

## Lecture 5 — September 15

Lecturer: Purnamrita Sarkar

Scribe: Ryan O'Donnell

**Disclaimer:** These scribe notes have been slightly proofread and may have typos etc.

**Note:** The latex template was borrowed from EECS, U.C. Berkeley.

## 5.1 Quick Note on MLE Existence

Two examples were provided to highlight the possibility that an MLE may not exist. Example 1: Let  $X \sim N(\mu, \sigma^2)$ . Let  $\theta = \langle \mu, \sigma \rangle$ ,  $\theta \in \mathbb{R} \times \mathbb{R}^+$ . The usual pdf for the normal is used here, which is equal to

$$f(x, \theta) = \frac{1}{\sqrt{2\pi} * \sigma} * \exp\left(\frac{-(x - \mu)^2}{2\sigma^2}\right)$$

$$\log(f(x, \theta)) = c - \frac{(x - \mu)^2}{2\sigma^2} - \log(\sigma)$$

So the MLE would be

$$0 = \frac{-(x - \mu)^2}{2\hat{\sigma}^3} - \frac{1}{\hat{\sigma}}$$

$$\hat{\sigma}^2 = (x - \mu)^2$$

But if  $x = \mu$ , then  $\sigma = 0$ , which is not allowed. So the MLE does not exist in this case.

Example 2: Let  $X_i \sim \text{uniform}(0, \theta)$ . If the interval included its boundary, then clearly the MLE would be  $\theta = \max[X_i]$ . But since this interval does not include its boundary, the MLE cannot be the maximum, and therefore an MLE does not exist.

## 5.2 Expounding on the Admissibility of Shrinkage Estimators

As was previously mentioned, it is somewhat difficult to intuitively understand why these particular shrinkage estimators are admissible over the MLE. The below begins with the Bayesian approach to the problem. Beginning with

$$\mathbf{X}|\theta \sim N(\theta, I), \theta \sim N(0, \tau^2 * I)$$

The posterior mean from empirical bayes is just

$$\mathbf{X} * \left(1 - \frac{1}{1 + \tau^2}\right)$$

With that, we can aim to show that the MSE of this posterior estimator is preferable to the MSE of the MLE. For the moment, we will assume we knew  $\tau$  exactly. This allows for an easier proof of the MSE decreasing. In reality, we could perhaps approximate it from the data, though the classical Bayesian approach would not allow for this, as it violates the idea of a prior distribution.

$$E[(\boldsymbol{\theta}_{post} - \boldsymbol{\theta})^T (\boldsymbol{\theta}_{post} - \boldsymbol{\theta})] = E\left[\left(\frac{\tau^2}{1 + \tau^2} \mathbf{X} - \boldsymbol{\theta}\right)^T * \left(\frac{\tau^2}{1 + \tau^2} \mathbf{X} - \boldsymbol{\theta}\right)\right] \quad (5.1)$$

For ease of notation, let

$$c = \frac{1}{1 + \tau^2} \quad (5.2)$$

Therefore, the above becomes

$$\begin{aligned} &= E[(\mathbf{X} - \boldsymbol{\theta} - c\mathbf{X})^T (\mathbf{X} - \boldsymbol{\theta} - c\mathbf{X})] \\ &= E[(\mathbf{X} - \boldsymbol{\theta})^T (\mathbf{X} - \boldsymbol{\theta})] + c^2 E[\mathbf{X}^T \mathbf{X}] - 2cE[\mathbf{X}^T (\mathbf{X} - \boldsymbol{\theta})] \\ &= MSE(\mathbf{X}) - (2c - c^2)E[\mathbf{X}^T \mathbf{X}] + 2cE[\mathbf{X}^T \boldsymbol{\theta}] \end{aligned}$$

To show that the above does indeed equal something smaller than the MSE, it is easiest to break it up into pieces. First, recall the law of iterated expectations. Using this law,

$$\begin{aligned} E[\mathbf{X}] &= E[E[\mathbf{X}|\boldsymbol{\theta}]] = E[\boldsymbol{\theta}] = 0 \\ \text{var}[\mathbf{X}] &= E[\text{var}(\mathbf{X}|\boldsymbol{\theta})] + \text{var}(E[\mathbf{X}|\boldsymbol{\theta}]) = E[I] + \text{var}(\boldsymbol{\theta}) = I(1 + \tau^2) \end{aligned}$$

