

Empirical Bayes

→ player quality

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

Model Observe Mookie's mid-season data $\{H, N\}$

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

p is an unbiased measure of his quality
successes (hits) in N trials (at-bats)

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

$$\hat{p}^{(MLE)} = \frac{H}{N} = BA = .300$$

$\frac{W}{W+L}$ $\frac{H}{N-H}$ $\frac{W+W'}{W+L+W'+L'}$

* Yesterday we used a Prior to stabilize the MLE early in the season. But here we assumed we had no other info, so we can't do that.

Q Suppose we know each player's batting average midway through the season. Using no info from previous seasons, predict each player's end-of-season batting average.

Data $\{H_i, N_i\}$ player i , $BA_i = \frac{H_i}{N_i}$
| |
hits # at bats

Model $\frac{H_i}{N_i} \sim \frac{1}{N_i} \text{Binomial}(N_i, p_i)$

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

observable measure of
batter's quality

Can we do better?

Prior idea

Mon: Ridge \rightarrow shrinkage \rightarrow
make your coeff. estimates
smaller so as not to
overfit

Tues: predict end from WP \rightarrow
add prior / fake data
to enhance our estimates

Q What are we going to shrink to?

towards the overall mean batting avg.
 \rightarrow Mookie is a baseball player
How to accomplish this?

\rightarrow PRIOR

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

Binomial math can be difficult.

Central Limit Theorem \rightarrow the sum of any of iid random variables is approx. normal

$$X_i = \frac{H_i}{N_i} = \frac{\text{Binomial}(N_i, P_i)}{N_i} = \frac{1}{N_i} \sum_{j=1}^{N_i} \begin{cases} 1 & \text{if hit} \\ 0 & \text{if not hit} \end{cases}$$

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

Math with Normal is easy!

$$\mathbb{E} X_i = \frac{1}{N_i} \mathbb{E} H_i = \frac{1}{N_i} \cdot \underbrace{N_i P_i}_{\mathbb{E} \text{Binomial}(N_i, P_i)} = P_i$$

$$\begin{aligned} \text{VAR}(X_i) &= \text{VAR}\left(\frac{H_i}{N_i}\right) = \frac{1}{N_i^2} \text{VAR}(H_i) \\ &= \frac{1}{N_i^2} \text{VAR}(\text{Binom}(N_i, P_i)) = \frac{N_i P_i (1-P_i)}{N_i^2} \end{aligned}$$

$$\text{VAR}(c X_i) = c^2 \text{VAR}(X_i) =$$

$$\begin{aligned} \text{VAR}(cX_i) &= \mathbb{E}(cX_i)^2 - [\mathbb{E}(cX_i)]^2 = \\ &= c^2 (\mathbb{E}X_i^2 - (\mathbb{E}X_i)^2) = c^2 \text{Var}(X_i) \end{aligned}$$

i^{th} batting avg. $X_i \approx \mathcal{N}(p_i, \frac{p_i(1-p_i)}{N_i})$

$$\rightarrow \begin{cases} X_i \sim \mathcal{N}(p_i, \sigma_i^2) \\ p_i \sim \mathcal{N}(p, \tau^2) \end{cases}$$

→ PRIOR:
batter i is
a baseball
player whose
quality is
drawn from
some dist.
across baseball
players

$$\hat{p}_i^{\text{(MLE)}} = X_i$$

Bayesian: MAP, posterior mean

posterior → compute
estimate $\hat{p}_i = \mathbb{E}(p_i | X_i)$

Posterior:

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

proportional to

$$\propto P(X_i | P_i) P(P_i)$$

$$= P(\mathcal{N}(P_i, \sigma_i^2) = X_i) \cdot P(\mathcal{N}(P, \tau^2) = P_i)$$

$$= \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2\sigma_i^2} (X_i - P_i)^2\right)$$

$$\frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{1}{2\tau^2} (P - P_i)^2\right)$$

$$\propto \exp\left(-\frac{1}{2} \frac{1}{\sigma_i^2} (X_i - P_i)^2 - \frac{1}{2} \frac{1}{\tau^2} (P - P_i)^2\right)$$

$$\exp(a) \cdot \exp(b) = e^a \cdot e^b = e^{a+b} = \exp(a+b)$$

$$= \exp \left(-\frac{1}{2} \left[\frac{x_i^2 - 2x_i p_i + p_i^2}{\sigma_i^2} + \frac{p^2 - 2p p_i + p_i^2}{\tau^2} \right] \right)$$

$$\alpha \exp \left(-\frac{1}{2} \left[p_i^2 \left(\frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right) - 2p_i \left(\frac{x_i}{\sigma_i^2} + \frac{p}{\tau^2} \right) + \left(\frac{p^2}{\sigma_i^2} + \frac{p^2}{\tau^2} \right) \right] \right)$$

