1998-05-31", "2017-01-14")
pat <- "[0-9]{4}-([0-9]{1,2})-([0-9]{1,2})"
repl <- "\\2/\\3/\\1"
str_view_all(x, pat)
str_replace_all(x, pat, repl)
```

如果某个分组仅想到分组作用但是不会提取具体的匹配内容也不会用该组内容做替换, 可以将该组变成“非捕获分组”, 办法是把表示分组开始左圆括号变成 (? : 三个字符。这在用分组表示优先级时比较有用, 如 "Jam(es|Ji)m" 可以写成 "Jam(?:es|Ji)m"。非捕获分组在向后引用和替换时不计入 \\1、\\2 这样的排列中。

比如, 把 1921-2020 之间的世纪号删去, 可以用

```
x <- c("1978", "2017", "2035")
pat <- "\\A(?:19|20)([0-9]{2})\\Z"
repl <- "\\1"
str_view_all(x, pat)
str_replace_all(x, pat, repl)
```

其中用了非捕获分组使得备择模式 19|20 优先匹配。注意模式并没有能保证日期在 1921-2020 之间。更周密的程序可以写成:

```
x <- c("1978", "2017", "2035")
pat1 <- "\\A19(2[1-9]|[3-9][0-9])\\Z"
pat2 <- "\\A20([01][0-9]|20)\\Z"
```

```
repl <- "\\1"
str_view_all(x, pat1)
str_view_all(x, pat2)
x %>%
 str_replace_all(pat1, repl) %>%
 str_replace_all(pat2, repl)
```

这里用了 `stringr` 包重新定义的管道运算，与 `magrittr` 包定义的管道运算作用相同，可以将函数的第一个自变量自动输入给管道的下一个处理层级。

## 35.4 stringr 包的正则表达式函数

### 35.4.1 str\_view() 函数

`str_view(string, pattern)` 在 RStudio 中打开 Viewer 窗格，显示 `pattern` 给出的正则表达式模式在 `string` 中的首个匹配。`string` 是输入的字符型向量。用 `str_view_all()` 显示所有匹配。

如果要匹配的是固定字符串，写成 `str_view(string, fixed(pattern))`。

如果要匹配的是单词等的边界，模式用 `boundary()` 函数表示，如 `str_view("a brown fox", boundary("word"))` 将匹配首个单词。

### 35.4.2 regex() 函数

`stringr` 包用到正则表达式模式的地方，实际上应该写成 `regex(pattern)`，只写模式本身是一种简写。`regex()` 函数可以指定 `ignore_case=TRUE` 要求不区分大小写，指定 `multi_line=TRUE` 使得 `^` 和 `$` 匹配用换行符分开的每行的开头和结尾，`dotall=TRUE` 使得 `.` 能够匹配换行符。`comment=TRUE` 使得模式可以写成多行，行尾的井号后面表示注释，这时空格不再原样匹配，为了匹配空格需要写在方括号内或者用反斜杠开头。

与 `regex()` 类似的表示模式的函数有 `fixed()`，`boundary()`，`coll()`。

### 35.4.3 检查那些元素能够匹配

`str_detect(string, pattern)` 返回字符型向量 `string` 的每个元素是否匹配 `pattern` 中的模式的逻辑型结果。与基本 R 的 `grepl()` 作用类似。如

```
x <- c("New theme", "Old times", "In the present theme")
str_view(x, "the")
str_detect(x, "the")
```

上例中的 `str_view()` 仅用作调试目的。

`str_which(string, pattern)` 返回字符型向量 `string` 的元素当中能匹配 `pattern` 中的模式的元素序号。与基本 R 的 `grep()` 作用类似。如

```
x <- c("New theme", "Old times", "In the present theme")
str_which(x, "the")
```

```
[1] 1 3
```

`str_count()` 则返回模式在每个元素中匹配的次数。如

```
str_count(c("123,456", "011"), "[[:digit:]]")
```

```
[1] 6 3
```

### 35.4.4 替换

`stringr` 包的 `str_replace_all(string, pattern, replacement)` 在字符型向量 `string` 的每个元素中查找模式 `pattern`，并将所有匹配按照 `replacement` 进行替换。在 `replacement` 可以用 `\1`、`\2` 中表示模式中的捕获，除此之外元字符没有特殊作用。

基本 R 中 `gsub()` 有类似功能。

如：

```
str_replace_all(c("123,456", "011"), ",", "")
```

```
[1] "123456" "011"
```

又如：

```
str_replace_all(c("123,456", "011"),
 "([[:digit:]]+),([[:digit:]]+)", "\\2,\\1")

[1] "456,123" "011"
```

注意源数据中第二个元素因为不能匹配所以就原样返回了，没有进行替换。

`str_replace()` 则仅对输入字符型向量的每个元素中模式的第一次出现进行替换，不如 `str_replace_all()` 常用。

### 35.4.5 返回匹配的元素

`str_subset(string, pattern)` 返回字符型向量中能匹配 `pattern` 的那些元素组成的子集，与基本 R 函数 `grep(pattern, string, value=TRUE)` 效果相同。注意，返回的是整个元素而不是匹配的子串。

比如，查找人名中间有空格的：

```
str_view_all(c("[马思聪]", "[李 明]"),
 "[[:alpha:]]+[[:space:]]+[[:alpha:]]+")
str_subset(c("[马思聪]", "[李 明]"),
 "[[:alpha:]]+[[:space:]]+[[:alpha:]]+")
```

注意上例中仅返回了有匹配的元素，而且是匹配元素的整个字符串而不是匹配的部分。

当要查找的内容是 `tibble` 的一列时，用 `filter()` 与 `str_detect()` 配合，可以进行行子集选择。比如，在数据框的人名中查找中间有空格的名字：

```
tibble(name=c(" 马思聪", " 李 明")) %>%
 filter(str_detect(name, "[[:alpha:]]+[[:space:]]+[[:alpha:]]+"))

A tibble: 1 x 1
name
<chr>
1 李 明
```

### 35.4.6 提取匹配内容

`str_subset()` 返回的是有匹配的源字符串，而不是匹配的部分子字符串。用 `str_extract(string, pattern)` 从源字符串中取出首次匹配的子串。

如

```
str_view_all("A falling ball", "all")
str_extract("A falling ball", "all")
```

`str_extract_all(string, pattern)` 取出所有匹配子串，结果是一个列表，列表的每个元素对应于字符型向量 `string` 的每个元素，结果列表的每个元素是一个字符型数组，存放所有匹配的子字符串。如：

```
x <- c("A falling ball", "Phone call.")
str_view_all(x, "all")
str_extract_all(x, "all")
```

`str_extract_all()` 可以加选项 `simplify=TRUE`，使得返回结果变成一个字符型矩阵，每行是原来一个元素中取出的各个子串，列数等于最大匹配次数，没有那么多匹配次数的填以空字符串。如果正常匹配结果不会出现空字符就可以用这种方法简化结果的保存和访问。如

```
x <- c("A falling ball", "Phone call.")
str_view_all(x, "all")
str_extract_all(x, "all", simplify=TRUE)
```

### 35.4.7 提取分组捕获内容

`str_subset()` 提取的是能匹配模式的元素子集，而不是匹配的模式或者捕获；`str_extract()` 和 `str_extract_all()` 提取的是每个元素的首次或者所有匹配的子字符串，而不是其中的捕获。

`str_match(string, pattern)` 提取每个元素的首次匹配内容以及其中各个捕获分组内容，结果是一个矩阵，每行对应于字符型向量 `string` 中的一个元素，结果矩阵的每行的第一个元素是匹配内容，其它元素是各个捕获，没有则为字符型缺失值（不是空字符串）。

比如，希望匹配中间有空格的人名并捕获空格前后部分：

```
str_match(c(" 马思聪", " 李 明"),
 "([[:alpha:]]+)[[:space:]]+([[:alpha:]]+)")
```

```
[,1] [,2] [,3]
[1,] NA NA NA
[2,] "李 明" "李" "明"
```

上例中源数据第一个元素没有匹配，所以结果都是缺失值 `NA`，第二个元素的结果在第二行，首先是整个匹配的子字符串，然后是捕获的两个部分。

`stringr::str_match_all(string, pattern)` 匹配每个字符串中所有出现位置，结果是一个列表，每个列表元素对应于输入的字符型向量 `string` 的每个元素，结果中每个列表元素是一个字符型矩阵，用来保存所有各个匹配以及匹配中的捕获，每行是一个匹配的结果，首先是匹配结果，其次是各个捕获。结果列表中每个作为列表元素的矩阵大小不一定相同。比如，模式为 `19xx` 或者 `20xx` 的年份，并将其分为前两位和后两位：

```
x <- c("1978-2000", "2011-2020-2099")
pat <- "\\b(19|20)([0-9]{2})\\b"
str_view_all(x, pat)
str_match_all(x, pat)
```

### 35.4.8 定位匹配位置

`str_locate(string, pattern)` 对输入字符型向量 `string` 的每个元素返回首次匹配 `pattern` 的开始和结束位置。输出结果是一个两列的矩阵，每行对应于输入的一个元素，每行的两个元素分别是首次匹配的开始和结束字符序号（按字符计算）。如

```
x <- c("A falling ball", "Phone call.")
str_view_all(x, "all")
str_locate(x, "all")
```

`str_locate_all(string, pattern)` 则可以返回每个元素中所有匹配的开始和结束位置，结果是一个列表，每个列表元素对应于输入字符型向量的每个元

素，结果中每个列表元素是一个两列的数值型矩阵，每行为一个匹配的的开始和结束字符序号。如

```
x <- c("A falling ball", "Phone call.")
str_view_all(x, "all")
str_locate_all(x, "all")
```

注意如果需要取出匹配的元素可以用 `str_subset()`，要取出匹配的子串可以用 `str_extract()` 和 `str_extract_all()`，取出匹配的子串以及分组捕获可以用 `str_match()` 和 `str_match_all()`。

## 35.5 利用基本 R 函数进行正则表达式处理

基本 R 函数 `grep`, `sub`, `gsub`, `regexpr`, `gregexpr`, `regexec` 中的 `pattern` 参数可以是正则表达式，这时应设参数 `fixed=FALSE`。`strsplit` 函数中的参数 `split` 也可以是正则表达式。`regmatches` 函数从 `regexpr`, `gregexpr`, `regexec` 的结果中提取匹配的字符串。

以原样匹配为例。

```
x <- c("New theme", "Old times", "In the present theme")
regexpr("the", x, perl=TRUE)
```

```
[1] 5 -1 4
attr(,"match.length")
[1] 3 -1 3
attr(,"index.type")
[1] "chars"
attr(,"useBytes")
[1] TRUE
```

这里使用了 `regexpr` 函数。`regexpr` 函数的一般用法为：

```
x <- c("New theme", "Old times", "In the present theme")
regexpr(pattern, text, ignore.case = FALSE, perl = FALSE,
 fixed = FALSE, useBytes = FALSE)
```



自变量为：

- `pattern` 是一个正则表达式，如果用了 `fixed=TRUE` 选项，则当作普通原样文本来匹配；
- `text` 是源字符串向量，要从其每个元素中查找 `pattern` 模式出现的位置；
- `ignore.case`：是否要忽略大小写匹配；
- `perl` 选择是否采用 perl 格式，如果不把 `pattern` 当作普通原样文本，应该选 `perl=TRUE`，perl 语言的正则表达式是事实上的标准，所以这样兼容性更好；
- `fixed` 当 `fixed=TRUE` 时 `pattern` 作为普通原样文本解释；
- `useBytes` 为 `TRUE` 时逐字节进行匹配，否则逐字符进行匹配。之所以有这样的区别，是因为有些编码中一个字符由多个字节构成，GBK 编码的汉字由两个字节组成，UTF-8 编码的汉字也是由两个字节构成。

`regexpr()` 函数返回一个整数值的向量，长度与 `text` 向量长度相同，结果的每个元素是在 `text` 的对应元素中 `pattern` 的首次匹配位置；没有匹配时结果元素取-1。结果会有一个 `match.length` 属性，表示每个匹配的长度，无匹配时取-1。

如果仅关心源字符串向量 `text` 中哪些元素能匹配 `pattern`，可以用 `grep` 函数，如

```
x <- c("New theme", "Old times", "In the present theme")
grep("the", x, perl=TRUE)
```

```
[1] 1 3
```

结果说明源字符串向量的三个元素中仅有第 1、第 3 号元素能匹配。如果都不匹配，返回 `integer(0)`。

`grep` 可以使用与 `regexpr` 相同的自变量，另外还可以加选项 `invert=TRUE`，这时返回的是不匹配的元素的下标。

`grep()` 如果添加选项 `value=TRUE`，则结果不是返回有匹配的元素的下标而是返回有匹配的元素本身（不是匹配的子串），如

```
x <- c("New theme", "Old times", "In the present theme")
grep("the", x, perl=TRUE, value=TRUE)
```

```
[1] "New theme" "In the present theme"
```

`grepl` 的作用与 `grep` 类似，但是其返回值是一个长度与源字符串向量 `text` 等长的逻辑型向量，每个元素的真假对应于源字符串向量中对应元素的匹配与否。如

```
x <- c("New theme", "Old times", "In the present theme")
grepl("the", x, perl=TRUE)
```

```
[1] TRUE FALSE TRUE
```

就像 `grep()` 与 `grepl()` 本质上给出相同的结果，只是结果的表示方式不同，`regexec()` 与 `regexpr()` 也给出仅在表示方式上有区别的结果。`regexpr()` 主要的结果是每个元素的匹配位置，用一个统一的属性返回各个匹配长度；`regexec()` 则返回一个与源字符串向量等长的列表，列表的每个元素为匹配的位置，并且列表的每个元素有匹配长度作为属性。所以，这两个函数只需要用其中一个就可以，下面仅使用 `regexpr()`。`regexec()` 的使用效果如

```
x <- c("New theme", "Old times", "In the present theme")
regexec("the", x, perl=TRUE)
```

```
[[1]]
[1] 5
attr(,"match.length")
[1] 3
attr(,"useBytes")
[1] TRUE
##
[[2]]
[1] -1
attr(,"match.length")
[1] -1
attr(,"useBytes")
[1] TRUE
##
[[3]]
[1] 4
```

```
attr(,"match.length")
[1] 3
attr(,"useBytes")
[1] TRUE
```

`grep()`, `grepl()`, `regexpr()`, `regexec()` 都只能找到源字符串向量的每个元素中模式的首次匹配, 不能找到所有匹配。`gregexpr()` 函数可以找到所有匹配。如

```
x <- c("New theme", "Old times", "In the present theme")
gregexpr("the", x, perl=TRUE)
```

```
[[1]]
[1] 5
attr(,"match.length")
[1] 3
attr(,"index.type")
[1] "chars"
attr(,"useBytes")
[1] TRUE
##
[[2]]
[1] -1
attr(,"match.length")
[1] -1
attr(,"index.type")
[1] "chars"
attr(,"useBytes")
[1] TRUE
##
[[3]]
[1] 4 16
attr(,"match.length")
[1] 3 3
attr(,"index.type")
```

```
[1] "chars"
attr(,"useBytes")
[1] TRUE
```

其结果是一个与源字符串向量等长的列表，格式与 `regexec()` 的结果格式类似，列表的每个元素对应于源字符串向量的相应元素，列表元素值为匹配的位置，并有属性 `match.length` 保存了匹配长度。匹配位置和匹配长度包含了所有的匹配，见上面例子中第三个元素的匹配结果。

函数 `grep`, `grep1` 结果仅给出每个元素能否匹配。`regexpr()`, `regexec()`, `gregexpr()` 则包含了匹配位置与匹配长度，这时，可以用 `regmatches()` 函数取出具体的匹配字符串。`regmatches()` 一般格式为

```
regmatches(x, m, invert = FALSE)
```

其中 `x` 是源字符串向量，`m` 是 `regexpr()`、`regexec()` 或 `gregexpr()` 的匹配结果。如

```
x <- c("New theme", "Old times", "In the present theme")
m <- regexpr("the", x, perl=TRUE)
regmatches(x, m)
```

```
[1] "the" "the"
```

可以看出，`regmatches()` 仅取出有匹配时的匹配内容，无匹配的内容被忽略。

取出多处匹配的例子如：

```
x <- c("New theme", "Old times", "In the present theme")
m <- gregexpr("the", x, perl=TRUE)
regmatches(x, m)
```

```
[[1]]
[1] "the"
##
[[2]]
character(0)
##
[[3]]
```

```
[1] "the" "the"
```

当 `regmatches()` 第二个自变量是 `gregexpr()` 的结果时，其输出结果变成一个列表，并且不再忽略无匹配的元素，无匹配元素对应的列表元素为 `character(0)`，即长度为零的字符型向量。对有匹配的元素，对应的列表元素为所有的匹配字符串组成的字符型向量。

实际上，如果 `pattern` 中没有正则表达式，`grep()`，`grepl()`，`regexpr()`，`gregexpr()` 中都可以用 `fixed=TRUE` 参数取代 `perl=TRUE` 参数，这时匹配总是解释为原样匹配，即使 `pattern` 中包含特殊字符也是进行原样匹配。

### 35.5.1 不区分大小写匹配

在基本 R 中，为了不区分大小写匹配，可以在 `grep` 等函数调用时加选项 `ignore.case=TRUE`；如

```
grep("Dr", c("Dr. Wang", "DR. WANG", "dR. W.R.))
```

```
[1] 1
```

```
grep("dr", c("Dr. Wang", "DR. WANG", "dR. W.R."), ignore.case=TRUE)
```

```
[1] 1 2 3
```

```
grep("(?i)dr", c("Dr. Wang", "DR. WANG", "dR. W.R.))
```

```
[1] 1 2 3
```

### 35.5.2 匹配单个字符

在模式中用 “.” 匹配任意一个字符（除了换行符“\n”，能否匹配此字符与选项有关）。如

