## VARIANCE EXPLAINED

<div>&#9776; menu</div>

# Understanding empirical Bayes estimation (using baseball statistics)

Which of these two proportions is higher: **4 out of 10**, or **300 out of 1000**? This sounds like a silly question. Obviously $4/10 = .4$, which is greater than $300/1000 = .3$.

But suppose you were a baseball recruiter, trying to decide which of two potential players is a better batter based on how many hits they get. One has achieved 4 hits in 10 chances, the other 300 hits in 1000 chances. While the first player has a higher proportion of hits, it's not a lot of evidence: a typical player tends to achieve a hit around 27% of the time, and this player's 4/10 could be due to luck. The second player, on the other hand, has a lot of evidence that he's an above-average batter.

This post isn't really about baseball, I'm just using it as an illustrative example. (I actually know very little about sabermetrics. If you want a more technical version of this post, check out this great paper). This post is, rather, about a very useful statistical method for estimating a large number of proportions, called **empirical Bayes estimation**. It's to help you with data that looks like this:

| Success | Total |
|---:|---:|
| 11 | 104 |
| 82 | 1351 |
| 2 | 26 |
| 0 | 40 |
| 1203 | 7592 |
| 5 | 166 |

A lot of data takes the form of these success/total counts, where you want to estimate a "proportion of success" for each instance. Each row might represent:

- **An ad you're running**: Which of your ads have higher clickthrough rates, and which have lower? (Note that I'm not talking about running an A/B test comparing two options, but rather about ranking and analyzing a large list of choices.)
- **A user on your site**: In my work at Stack Overflow, I might look at what fraction of a user's visits are to Javascript questions, to guess whether they are a web developer. In another application, you might consider how often a user decides to read an article they browse over, or to purchase a product they've clicked on.

When you work with pairs of successes/totals like this, you tend to get tripped up by the uncertainty in low counts. 1/2 does not mean the same thing as 50/100; nor does 0/1 mean the same thing as 0/1000. One approach is to filter out all cases that don't meet some minimum, but this isn't always an option: you're throwing away useful information.

I previously wrote a post about one approach to this problem, using the same analogy: Understanding the beta distribution (using baseball statistics). Using the beta distribution to represent your *prior expectations*, and *updating* based on the new evidence, can help make your estimate more accurate and practical. Now I'll demonstrate the related method of empirical Bayes estimation, where the beta distribution is used to improve a large set of estimates. What's great about this method is that as long as you have a lot of examples, *you don't need to bring in prior expectations*.

Here I'll apply empirical Bayes estimation to a baseball dataset, with the goal of improving our estimate of each player's batting average. I'll focus on the intuition of this approach, but will also show the R code for running this analysis yourself. (So that the post doesn't get cluttered, I don't show the code for the graphs and tables, only the estimation itself. But you can find *all* this post's code here).

## Working with batting averages

In my original post about the beta distribution, I made some vague guesses about the distribution of batting averages across history, but here we'll work with real data. We'll use the `Batting` dataset from the excellent Lahman package. We'll prepare and clean the data a little first, using dplyr and tidyr:

```r
library(dplyr)
library(tidyr)
library(Lahman)

career <- Batting %>%
  filter(AB > 0) %>%
  anti_join(Pitching, by = "playerID") %>%
  group_by(playerID) %>%
  summarize(H = sum(H), AB = sum(AB)) %>%
  mutate(average = H / AB)

# use names along with the player IDs
career <- Master %>%
  tbl_df() %>%
  select(playerID, nameFirst, nameLast) %>%
  unite(name, nameFirst, nameLast, sep = " ") %>%
  inner_join(career, by = "playerID") %>%
  select(-playerID)
```

Above, we filtered out pitchers (generally the weakest batters, who should be analyzed separately). We summarized each player across multiple years to get their *career* Hits

(H) and At Bats (AB), and batting average. Finally, we added first and last names to the dataset, so we could work with them rather than an identifier:

```
career
```

```
## Source: local data frame [9,256 x 4]
##
##                  name    H    AB average
##                 (chr) (int) (int)  (dbl)
## 1         Hank Aaron  3771 12364  0.3050
## 2       Tommie Aaron   216   944  0.2288
## 3          Andy Abad     2    21  0.0952
## 4        John Abadie    11    49  0.2245
## 5     Ed Abbaticchio   772  3044  0.2536
## 6        Fred Abbott   107   513  0.2086
## 7        Jeff Abbott   157   596  0.2634
## 8        Kurt Abbott   523  2044  0.2559
## 9         Ody Abbott    13    70  0.1857
## 10 Frank Abercrombie     0     4  0.0000
## ..                  ...   ...   ...    ...
```

I wonder who the best batters in history were. Well, here are the ones with the highest batting average:

