

## Lecture 13: Shrinkage Estimation

*Instructor: Jonathan Pipping**Author: Ryan Brill*

## 13.1 Problem Setup

### 13.1.1 Revisiting the Problem

Recall our question from the previous lecture:

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

We reduced this problem to estimating the parameters of the following model:

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

where  $X_i$  is the observed mid-season BA for player  $i$ ,  $\mu_i$  is the latent "true quality" of player  $i$ ,  $\sigma_i^2$  is a known variance which depends on  $N_i$ , the number of at-bats for player  $i$ , and  $\mu$  and  $\tau^2$  are unknown hyperparameters representing the overall mean and variance of the population of players.

We solved this previously with Empirical Bayes, and will now consider another perspective to understand why this works.

### 13.1.2 Revising the Model

Since  $\sigma_i^2$  is known, we can divide by  $\sigma_i^2$  to remove some parameters:

- $X_i \leftarrow X_i / \sigma_i^2$
- $\theta_i \leftarrow \mu_i / \sigma_i^2$

Now suppose we are taking a frequentist approach, that is, we will think of each parameter as unknown fixed constants rather than random variables (i.e. the unknown "true" quality of each player). Then we are left with the model:

$$X_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\theta_i, 1), \quad i = 1, \dots, k$$

and our task is to estimate the fixed unknown constants  $\theta_i$  (the normal means) given the data  $\{X_i\}_{i=1}^k$  to optimize the composite loss function

$$L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) = \sum_{i=1}^k \left( \theta_i - \hat{\theta}_i \right)^2$$

where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$  are the true parameters to be estimated, and  $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$  are our estimates.

The performance of the joint estimator  $\hat{\theta}$  is judged by the risk function, or the expected loss:

$$R(\theta, \hat{\theta}) = \mathbb{E}L(\theta, \hat{\theta})$$

This simple setup leads to one of the most significant results in mathematical statistics: Stein's Paradox and Shrinkage Estimation.

## 13.2 Shrinkage Estimation

### 13.2.1 Starting Point: The MLE

The estimation problem above involves pairs of values  $X_i, \theta_i$ ,  $i = 1, \dots, k$  where one element of each pair is known ( $X_i$ ) and one ( $\theta_i$ ) is unknown. The "obvious" or ordinary estimator is just  $\hat{\theta}_i = X_i$ , which is the **Maximum Likelihood Estimator (MLE)** we're already familiar with! This estimator maximizes the probability of observing the data we did.

$$\begin{aligned}\hat{\theta}^{(MLE)} &= \arg \max_{\theta} \mathbb{P}(\text{data} \mid \theta) \\ &= \arg \max_{\theta} \mathbb{P}(X_1, \dots, X_k \mid \theta_1, \dots, \theta_k)\end{aligned}$$

By assumed independence between the  $X_i$  and the monotonicity of the log function, we have

$$\begin{aligned}\hat{\theta}^{(MLE)} &= \arg \max_{\theta} \prod_{i=1}^k \mathbb{P}(X_i \mid \theta_i) \\ &= \arg \max_{\theta} \sum_{i=1}^k \log \mathbb{P}(X_i \mid \theta_i)\end{aligned}$$

Then from our model, we have

$$\begin{aligned}\hat{\theta}^{(MLE)} &= \arg \max_{\theta} \sum_{i=1}^k \log \mathcal{N}(X_i; \theta_i, 1) \\ &= \arg \max_{\theta} \sum_{i=1}^k \log \left[ \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{1}{2}(X_i - \theta_i)^2 \right) \right] \\ &= \arg \max_{\theta} \sum_{i=1}^k -\frac{1}{2}(X_i - \theta_i)^2 \\ &= \mathbf{X}\end{aligned}$$

where  $\mathbf{X} = (X_1, \dots, X_k)$ .

In terms of our baseball example, this means just predict using each player's mid-season batting average. However, as we saw in the last lecture, these predictions are **terrible**. But here  $\hat{\theta}$  are unknown fixed constants and a prior is mis-specified. Why does it work?

### 13.2.2 Visualizing the Problem

Recall that this estimation problem involves pairs of values  $X_i, \theta_i$ ,  $i = 1, \dots, k$  where one element of each pair is known ( $X_i$ ) and one ( $\theta_i$ ) is unknown. Since the  $\theta_i$ 's are unknown, we cannot plot the pairs, but we

can imagine what such a plot would look like to understand the problem. Such a plot is included in Figure 13.1.