```
s <- c("abc", "cabs", "lab")
mres <- regexpr("ab.", s, perl=TRUE); mres
```

```
[1] 1 2 -1
```

```
attr(,"match.length")
```

```
[1] 3 3 -1
attr(,"index.type")
[1] "chars"
attr(,"useBytes")
[1] TRUE
```

`regexpr` 仅给出每个元素中模式的首次匹配位置而不是给出匹配的内容。`regmatches` 函数以原始字符型向量和匹配结果为输入，结果返回每个元素中匹配的各个子字符串（不是整个元素），如：

```
regmatches(s, mres)
```

```
[1] "abc" "abs"
```

注意返回结果和输入字符型向量元素不是一一对应的，仅返回有匹配的结果。

像句点这样的字符称为元字符（meta characters），在正则表达式中有特殊函数。如果需要匹配句点本身，用“`[.]`”或者“`\.`”表示。比如，要匹配 `a.txt` 这个文件名，如下做法有错误：

```
grep("a.txt", c("a.txt", "a0txt"), perl=TRUE)
```

```
[1] 1 2
```

结果连 `a0txt` 也匹配了。用“`[.]`”表示句点则将句点不做特殊解释：

```
grep("a[.]txt", c("a.txt", "a0txt"), perl=TRUE)
```

```
[1] 1
```

```
grep("a\\.txt", c("a.txt", "a0txt"), perl=TRUE)
```

```
[1] 1
```

注意在 R 语言字符型常量中一个 `\` 需要写成两个。

如果仅需按照原样进行查找，也可以在 `grep()`、`grep1()`、`regexpr()`、`gregexpr()` 等函数中加选项 `fixed=TRUE`，这时不要再用 `perl=TRUE` 选项。如

```
grep("a.txt", c("a.txt", "a0txt"), fixed=TRUE)
```

```
[1] 1
```

### 35.5.3 匹配一组字符中的某一个

模式 “[ns]a.[.]xls” 表示匹配的第一个字符是 n 或 s，第二个字符是 a，第三个字符任意，第四个字符是句点，然后是 xls。例：

```
regexpr("[ns]a.[.]xls", c("sa1.xls", "dna2.xlss", "na3.xls"), perl=T)
```

```
[1] 1 2 1
attr(,"match.length")
[1] 7 7 7
attr(,"index.type")
[1] "chars"
attr(,"useBytes")
[1] TRUE
```

### 35.5.4 原样匹配元字符

例：

```
grep("int x\\[5\\]", c("int x;", "int x[5]"), perl=TRUE)
```

```
[1] 2
```

也可以用 “[[]” 表示 “[”，用 “[”]” 表示 “[”]”，如

```
grep("int x[[5]]", c("int x;", "int x[5]"), perl=TRUE)
```

```
[1] 2
```

### 35.5.5 匹配数字

例：

```
grep("n\\d[.]xls", c("n1.xls", "na.xls"), perl=TRUE)
```

```
[1] 1
```

### 35.5.6 匹配开头和末尾

例:

```
grep("^n\\d[.]xls$", c("n1.xls", "na.xls", "cn1.xls", "n1.xlsx"), perl=TRUE)

[1] 1

grep("\\An\\d[.]xls\\Z", c("n1.xls", "na.xls", "cn1.xls", "n1.xlsx"), perl=TRUE)

[1] 1
```

只匹配了第一个输入字符串。

### 35.5.7 匹配字母、数字、下划线

例:

```
m <- regexpr("s\\w[.]", c("file-s1.xls", "s#.xls"), perl=TRUE)
regmatches(c("file-s1.xls", "s#.xls"), m)

[1] "s1."
```

可以看出，模式匹配了 `s1.` 而没有匹配 `s#.`。

### 35.5.8 十六进制和八进制数

例如

```
gregexpr("\\x0A", "abc\\nefg\\n")[[1]]

[1] 4 8
attr(,"match.length")
[1] 1 1
attr(,"index.type")
[1] "chars"
attr(,"useBytes")
[1] TRUE
```



匹配了两个换行符。

### 35.5.9 POSIX 字符类

例如：

```
grep("[[:alpha:]_][[:alnum:]_]", c("x1", "_x", ".x", ".1"))
```

```
[1] 1 2 3 4
```

#### 35.5.10 加号重复匹配

例

```
s <- c("sa1", "dsa123")
mres <- regexpr("sa[[:digit:]]+", s, perl=TRUE)
regmatches(s, mres)
```

```
[1] "sa1" "sa123"
```

例如：

```
p <- "^[[[:alnum:]]_]+@[[:alnum:]]_+[.][[:alnum:]]_+ $"
x <- "abc123@efg.com"
m <- regexpr(p, x, perl=TRUE)
regmatches(x, m)
```

```
[1] "abc123@efg.com"
```

匹配的电子邮件地址在 @ 前面可以使用任意多个字母、数字、下划线，在 @ 后面由小数点分成两段，每段可以使用任意多个字母、数字、下划线。这里用了 ^ 和 \$ 表示全字符串匹配。

#### 35.5.11 星号和问号重复匹配

`^https?:/[[:alnum:]]./]+$` 可以匹配 http 或 https 开始的网址。如

```
s <- c("http://www.163.net", "https://123.456.")
grep("^https?:/[[:alnum:]_./]+$", s, perl=TRUE)
```

```
[1] 1 2
```

(注意第二个字符串不是合法网址但是按这个正则表达式也能匹配)

### 35.5.12 计数重复

例:

```
grep("[[:digit:]]{3}", c("1", "12", "123", "1234"))
```

```
[1] 3 4
```

模式匹配的是三位的数字。

日期匹配例:

```
pat <- paste(
 c("[[:digit:]]{1,2}[-/]",
 "[[:digit:]]{1,2}[-/]",
 "[[:digit:]]{2,4}"), collapse="")
grep(pat, c("2/4/1998", "13/15/198"))
```

```
[1] 1 2
```

### 35.5.13 贪婪匹配和懒惰匹配

例如:

```
s <- "1st other 2nd"
p1 <- "<[Bb]>. *</[Bb]>"
m1 <- regexpr(p1, s, perl=TRUE)
regmatches(s, m1)[[1]]
```

```
[1] "1st other 2nd"
```

我们本来期望的是提取第一个 “<B>.....</B>” 组合，不料提取了两个 “<B>.....</B>” 组合以及中间的部分。

比如，上例中模式修改后得到了期望的结果：

```
s <- "1st other 2nd"
p2 <- "<[Bb]>.*?</[Bb]>"
m2 <- regexpr(p2, s, perl=TRUE)
regmatches(s, m2)[[1]]

[1] "1st"
```

### 35.5.14 单词边界

例：

```
grep("\\bcat\\b", c("a cat meaos", "the category"))

[1] 1
```

### 35.5.15 句点全匹配与多行模式

句点通配符一般不能匹配换行，如

```
s <- "1st\n\n"
grep("<[Bb]>.*?</[Bb]>", s, perl=TRUE)

integer(0)
```

跨行匹配失败。而在 HTML 的规范中换行是正常的。一种办法是预先用 `gsub` 把所有换行符替换为空格。但是这只能解决部分问题。

另一方法是在 Perl 正则表达式开头添加 `(?s)` 选项，这个选项使得句点通配符可以匹配换行符。如

```
s <- "1st\n\n"
mres <- regexpr("(?s)<[Bb]>.*?</[Bb]>", s, perl=TRUE)
regmatches(s, mres)
```

```
[1] "1st\n"
```

多行模式例:

```
s <- "1st\n\n"
mres1 <- gregexpr("^<.+?>", s, perl=TRUE)
mres2 <- gregexpr("(?m)^<.+?>", s, perl=TRUE)
regmatches(s, mres1)[[1]]
```

```
[1] ""
```

```
regmatches(s, mres2)[[1]]
```

```
[1] "" ""
```

字符串 `s` 包含两行内容，中间用 `\n` 分隔。`mres1` 的匹配模式没有打开多行选项，所以模式中的 `^` 只能匹配 `s` 中整个字符串开头。`mres2` 的匹配模式打开了多行选项，所以模式中的 `^` 可以匹配 `s` 中每行的开头。

### 35.5.16 备择模式

例如，某人的名字用 James 和 Jim 都可以，表示为 `James|Jim`，如

```
s <- c("James, Bond", "Jim boy")
pat <- "James|Jim"
mres <- gregexpr(pat, s, perl=TRUE)
regmatches(s, mres)
```

```
[[1]]
```

```
[1] "James"
```

```
##
```

```
[[2]]
```

```
[1] "Jim"
```

### 35.5.17 分组与捕获

例：希望把 “<B>.....</B>” 两边的 “<B>” 和 “</B>” 删除，可以用如下的替换方法：

```
x <- "1st other 2nd"
pat <- "(?s)<[Bb]>(.*?)</[Bb]>"
repl <- "\\1"
gsub(pat, repl, x, perl=TRUE)
```

```
[1] "1st other 2nd"
```

替换模式中的\1(写成 R 字符型常量时\要写成\\) 表示第一个圆括号匹配的内容，但是表示选项的圆括号 ((?s)) 不算在内。

例：希望把带有前导零的数字的前导零删除，可以用如

```
x <- c("123", "0123", "00123")
pat <- "\\b0+([1-9][0-9]*)\\b"
repl <- "\\1"
gsub(pat, repl, x, perl=TRUE)
```

```
[1] "123" "123" "123"
```

其中的\b 模式表示单词边界，这可以排除在一个没有用空格或标点分隔的字符串内部拆分出数字的情况。

例：为了交换纵横坐标，可以用如下替换

```
s <- "1st: (5,3.6), 2nd: (2.5, 1.1)"
pat <- paste0(
 "[]([[:digit:]]+),",
 "[[:space:]]*([[:digit:]]+) []")
repl <- "(\\2, \\1)"
gsub(pat, repl, s, perl=TRUE)
```

```
[1] "1st: (3.6, 5), 2nd: (1.1, 2.5)"
```

例如，要匹配 yyyy-mm-dd 这样的日期，并将其改写为 mm/dd/yyyy，就可以用这样的替换模式：

```
pat <- "[0-9]{4}-([0-9]{1,2})-([0-9]{1,2})"
repl <- "\\2/\\3/\\1"
gsub(pat, repl, c("1998-05-31", "2017-01-14"))
```

```
[1] "05/31/1998" "01/14/2017"
```

分组除了可以做替换外，还可以用来表示模式中的重复出现内容。例如，`[a-z]{3}\\1` 这样的模式可以匹配如 `abcabc`, `uxzuxz` 这样的三字母重复。如

```
grep("[a-z]{3}\\1", c("abcabc", "aabbcc"))
```

```
[1] 1
```

又例如，下面的程序找出了年（后两位）、月、日数字相同的日期：

```
x <- c("2008-08-08", "2017-01-18")
m <- regexpr("\\d\\d(\\d\\d)-\\1-\\1", x)
regmatches(x, m)
```

```
[1] "2008-08-08"
```

下面是一个非捕获分组示例。设需要把 1921-2020 之间的世纪号删去，可以用

```
pat <- "\\A(?:19|20)([0-9]{2})\\Z"
repl <- "\\1"
x <- c("1978", "2017", "2035")
gsub(pat, repl, x, perl=TRUE)
```

```
[1] "78" "17" "35"
```

其中用了非捕获分组使得备择模式 `19|20` 优先匹配。注意模式并没有能保证日期在 1921-2020 之间。更周密的程序可以写成：

```
x <- c("1978", "2017", "2035")
p1 <- "\\A19(2[1-9] | [3-9][0-9])\\Z"
r1 <- "\\1"
p2 <- "\\A20([01][0-9] | 20)\\Z"
x <- gsub(p1, r1, x, perl=TRUE)
x <- gsub(p2, r1, x, perl=TRUE)
```

```
x
[1] "78" "17" "2035"
```

## 35.6 正则表达式应用例子

### 35.6.1 数据预处理

在原始数据中，经常需要审核数据是否合法，已经把一些常见错误输入自动更正。这都可以用正则表达式实现。

#### 35.6.1.1 除去字符串开头和结尾的空格

函数 `stringr::str_trim()` 和 `trimws()` 可以除去字符串开头与结尾的空格，也可以仅除去开头或仅除去结尾的空格。

这个任务如果用正则表达式字符串替换函数来编写，可以写成：

```
把字符串向量 x 的元素去除首尾的空白。
strip <- function(x){
 x <- str_replace_all(x, "^[:space:]+", "")
 x <- str_replace_all(x, "[:space:]+$", "")
 x
}
```

或者

```
把字符串向量 x 的元素去除首尾的空白。
strip <- function(x){
 x <- gsub("^[:space:]+", "", x, perl=TRUE)
 x <- gsub("[:space:]+$", "", x, perl=TRUE)
 x
}
```

这个版本可以除去包括空格在内的所有首尾空白字符。

### 35.6.1.2 除去字符串向量每个元素中所有空格

```
compress <- function(x){
 str_replace_all(x, " ", "")
}
```

或者

```
compress <- function(x){
 gsub(" ", "", x, fixed=TRUE)
}
```

这可以解决"李明"与"李 明"不相等这样的问题。类似的程序也可以用来把中文的标点替换成英文的标点。

### 35.6.1.3 判断日期是否合法

设日期必须为 yyyy-mm-dd 格式，年的数字可以是两位、三位、四位，程序为：

```
is_yyyymmdd <- function(x){
 pyear <- "[0-9]{2}|[1-9][0-9]{2}|[1-2][0-9]{3}"
 pmon <- "[1-9]|0[1-9]|1[0-2]"
 pday <- "[1-9]|0[1-9]|1[0-9]|2[0-9]|3[01]"
 pat <- paste("\\A", pyear, "-", pmon, "-", pday, "\\Z", sep="")
 str_detect(x, pat)
}
```

或

```
is.yyyymmdd <- function(x){
 pyear <- "[0-9]{2}|[1-9][0-9]{2}|[1-2][0-9]{3}"
 pmon <- "[1-9]|0[1-9]|1[0-2]"
 pday <- "[1-9]|0[1-9]|1[0-9]|2[0-9]|3[01]"
 pat <- paste("\\A", pyear, "-", pmon, "-", pday, "\\Z", sep="")
 grepl(pat, x, perl=TRUE)
}
```



这样的规则还没有排除诸如 9 月 31 号、2 月 30 号这样的错误。

例:

```
x <- c("49-10-1", "1949-10-01", "532-3-15", "2015-6-1",
 "2017-02-30", "2017-13-11", "2017-1-32")
is_yyyymmdd(x)
```

```
[1] TRUE TRUE TRUE TRUE TRUE FALSE FALSE
```

注意错误的 2 月 30 号没有识别出来。

#### 35.6.1.4 把字符型日期变成 yyyy-mm-dd 格式。

```
make_date <- function(x){
 x %>%
 str_trim() %>%
 str_replace_all("[[:space:]]+", "-") %>%
 str_replace_all("/", "-") %>%
 str_replace_all("[.]", "-") %>%
 str_replace_all("^([0-9]{2})(-[0-9]{1,2})-([0-9]{1,2})$", "20\\1\\2") %>%
 str_replace_all("^([0-9]{4})-([0-9])-([0-9]{1,2})$", "\\1-0\\2-\\3") %>%
 str_replace_all("^([0-9]{4})-([0-9]{2})-([0-9])$", "\\1-0\\2")
}
```

或者

```
make.date <- function(x){
 x <- trimws(x)
 x <- gsub("[[:space:]]+", "-", x)
 x <- gsub("/", "-", x)
 x <- gsub("[.]", "-", x)
 x <- gsub("^([0-9]{2})(-[0-9]{1,2})-([0-9]{1,2})$", "20\\1\\2", x)
 x <- gsub("^([0-9]{4})-([0-9])-([0-9]{1,2})$", "\\1-0\\2-\\3", x)
 x <- gsub("^([0-9]{4})-([0-9]{2})-([0-9])$", "\\1-0\\2", x)
}
```

```
x
}
```

另一办法是用 `strsplit()` 拆分成三个部分，转换为整数，再转换回字符型。

测试:

```
x <- c("49/10/1", "1949.10.01", "532 3 15", "2015/6.1",
 "20170230", "2017.13/11", "2017 1 32")
make_date(x)
```

```
[1] "2049-10-01" "1949-10-01" "532-3-15" "2015-06-01" "20170230"
[6] "2017-13-11" "2017-01-32"
```

目前的函数还不能处理没有分隔符的情况，也不能验证日期的合法性。

### 35.6.1.5 合并段落为一行

在某些纯文本格式中，各段之间用空行分隔，没有用空行分隔的各行看成同一段。如下的函数把其中的不表示分段的换行删去从而合并这些段落。函数以一个文件名作为输入，合并段落存回原文件。注意，这样修改文件的函数在调试时，应该注意先备份文件，等程序没有任何错误以后才可以忽略备份。

```
combine_paragraph <- function(fname){
 lines <- readLines(fname)

 s <- str_flatten(lines)
 s <- s %>%
 str_replace_all("^[:space:]+\n", "\n") %>%
 str_replace_all("(^[^\n]+\n)", "\\1 ") %>%
 str_replace_all("(^[^\n]+\n)", "\\1\n\n")

 writeLines(str_split(s, "\n")[[1]],
 con=fname)
}
```

函数首先把仅有空格的行中的空格删除，将有内容的行的行尾换行符替换成一

个空格，再把剩余的有内容的行的行尾换行符多加一个换行符。

上面的程序中特意用了基本 R 的 `readLines()` 函数而不是 `readr` 包的 `read_lines()` 函数，因为 `readLines()` 使用操作系统的默认中文编码，而 `read_lines()` 默认使用 UTF-8 编码，需要用选项 `locale=locale(encoding="GB18030")` 才能在中文 MS Windows 中正确读取中文文件。

下面的版本不使用 `stringr`:

```
combine.paragraph <- function(fname){
 lines <- readLines(fname)
 s <- paste(lines, collapse="\n")
 s <- gsub("[[:space:]]+\n", "\n", s, perl=TRUE)
 s <- gsub("([^\n]+\n)", "\\1 ", s, perl=TRUE)
 s <- gsub("([^\n]+\n)", "\\1\n\n", s, perl=TRUE)
 writeLines(strsplit(s, "\n", fixed=TRUE)[[1]],
 con=fname)
}
```

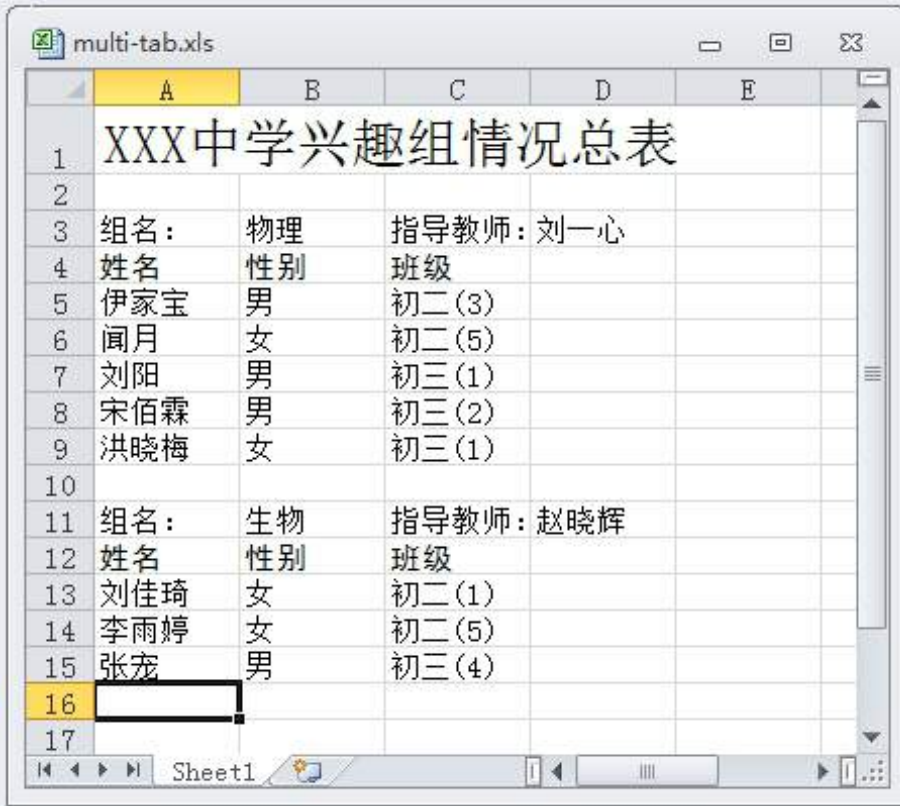
### 35.6.2 不规则 Excel 文件处理

- 作为字符型数据处理示例，考察如下的一个 Excel 表格数据。

假设一个中学把所有课外小组的信息汇总到了 Excel 表的一个工作簿中。每个课外小组占一块区域，各小组上下排列，但不能作为一个数据框读取。下图为这样的文件的一个简化样例：

实际数据可能有很多个小组，而且数据是随时更新的，所以复制粘贴另存的方法不太可行，需要用一个通用的程序处理。Excel 文件 (.xls 后缀或.xlsx 后缀) 不是文本型数据。在 Excel 中，用“另存为”把文件保存为 CSV 格式，内容如下：

```
XXX中学兴趣组情况总表,,
,,
组名: ,物理,指导教师: ,刘一心
姓名,性别,班级,
```



	A	B	C	D	E
1	XXX中学兴趣组情况总表				
2					
3	组名:	物理	指导教师:	刘一心	
4	姓名	性别	班级		
5	伊家宝	男	初二(3)		
6	闻月	女	初二(5)		
7	刘阳	男	初三(1)		
8	宋佰霖	男	初三(2)		
9	洪晓梅	女	初三(1)		
10					
11	组名:	生物	指导教师:	赵晓辉	
12	姓名	性别	班级		
13	刘佳琦	女	初二(1)		
14	李雨婷	女	初二(5)		
15	张宠	男	初三(4)		
16					
17					

图 35.1: 不规则 Excel 文件样例图形

伊家宝,男,初二(3),  
 闻月,女,初二(5),  
 刘阳,男,初三(1),  
 宋佰霖,男,初三(2),  
 洪晓梅,女,初三(1),  
 ,,,  
 组名: ,生物,指导教师: ,赵晓辉  
 姓名,性别,班级,  
 刘佳琦,女,初二(1),  
 李雨婷,女,初二(5),  
 张宠,男,初三(4),

生成测试用的数据文件:

```
demo.multitab.data <- function(){
 s <- "
 XXX 中学兴趣组情况总表,,,
 ,,,
 组名: , 物理, 指导教师: , 刘一心
 姓名, 性别, 班级,
 伊家宝, 男, 初二 (3),
 闻月, 女, 初二 (5),
 刘阳, 男, 初三 (1),
 宋佰霖, 男, 初三 (2),
 洪晓梅, 女, 初三 (1),
 ,,,
 组名: , 生物, 指导教师: , 赵晓辉
 姓名, 性别, 班级,
 刘佳琦, 女, 初二 (1),
 李雨婷, 女, 初二 (5),
 张宠, 男, 初三 (4),
 "
 writeLines(s, "multitab.csv")
}
demo.multitab.data()
```

读入测试用的数据，转换为一整个数据框：

```
demo_multitab <- function(){
 ## 读入所有行
 lines <- readLines("multitab.csv")

 ## 去掉首尾空格
 lines <- str_trim(lines)

 ## 删去所有空行和只有逗号的行
 empty <- str_detect(lines, "^[:,space:],*$")
 if(length(empty)>0){
 lines <- lines[!empty]
 }

 ## 找到所有包含“组名:”的行对应的行号
 heads <- str_which(lines, " 组名: ")

 ## 定位每个表的开始行和结束行（不包括组名和表头所在的行）
 start <- heads + 2
 end <- c(heads[-1]-1, length(lines))
 ngroups <- length(heads)

 ## 先把数据读入一个列表。
 for(ii in seq(ngroups)){
 ## 组名和指导教师所在行:
 line <- lines[heads[ii]]
 v <- str_split(line, ",")[[1]]
 ## 组名: v[2] 指导教师: v[4]

 ## 将表格内容各行合并成一个用换行符分隔的长字符串,
 ## 然后变成可读取的文件
 s <- str_flatten(lines[start[ii]:end[ii]], collapse="\n")
 }
}
```

```

con <- textConnection(s, "rt")
da1 <- read.csv(
 con, header=FALSE,
 colClasses=c("character", "character", "character", "NULL"))
close(con)
names(da1) <- c("姓名", "性别", "班级")
da1 <- cbind("组名"=v[2], "指导教师"=v[4], da1)

if(ii==1) {
 da <- da1
} else {
 da <- rbind(da, da1)
}
}

da
}
da <- demo_multitab()
knitr::kable(da)

```

组名	指导教师	姓名	性别	班级
物理	刘一心	伊家宝	男	初二 (3)
物理	刘一心	闻月	女	初二 (5)
物理	刘一心	刘阳	男	初三 (1)
物理	刘一心	宋佰霖	男	初三 (2)
物理	刘一心	洪晓梅	女	初三 (1)
生物	赵晓辉	刘佳琦	女	初二 (1)
生物	赵晓辉	李雨婷	女	初二 (5)
生物	赵晓辉	张宠	男	初三 (4)

不使用 stringr 的版本:

```

demo_multitab <- function(){
 ## 读入所有行
 lines <- readLines("multitab.csv")

```

```
去掉首尾空格
lines <- trimws(lines)

删去所有空行和只有逗号的行
(1) 不用正则表达式做法
empty <- which(lines == "" | substring(lines, 1, 3)=="",",",")
(2) 用正则表达式做法:
empty <- grep("^[:space:]*$", lines)
if(length(empty)>0){
 lines <- lines[-empty]
}

找到所有包含 “组名:” 的行对应的行号
heads <- grep(" 组名: ", lines, fixed=TRUE)

定位每个表的开始行和结束行 (不包括组名和表头所在的行)
start <- heads + 2
end <- c(heads[-1]-1, length(lines))
ngroups <- length(heads)

先把数据读入一个列表。
for(ii in seq(ngroups)){
 ## 组名和指导教师所在行:
 line <- lines[heads[ii]]
 v <- strsplit(line, ",")[[1]]
 ## 组名: v[2] 指导教师: v[4]

 ## 将表格内容各行合并成一个用换行符分隔的长字符串,
 ## 然后变成可读取的文件
 s <- paste(lines[start[ii]:end[ii]], collapse="\n")
 con <- textConnection(s, "rt")
 da1 <- read.csv(
```



```

 con, header=FALSE,
 colClasses=c("character", "character", "character", "NULL"))
close(con)
names(da1) <- c(" 姓名", " 性别", " 班级")
da1 <- cbind(" 组名"=v[2], " 指导教师"=v[4], da1)

if(ii==1) {
 da <- da1
} else {
 da <- rbind(da, da1)
}
}

da
}
da <- demo.multitab()
da

```

在程序中，用 `readLines` 函数读取文本文件各行到一个字符型向量。用 `grep` 可以找到每个小组开头的行（有“组名：”的行）。然后可以找出每个小组学生名单的开始行号和结束行号。各小组循环处理，读入后每个小组并入结果数据框中。用 `strsplit` 函数拆分用逗号分开的数据项。用 `textConnection` 函数可以把一个字符串当作文件读取，这样 `read.csv` 函数可以从一个字符串读入数据。

### 35.6.3 网站数据获取

很多网站定期频繁发布数据，所以传统的手工复制粘贴整理是不现实的。

但是，这些数据网页往往有固定模式，如果网页不是依赖 JavaScript 来展示的话，可以读取网页然后通过字符型数据处理方法获得数据。

`url` 可以打开一个网址链接然后用 `readLines` 或者 `readr::read_lines()` 函数读取。用 `str_replace_all()` 或者 `gsub()` 去掉不需要的成分。用 `str_which()` 或者 `grep` 查找关键行。具体例子略。

### 35.6.4 字频统计

正则表达式中的字符类 `[:alpha:]` 指的是当前系统中的字母，所以在中文环境中的中文字也是字母，但中文标点不算。下面是《红楼梦》中“秋窗风雨夕”的文本：

秋花惨淡秋草黄，耿耿秋灯秋夜长。  
已觉秋窗秋不尽，那堪风雨助凄凉！  
助秋风雨来何速！惊破秋窗秋梦绿。  
抱得秋情不忍眠，自向秋屏移泪烛。  
泪烛摇摇蕙短檠，牵愁照恨动离情。  
谁家秋院无风入？何处秋窗无雨声？  
罗衾不奈秋风力，残漏声催秋雨急。  
连宵脉脉复飏飏，灯前似伴离人泣。  
寒烟小院转萧条，疏竹虚窗时滴沥。  
不知风雨几时休，已教泪洒窗纱湿。

希望统计每个字的出现次数，并显示频数前十的字。设变量 `poem_autumnwindow` 中包含了上述诗词的文本。

首先用 `str_extract_all()` 提取每个中文字，组成一个字符型向量：

```
words_vec <- str_extract_all(poem_autumnwindow, "[[:alpha:]]")[[1]]
head(words_vec)
```

```
[1] "秋" "花" "惨" "淡" "秋" "草"
```

用 `table()` 函数计算频数，并按频数排序，输出前 10 结果：

```
words_freq <- sort(table(words_vec), decreasing=TRUE)
knitr::kable(head(words_freq, 10))
```

words_vec	Freq
秋	15
窗	5
风	5
雨	5
不	4
泪	3
灯	2
耿	2
何	2
离	2

### 35.6.5 数字验证

#### 35.6.5.1 整数

字符串完全为十进制正整数的模式，写成 R 字符型常量：

```
"\\A[0-9]+\\Z"
```

这个模式也允许正整数以 0 开始，如果不允许以零开始，可以写成

```
"\\A[1-9][0-9]*\\Z"
```

对于一般的整数，字符串完全为十进制整数，但是允许前后有空格，正负号与数字之间允许有空格，模式可以写成：

```
"\\A[]*[+-]?[]*[1-9][0-9]*\\Z"
```

#### 35.6.5.2 十六进制数字

字符串仅有十六进制数字，模式写成 R 字符型常量为

```
"\\A[0-9A-Fa-f]+\\Z"
```

在文中匹配带有 0x 前缀的十六进制数字，模式为

```
"\\b0x[0-9A-Fa-f]+\\b"
```

### 35.6.5.3 二进制数字

为了在文中匹配一个以 B 或 b 结尾的二进制非负整数，可以用

```
"\\b[01]+[Bb]\\b"
```

### 35.6.5.4 有范围的整数

1-12:

```
"\\b1[012] | [1-9]\\b"
```

1-24:

```
"\\b2[0-4] | 1[0-9] | [1-9]\\b"
```

1-31:

```
"\\b3[01] | [12][0-9] | [1-9]\\b"
```

1900-2099:

```
"\\b(?:19|20)[0-9]{2}\\b"
```

这里的分组仅用于在 19 和 20 之间选择，不需要捕获，所以用了 (?: 的非捕获分组格式。

### 35.6.5.5 判断字符型向量每个元素是否数值

如下的 R 函数用了多种数字的正则表达式来判断字符型向量每个元素是否合法数值。

```
all_numbers <- function(x){
 x <- x %>%
 str_replace_all("\\A[]+", "") %>%
 str_replace_all("[]+\\Z", "")
```





## Part IX

# 用 Rcpp 连接 C++ 代码





## Chapter 36

# Rcpp 介绍

为了提高 R 程序的运行效率，可以尽量使用向量化编程，减少循环，尽量使用内建函数。对于效率的瓶颈，尤其是设计迭代算法时，可以采用编译代码，而 Rcpp 扩展包可以很容易地将 C++ 代码连接到 R 程序中，并且支持在 C++ 中使用类似于 R 的数据类型。

没有学过 C++ 语言的读者，如果需要编写比较独立的不太依赖于 R 的已有功能的算法，可以考虑学习使用 Julia 语言编写。Julia 语言是最近几年才发明的一种比 R 更现代、理念更先进的程序语言，其运行效率一般比 R 高得多，经常接近编译代码的效率。

Rcpp 可以很容易地把 C++ 代码与 R 程序连接在一起，可以从 R 中直接调用 C++ 代码而不需要用户关心那些繁琐的编译、链接、接口问题。可以在 R 数据类型和 C++ 数据类型之间容易地转换。

因为涉及到编译，所以 Rcpp 比一般的扩展包有更多的安装要求：除了要安装 Rcpp 包之外，MS Windows 用户还需要安装 RTools 包，这是用于 C, C++, Fortran 程序编译链接的开发工具包，是自由软件。用户的应用程序路径 (PATH) 中必须有 RTools 包可执行程序的路径 (安装 RTools 可以自动设置)。如果 Rcpp 不能找到编译器，可以把编译器安装到 Rcpp 默认的位置。Mac 操作系统和 Linux 操作系统中可以用操作系统自带的编译器。

在 MS Windows 系统中，从CRAN网站下载 RTools 工具包，并将其安装到默认位置 `C:\RTools` 中，否则 RStudio 中使用 Rcpp 可能会出错。注意，RStudio

自己的自动安装 RTools 的功能有问题，安装后不能正常编译含有 Rcpp 的 Rmd 文件。

Rcpp 支持把 C++ 代码写在 R 源程序文件内，执行时自动编译连接调用；也支持把 C++ 代码保存在单独的源文件中，执行 R 程序时自动编译连接调用；对较复杂的问题，应制作 R 扩展包，利用构建 R 扩展包的方法实现 C++ 代码的编译连接，这时接口部分也可以借助 Rcpp 属性功能或模块功能完成。

## 36.1 Rcpp 的用途

- 把已经用 R 代码完成的程序中运行速度瓶颈部分改写成 C++ 代码，提高运行效率。
- 对于 C++ 或 C 程序源代码或二进制代码提供的函数库，可以用 Rcpp 编写 C++ 界面程序进行 R 与 C++ 程序的输入、输出的传送，并在 C++ 界面程序中调用外来的函数库。
- 注意，用 Rcpp 编写 C++ 程序，不利于把程序脱离 R 运行或被其他的 C++ 程序调用。当然，可以只把 Rcpp 作为界面，主要的算法引擎完全不用 Rpp 的数据类型。
- RInside 扩展包支持把 R 嵌入到 C++ 主程序中。

## 36.2 Rcpp 入门样例

### 36.2.1 用 `cppFunction()` 转换简单的 C++ 函数—Fibonacci 例子

考虑用 C++ 程序计算 Fibonacci 数的问题。Fibonacci 数满足  $f_0 = 0, f_1 = 1, f_t = f_{t-1} + f_{t-2}$ 。

可以使用如下 R 代码，其中有一部分 C++ 代码，用 `cppFunction` 转换成了 R 可以调用的同名 R 函数。

```
library(Rcpp)
cppFunction(code='
 int fibonacci(const int x){
 if(x < 2) return x;
 else
 return (fibonacci(x-1) + fibonacci(x-2));
 }
')
print(fibonacci(5))
[1] 5
```

编译、链接、导入是在后台由 Rcpp 控制自动进行的，不需要用户去设置编译环境，也不需要用户执行编译、链接、导入 R 的工作。在没有修改 C++ 程序时，同一 R 会话期间重新运行不必重新编译。

上面的 Fibonacci 函数仅接受标量数值作为输入，不允许向量输入。

从算法角度评价，这个算法是极其低效的，其算法规模是  $O(2^n)$ ， $n$  是自变量值。

### 36.2.2 用 sourceCpp() 转换 C++ 程序——正负交替迭代例子

设  $x_t, t = 1, 2, \dots, n$  保存在 R 向量  $x$  中，令

$$y_1 = x_1$$

$$y_t = (-1)^t y_{t-1} + x_t, \quad t = 2, \dots, n$$

希望用 C++ 函数对输入序列  $x$  计算输出  $y$ ，并用 R 调用这样的函数。

下面的程序用 R 函数 `sourceCpp()` 把保存在 R 字符串中的 C++ 代码编译并转换为同名的 R 函数。

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
```

```
NumericVector iters(NumericVector x){
 int n = x.size();
 NumericVector y(n);

 y[0] = x[0];
 double sign=-1;
 for(int t=1; t<n; t++){
 sign *= -1;
 y[t] = sign*y[t-1] + x[t];
 }

 return y;
}
')
print(iters(1:5))
[1] 1 3 0 4 1
```