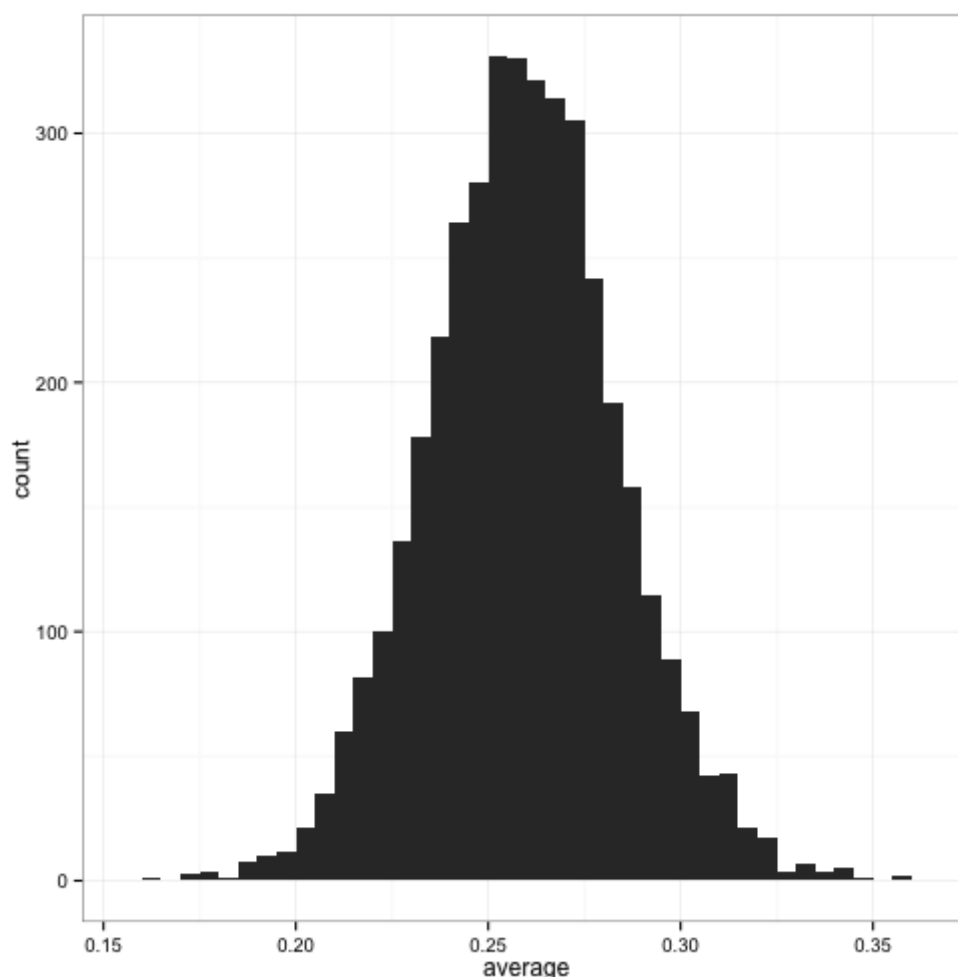| name | H | AB | average |
|---|---|---|---|
| Jeff Banister | 1 | 1 | 1 |
| Doc Bass | 1 | 1 | 1 |
| Steve Biras | 2 | 2 | 1 |
| C. B. Burns | 1 | 1 | 1 |
| Jackie Gallagher | 1 | 1 | 1 |

Err, that's not really what I was looking for. These aren't the best batters, they're just the batters who went up once or twice and got lucky. How about the worst batters?

| name | H | AB | average |
|---|---|---|---|
| Frank Abercrombie | 0 | 4 | 0 |
| Horace Allen | 0 | 7 | 0 |
| Pete Allen | 0 | 4 | 0 |
| Walter Alston | 0 | 1 | 0 |
| Bill Andrus | 0 | 9 | 0 |

Also not what I was looking for. That "average" is a really crummy estimate. **Let's make a better one.**

## Step 1: Estimate a prior from all your data

Let's look at the distribution of batting averages across players.



(For the sake of estimating the prior distribution, I've filtered out all players that have fewer than 500 at-bats, since we'll get a better estimate from the less noisy cases. I show a more principled approach in the Appendix).

The first step of empirical Bayes estimation is to estimate a beta prior using this data. Estimating priors from the data you're currently analyzing is not the typical Bayesian approach- usually you decide on your priors ahead of time. There's a lot of debate and discussion about when and where it's appropriate to use empirical Bayesian methods, but it basically comes down to how many observations we have: if we have a lot, we can get a good estimate that doesn't depend much on any one individual. Empirical Bayes is an **approximation** to more exact Bayesian methods- and with the amount of data we have, it's a very good approximation.

So far, a beta distribution looks like a pretty appropriate choice based on the above histogram. (What would make it a bad choice? Well, suppose the histogram had two peaks, or three, instead of one. Then we might need a mixture of Betas, or an even more complicated model). So we know we want to fit the following model:

$$X \sim \text{Beta}(\alpha_0, \beta_0)$$

We just need to pick $\alpha_0$ and $\beta_0$, which we call "hyper-parameters" of our model. There are many methods in R for fitting a probability distribution to data ( `optim` , `mle` , `bbmle` , etc). You don't even have to use maximum likelihood: you could use the mean
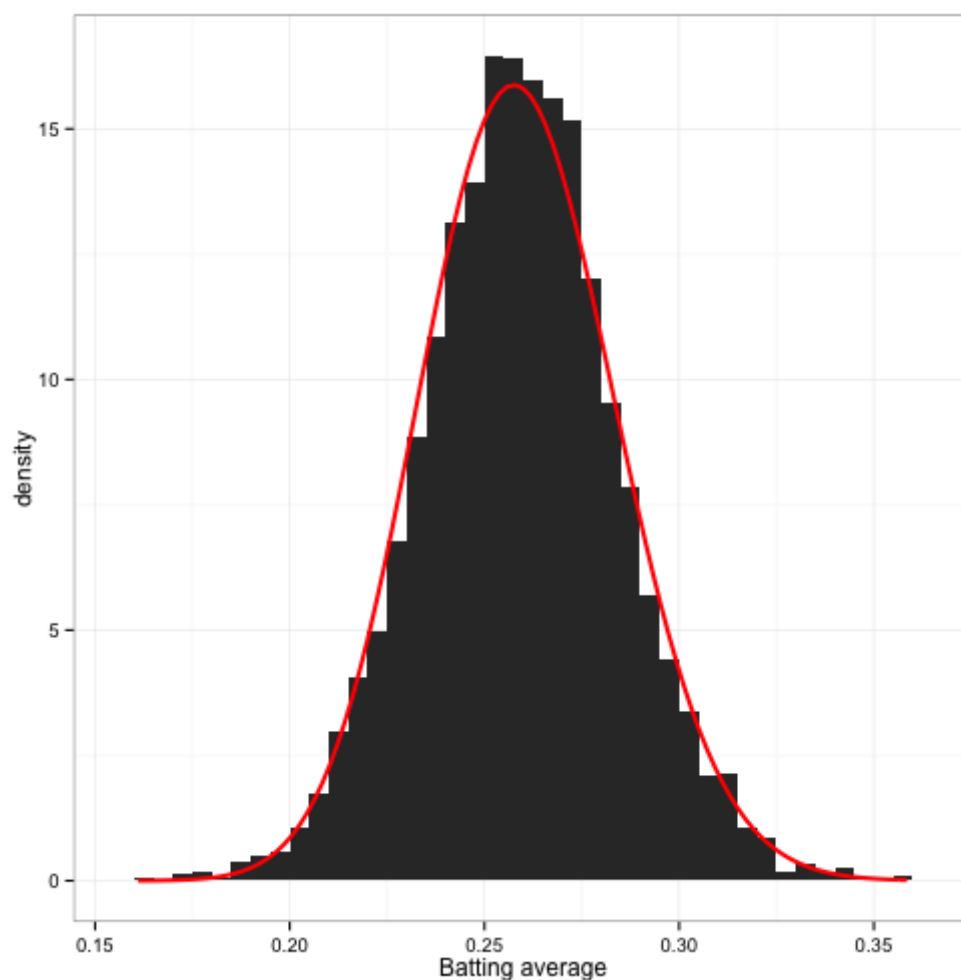
and variance, called the "method of moments". But we'll use the fitdistr function from MASS.

```
# just like the graph, we have to filter for the players we actually
# have a decent estimate of
career_filtered <- career %>%
    filter(AB >= 500)

m <- MASS::fitdistr(career_filtered$average, dbeta,
                    start = list(shape1 = 1, shape2 = 10))

alpha0 <- m$estimate[1]
beta0 <- m$estimate[2]
```

This comes up with $\alpha_0 = 78.661$ and $\beta_0 = 224.875$. How well does this fit the data?



Not bad! Not perfect, but something we can work with.

## Step 2: Use that distribution as a prior for each individual estimate

Now when we look at any individual to estimate their batting average, we'll start with our overall prior, and update based on the individual evidence. I went over this process

in detail in the <u>original Beta distribution post</u>: it's as simple as adding $\alpha_0$ to the number of hits, and $\alpha_0 + \beta_0$ to the total number of at-bats.

For example, consider our hypothetical batter from the introduction that went up 1000 times, and got 300 hits. We would estimate his batting average as:

$$\frac{300 + \alpha_0}{1000 + \alpha_0 + \beta_0} = \frac{300 + 78.7}{1000 + 78.7 + 224.9} = 0.29$$

How about the batter who went up only 10 times, and got 4 hits. We would estimate his batting average as:

$$\frac{4 + \alpha_0}{10 + \alpha_0 + \beta_0} = \frac{4 + 78.7}{10 + 78.7 + 224.9} = 0.264$$

Thus, even though $\frac{4}{10} > \frac{300}{1000}$, we would guess that the $\frac{300}{1000}$ batter is better than the $\frac{4}{10}$ batter!

