

# Modelling at-bats in baseball using the generalised Pareto distribution

A report created using bibat version 0.0.7

Teddy Groves

I used to do a lot of statistical analyses of sports data where there was a latent parameter for the player's ability. You can see an example [here](#).

It was a natural choice to use a Bayesian hierarchical model where the abilities have a location/scale type distribution with support on the whole real line, and in fact this worked pretty well! You would see the kind of players at the top and bottom of the ability scale that you would expect.

Still, there were a few problems. In particular the data were typically unbalanced because better players tended to produce more data than worse players. The result of this was that my models would often inappropriately think the bad players were like the good players: they would not only tend to be too certain about the abilities of low-data players, but also be biased, thinking that these players are probably a bit better than they actually are. I never came up with a good way to solve this problem, despite trying a lot of things!

Even though I don't work with sports data very much any more, the problem still haunted me, so when I read [this great case study](#) about geomagnetic storms it gave me an idea for yet another potential solution.

The idea was this: just as data about intense storms that cause electricity problems tell us about a tail of the bigger solar mag-

netism distribution, maybe data about professional sportspeople is best thought about as coming from a tail of the general sports ability distribution. If so, maybe something like the [generalised pareto distribution](#) might be better than the bog standard normal distribution for describing the pros' abilities.

I thought I'd test this out with some sports data, and luckily there is a really nice baseball example on [the Stan case studies website](#), complete with [data from the 2006 Major league season])<https://github.com/stan-dev/example-models/blob/master/knitr/pool-binary-trials/baseball-hits-2006.csv>). After I posted some early results [on the Stan discourse forum](#), other users suggested that it might be interesting to model similar data from different seasons. This data can now be found quite easily on the [baseball databank](#).

With a few datasets and at least two different statistical models to consider, the full analysis looked like it would be a little too big to fit in a trivial file structure, so it seemed like a good opportunity to showcase my batteries-included Bayesian analysis template package [bibat](#).

The rest of this vignette describes how I used bibat to (relatively) painlessly see if my generalised Pareto distribution idea would work.

Check out the analysis's [github repository](#) for all the details.

## Setup

First I installed bibat in my global Python 3.11 environment with this command:

```
$ pip install bibat
```

Next I started bibat's wizard like this:

```
bibat
```

After I answered the wizard's questions bibat created a new folder called `baseball` that looked like this:

```
baseball
  CODE_OF_CONDUCT.md
  LICENSE
  Makefile
  README.md
  bibat_version.txt
  data
    raw
      raw_measurements.csv
      readme.md
  docs
    bibliography.bib
    img
      example.png
      readme.md
    report.qmd
  inferences
    fake_interaction
      config.toml
    interaction
      config.toml
    no_interaction
      config.toml
  investigate.ipynb
  plots
    posterior_ll_comparison.png
    posterior_predictive_comparison.png
  prepare_data.py
  pyproject.toml
  requirements-tooling.txt
  requirements.txt
  sample.py
  src
    __init__.py
    data_preparation_functions.py
    inference_configuration.py
    prepared_data.py
    readme.md
    stan
      custom_functions.stan
```

```

    multilevel-linear-regression.stan
    readme.md
    stan_input_functions.py
    util.py
tests
    test_integration
        test_data_preparation.py
    test_unit
        test_inference_configuration.py
        test_util.py

```

This folder implements bibat's example analysis - a comparison of linear regression with two different design matrices. To check that everything was working correctly I ran the analysis like this:

```

$ cd baseball
$ make analysis

```

Some cogs turned, some new lines appeared in my terminal and some new files were created - nice, the setup worked!

## Getting raw data

To fetch raw data from the internet, I wrote a new script called `fetch_data.py`:

```

"""Fetch raw data from the internet."""

import os

import pandas as pd

URLS = {
    "2006": "https://raw.githubusercontent.com/stan-dev/"
    "example-models/master/knitr/pool-binary-trials/baseball-hits-2006.csv",
    "bdb-main": "https://raw.githubusercontent.com/chadwickbureau/"

