

МС-26 Теоретический материал
Критерии независимости и однородности χ^2 -Пирсона,
критерии однородности Колмогорова-Смирнова.
Сравнение генеральных средних и генеральных дисперсий
двух и более нормальных совокупностей

Критерий независимости χ^2

Имеются две дискретные случайные величины X и Y .

Требуется проверить **гипотезу об их независимости**.

Пусть различаются r значений случайной величины X (обозначим их x_1, \dots, x_r) и s значений случайной величины Y (обозначим их y_1, \dots, y_s).

Через k_{ij} обозначим общее количество таких элементов выборки, в которых X принимает значение x_i , а Y — значение y_j .

Тогда

$$\sum_{i=1}^r \sum_{j=1}^s k_{ij} = n,$$

где n — объем выборки. Введем также обозначения:

$$v_i = \sum_{j=1}^s k_{ij}; \mu_j = \sum_{i=1}^r k_{ij}.$$

В рассматриваемом случае результаты наблюдений удобно оформлять в виде таблицы, называемой **таблицей сопряженности признаков**:

X/Y	y_1	y_2	...	y_s	
x_1	k_{11}	k_{12}	...	k_{1s}	v_1
x_2	k_{21}	k_{22}	...	k_{2s}	v_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
x_r	k_{r1}	k_{r2}	...	k_{rs}	v_r
	μ_1	μ_2	...	μ_s	n

Пусть далее $p_{ij} = P(\{X = x_i, Y = y_j\})$; $p_i = P(\{X = x_i\})$; $q_j = P(\{Y = y_j\})$.

$H_0: p_{ij} = p_i q_j, i = 1, \dots, r, j = 1, \dots, s.$

$H_1: p_{ij} \neq p_i q_j$ для некоторых $i = 1, \dots, r, j = 1, \dots, s.$

Рассматривается следующая статистика:

$$\chi^2 = n \left(\sum_{i=1}^r \sum_{j=1}^s \frac{k_{ij}^2}{v_i \mu_j} - 1 \right).$$

Гипотеза H_0 отклоняется, если вычисленное по выборочным данным значение статистики $\chi_{\text{набл}}^2$ удовлетворяет неравенству:

$$\chi_{\text{набл}}^2 > \chi_{\alpha; (r-1)(s-1)}^2.$$

Критерий однородности χ^2

Проверяется гипотеза о том, что две выборки принадлежат одной генеральной совокупности.

Данные должны быть представлены в виде интервального статистического ряда.

Имеются выборка объема n_1 из генеральной совокупности X_1 и выборка объема n_2 из генеральной совокупности X_2 ; l — количество интервалов группировки (одинаковое для обеих выборок); μ_i и ν_i — количество попаданий в i -й интервал группирования, соответственно, первой и второй выборок, $i = 1, 2, \dots, l$; уровень значимости α .

Пусть $F_j(x)$ — функция распределения случайной величины X_j , $j = 1, 2$.

Проверяется гипотеза

$$H_0: F_1(x) = F_2(x), x \in \mathbb{R},$$

$$H_1: F_1(x) \neq F_2(x), \text{ для некоторых } x \in \mathbb{R}.$$

В случае совпадения объемов выборок: $n_1 = n_2 = n$

статистика вычисляется по формуле

$$\chi^2 = \sum_{i=1}^l \frac{(\mu_i - \nu_i)^2}{\mu_i + \nu_i}.$$

Критическое значение статистики: $\chi_{\alpha; l-1}^2$.

Гипотеза H_0 отклоняется, если вычисленное по выборочным данным значение статистики $\chi_{\text{набл}}^2$ удовлетворяет неравенству:

$$\chi_{\text{набл}}^2 > \chi_{\alpha; l-1}^2.$$

Статистика критерия имеет следующий вид:

$$\chi^2 = n_1 n_2 \sum_{i=1}^l \frac{\left(\frac{\mu_i}{n_1} - \frac{\nu_i}{n_2}\right)^2}{\mu_i + \nu_i}.$$

Критерий однородности Колмогорова-Смирнова

Имеются две выборки — объема n_1 из генеральной совокупности X_1 и объема n_2 из генеральной совокупности X_2 .

