

- Previously, we covered ridge regression, which is a popular shrinkage estimator.
- OLS / MLE
Ridge.
-
- \rightarrow Ridge is shrinking the coefficients towards a common value (zero).

- Ridge is not the only shrinkage estimator. Today we'll cover the first shrinkage estimator probably that came to be:

The James - Stein estimator

I once met Stein! G!

- Suppose we wish to estimate a single parameter μ from observation x in Bayesian context:
- $\mu \sim N(M, A)$
- Normal Prior.
- $x | \mu \sim N(\mu, 1)$
- Normal $(\mu, 1)$ likelihood
- In this case $\mu | x$ has posterior distribution

$$\mu | x \sim N(M + B(x - M), B)$$

where

$$B = \frac{A}{A+1}$$

Example:

→ If this was MLE and $x = 10$

$$\Rightarrow \hat{\mu}^{\text{MLE}} = \frac{10}{1} = 10$$

• However, in this context, say we had $M = 9$,

then

$$\hat{\mu}^{\text{Bayes}} = M + \left(\frac{A}{A+1} \right) \cdot (x - M)$$

$$\text{Always } = 9 + \left(\frac{A}{A+1} \right) (10 - 9)$$

the same.

$$0 \leq \frac{A}{A+1} \leq 1$$

• Let's now compare the overall expected error of either estimator:

$$E \{ (\hat{\mu}^{\text{MLE}} - \mu)^2 \} = 1$$

$$\text{If we had } x_1 = 10 \rightarrow \hat{\mu}^{\text{MLE}} = 10$$

$$x_2 = 10.5 \rightarrow \hat{\mu}^{\text{MLE}} = 10.5$$

⋮

$$\text{Var} = 1$$

$$E \left\{ (\hat{\mu}^{\text{Bayes}} - \mu)^2 \right\} = B^2$$

$$\begin{aligned} V(\hat{\mu}^{\text{Bayes}}) &= V(M + B(\bar{x} - M)) \\ &= V(B\bar{x} - BM) \\ &= B^2 V(\bar{x}) = B^2 \end{aligned}$$

This variance is smaller than that of MLE.

Now, if we have N independent observations,

say ...

$$\boldsymbol{\mu} = [\mu_1, \mu_2, \dots, \mu_N]^T$$

and

$$\underline{x} = [x_1, x_2, \dots, x_N]^T$$

The Data Generation Process (DGP)
of each x_i is not the same

Different μ_i 's for each observation.

$$\mu_i \sim N(M, A)$$

and $x_i | \mu_i \sim N(\mu_i, I)$

We're assuming each μ_i has the same prior.

x_i 's have the same variance around their different means.

Then,

$$\hat{\boldsymbol{\mu}}^{\text{Bayes}} = \underline{M} + B(\underline{x} - \underline{M})$$

$$\begin{matrix} \text{N terms} \\ \downarrow \end{matrix} \quad \left[\begin{matrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_N \end{matrix} \right] \quad \left[\begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{matrix} \right]$$

And $\hat{\mu}^{\text{MLE}} = \bar{x}$

Thus, $E\left\{\|\hat{\mu}^{\text{MLE}} - \mu\|^2\right\} = E\left\{\sum_{i=1}^N (\hat{\mu}^{\text{MLE}} - \mu)^2\right\} = N$

$$E\left\{\|\hat{\mu}^{\text{Bayes}} - \mu\|^2\right\} = N B^2$$

Again, the overall error of $\hat{\mu}^{\text{Bayes}}$ is smaller vs. $\hat{\mu}^{\text{MLE}}$.