这个例子说明 C++ 程序可以直接写在 R 程序文件内 (保存为 R 多行字符串), 用 `sourceCpp()` 函数编译。

Rcpp 包为 C++ 提供了一个 `NumericVector` 数据类型, 用来存储数值型向量。用成员函数 `size()` 访问其大小, 用方括号下标访问其元素。

C 程序中定义的函数可以返回 `NumericVector` 数据类型, 将自动转换为 R 的数值型向量。

特殊的注释 `//[[Rcpp::export]]` 用来指定哪些 C++ 函数是要转换为 R 函数的。这叫做 Rcpp 属性 (attributes) 功能。

### 36.2.3 用 `sourceCpp()` 转换 C++ 源文件中的程序—正负交替迭代例子

直接把 C++ 代码写在 R 源程序内部的好处是不用管理多个源文件, 缺点是当 C++ 代码较长时, 不能利用专用 C++ 编辑环境和调试环境, 出错时显示的错误行号不好定位, 而且把代码保存在 R 字符串内, C++ 代码中用到字符时需

要特殊语法。所以，稍复杂的 C++ 代码应该放在单独的 C++ 源文件内。

假设上面的 `iters` 函数的 C++ 代码单独存入了一个 `iters.cpp` 源文件中，内容如下：

```
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector iters(NumericVector x){
 int n = x.size();
 NumericVector y(n);

 y[0] = x[0];
 double sign=-1;
 for(int t=1; t<n; t++){
 sign *= -1;
 y[t] = sign*y[t-1] + x[t];
 }

 return y;
}
```

用如下的 `sourceCpp()` 函数把 C++ 源文件中代码编译并转换为 R 可访问的同名函数，测试调用：

```
sourceCpp(file='iters.cpp')
print(iters(1:5))
[1] 1 3 0 4 1
```

### 36.2.4 用 sourceCpp() 转换 C++ 源程序文件—卷积例子

考虑向量  $\mathbf{x} = (x_0, x_1, \dots, x_{n-1})$ ,  $\mathbf{y} = (y_0, y_1, \dots, y_{m-1})$ , 定义  $\mathbf{x}$  与  $\mathbf{y}$  的离散卷积  $\mathbf{z} = (z_0, z_1, \dots, z_{n+m-2})$  为

$$z_k = \sum_{(i,j): i+j=k} x_i y_j = \sum_{i=\max(0, k-m+1)}^{\min(k, n)} x_i y_{k-i}.$$

假设如下的 C++ 程序保存到了源文件 conv.cpp 中。

```
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector convolveCpp(
 NumericVector a, NumericVector b){
 int na = a.size(), nb = b.size();
 int nab = na + nb - 1;
 NumericVector xab(nab);

 for(int i=0; i < na; i++)
 for(int j=0; j < nb; j++)
 xab[i+j] += a[i] * b[j];

 return xab;
}
```

如下的代码用 sourceCpp() 函数把上面的 C++ 源文件 conv.cpp 自动编译并转换为同名的 R 函数, 进行测试:

```
sourceCpp(file='conv.cpp')
print(convolveCpp(1:5, 1:3))
[1] 1 4 10 16 22 22 15
```

### 36.2.5 在 Rmd 文件中使用 C++ 源程序文件

在 Rmd 文件中，可以使用特殊的 Rcpp 代码块，其中包含 C++ 源代码，可以直接起到对其调用 `sourceCpp()` 的作用。

但是，目前在 MS Windows 系统下仍有错误，需要将 RTools 安装在默认路径，即 `C:\RTools\bin` 与 `C:\RTools\mingw_64\bin` 中，并可能需要将这两个路径人为地添加到系统的 PATH 变量中，办法是“Windows 程序–Windows 系统–控制面板–系统和安全–系统–高级系统设置–高级–环境变量–系统变量”，其中的 PATH 变量是用分号分开的可执行程序搜索路径，将其中添加 RTools 包的 `C:\RTools\bin` 与 `C:\RTools\mingw_64\bin`。

设置好以后，在 Rmd 文件中输入如下的代码块：

```
```{Rcpp}
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector convolveCpp(
    NumericVector a, NumericVector b){
    int na = a.size(), nb = b.size();
    int nab = na + nb - 1;
    NumericVector xab(nab);

    for(int i=0; i < na; i++)
        for(int j=0; j < nb; j++)
            xab[i+j] += a[i] * b[j];

    return xab;
}
...
```
```

可以使得 C++ 函数 `convolveCpp()` 被编译并转换成 R 可以调用的同名函数，调用如：

```
convolveCpp(1:5, 1:3)
[1] 1 4 10 16 22 22 15
```

Rcpp 代码块也可以加 `cache=TRUE` 选项，使得 C++ 程序在没有修改时不必每次重新编译。



## Chapter 37

# R 与 C++ 的类型转换

R 程序与由 Rcpp 支持的 C++ 程序之间需要传递数据，就需要将 R 的数据类型经过转换后传递给 C++ 函数，将 C++ 函数的结果经过转换后传递给 R。

### 37.1 用 wrap() 把 C++ 变量返回到 R 中

在 R API 中用 .Call() 函数调用 C 程序库函数时，R 对象的数据类型一般是 SEXP。Rcpp 提供了模板化的 wrap() 函数把 C++ 的函数返回值转换成 R 的 SEXP 数据类型。此函数的声明为

```
template <typename T> SEXP wrap(const T& object);
```

wrap() 能转换的类型包括：

- 把 int, double, bool 等基本类型转换为 R 的原子向量类型（所有元素数据类型相同的向量）；
- 把 std::string 转换为 R 的字符型向量；
- 把 STL 容器如 std::vector<T> 或 std::map<T> 转换成基本类型为 T 的向量，条件是 T 能够转换；
- 把 STL 的映射 std::map<std::string, T> 转换为基本类型为 T 的有名向量，条件是 T 能够转换；

- 可以转换定义了 `operator SEXP()` 的 C++ 类的对象;
- 可以转换专门化过 `wrap` 模板的 C++ 对象。

是否可用 `wrap()` 转换是在编译时确定的。

## 37.2 用 `as()` 函数把 R 变量转换为 C++ 类型

Rcpp 提供了模板化的 `as()` 用来把 `SEXP` 类型转换成适当的 C++ 类型。`as()` 函数的声明为:

```
template <typename T> T as(SEXP x);
```

`as()` 可以把 R 对象转换为为基础的类型如 `int`, `double`, `bool`, `std::string` 等, 可以转换到元素为基础类型的 STL 向量如 `std::vector` 等。如果 C++ 类定义了以 `SEXP` 为输入的构造函数也可以利用 `as()` 来转换。`as()` 可以针对用户自定义类作专门化。

## 37.3 `as()` 和 `wrap()` 的隐含调用

当 C++ 中赋值运算的右侧表达式是一个 R 对象或 R 对象的部分内容时, 可以隐含地调用 `as()` 将其转换成左侧的 C++ 类型。

当 C++ 中赋值运算的左侧表达式是一个 R 对象或其部分内容时, 可以隐含地调用 `wrap()` 将右侧的 C++ 类型转换成 R 类型。

在用 Rcpp 属性 (Rcpp 属性见下一节) 声明的 C++ 函数中, 可以直接以 `IntegerVector`, `NumericVector`, `CharacterVector`, `Function` 等作为自变量类型或返回值, 可以与 R 中相应的类型直接对应。

能自动转换到 R 中的缺省值类型还包括:

- 用双撇号界定的字符串常量;
- 十进数值如 10, 4.5;
- 预定义的常数如 `true`, `false`, `R_NilValue`, `NA_STRING`, `NA_INTEGER`, `NA_REAL`, `NA_LOGICAL`;

# Chapter 38

## Rcpp 属性

### 38.1 Rcpp 属性介绍

Rcpp 属性 (attributes) 用来简化把 C++ 函数变成 R 函数的过程。做法是在 C++ 源程序中加入一些特殊注释，利用其指示自动生成 C++ 与 R 的接口程序。属性是 C++11 标准的内容，现在的编译器支持还不多，所以在 Rcpp 支持的 C++ 程序中写成了特殊格式的注释。

Rcpp 属性有如下优点：

- 降低了同时使用 R 与 C++ 的学习难度；
- 取消了很多繁复的接口代码；
- 可以在 R 会话中很简单地调用 C++ 代码，不需要用户自己考虑编译、连接、接口问题；
- 可以先交互地调用 C++，成熟后改编为 R 扩展包而不需要修改界面代码。

Rcpp 属性的主要组成部分如下：

- 在 C++ 中，提供 `Rcpp::export` 标注要输出到 R 中的 C++ 函数。
- 在 R 中，提供 `sourceCpp()`，用来自动编译连接保存在文件或 R 字符串中的 C++ 代码，并自动生成界面程序把 C++ 函数转换为 R 函数。

- 在 R 中，提供 `cppFunction()` 函数，用来把保存在 R 字符串中的 C++ 函数自动编译连接并转换成 R 函数。提供 `evalCpp()` 函数，用来把保存在 R 字符串中的 C++ 代码片段自动编译连接并执行。
- 在 C++ 中，提供 `Rcpp::depends` 标注，说明编译连接时需要的外部头文件和库的位置。
- 在构建 R 扩展包时，提供 `compileAttributes()` R 函数，自动给 C++ 函数生成相应的 `extern C` 声明和 `.Call` 接口代码。

## 38.2 在 C++ 源程序中指定要导出的 C++ 函数

用特殊注释 `//[[Rcpp::export]]` 说明某 C++ 函数需要在编译成动态链接库时，把这个函数导出到链接库的对外可见部分。

例如

```
//[[Rcpp::export]]
NumericVector convolveCpp(
 NumericVector a, NumericVector b){

}
```

具体程序参见前面“用 `sourceCpp()` 转换 C++ 源程序文件—卷积例子”。假设此 C++ 源程序保存到了当前工作目录的 `conv.cpp` 源文件中，为了在 R 中调用此 C++ 程序，只要用如：

```
sourceCpp(file='conv.cpp')
convolveCpp(1:3, 1:5)
```

注意 `sourceCpp()` 把 C++ 源程序自动进行了编译链接并转换成了同名的 R 函数。在同一 R 会话内，如果源程序和其依赖资源没有变化（根据文件更新时间判断），就不重新编译 C++ 源代码。

在用特殊注释说明要导出的 C++ 函数时，可以用特殊的 `name=` 参数指定函数导出到 R 中的 R 函数名。如果不指定，R 函数名和 C++ 函数名是相同的。

例如

```

//[[Rcpp::export(name="conv")]]
NumericVector convolveCpp(
 NumericVector a, NumericVector b){

}

```

则 C++ 函数 `convolveCpp` 导入到 R 中后，改名为 “conv”。

对于要导出的 C++ 函数，必须在全局名字空间中定义，而不能在某个 C++ 名字空间声明内定义。自变量必须能够用 `Rcpp::as` 转换成 C++ 类型，返回值必须是空值或者能够用 `Rcpp::wrap` 转换成 R 类型。在自变量和返回值类型说明中，必须使用完整的类型，比如 `std::string` 不能简写成 `string`。Rcpp 提供的类型如 `NumericVector` 可以不必用 `Rcpp::` 修饰。

### 38.3 在 R 中编译链接 C++ 代码

`sourceCpp()` 函数可以用 `code=` 指定一个 R 字符串，字符串的内容是 C++ 源程序，其中还是用特殊注释 `//[[Rcpp::export]]` 标识要导出的 C++ 函数。如

```

sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector convolveCpp(
 NumericVector a, NumericVector b){

}
')
convolveCpp(1:3, 1:5)

```

对于比较简单的单个 C++ 函数，可以用 `cppFunction()` 函数的 `code=` 指定一个 R 字符串，字符串的内容是一个 C++ 函数定义，转换为一个 R 函数。例

如

```
cppFunction(code='
 int fibonacci(const int x){
 if(x < 2) return x;
 else
 return (fibonacci(x-1) + fibonacci(x-2));
 }
')
print(fibonacci(5))
```

为了在 R 中计算一个简单的 C++ 表达式, 可以用 `evalCpp(' C++ 表达式内容')`, 如

```
evalCpp('std::numeric_limits<double>::max()')
```

函数将返回该 C++ 表达式的值。

在 `cppFunction()` 和 `evalCpp()` 中, 可以用 `depends=` 参数指定要链接的其它库, 如

```
sourceCpp(depends='RcppArmadillo', code='.....')
```

在编译代码时与 `RcppArmadillo` 的动态连接库连接。

也可以把这样的链接依赖关系写在特殊的 C++ 注释中, 如

```
///[[Rcpp::depends(RcppArmadillo)]]
```

这样的注释仅对 `sourceCpp()` 和 `cppFunction()` 有效, 在编译 R 扩展包时, 仍需要把依赖的包列在 `DESCRIPTION` 文件的 `Imports` 中, 把要链接的包列在 `LinkingTo` 中。

## 38.4 Rcpp 属性的其它功能

### 38.4.1 自变量有缺省值的函数

借助于 Rcpp, 自变量有缺省值的 C++ 函数可以自动转换成自变量有缺省值的 R 函数。定义时要符合 C++ 语法, 比如带缺省值的自变量都要在不带缺省值的自变量的后面, 缺省值不能有变量。

例如

```
DataFrame readData(
 CharacterVector file,
 CharacterVector colNames = CharacterVector::create(),
 std::string comment = "#",
 bool header = true){ ... }
```

转换到 R 中, 相当于

```
function(file,
 colNames=character(),
 comment="#",
 header=TRUE)
```

### 38.4.2 允许用户中断

在 C++ 代码中进行长时间的计算时, 应该允许用户可以中断计算。Rcpp 的办法是在 C++ 计算过程中每隔若干步循环就插入一个 `Rcpp::checkUserInterrupt()`; 语句。

### 38.4.3 把 R 代码写在 C++ 源文件中

正常情况下, 应该把 R 代码和 C++ 代码写在分别的源程序中, 当 C++ 代码比较短时, 也可以把 C++ 代码写在 R 源程序中作为一个字符串。

Rcpp 允许把 C++ 代码和 R 代码都写在一个 C++ 源文件中, R 代码作为特殊的注释, 以 `/** R` 行开头, 以正常的 `*/` 结束。在 R 中用 `sourceCpp()` 调

用这个 C++ 源文件，就可以编译 C++ 后执行其中特殊注释内的 R 代码。这样的特殊注释可以有多个。

例如，下述内容保存在文件 `fibonacci.cpp` 中：

```
[[Rcpp::export]]
int fibonacci(const int x){
 if(x < 2) return x;
 else
 return (fibonacci(x-1) + fibonacci(x-2));
}

/** R
 * 调用 C++ 中的 fibonacci() 函数
 * print(fibonacci(10))
 */
```

只要在 R 中运行

```
sourceCpp(file='fibonacci.cpp')
```

就可以编译连接此 C++ 文件，把其中用 `[[Rcpp::export]]` 标识的函数转换为 R 函数，并在 R 中执行源文件内特殊注释中的 R 代码。

#### 38.4.4 在 C++ 中调用 R 的随机数发生器

在 C 或 C++ 中调用 R 的随机数发生器，需要能够同步地更新随机数发生器状态。如果利用 `Rcpp` 属性编译 C++ 源程序，则 `Rcpp` 属性会自动添加一个 `RNGScope` 实例进行随机数发生器状态的同步。



# Chapter 39

## Rcpp 提供的 C++ 数据类型

### 39.1 RObject 类

Rcpp 包为 C++ 定义了 `NumericVector`, `IntegerVector`, `CharacterVector`, `Matrix` 等新数据类型, 可以直接与 R 的 `numeric`, `charactor`, `matrix` 对应。

Rcpp 最基础的 R 数据类型是 `RObject`, 这是 `NumericVector`, `IntegerVector` 等的基类, 通常不直接使用。`RObject` 包裹了原来 R 的 C API 的 `SEXP` 数据结构, 并且提供了自动的内存管理, 不再需要用户自己处理建立内存和消除内存的工作。`RObject` 存储数据完全利用 R C API 的 `SEXP` 数据结构, 不进行额外的复制。

因为 `RObject` 类是基类, 所以其成员函数也适用于 `NumericVector` 等类。`isNull`, `isObject`, `isS4` 可以查询是否 `NULL`, 是否对象, 是否 `S4` 对象。`inherits` 可以查询是否继承自某个特定类。用 `attributeNames`, `hasAttribute`, `attr` 可以访问对象的属性。用 `hasSlot`, `slot` 可以访问 `S4` 对象的插口 (slot)。

`RObject` 有如下导出类:

- `IntegerVector`: 整数向量;
- `NumericVector`: 数值向量;

- LogicalVector: 逻辑向量;
- CharacterVector: 字符型向量;
- GenericVector: 列表;
- ExpressionVector: 表达式向量;
- RawVector: 元素为 raw 类型的向量。
- IntegerMatrix, NumericMatrix: 整数值或数值矩阵。

在 R 向量中, 如果其元素都是同一类型 (如整数、双精度数、逻辑、字符型), 则称为原子向量。Rcpp 提供了 IntegerVector, NumericVector, LogicalVector, CharacterVector 等数据类型与 R 的原子向量类型对应。在 C++ 中可以用 [] 运算符存取向量元素, 也可以用 STL 的迭代器。用 .begin(), .end() 等界定范围, 用循环或或者 accumulate 等 STL 算法处理整个向量。