Предполагается, что случайные величины X_j — непрерывные с функциями распределения $F_j(x)$, $j = 1, 2$.

$$H_0: F_1(x) = F_2(x), x \in \mathbb{R},$$

$$H_1: F_1(x) \neq F_2(x), \text{ для некоторых } x \in \mathbb{R}.$$

Проверка гипотезы производится по следующей схеме:

1. По имеющимся выборкам находятся эмпирические функции распределения $F_1^*(x)$ и $F_2^*(x)$.
2. Рассматривается статистика следующего вида:

$$D = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} \cdot \max_x |F_1^*(x) - F_2^*(x)|.$$

3. По таблицам распределения Колмогорова определяется величина k_α — 100α -процентная точка распределения Колмогорова уровня α .

4. Гипотеза H_0 отклоняется на уровне значимости α , если вычисленное по выборочным данным значение статистики $D_{\text{набл}}$ удовлетворяет неравенству:

$$D_{\text{набл}} > k_\alpha.$$

Замечание. Критерий Колмогорова–Смирнова применяется при $n_1, n_2 \geq 50$.

Дисперсионный анализ — это статистический метод анализа результатов наблюдений, зависящих от различных одновременно действующих факторов, основанный на сравнении оценок дисперсий соответствующих групп выборочных данных.

Под **фактором** понимают различные, независимые, качественные показатели, влияющие на изучаемые признаки.

Признаки, изменяющиеся под воздействием тех или иных факторов, называют **результативными**.

Сущность метода дисперсионного анализа заключается в измерении отдельных дисперсий (общая, факторная, остаточная), и дальнейшем определении силы влияния изучаемых факторов (оценки роли каждого из факторов, либо их совместного влияния) на результативный признак.

Классификация методов дисперсионного анализа



***ANOVA** - Analysis of Variance

***MANOVA** – Multivariate Analysis of Variance

Предположения для использования однофакторного дисперсионного анализа:

Чтобы результаты однофакторного дисперсионного анализа были достоверными, должны выполняться следующие допущения:

- 1. Нормальность.** Каждая выборка была взята из нормально распределенной популяции.
- 2. Равные дисперсии** — дисперсии совокупностей, из которых взяты выборки, равны (можно использовать **тест Бартлетта**, чтобы проверить это предположение).
- 3. Независимость.** Наблюдения в каждой группе независимы друг от друга, а наблюдения внутри групп были получены путем случайной выборки.

Однофакторный дисперсионный анализ основан на том, что сумму квадратов отклонений статистического комплекса возможно разделить на компоненты:

$$SST = SSA + SSW,$$

где SST - общая сумма квадратов отклонений, SSA - объяснённая влиянием фактора А сумма квадратов отклонений, SSW - необъяснённая сумма квадратов отклонений или сумма квадратов отклонений ошибки.

- 1. Межгрупповая вариация SSA** – вариация между средним каждой группы и общим средним значением всей выборки
- 2. Внутригрупповая вариация SSW** – вариация между каждым объектом исследования группы и средним значением соответствующей группы

Компоненты дисперсии	Сумма квадратов	Число степеней свободы	Оценки дисперсии
Межгрупповая (SSA) Sum of Squares Among Groups	$SSA = \sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2$	$k - 1$	$MSA = \frac{SSA}{k - 1}$ (Факторная дисперсия)
Внутригрупповая (SSW) Sum of Squares Within Groups	$SSW = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$	$n - k$	$MSW = \frac{SSW}{n - k}$ (Остаточная дисперсия)
Общая ($SST = SSA + SSW$) Total Sum of Squares	$SST = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2$	$n - 1$	$MST = \frac{SST}{n - 1}$

Чтобы провести однофакторный дисперсионный анализ данных статистического комплекса, нужно найти **фактическое отношение Фишера** - отношение дисперсии, объяснённой влиянием фактора (**межгрупповой, факторной**), и необъяснённой дисперсии (**внутригрупповой, остаточной**).

Тогда **дисперсионное отношение** это $F_{\text{набл}} = \frac{MSA}{MSW}$.