Performing this calculation for all the batters is simple enough:

```
career_eb <- career %>%
    mutate(eb_estimate = (H + alpha0) / (AB + alpha0 + beta0))
```

## Results

Now we can ask: who are the best batters by this improved estimate?

| name | H | AB | average | eb_estimate |
|---|---|---|---|---|
| Rogers Hornsby | 2930 | 8173 | 0.358 | 0.355 |
| Shoeless Joe Jackson | 1772 | 4981 | 0.356 | 0.350 |
| Ed Delahanty | 2596 | 7505 | 0.346 | 0.343 |
| Billy Hamilton | 2158 | 6268 | 0.344 | 0.340 |
| Harry Heilmann | 2660 | 7787 | 0.342 | 0.339 |

Who are the *worst* batters?

| name | H | AB | average | eb_estimate |
|---|---|---|---|---|
| Bill Bergen | 516 | 3028 | 0.170 | 0.178 |
| Ray Oyler | 221 | 1265 | 0.175 | 0.191 |
| John Vukovich | 90 | 559 | 0.161 | 0.196 |
| John Humphries | 52 | 364 | 0.143 | 0.196 |
| George Baker | 74 | 474 | 0.156 | 0.196 |

Notice that in each of these cases, empirical Bayes didn't simply pick the players who had 1 or 2 at-bats. It found players who batted well, or poorly, across a long career.

What a load off our minds: we can start using these empirical Bayes estimates in downstream analyses and algorithms, and not worry that we're accidentally letting 0/1 or 1/1 cases ruin everything.

Overall, let's see how empirical Bayes changed all of the batting average estimates:



The horizontal dashed red line marks $y = \frac{\alpha_0}{\alpha_0 + \beta_0} = 0.259$ - that's what we would guess someone's batting average was if we had *no* evidence at all. Notice that points above that line tend to move down towards it, while points below it move up.

The diagonal red line marks $x = y$. Points that lie close to it are the ones that didn't get shrunk at all by empirical Bayes. Notice that they're the ones with the highest number of at-bats (the brightest blue): they have enough evidence that we're willing to believe the naive batting average estimate.

This is why this process is sometimes called *shrinkage*: we've moved all our estimates towards the average. How much it moves these estimates depends on how much evidence we have: if we have very little evidence (4 hits out of 10) we move it a lot, if we have a lot of evidence (300 hits out of 1000) we move it only a little. That's shrinkage in a nutshell: *Extraordinary outliers require extraordinary evidence*.

## Conclusion: So easy it feels like cheating

Recall that there were two steps in empirical Bayes estimation:

1. Estimate the overall distribution of your data.

2. Use that distribution as your prior for estimating each average.

Step 1 can be done once, "offline"- analyze all your data and come up with some estimates of your overall distribution. Step 2 is done for each new observation you're considering. You might be estimating the success of a post or an ad, or classifying the behavior of a user in terms of how often they make a particular choice.

And because we're using the beta and the binomial, consider how *easy* that second step is. All we did was add one number to the successes, and add another number to the total. You can build that into your production system with a single line of code that takes nanoseconds to run.

```
// We hired a Data Scientist to analyze our Big Data

// and all we got was this lousy line of code.

float estimate = (successes + 78.7) / (total + 303.5);
```
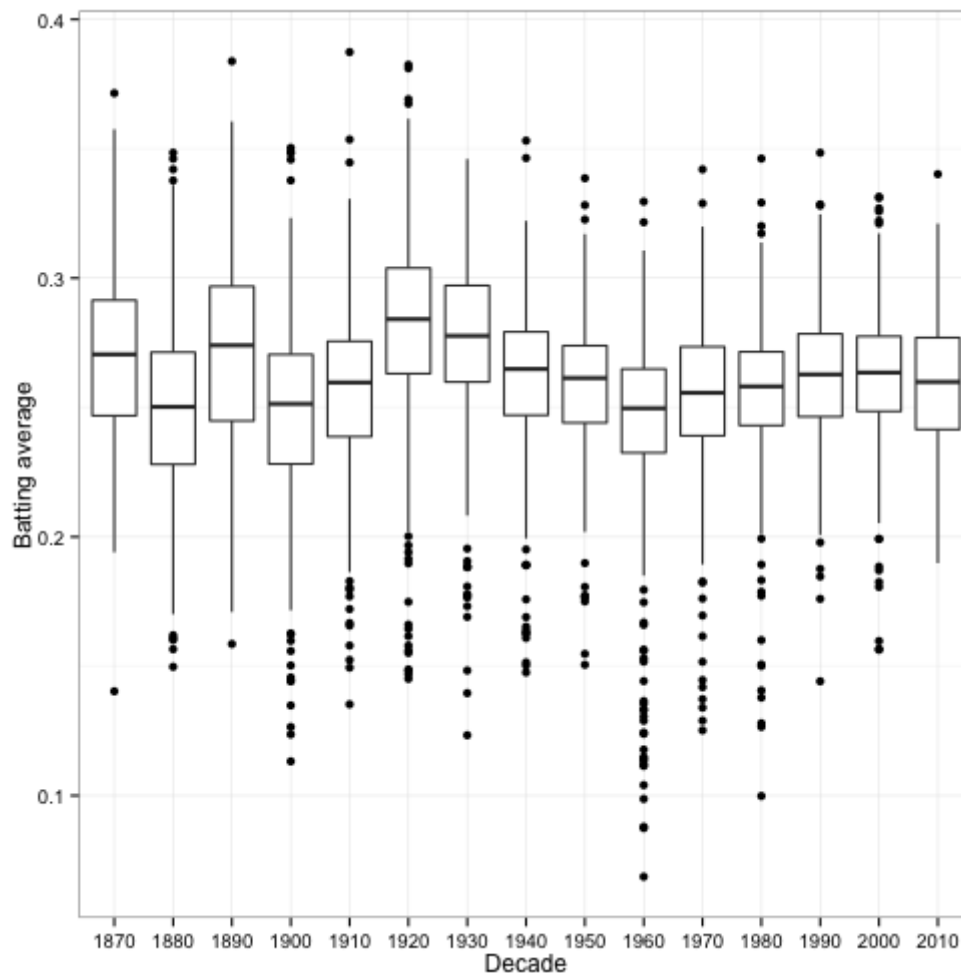
That really is so simple that it feels like cheating- like the kind of "fudge factor" you might throw into your code, with the intention of coming back to it later to do some real Machine Learning.

