

# **Bayesian Method and Hierarchical Modeling for Baseball**

**GR5224**

**Jiaqian Yu(jy2880)**

**Hanying Ji(hj2473)**

**3rd,May,2018**

## ● Introduction

New technology is changing the sports industry. Sports data analysis is a quickly developing field in United State, which combines the way of data analysis and the deep understanding of a specific sports to help league hiring good players in a quantitative way and make fans understanding their favorite team and players. As the origin of sports data analysis, Sabermetrics—the empirical analysis of baseball, evaluate past performance and predicting future performance to determine a player's contributions to his team. For a batter, the traditional measure of batting performance of a batter is considered to be the batting average - simply the number of times a player gets a base hit divided by the number of times he goes up at bat. We can evaluate a batter by calculating the Batting Average and rank them. But is that enough? Or are we confident to evaluate a batter simply apply this standard? Are there any other factors that can be consider to correctly get the Batting Average, or saying performance of a batter?

Being interested in this typical American sport, in this project, we are going to use the idea of Bayesian method which using more information as prior and the idea of hierarchical modeling to analyze the performance of baseball batters. The dataset comes from the amazing package Lahman in R, providing the tables from the 'Sean Lahman Baseball Database' as a set of R data frames.

## ● Description

The main purpose of this project is to see who are the Top Five players if we have their historical batting data. In order to simplify the problem, we make the assumption that a batter's batting performance doesn't change with time. A full bayesian treatment of the baseball batter hierarchical model will be used.

## ● Data

Our dataset contains information of 2232 baseball batters from year 2000 to 2016, which including playerID, name, bats(left/right/both hand), Hits(amount of Hits in total), AtBats(amount of AtBats in total), BattingAvg(batting accuracy).

- **Modeling**

To clarify, we use  $BattingAverage = Hits/AtBats$  to define the performance of a batter. Since *Hit* is a series of success and failure, which can be represented with a Binomial Random Variable, or  $Binomial(N, \theta)$ , where  $N$  is the observable total number of Hits in a season, and  $\theta$  is the true batting accuracy for every batter, the unknown parameter that we are interested in.

- **Model 1: frequentist parameter estimation**

In the general idea of frequentist, we simply use the sample mean of *BattingAverage* of each batter as MLE to estimate the parameter  $\theta$ . Let's look at the data:

playerID <chr>	name <chr>	bats <fctr>	Hits <int>	AtBats <int>	BattingAvg <dbl>
adamsla01	Lane Adams	R	0	3	0
bantzbr01	Brandon Bantz	R	0	2	0
barkese01	Sean Barker	R	0	2	0
baronst01	Steve Baron	R	0	11	0
barteki01	Kimera Bartee	B	0	19	0
barthji01	Jimmy Barthmaier	R	0	3	0

playerID <chr>	name <chr>	bats <fctr>	Hits <int>	AtBats <int>	BattingAvg <dbl>
davidda01	Dave Davidson	L	1	1	1.0000000
ohmeke01	Kevin Ohme	L	1	1	1.0000000
roachja01	Jason Roach	R	2	2	1.0000000
tupmama01	Matt Tupman	L	1	1	1.0000000
mantoje01	Jeff Manto	R	4	5	0.8000000
brettry01	Ryan Brett	B	2	3	0.6666667

As we can see, players with the lowest *BattingAverage* or highest *BattingAverage* are not either the worst player or the best player. The reason that let them get these *BattingAverage* is just because of lack of game experience(low *AtBats*). That "average" is a really crummy estimate. That's the defect of classical probability or Frequentist for this specific analysis. Let's make a better one.

- **Model 2: simple Bayesian method with different prior**

From the former result, we cannot just make conclusion that if we observe a new player goes up to bat once and gets a single, then his *BattingAverage* is briefly 100%, or if goes up to bat once and gets none, then he'll never get a hit all season. The best way to represent these prior expectations is with the beta distribution- it's saying, before we've seen the player taking his first

swing, what we roughly expect his *BattingAverage* to be. The reason why we choose Beta distribution is because it is the conjugate family of Binomial distribution and the domain of Beta distribution is (0,1), just like a probability, so we already know we're on the right track, but the appropriateness of the Beta for this task goes far beyond that. Since we have the historical data of each batter, we can use that data to help us make this estimate, in a Bayesian manner, by calculating the posterior for  $\theta$ .

For batters  $i$ ,

$$p(y_i|\theta_i) = \text{Bin}(N, \theta_i) \text{ where } i \in \{1, 2, \dots, 2232\}$$

We need to choose the prior. Here comes the idea of hierarchy, which imply a perspective from the whole to an individual.

In this part, we assume all batters have the same prior distribution  $\text{Beta}(\alpha_0, \beta_0)$ . Then the posterior for  $\theta$  is just  $\text{Beta}(\alpha_0 + \text{Hits}, \beta_0 + N - \text{Hits})$ . And our estimate for batter's true accuracy is  $\frac{\alpha_0 + \text{Hits}}{\alpha_0 + \beta_0 + N}$ .