```

```

    "baseballdatabank/master/core/Batting.csv",
    "bdb-post": "https://raw.githubusercontent.com/chadwickbureau/"
    "baseballdatabank/master/core/BattingPost.csv",
    "bdb-apps": "https://raw.githubusercontent.com/chadwickbureau/"
    "baseballdatabank/master/core/Appearances.csv"
}
OUT_FILES = {
    "2006": os.path.join("data", "raw", "2006.csv"),
    "bdb-main": os.path.join("data", "raw", "bdb-main.csv"),
    "bdb-post": os.path.join("data", "raw", "bdb-post.csv"),
    "bdb-apps": os.path.join("data", "raw", "bdb-apps.csv"),
}

if __name__ == "__main__":
    for name, url in URLS.items():
        print(f"Fetching {name} data from {url}")
        data = pd.read_csv(url, comment="#")
        print(f"Writing {name} data to {OUT_FILES[name]}")
        data.to_csv(OUT_FILES[name])

```

To get the files I ran the script:

```

$ source .venv/bin/activate
$ python fetch_data.py

```

Since this worked, I added a new makefile target for the raw data files:

```

RAW_DATA = data/raw/2006.csv data/raw/bdb-main.csv data/raw/bdb-post.csv data/raw/bdb-apps.c
...
$(RAW_DATA): $(ACTIVATE_VENV) && python fetch_data.py

```

Finally I removed the example analysis's raw data:

```

$ rm data/raw/raw_measurements.csv

```

## Preparing the data

The first step in preparing data is to decide what prepared data looks like for the purposes of our analysis. Bibat provides dataclasses called `PreparedData` and `MeasurementsDF` in the file `src/prepared_data.py` which can help get us started with this.

As it happens, prepared data looks very similar in our analysis and the example. All we need to do is change the `MeasurementsDF` definition a little<sup>1</sup>:

<sup>1</sup> note that this class uses [pandera](#), a handy library for defining what a pandas dataframe should look like

```
from typing import Optional
import pandera as pa

# ...

class MeasurementsDF(pa.SchemaModel):
    """A PreparedData should have a measurements dataframe like this.

    Other columns are also allowed!
    """

    player_season: pa.typing.Series[str]
    season: pa.typing.Series[str]
    n_attempt: pa.typing.Series[int] = pa.Field(ge = 1)
    n_success: pa.typing.Series[int] = pa.Field(ge = 0)
```

Next we can write some functions that create `PreparedData` objects. These should live in the file `data_preparation_functions.py`. In this case we write a couple of data preparation functions: `prepare_data_2006` and `prepare_data_bdb`:

```
"""Provides functions prepare_data_x.

These functions should take in a dataframe of measurements and return a
PreparedData object.

"""
```

```

import pandas as pd
from src.prepared_data import PreparedData

def prepare_data_2006(measurements_raw: pd.DataFrame) -> PreparedData:
    """Prepare the 2006 data."""
    measurements = measurements_raw.rename(
        columns={"K": "n_attempt", "y": "n_success"}
    ).assign(
        season="2006",
        player_season=lambda df: [f"2006-player-{i+1}" for i in range(len(df))],
    )
    return PreparedData(
        name="2006",
        coords={
            "player_season": measurements["player_season"].tolist(),
            "season": measurements["season"].tolist(),
        },
        measurements=measurements,
    )

def prepare_data_bdb(
    measurements_main: pd.DataFrame,
    measurements_post: pd.DataFrame,
    appearances: pd.DataFrame,
) -> PreparedData:
    """Prepare the baseball-databank data.

```

There are a few substantive data choices here.

First, the function excludes players who have a '1' in their position as these are likely pitchers, as well as players with fewer than 20 at bats.

Second, the function defines a successes and attempts according to the 'on-base percentage' metric, so a success is a time when a player got a hit, a base on ball/walk or a hit-by-pitch and an attempt is an at-bat or a base-on-ball/walk or a hit-by-pitch or a sacrifice fly. This could have alternatively been calculated as just hits divided by at-bats, but my

