

## WEEK 5: DEC. 1, 2022

Note Writer: Yu-Chieh Kuo<sup>†</sup> and Collaborators<sup>1</sup>: Whiney Yu, Tzu-Yue Huang<sup>\*</sup>

<sup>†</sup>Department of Information Management, National Taiwan University

<sup>\*</sup>Department of Economics, National Taiwan University

### Summary of Consistent Estimators

**Least Squares:**  $Q_n(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \xrightarrow{p} Q_\infty(\theta) = \mathbb{E}[y_i - \hat{y}_i]^2$ , where  $\hat{\theta} \equiv \arg \min Q_n(\theta)$ .

**Maximum Likelihood:**  $Q_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(y_i, x_i, \theta) \xrightarrow{p} Q_\infty(\theta) = \mathbb{E}[\log f(y_i, x_i, \theta)]$ , where  $\hat{\theta} \equiv \arg \max Q_n(\theta)$ .

**GMM, Minimum Distance Estimators:** We have  $\ell$  equations satisfying  $\mathbb{E}[g(y_i, x_i, z_i, \theta)] = 0$  such that  $\bar{g}_n \equiv \frac{1}{n} \sum_{i=1}^n g_i$  and

$$Q_n(\theta) = \bar{g}_n' \hat{W} \bar{g}_n \xrightarrow{p} Q_\infty(\theta) = \mathbb{E}[g_i]' W \mathbb{E}[g_i],$$

where  $\hat{\theta} \equiv \arg \min Q_n(\theta)$ .

### Restricted Estimation

(This section refers to Hansen's textbook, CH8.)

Given  $y_i = x_i' \beta + e_i$  and  $\mathbb{E} x_i e_i = 0$ , we have  $q$  linear constraints such that

$$\underbrace{R'}_{q \times k} \underbrace{\beta}_{k \times 1} = \underbrace{C}_{q \times 1}.$$

Note that the constraint is on the population (parameter space).

### High-dimensional / regularized estimators

The objective function here might be

$$\min_{\hat{\beta}_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^k \hat{\beta}_j^2,$$

where the last term  $\lambda \sum_{j=1}^k \hat{\beta}_j^2$  is the Lagrange multiplier corresponding to  $\sum_{j=1}^k \hat{\beta}_j^2 \leq C$ . It is called **ridge regression**.

In addition, the objective function can be also in the form

$$\min_{\hat{\beta}_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^k |\hat{\beta}_j|,$$

where the last term  $\lambda \sum_{j=1}^k |\hat{\beta}_j|$  is the Lagrange multiplier corresponding to  $\sum_{j=1}^k |\hat{\beta}_j| \leq C$ . It is called **LASSO**.

<sup>1</sup>Yu-Chieh thanks their supports to take photo and provide notes.