(1) Noninformative prior:  $\text{Beta}(0,0)$

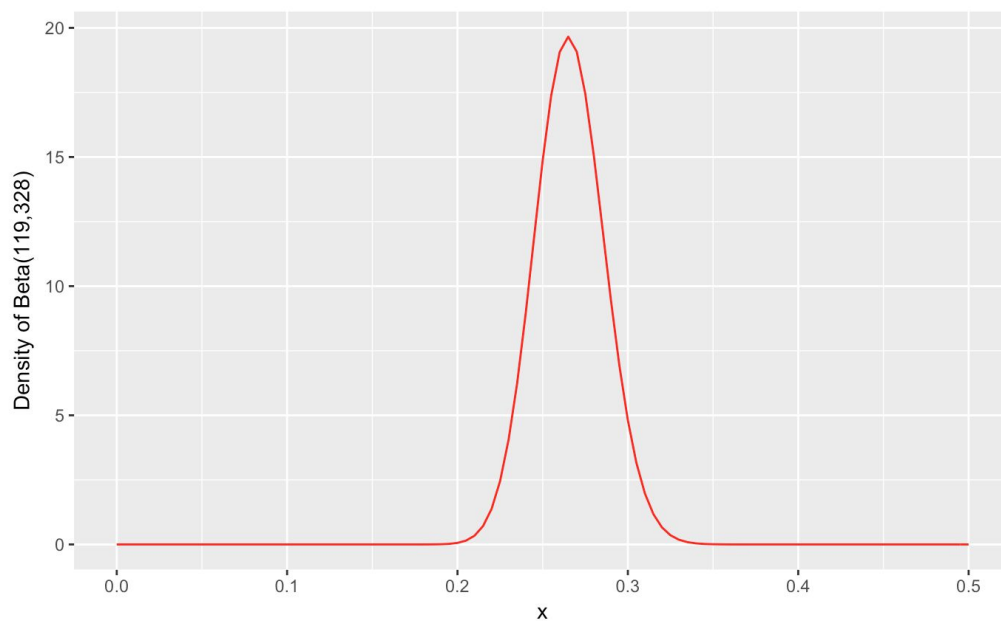
If we have no prior information about a parameter, we just use noninformative prior. The posterior mean are used to estimate everyone's true batting average, so the result is the same with simply calculating it by  $\text{BattingAvg} = \text{Hits}/\text{AtBats}$ , which is not desired.

(2) prior information:  $0.266$  is in general considered an average batting average, while  $0.300$  is considered an excellent one.

Since we get the prior information, letting  $\frac{\alpha_0}{\alpha_0 + \beta_0} = 0.266$  and the 95% quantile of beta distribution equals 0.3, we can get  $\alpha_0$  and  $\beta_0$  from the following list

[1]	0.6262650	0.6290661	0.6377271	0.6471558	0.6564142	0.6652789	0.6737082	0.6817149	0.6893279
[10]	0.6965793	0.7035000	0.7101178	0.7164580	0.7225426	0.7283913	0.7340215	0.7394487	0.7446866
[19]	0.7497474	0.7546424	0.7593813	0.7639733	0.7684266	0.7727485	0.7769459	0.7810252	0.7849919
[28]	0.7888516	0.7926090	0.7962687	0.7998349	0.8033117	0.8067026	0.8100112	0.8132405	0.8163938
[37]	0.8194737	0.8224830	0.8254242	0.8282997	0.8311117	0.8338624	0.8365538	0.8391879	0.8417663
[46]	0.8442910	0.8467635	0.8491854	0.8515582	0.8538834	0.8561623	0.8583963	0.8605866	0.8627344
[55]	0.8648409	0.8669071	0.8689342	0.8709232	0.8728751	0.8747908	0.8766712	0.8785172	0.8803297
[64]	0.8821094	0.8838572	0.8855739	0.8872601	0.8889167	0.8905441	0.8921433	0.8937147	0.8952590
[73]	0.8967768	0.8982688	0.8997354	0.9011772	0.9025948	0.9039887	0.9053593	0.9067072	0.9080329
[82]	0.9093367	0.9106193	0.9118809	0.9131220	0.9143431	0.9155444	0.9167266	0.9178898	0.9190345
[91]	0.9201610	0.9212697	0.9223610	0.9234351	0.9244924	0.9255332	0.9265579	0.9275666	0.9285598
[100]	0.9295377	0.9305006	0.9314487	0.9323824	0.9333019	0.9342075	0.9350993	0.9359777	0.9368429
[109]	0.9376951	0.9385345	0.9393614	0.9401760	0.9409786	0.9417692	0.9425482	0.9433156	0.9440718
[118]	0.9448170	0.9455512	0.9462746	0.9469876	0.9476902	0.9483826	0.9490650	0.9497375	0.9504004
[127]	0.9510537	0.9516976	0.9523324	0.9529580	0.9535747	0.9541827	0.9547820	0.9553728	0.9559553
[136]	0.9565295	0.9570957	0.9576538	0.9582042	0.9587468	0.9592818	0.9598093	0.9603295	0.9608424
[145]	0.9613482	0.9618470	0.9623388	0.9628239	0.9633023	0.9637740	0.9642393	0.9646981	0.9651507
[154]	0.9655971	0.9660373	0.9664715	0.9668998	0.9673223	0.9677390	0.9681501	0.9685555	0.9689555
[163]	0.9693501	0.9697393	0.9701232	0.9705020	0.9708757	0.9712444	0.9716081	0.9719669	0.9723210
[172]	0.9726703	0.9730149	0.9733549	0.9736904	0.9740214	0.9743480	0.9746703	0.9749883	0.9753021
[181]	0.9756117	0.9759173	0.9762188	0.9765163	0.9768099	0.9770996	0.9773856	0.9776677	0.9779462
[190]	0.9782210	0.9784922	0.9787599	0.9790240	0.9792847	0.9795420	0.9797960	0.9800466	0.9802940
[199]	0.9805382	0.9807792							

Then, we roughly choose  $\alpha_0=126$  and  $\beta_0=347.6842$  as hyperparameter and following is the prior distribution.