understanding is that this method underrates players who are good at getting walks.

```
"""
pitchers = appearances.loc[
    lambda df: df["G_p"] == df["G_all"], "playerID"
].unique()

def filter_batters(df: pd.DataFrame):
    return (
        (df["AB"] >= 20)
        & (df["season"].ge(2017))
        & (~df["player"].isin(pitchers))
    )

measurements_main, measurements_post = (
    m.rename(columns={"yearID": "season", "playerID": "player"})
    .assign(
        player_season=lambda df: df["player"].str.cat(
            df["season"].astype(str)
        ),
        n_attempt=lambda df: df[["AB", "BB", "HBP", "SF"]]
        .fillna(0)
        .sum(axis=1)
        .astype(int),
        n_success=lambda df: (
            df[["H", "BB", "HBP"]].fillna(0).sum(axis=1).astype(int)
        ),
    )
    .loc[
        filter_batters,
        ["player_season", "season", "n_attempt", "n_success"],
    ]
    .copy()
    for m in [measurements_main, measurements_post]
)

measurements = (
    pd.concat([measurements_main, measurements_post])
    .groupby(["player_season", "season"])
    .sum()
)
```



```

        .reset_index()
    )
    return PreparedData(
        name="bdb",
        coords={
            "player_season": measurements["player_season"].tolist(),
            "season": measurements["season"].tolist(),
        },
        measurements=measurements,
    )

```

Finally we need to update `prepare_data.py`, the script that runs the data preparation functions. Again there isn't much to change from the example analysis.

```

"""Read the data in RAW_DIR and save prepared data to PREPARED_DIR."""

import os

import pandas as pd
from src import data_preparation_functions
from src.prepared_data import write_prepared_data

DATA_PREPARATION_FUNCTIONS_TO_RUN = {
    "2006": data_preparation_functions.prepare_data_2006,
    "bdb": data_preparation_functions.prepare_data_bdb,
}

DATA_DIR = os.path.join(os.path.dirname(__file__), "data")
RAW_DIR = os.path.join(DATA_DIR, "raw")
PREPARED_DIR = os.path.join(DATA_DIR, "prepared")
RAW_DATA_FILES = {
    "2006": [os.path.join(RAW_DIR, "baseball-hits-2006.csv")],
    "bdb": [
        os.path.join(RAW_DIR, "bdb-main.csv"),
        os.path.join(RAW_DIR, "bdb-post.csv"),
        os.path.join(RAW_DIR, "bdb-apps.csv"),
    ]
}

```

```

def main():
    """Save prepared data in the PREPARED_DIR."""
    print("Reading raw data...")
    raw_data = {
        k: [pd.read_csv(f, index_col=None) for f in v]
        for k, v in RAW_DATA_FILES.items()
    }
    print("Preparing data...")
    for name, dpf in DATA_PREPARATION_FUNCTIONS_TO_RUN.items():
        print(f"Running data preparation function {dpf.__name__}...")
        prepared_data = dpf(*raw_data[name])
        output_dir = os.path.join(PREPARED_DIR, prepared_data.name)
        print(f"\twriting files to {output_dir}")
        if not os.path.exists(PREPARED_DIR):
            os.mkdir(PREPARED_DIR)
        write_prepared_data(prepared_data, output_dir)

if __name__ == "__main__":
    main()

```

To check that all this works, we can run the script `prepare_data.py` manually or using `make analysis`. Now if we look in the file `data/prepared/bdb/measurements.csv` we should see some lines that look like this:

```

,player_season,season,n_attempt,n_success
0,abreujo02,2018,553,180
1,acunaro01,2018,487,178
3,adamewi01,2018,322,112
6,adamsla01,2018,29,10
7,adamsma01,2018,337,104
8,adamsma01,2018,277,92
9,adamsma01,2018,60,12

```

## Specifying statistical models

I wanted to test two statistical models: one with the modelling the distribution of per-player logit-scale at-bat success rates as a normal distribution with unknown mean and standard deviation, and another where the same logit-scale rates have a generalised pareto distribution.

So, given a table of  $N$  player profiles, with each player has  $y$  successes out of  $K$  at-bats and an unknown latent success rate  $\alpha$ , I wanted to use this measurement model:

$$y \sim \text{binomial logit}(K, \alpha)$$

In the generalised pareto model I would give the  $\alpha$ s this prior model, with the hyperparameter  $\min \alpha$  assumed to be known exactly and  $k$  and  $\sigma$  given prior distributions that put the  $\alpha$ s in the generally plausible range of between roughly 0.1 and 0.4.

$$\alpha \sim GPareto(\min \alpha, k, \sigma)$$

In the normal model I would use a standard hierarchical regression model with an effect for the log-scale number of at-bats to attempt to explicitly model the tendency of players with more appearances to be better:

$$\alpha \sim Normal(\mu + b_K \cdot \ln K, \tau)$$

Again I would choose priors for the hyperparameters that put most of the alphas between 0.1 and 0.4.