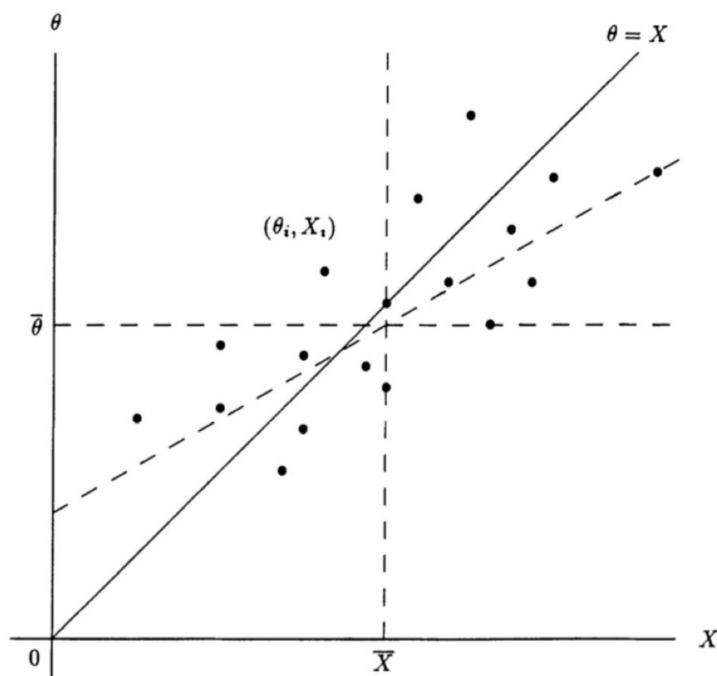


Figure 13.1: A hypothetical bivariate plot of the pairs  $X_i, \theta_i$ ,  $i = 1, \dots, k$ .

Since  $\mathbf{X}$  is  $\mathcal{N}(\boldsymbol{\theta}, 1)$ , we can think of the  $X_i$  as being generated by  $\mathcal{N}(0, 1)$  "errors" to the given  $\theta_i$ 's. So the horizontal deviates of the  $X_i$ 's from the 45° line  $\theta = X$  are independent  $\mathcal{N}(0, 1)$  random variables. Our goal is to estimate all the  $\theta_i$ 's given all of the  $X_i$ 's with **no assumptions about a possible distributional structure for the  $\theta_i$ 's**: they are simply to be viewed as unknown constants.

### 13.2.3 The Galtonian Perspective

**Q:** Why should we expect that the ordinary estimator  $\hat{\boldsymbol{\theta}} = \mathbf{X}$  can be improved upon?

Well, if the  $\theta_i$ 's and hence the pair  $(X_i, \theta_i)$  had a known joint distribution, a natural method of proceeding is  $\hat{\boldsymbol{\theta}} = \mathbb{E}[\boldsymbol{\theta} | \mathbf{X}]$  and use this, the theoretical regression function of  $\boldsymbol{\theta}$  on  $\mathbf{X}$ , to generate estimates of the  $\theta_i$ 's by evaluating it for each  $X_i$ . This is an unattainable ideal though, because we do **not** know the conditional distribution of  $\boldsymbol{\theta}$  given  $\mathbf{X}$ .

Moreover, we don't assume that our uncertainty about the unknown constants  $\theta_i$  can be described by a probability distribution at all; we are relying purely on frequentist principles. We **do** however know the conditional distribution of  $\mathbf{X}$  given  $\boldsymbol{\theta}$ ,  $\mathcal{N}(\boldsymbol{\theta}, 1)$ , and we can calculate  $\mathbb{E}[\mathbf{X} | \boldsymbol{\theta}]$ . Indeed, this theoretical regression line corresponds to the 45° line  $\theta = X$ , and this line yields the ordinary estimators  $\hat{\boldsymbol{\theta}}^{(MLE)} = \mathbf{X}$ .

So the ordinary estimator may be viewed as being based on the "wrong" regression line, on  $\mathbb{E}[\mathbf{X} | \boldsymbol{\theta}]$  instead of  $\mathbb{E}[\boldsymbol{\theta} | \mathbf{X}]$ .

As Francis Galton knew in the 1880s, the regressions of  $\mathbf{X}$  on  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}$  on  $\mathbf{X}$  can be markedly different, as

you'll show in the [Homework](#) . He suggested that this ordinary estimator could be improved upon, and even suggested a method for doing so: by attempting to approximate  $\mathbb{E}[\boldsymbol{\theta} \mid \mathbf{X}]$  in a setting where the  $\theta_i$ 's do not have a distribution.