## 39.2 IntegerVector 类

在 R 中通常不严格区分整数与浮点实数, 但是在与 C++ 交互时, C++ 对整数与实数严格区分, 所以 RCpp 中整数向量与数值向量是区分的。

在 R 中, 如果定义了一个仅有整数的向量, 其类型是整数 (integer) 的, 否则是数值型 (numeric) 的, 如:

```
x <- 1:5
class(x)
[1] "integer"
y <- c(0, 0.5, 1)
class(y)
[1] "numeric"
```

用 as.integer() 和 as.numeric() 函数可以显式地确保其自变量转为需要的整数型或数值型。

RCpp 可以把 R 的整数向量传递到 C++ 的 IntegerVector 中, 也可以把 C++ 的 IntegerVector 函数结果传递回 R 中变成一个整数向量。也可以在 C++ 中生成 IntegerVector 向量, 填入整数值。

### 39.2.1 IntegerVector 示例 1: 返回完全数

如果一个正整数等于它所有的除本身以外的因子的和，称这个数为完全数。如

$$6 = 1 + 2 + 3$$

$$28 = 1 + 2 + 4 + 7 + 14$$

是完全数。

任务：用 C++ 程序输入前 4 个完全数偶数，返回到 R 中。这 4 个数为 6, 28, 496, 8182。

程序：

```
library(Rcpp)
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
IntegerVector wholeNumberA(){
 IntegerVector epn(4);
 epn[0] = 6; epn[1] = 28;
 epn[2] = 496; epn[3] = 8182;

 return epn;
}
')
```

```
print(wholeNumberA())
[1] 6 28 496 8182
```

从例子看出，可以在 C++ 中建立一个 IntegerVector，需指定大小。可以逐个填入数值。直接返回 IntegerVector 到 R 即可，不需用 wrap() 显式地转换。

### 39.2.2 IntegerVector 示例 2: 输入整数向量

任务: 用 C++ 编写函数, 从 R 中输入整数向量, 计算其元素乘积 (与 R 的 `prod()` 函数功能类似)。程序如下:

```
require(Rcpp)
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
IntegerVector prod1(IntegerVector x){
 int prod = 1;
 for(int i=0; i < x.size(); i++){
 prod *= x[i];
 }
 return wrap(prod);
}
')
print(prod1(1:5))
```

从程序看出, 用 `IntegerVector` 从 R 中接受一个整数值向量时, 不需要显式地转换。把一个 C++ 整数值返回给 R 时, 必须用 `IntegerVector` 返回, 因为返回值是一个 C++ 的 `int` 类型, 所以需要用 `Rcpp::wrap()` 转换一下。在 `sourceCpp` 中可以省略 `Rcpp::` 部分。

还可以用 C++ STL 的算法库进行这样的累计乘积计算, `std::accumulate()` 可以对指定范围进行遍历累计运算。前两个参数是一个范围, 用迭代器 (iterators) 表示开始和结束, 第三个参数是初值, 第四个参数是对每个元素累计的计算。程序如下:

```
require(Rcpp)
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;
```

```
[[Rcpp::export]]
IntegerVector prod2(IntegerVector x){
 int prod = std::accumulate(
 x.begin(), x.end(),
 1, std::multiplies<int>());
 return wrap(prod);
}
')
```

```
print(prod2(1:5))
```

在以上的输入 `IntegerVector` 的 C++ 程序中，如果从 R 中输入了实数型的向量，则元素被转换成整数型。比如 `prod2(seq(1,1.9,by=0.1))` 结果将等于 1。如果输入了无法转换为整数型向量的内容，比如 `prod2(c(' a' , ' b' , ' c' ))`，程序会报错。

## 39.3 NumericVector 类

`NumericVector` 类在 C++ 中保存双精度型一维数组，可以与 R 的实数型向量 (class 为 `numeric`) 相互转换。这是自己用 C++ 程序与 R 交互时最常用的数据类型。

### 39.3.1 示例 1: 计算元素 $p$ 次方的和

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

[[Rcpp::export]]
NumericVector ssp(
 NumericVector vec, double p){
 double sum = 0.0;
 for(int i=0; i < vec.size(); i++){
```

```

 sum += pow(vec[i], p);
 }
 return(wrap(sum));
}
')
ssp(1:4, 2)
[1] 30
ssp((1:4)/10, 2.2)
[1] 0.2392496
sum(((1:4)/10)^2.2)
[1] 0.2392496

```

从程序看出, 用 Rcpp 属性编译时, C++ 函数的输入与返回值类型转换有如下规则: R 中数值型向量在 C++ 中可以用 NumericVector 接收; R 中单个的实数在 C++ 中可以用 double 来接收; 为了返回单个的实数, 在 C++ 中需要以 NumericVector 为返回类型, 可以从一个 double 型用 wrap() 转换。C++ 中如果返回 NumericVector, 在 R 中转换为数值型向量。

### 39.3.2 示例 2: clone 函数

在自定义 R 函数时, 输入的自变量的值不会被改变, 相当于自变量都是局部变量。如果在自定义函数中修改了自变量的值, 实际上只能修改自变量的一个副本的值。如

```

x <- 100
f <- function(x){
 print(x); x <- 99; print(x)
}
c(f(x), x)
[1] 100
[1] 99
[1] 99 100

```

但是, 在用 Rcpp 编写 R 函数时, 因为 RObject 传递的是指针, 并不自动复

制自变量值，所以修改自变量值会真的修改原始的自变量变量值。如：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector f2(NumericVector x){
 x[0] = 99.0;
 return(x);
}
')
x <- 100
c(f2(x), x)
[1] 99 99
```

可见自变量的值被修改了。当然，对这个问题而言，因为输入的是一个标量，只要函数自变量不是 `NumericVector` 类型而是用 `double` 类型，则自变量值会被复制，达到值传递的效果，自变量值也就不会被真的修改。如

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector f3(double x){
 x = 99.0;
 return(wrap(x));
}
')
x <- 100
c(f3(x), x)
[1] 99 100
```

下面的程序把输入向量每个元素平方后返回，为了不修改输入自变量的值而是返回一个修改后的副本，使用了 `Rcpp` 的 `clone` 函数：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector square(NumericVector x){
 NumericVector y = clone(x);
 for(int i=0; i < x.size(); i++){
 y[i] = x[i]*x[i];
 }
 return(y);
}
')
x <- c(2, 7)
cbind(square(x), x)
x
[1,] 4 2
[2,] 49 7
```

### 39.3.3 示例 3: 把输入矩阵制作副本计算元素平方根

`NumericMatrix` 是 `Rcpp` 提供的元素为双精度型的矩阵。下面的例子输入一个 `R` 矩阵，输出其元素的平方根，用了 `clone` 函数来避免对输入的直接修改。

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericMatrix matSqrt(NumericMatrix x){
 NumericMatrix y = clone(x);
 std::transform(y.begin(), y.end(),
 y.begin(), ::sqrt);
}
```



```

 return(y);
}
')
x <- rbind(c(1,2), c(3,4))
cbind(matSqrt(x), x)
[,1] [,2] [,3] [,4]
[1,] 1.000000 1.414214 1 2
[2,] 1.732051 2.000000 3 4

```

在上面的 C++ 程序中, NumericMatrix 看成了一维数组, 用 STL 的 iterater 遍历, 用 STL 的 transform 对每个元素计算变换。

## 39.4 Rcpp 的其它向量类

### 39.4.1 Rcpp 的 LogicalVector 类

LogicalVector 类可以存储 C++ 值 true, false, 还可以保存缺失值 NA\_REAL, R\_NaN, R\_PosInf, 但是这些不同的缺失值转换到 R 中都变成 NA。

如:

```

sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
LogicalVector f4(){
 LogicalVector x(5);
 x[0] = false; x[1] = true;
 x[2] = NA_REAL;
 x[3] = R_NaN; x[4] = R_PosInf;
 return(x);
}
')
```

```
f()
[1] FALSE TRUE NA NA NA
identical(f(), c(FALSE, TRUE, rep(NA,3)))
[1] TRUE
```

### 39.4.2 Rcpp 的 CharacterVector 类型

CharacterVector 类型可以与 R 的字符型向量相互交换信息，在 C++ 中其元素为字符串。字符型缺失值在 C++ 中为 R\_NAString。R 的字符型向量也可以转换为 std::vector。

如：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
CharacterVector f5(){
 CharacterVector x(3);
 x[0] = "This is a string";
 x[1] = "Test";
 x[2] = R_NAString;
 return(x);
}
')
f5()
[1] "This is a string" "Test" NA
```

## 39.5 Rcpp 提供的其它数据类型

### 39.5.1 Named 类型

R 中的向量、矩阵、数据框可以有元素名、列名、行名。这些名字可以借助 Named 类添加。

例如：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector f6(){
 NumericVector x = NumericVector::create(
 Named("math") = 82,
 Named("chinese") = 95,
 Named("English") = 60);
 return(x);
}
')
```

**f6()**

| ## | <i>math</i> | <i>chinese</i> | <i>English</i> |
|----|-------------|----------------|----------------|
| ## | 82          | 95             | 60             |

“Named(元素名)”可以简写成“\_\_元素名”。如：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector f6b(){
 NumericVector x = NumericVector::create(
 _["math"] = 82,
```

```
 _["chinese"] = 95,
 _["English"] = 60);
 return(x);
}
')
```

### 39.5.2 List 类型

Rcpp 提供的 List 类型对应于 R 的 list(列表) 类型，在 C++ 中也可以写成 GenericVector 类型。其元素可以不是同一类型，在 C++ 中可以用方括号和字符串下标的格式访问其元素。

例如，下面的函数输入一个列表，列表元素 vec 是数值型向量，列表元素 multiplier 是数值型标量，返回一个列表，列表元素 sum 为 vec 元素和，列表元素 dsum 为 vec 元素和乘以 multiplier 的结果：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
List f7(List x){
 NumericVector vec = as<NumericVector>(x["vec"]);
 double multiplier = as<double>(x["multiplier"]);
 double y = 0.0, y2;
 for(int i=0; i<vec.length(); i++){
 y += vec[i];
 }
 y2 = y*multiplier;
 return(List::create(Named("sum")=y,
 Named("dsum")=y2));
}
')
f7(list(vec=1:5, multiplier=10))
```

```
$sum
[1] 15
##
$dsum
[1] 150
```

上面的程序用了 `Rcpp::List::create()` 当场生成 `List` 类型，因为用 `Rcpp` 属性功能编译所以可以略写 `Rcpp::`。也可以在程序中预先生成指定大小的列表，然后再给每个元素赋值，元素值可以是任意能够转化为 `SEXP` 的类型，如：

```
.....
List gv(2);
gv[0] = "abc";
gv[1] = 123;
```

可以用 `List` 的 `reserve` 函数为列表指定元素个数。

### 39.5.3 Rcpp 的 DataFrame 类

`Rcpp` 的 `DataFrame` 类用来与 `R` 的 `data.frame` 交换信息。

示例如：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
DataFrame f8(){
 IntegerVector vec =
 IntegerVector::create(7,8,9);
 std::vector<std::string> s(3);
 s[0] = "abc"; s[1] = "ABC"; s[2] = "123";
 return(DataFrame::create(
 Named("x") = vec,
 Named("s") = s));
}
```

```

}
')
f8()
x s
1 7 abc
2 8 ABC
3 9 123

```

### 39.5.4 Rcpp 的 Function 类

Rcpp 的 Function 类用来接收一个 R 函数，并且可以在 C++ 中调用这样的函数。

示例如：

```

sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
SEXP dsort(Function sortFun, SEXP x){
 return sortFun(x, Named("decreasing", true));
}
')
dsort(sort, c('bb', 'ab', 'ca'))
[1] "ca" "bb" "ab"
dsort(sort, c(2,1,3))
[1] 3 2 1

```

程序用 Function 对象 sortFun 接收从 R 中传递过来的排序函数，实际调用时传递过来的是 R 的 sort 函数。

在 C++ 中调用 R 函数时，有名的自变量用 “Named(自变量字符串, 自变量值)” 的格式给出。

程序中的待排序的向量与排序后的向量都用了 SEXP 来说明，即直接传送原始

R API 指针，这样可以不管原来类型是什么，在 C++ 中完全不进行类型转换或复制。

从运行例子看出数值和字符串都正确地按照降序排序后输出了。

R 函数可以不作为 C++ 的函数自变量传递进来，而是直接调用 R 的函数。

下面的函数直接调用 `rt` 函数生成 2 个自由度为 3 的 t 分布随机数：

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector rt23(){
 Function rt("rt");
 return rt(2, 3);
}
')
set.seed(1)
rt23()
[1] -0.7027211 -0.5693196
rt23()
[1] 0.6842766 -0.3620012
```

使用了 `Rcpp` 属性时，生成的界面程序会自动生成一个 `RNGScope` 的实例，用来保存当前的随机数发生器状态，在解构时将自动更新随机数发生器状态。

用 `rt` 这个 R 函数名初始化了一个 C++ 的 `Function` 对象 `rt`，然后就可以调用这个 C++ 函数了。从程序可以看出，连续两次调用的结果不同。

### 39.5.5 Rcpp 的 Environment 类

R 的环境是分层的，可以逐层查找变量名对应的内容。`Rcpp` 的 `Environment` 类用来与 R 环境对应。可以利用 `Environment` 来定位 R 的扩展包中的函数或数据，例如下面的程序片段在 C++ 中定位了 `stats` 扩展包中的 `rnorm` 函数并进行了调用：

```
Environment stats("package:stats");
Function rnorm = stats["rnorm"];
return rnorm(3, Named("sd", 100.0));
```

当然，也可以用 Function 对象直接从各环境中搜索 rnorm 函数名，但是这样指定环境更可靠。

下面的例子访问了 R 全局环境，取出了全局变量 x 的值存入 C++ 的双精度型 STL 向量中，并把一个 C++ 的 STL map 型数据转换成有名字符型向量存到了全局变量 y 中。

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
void f9(){
 Environment global = Environment::global_env();
 std::vector<double> vx = global["x"];
 std::map<std::string, std::string> ma;
 ma["foo"] = "abc";
 ma["bar"] = "123";
 global["y"] = ma;
 return;
}
')
x <- c(1,5)
f9()
y
bar foo
"123" "abc"
```



# Chapter 40

## Rcpp 糖

在 C++ 中，向量和矩阵的运算通常需要逐个元素进行，或者调用相应的函数。Rcpp 通过 C++ 的表达式模板 (expression template) 功能，可以在 C++ 中写出像 R 中对向量和矩阵运算那样的表达式。这称为 Rcpp 糖 (sugar)。

R 中的很多函数如 `sin` 等是向量化的，Rcpp 糖也提供了这样的功能。Rcpp 糖提供了一些向量化的函数如 `ifelse`, `sapply` 等。

比如，两个向量相加可以直接写成 `x + y` 而不是用循环或迭代器 (iterator) 逐元素计算；若 `x` 是一个 `NumericVector`，用 `sin(x)` 可以返回由 `x` 每个元素的正弦值组成的 `NumericVector`。

Rcpp 糖不仅简化了程序，还提高了运行效率。

### 40.1 简单示例

比如，函数

$$f(x, y) = \begin{cases} x^2 & x < y, \\ -y^2 & x \geq y \end{cases}$$

如下的程序可以在 C++ 中定义一个向量化的版本：

```

sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector f11(NumericVector x, NumericVector y){
 return ifelse(x < y, x*x, -(y*y));
}
')
f11(c(1, 3), c(4,2))
[1] 1 -4

```

上面简单例子中,  $x < y$  是向量化的,  $x * x$ ,  $y * y$ ,  $-(y * y)$  都是向量化的。  
`ifelse` 也是向量化的。

## 40.2 向量化的运算符

### 40.2.1 向量化的四则运算

Rcpp 糖向量化了  $+$ ,  $-$ ,  $*$ ,  $/$ 。设  $x$ ,  $y$  都是 `NumericVector`, 则  $x + y$ ,  $x - y$ ,  $x * y$ ,  $x / y$  将返回对应元素进行四则运算后的 `NumericVector` 变量。

向量与标量运算, 如  $x + 2.0$ ,  $2.0 - x$ ,  $x * 2.0$ ,  $2.0 / x$  将返回标量与向量每个元素进行四则运算后的 `NumericVector` 变量。

还可以进行混合四则运算, 如:

```

NumericVector res = x * y + y / 2.0;
NumericVector res = x * (y - 2.0);
NumericVector res = x / (y * y);

```

参加四则运算的或者都是同基本类型的向量而且长度相等, 或者一边是向量, 一边是同类型的标量。

注意: 对向量整体赋一个相同的值, 不能简单地写成如  $x=0$ ; 这样的赋值, 需要用循环或 STL 的 `fill` 算法, 如 `std::fill(x.begin(), x.end(), 0);`。

### 40.2.2 向量化的二元逻辑运算

Rcpp 糖扩展了两个元素的比较到两个等长向量的比较，以及一个向量与一个同类型标量的比较，结果是一个同长度的 LogicalVector。比较运算符包括 `<`, `>`, `<=`, `>=`, `==`, `!=`。

也可以使用嵌套的表达式，比如，设 `x, y` 是两个 NumericVector，可以写：

```
LogicalVector res = (x + y) < (x * x);
```

### 40.2.3 向量化的一元运算符

对数值型向量或返回数值型向量的表达式前面加负号，可以返回每个元素取相反数的结果。比如，设 `x` 是 NumericVector，可以写：

```
NumericVector res = -x;
NumericVector res = -x * (x + 2.0);
```

对逻辑型向量或返回逻辑型向量的表达式前面加叹号，可以返回每个元素取反的结果，如：

```
NumericVector y, z;
NumericVector res = ! (y < z);
```

## 40.3 用 Rcpp 访问数学函数

在 C 和 C++ 代码中可以调用 R 中的函数，但是格式比较复杂，而且原来不支持向量化。Rcpp 糖则使得从 C++ 中调用 R 函数变得和在 R 调用函数格式类似。