To implement these models using Stan I first added the following function to the file `custom_functions.stan`. This was simply copied from [the relevant part of the geomagnetic storms case study](#).

```

real gpareto_lpdf(vector y, real ymin, real k, real sigma) {
  // generalised Pareto log pdf
  int N = rows(y);
  real inv_k = inv(k);
  if (k<0 && max(y-ymin)/sigma > -inv_k)
    reject("k<0 and max(y-ymin)/sigma > -1/k; found k, sigma =", k, ", ", sigma);
  if (sigma<=0)
    reject("sigma<=0; found sigma =", sigma);
  if (fabs(k) > 1e-15)
    return -(1+inv_k)*sum(log1p((y-ymin) * (k/sigma))) -N*log(sigma);
  else
    return -sum(y-ymin)/sigma -N*log(sigma); // limit k->0
}

```

Next I wrote a file `gpareto.stan`:

```

functions {
#include custom_functions.stan
}
data {
  int<lower=0> N; // items
  array[N] int<lower=0> K; // trials
  array[N] int<lower=0> y; // successes
  real min_alpha; // noone worse than this would be in the dataset
  real max_alpha;
  array[2] real prior_sigma;
  array[2] real prior_k;
  int<lower=0,upper=1> likelihood;
}
parameters {
  real<lower=0.001> sigma; // scale parameter of generalised pareto distribution
  real<lower=-sigma/(max_alpha-min_alpha)> k; // shape parameter of generalised pareto distr
  vector<lower=min_alpha,upper=max_alpha>[N] alpha; // success log-odds
}
model {
  sigma ~ normal(prior_sigma[1], prior_sigma[2]);
  k ~ normal(prior_k[1], prior_k[2]);
  alpha ~ gpareto(min_alpha, k, sigma);
  if (likelihood){

```

```

    y ~ binomial_logit(K, alpha);
  }
}
generated quantities {
  vector[N] yrep;
  vector[N] llik;
  for (n in 1:N){
    yrep[n] = binomial_rng(K[n], inv_logit(alpha[n]));
    llik[n] = binomial_logit_lpmf(y[n] | K[n], alpha[n]);
  }
}

```

Finally I wrote a file `normal.stan`:

```

data {
  int<lower=0> N; // items
  array[N] int<lower=0> K; // trials
  array[N] int<lower=0> y; // successes
  array[2] real prior_mu;
  array[2] real prior_tau;
  array[2] real prior_b_K;
  int<lower=0,upper=1> likelihood;
}
transformed data {
  vector[N] log_K = log(to_vector(K));
  vector[N] log_K_std = (log_K - mean(log_K)) / sd(log_K);
}
parameters {
  real mu; // population mean of success log-odds
  real<lower=0> tau; // population sd of success log-odds
  real b_K;
  vector[N] alpha_std; // success log-odds (standardized)
}
model {
  b_K ~ normal(prior_b_K[1], prior_b_K[2]);
  mu ~ normal(prior_mu[1], prior_mu[2]);
  tau ~ normal(prior_tau[1], prior_tau[2]);
  alpha_std ~ normal(0, 1);
  if (likelihood){

```

```

    y ~ binomial_logit(K, mu + b_K * log_K_std + tau * alpha_std);
  }
}
generated quantities {
  vector[N] alpha = mu + b_K * log_K_std + tau * alpha_std;
  vector[N] yrep;
  vector[N] llik;
  for (n in 1:N){
    yrep[n] = binomial_rng(K[n], inv_logit(alpha[n]));
    llik[n] = binomial_logit_lpmf(y[n] | K[n], alpha[n]);
  }
}

```

## Generating Stan inputs

Next we need to tell our analysis how to turn some prepared data into a dictionary that can be used as input for Stan. Bibat assumes that this task is handled by functions that live in the file `src/stan_input_functions.py`, each of which takes in a `PreparedData` and returns a Python dictionary. You can write as many Stan input functions as you like and choose which one to run for any given inference.