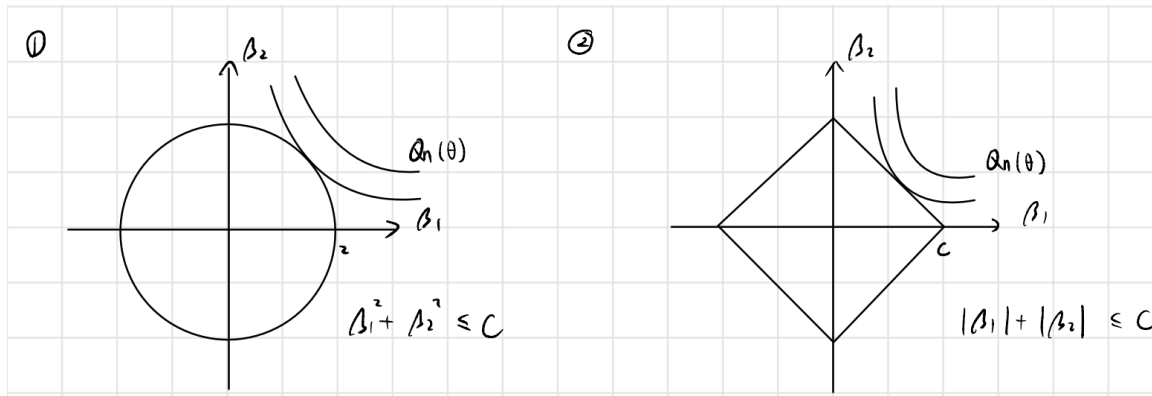


Figure 1: Visualization of the ridge regression and LASSO

### Lagrange function

First, we define the sum of squared errors (SSE) as

$$\begin{aligned}
 SSE &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\
 &= \sum_{i=1}^n (y_i - x_i' \beta)^2 \\
 &= (Y - X\beta)'(Y - X\beta) \\
 &= Y'Y - 2Y'X\beta + \beta'X'X\beta.
 \end{aligned}$$

Be careful about the dimension issues of each matrix:  $Y$  is  $n \times 1$ ,  $X$  is  $n \times k$ , and  $\beta$  is  $k \times 1$ .

Combining SSE, restricted and regularized estimations, we can define the Lagrange function as

$$\mathcal{L} = \frac{1}{2}(Y'Y - 2Y'X\beta + \beta'X'X\beta) + \underbrace{\lambda'(R'\beta - C)}_{1 \times g \text{ and } g \times 1},$$

where the fraction  $\frac{1}{2}$  in the first term is used to cancel the left coefficient 2 after the derivation, and the second term is the Lagrangem multiplier.

Hence, the first partial derivative of the Lagrange function w.r.t.  $\beta$  and  $\lambda$  is

$$\frac{\partial \mathcal{L}}{\partial \beta} = -X'Y + X'X\tilde{\beta} + R\tilde{\lambda} = 0 \tag{1}$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = R'\tilde{\beta} - C = 0. \tag{2}$$

By solving the system to obtain  $\tilde{\beta}$  and  $\tilde{\lambda}$ , we can use  $\tilde{\beta}$  and  $\tilde{\lambda}$  to denote the solutions to restricted estimation problem.

To solve the system, we first pre-multiply (1) by  $R'(X'X)^{-1}$ :

$$\begin{aligned}
 &\underbrace{-R'(X'X)^{-1}X'Y + R'(X'X)^{-1}X'X\tilde{\beta} + R'(X'X)^{-1}R\tilde{\lambda}}_{\hat{\beta}} = 0 \\
 \iff &-R'\hat{\beta} + R'\tilde{\beta} + R'(X'X)^{-1}R\tilde{\lambda} = 0 \\
 \iff &R'\tilde{\beta} = R'\hat{\beta} + R'(X'X)^{-1}R\tilde{\lambda}.
 \end{aligned}$$

Next, we substitute  $R'\hat{\beta} + R'(X'X)^{-1}R\tilde{\lambda}$  for  $R'\tilde{\beta}$  in (2) to solve  $\tilde{\lambda}$ :

$$\begin{aligned} R'\tilde{\beta} &= C \\ \iff R'\hat{\beta} + R'(X'X)^{-1}R\tilde{\lambda} &= C \\ \iff \tilde{\lambda} &= (R'(X'X)^{-1}R)^{-1}(R'\hat{\beta} - C). \end{aligned}$$

Lastly, we substitute  $(R'(X'X)^{-1}R)^{-1}(R'\hat{\beta} - C)$  for  $\tilde{\lambda}$  in (1) to solve  $\tilde{\beta}$ :

$$\begin{aligned} -X'Y + X'X\tilde{\beta} + R(R'(X'X)^{-1}R)^{-1}(R'\hat{\beta} - C) &= 0 \\ \iff \tilde{\beta} &= (X'X)^{-1}X'Y - (X'X)^{-1}R(R'(X'X)^{-1}R)^{-1}(R'\hat{\beta} - C) \\ \iff \tilde{\beta} &= \hat{\beta} - R(R'(X'X)^{-1}R)^{-1}(R'\hat{\beta} - C). \end{aligned}$$

Note that  $R$  is  $k \times q$ .

