

# Empirical Bayes

Q1 Suppose Mookie Betts' batting average midway through the season is .300. Using no other information, predict his end-of-season batting average.

Model mid-season batting average is

$$\frac{H}{N} \sim \frac{1}{N} \cdot \text{Binomial}(N, p)$$

$H = \# \text{ hits}$ ,  $N = \# \text{ at-bats}$

As discussed previously, the MLE of a binomial is

$$\hat{p}_{\text{MLE}} = \frac{H}{N}$$

and this is our prediction.

Concretely, we predict Mookie's end-of-season batting average to be .300.

Yesterday we used a prior, but we can't do that here since we have no other information.

Q2 Suppose we know each player's batting average midway through the 2023 season. Using no information from any previous season, i.e. only using these 2023 mid-season batting averages, predict each player's end-of-season batting average.

One approach:

Notation player index  $i$   
 Player  $i$ 's # hits  $H_i$   
 Player  $i$ 's # at-bats  $N_i$   
 Player  $i$ 's batting average  $\frac{H_i}{N_i}$

Model  $\frac{H_i}{N_i} \sim \text{Binomial}(N_i, p_i)$

$$\underline{\text{MLE}} \quad \hat{p}_i^{(\text{MLE})} = \frac{H_i}{N_i}$$

Can we do better??

Previous idea : Shrinkage

But, we only have access to batting averages from 2023 and nothing else.

So, we can't shrink to prior information from previous seasons.

What can we shrink to??

{ Well, we have the midseason batting average of each baseball player, so perhaps we can pool information and shrink to the overall mean batting average.

Idea: to predict the batting average of Mookie Betts, use the batting averages of every other player!

Insight: Mookie Betts is a baseball player  
use that information!

Notation player index  $i$

player  $i$ 's # hits  $H_i$

player  $i$ 's # at-bats  $N_i$

player  $i$ 's batting average  $X_i = \frac{H_i}{N_i}$

Model  $X_i \sim \frac{1}{N_i} \text{Binomial}(N_i, P_i)$

Remove players from the dataset with small  $N_i$  (say,  $N_i < 25$ ).  
Since  $N_i$  large,

$$X_i \approx N\left(P_i, \frac{P_i(1-P_i)}{N_i}\right)$$

by Central Limit Theorem.

Simplification

$P_i$  unknown, so variance  $\sigma_i^2 = \frac{P_i(1-P_i)}{N_i}$  unknown

Much easier to work with known variance

so let's just assume for simplicity

that  $\sigma_i^2 = \frac{C}{N_i}$  for some known

constant  $C$ .

We'll use  $C = 0.035$ , treat this as a tuning hyperparameter... (one using a previous season).

In practice, we use a variance stabilizing transform  $h(X)$  which transforms the batting average so that it has a known variance... .

Now,  $X_i \sim N(\mu, \sigma_i^2)$   $\sigma_i^2$  known... .

## Parametric Bayesian Model

$$\text{Prior: } \begin{cases} X_i \sim N(\mu_i, \sigma_i^2) \\ \mu_i \sim N(\mu, \tau^2) \end{cases}$$

insight: each player  $i$  is a baseball player  
so, let's pool information!  
there is a mean  $\mu$  across all baseball players  $\mu_i$   
and some s.d.  $\tau$

$$\mu_i = \mu; \text{ unknown} \quad \sigma_i = \frac{\sqrt{C}}{\sqrt{N_i}} \text{ known}$$
$$\mu, \tau \text{ unknown}$$

MLE: ignore the prior, only use the information relevant to player i

$$\hat{\mu}_i^{(MLE)} = X_i \quad \text{observed mid-season batting average}$$

Bayesians: use the prior!