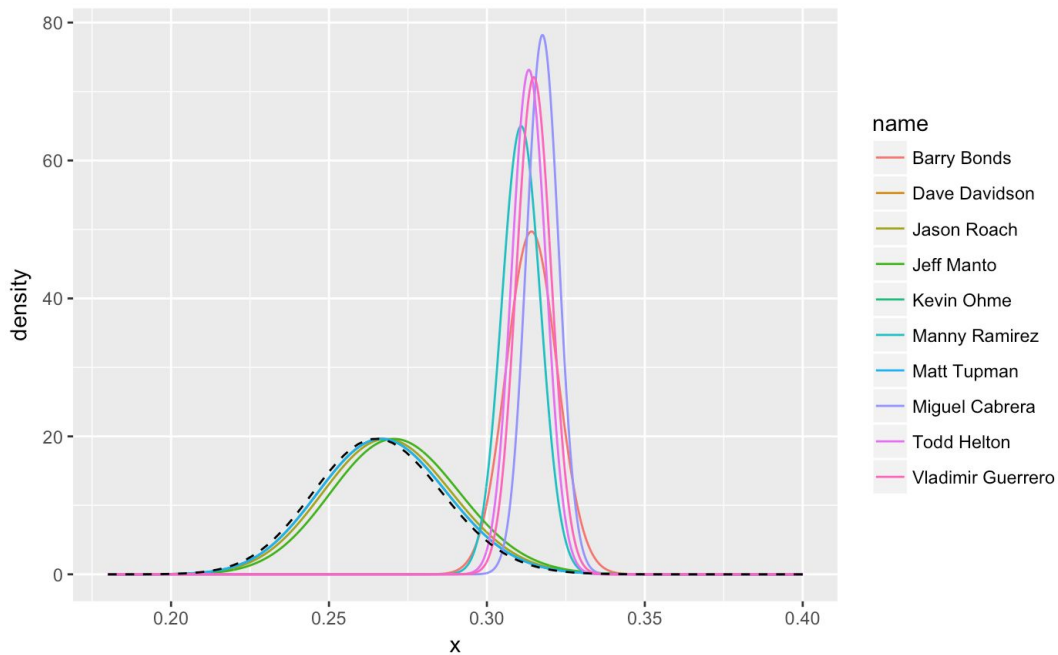


After using data to update the prior, we can get the posterior mean of each batters. Let's check the new Top 5 and Tail 5.

playerID <chr>	name <chr>	bats <fctr>	Hits <int>	AtBats <int>	BattingAvg <dbl>	alpha1 <dbl>	beta1 <dbl>	BattingAvg_post <dbl>
cabremi01	Miguel Cabrera	R	2519	7853	0.3207691	2645	5681.684	0.3176535
guerrvl01	Vladimir Guerrero	R	2092	6570	0.3184170	2218	4825.684	0.3148920
bondsba01	Barry Bonds	L	925	2871	0.3221874	1051	2293.684	0.3142300
heltoto01	Todd Helton	L	2141	6761	0.3166691	2267	4967.684	0.3133516
ramirma02	Manny Ramirez	R	1642	5213	0.3149818	1768	3918.684	0.3109017
suzukic01	Ichiro Suzuki	L	3030	9689	0.3127258	3156	7006.684	0.3105479

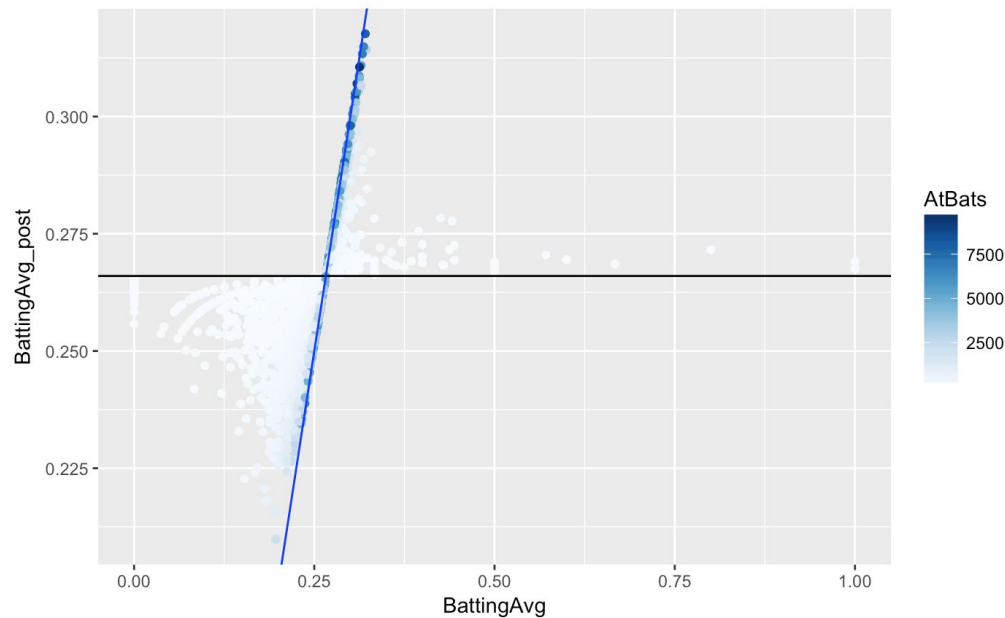
playerID <chr>	name <chr>	bats <fctr>	Hits <int>	AtBats <int>	BattingAvg <dbl>	alpha1 <dbl>	beta1 <dbl>	BattingAvg_post <dbl>
mathije01	Jeff Mathis	R	401	2038	0.1967615	527	1984.6842	0.2098194
zuninmi01	Mike Zunino	R	219	1125	0.1946667	345	1253.6842	0.2158025
cashke01	Kevin Cash	R	117	641	0.1825273	243	871.6842	0.2179990
woodbr01	Brandon Wood	R	130	700	0.1857143	256	917.6842	0.2181166
sadledo01	Donnie Sadler	R	97	537	0.1806331	223	787.6842	0.2206426
hicksbr01	Brandon Hicks	R	45	294	0.1530612	171	596.6842	0.2227478

As we can see, those batters with low *AtBats* are not taken into consideration thanks to the property of conjugate prior of Binomial distribution. The posterior mean is the weighted average of prior mean and data, larger the size of data(in here which means the amount of *AtBats* ), more weight is given to data.



The dash line is our prior, and those lines that are close to the dash line are the posterior distribution of the top 5 batters from model 1, they used to have *BattingAverage* close to 1, however by applying prior information, their posterior mean shrinkage to prior mean. That is because the small size of data and we prefer they are just average level; However for those Top 5

batters from the second model, due to their large amount of *AtBats*, the posterior mean are more close to their own sample mean. This is the effect called shrinkage.