With no distributional assumptions on the  $\theta_i$ 's, we are of course prevented from looking at an optimal estimator of  $\mathbb{E}[\boldsymbol{\theta} \mid \mathbf{X}]$ . Instead we note that  $\hat{\boldsymbol{\theta}}^{(MLE)} = \mathbf{X}$  is a linear function of  $\mathbf{X}$ , and we can look for a best linear estimator of the form  $\hat{\boldsymbol{\theta}} = \mathbf{a} + b\mathbf{X}$  so as to minimize the composite loss function

$$\begin{aligned} L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) &= \sum_{i=1}^k (\theta_i - \hat{\theta}_i)^2 \\ &= \sum_{i=1}^k (\theta_i - (a + bX_i))^2 \end{aligned}$$

If the  $\theta_i$ 's are known, we would have a standard [simple linear regression problem](#) with the best linear estimator given by the regression line

$$\hat{\boldsymbol{\theta}} = \bar{\boldsymbol{\theta}} + \hat{\beta}(\mathbf{X} - \bar{\mathbf{X}})$$

where

$$\hat{\beta} = \frac{\sum_{i=1}^k (X_i - \bar{X})(\theta_i - \bar{\theta})}{\sum_{i=1}^k (X_i - \bar{X})^2}$$

These  $\theta_i$ 's are unknown, but if we could estimate the functions of these unknown parameters ( $\bar{\theta}$  and  $\hat{\beta}$ ), we could estimate the regression line of  $\boldsymbol{\theta}$  on  $\mathbf{X}$ . What can we use to estimate these functions?

### 13.2.4 Deriving the Best Linear Shrinkage Estimator

First, we can use the sample mean  $\bar{X}$  to estimate  $\bar{\theta}$ . Then for the slope  $\hat{\beta}$ :

$$\begin{aligned} \hat{\beta} &= \frac{\sum_{i=1}^k (X_i - \bar{X})(\theta_i - \bar{\theta})}{\sum_{i=1}^k (X_i - \bar{X})^2} \\ &= \frac{S_{X\theta}}{S_X^2} \end{aligned}$$

where

$$S_X^2 = \sum_{i=1}^k (X_i - \bar{X})^2$$

is the sample variance of the  $X_i$ 's, an unbiased estimator of  $\text{Var}(\mathbf{X})$  and

$$S_{X\theta} = \sum_{i=1}^k (X_i - \bar{X})(\theta_i - \bar{\theta})$$

is the sample covariance of  $\mathbf{X}$  and  $\boldsymbol{\theta}$ , an unbiased estimate of  $\text{Cov}(\mathbf{X}, \boldsymbol{\theta})$ . Of these two,  $S_{X\theta}$  is unknown, so we'll need to estimate it from the data. Since  $\mathbf{X} = \boldsymbol{\theta} + \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, 1)$  and  $\text{Var}(\mathbf{X}) = \text{Var}(\boldsymbol{\theta}) + \text{Var}(\boldsymbol{\epsilon}) = \text{Var}(\boldsymbol{\theta}) + 1$ , we have that

$$\begin{aligned} \text{Cov}(\mathbf{X}, \boldsymbol{\theta}) &= \text{Cov}(\boldsymbol{\theta} + \boldsymbol{\epsilon}, \boldsymbol{\theta}) \\ &= \text{Var}(\boldsymbol{\theta}) + \text{Cov}(\boldsymbol{\epsilon}, \boldsymbol{\theta}) \end{aligned}$$

And since  $\epsilon$  is independent of  $\theta$ ,  $\text{Cov}(\epsilon, \theta) = 0$  and we have that

$$\text{Cov}(\mathbf{X}, \theta) = \text{Var}(\theta) = \text{Var}(\mathbf{X}) - 1$$

so  $S_X^2 - 1$  and  $S_{X\theta}$  have the same expectation, and

$$\begin{aligned}\hat{\beta} &= \frac{S_{X\theta}}{S_X^2} \\ &= \frac{\sum_{i=1}^k (X_i - \bar{X})(\theta_i - \bar{\theta})}{\sum_{i=1}^k (X_i - \bar{X})^2} \\ &\approx \frac{\sum_{i=1}^k (X_i - \bar{X})^2 - (k-1)}{\sum_{i=1}^k (X_i - \bar{X})^2} \\ &= 1 - \frac{k-1}{\sum_{i=1}^k (X_i - \bar{X})^2}\end{aligned}$$

This leads to the **Efron-Morris estimated least-squares line**:

$$\hat{\theta}^{(EM)} = \bar{X} + \left(1 - \frac{k-1}{S_X^2}\right) (\mathbf{X} - \bar{X}), \text{ where } S_X^2 = \sum_{i=1}^k (X_i - \bar{X})^2$$