**Remark.**

1. If  $R'\hat{\beta} - C = 0$ , then  $\tilde{\beta} = \hat{\beta}$ .
2.  $R'(X'X)^{-1}R$  is invertible **only if**  $\text{rank}(R) = q$ .

□

## Consistency

Now we discuss the consistency of the restricted estimation. If it is given  $R'B = C$  and  $\hat{\beta} \xrightarrow{p} \beta_0$  (true  $\beta$ ), then we have

$$R'\hat{\beta} - C \xrightarrow{p} 0 \text{ and } \hat{\beta} \xrightarrow{p} \beta_0 \implies \tilde{\beta} \xrightarrow{p} \beta.$$

## Asymptotic normality

$$\begin{aligned} \sqrt{n}(\tilde{\beta} - \beta) &= \underbrace{\left( \begin{array}{c} \underbrace{\underbrace{\underbrace{0}_{\equiv Q_{XX}^{-1}}} \underbrace{\underbrace{E[x_i x_i' e_i^2]}_{\equiv \Omega}} \underbrace{\underbrace{0}_{\equiv Q_{XX}^{-1}}} \\ \underbrace{\underbrace{(Ex_i x_i')^{-1}} \underbrace{E[x_i x_i' e_i^2]} \underbrace{(Ex_i x_i')^{-1}} \end{array} \right)}_{\stackrel{d}{\sim} \mathcal{N}(0, Cov)} \underbrace{\underbrace{\underbrace{- (X'X)^{-1} R (R'(X'X)^{-1} R)^{-1}}_{\equiv \hat{M} \xrightarrow{p} M}} \underbrace{\underbrace{\sqrt{n}(\hat{\beta} - \beta)}}_{\text{since } C=R'\beta}}_{\sqrt{n}(\hat{\beta} - \beta)} \\ &\stackrel{d}{\sim} \mathcal{N}(0, Cov), \end{aligned}$$

where  $Cov$  is derived by  $\sqrt{n}(\tilde{\beta} - \beta) \sqrt{n}(\tilde{\beta} - \beta)'$ . Clearly,

$$\begin{aligned} Cov &= \sqrt{n}(\tilde{\beta} - \beta) \sqrt{n}(\tilde{\beta} - \beta)' \\ &= n(\hat{\beta} - \beta)(\hat{\beta} - \beta)' - n\hat{M}(\hat{\beta} - \beta)(\hat{\beta} - \beta)' - n(\hat{\beta} - \beta)(\hat{\beta} - \beta)'\hat{M}' + n\hat{M}(\hat{\beta} - \beta)(\hat{\beta} - \beta)'\hat{M}' \\ &\xrightarrow{p} V_\beta - MV_\beta - V_\beta M' + MV_\beta M' \\ &= V_\beta - Q_{XX}^{-1}R(R'Q_{XX}^{-1}R)^{-1}R'V_\beta - V_\beta R(R'Q_{XX}^{-1}R)^{-1}R'Q_{XX}^{-1} + Q_{XX}^{-1}R(R'Q_{XX}^{-1}R)^{-1}R'V_\beta. \end{aligned}$$

### Can we do better?

The answer is **yes**. We may set up a minimum distance estimation as

$$\begin{aligned} \min_{\beta} \quad & \mathcal{J}(\beta) = n(\hat{\beta} - \beta)' \hat{W}(\hat{\beta} - \beta) \\ \text{s.t.} \quad & \text{Constraints,} \end{aligned}$$

where  $\hat{\beta}$  is an OLS estimator (treated as given). Note that  $\beta$  here is a choice variable, not the true parameter.