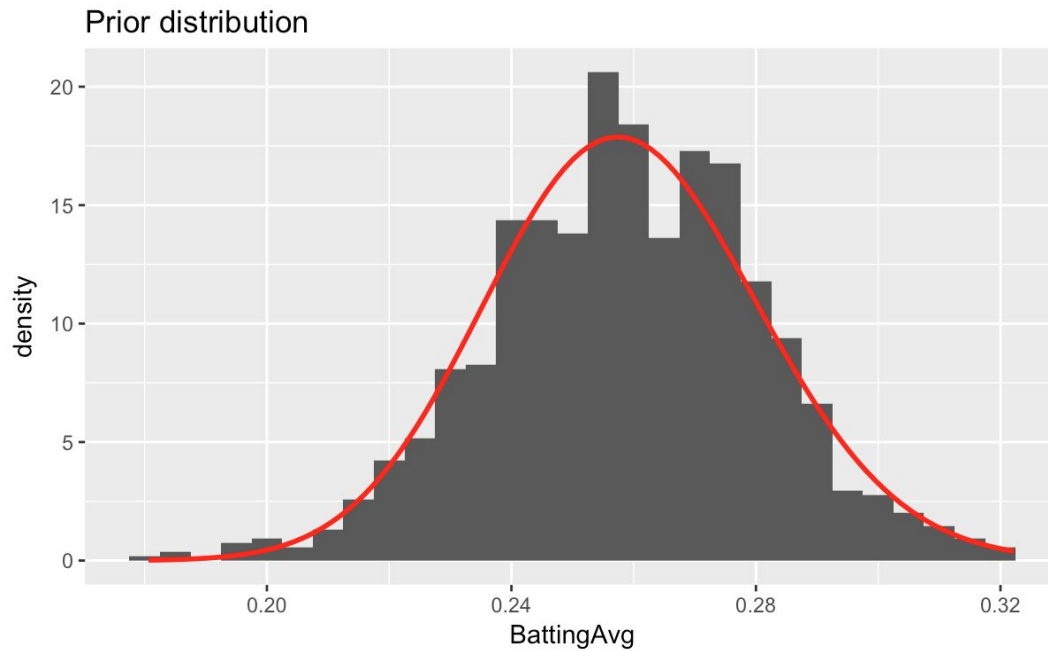
(3) Empirical Bayesian Method: When estimating a parameter of an individual, then the distribution of the whole group in which the individual belongs to can be treat as a prior.

Empirical Bayes methods are procedures for statistical inference in which the prior distribution is estimated from the data. This approach stands in contrast to standard Bayesian methods, for which the prior distribution is fixed before any data are observed. Despite this difference in perspective, empirical Bayes may be viewed as an approximation to a fully Bayesian treatment of a hierarchical model wherein the parameters at the highest level of the hierarchy are set to their most likely values, instead of being integrated out. Empirical Bayes, also known as maximum marginal likelihood, represents one approach for setting hyperparameters. It is an approximation to more exact Bayesian methods and with the amount of data we have, it's a very good approximation. To put it simply, we get the prior  $Beta(\alpha_0, \beta_0)$  from our data and use Maximum likelihood to estimate the two parameter  $\alpha_0$  and  $\beta_0$ . In order to estimate a better prior distribution of all players, we filtered out all the players that have fewer than 500 *AtBats* to get the more representative players with less noise to calculate the prior.

```

> alpha0_eb
shape1
99.4927
> beta0_eb
shape2
285.158

```



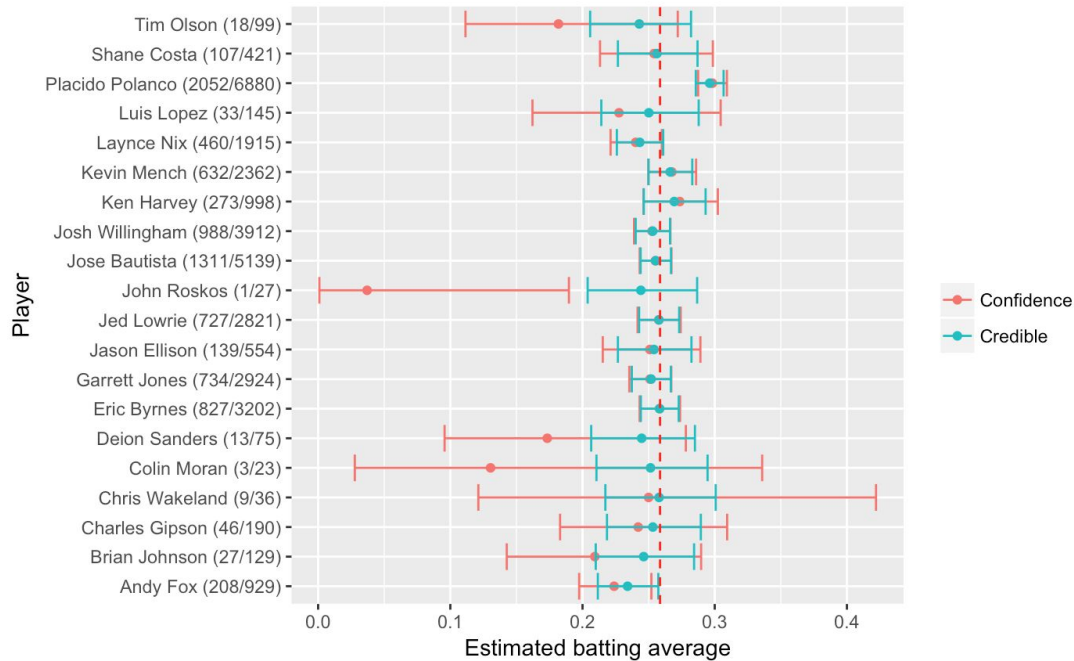
playerID <chr>	name <chr>	bats <fctr>	Hits <int>	AtBats <int>	BattingAvg <dbl>	alpha1 <dbl>	beta1 <dbl>	BattingAvg_post <dbl>	EB_BattingAvg_post <dbl>
cabremi01	Miguel Cabrera	R	2519	7853	0.3207691	2645	5681.684	0.3176535	0.3178689
guerrvl01	Vladimir Guerrero	R	2092	6570	0.3184170	2218	4825.684	0.3148920	0.3151118
bondsba01	Barry Bonds	L	925	2871	0.3221874	1051	2293.684	0.3142300	0.3146814
heltoto01	Todd Helton	L	2141	6761	0.3166691	2267	4967.684	0.3133516	0.3135464
ramirma02	Manny Ramirez	R	1642	5213	0.3149818	1768	3918.684	0.3109017	0.3111114
suzukic01	Ichiro Suzuki	L	3030	9689	0.3127258	3156	7006.684	0.3105479	0.3106612