I bring this up to disprove the notion that statistical sophistication necessarily means dealing with complicated, burdensome algorithms. This Bayesian approach is based on sound principles, but it's still easy to implement. Conversely, next time you think "I only have time to implement a dumb hack," remember that you can use methods like these: it's a way to choose your fudge factor. Some dumb hacks are better than others!

But when anyone asks what you did, remember to call it "empirical Bayesian shrinkage towards a Beta prior." We statisticians have to keep up appearances.

## Appendix: How could we make this more complicated?

We've made some enormous simplifications in this post. For one thing, we assumed all batting averages are drawn from a single distribution. In reality, we'd expect that it depends on some known factors. For instance, the distribution of batting averages has changed over time:

Ideally, we'd want to estimate a different Beta prior for each decade. Similarly, we could estimate separate priors for each team, a separate prior for pitchers, and so on. One useful approach to this is Bayesian hierarchical modeling (as used in, for example, this study of SAT scores across different schools).

Also, as alluded to above, we shouldn't be estimating the distribution of batting averages using only the ones with more than 500 at-bats. Really, we should use all of our data to estimate the distribution, but give *more consideration* to the players with a higher number of at-bats. This can be done by fitting a beta-binomial distribution. For instance, we can use the dbetabinom.ab function from VGAM, and the `mle` function:

```
library(VGAM)

# negative log likelihood of data given alpha; beta
ll <- function(alpha, beta) {
  -sum(dbetabinom.ab(career$H, career$AB, alpha, beta, log = TRUE))
}

m <- mle(ll, start = list(alpha = 1, beta = 10), method = "L-BFGS-B")
coef(m)
```

```
## alpha  beta
##    75   222
```

We end up getting almost the same prior, which is reassuring!

---

## David Robinson

*Data Insights Engineering Manager at Flatiron Health, works in R and Python.*

✉ Email    🐦 Twitter    🜲 Github    ✍ Stack Overflow

### Subscribe

> Your email

Subscribe to this blog

### Recommended Blogs

- DataCamp
- R Bloggers
- RStudio Blog
- R4Stats
- Simply Statistics

**Understanding empirical Bayes estimation (using baseball statistics)** was published on October 01, 2015.

---

**Comments**    **Community**                              ⬤    Ava

♡ Recommend  20          🐦 Tweet        f Share        Sort by Best ▾

> Join the discussion…

**Freddy** · 4 years ago · **edited**
Hi David,

This is a very informative post. I'm wondering: can this be applied to a case where you're not estimating a fraction?

As in, if you wanted to estimate a score (s) given a sample size (n), but aren't looking to calculate s/n.

Sticking to your example it might be something like hits per game (s) over a course of (n) games. How would that influence the formula of (hits + a0)/at_bat+a0 + b0)

Thank you,

Freddy

2 ∧ | ∨ · Reply · Share ›

**David Robinson** Mod → Freddy · 4 years ago

This is a great question. If you are able to assume the number of hits per game is Poisson distributed for each player, you can use a gamma as the prior instead of the beta. I use that approach used in my post here:

https://github.com/dgrtwo/d...

(Look for the use of "dgamma"). I may make that into a post sometime!

1 ∧ | ∨ · Reply · Share ›

**Himanshu Verma** → David Robinson
· 3 years ago

Hi David, thanks for your awesome posts. I was wondering if you already wrote the post about the use of gamma distribution that you mentioned in the comment above.

∧ | ∨ · Reply · Share ›

**David Robinson** Mod →
Himanshu Verma · 3 years ago

Unfortunately I never got around to it, since I couldn't find a dataset within baseball that I was comfortable modeling with the Poisson! (Many reasonable guesses, like hits per game, were overdispersed). Instead, in the e-book version of this series of posts (see here: http://varianceexplained.or... ) I included a chapter on the Multinomial/Dirichlet.

∧ | ∨ · Reply · Share ›

**Himanshu Verma** → David
Robinson · 3 years ago

I got the book. Thanks a lot :)

∧ | ∨ · Reply · Share ›

**Александр Черкасов** · 3 years ago

Hi David!

Thanks a lot, I really enjoy to read your posts!

I would like to ask you for a hint. I have a similar problem, but my distribution is EXTREMELY skewed to the left. To make it clear, I have ~110.000 observations and among them ~80.000 averages are equal to 0 (mean is about 0.002).

Obviously, I can't filter all them out, for example fitting the log(average>0).

Function fitdistr don't work, I guess because it uses log, so gets Inf at process and can't optimize it.