$$= \exp \left(-\frac{1}{2} \left(\frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right) \left[p_i^2 - 2p_i \left(\frac{x_i}{\sigma_i^2} + \frac{p}{\tau^2} \right) \right] \right)$$

$$\alpha \exp \left(-\frac{1}{2} \left(\frac{1}{\sigma_i^2} + \frac{1}{\tau^2} \right) \cdot \left(p_i - \left(\frac{x_i}{\sigma_i^2} + \frac{p}{\tau^2} \right) \right)^2 \right)$$

$$(p_i - a)^2 = p_i^2 - 2p_i a + \cancel{a^2}$$

$$P(P_i | X_i)$$

$$\propto \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}\right) \cdot \left(P_i - \left(\frac{\frac{X_i}{\sigma_i^2} + \frac{\mu}{\tau^2}}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}}\right)^2\right)\right)$$

$$P(P_i | N(\mu, \nu)) = \frac{1}{\sqrt{2\pi\nu^2}} \exp\left(-\frac{1}{2} \frac{1}{\nu} (P_i - \mu)^2\right)$$

$$\Rightarrow P_i | X_i \sim N\left(\left(\frac{\frac{X_i}{\sigma_i^2} + \frac{\mu}{\tau^2}}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}}\right), \frac{1}{\left(\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}\right)}\right)$$

We now have the post. dist $P_i | X_i$

Prior $P_i \sim N(\mu, \tau^2)$

Posterior: updated our beliefs about P_i after having seen the data X_i .

Posterior mean:

X_i = player's obs. batting avg
 μ = overall mean player batting
 τ^2 = variability across batters
(variance of μ_i overall batters)
 σ_i^2 = variability within a batter

$$\hat{p}_i^{(MLE)} = X_i$$

$$\hat{p}_i^{(Bayes)} = \mathbb{E}(p_i | X_i) = \left(\frac{X_i \frac{1}{\sigma_i^2} + \frac{\mu}{\tau^2}}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}} \right)$$
$$= \mu + \frac{\tau^2}{\tau^2 + \sigma_i^2} (X_i - \mu)$$

$$\left(\frac{X_i \frac{1}{\sigma_i^2} + \frac{\mu}{\tau^2}}{\frac{1}{\sigma_i^2} + \frac{1}{\tau^2}} \right) \frac{\sigma_i^2 \tau^2}{\sigma_i^2 \tau^2} = \frac{X_i \tau^2 + \mu \sigma_i^2}{\sigma_i^2 + \tau^2}$$

$$= \mu \left(\frac{\sigma_i^2 + \tau^2 - \tau^2}{\sigma_i^2 + \tau^2} \right) + X_i \left(\frac{\tau^2}{\sigma_i^2 + \tau^2} \right)$$

$$= \mu + \mu \left(\frac{-\tau^2}{\sigma_i^2 + \tau^2} \right) + X_i \left(\frac{\tau^2}{\sigma_i^2 + \tau^2} \right)$$

If $\frac{\tau^2}{\tau^2 + \sigma_i^2} = 1$, then $\hat{p}_i^{\text{Bayes}} = \hat{p}_i^{\text{MLE}} = X_i$
 but its < 1 , so \hat{p}_i^{Bayes} is closer to p than \hat{p}_i^{MLE}

$$\begin{cases} X_i \sim \mathcal{N}(p_i, \sigma_i^2) \\ p_i \sim \mathcal{N}(p, \tau^2) \end{cases}$$

Model

$$\hat{p}_i^{\text{(Bayes)}} = \frac{\frac{X_i}{\hat{\sigma}_i^2} + \frac{\hat{p}}{\hat{\tau}^2}}{\frac{1}{\hat{\sigma}_i^2} + \frac{1}{\hat{\tau}^2}}$$

observed X_i
 estimate p_i

Problem: we never observed p, τ^2, σ_i^2

Empirical Bayes: estimate these other hyperparameters $\hat{p}, \hat{\tau}^2, \hat{\sigma}_i^2$ and plug them in!

$$\left\{ \begin{aligned} \sigma_i^2 &= \frac{P_i(1-P_i)}{N_i} = \frac{1}{N_i} \text{var}(\text{Binom}(N_i, P_i)) \end{aligned} \right.$$

Simplification: $\sigma_i^2 = \frac{C}{N_i}$ C : some constant

$$X_i = \frac{H_i}{N_i} \quad P_i(1-P_i) = \frac{1}{4} \cdot \frac{3}{4} = \frac{3}{16}$$

$$\hat{P}_i^{\text{Bayes}} = \frac{\frac{H_i/N_i}{C/N_i} + \frac{\hat{P}}{\tau^2}}{\frac{1}{C/N_i} + \frac{1}{\tau^2}}$$

$$\hat{P}_i^{(\text{MLE})} = \frac{H_i}{N_i}$$

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

$$\tau^2 = 0 \rightarrow \hat{P}$$

$$\frac{\hat{\sigma}^2}{\sigma^2} = \infty \rightarrow \frac{H_i}{N_i}$$

$$\widehat{WP}^{MLE} = \frac{W}{N}, \quad \widehat{WP}^{Bayes} = \frac{W + W'}{N + N'}$$

$$\text{as } N_i \uparrow, \quad \hat{p}_i^{(Bayes)} \rightarrow \hat{p}_i^{(MLE)}$$