playerID <chr>	name <chr>	bats <fctr>	Hits <int>	AtBats <int>	BattingAvg <dbl>	alpha1 <dbl>	beta1 <dbl>	BattingAvg_post <dbl>	EB_BattingAvg_post <dbl>
mathije01	Jeff Mathis	R	401	2038	0.1967615	527	1984.6842	0.2098194	0.2065889
zuninmi01	Mike Zunino	R	219	1125	0.1946667	345	1253.6842	0.2158025	0.2109711
cashke01	Kevin Cash	R	117	641	0.1825273	243	871.6842	0.2179990	0.2110784
woodbr01	Brandon Wood	R	130	700	0.1857143	256	917.6842	0.2181166	0.2115821
hicksbr01	Brandon Hicks	R	45	294	0.1530612	171	596.6842	0.2227478	0.2129118
sadledo01	Donnie Sadler	R	97	537	0.1806331	223	787.6842	0.2206426	0.2131965

As we can see, it gives the same Top5 and Tail 5 ranking result as the informative prior.

Then We randomly sample 20 players to see their confidence interval and credible interval





The main difference between credible interval and confidence interval is that the former one takes prior into consideration and then realizes shrinkage.

- Model 3: Bayesian hierarchical Modeling

From the above model, we set different priors based on our general information about baseball or based on empirical Bayesian Method which actually gather information from data and then get the estimate for posterior distribution to choose batters. But we always assume a same prior for every batters in each model.

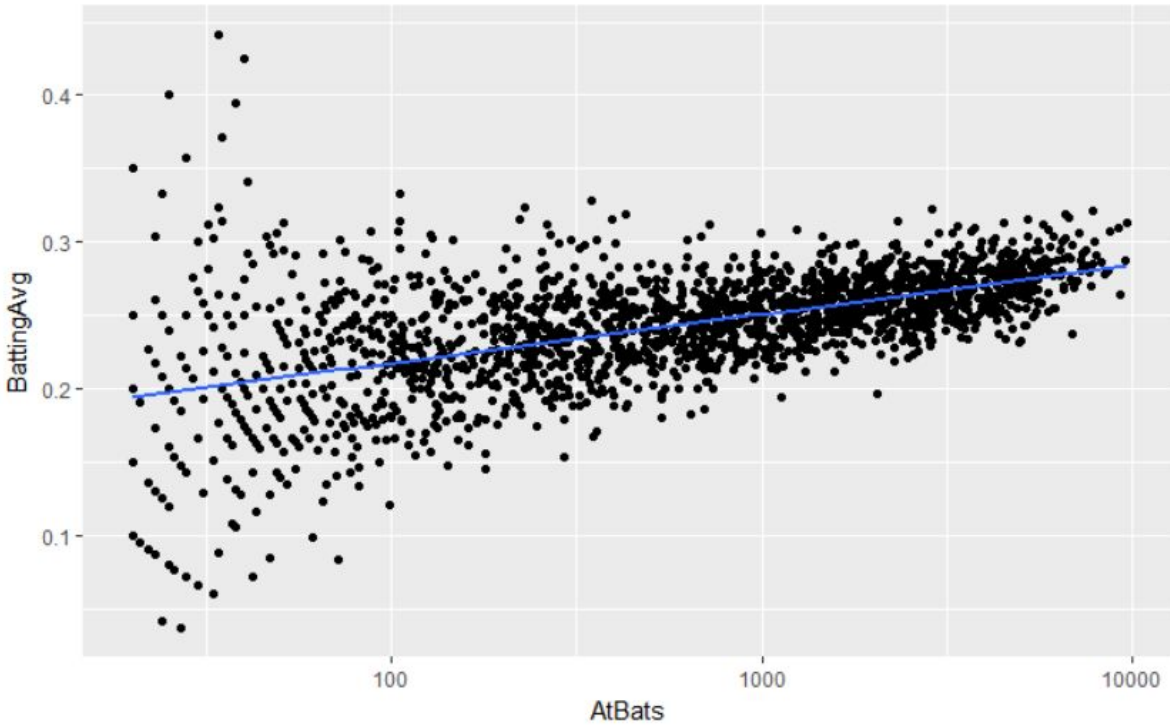
Our next question is: will some variables affect our priors? So in this model, we try to assign different prior distributions for batters with different characteristics. For example, here we typically choose different *AtBats* and different dominant hands which is listed as “*Bats*” in our dataset. This is called Bayesian hierarchical modeling, means the prior parameter for each batter are not fixed and will depend on some other variables.

(1) Experienced Batters tend to play better?

First, we will take *AtBats* into consideration.

Due to the shrinkage property of Bayesian models (estimate for true accuracy is  $\frac{\alpha_0 + Hits}{\alpha_0 + \beta_0 + N}$ ), the estimate of parameter  $\theta$  for batters with high *AtBats* is close to their previous performance recorded in the dataset as *BattingAverage*, and the estimate for batters with low *AtBats* is close to the mean of prior distribution which is the overall average as we set in model 2(3). This is the

result for shrinkage property of Bayesian models, but there is a problem. As we can see from the plot below, when the number *AtBats* increases, the *BattingAverage* also increases.



This is reasonable, because when the batters have more *AtBats*, they have more practical experience and they are more likely to hit during the game.

So, for the batters with low *AtBats*, our previous model tend to overestimate their true batting accuracy, but for the batters with high *AtBats*, it seems we are having a more reasonable estimate. How to deal with this problem?

One way to improve our previous model is to add the influence of *AtBats* in our model, which means our prior will be affected by the number of *AtBats*, thus every batters will have different priors based on their *AtBats*.

First, as is mentioned above in the previous model,

for each batter  $i$ , we assume  $p(y_i|\theta_i) = \text{Bin}(N, \theta_i)$  where  $i \in \{1, 2, \dots, 2232\}$

Also, we assume our prior is a Beta distribution:  $p(\theta_i) = \text{Beta}(\alpha_0, \beta_0)$ .