The Bayesian estimate is the posterior mean, i.e. the mean of the dist of  $\mu_i$  after seeing the data  $\{X_i\}$ ,

$$\hat{\mu}_i^{(\text{Bayes})} = \mathbb{E}(\mu_i | X_i)$$

need the posterior dist  $\mu_i | X_i$ ,

$$P(\mu_i | X_i) = \frac{P(X_i | \mu_i) P(\mu_i)}{P(X_i)} \quad \text{Bayes Rule}$$

$$\propto \underbrace{P(X_i | \mu_i)}_{N(\mu_i; \sigma_i^2)} \underbrace{P(\mu_i)}_{N(\mu, \tau^2)}$$

since  $P(X_i)$  has no  $\mu_i$  term

$$= P(N(\mu_i, \sigma_i^2) = x_i) \cdot P(N(\mu, \tau^2) = \mu_i)$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right) \cdot \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{1}{2} \frac{(\mu_i - \mu)^2}{\tau^2}\right)$$

$$\propto \exp\left(-\frac{1}{2} \left[ \frac{x_i^2}{\sigma_i^2} - 2\mu_i x_i + \frac{\mu_i^2}{\sigma_i^2} + \frac{\mu_i^2}{\tau^2} - 2\mu_i \mu + \frac{\mu^2}{\tau^2} \right]\right)$$

$$\propto \exp\left(-\frac{1}{2} \left[ \mu_i^2 \left( \frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right) - 2\mu_i \left( \frac{x_i}{\sigma_i^2} + \frac{\mu}{\tau^2} \right) \right]\right)$$

$$\propto \exp\left(-\frac{1}{2} \left( \frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right) \left[ \mu_i^2 - 2\mu_i \frac{\left( \frac{x_i}{\sigma_i^2} + \frac{\mu}{\tau^2} \right)}{\left( \frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right)} \right]\right)$$

$$\propto \exp\left(-\frac{1}{2} \left( \frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right) \left[ \mu_i - \frac{\left( \frac{x_i}{\sigma_i^2} + \frac{\mu}{\tau^2} \right)}{\left( \frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right)} \right]^2\right)$$

$$\implies \mu_i / x_i \sim N\left(\frac{\left( \frac{x_i}{\sigma_i^2} + \frac{\mu}{\tau^2} \right)}{\left( \frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right)}, \frac{1}{\left( \frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right)}\right)$$

$\Rightarrow$  Posterior Mean

$$\hat{\mu}_i^{(\text{Bayes})} = \frac{\bar{x}_i + \frac{\mu}{\tau^2}}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}} = \mu + \frac{\tau^2}{\tau^2 + \sigma_i^2} (\bar{x}_i - \mu)$$

$$\bar{x}_i = \frac{H_i}{N_i}, \quad \sigma_i^2 = \frac{C}{N_i}$$

$$\Rightarrow \hat{\mu}_i^{(\text{Bayes})} = \frac{\frac{H_i}{C} + \frac{\mu}{\tau^2}}{\frac{N_i}{C} + \frac{1}{\tau^2}}$$

\* Looks a lot like  $\frac{W+W'}{W+L+W'+L}$  from last lecture!

Problem:  $\mu$  and  $\tau$  are unknown quantities

Empirical Bayes: Estimate  $\mu$  and  $\tau$  from observed data!

$$\hat{\mu}_i^{(EB)} = \frac{\frac{H_i}{C} + \frac{\hat{\mu}}{\frac{1}{\tau^2}}}{\frac{N_i}{C} + \frac{1}{\frac{1}{\tau^2}}}$$

$$\left. \begin{array}{l} \text{If } \tau^2 = 0, \quad \mu_i \sim N(\mu, \tau^2) \approx N(\mu, 0) = \mu \\ \Rightarrow \hat{\mu}_i = \bar{\mu} \end{array} \right\}$$

$$\left. \begin{array}{l} \text{If } \tau^2 = \infty, \quad \mu_i \sim N(\mu, \infty) \approx \text{Uniform}(-\infty, \infty) \\ \text{and } \hat{\mu}_i = \frac{H_i}{N_i} = \bar{x}_i = \hat{\mu}_i^{(\text{MLE})} \end{array} \right\}$$

$$\left. \begin{array}{l} \text{Else, } \hat{\mu}_i = \bar{\mu} + \frac{\tau^2}{\tau^2 + \sigma_i^2} (\bar{x}_i - \bar{\mu}) \\ \text{is closer to } \bar{\mu} \text{ if } \frac{\tau^2}{\tau^2 + \sigma_i^2} \text{ is small} \\ \text{is closer to } \bar{x}_i \text{ if it is large} \end{array} \right\} \begin{array}{l} \text{large } \sigma_i^2 \\ \text{small } N_i \end{array} \begin{array}{l} \text{small } \sigma_i^2 \\ \text{large } N_i \end{array}$$