As it turns out,  $\mathbf{X}$  also has a normal distribution whose parameters using the above derivation is:

$$\mathbf{X} \sim N(0, (1 + \tau^2)I)$$

This is useful because it implies that

$$\frac{\mathbf{X}^T \mathbf{X}}{1 + \tau^2}$$

is a chi-squared distribution with degrees of freedom  $p$ . So, by properties of the chi squared distribution,

$$E[\mathbf{X}^T \mathbf{X}] = (1 + \tau^2)p \quad (5.3)$$

Combining this with the original definition for  $c$  shows that:

$$\begin{aligned} &= (2c - c^2)E[\mathbf{X}^T \mathbf{X}] \\ &= c(2 - c) * E[\mathbf{X}^T \mathbf{X}] \\ &= \frac{1 + 2\tau^2}{(1 + \tau^2)^2} * p \end{aligned}$$

For the next part, the law of iterated expectation and the chi squared distributions are again very useful. The bulk of the work comes from simply implementing the law.

$$E[\boldsymbol{\theta}^T \mathbf{X}] = E[E[\boldsymbol{\theta}^T \mathbf{X} | \theta]] = E[\theta^T E[\mathbf{X} | \theta]] = E[\theta^T \theta] = \sum_i E[\theta_i^2] = \tau^2 p$$

Combining this with the original definition for  $c$  shows that:

$$\begin{aligned} &= 2c * E[\boldsymbol{\theta}^T \mathbf{X}] \\ &= \frac{2\tau^2}{1 + \tau^2} * p \end{aligned}$$

Now, if we combine these two facts with the original definition of  $c$ , we can simplify our original expression for the MSE.

$$\begin{aligned} MSE(\boldsymbol{\theta}_{post}) &= MSE(\mathbf{X}) - (2c - c^2)E[\mathbf{X}^T \mathbf{X}] + 2cE[\mathbf{X}^T \boldsymbol{\theta}] \\ &= MSE(\mathbf{X}) - \frac{1 + 2\tau^2}{1 + \tau^2} * p + \frac{2\tau^2}{1 + \tau^2} * p \\ &= MSE(\mathbf{X}) - \frac{1}{1 + \tau^2} * p \end{aligned}$$

So, as long as we know  $\tau$ , we have found a way to create a shrinkage estimator that is uniformly better than MLE in terms of its MSE. Also, this posterior mean approach creates something that is similar to the James-Stein Estimator. However, this example was not entirely realistic. What if we did not know  $\tau$ ? Would we still do better than the MLE? It turns out that if we use some  $y$  to estimate  $\tau$ , we arrive at the James Stein Estimator.

Recall the following:  $\hat{\boldsymbol{\theta}}_{post} = (1 - \frac{1}{1+\tau^2}) * \mathbf{X}$ .

If we don't know  $\tau$ , we must estimate it. Consider a random variable  $Y$  s.t.

$$E[Y] = \frac{1}{1 + \tau^2}$$

Now, let  $V = \frac{\mathbf{X}^T \mathbf{X}}{1 + \tau^2}$ . By definition,  $V$  is a chi-squared distribution, as it is equal to  $\Sigma(\frac{\mathbf{X}}{\sqrt{1 + \tau^2}})^2$ . Now, take  $\frac{1}{V}$ . This has the inverse chi squared distribution. By properties of the inverse chi squared,  $E[\frac{1}{V}] = \frac{1}{p-2} = E[\frac{1 + \tau^2}{\mathbf{X}^T \mathbf{X}}]$ . Now, notice the following:

$$E[\frac{p-2}{\mathbf{X}^T \mathbf{X}}] = (p-2)E[\frac{1}{\mathbf{X}^T \mathbf{X}}] = (p-2)(\frac{1}{(p-2)(1 + \tau^2)}) = \frac{1}{1 + \tau^2}$$

Therefore, since this yields the desired expectation,  $Y = \frac{p-2}{\mathbf{X}^T \mathbf{X}}$ . Now, using this value of  $y$  as an estimator for  $1 - \frac{1}{1 + \tau^2}$  yields the following, which is equivalent to the James Stein Estimator:

$$\hat{\boldsymbol{\theta}}_{\text{empirical bayes}} = (1 - \frac{p-2}{\mathbf{X}^T \mathbf{X}}) * \mathbf{X}$$

We call this “empirical bayes” since here we used a Bayesian model and then played frequentist by estimating the hyperparameter using the data.

