

Brown 2008: Prediction of Batting Averages

Tues go through 1st 7 pages of those notes (the context)

Original data $\{H_i, N_i\}$ # hits for player i
(1st half of season) # at bats for player i
 p players

Original model $H_i \sim \text{Binomial}(N_i, P_i)$

unknown success/hit rate parameter P_i

$$\frac{H_i}{N_i} = \frac{\text{Bin}(N_i, P_i)}{N_i} \xrightarrow{d} \mathcal{N}\left(P_i, \frac{P_i(1-P_i)}{N_i}\right)$$

with this didn't depend on the unknown P_i

CLT $\frac{1}{\sqrt{n}} \sum_{k=1}^n X_k \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ if $\{X_k\}$ iid mean 0
variance $\sigma^2 < \infty$.

$$\frac{1}{\sqrt{N_i}} \left(\sum_{k=1}^{N_i} \text{Bernoulli}(P_i) - N_i P_i \right) \xrightarrow{d} \mathcal{N}(0, P_i(1-P_i))$$

$$\sqrt{N_i} \left(\frac{\text{Bin}(N_i, P_i)}{N_i} - P_i \right) \xrightarrow{d} \mathcal{N}(0, P_i(1-P_i))$$

$$\frac{H_i}{N_i} \xrightarrow{d} \mathcal{N}\left(P_i, \frac{P_i(1-P_i)}{N_i}\right)$$

2. Variance stabilizing Transformation

$$X_i = \arcsin \sqrt{\frac{H_i + 1/4}{N_i + 1/2}}$$

transformed data $\{X_i\}$

Where does the form $\arcsin \sqrt{\frac{H}{N}}$ come from?

$H \sim \text{Binomial}(N, p)$.

$$\mu = E\left(\frac{H}{n}\right) = p, \quad \text{var}\left(\frac{H}{n}\right) = \frac{p(1-p)}{n} =: V(p).$$

Find transformation $T(\cdot)$ so that $\text{var}(T(\frac{H}{n}))$ is constant \rightarrow not depending on p .

1st order Taylor APPROX. of T about mean p at point $\frac{H}{N}$:

$$T\left(\frac{H}{N}\right) \approx T(p) + T'(p)\left(\frac{H}{N} - p\right)$$

$$\begin{aligned} \text{var}\left(T\left(\frac{H}{N}\right)\right) &= (T'(p))^2 \text{var}\left(\frac{H}{N}\right) \\ &= (T'(p))^2 V(p) \stackrel{\text{want}}{=} C \end{aligned}$$

$$\Rightarrow \text{solve } T'(p) = \sqrt{\frac{C}{V(p)}}$$

$$\begin{aligned} \Rightarrow T(p) &= \int_0^p \frac{\sqrt{C}}{\sqrt{V(t)}} dt = \sqrt{C} \int_0^p \frac{dt}{\sqrt{t(1-t)}} \\ &= \sqrt{C} \int_0^{\arcsin p} \frac{2 \sin \theta \cos \theta d\theta}{\sqrt{\sin^2 \theta (1 - \cos^2 \theta)}} \quad \text{letting } t = \sin \theta \end{aligned}$$

$$= \underbrace{2 \sqrt{nc}}_{\text{constant}} \arcsin \sqrt{p}$$

Why $X_i = \arcsin \sqrt{\frac{H_i + 1/4}{N_i + 1/2}}$? (Significance of the constants $1/4, 1/2$)

- Constant variance $\text{var}(X) = \frac{1}{4N} + \mathcal{O}(\frac{1}{N^2})$
 - Nice Expectation $EX = \arcsin \sqrt{p} + \mathcal{O}(\frac{1}{N^2})$
 - $X_i \xrightarrow{d} \mathcal{N}(\theta_i, \delta_i^2)$
- + term \nearrow

will be estimating this parameter

$\theta_i = \arcsin \sqrt{p}$	unknown
$\delta_i^2 = 1/4N_i$	<u>known</u>

Why Lots of Algebra & Taylor Series manipulations.
I have a (poorly written) reference from Brown on this.

③ Measuring the Efficacy of an Estimator

Transformed Data $X_i \sim \mathcal{N}(\theta_i, \sigma_i^2)$

Goal Estimate $\{\theta_i\}$ from the data $\{X_i\}$,
which is equivalent to estimating the i^{th}
player's Batting Average

Estimator of θ_i is a function of the data

Notation: $\hat{\theta}_i = \delta(X_i)$

Total Squared Error of the Estimator δ

$$\text{TSE}(\delta) := \sum_i (\delta_i - \theta_i)^2 \quad \text{if we know the } \theta_i$$

$$= \sum_i (\delta_i - X_i)^2 + \sum_i (X_i - \theta_i)^2 - \sum_i 2(\delta_i - X_i)(X_i - \theta_i)$$

Estimate of TSE

$$\widehat{\text{TSE}}(\delta) := \sum_i (\delta_i - X_i)^2 + \mathbb{E} \left[\sum_i (X_i - \theta_i)^2 - \sum_i 2(\delta_i - X_i)(X_i - \theta_i) \right]$$

$$\widehat{\text{TSE}}(\delta) = \sum_i (\delta_i - X_i)^2 - \sum_i \frac{1}{4N_i}$$

\mathbb{E} over $X \sim \mathcal{N}(\theta_i, \sigma_i^2)$
 $\sigma_i^2 = \frac{1}{4N_i}$ known

$\left\{ \begin{array}{l} X_i \text{ known} \\ \delta_i = \delta_i(X_i) \\ N_i \text{ known} \end{array} \right\}$ known

so, we can calculate $\widehat{TSE}(\delta)$ to evaluate our estimator $\delta_i(X_i)$

Naive Estimator

Normalization

of θ_i

$$\delta_n(x) = x$$

Identity: guess future bathing coverage will be the exact same as the previous BA

want to compare

estimator δ

to δ_0 .

$$\widehat{TSE}^*(\delta) := \frac{\widehat{TSE}(\delta)}{\widehat{TSE}(\delta_0)}$$

$$\widehat{TSE}^*(\delta) < 1 \Rightarrow \delta \text{ better than } \delta_0$$

4. Choosing our Estimators $\hat{\delta}(x)$ for Θ

- Naive Estimator
- Overall Mean
- Parametric Empirical Bayes
 - Method of Moments
 - MLE
- Nonparametric Empirical Bayes
- Harmonic Bayes Estimator
- James Stein Estimator

5 nontrivial
estimators

Next Week: Lets' try to understand these estimators.

- Results
Page 133 of Brown 2008, Table 2, TSE^* column.

Naive estimator did the worst!

EBCMM, NP EB, JS did the best! $TSE^* \approx 0.5$

- Explanation

1. $\{X_i\}$ not normal (moderately)
2. Sample values for $\{N_i\}$ and $\{X_i\}$ are moderately correlated. (When separate pitchers & nonpitchers, this correlation disappears, and the other estimators do way better!)