Here,  $N$  is the total *AtBats* for every single batter. For this project, our goal is to estimate every batter's true accuracy, thus we are interested in the Bayesian estimator of the distribution. For

Beta distribution  $Beta(\alpha_0, \beta_0)$  under squared loss, the Bayesian estimator is just  $\frac{\alpha_0}{\alpha_0 + \beta_0}$ . To make it clear, we try to transform parameters  $\alpha_0$  and  $\beta_0$  to  $\mu_0$  and  $\sigma_0$  by letting:

$$\mu_0 = \frac{\alpha_0}{\alpha_0 + \beta_0} \text{ and } \sigma_0 = \frac{1}{\alpha_0 + \beta_0}.$$

Here,  $\mu_0$  is the Bayesian estimator of prior distribution and is the parameter we try to obtain after we have real data. Then our initial model becomes:

$$p(y_i | \theta_i) = \text{Bin}(AtBats, \theta_i) \text{ where } i \in \{1, 2, \dots, 2232\}$$

$$p(\theta_i) = \text{Beta}\left(\frac{\mu_0}{\sigma_0}, \frac{1 - \mu_0}{\sigma_0}\right).$$

Next, we will contain the influence of *AtBats* to the prior distribution in our model in order to provide different prior distributions for every different batters. Basically, we just want the Bayesian estimator of prior distribution for every different batters  $\mu_i$ ,  $i \in \{1, 2, \dots, 2232\}$  to be related to *AtBats*. Notice that the range of *AtBats* is from 1 to 9689 in our dataset which affected the result too much if we contain it as a linear term in our model, so we simply take the log of *AtBats* and contain it in our linear model. We set:

$$\mu_i = \mu_0 + \lambda \times \log(AtBats)$$

Thus we can get our updated prior distribution for every batter which is  $p(\theta_i) = \text{Beta}\left(\frac{\mu_i}{\sigma_0}, \frac{1 - \mu_i}{\sigma_0}\right)$  (\*)

(This particular model is called beta-binomial regression.) This is a little bit strange because not only does the prior distribution influences the likelihood of the real data, but the real data also influences the prior distribution.

Next, we try to fit parameters in our prior distribution (\*) which are  $\mu_0, \lambda, \sigma_0$  using maximum likelihood. Here, we use “gamlss” package to fit the beta-binomial regression model using the maximum likelihood.

The coefficients results are:

parameter <fctr>	term <chr>	estimate <dbl>	std.error <dbl>	statistic <dbl>	p.value <dbl>
mu	(Intercept)	0.15170322	0.0031952427	47.47784	0.000000e+00
mu	log(AtBats)	0.01432317	0.0004368415	32.78803	8.850021e-193
sigma	(Intercept)	-6.72320574	0.0592102640	-113.54798	0.000000e+00

So we can see the coefficients for our three parameters in our regression model:

$$\mu_0 = 0.1517$$

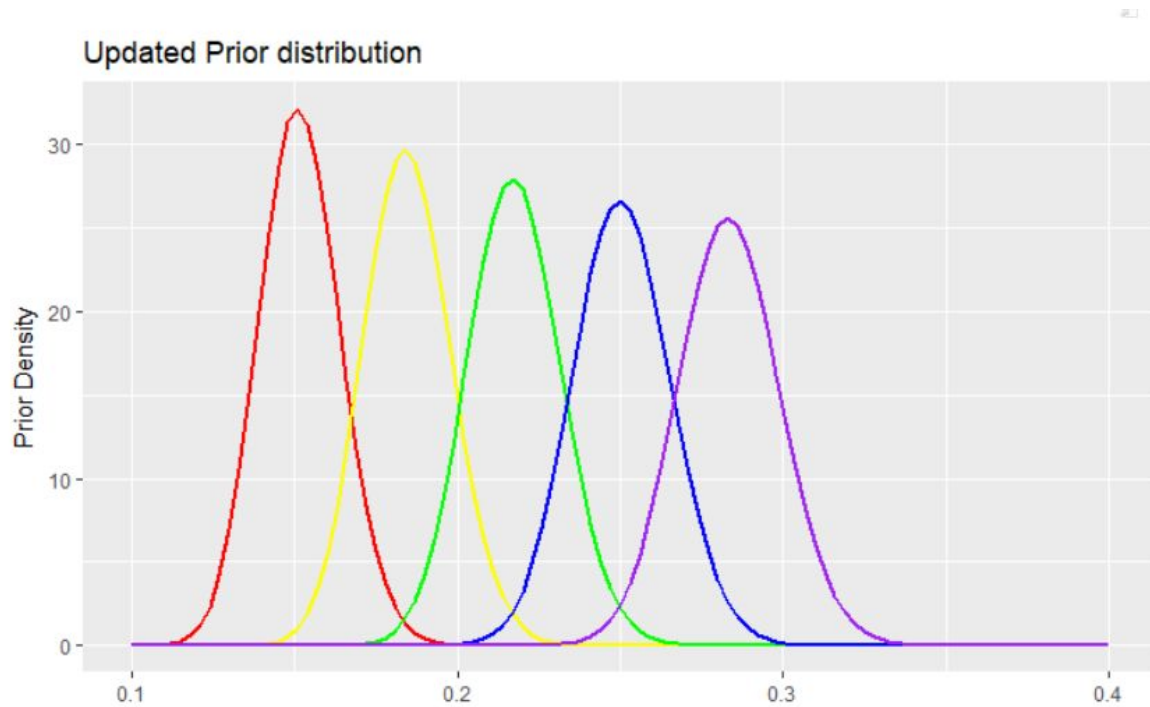
$$\lambda = 0.0143$$

$$\log(\sigma_0) = -6.723 \Rightarrow \sigma_0 = e^{-6.723} = 0.0012$$

Also, we can see that the  $p$ -value for  $\log(AtBats)$  is really small, which means taking  $AtBats$  into consideration is reasonable.

Then we plug in the  $AtBats$  in our regression models and get prior distribution  $p(\theta_i) = \text{Beta}\left(\frac{\mu_i}{\sigma_0}, \frac{1-\mu_i}{\sigma_0}\right)$  for every batters.

We choose some prior distribution with different  $AtBats$  and the density functions are shown below.