R 源程序中提供了许多数学和统计相关的函数，这些函数可以编译成一个独立的函数库，供其它程序链接使用，函数内容在 R 的 `Rmath.h` 头文件中有详细列表。

Rcpp 糖把 R 中的数学函数在 C++ 中向量化了，输入一个向量，结果是对应元素为函数值的向量。自变量类型可以是数值型或整数型。这些数学函数包括 `abs`, `exp`, `floor`, `ceil`, `pow` 等。在 Rcpp 中支持的 C++ 源程序中调用这些函

数, 最简单的一种方法是直接在 Rcpp 名字空间中使用和 R 中相同的函数名和调用方法, 类似 sqrt 这样的函数允许输入一个向量, 对向量的每个元素计算。

比如, 下面的程序输入一个向量, 计算相应的标准正态分布函数值:

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
NumericVector f10(NumericVector x){
 NumericVector y(x.size());
 y = pnorm(x);
 return y;
}
')
f10(c(0, 1.96, 2.58))
[1] 0.5000000 0.9750021 0.9950600
```

如果需要对单个的数值计算, 可以使用 Rmath.h 中定义的带有 Rf\_ 的版本, 如:

```
y = ::Rf_pnorm5(x, 0.0, 1.0, 1, 0);
```

这里用:: 使用了缺省的名字空间。注意所有自变量都不能省略。

Rcpp 还提供了一个 R 名字空间, 可以用不带 Rf\_ 前缀的函数, 但是自变量也不能省略。如

```
y = R::pnorm(x, 0.0, 1.0, 1, 0);
```

在 Rcpp 的 R 名字空间中有许多的数学和统计相关的函数, 各函数的自变量参见 Rcpp 的 Rmath.h 文件。

R 中提供了许多分布的分布密度 (或概率质量函数)、分布函数、分位数函数, 分布密度函数和概率质量函数命名类似 dxxxx, 分布函数命名类似 pxxxx, 分位数函数命名类似 qxxxx。在 Rcpp 糖支持下, 这些函数可以直接用类似 R 中的格式调用。如

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
List f14(){
 NumericVector x =
 NumericVector::create(0, 1.96, 2.58);
 NumericVector p =
 NumericVector::create(0.95, 0.975, 0.995);
 NumericVector y1 = dnorm(x, 0.0, 1.0);
 NumericVector y2 = pnorm(x, 0.0, 1.0);
 NumericVector y3 = qnorm(p, 0.0, 1.0);
 return List::create(Named("y1")=y1,
 Named("y2")=y2, Named("y3")=y3);
}
')
```

```
f14()
$y1
[1] 0.39894228 0.05844094 0.01430511
##
$y2
[1] 0.5000000 0.9750021 0.9950600
##
$y3
[1] 1.644854 1.959964 2.575829
##
```

R 中的 `rxxxx` 类的函数可以产生各种分布的随机数向量，随机数向量与当前种子有关。Rcpp 属性会自动地维护随机数发生器的状态使其与 R 同步。如

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;
```

```

//[[Rcpp::export]]
NumericVector f15(){
 NumericVector x = rnorm(10, 0.0, 1.0);
 return x;
}
')
round(f15(), 2)
[1] 0.62 -0.06 -0.16 -1.47 -0.48
[6] 0.42 1.36 -0.10 0.39 -0.05

```

## 40.4 返回单一逻辑值的函数

在 R 中，`any()` 和 `all()` 对一个逻辑向量分别判断是否有任何真值，以及所有元素为真值。Rcpp 糖在 C++ 中也提供了这样的 `any()` 和 `all()` 函数。如

```

sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
List f12(){
 IntegerVector x = seq_len(1000);
 LogicalVector res1 = any(x*x < 3);
 LogicalVector res2 = all(x*x < 3);
 return List::create(Named("any")=res1,
 Named("all")=res2);
}
')
f12()
$any
[1] TRUE
##

```

```
$all
[1] FALSE
```

`any()` 和 `all()` 不是直接返回逻辑值，而是返回一个类对象，该类定义了 `is_true`, `is_false`, `is_na` 方法与向 SEXP 转换的运算符。此种结果可以保存到 `LogicalVector` 中，但不能赋值到 `bool` 类型，因为可能有缺失值。

逻辑型糖表达式结果可以用 `is_true`, `is_false`, `is_na` 来判断结果。如

```
bool res1 = is_true(any(x < y));
bool res2 = is_na(all(x < y));
```

在求 `any()` 结果时，一旦遇到一个真值结果就为真值，即使后面有缺失值也没有关系。在求 `all()` 结果时，一旦遇到一个假值结果就为假值，即使后面有缺失值也没有关系。

## 40.5 返回糖表达式的函数

`is_na` 以任何糖表达式为输入，输出一个逻辑类型的元素个数相同的糖表达式。结果每个元素当输入中对应元素缺失时为 `TRUE`，否则为 `FALSE`。如

```
IntegerVector x = IntegerVector::create(0, 1, NA_INTEGER, 3);
LogicalVector res1 is_na(x);
LogicalVector res2 = all(is_na(x));
if(is_true(any(! is_na(x)))) ...
```

`seq_along` 输入一个向量，输出一个元素为该向量的各个下标值的糖表达式。如

```
IntegerVector x = IntegerVector::create(0, 1, NA_INTEGER, 3);
seq_along(x);
seq_along(x*x*x*x*x);
```

注意上述程序不会计算 `x*x*x*x*x` 的值而只利用其结果的元素个数。所以两次调用的计算量是一样的。

`seq_len` 自变量为个数，返回一个元素值为正数的元素个数等于自变量值的糖

表达式，第  $i$  个元素等于  $i$ 。常可与 `sapply`, `lapply` 配合使用。如

```
sourceCpp(code='
#include <Rcpp.h>
using namespace Rcpp;

//[[Rcpp::export]]
List f13(){
 IntegerVector x =seq_len(3);
 List y = lapply(seq_len(3), seq_len);
 return List::create(Named("x")=x,
 Named("y")=y);
}
')
```

```
f13()
$x
[1] 1 2 3
##
$y
$y[[1]]
[1] 1
##
$y[[2]]
[1] 1 2
##
$y[[3]]
[1] 1 2 3
##
```

`pmin` 和 `pmax` 自变量为两个等长的向量或糖表达式，返回长度相同的结果，结果元素只等于对应元素的最小值或最大值。

自变量也可以取一个向量一个标量，向量的每个元素与标量比较得到最小值或最大值。如



```
IntegerVector x =seq_len(10);
pmin(x, x*x);
pmin(x*x, 2);
```

`ifelse` 函数有三个自变量，第一自变量是一个逻辑型向量值的糖表达式，第二和第三自变量或者是两个同类型并与第一自变量等长的糖表达式，或者其中一个为同类型标量，结果仍为向量型糖表达式，第  $i$  元素当第一自变量第  $i$  元素为真时取第一自变量的第  $i$  元素，当第一自变量第  $i$  元素为假时取第二自变量的第  $i$  元素，当第一自变量第  $i$  元素为缺失值时取缺失值。如

```
IntegerVector x, y;
IntegerVector res1 = ifelse(x < y, x, (x+y)*y);
IntegerVector res2 = ifelse(x > y, x, 2);
```

`sapply` 函数第一自变量是一个向量或列表，第二自变量是一个函数。返回值类型在编译时从函数的结果类型导出。第二自变量可以是任意的 C++ 函数，比如，可以是如下的重载的模板化函数：

```
template <typename T>
T square(const T& x){
 return x * x;
}
sapply(seq_len(4), square<int>);
```

下面的例子中 `sapply` 第二自变量使用了 functor，这是一种能够产生函数的函数：

```
template <typename T>
struct square : std::unary_function<T,T> {
 T operator() (const T& x) {
 return x * x;
 }
}
sapply(seq_len(4), square<int>());
```

`lapply` 函数与 `sapply` 函数基本相同，只不过 `lapply` 函数总是返回列表，列表在 Rcpp 中为 `List` 或 `GenericVector`，在 R API 中类型为 `VECSXP`。

`sign` 函数输入一个数值型或整型表达式，返回各元素值在  $\{1, 0, -1, NA\}$  中取值的糖表达式，可以保存到 `IntegerVector` 中。结果各元素的取值表示输入中对应元素为正、零、负和缺失。如

```
IntegerVector x;
IntegerVector res1 = sign(x);
IntegerVector res2 = sign(x*x);
```

`diff` 函数自变量为数值型或整数型向量或有这样结果的糖表达式，输出后一元素减去前一元素的结果，结果长度比输入长度少一。如

```
IntegerVector res = diff(seq_len(5));
```

## 40.6 R 与 Rcpp 不同语法示例

Rcpp 糖通过一些现代的 C++ 技术，支持部分的向量化运算，比如，`NumericVector` 与标量相加，两个等长 `NumericVector` 相加，对 `NumericVector` 计算如绝对值这样的数学函数值等。但是，C++ 毕竟和 R 有很大差别，要区分 Rcpp 能做的和不能做的。

例如，`NumericVector` 为了把向量所有元素都改成同一值，不能直接用等于赋值。可以用 `std::fill(y.begin(), y.end(), 99.99)` 这样的做法。

## Chapter 41

# 用 Rcpp 帮助制作 R 扩展包

R 扩展包是把解决某种问题的可复用代码、文档整合在一起的最好的方法。写成 R 扩展包后，可以自己用，也可以利用 CRAN 分发。扩展包用户一般不用自己编译。

使用扩展包来组织程序，多个源程序、头文件之间的依赖关系可以自动得到处理。

扩展包提供了测试、文档和一致性检查的统一框架。

扩展包中代码可以仅有 R 程序，也可以包括 C 程序、C++ 程序、Fortran 程序。如果仅有 R 代码，就不需要借助于 Rcpp，可以使用 `package.skeleton()` 函数生成一个扩展包框架。如果有 C++ 代码，就可以用 Rcpp 作为接口，并用 Rcpp 提供的 `Rcpp.package.skeleton()` 函数制作扩展包框架。

Rcpp 属性的 `Exports` 注释仍可在制作扩展包时指定如何输出 C++ 中定义的函数使其在 R 中可调用。

### 41.1 不用扩展包共享 C++ 代码的方法

Rcpp 属性的 `sourceCpp()` 通常只适用于写在 R 程序内部的简短 C++ 代码，或者写在一个单独 C++ 文件中，不依赖于其它 C++ 程序的单独代码。如果有多个 C++ 源程序、头文件，彼此有依赖关系，最好使用扩展包。

在多个单独的 C++ 文件共享某些简单的代码，彼此不互相依赖时，可以用 C++ 的预处理 `include` 命令共享这些代码。

比如，有多个 C++ 源程序都用到如下的代码：

```
#ifndef __UTILITIES__
#define __UTILITIES__
inline double timesTwo(double x) {
 return x * 2;
}
#endif // __UTILITIES__
```

假设这段代码保存到目前子目录的“`utilities.hpp`”文件中。则在每个需要用到这段代码的 C++ 源程序中，插入如：

```
#include "utilities.hpp"
//[[Rcpp::export]]
double transformValue(double x){
 return timesTwo(x) * 10;
}
```

## 41.2 生成扩展包

### 41.2.1 利用已有基于 Rcpp 属性的源程序制作扩展包

假设在当前目录中有了若干个 C++ 文件，其中需要转换到 R 中的 C++ 函数已经用 `Rcpp::export` 声明过。其中一个是 `conv1.cpp`。

从当前目录启动 R，运行

```
Rcpp.package.skeleton("testpack",
 example_code=FALSE,
 attributes=TRUE,
 cpp_files=c("conv1.cpp"))
```

运行显示：

```
Creating directories ...
Creating DESCRIPTION ...
Creating NAMESPACE ...
Creating Read-and-delete-me ...
Saving functions and data ...
Making help files ...
Done.
Further steps are described in './testpack/Read-and-delete-me'.
```

#### Adding Rcpp settings

```
>> added Imports: Rcpp
>> added LinkingTo: Rcpp
>> added useDynLib directive to NAMESPACE
>> added importFrom(Rcpp, evalCpp) directive to NAMESPACE
>> copied conv1.cpp to src directory
```

运行完后，在当前目录生成了一个 `testpack` 子目录，这是要制作的扩展包的名字。在 `testpack` 子目录中，有文件 `DESCRIPTION`, `NAMESPACE`, `Read-and-delete-me`, 有子目录 `src`, `R`, `man`。

子目录 `src` 中为 C++ 和 C 源程序、头文件。子目录 `R` 中为 Rcpp 从 C++ 程序转换过来的 R 接口程序，用户自己的 R 程序也可以放在这里。子目录 `man` 是特殊格式的文档，其格式类似 LaTeX。

### 41.2.2 DESCRIPTION 文件

在 `DESCRIPTION` 文件中，有扩展包名称、版本、日期、作者姓名、维护者姓名和联系方式、简单描述、授权，还有 `Imports` 和 `LinkingTo` 两项。除此之外，还有许多可选的域，如 `Depends`。

`Imports` 给出本软件包要调用的其它扩展包，但是这些扩展包并不随本扩展包一起调入，仅是会调入其名字空间。这里的值为

```
Imports: Rcpp (>= 0.12.3)
```

`LinkingTo` 指定在编译本软件包的 C、C++、Fortran 等源程序时，会用到哪

些其它扩展包的头文件。这里的值为

**LinkingTo:** Rcpp

指定的这些扩展包一般是编译时才有用的，所以一般不会出现在 `Depends` 和 `Imports` 域中。

`LinkingTo` 只解决了头文件的问题，要链接除了 Rcpp 之外的二进制库文件，还需要手工编辑 `src/Makevars` 和 `src/Makevars.win` 文件。

和 `Imports` 有些相像的 `DESCRIPTION` 域是 `Depends`，指定调入本扩展包时必须预先调入的软件包。这里“调入”是指用 `library()` 或 `require()` 调入扩展包。多个扩展包名用逗号分开，可以在扩展包名字后面加圆括号，在圆括号内写上 `>=` 某个版本号，如“`MASS(>=3.1-20)`”。

`Depends` 也可以指定依赖于某个 R 版本之后，如“`R(>=2.14.0)`”。

`DESCRIPTION` 文件中的 `Suggests` 与 `Depends` 域类似，但不是本扩展包必须的，比如仅用在某个例子中或测试中、仅用来编译 vignettes。

### 41.2.3 NAMESPACE 文件

示例 `NAMESPACE` 文件如下：

```
useDynLib(testpack)
exportPattern("^[:alpha:]+")
importFrom(Rcpp, evalCpp)
```

其中第一行指定调用本软件包时，需要调入的本扩展包的动态链接库。第二行指定扩展包需要对外部可见的 R 函数是所有函数名字以字母开头的 R 函数。用户可以自己指定其它的模式或者指定固定的若干个函数。第三行说明了需要从 Rcpp 包导入 `evalCpp` 函数。

## 41.3 重新编译

修改了扩展包中的 C++ 源程序后，需要重新编译。只要在 R 中把工作目录设为软件包的子目录内，运行

```
compileAttributes()
```

这会自动生成两个文件，一个是 `src/RcppExports.cpp`，是 C++ 程序的接口函数。另一个是 `R/RcpExports.R`，用 `.Call` 来调用 C++ 接口函数，转换成 R 函数。这两个文件不要自己修改。

## 41.4 建立 C++ 用的接口界面

利用了 `Rcpp` 属性可以指定要输出到 R 中的函数。

在 C++ 源程序中加入特殊注释

```
//[[Rcpp::interfaces(r, cpp)]]
```

则软件包在编译时也会生成该源程序文件中函数的外部可访问的接口，这些接口的界面会在安装后的包的 `include` 子目录中出现，在开发时出现在 `inst/include` 子目录中。

设要生成的扩展包名为 `testpack`，则界面文件包括 `include` 子目录中的 `testpack_RcppExports.h` 文件和 `testpack.h` 文件，`testpack.h` 文件仅用来包含入 `testpack_RcppExports.h` 文件。

如果需要添加自己的一些界面程序，可以修改 `testpack.h` 文件，这时需要去掉文件开始的自动生成标记，并且保留对 `testpack_RcppExports.h` 文件的包含。

导出的 C++ 界面都在与制作的扩展包同名的名字空间中，比如，如果制作的软件包名为 `testpack`，其中导出的一个 C++ 函数为 `convolveCpp`，则在别的包的 C++ 源程序中调用时，应该包含 `testpack.h` 文件，并用 `testpack::convolveCpp()` 格式调用。

如果自己本扩展包需要在编译时包含这些头文件，需要自己编辑 `src` 子目录中的 `Makevars` 文件和 `Makevars.win` 文件，添加行：

```
PKG_CPPFLAGS += -I../inst/include/
```





## Part X

# R 编程例子



# Chapter 42

## R 编程例子

因为这些例子要用作学生习题，所以这里只有问题，没有实际内容。

### 42.1 R 语言

#### 42.1.1 用向量作逆变换

设向量  $x$  长度为  $n$ ，其中保存了 1 到  $n$  的正整数的一个排列。把  $x$  看成是在集合  $\{1, 2, \dots, n\}$  上的一个一一变换，求向量  $y$  使得  $y$  能够表示上述变换的逆变换。即任给长度为  $n$  的向量  $z$ ， $z[x]$  表示按照  $x$  的次序重新排列  $z$  的元素，而  $z[x][y]$  则应该恢复为  $z$ 。

#### 42.1.2 斐波那契数列计算

设数列  $x_0 = 0, x_1 = 1$ ，后续值按如下公式递推计算：

$$x_n = x_{n-2} + x_{n-1}, \quad n = 2, 3, \dots$$

这样的数列叫做斐波那契数列。希望编写 R 函数，输入  $n$ ，返回计算的  $x_n$  的值。

### 42.1.3 穷举所有排列

设向量  $\mathbf{x}$  的各个元素为某个集合的元素。想要列出  $\mathbf{x}$  的元素的所有不同排列。比如, 如果  $\mathbf{x} = 1:3$ , 所有排列为