**James-Stein's original shrinkage estimator** can be derived by considering the class of estimators that are linear in  $\mathbf{X}$  with 0 intercept,

$$\hat{\theta} = b\mathbf{X}$$

where  $b$  is a constant. The least squares estimator has

$$\hat{\beta} = \frac{\sum_{i=1}^k \theta_i X_i}{\sum_{i=1}^k X_i^2}$$

and  $\theta_i X_i$  has the same expectation as  $\sum_{i=1}^k X_i^2 - k$ , which yields the **James-Stein shrinkage estimator**:

$$\hat{\theta}^{(JS)} = \left(1 - \frac{k}{S_X^2}\right) \mathbf{X}, \text{ where } S_X^2 = \sum_{i=1}^k (X_i - \bar{X})^2$$

### 13.2.5 Comparing the Estimators

So the ordinary estimators  $\hat{\theta}^{(MLE)} = \mathbf{X}$  are derived from the theoretical regression line of  $\mathbf{X}$  on  $\theta$ , which is useful if our goal is to predict  $\mathbf{X}$  from  $\theta$ . But our goal is the opposite, to predict  $\theta$  from  $\mathbf{X}$  with the sum of squares criterion:

$$\sum_{i=1}^k (\theta_i - \hat{\theta}_i)^2$$

So the optimal estimator is the least squares regression line of  $\theta$  on  $\mathbf{X}$ , and the James-Stein and Efron-Morris estimators are themselves approximations of this correct regression line. It turns out that both of these estimates are better (have less risk) than the ordinary estimate (the MLE), but we need to **show this**.

Let  $\hat{\boldsymbol{\theta}}^{(b)} = b\mathbf{X}$  represent the class of linear estimators with zero intercept. The Risk is:

$$\begin{aligned} R(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(b)}) &= \mathbb{E} L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(b)}) \\ &= \mathbb{E} \sum_{i=1}^k \left( \theta_i - \hat{\theta}_i \right)^2 \\ &= \mathbb{E} \sum_{i=1}^k \left( \theta_i - \hat{\theta}_i^{(LS)} + \hat{\theta}_i^{(LS)} - \hat{\theta}_i^{(b)} \right)^2 \end{aligned}$$

where  $\hat{\theta}_i^{(LS)} = \hat{\beta}X_i$  and  $\hat{\beta} = \frac{\sum_{i=1}^k \theta_i X_i}{\sum_{i=1}^k X_i^2}$ . Then

$$R(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(b)}) = \mathbb{E} \left[ \sum_{i=1}^k \left( \theta_i - \hat{\theta}_i^{(LS)} \right)^2 + 2 \sum_{i=1}^k \left( \theta_i - \hat{\theta}_i^{(LS)} \right) \left( \hat{\theta}_i^{(LS)} - \hat{\theta}_i^{(b)} \right) + \sum_{i=1}^k \left( \hat{\theta}_i^{(LS)} - \hat{\theta}_i^{(b)} \right)^2 \right]$$

Since the residuals  $\theta_i - \hat{\theta}_i^{(LS)}$  are orthogonal to any linear function of  $X_i$ , the second term equals 0 by the orthogonality principle, and

$$\begin{aligned} R(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(b)}) &= \mathbb{E} \left[ \sum_{i=1}^k \left( \theta_i - \hat{\theta}_i^{(LS)} \right)^2 + \sum_{i=1}^k \left( \hat{\theta}_i^{(LS)} - \hat{\theta}_i^{(b)} \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{i=1}^k \left( \theta_i - \hat{\theta}_i^{(LS)} \right)^2 \right] + \mathbb{E} \left[ \sum_{i=1}^k \left( \hat{\theta}_i^{(LS)} - \hat{\theta}_i^{(b)} \right)^2 \right] \\ &= R(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(LS)}) + \mathbb{E} \left[ \sum_{i=1}^k \left( \hat{\beta}X_i - bX_i \right)^2 \right] \\ &= R(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(LS)}) + \mathbb{E} \left[ \left( \hat{\beta} - b \right)^2 S_X^2 \right] \text{ where } S_X^2 = \sum_{i=1}^k X_i^2 \end{aligned}$$

So a James-Stein estimator will improve on the ordinary estimator if and only if

$$\mathbb{E} \left[ \left( \hat{\beta} - b \right)^2 S_X^2 \right] < \mathbb{E} \left[ \left( \hat{\beta} - 1 \right)^2 S_X^2 \right] \text{ for all } \boldsymbol{\theta}$$

Since  $\hat{\beta}$  is a "reasonable" estimator of  $\beta$  but the constant 1 is not, we can expect the James-Stein estimator to dominate the MLE. This leads us into a landmark result in statistics: Stein's Paradox.