Указанное дисперсионное отношение **тем больше, чем сильнее** разнятся уровни фактора по воздействию на групповые средние.

Проверка гипотезы о совпадении нескольких генеральных средних методом дисперсионного анализа

Пусть $\vec{X}_i = (X_{i1}, \dots, X_{in_i})$ — выборка объема n_i из $N(\mu_i, \sigma^2)$, где $i = 1, \dots, k$.

Предположим также, что $n = n_1 + \dots + n_k$ случайных величин

$$X_{11}, \dots, X_{1n_1}; X_{21}, \dots, X_{2n_2}; \dots; X_{k1}, \dots, X_{kn_k}$$

независимы в совокупности. Таким образом, выборки $\vec{X}_1, \dots, \vec{X}_k$ **независимы** и получены из нормальных распределений с одинаковой дисперсией σ^2 и, возможно, различными средними μ_1, \dots, μ_k . **Гипотеза о равенстве всех средних** одновременно записывается как

$$H_0: \mu_1 = \dots = \mu_k, H_1: (\exists i, j) \mu_i \neq \mu_j.$$

Заметим, что при верной H_0 генеральные распределения совпадают:

$$N(\mu_1, \sigma^2) = \dots = N(\mu_k, \sigma^2).$$

Рассмотрим объединенную выборку объема $n = n_1 + \dots + n_k$

$$\vec{X} = (X_{11}, \dots, X_{1n_1}; X_{21}, \dots, X_{2n_2}; \dots; X_{k1}, \dots, X_{kn_k}).$$

Интерпретируя выборки $\vec{X}_1, \dots, \vec{X}_k$ как группы, на которые разбита совокупность \vec{X} , введем обозначения, аналогичные тем, что использовались при изучении межгрупповой дисперсии:

$$\bar{X}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij} — \text{выборочное среднее в } i\text{-й совокупности};$$

$$\hat{\sigma}_i^2 = \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 — \text{выборочная дисперсия в той же выборке};$$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^k \bar{X}_i n_i = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} X_{ij} — \text{выборочное среднее в объединенной выборке } \vec{X};$$

$$\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^k \hat{\sigma}_i^2 n_i — \text{средняя групповая дисперсия } (SSW/n);$$

$$\delta^2 = \frac{1}{n} \sum_{i=1}^k (\bar{X}_i - \bar{X})^2 n_i — \text{межгрупповая дисперсия } (SSA/n);$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2 — \text{выборочная дисперсия признака в объединенной выборке } \vec{X}.$$

Известно, что выборочную дисперсию $\hat{\sigma}^2$ можно представить в виде суммы

$$\hat{\sigma}^2 = \bar{\sigma}^2 + \delta^2 = (SSA + SSW)/n = SST/n,$$

где слагаемое $\bar{\sigma}^2$ характеризует среднюю изменчивость признака в каждой выборке $\vec{X}_1, \dots, \vec{X}_k$, а слагаемое δ^2 характеризует разброс выборочных средних $\bar{X}_1, \dots, \bar{X}_k$.

Критерий проверки H_0 против H_1 использует так называемое **отношение Фишера**:

$$F = \frac{\frac{1}{k-1} \sum_{i=1}^k (\bar{X}_i - \bar{X})^2 n_i}{\frac{1}{n-k} \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2} = \frac{\frac{1}{k-1} n \delta^2 / s^2}{\frac{1}{n-k} n \bar{\sigma}^2 / s^2} = \frac{MSA}{MSW}.$$

Можно доказать, что $F \sim F(k-1, n-k)$, где $F(k-1, n-k)$ — **распределение Фишера с $k-1$ и $n-k$ степенями свободы**.

Для проверки H_0 с уровнем значимости α применяется критерий с критической областью $F > F_\alpha(k-1, n-k)$, где $F_\alpha(k-1, n-k)$ — верхняя процентная точка распределения $F(k-1, n-k)$.

Замечание. Если $F < 1$, то следует сразу принять гипотезу H_0 , поскольку $F_{кр}$ всегда больше единицы.