But I thought that your solution with VGAM library must work, because you use not filtered data frame there, so you have zeros too. But I get the same exact error:

> m <- mle(ll, start = list(alpha = 1, beta = 10), method = "L-BFGS-B")
Error in optim(start, f, method = method, hessian = TRUE, ...) :
L-BFGS-B needs finite values of 'fn'
(and a couple of additional warnings:
In addition: Warning messages:
1: In lbeta(shape1[okk] + x[okk], shape2[okk] + size[okk] - x[okk]) :
NaNs produced
2: In lbeta(shape1[okk], shape2[okk]) : NaNs produced)

I'm a bit embarrassed, looks like I don't fully understand how it works. Do you think problem is in my starting alpha and beta?

Sorry for my poor English and thanks a lot in advance for your insights!

1 ∧ | ∨ · Reply · Share ›

**Matt F.** · 4 years ago

Thank you for the excellent post. I'm new to Bayesian analysis and you've made it very clear and applicable to a problem I routinely face. I have two questions for you:

1. Can you comment on the applicability of the Beta distribution when the empirical distribution I'm fitting it to is HIGHLY skewed? I'm performing an analysis much like the one presented here but on a very large data set where the vast majority of the probabilities are around .01.

2. What approach to calculating alpha and beta would you recommend on a large dataset (~6 million records). MLE seems to be quite slow (at least using Python's

scipy.stats.beta.fit).

Thanks!

1 ∧ | ∨ · Reply · Share ›

David Robinson  Mod  → Matt F.
· 4 years ago · edited

1. Betas can be very heavily skewed when beta is
much greater than alpha. I often use it to model
probabilities around 10^-4. The thing to watch out
for is that you need much higher n's to get an
estimate that is even reasonable.

2. I would recommend first filtering for cases with
higher n's. Presumably some of the n's are much
higher, and those will give you much more
information and less noise. Once you have it down
to a reasonable number of records try the MLE
again.

You can also use the method of moments as
described here: https://en.wikipedia.org/wi..., for
which you need only the mean and variance. But
this will also be very skewed- indeed, much more
so- if you do not first filter for the highest n's.

∧ | ∨ · Reply · Share ›

Michael Wiebe · 4 months ago

Your AB vs eb graph seems off: the y-axis has 0.2 twice,
and the x-axis spacing is weird.

∧ | ∨ · Reply · Share ›

Aman Jain · 5 months ago

Hello David.

Thanks for a really nice explanation of Bayes
approximation.
Question: How can we implement this method for Query
classification problem. Given a query and we know the
frequency of clicks on each category/ class. How to find the
probability of each category for each Query.

∧ | ∨ · Reply · Share ›

Jake · 8 months ago

Hi David,

I had a quick question, do you have any code for checking
goodness-of-fit for the beta-binomial distributed batting

averages?

Much appreciated,
Jake

ᴧ | ᴠ · Reply · Share ›

**Karl Gierach** · a year ago

Hi David,

Regarding the paper you cited; "Empirical Bayes in-season prediction of
baseball batting averages", by Jiang & Zhang, which specific techniques from the paper are implemented in this post? MLE and MM? It seems they are in general working with much more complicated models.

Thanks for the post!

Karl

ᴧ | ᴠ · Reply · Share ›

**David Robinson** Mod → Karl Gierach · a year ago

By "for a more technical version" above I really mean "for a more advanced/accurate version of using empirical Bayes to estimate batting averages."

This post uses the example of baseball to teach a simple application of empirical Bayes, but it doesn't handle the many complicated issues specific to baseball (field effects, changes over time, how more AB is confounded with better batters, etc). Jiang and Zhang is an example of a paper that takes those very seriously and is where one should go if you really want to estimate batting averages during a season.

Having said that, the theoretical underpinnings of this post and their empirical Bayes methods (estimate a prior, use that to calculate posteriors) are still very similar- though they bear more similarity to later posts in my series, like this one: http://varianceexplained.or...

10 ᴧ | ᴠ · Reply · Share ›

**Babas** · 2 years ago · edited

Hi David,

I am very happy to have found your blog, I believe your way to teach with simple examples is great. I still have a

question, in Step 1 "Estimate a prior from all your data"
When you fit the data using the beta distribution,
obviously you are using the beta distribution because it's
the topic of the post but you could have fit the data using a
normal distribution or a gamma distribution. My question
is when to use one kind of distribution rather than another
and, what are the advantages of using a distribution and
not another.

Best,

Babas

˄ | ˅ · **Reply** · **Share ›**

**Anicet Ebou** · 2 years ago

Hello David,

I have learned a lot, understood key principles in this post.
I can just say thanks you so much.

˄ | ˅ · **Reply** · **Share ›**

**Michael Chirico** · 2 years ago

Just want to emphasize how much I love this post --
successfully implemented on real data for the first time
just now!! I was just about to start excluding small-sample
entities when looking at %ages and said... nah, EB time!

A question -- why bother using the "distribution fitter" to
get a0 and b0? It's pretty straightforward to get them for
beta distribution:

(m = sample mean, s = sample variance)
denominator = m*(1-m)/s - 1
numerator = m*denominator

(note we don't even need to calculate b0 directly, though
the formula for that is easy enough -- (1-m)*denominator)

I peaked a little at MASS::fitdistr and it seems it has the
MLE shortcuts hardwired for normal, lognormal, a few
others, but not for beta, so this may be slow for large data
sets.