**Remark.** The Constrained Least Squares (CLS) is a special case where  $\hat{W} = Q_{XX}$ . □

Now, consider the SSE (what is  $\beta$  below. choice variable or true para?)

$$\begin{aligned} SSE &= \sum_{i=1}^n (y_i - x_i' \beta)^2 \\ &= \sum_{i=1}^n (x_i' \hat{\beta} + \hat{e}_i - x_i' \beta)^2 \\ &= \sum_{i=1}^n (\hat{e}_i + x_i' (\hat{\beta} - \beta))^2 \\ &= \sum_{i=1}^n \hat{e}_i^2 + (\hat{\beta} - \beta)' \left( \sum_{i=1}^n x_i x_i' \right) (\hat{\beta} - \beta) + 2 \sum_{i=1}^n \hat{e}_i x_i' (\hat{\beta} - \beta) \\ &= \sum_{i=1}^n \hat{e}_i^2 + (\hat{\beta} - \beta)' \left( \sum_{i=1}^n x_i x_i' \right) (\hat{\beta} - \beta), \end{aligned}$$

where we define  $\mathcal{J}(\beta)$  as the last term  $(\hat{\beta} - \beta)' \left( \sum_{i=1}^n x_i x_i' \right) (\hat{\beta} - \beta)$  with  $\hat{W} = \sum_{i=1}^n x_i x_i'$ .

After obtaining  $\mathcal{J}(\beta)$ , we want to conduct the minimum distance estimation. That is, we solve the system

$$\begin{aligned} \min_{\beta} \quad & \mathcal{J}(\beta) \\ \text{s.t.} \quad & R' \beta = C \quad (\text{Note that } R' \beta_0 = C). \end{aligned}$$

The corresponding Lagrange function is

$$\begin{aligned} \mathcal{L} &= \frac{1}{2} \mathcal{J}(\beta, \hat{W}) + \lambda'(R' \beta - C) \\ &= \frac{n}{2} (\hat{\beta} - \beta)' \hat{W} (\hat{\beta} - \beta) + \lambda'(R' \beta - C), \end{aligned}$$

and FOC w.r.t.  $\beta$  and  $\lambda$  yields

$$\frac{\partial \mathcal{L}}{\partial \beta} = -n \hat{W} (\hat{\beta} - \tilde{\beta}) + R \tilde{\lambda} = 0 \tag{3}$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = R' \tilde{\beta} - C = 0. \tag{4}$$

Extending (3) solves  $\tilde{\beta}$ :

$$\tilde{\beta} = \hat{\beta} - \frac{1}{n} \hat{W}^{-1} R \tilde{\lambda}.$$

Substituting (4) for (3) gives

$$R' \left( \hat{\beta} - \frac{1}{n} \hat{W}^{-1} R \tilde{\lambda} \right) - C = 0 \iff \tilde{\lambda} = n (R' \hat{W}^{-1} R)^{-1} (R' \hat{\beta} - C).$$

Lastly, we use  $\tilde{\lambda}$  to solve  $\tilde{\beta}$  in (3):

$$\begin{aligned} -n\hat{W}(\hat{\beta} - \tilde{\beta}) + nR(R'\hat{W}^{-1}R)^{-1}(R'\hat{\beta} - C) &= 0 \\ \iff \tilde{\beta} &= \hat{\beta} - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}(R'\hat{\beta} - C). \end{aligned}$$

### Consistency

Given  $\hat{\beta} \xrightarrow{p} \beta_0$  (true parameter) and  $R'\hat{\beta} - C \xrightarrow{p} 0$ , we obtain  $\tilde{\beta} \xrightarrow{p} \beta_0$ .