Если гипотеза о равенстве средних не подтверждается, имеет смысл оценивать величины μ_i по отдельности. Получаем для них следующие **доверительные интервалы (с надежностью γ)**:

$$\bar{X}_i - \delta < \mu_i < \bar{X}_i + \delta,$$

где $\delta = \sqrt{\frac{MSW}{n_i}} t_\gamma$, t_γ – критическая точка распределения Стьюдента с $n - k$ степенями свободы (для двусторонней области)

scipy.stats.t.ppf((1+gamma)/2,n-k)

Одним из условий применения дисперсионного анализа является равенство генеральных групповых дисперсий $\sigma_i^2 = \sigma_0^2, i = 1, \dots, k$.

Сравнение нескольких дисперсий нормальных генеральных совокупностей по выборкам различного объема. Критерий Бартлетта

Пусть генеральные совокупности X_1, X_2, \dots, X_l распределены нормально. Из этих совокупностей извлечены независимые выборки, вообще говоря различных объемов n_i (некоторые n_i могут быть одинаковыми; если все выборки имеют одинаковые объём, то предпочтительнее пользоваться **критерием Кочрена**).

$$H_0: \sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2$$

1. Найдем исправленные выборочные дисперсии $s_i^2, i = 1, \dots, l$.
2. Найдем $k_i = n_i - 1$ – число степеней свободы дисперсии $s_i^2, i = 1, \dots, l$ и сумму степеней свободы $k = \sum_{i=1}^l k_i$.
3. Найдем среднюю арифметическую исправленных дисперсий, взвешенную по числам степеней свободы:

$$\overline{s^2} = \frac{\sum_{i=1}^l k_i s_i^2}{k}.$$

4. Случайная величина (**критерий Бартлетта**) $B = \frac{V}{C}$,

где

$$V = 2,303 \left[k \lg \overline{s^2} - \sum_{i=1}^l k_i \lg s_i^2 \right]; C = 1 + \frac{1}{3(l-1)} \left[\sum_{i=1}^l \frac{1}{k_i} - \frac{1}{k} \right],$$

которая при условии справедливости гипотезы об однородности дисперсий распределена приближенно как χ^2 с $l - 1$ степенями свободы, если объем каждой выборки $n_i > 3, i = 1, \dots, l$.

Правило. Для того чтобы при заданном уровне значимости α проверить **нулевую гипотезу об однородности дисперсий нормальных совокупностей**, надо вычислить наблюдаемое значение критерия Бартлетта $B_{набл.} = V/C$, и критическую точку $\chi_{кр.}^2(\alpha; l - 1)$ правосторонней критической области. Если $B_{набл.} < \chi_{кр.}^2$ — нет оснований отвергнуть нулевую гипотезу. Если $B_{набл.} > \chi_{кр.}^2$ — нулевую гипотезу отвергают.

Замечание 1. Не следует торопиться вычислять постоянную C . Сначала надо найти V и сравнить с $\chi^2_{кр.}$. если окажется, что $V < \chi^2_{кр.}$, то подалвно (так как $C > 1$) $B_{набл.} = V/C < \chi^2_{кр.}$ и, следовательно, C вычислять не нужно. Если же $V > \chi^2_{кр.}$, то надо вычислить C и затем сравнить $B_{набл.}$ с $\chi^2_{кр.}$.

Замечание 2. Критерий Бартлетта чувствителен к отклонениям распределений от нормального, поэтому к выводам следует относиться с осторожностью.

Замечание 3. При условии однородности дисперсий **в качестве оценки генеральной дисперсии** принимают среднюю арифметическую исправленных дисперсий, взвешенную по числам степеней свободы:

Критерий Бартлетта также называют **гипотезой об однородности дисперсий**.

<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.bartlett.html>

Предположим, что фактор A влияет на результативный признак X . Для измерения степени этого влияния используют **выборочный коэффициент детерминации**, равный

$$\eta^2 = \frac{SSA}{SST},$$

который показывает, какую долю выборочной дисперсии составляет дисперсия групповых средних, иначе говоря, какая доля общей дисперсии объясняется зависимостью результативного признака X от фактора A .

Для выполнения **однофакторного дисперсионного анализа в Python** применяют функцию

f_oneway()

из библиотеки **SciPy**.

https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.f_oneway.html