˄ | ˅ · **Reply** · **Share ›**

**David Robinson** Mod → Michael Chirico
· 2 years ago

Thanks! I really recommend you check out my e-
book; Introduction to Empirical Bayes, where I

expand this to include many other empirical bayes
methods and analyses. (I also introduce the ebbr
package for applying this to real data!).

http://varianceexplained.or...

You're referring to the method of moments. That's
pretty accurate for beta and almost certainly good
enough for this problem, but I typically go with the
MLE (honestly I can't remember why I did for this
post, and I don't know why I assume it's more
accurate). In later posts, and more importantly in
the book, I usually fit a beta-binomial, which will be
more accurate because it knows to discount low-n
observations. The ebbr package lets you choose
between that and the method of moments!

1 ⌃ | ⌄ · Reply · Share ›

**Michael Chirico** ↗ David Robinson
· 2 years ago

excellent... goes straight to my other Q about
implementation (unequal weighting of the
input %s to the prior).

purchased and downloaded, cheers! 🤓

⌃ | ⌄ · Reply · Share ›

**Michael Chirico** · 2 years ago · edited

why does .2 show up twice on the y axis for the true vs. EB
avg figure? that whole y axis looks all sorts of messed up...
EB values don't seem to correspond to the "top players"
table above. maybe rounding?

⌃ | ⌄ · Reply · Share ›

**Charles Neiswender** · 2 years ago

Hi - I liked this blog post so much I recreated the work in
python. I'm sure I'm not the only one who has done it, but
I thought I'd share. Thanks for combining two of my
favorite things - baseball and data science!
http://coldfashioned.net/20...

⌃ | ⌄ · Reply · Share ›

**Rich Pauloo** · 2 years ago · edited

Nice post David. I've been following along in R.

I have a quick question. When you solve for the alpha and
beta parameters for the beta distribution in part 1, what is
the code you use to plot the beta distribution?

∧ | ∨ · **Reply** · **Share ›**

**Avatar**  This comment was deleted.

**Rich Pauloo** ➤ Guest · 2 years ago
Following up on this, I actually asked this
question on SO, and got this response:
https://stackoverflow.com/q...

∧ | ∨ · **Reply** · **Share ›**

**Andrew Zmeul** · 2 years ago
Hi!
I have a problem with the first code part it gives me an
error just after the first %>% part
Error in is_null(nms) : object 'rlang_is_null' not found

∧ | ∨ · **Reply** · **Share ›**

**Charles Neiswender** · 2 years ago
Really great post. Apart from not knowing R, the content
was well-structured and straightforward to follow. Add to
that the use of baseball stats, and I thoroughly enjoyed
reading. I've been brainstorming ways to cook up a Data
Science project related to baseball, and this really helps.
This will be a great jumping off point for my project.

∧ | ∨ · **Reply** · **Share ›**

**Professorwiki** · 3 years ago
What if I am just estimating a mean? What formula akin to
(H + alpha0) / (AB + alpha0 + beta0)) would I use? To be
more specific, I'm working with a poisson distribution
trying to estimate tornado frequency counts over the years.
The number of tornado reports has naturally increased
over time due to better tech and storm chasers, so I want to
get an empirical Bayes estimate for the earlier years in the
data

∧ | ∨ · **Reply** · **Share ›**

**Derek Thomas** · 3 years ago
Hi David,

Would it be even better to look at the distribution +/- 5
years from the midpoint of when the batter was batting?

Thanks,
Derek

∧ | ∨ · **Reply** · **Share ›**

**prateek** · 3 years ago

This was really helpful! Thanks a lot for the lucid
explanation

⌃ | ⌄ · **Reply** · **Share ›**

**HistorySquared** · 4 years ago

Great stuff, simple examples trump derived formulas.
Question, are there any particular packages or processes
anyone could recommend to fit a power law to a bayesian
distribution using R. Second, and just a request, but if if
you know how to conduct and interpret Bayesian Model
Averaging, that'd be great also. Regardless, thanks for the
post.

⌃ | ⌄ · **Reply** · **Share ›**
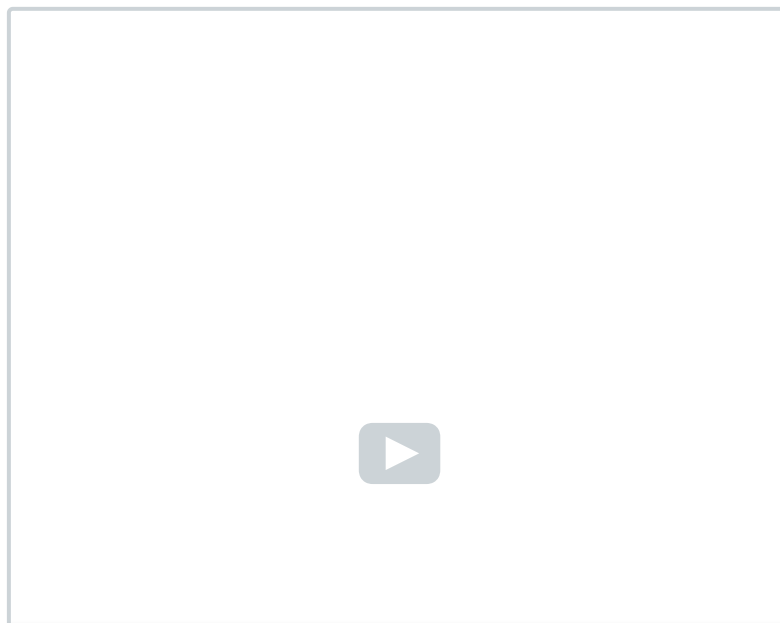
**Albert Y. Kim** · 4 years ago