We can start by defining some Stan input functions that pass arbitrary prepared data on to each of the models:

```

"""General function for creating a Stan input."""
return {
    "N": len(ppd.measurements),
    "K": ppd.measurements["n_attempt"].values,
    "y": ppd.measurements["n_success"].values,
    "prior_mu": [logit(0.25), 0.2],
    "prior_tau": [0.2, 0.1],
    "prior_b_K": [0, 0.03],
}

def get_stan_input_gpareto(ppd: PreparedData) -> Dict:

```

```

"""General function for creating a Stan input."""
return {
    "N": len(ppd.measurements),
    "K": ppd.measurements["n_attempt"].values,
    "y": ppd.measurements["n_success"].values,
    "min_alpha": logit(0.07),
    "max_alpha": logit(0.5),
    "prior_sigma": [1.5, 0.4],
    "prior_k": [-0.5, 1],
}

```

But why stop there? It can also be useful to generate Stan inputs consistently with a model, based on some hardcoded hyperparameter values. Here are some functions that do this for both of our models:

```

"""Generate fake Stan input consistent with the normal model."""
N = len(ppd.measurements)
rng = np.random.default_rng(seed=1234)
true_param_values = {
    "mu": logit(0.25),
    "tau": 0.18, # 2sds is 0.19 to 0.32 batting average
    "b_K": 0.04, # slight effect of more attempts
    "alpha_std": rng.random.normal(loc=0, scale=1, size=N),
}
K = ppd.measurements["n_attempt"].values
log_K_std = (np.log(K) - np.log(K).mean()) / np.log(K).std()
alpha = (
    true_param_values["mu"]
    + true_param_values["b_K"] * log_K_std
    + true_param_values["tau"] * true_param_values["alpha_std"]
)
y = rng.random.binomial(K, expit(alpha))
return {"N": N, "K": K, "y": y} | true_param_values

def get_stan_input_gpareto_fake(ppd: PreparedData) -> Dict:
    """Generate fake Stan input consistent with the gpareto model."""
    N = len(ppd.measurements)

```

```

K = ppd.measurements["n_attempt"].values
min_alpha = 0.1
rng = np.random.default_rng(seed=1234)
true_param_values = {"sigma": -1.098, "k": 0.18}
true_param_values["alpha"] = gpareto_rvs(
    rng,
    N,
    min_alpha,
    true_param_values["k"],
    true_param_values["sigma"],
)
y = rng.binomial(K, expit(true_param_values["alpha"]))
return {"N": N, "K": K, "y": y, "min_alpha": min_alpha} | true_param_values

def gpareto_rvs(
    rng: np.random.Generator, size: int, mu: float, k: float, sigma: float
):
    """Generate random numbers from a generalised pareto distribution.

    See https://en.wikipedia.org/wiki/Generalized\_Pareto\_distribution for
    source.

    """
    U = rng.uniform(size)
    if k == 0:
        return mu - sigma * np.log(U)
    else:
        return mu + (sigma * (U**-k) - 1) / sigma

```

## Specifying inferences

Now all the building blocks for making statistical inferences - raw data, data preparation rules, statistical models and recipes for turning prepared data into model inputs - are in place. The last step before actually running Stan is to write down how put these blocks together. Bibat has another concept for this, called ‘inferences’.



An inference in bibat is a folder containing a special file called `config.toml`. This file sets out what inferences you want to make: which statistical model, which prepared data function, which Stan input function, which parameters have which dimensions, which sampling modes to use and how to configure the sampler. The folder will later be filled up with the results of performing the specified inferences.

I started by deleting the example inferences and creating two fresh folders, leaving me with an `inferences` folder looking like this:

```
.
├── gpareto2006
│   └── config.toml
└── normal2006
    └── config.toml
```

Here is the file `inferences/gpareto2006/config.toml`:

```
name = "gpareto2006"
stan_file = "gpareto.stan"
prepared_data_dir = "2006"
stan_input_function = "get_stan_input_gpareto"
modes = ["prior", "posterior"]

[dims]
alpha = ["player"]

[stanc_options]
warn-pedantic = true

[sample_kwargs]
save_warmup = false
iter_warmup = 2000
iter_sampling = 2000

[sample_kwargs.prior]
chains = 2
iter_warmup = 1000
iter_sampling = 1000
```

Here is the file `inferences/normal2006/config.toml`:

```
name = "normal2006"
stan_file = "normal.stan"
prepared_data_dir = "2006"
stan_input_function = "get_stan_input_normal"
modes = ["prior", "posterior"]

[dims]
alpha = ["player"]

[stanc_options]
warn-pedantic = true

[sample_kwargs]
save_warmup = false
iter_warmup = 2000
iter_sampling = 2000
```

Note that: \* The Stan file, prepared data folder and stan input function are referred to by strings. The analysis should raise an error if you enter a non-existing value. \* Both inferences are set to run in “prior” and “posterior” modes - the other pre-defined mode is “kfold”, but you can also write your own! \* You can enter arbitrary arguments to `cmdstanpy's CmdStanModel.sample` method in the `[sample_kwargs]` table. \* You can enter mode-specific overrides in `[sample_kwargs.<MODE>]`. This can be handy if you want to run more or fewer iterations for a certain mode.