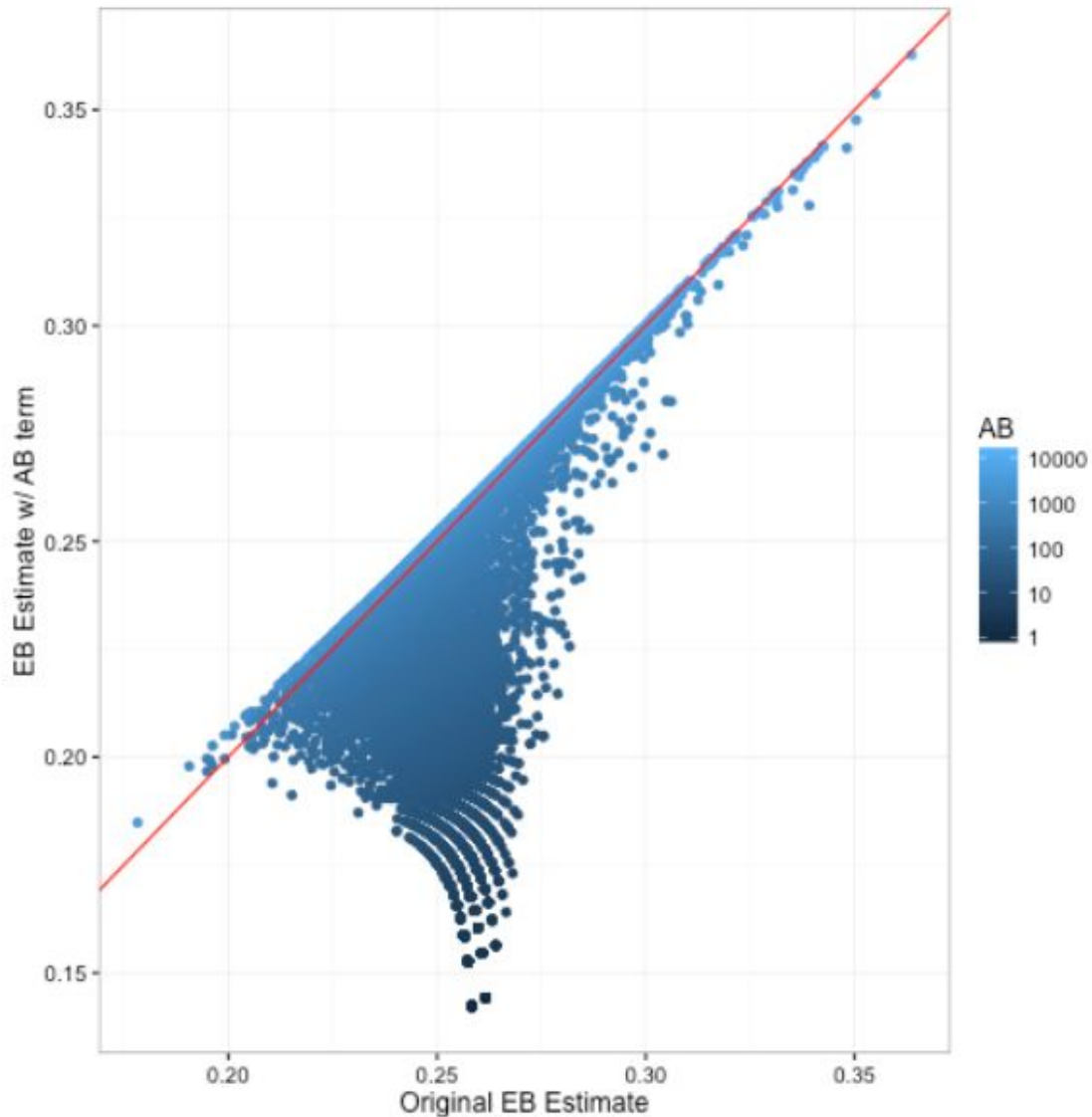
Red:  $AtBats=1$ , Yellow:  $AtBats=10$ , Green:  $AtBats=100$ ,

Blue:  $AtBats=1000$ , Purple:  $AtBats=10000$

Then, based on the prior distributions for batters with different  $AtBats$ , we can calculate their posterior distribution given the real data for them.

Then we calculate the Bayesian estimate for posterior distribution under squared error loss which is also the mean of the posterior distribution.

The results for the estimates are shown below.



(2) Left-handed or right-handed?

Now we want to go further by trying to find more structure of our data. We notice that in our dataset, there is a variable called “bat” which list the batters’ dominant hand. Suppose there are two batters who share the same *AtBats* and Hits which means the estimate for their true accuracy will be exactly the same based on all of our previous model. But the only difference is that one player is left-handed and another is right-handed. And we hope to see if this will affect our estimate for their accuracy. If so, we will also add the influence of “*bats*” in our model.

Same as the model with different *AtBats* above, we just add another variable *bats* in our regression model. Set right-handed batters as the baseline and use “gamlss” function to fit the maximum likelihood estimator, we just get the results below:

parameter <fctr>	term <chr>	estimate <dbl>	std.error <dbl>	statistic <dbl>	p.value <dbl>
mu	(Intercept)	0.1514003660	0.0032169312	47.063601	0.000000e+00
mu	log(AtBats)	0.0143041068	0.0004367956	32.747826	2.306952e-192
mu	batsB	-0.0006405712	0.0015137873	-0.423158	6.722208e-01
mu	batsL	0.0016889352	0.0011460300	1.473727	1.406964e-01
sigma	(Intercept)	-6.7260545279	0.0593051216	-113.414396	0.000000e+00

From the coefficients, again we can see that the  $p$  – value for  $\log(\text{AtBats})$  is really small, which means taking *AtBats* into consideration is quite reasonable. We could also see that the  $p$  – value for left-handed batters is 0.14 which is not so low, under significance level 0.1, it is not statistically significant. But it can still have a really small effect on our choice of prior distribution. Meanwhile, the  $p$  – value for both-handed batters is 0.67 which is high, so we decide not to take both-handed batters into consideration.