Recall

$$V(\hat{\theta}) = E[(\hat{\theta} - E(\hat{\theta}))^2]$$

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

$$\text{Bias}^2 = [E(\hat{\theta}) - \theta]^2$$

If true θ was 10.1, then

$$\text{Bias}(\hat{\theta}) = -0.1$$

$$\text{Bias}^2(\hat{\theta}) = 0.01$$

$$E(\hat{\theta}) = 10$$

Sample 1
 $\hat{\theta}_1 = 10.5$

Sample 2
 $\hat{\theta}_2 = 10$
 \vdots

Sample n
 $\hat{\theta}_n = \overline{a.s}$

$$V(\hat{\theta}) = \frac{(10.5 - 10)^2 + (10 - 10)^2 + \dots + (9.5 - 10)^2}{n-1}$$

$V(\hat{\theta}) + \text{Bias}^2(\hat{\theta})$ to check how estimator performs

Overall Error: $\frac{(\hat{\theta}_1 - \text{True } \theta)^2 + (\hat{\theta}_2 - \text{True } \theta)^2 + \dots + (\hat{\theta}_n - \text{True } \theta)^2}{n}$

→ If we know m and A , then we're in great shape with Bayes. However, we usually don't have prior information. what to do then?

- Estimate m !

$$\hat{m} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\hat{\beta} = 1 - \frac{(N-3)}{\sum_{i=1}^{n-3} (x_i - \bar{x})^2}$$

Then, $\hat{\mu}_i^{\text{JS}} = \hat{m} + \hat{\beta}(x_i - \hat{m})$ for $i = 1, \dots, n$.

James - Stein Estimator

It can be shown that

$$E[\|\hat{\mu}^{\text{JS}} - \mu\|^2] = N\beta^2 + 3(1-\beta)$$

James - Stein Theorem

Assumption

normally. Different means
Same variance

Suppose that $x_i | \mu \sim N(\mu, 1)$
 independently for $i = 1, 2, \dots, N$ with $N \geq 4$.
 (Regardless of the prior).

Statement

$$\text{Then, } E\left\{\|\hat{\mu}^{\text{JS}} - \mu\|^2\right\} \leq N = E\left\{\|\hat{\mu}^{\text{MLE}} - \mu\|^2\right\}$$

for all choices of $\mu \in \mathbb{R}^N$

Overall error of JS estimator

Overall error of MLE.

Baseball Players

- Same season of MLB
- Each of them batted a # of hits over a number of at-bats.
- If we divide $\frac{\# \text{ of hits}}{\# \text{ of at-bats}} = \text{batting avg.}$

→ Let's get the JS for these batting avg's and compare it to the MLE.

Example: player 1: $\frac{31}{90} = \underline{\underline{.345}} = \hat{\mu}_1^{\text{MLE}}$

$$\text{player 18: } \frac{13}{90} = .145 = \hat{\mu}_{18}^{\text{MLE}}$$

Long-run
behavior of the player
avg. of the player

probabilities,
 $\underline{p_i}$'s

$$p_i \sim B_i(\underline{q_0}, p_i)$$

These are
like the x_i 's

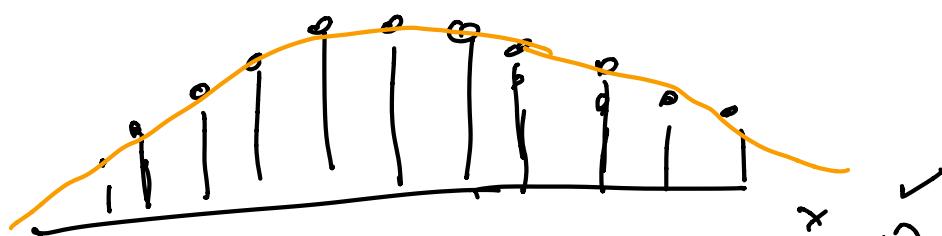
$$V(p_i) = \frac{1}{n^2} np(1-p) = \frac{p(1-p)}{n}$$

These are like the p_i 's

Issue: this is Binomial! not normal!!.

Good News is that we can approximate
the Binomial w/ a Normal!

$$\hat{p}_i \approx N(p_i, \sigma_0^2)$$



$$\sigma_0^2 = \frac{\bar{p}(1-\bar{p})}{90}$$

$$z = \frac{\hat{p}_i - p_i}{\sigma_0}, z \sim N(0, 1)$$

$$x \sim N(p_i, 1)$$

Now, let $x_i = \frac{\hat{p}_i}{\sigma_0}$ why? $\Rightarrow x_i \sim N(p_i, 1)$

and apply JS estimator

$$\hat{\mu}^{ss} = \hat{m} + \hat{B}(x; -\hat{m})$$