Thanks for this.

Can you recommend a good reference detailing the
conditions under which it is appropriate to use an
empirical Bayes vs a pure Bayes approach? You hinted at it
in your paper in a slightly hand-wavy way (which is
appropriate for a high level post): "if you have lots of
observations."

⌃ | ⌄ · **Reply** · **Share ›**

**Aurélien** · 4 years ago

Thanks David for this post, I recently run into. I have a
question about correction for multiple comparisons. In this
video

▶

see more

see more

∧ | ∨ · Reply · Share ›

**David Robinson** Mod → Aurélien · 4 years ago

I address exactly the matter of multiple testing (specifically FDR control) on this same example in this post:

http://varianceexplained.or...

The reason it doesn't come up in the above post is that I'm not doing any hypothesis testing, we're just doing estimation.

1 ∧ | ∨ · Reply · Share ›

**Aurélien** → David Robinson · 4 years ago

ah! great. Thanks again. This means I was doing Bayesian statistics while using library("qvalue") without even noticing. Shame on me

∧ | ∨ · Reply · Share ›

**David Robinson** Mod → Aurélien · 4 years ago

Sort of! q-value has both frequentist and Bayesian interpretations. See https://projecteuclid.org/d... and http://noble.gs.washington....

∧ | ∨ · Reply · Share ›

**Aurélien** → David Robinson · 4 years ago

thanks! will need to dive into those papers. But, clearly my comprehension took off thank to your posts. Such a good teacher!

∧ | ∨ · Reply · Share ›

**S_Bushmanov** · 4 years ago · edited

David, thanks for the enlightening post! Really learnt something new about baysean statistics and learning from data.

Perhaps a dumb question, in the context of your post, more related to the base R.

After I follow your post and create `career` df, I try to explore it with `head(career[order(career$average, dec

=T),])`, but this statement fails due to shortened `decreasing` argument. I vaguely remember from R docs that I can shorten arguments in R functions and R will search for ones that match the shortened versions.

Do you have an idea why shortening `decreasing` does not work???

∧ | ∨ · **Reply** · **Share ›**

> **David Robinson** Mod ➔ S_Bushmanov
> · 4 years ago
>
> Yes- it's because there's an ellipsis (one of these: "...") before that argument. The ellipsis grabs any of the arguments that don't match exactly the keywords that follow. Take a look here for the formal explanation:
>
> https://cran.r-project.org/...
>
> If the ellipsis had been *after* the decreasing argument, it would have accepted partial matching. (Which can be a problem; see here: http://stackoverflow.com/qu.... But if that had been before the ellipsis, you wouldn't be able to break ties using order, e.g.:
>
> career[order(career$H, career$AB), ]
>
> Finally, note that a great way to do this sorting with dplyr is:
>
> career %>% arrange(desc(average))
>
> ∧ | ∨ · **Reply** · **Share ›**

**John Fries** · 4 years ago · edited
Hi David,

I noticed that you used the mean of the beta posterior distribution to sort your batters. I was wondering if you could comment on Cameron Davidson-Pilon's recommendation to use the 95% least plausible value? http://nbviewer.ipython.org...

Alternatively, I would propose that the comparison function for your sorting algorithm should be Prob(A > B), where A & B are batters.
So, A > B if Prob(A > B) > 0.5

Given that you have the posterior distribution for every batter, there's no reason you can't directly calculate this for

batter, there's no reason you can't directly calculate this for every pair of batters (or nlogn batters with memoization if computational resources are a concern). I've tried this technique for ranking random variables in a similar context, and it gives somewhat different answers than the other methods (mode, mean, median, least plausible value, etc). My only concern is that I have not been able to prove the transitivity of this comparison function (although I admittedly haven't been able to devote as much time to it as I would like).

Thank you,
John Fries

note: for my proposal, I ended up forcing a total ordering on the batters by sorting them lexically for the cases where there batting averages were the same, to avoid having to think through the equality case.

⌃ | ⌄  ·  Reply  ·  Share ›

**Julia Silge** · 4 years ago
This is such a clear, helpful post! I feel like this is the closest I've come to really grasping the math behind

**YOU MIGHT ALSO ENJOY**                                              (VIEW ALL POSTS)

- The 'largest stock profit or loss' puzzle: efficient computation in R
- The 'knight on an infinite chessboard' puzzle: efficient simulation in R
- Exploring college major and income: a live data analysis in R