Our new beta-binomial regression is:

$$\mu_i = \mu_0 + \lambda \times \log(\text{AtBats}) + \psi \times I(\text{Bats} = L)$$

And the coefficients are:

$$\mu_0 = 0.1514$$

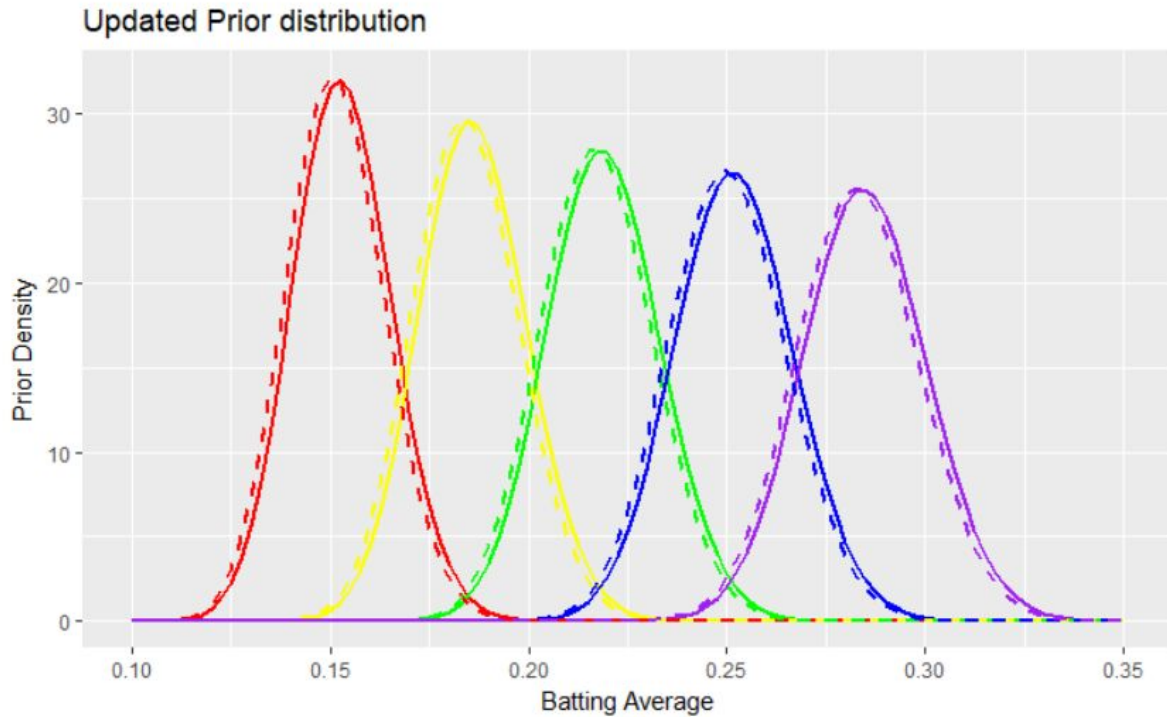
$$\lambda = 0.0143$$

$$\psi = 0.0017$$

$$\log(\sigma_0) = -6.726 \Rightarrow \sigma_0 = e^{-6.726} = 0.0012$$

Thus we can get our updated prior distribution for batters with different dominant hands and different *AtBats* which is  $p(\theta_i) = \text{Beta}\left(\frac{\mu_i}{\sigma_0}, \frac{1-\mu_i}{\sigma_0}\right)$

We choose some prior distribution with different *AtBats* and different dominant hands. The density functions are shown below.



Dashed lines: right-handed batters, Full lines: left-handed batters

Red:  $AtBats=1$ , Yellow:  $AtBats=10$ , Green:  $AtBats=100$ ,

Blue:  $AtBats=1000$ , Purple:  $AtBats=10000$

We use these priors to improve our estimates for every batters. As we can see from the plot above, for the batters with same  $AtBats$  , we will give left-handed batters a higher prior. Then we can calculate their posterior distribution with their own priors and their real data.

Then we calculate the Bayesian estimate for posterior distribution under squared error loss which is also the mean of the posterior distribution.

## ● Conclusion

In this project, we analyze the baseball dataset as well as showing the process of how we start to think problem from frequentist to Bayesian and also the process of the Bayesian hierarchical modeling, combining what we have learnt in class with some new knowledge such as empirical Bayesian method and Beta-Binomial Regression. From Model 1 and 2, The Top Five Players are Miguel Cabrera, Vladimir Guerrero, Barry Bonds, Todd Helton, Manny Ramirez, and as what we can get from the result is that different ways of setting prior all come to the conclusion that they will prefer experienced player more. From Model 3 and Model 4, we found some structures inside the dataset. After running the regression, we found out it is reasonable to conduct a Bayesian hierarchical model by taking *AtBats* and dominant-hand into consideration. So we tent to give batters with higher *AtBats* and left-handed batters a higher prior. Also, by adding more layer of information to the original data, we can get a more informative and confident result.

## ● Discussion

For all the analysis above, we just ignore the influence of time and assume that batters tend to have a stable performance during all the years in their career. But if we can consider the time variable, we may able to find some more interesting structure of the data. Also for future study, we can go deep in to the empirical Bayesian method to see what's the advantage and disadvantage of using the method that seems quite different from what we learnt in class.



# Appendix

## Bibliography

[1]Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, Donald B. Rubin(2014)*Bayesian Data Analysis*. CRC Press;Taylor & Francis Group

[2]Peter D. Hoff(2009)*A first course in Bayesian Statistical Methods*. Springer

[3]David Robinson(2017)*Understanding empirical Bayes estimation (using baseball statistics)*

Retrieved from [http://varianceexplained.org/r/empirical\\_bayes\\_baseball/](http://varianceexplained.org/r/empirical_bayes_baseball/)

[4]David Robinson(2017)*Understanding empirical Bayes hierarchical modeling (using baseball statistics)*

Retrieved from [http://varianceexplained.org/r/empirical\\_bayes\\_baseball/](http://varianceexplained.org/r/empirical_bayes_baseball/)

## Contribution

Hanying Ji: Mainly taking responsibility of data processing; Analyzing the data by using the method that setting the same prior for all batters in different ways.

Jiaqian Yu: Analyzing the data by using Bayesian hierarchical modeling with different *AtBats* and different dominant hands.