\* use  $\hat{\mu}^{(\text{MLE})}$ ,  $\hat{\tau}^2_{\text{MLE}}$  to estimate  $\mu, \tau^2$

Model  $\begin{cases} X_i \sim N(\mu_i, \sigma_i^2) \\ \mu_i \sim N(\mu, \tau^2) \end{cases}$

Get Rid of  $\mu_i$

Marginal distribution  $X_i \sim N(\mu, \tau^2 + \sigma_i^2)$   
By Bayes Rule.

Log-Likelihood

$$\begin{aligned} L(\mu, \tau^2) &= \log P(X | \mu, \tau^2) \\ &= \log \prod_{i=1}^n P(X_i | \mu, \tau^2) \quad \text{by independence} \\ &= \sum_{i=1}^n \log P(N(\mu, \tau^2 + \sigma_i^2) = X_i) \\ &= \sum_{i=1}^n \log \left[ \frac{1}{\sqrt{2\pi(\tau^2 + \sigma_i^2)}} e^{-\frac{(X_i - \mu)^2}{2(\tau^2 + \sigma_i^2)}} \right] \\ &= \sum_{i=1}^n \left\{ \log \left( \frac{1}{\sqrt{2\pi(\tau^2 + \sigma_i^2)}} \right) - \frac{(X_i - \mu)^2}{2(\tau^2 + \sigma_i^2)} \right\} \\ &\propto -\frac{1}{2} \sum_i \log(\tau^2 + \sigma_i^2) - \frac{1}{2} \sum_i \frac{(X_i - \mu)^2}{\tau^2 + \sigma_i^2} \end{aligned}$$

$$\text{MLE } (\hat{\mu}, \hat{\tau}^2) = \underset{\mu, \tau^2}{\operatorname{Argmin}} L(\mu, \tau^2) \Rightarrow \text{solve } \frac{dL}{d\mu} = 0, \frac{dL}{d\tau^2} = 0.$$

$$\frac{\partial L}{\partial \mu} = \frac{1}{2} \sum_i \frac{2(X_i - \mu)}{\tau^2 + \sigma_i^2} = 0$$

$$\Rightarrow \mu = \frac{\sum_i X_i / (\tau^2 + \sigma_i^2)}{\sum_i 1 / (\tau^2 + \sigma_i^2)}$$

$$\frac{\partial L}{\partial \tau^2} = -\frac{1}{2} \sum_i \frac{1}{\tau^2 + b_i^2} + \frac{1}{2} \sum_i \frac{(x_i - \mu)^2}{(\tau^2 + b_i^2)^2} = 0$$

$$\Rightarrow \sum_i \frac{(x_i - \mu)^2}{(\tau^2 + b_i^2)^2} = \sum_i \frac{1}{\tau^2 + b_i^2}$$

Problem  $\mu$  is in terms of  $\tau^2$  and  $\tau^2$  is in terms of  $\mu$ , but we need both...

Iteratively Solve for  $\hat{\mu}, \hat{\tau}^2$ :

$$\hat{\mu} = \frac{\sum x_i / (\hat{\tau}^2 + b_i^2)}{\sum 1 / (\hat{\tau}^2 + b_i^2)}$$

$$\sum \frac{(x_i - \hat{\mu})^2}{(\hat{\tau}^2 + b_i^2)^2} = \sum \frac{1}{\hat{\tau}^2 + b_i^2}$$