## 5.3 Linear Regression

### 5.3.1 Model and MLE

Here is a linear model for linear regression. Lets first do it for one pair of data points  $(x, y)$ .

$$y = \beta_0 + \beta_1 \cdot x_1 + \dots + \beta_p \cdot x_p + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

Now, for  $n$  data-points  $(x_i, y_i)$ , where  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$  we can write it in matrix notation as follows:

We can write this in matrix form by stacking the datapoints as the rows of a matrix  $\mathbf{X}$  so that  $x_{ij}$  is the  $j$ -th feature of the  $i$ -th datapoint. Then writing  $Y$ ,  $\beta$  and  $\epsilon$  as column vectors, we can write the matrix form of the linear regression model as:

$$\mathbf{y} = \mathbf{X}\beta + \epsilon$$

where:

$$\mathbf{y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \text{ and } \mathbf{X} = \begin{bmatrix} 1 & x_{12} & \dots & x_{1p} \\ 1 & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \dots & x_{np} \end{bmatrix}$$

Assume that  $\epsilon_i$  is normally distributed with variance  $\sigma^2$ . And so  $\epsilon$ . We will now calculate the MLE  $\hat{\beta}$  of  $\beta$ .

We are using the notation where smaller case bold letters denote vectors, capital bold letters denote matrices.

$$f(\mathbf{y}, \beta) \propto \exp\left(\frac{-(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)}{2\sigma^2}\right)$$

Take Log, we can get:

$$\frac{-(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)}{2\sigma^2} \tag{5.4}$$

Same drill— differentiate and set it to zero.

$$-\mathbf{X}^T(\mathbf{y} - \mathbf{X}\hat{\beta}) = 0 \rightarrow \mathbf{X}^T\mathbf{X}\hat{\beta} = \mathbf{X}^T\mathbf{y} \rightarrow \hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

### 5.3.2 Relation to least squares

Lets say I wanted to calculate an estimate that minimized the residual sum of squares (RSS).

$$\beta_{LS} = \min_{\beta'} RSS(\beta') := \min_{\beta'} \sum_i (y_i - \mathbf{x}_i^T \beta')^2$$

As it turns out,  $RSS(\beta')$  is none other than  $(\mathbf{y} - \mathbf{X}\beta')^T(\mathbf{y} - \mathbf{X}\beta')$ . But remember, because the noises are all independently drawn from the same mean zero normal distribution, maximizing log likelihood boils down to minimizing the RSS. And in this special case, the least squares estimate is identical to the MLE.

### 5.3.3 Expectation and Variance of $\hat{\beta}$

Now, we want to find the  $E[\hat{\beta}]$ ,  $Var[\hat{\beta}]$ . Lets put down some ground rules for taking expectations of vector valued random variables. Say  $\mathbf{z} = A\mathbf{y}$  where  $A$  is a fixed matrix.  $E[\mathbf{z}] = AE[\mathbf{y}]$  and  $var(\mathbf{z}) = Avar(\mathbf{y})A^T$ . Recall that  $E[\mathbf{y}] = \mathbf{X}\beta$  and  $var(\mathbf{y}) = \sigma^2\mathbf{I}$

$$E[\hat{\beta}] = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^TE[\mathbf{y}] = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\beta$$
$$var[\hat{\beta}] = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^Tvar[\mathbf{y}]\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$$

Conclusion:  $\hat{\beta} \sim N(\beta, \sigma^2(\mathbf{X}^T\mathbf{X})^{-1})$ . Note: this is not approximate, but exact!