### 13.3 Stein's Paradox

**Theorem 13.1** (Stein's Paradox). *Suppose  $\{X_i\}_{i=1}^k$  are drawn independently by  $X_i \sim \mathcal{N}(\theta_i, 1)$ . Then the James-Stein estimator of  $\boldsymbol{\theta}$ ,*

$$\hat{\boldsymbol{\theta}}^{(JS)} = \left( 1 - \frac{C}{S_X^2} \right) \mathbf{X}, \text{ where } S_X^2 = \sum_{i=1}^k (X_i - \bar{X})^2, k \geq 3, \text{ and } 0 < C \leq 2(k-2)$$

and the Efron-Morris estimator of  $\theta$ ,

$$\hat{\theta}^{(EM)} = \bar{X} + \left(1 - \frac{C}{S_X^2}\right)(\mathbf{X} - \bar{X}), \text{ where } S_X^2 = \sum_{i=1}^k (X_i - \bar{X})^2, k \geq 4, \text{ and } 0 < C \leq 2(k-3)$$

both uniformly dominate (i.e. have uniformly-lower squared error risk) the "obvious" maximum likelihood estimator of  $\theta$ ,

$$\hat{\theta}^{(MLE)} = \mathbf{X}$$

Concisely,

$$\begin{aligned} R(\theta, \hat{\theta}^{(JS)}) &< R(\theta, \hat{\theta}^{(MLE)}) \quad \forall \theta \\ R(\theta, \hat{\theta}^{(EM)}) &< R(\theta, \hat{\theta}^{(MLE)}) \quad \forall \theta \end{aligned}$$

See [SS] for a formal proof.

### 13.3.1 Shrinkage

Why do we call these estimators "shrinkage" estimators?

- The James-Stein estimator  $\hat{\theta}^{(JS)}$  is the weighted average of 0 and  $\mathbf{X}$ , and so **shrinks** the ordinary estimator  $\hat{\theta}^{(MLE)} = \mathbf{X}$  towards 0.
- The Efron-Morris estimator  $\hat{\theta}^{(EM)}$  is the weighted average of  $\bar{X}$  and  $\mathbf{X}$ , and so **shrinks** the ordinary estimator  $\hat{\theta}^{(MLE)} = \mathbf{X}$  towards  $\bar{X}$ .

These shrinkage estimators dominate the ordinary estimator  $\hat{\theta}^{(MLE)} = \mathbf{X}$  as long as  $k \geq 3$  for the James-Stein estimator and  $k \geq 4$  for the Efron-Morris estimator.

### 13.3.2 Connection to Empirical Bayes

Recall the Empirical Bayes estimator of  $\theta$  is given by

$$\hat{\theta}^{(EB)} = \bar{X} + \left(\frac{\tau^2}{\tau^2 + 1}\right)(\mathbf{X} - \bar{X})$$

This estimator has the same form as the Efron-Morris estimator, shrinking towards the overall mean  $\bar{X}$ . Here, the prior variance  $\tau^2$  (where  $\theta_i \sim \mathcal{N}(\theta, \tau^2)$ ) determines how much the estimator shrinks towards the overall mean.

### 13.3.3 The Paradox

To estimate  $\theta_i$ , one of our parameters, it is optimal to use information from all the other observations  $\{X_j\}_{j \neq i}$ , via  $\bar{X}$  and  $S_X^2$ , even though the  $X_i$  are drawn **independently** and are all unrelated in the sense that each has its own mean  $\theta_i$ . This seems preposterous!

How can information about player  $A$ 's batting average and player  $B$ 's batting average help us improve our estimate of player  $C$ 's batting average? How could information about the price of apples in Washington and the price of oranges in Florida help us improve our estimate of the price of French wine, when it's assumed they're unrelated? Herein lies the paradox.

### 13.3.4 Consultant's Dilemma

Suppose that in the middle of the season, an MLB general manager asks you to predict the end-of-season batting average of one player on his team, player  $A$ , using only that season's available data.

The estimator that is best on average across all players (a shrinkage estimator) is different than the estimator that is best for one specific individual player (the MLE).

Optimizing for the squared error aggregated across **all players** is not the same as optimizing for the errors of **separate estimators** of the individual parameters. A combined shrinkage estimator should be used to optimize a **combined loss**, but this combined estimator is **worse** if we want to estimate just one individual parameter.

So... which estimator should we use?

## References

- [SS] Stigler, S. M., *The 1988 Neyman Memorial Lecture: A Galtonian Perspective on Shrinkage Estimators*, Statistical Science, 1990.