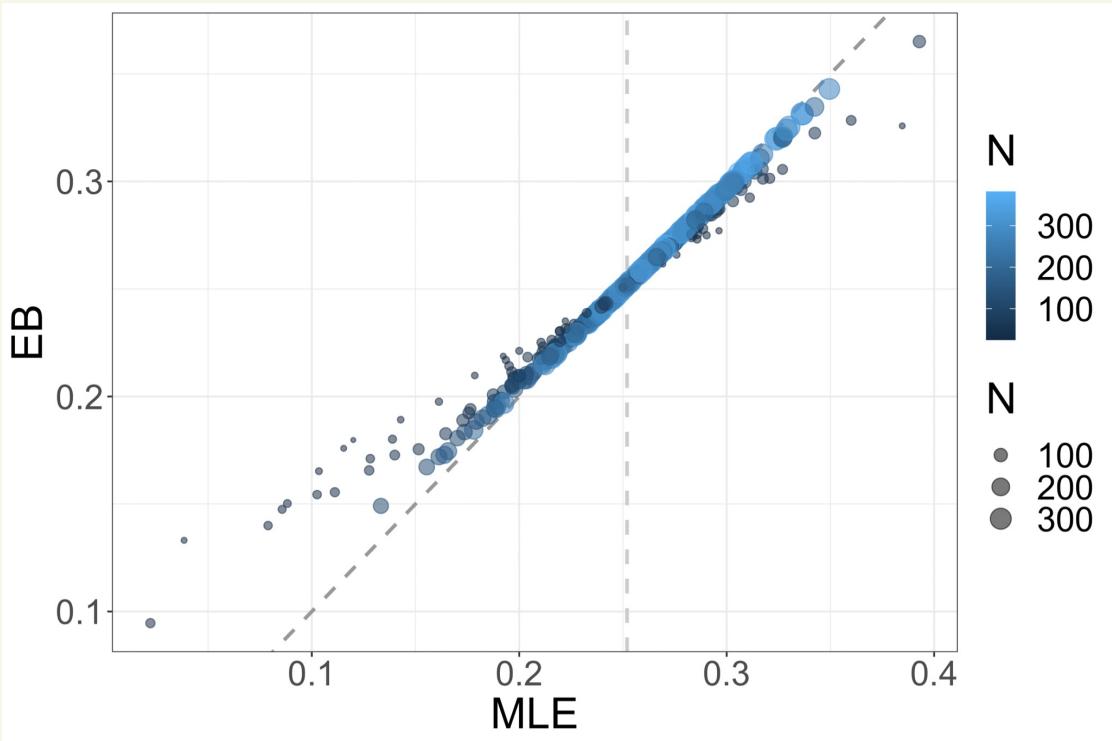
$\mu^{(0)}$  = educated guess

$\tau^{2(0)}$  = educated guess

while not  $|\mu^{(t)} - \mu^{(t-1)}| \leq \delta$  and  $|\tau^{2(t)} - \tau^{2(t-1)}| < \delta$ :

$$\mu^{(t)} \leftarrow \frac{\sum_i x_i / (\tau^{2(t)} + b_i^2)}{\sum_i 1 / (\tau^{2(t)} + b_i^2)}$$

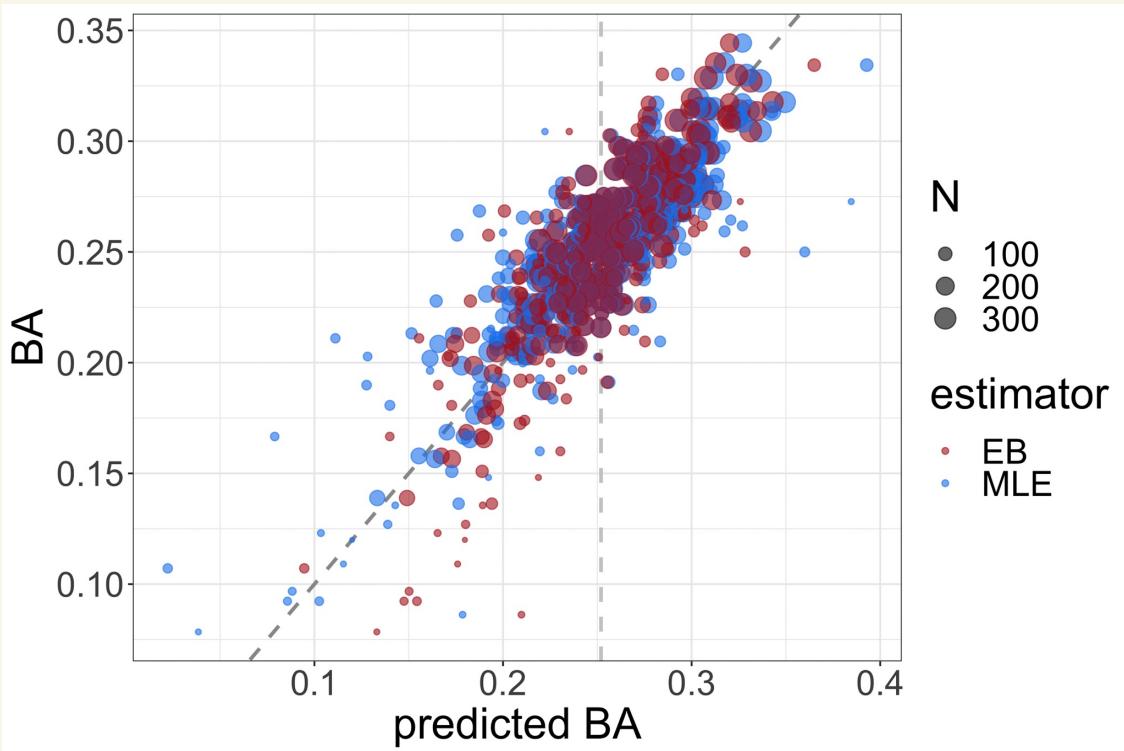
$\tau^{2(t)}$  ← the solution of  $\sum_i \frac{(x_i - \mu)^2}{(\tau^2 + b_i^2)^2} = \sum_i \frac{1}{\tau^2 + b_i^2}$   
(use `uniRoot in R`)