Now when I ran `make analysis` again, I saw messages indicating that Stan had run, and found that the `inferences` subfolders had been populated:

```
inferences
  gpareto2006
    config.toml
    idata.json
  normal2006
    config.toml
```

```
idata.json
```

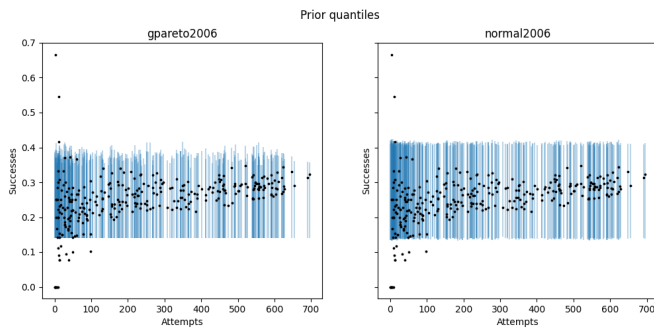
## Investigating the inferences

Now that the inferences are ready it's time to check them out. Bibat provides a jupyter notebook called `investigate.ipynb` for exactly this purpose.

A lot of code from the example analysis's notebook was reusable, so I largely followed its structure, with a few tweaks.

## Choosing priors using push-forward calibration

The trickiest thing about my analysis was setting prior distributions for the parameters  $k$  and  $\sigma$  in the generalised Pareto models. To choose some more or less plausible values I did a few prior-mode model runs and checked the distributions of the resulting `alpha` variables. I wanted to make sure that they all lay in the range corresponding to batting averages between about 0.15 and a little over 0.4. Here is a graph that shows the 1% to 99% prior quantiles for each player's latent success percentage in the 2006 dataset alongside their actually realised success rate.



## Extending the analysis to the baseball databank data

To model the more recent data, all I had to do was create some new inference folders with appropriate `prepared_data_dir` fields in their `config.toml` files. For example, here is the `config.toml` file for the `gparetobdb` inference:

```
name = "gparetobdb"
stan_file = "gpareto.stan"
prepared_data_dir = "bdb"
stan_input_function = "get_stan_input_gpareto"
modes = ["prior", "posterior"]

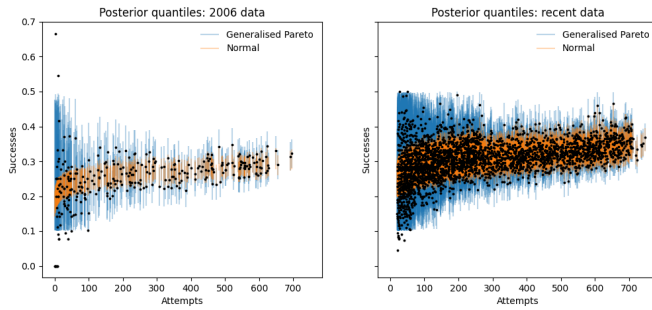
[dims]
alpha = ["player"]

[stanc_options]
warn-pedantic = true

[sample_kwargs]
save_warmup = false
iter_warmup = 2000
iter_sampling = 2000

[sample_kwargs.prior]
chains = 2
iter_warmup = 2000
iter_sampling = 1000
```

After running `make analysis` one more time, I went back to the notebook `investigate.ipynb` and made plots of both models' posterior 1%-99% success rate intervals for both datasets:



I think this is very interesting. Both models' prior distributions had similar regularisation levels, and they more or less agree about the abilities of the players with the most at-bats, both in terms of locations and the widths of the plausible intervals for the true success rate. However, the models ended up with dramatically different certainty levels about the abilities of players with few at-bats. This pattern was true both for the small 2006 dataset and the much larger baseball-databank dataset.