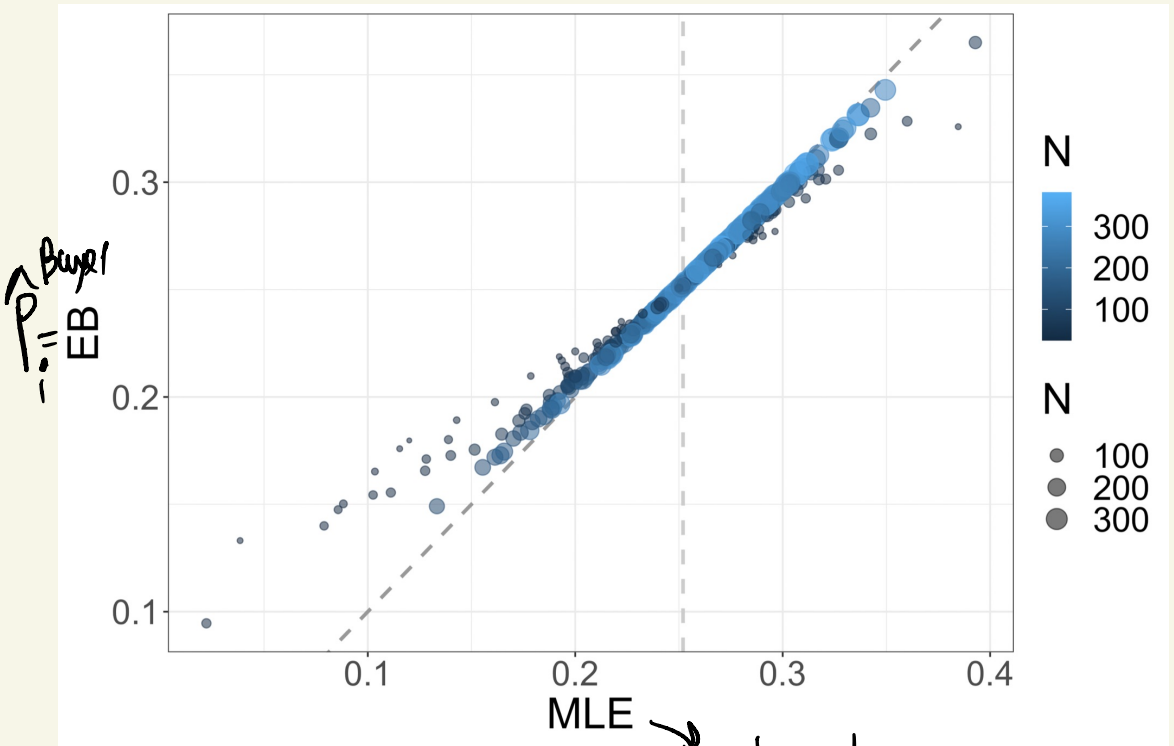
$$\hat{p}_i^{(Bayes)} = \frac{H_i + \hat{p} \frac{\hat{C}}{\hat{\sigma}^2}}{N_i + 1 \cdot \frac{\hat{C}}{\hat{\sigma}^2}}$$

- estimate $\hat{p}, \hat{\sigma}^2$ from data $\{X_i\}$ in a smart way

$$\text{e.g. } \hat{p} = \text{mean}(X_i)$$

- in practice, there is no C b.c. that was a simplification

- treated c as a tuning parameter
 → chose the \hat{c} which had best pred. perf. (lowest RMSE of \hat{p}_i to p_i)



obs. batch, average

rmse_MLE	rmse_EB
0.02629828	0.02383808

Taxonomy

- Shrinkage — Shrank obs. MLE bathy avg. to overall mean, shrinking more if have fewer observations, to stabilize estimates and improve predictive performance
- Bayesian modeling — used prior to encode additional information

Consultant's Dilemma (Stein's Paradox)

James Stein's Theorem

Client index i , (unobserved) μ_i each client,
data $X_i \stackrel{iid}{\sim} N(\mu_i, 1)$ (known variance)

the JS estimator $\hat{\mu}_i^{(JS)}$
dominates the MLE,

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

Leaderboard:

predict $\hat{\mu}_i$

observe μ_i

$$\text{score} = \sum_{i=1}^n (\hat{\mu}_i - \mu_i)^2$$

$$\hat{\mu}_i^{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$