\* Players with smaller  $N$  have  $\hat{\mu}^{EB}$  shrunk towards the overall mean, while players with larger  $N$  have  $\hat{\mu}^{(EB)} = \hat{\mu}^{(MLE)}$  their end-of-season BA

rmse_MLE	rmse_EB
0.02629828	0.02383808

\*  $\hat{\mu}^{(EB)}$  predict better the actual end-of-stn BA!



\* There really isn't a huge difference b/t the predictions though..

### Takeaway

- shrinkage towards the overall mean helps prediction when have smaller sample sizes
- sharing information helps!

# Consultants Dilemma (Stein's Paradox)

Billy Beane asks you to project the end-of-season performance of Elvis Andrus. Using just Andrus' at bats alone, the MLE is best.

Using all MLB at bats, a shrinkage estimator is better than the MLE on average across all batters.

But, we only need a good prediction for one player, Andrus, not a good prediction on average over all players.

Which estimator do we use ???

An intuitive explanation is that optimizing for the mean-squared error of a *combined* estimator is not the same as optimizing for the errors of separate estimators of the individual parameters. In practical terms, if the combined error is in fact of interest, then a combined estimator should be used, even if the underlying parameters are independent. If one is instead interested in estimating an individual parameter, then using a combined estimator does not help and is in fact worse.

James Stein's Theorem A similar Empirical Bayes shrinkage estimator, the James Stein estimator  $\hat{\mu}_i^{(JS)}$ , everywhere dominates the MLE,

$$E_{\mu} \left[ \sum_{i=1}^n (\hat{\mu}_i^{(JS)} - \mu_i)^2 \right] < E_{\mu} \left[ \sum_{i=1}^n (\hat{\mu}^{(MLE)} - \mu_i)^2 \right],$$

i.e. it has smaller expected mean squared error.  
 Expected because  $\hat{\mu}_i$  is random because  $X_i$  is random.

# Shrinkage, with Access to Previous Season's Data

Parametric Empirical Bayes Estimator:

$$\hat{\mu}_i^{(PEB)} = \hat{\mu} + \frac{\tau^2}{\tau^2 + \sigma_i^2} (x_i - \hat{\mu})$$

General Shrinkage Estimator:

$$\hat{\mu}_i = \hat{\mu} + \beta (x_i - \hat{\mu})$$

Procedure  $\{x_i\}$  data

$\hat{\mu}$  = overall mean  
(sample mean of data  
from the season of  
datapoint  $i$ )

$\hat{\mu}_i$  = response column  
= known end-of-season  
battery avg of datapoint  $i$   
(player  $i$ ) from a previous  
season

Estimate  $\beta$  using Regression,  
call it  $\hat{\beta}$

Then, for future prediction,

$$\hat{\mu}_i = \hat{\mu} + \hat{\beta} (x_i - \hat{\mu}).$$