```
1 2 3
1 3 2
2 1 3
2 3 1
3 1 2
3 2 1
```

共  $3! = 6$  种不同的排列。

### 42.1.4 可重复分组方式穷举

设有  $n$  个编号卡片, 分别有号码  $1, 2, \dots, n$ 。从中有放回地抽取  $m$  个并记录每次的号码, 穷举  $m$  个号码中多少个 1, 多少个 2,  $\dots$ , 多少个  $n$  这样的结果。

例如, 有 3 个编号卡片, 随机有放回地抽取 2 次。用  $(x_1, x_2, x_3)$  表示每一种个数组合,  $x_1$  表示 2 次抽取中号码 1 的个数,  $x_2$  表示 2 次抽取中号码 2 的个数,  $x_3$  表示 2 次抽取中号码 3 的个数。问题就是列出  $(x_1, x_2, x_3)$  的所有不同值。

## 42.2 概率

### 42.2.1 智者千虑必有一失

成语说: “智者千虑, 必有一失; 愚者千虑, 必有一得”。设智者作判断的准确率为  $p_1 = 0.99$ , 愚者作判断的准确率为  $p_2 = 0.01$ , 计算智者做 1000 次独立的判断至少犯一次错误的概率, 与愚者做 1000 次独立判断至少对一次的概率。

### 42.2.2 圆桌夫妇座位问题

在一张圆桌上用餐时,  $n$  对夫妇随机入座。计算没有任何一位妻子和她的丈夫相邻的概率。通过推导可得此概率为

$$p_n = 1 + \sum_{k=1}^n (-1)^k C_n^k \frac{(2n-k-1)! 2^k}{(2n-1)!}. \quad (1)$$

例如  $p_2 = 1/3$ ,  $p_3 = 4/15 = 0.2667$ ,  $p_4 = 0.2952$ ,  $p_5 = 0.3101$ 。

分别用上面的理论公式以及直接穷举验证的方法, 对  $n = 2, 3, 4, 5$  的情形进行验证。

## 42.3 科学计算

### 42.3.1 城市间最短路径

假设有  $n$  个城市, 编号为  $1, 2, \dots, n$ 。已知其中的部分城市之间有高速公路连通, 每对连通城市记为  $(F_i, T_i)$ , 其中  $F_i, T_i \in \{1, 2, \dots, n\}$  且  $F_i < T_i$ ,  $i = 1, 2, \dots, m$ 。除了这些直接连通的城市以外, 其它的任意两个城市只能途经别的城市连通, 或者根本不能靠高速公路连通。用一个  $m \times 2$  的 R 矩阵  $M$  可以输入这些连通情况, 矩阵的每行是一对  $(F_i, T_i)$  值。

要求编写一个 R 函数, 输入直接连通情况  $M$  后, 输出一个  $n \times n$  矩阵  $R$ ,  $R[i,i]=0$ ,  $R[i,j]=1$  表示直接相连,  $R[i,j]=k$  ( $k \geq 2$ ) 表示城市  $i$  与城市  $j$  至少需要经过  $k$  段高速公路连通,  $R[i,j]=\text{Inf}$  表示城市  $i$  与城市  $j$  不能靠高速公路连通。 $R$  的元素值仅考虑途经的高速公路段数而不考虑具体里程。如果从一个城市通过高速公路移动到直接相连城市叫做移动一步,  $R$  的  $(i, j)$  元素是从第  $i$  城市通过高速公路到第  $j$  城市需要移动的步数。

这个问题也可以作为“相识”问题的模型。设  $n$  个人中有些人是直接相识的, 如果两个不相识的人想认识, 假设必须经过相识的人引荐, 问最少需要多少个引荐人。

本问题必须使用循环, 很难向量化, 属于 R 比较不擅长的问题。可以考虑使用 Rcpp 把程序用 C++ 语言实现, 可以大大加速。当然, 如果这个程序仅需要执行不多的次数, 用 R 就足够了。

### 42.3.2 Daubechies 小波函数计算

这个例子主要展示了不用每次计算函数值而是尽可能从已经计算并储存的函数值中查找的技巧。程序中用了比较多的循环，如果有需要，可考虑用 Rcpp 转换成 C++ 代码以提高效率。

小波是重要数学工具，在图像处理、信号处理等方面有广泛应用。小波中一个重要的函数叫做尺度函数 (scale function)，它满足所谓双尺度方程：

$$\phi(x) = \sqrt{2} \sum_k h_k \phi(2x - k)$$

一种特殊的尺度函数是只在有限区间上非零的，叫做紧支集的。紧支集尺度函数可以在给定  $\{h_k\}$  后用以下迭代公式生成：

$$\begin{aligned} \eta_0(x) &= I_{[-0.5, 0.5]}(x) \\ \eta_{n+1}(x) &= \sqrt{2} \sum_{k=0}^{2N-1} h_k \eta_n(2x - k) \end{aligned}$$

其中  $N$  是正整数， $N=2$  时  $h_0=0.482962913145$ ,  $h_1=0.836516303738$ ,  $h_2=0.224143868042$ ,  $h_3 = -0.129409522551$ 。已知  $\phi(x)$  的支集 (不为零的区间) 为  $[0, 2N - 1]$ ,  $\eta_n(x)$  的支集包含于  $[-0.5, 2N - 1]$  中。

任务：编写计算  $\phi(x)$  的 R 程序，通过 20 次迭代计算，输出  $\phi(x)$  在  $[0, 2N - 1]$  区间的 256 个等间隔点上的函数值并作图。在迭代过程中，应不重新结算函数格子点的值，仅计算新加入的格子点的值。不要利用递归函数，递归函数需要每次重新计算已经计算过的函数值。

### 42.3.3 房间加热温度变化

某个房屋带有天花板，天花板与屋顶之间有一定的空间。用壁炉保持房间温度。为了研究房间内与天花板上方的温度变化，建立了如下的微分方程组：

$$\begin{aligned} \frac{dx_1}{dt} &= 0.35 \left( -9.7 \sin \frac{(t+3)\pi}{12} + 8.3 - x_1(t) \right) + 0.46(x_2(t) - x_1(t)) + 11.1 \\ \frac{dx_2}{dt} &= 0.28 \left( -9.7 \sin \frac{(t+3)\pi}{12} + 8.3 - x_2(t) \right) + 0.46(x_1(t) - x_2(t)) \end{aligned}$$

其中  $x_1(t)$  是房间在  $t$  时刻的温度 (单位：摄氏度)， $x_2(t)$  是天花板上方在  $t$  时刻的温度， $t$  是单位为小时的时间。

设  $t = 0$  时  $x_1(t) = x_2(t) = 4$ , 用每一秒钟重新计算的方法计算 24 小时内的房间温度与天花板上温度, 绘图。计算 7 天的温度, 查看周期性。当温度循环变化差距小于 0.05 度时认为开始周期变化了。

## 42.4 统计计算

### 42.4.1 核回归与核密度估计

考虑核回归问题。核回归是非参数回归的一种, 假设变量  $Y$  与变量  $X$  之间的关系为:

$$Y = f(X) + \varepsilon$$

其中函数  $f$  未知。观测到  $X$  和  $Y$  的一组样本  $X_i, Y_i, i=1, \dots, n$  后, 对  $f$  的一种估计为:

$$\hat{f}(x) = \frac{\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) Y_i}{\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right)}$$

其中  $h > 0$  称为窗宽, 窗宽越大, 得到的密度估计越平滑。 $K$  叫做核函数, 一般是一个非负的偶函数, 原点处的函数值最大, 在两侧迅速趋于零。例如正态密度函数, 或所谓双三次函数核:

$$K(x) = \begin{cases} (1 - |x|^3)^3 & |x| \leq 1 \\ 0 & \text{其它} \end{cases}$$

与核回归类似, 可以用核平滑方法估计总体分布密度。设样本为  $Y_1, Y_2, \dots, Y_n$ , 密度估计公式为

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - Y_i}{h}\right)$$

其中  $K(x)$  满足  $\int_{-\infty}^{\infty} K(x) dx = 1$ 。 $h$  是窗宽, 窗宽越大, 估计的曲线越光滑。 $h$  的一种建议公式为  $h = 1.06Sn^{-1/5}$ ,  $S$  为样本标准差。

对以上两个问题进行编程, 其中窗宽  $h$  由用户输入。

## 42.4.2 二维随机模拟积分

设二元函数  $f(x, y)$  定义如下

$$\begin{aligned} f(x, y) &= \exp\{-45(x+0.4)^2 - 60(y-0.5)^2\} \\ &\quad + 0.5 \exp\{-90(x-0.5)^2 - 45(y+0.1)^4\} \end{aligned}$$

(注意其中有一个 4 次方。) 求如下二重定积分

$$I = \int_{-1}^1 \int_{-1}^1 f(x, y) dx dy$$

$f(x, y)$  有两个分别以  $(-0.4, 0.5)$  和  $(0.5, -0.1)$  为中心的峰, 对积分有贡献的区域主要集中在  $(-0.4, 0.5)$  和  $(0.5, -0.1)$  附近, 在其他地方函数值很小, 对积分贡献很小。可以用随机模拟方法估计  $I$  的值。

第一种估计方法是平均值法。设  $X_i, i = 1, 2, \dots, N$  是均匀分布  $U(-1, 1)$  的随机数,  $Y_i, i = 1, 2, \dots, N$  也是均匀分布  $U(-1, 1)$  的随机数, 两者独立, 可估计  $I$  为

$$\hat{I}_1 = \frac{4}{N} \sum_{i=1}^N f(X_i, Y_i).$$

第二种估计方法是重要抽样法。设正态分布  $N(\mu, \sigma^2)$  的密度记为  $p(x; \mu, \sigma^2)$ , 令

$$\begin{aligned} g(x, y) &= 0.5358984p(x; -0.4, 90^{-1})p(y; 0.5, 120^{-1}) \\ &\quad + 0.4641016p(x; 0.5, 180^{-1})p(y; -0.1, 20^{-1}), \\ &\quad -\infty < x < \infty, -\infty < y < \infty, \end{aligned}$$

这是一个二元随机向量的密度, 是两个二元正态密度的混合分布。设  $K_i, i = 1, 2, \dots, N$  是取值于  $\{1, 2\}$  的随机数,  $P(K_i = 1) = 0.5358984$ ,  $P(K_i = 2) = 1 - P(K_i = 1)$ 。当  $K_i = 1$  时, 取  $X_i$  为  $N(-0.4, 90^{-1})$  随机数,  $Y_i$  为  $N(0.5, 120^{-1})$  随机数; 当  $K_i = 2$  时, 取  $X_i$  为  $N(0.5, 180^{-1})$  随机数,  $Y_i$  为  $N(-0.1, 20^{-1})$  随机数, 这样得到的  $(X_i, Y_i), i = 1, 2, \dots, N$  是  $g(x, y)$  的随机数。令

$$\hat{I}_2 = \frac{1}{N} \sum_{i=1}^n \frac{f(X_i, Y_i)}{g(X_i, Y_i)},$$

称为  $I$  的重要抽样法估计。



### 42.4.2.1 编程任务

1. 编写 R 函数估计计算  $\hat{I}_1$ 。重复模拟  $B = 100$  次，得到  $\hat{I}_1$  的  $B$  个值，计算这些值的平均值和标准差。
2. 编写 R 函数估计计算  $\hat{I}_2$ 。重复模拟  $B = 100$  次，得到  $\hat{I}_2$  的  $B$  个值，计算这些值的平均值和标准差。
3. 比较模拟得到的平均值和标准差，验证两者是否基本一致，通过标准差大小比较两种方法的精度。

注意尽量用向量化编程。

### 42.4.3 潜周期估计

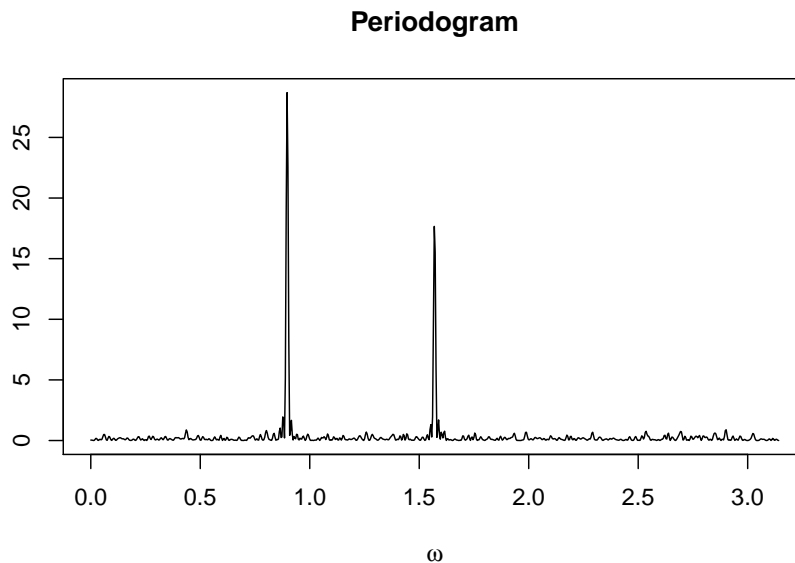
设时间序列  $\{y_t\}$  有如下模型:

$$y_t = \sum_{k=1}^m A_k \cos(\lambda_k t + \phi_k) + x_t, \quad t = 1, 2, \dots$$

其中  $x_t$  为线性平稳时间序列， $\lambda_k \in (0, \pi)$ ,  $k = 1, 2, \dots, m$ 。这样的模型称为潜周期模型。如果有  $\{y_t\}$  的一组样本  $y_1, y_2, \dots, y_n$ ，可以定义周期图函数

$$P(\omega) = \frac{1}{2\pi n} \left| \sum_{t=1}^n y_t e^{-it\omega} \right|^2, \quad \omega \in [0, \pi].$$

这里  $\omega$  是角频率。对于潜周期数据，在  $\lambda_j$  的对应位置  $P(\omega)$  会有尖峰，而且当  $n \rightarrow \infty$  时尖峰高度趋于无穷。下面是一个样例图形。



如下算法可以在  $n$  较大时估计  $m$  和  $\{\lambda_k\}$ : 首先, 对  $\omega_j = \pi j/n, j = 1, 2, \dots, n$  计算  $h_j = P(\omega_j)$ , 求  $\{h_j, j = 1, 2, \dots, n\}$  的 3/4 分位数记为  $q$ 。令  $C = qn^{0.25}$ , 以  $C$  作为分界线, 设  $\{h_j\}$  中大于  $C$  的下标  $j$  的集合为  $J$ , 当  $J$  非空时, 把  $J$  中相邻点分入一组, 但是当两个下标的差大于等于  $n^{0.6}$  时就把后一个点归入新的一组。在每组中, 以该组的  $h_j$  的最大值点对应的角频率  $j\pi/n$  作为潜频率  $\{\lambda_k\}$  中的一个的估计。

用如下 R 程序可以模拟生成一组  $\{y_t\}$  的观测数据:

```
set.seed(1); n <- 500; tt <- seq(n)
m <- 2; lam <- 2*pi/c(4, 7); A <- c(1, 1.2)
y <- A[1]*cos(lam[1]*tt) + A[2]*cos(lam[2]*tt) + rnorm(n)
```

编写 R 程序:

- (1) 编写计算  $P(\omega)$  的函数, 输入  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  和  $\boldsymbol{\omega} = (\omega_1, \omega_2, \dots, \omega_s)^T$ , 输出  $(P(\omega_1), P(\omega_2), \dots, P(\omega_s))$ 。
- (2) 对输入的时间序列样本  $y_1, y_2, \dots, y_n$ , 编写函数用以上描述的算法估计  $m$  和  $\{\lambda_j, j = 1, 2, \dots, m\}$ 。
- (3) 用上述模拟数据测试编写的算法程序。

(4) 进一步地, 用 R 函数 `fft()` 计算  $h_j = P(\pi j/n), j = 1, 2, \dots, n$ 。

(5) 把整个算法用 Rcpp 和 C++ 程序实现。

#### 42.4.4 ARMA(1,1) 模型估计

考虑如下的零均值高斯 ARMA(1,1) 模型:

$$X_t = aX_{t-1} + e_t + be_{t-1}, t = 1, 2, \dots, T$$

其中  $0 < |a| < 1, 0 < |b| < 1, -a \neq b, \{e_t\}$  独立同  $N(0, \sigma_e^2)$  分布。

- (1) 给定观测  $x_1, x_2, \dots, x_T$ , 写出在  $x_0 = 0$  且  $e_0 = 0$  条件下的条件对数似然函数。
- (2) 编写 R 程序, 用条件最大似然估计方法估计参数  $a, b, \sigma_e^2$ 。可使用数值优化程序如 `nlm()` 或 `optim()`。
- (3) 对  $a = 0.5, b = 0.7, \sigma_e = 1, T = 100$ , 模拟生成  $x_t$  的样本并重复模拟  $B = 1000$  次, 据此评估  $\hat{a}, \hat{b}, \hat{\sigma}_e$  的估计精度。模拟可使用 R 的 `arima.sim()` 函数。

若  $x_0, e_0$  已知, 从  $x_1, \dots, x_T$  可递推地计算

$$e_t = x_t - ax_{t-1} - be_{t-1}, t = 1, 2, \dots, T$$

记  $\boldsymbol{\theta} = (a, b, \sigma_e, x_0, e_0)^T$ , 有  $X_t | x_{t-1}, \dots, x_1, \boldsymbol{\theta}$  服从  $N(ax_{t-1} + be_{t-1}, \sigma_e^2)$  条件分布。于是用联合密度的乘积公式可得

$$\begin{aligned} f(x_1, \dots, x_T | \boldsymbol{\theta}) &= f(x_1 | \boldsymbol{\theta}) f(x_2 | x_1, \boldsymbol{\theta}) \dots f(x_T | x_{T-1}, \dots, x_1, \boldsymbol{\theta}) \\ &= \prod_{t=1}^T \text{dnorm}(x_t, ax_{t-1} + be_{t-1}, \sigma_e) \\ &= \prod_{t=1}^T (2\pi)^{-1/2} \sigma_e^{-1} \exp \left\{ -\frac{1}{2} \frac{1}{\sigma_e^2} (x_t - ax_{t-1} - be_{t-1})^2 \right\} \\ &= (2\pi)^{-T/2} \sigma_e^{-T} \exp \left\{ -\frac{1}{2} \frac{1}{\sigma_e^2} \sum_{t=1}^T e_t^2 \right\} \end{aligned}$$

其中  $\text{dnorm}(x, \mu, \sigma)$  表示  $N(\mu, \sigma^2)$  的分布密度。

取  $x_0 = 0, e_0 = 0$ , 给定  $a, b$  后递推计算  $\{e_t\}$  序列, 对数似然函数为

$$l(a, b, \sigma_e) = -T \ln \sigma_e - \frac{1}{2} \frac{1}{\sigma_e^2} \sum_{t=1}^T e_t^2$$

注意其中的  $\{e_t\}$  也与待估参数  $a, b$  有关。

可以用数值优化算法估计求如上的问题的最大值点。显然  $a, b$  的最大值点不依赖于  $\sigma_e$  的取值, 所以可以先求  $\sum_{t=1}^T e_t^2$  的最小值点。这基本是一个最小二乘估计问题, 但是  $e_{t-1}$  也依赖于  $a, b$  的值所以不能直接用线性最小二乘求解。

#### 42.4.5 VAR 模型平稳性

称  $k$  元时间序列  $\mathbf{r}_t$  服从一个 VAR( $p$ ) 模型, 如果

$$\mathbf{r}_t = \boldsymbol{\phi}_0 + \boldsymbol{\Phi}_1 \mathbf{r}_{t-1} + \cdots + \boldsymbol{\Phi}_p \mathbf{r}_{t-p} + \mathbf{a}_t \quad (42.1)$$

其中  $\boldsymbol{\phi}_0$  和  $\{\mathbf{a}_t\}$  同 VAR(1) 的规定,  $\boldsymbol{\Phi}_j$  是  $k$  阶方阵 ( $k = 1, 2, \dots, p$ )。利用向后推移算子 (滞后算子)  $B$  可以将模型写成

$$(\mathbf{I} - \boldsymbol{\Phi}_1 B - \cdots - \boldsymbol{\Phi}_p B^p) \mathbf{r}_t = \boldsymbol{\phi}_0 + \mathbf{a}_t$$

记

$$P(z) = \mathbf{I} - \boldsymbol{\Phi}_1 z - \cdots - \boldsymbol{\Phi}_p z^p \quad (42.2)$$

这是一个从复数  $z$  到  $k$  阶方阵  $P(z)$  的变换,  $P(z)$  的每个元素是关于  $z$  的阶数不超过  $p$  的多项式。称一元多项式函数  $\det(P(z))$  (或记为  $|P(z)|$ ) 为模型的 (逆序) 特征多项式。

如果  $|P(z)| \neq 0, \forall |z| < 1$  ( $z$  为复数), 则模型(42.1)是平稳的。

编写 R 函数, 输入三维 (三个下标) 的系数矩阵数组 `arrcoef`, 返回特征多项式的所有复根。`arrcoef[, , 1]` 保存  $\boldsymbol{\Phi}_1$ , `arrcoef[, , 2]` 保存  $\boldsymbol{\Phi}_2$ , 等等。

作为例子, 考虑如下的三个模型, 对每个计算特征多项式的复根并判断是否平稳:

$$\begin{pmatrix} r_{1t} \\ r_{2t} \end{pmatrix} = \begin{pmatrix} 0.2 & 0.3 \\ -0.6 & 1.1 \end{pmatrix} \begin{pmatrix} r_{1,t-1} \\ r_{2,t-1} \end{pmatrix} + \begin{pmatrix} a_{1t} \\ a_{2t} \end{pmatrix}$$

$$\begin{pmatrix} x_{1t} \\ x_{2t} \end{pmatrix} = \begin{pmatrix} 0.5 & -1.0 \\ -0.25 & 0.5 \end{pmatrix} \begin{pmatrix} x_{1,t-1} \\ x_{2,t-1} \end{pmatrix} + \begin{pmatrix} a_{1t} \\ a_{2t} \end{pmatrix}$$

$$\mathbf{r}_t = \begin{pmatrix} 0.39 & 0.10 & 0.05 \\ 0.35 & 0.34 & 0.47 \\ 0.49 & 0.24 & 0.24 \end{pmatrix} \mathbf{r}_{t-1} + \begin{pmatrix} 0.06 & 0.11 & 0.02 \\ -0.19 & -0.18 & -0.01 \\ -0.31 & -0.13 & 0.09 \end{pmatrix} \mathbf{r}_{t-2} + \mathbf{a}_t$$

#### 42.4.6 贮存可靠性评估

设某种设备的贮存寿命为随机变量  $X$ , 服从指数分布  $\text{Exp}(\theta)$ ,  $EX = \theta > 0$ . 假设有  $n$  台此种设备分别在时间  $t_1, t_2, \dots, t_n$  进行了试验, 第  $i$  台试验成功用  $\delta_i = 1$  表示, 试验失效用  $\delta_i = 0$  表示. 把  $n$  次试验的结果写成

$$Z_n = \begin{pmatrix} t_1 & \delta_1 \\ t_2 & \delta_2 \\ \vdots & \vdots \\ t_n & \delta_n \end{pmatrix} \quad (1)$$

这里  $\delta_1, \delta_2, \dots, \delta_n$  认为是随机变量,  $t_1, t_2, \dots, t_n$  是固定的时间.  $\delta_i$  服从两点分布  $B(1, \exp(-t_i/\theta))$ .  $Z_n$  取某个特定组合的概率为

$$P(Z_n) = \prod_{i=1}^n \left\{ [\exp(-t_i/\theta)]^{\delta_i} [1 - \exp(-t_i/\theta)]^{1-\delta_i} \right\}. \quad (2)$$

要利用观测数据  $Z_n$  估计  $\theta$  的置信度为  $1 - \alpha$  的置信下限, 可以用陈家鼎《生存分析与可靠性》中的样本空间排序法. 注意到  $Z_n$  中每个  $\delta_i$  可以取 1 或 0, 所以  $Z_n$  的所有不同取值共有  $2^n$  个, 记这些所有不同取值为  $\mathcal{D}$ , 在  $\mathcal{D}$  中定义如下的序: 设  $Z'_n$  与  $Z_n$  都是  $\mathcal{D}$  中的试验结果, 称  $Z'_n$  不次于  $Z_n$ , 并记作  $Z'_n \succeq Z_n$ , 如果如下两个条件之一成立:

- (1)  $\sum_{i=1}^n \delta'_i > \sum_{i=1}^n \delta_i$ ;
- (2)  $\sum_{i=1}^n \delta'_i = \sum_{i=1}^n \delta_i$ , 但  $\sum_{i=1}^n \delta'_i t_i \geq \sum_{i=1}^n \delta_i t_i$ .

记

$$G(\theta) = \sum_{Z'_n \succeq Z_n} P_\theta(Z'_n) = \sum_{Z'_n \succeq Z_n} \prod_{i=1}^n \left\{ [\exp(-t_i/\theta)]^{\delta'_i} [1 - \exp(-t_i/\theta)]^{1-\delta'_i} \right\}. \quad (3)$$

这是  $\theta$  的严格单调增函数，求解如下的方程

$$G(\theta) - \alpha = 0 \quad (4)$$

得到解  $\underline{\theta}$  是  $\theta$  的置信度为  $1 - \alpha$  的置信下限。

当所有  $n$  次试验都失效时，恒有  $G(\theta) = 1$ ，方程 (4) 无解，取  $\underline{\theta} = 0$ 。

当  $n$  次试验都没有失效，即失效数  $n - \sum_{i=1}^n \delta_i = 0$  时，方程为

$$\exp\left(-\frac{1}{\theta} \sum_{i=1}^n t_i\right) - \alpha = 0,$$

解得

$$\underline{\theta} = \frac{\sum_{i=1}^n t_i}{\ln \frac{1}{\alpha}}.$$

当失效数为 1 时，最多仅有  $n$  个结果不次于  $Z_n$ ，很容易可以计算  $G(\theta)$  的值。在  $n$  较大而且  $Z_n$  中失效数较多时， $G(\theta)$  的求和 (3) 中项数很多，计算量很大。当  $n = 10$  时， $\mathcal{D}$  约有一千项，当  $n = 20$  时， $\mathcal{D}$  约有一百万项，还在可以穷举计算的范围内。但是，当  $n \geq 30$  时， $\mathcal{D}$  就有 10 亿项，每计算一次  $G(\theta)$  都需要很长时间。

针对  $n$  不太大的情形，可以用精确的公式 (3) 计算  $G(\theta)$  并用 (4) 求解  $\underline{\theta}$ ，编写这个问题的纯 R 程序版本，输入一组  $Z_n$  值后（包括试验时间和试验结果），输出用 (3) 和 (4) 求解  $\underline{\theta}$  得到的置信下限  $\underline{\theta}$  的值。在文件 `store-reliab-data.csv` 中已经生成了 10 组模拟数据，对这 10 组模拟数据计算相应的  $\underline{\theta}$  的值。在此数据文件中，`testid` 相同的行属于同一次试验的不同设备的时间和结果。

这个程序中用到比较多的循环，考虑把主要计算部分用 Rcpp 包和 C++ 代码实现，看能够提高效率多少倍。

另一种计算  $G(\theta)$  的方法是用随机模拟方法估计  $G(\theta)$  的值然后求解 (4)。随机模拟方法计算简单而且不受失效个数多少的影响，但是有随机误差，在求解 (4) 要考虑到随机误差的影响。随机模拟方法如下。取模拟次数  $N$ ，对给定  $\theta$ ，为了计算  $G(\theta)$ ，模拟生成  $N$  组独立的寿命  $(X_1^{(i)}, X_2^{(i)}, \dots, X_n^{(i)})$ ， $i = 1, 2, \dots, N$ ，

其中每个  $X_j^{(i)}$  服从  $\text{Exp}(\theta)$  分布, 各分量相互独立。计算  $\delta_j^{(i)} = I_{\{X_j^{(i)} > t_j\}}$ , 记

$$Z_n^{(i)} = \begin{pmatrix} t_1 & \delta_1^{(i)} \\ t_2 & \delta_2^{(i)} \\ \vdots & \vdots \\ t_n & \delta_n^{(i)} \end{pmatrix}$$

用

$$\frac{1}{N} \sum_{i=1}^N I_{\{Z_n^{(i)} \succeq Z_n\}}$$

来估计  $G(\theta)$ 。

## 42.5 数据处理

### 42.5.1 小题分题型分数汇总

考虑中学某科的一次考试, 考卷各小题情况汇总在如下的用逗号符分隔的文本文件 `subscore-subtype.csv` 中:

序号, 题型

1, 选择题

2, 选择题

3, 选择题

4, 选择题

5, 选择题

6, 选择题

7, 选择题

8, 选择题

9, 选择题

10, 选择题

11, 简答题

12, 简答题

13, 填空题

14, 简答题

15, 简答题

16, 简答题

17, 简答题

18, 写作

读入此数据为 R 数据框，只要用如下程序：

```
dm <- read.csv('subscore-subtype.csv', header=TRUE,
 stringsAsFactors=FALSE)
```

结果显示如下：

```
knitr::kable(dm)
```

| 序号 | 题型  |
|----|-----|
| 1  | 选择题 |
| 2  | 选择题 |
| 3  | 选择题 |
| 4  | 选择题 |
| 5  | 选择题 |
| 6  | 选择题 |
| 7  | 选择题 |
| 8  | 选择题 |
| 9  | 选择题 |
| 10 | 选择题 |
| 11 | 简答题 |
| 12 | 简答题 |
| 13 | 填空题 |
| 14 | 简答题 |
| 15 | 简答题 |
| 16 | 简答题 |
| 17 | 简答题 |
| 18 | 写作  |

设部分学生的小题分录入到如下的用逗号分隔的文本文件 `subscore-subscore.csv` 中：

学号, Y1, Y2, Y3, Y4, Y5, Y6, Y7, Y8, Y9, Y10, Y11, Y12, Y13, Y14, Y15, Y16, Y17, Y18



```
1138010104,3,3,3,3,0,3,7.5,2.5,3.5,6,5,4,2.5,5.5,0,4,2,45.5
1138010108,3,0,3,0,3,0,6,3,4,4,2,5,2,5,6,4,2,48.5
1138010114,3,3,3,3,0,3,5.5,3.5,4,6,2,5,2,5.5,6,1,2,44.5
1138010128,3,3,3,0,0,0,8.5,3.5,2.5,4,5,5,0,4.5,6,4,2,45
1138010126,3,3,3,3,3,3,7,0,4,6,5,5,2.5,4.5,6,4.5,2,44
```

用如下 R 程序读入小题分数数据为 R 数据框:

```
ds <- read.csv('subscore-subscore.csv', header=TRUE,
 stringsAsFactors=FALSE)
```

结果显示如下:

```
knitr::kable(ds)
```

| 学号         | Y1 | Y2 | Y3 | Y4 | Y5 | Y6 | Y7  | Y8  | Y9  | Y10 | Y11 | Y12 | Y13 | Y14 | Y15 | Y16 |
|------------|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 1138010104 | 3  | 3  | 3  | 3  | 0  | 3  | 7.5 | 2.5 | 3.5 | 6   | 5   | 4   | 2.5 | 5.5 | 0   | 4.0 |
| 1138010108 | 3  | 0  | 3  | 0  | 3  | 0  | 6.0 | 3.0 | 4.0 | 4   | 2   | 5   | 2.0 | 5.0 | 6   | 4.0 |
| 1138010114 | 3  | 3  | 3  | 3  | 0  | 3  | 5.5 | 3.5 | 4.0 | 6   | 2   | 5   | 2.0 | 5.5 | 6   | 1.0 |
| 1138010128 | 3  | 3  | 3  | 0  | 0  | 0  | 8.5 | 3.5 | 2.5 | 4   | 5   | 5   | 0.0 | 4.5 | 6   | 4.0 |
| 1138010126 | 3  | 3  | 3  | 3  | 3  | 3  | 7.0 | 0.0 | 4.0 | 6   | 5   | 5   | 2.5 | 4.5 | 6   | 4.5 |

在数据框 dm 中有每个小题的题型信息, 在数据框 ds 中有每个学生的每个小题的分数。从这两个数据框, 汇总计算每个学生的题型分, 即每个学生选择题共考多少分, 简答题共考多少分, 等等。

要注意的是, 最终的程序应该写成一个函数, 其中的计算不依赖于具体的题型名称、小题个数、小题与题型如何对应, 只要输入 dm 和 ds 两个数据框就可以进行统计。

如果没有这样的通用性要求, 这个问题就可以这样简单解决:

```
resm <- data.frame(
 '学号'=ds[, '学号'],
 '选择题'=rowSums(ds[, paste('Y', 1:10, sep='')]),
 '简答题'=rowSums(ds[, paste('Y', c(11,12,14:17), sep='')]),
 '填空题'=ds[, 'Y13'],
 '作文'=ds[, 'Y18']
)
```

```
knitr::kable(resm[order(resm[, '学号']),], row.names=FALSE)
```

| 学号         | 选择题  | 简答题  | 填空题 | 作文   |
|------------|------|------|-----|------|
| 1138010104 | 34.5 | 20.5 | 2.5 | 45.5 |
| 1138010108 | 26.0 | 24.0 | 2.0 | 48.5 |
| 1138010114 | 34.0 | 21.5 | 2.0 | 44.5 |
| 1138010126 | 35.0 | 27.0 | 2.5 | 44.0 |
| 1138010128 | 27.5 | 26.5 | 0.0 | 45.0 |

## 42.6 文本处理

### 42.6.1 用 R 语言下载处理《红楼梦》htm 文件

网上许多资源是 html 格式的文本文件。比如，《红楼梦》在许多网站可以浏览，是在浏览器中按章节浏览。我们希望将其下载到本地，并转换为 txt 格式，在手机或者电子书阅读器中阅读。

这个任务涉及到 R 的文件访问，字符型连接，正则表达式，中文编码问题。

设下载网站主页是 <http://www.xiexingcun.net/honglouloumeng/index.html>，各个章节的文件名是 01.htm 到 99.htm，100.htm 到 120.htm。

已下载的文件在如下链接中：[honglouloumeng.zip](#)

将这些文件转换成 txt 格式，然后合并成一个 txt 文件。

要求：每一段仅用一行，中间不换行；不同段落之间用换行分开；每一回目开始有回目标题，标题前后空行，标题格式为“第 xxX 回标题内容”。

## 参考文献

- Becker, R. A. and Chambers, J. M. (1984). *S: An Interactive Environment for Data Analysis and Graphics*. Wadsworth Advanced Books Program, Belmont CA.
- Becker, R. A., Chambers, J. M., and Wilks, A. R. (1988). *The New S Language*. Chapman and Hall, New York.
- Chambers, J. M. (2008). *Software for Data Analysis: Programming with R*. Springer.
- Chambers, J. M. and Hastie, T. (1992). *Statistical Models in S*. Chapman and Hall, New York.
- Cleveland, W. S. (1979). Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, 74:829–36.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning with Applications in R*. Springer.
- Kabacoff, R. I. (2012). *R 语言实战*. 人民邮电出版社.
- Knuth, D. E. (1984). Literate programming. *The Computer Journal*, 27(2):97–111.
- Venables, W. N. and Ripley, B. D. (2002). *Modern Applied Statistics with S*. Springer, 4th ed. edition.
- Wilke, C. O. (2019). *Fundamentals of Data Visualization: A Primer on Making Informative and Compelling Figures*. O’Reilly Media.

- Xie, Y. (2017). *bookdown: Authoring Books and Technical Documents with R Markdown*.
- Xie, Y., Allaire, J. J., and Golemund, G. (2019). *R Markdown: The Definitive Guide*. CRC Press.
- Yihui Xie, Amber Thomas, A. P. H. (2017). *blogdown: Creating Websites with R Markdown*.