### Asymptotic normality

$$\begin{aligned} \sqrt{n}(\tilde{\beta} - \beta_0) &= \sqrt{n}(\hat{\beta} - \beta_0) - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}R' \overbrace{\sqrt{n}(\hat{\beta} - \beta_0)}^{\xrightarrow{d} \mathcal{N}(0, V_\beta)} \\ &\xrightarrow{d} \mathcal{N}(0, Cov), \end{aligned}$$

where  $Cov$  is

$$\begin{aligned} Cov &= V_\beta - W^{-1}R(R'W^{-1}R)^{-1}R'V_\beta - V_\beta R(R'W^{-1}R)^{-1}R'W^{-1} \\ &\quad + W^{-1}R(R'W^{-1}R)^{-1}R'V_\beta R(R'W^{-1}R)^{-1}R'W^{-1}. \end{aligned}$$

It shows that the most efficient choice of  $W$  is  $V_\beta^{-1}$ . Therefore, the covariance matrix alters to

$$Cov = V_\beta - V_\beta R(R'V_\beta^{-1}R)^{-1}R'V_\beta.$$

In general,  $\tilde{\beta}_{MD}$  (minimum distance) is more efficient than  $\tilde{\beta}_{CLS}$ .

### Short summary

**CLS:** We solve  $\min_\beta \sum_{i=1}^n (y_i - x_i'\beta)^2$  s.t.  $R'\beta = C \implies \tilde{\beta}_{CLS}$ .

**MD Estimation:** We solve  $\min_\beta (\hat{\beta} - \beta)' \hat{W}(\hat{\beta} - \beta)$  s.t.  $R'\beta = C \implies \tilde{\beta}_{MD}$ .

Note that CLS is a special case where

$$\hat{W} = \frac{1}{n} \sum_{i=1}^n x_i x_i' \xrightarrow{p} W = \mathbb{E}[x_i x_i'],$$

but the efficient weight matrix is

$$\hat{W} = \left( \frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n x_i x_i' \hat{e}_i^2 \right) \left( \frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \xrightarrow{p} W = \mathbb{E}[x_i x_i']^{-1} \mathbb{E}[x_i x_i' \hat{e}_i^2] \mathbb{E}[x_i x_i']^{-1} = V_\beta^{-1}.$$

Consequently,  $\tilde{\beta}_{MD}$  is more efficient than  $\tilde{\beta}_{CLS}$ .

**Example.** Given a regression  $y_i = x_{i1}'\beta_1 + x_{i2}'\beta_2 + e_i$  with a constraint  $\beta_2 = 0$ , we can show that the estimator from regression without  $x_{2i}$  is identical with the CLS estimator with  $\beta_2 = 0$ .

Another example can be found at Page. 269 in Hansen's textbook.  $\square$

## Misspecification

(This section refers to Hansen's textbook, CH8.13.)

In the case that  $R'\beta = C^* \neq C$ , the MD estimator alters to

$$\tilde{\beta}_{MD} = \hat{\beta} - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}(R'\hat{\beta} - C) \xrightarrow{p} \beta - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}(C^* - C) \equiv \beta_n^*.$$

The asymptotic normality becomes

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_{MD} - \beta_n^*) &= \sqrt{n}(\hat{\beta} - \beta) - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1} \sqrt{n}(R'\hat{\beta} - C^*) \\ &= \sqrt{n}(\hat{\beta} - \beta) - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1} \sqrt{n}(R'\hat{\beta} - R'\beta) \\ &= (I - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}R) \sqrt{n}(\hat{\beta} - \beta) \\ &\xrightarrow{d} \mathcal{N}(0, V_\beta(W)), \end{aligned}$$

where  $V_\beta(W)$  is the same asymptotic covariance in the case without misspecification. **why???????**

Another case for the misspecification issue might be in the form of  $R'\beta_n = C + \delta\sqrt{n}$ . In this case,  $R'\hat{\beta} - C = R'(\hat{\beta} - \beta_n) + \delta\sqrt{n}$ , and the MD estimator is

$$\begin{aligned} \tilde{\beta}_{MD} &= \hat{\beta} - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}(R'\hat{\beta} - C) \\ &= \hat{\beta} - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}R'(\hat{\beta} - \beta_n) - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}R'\delta\sqrt{n}. \end{aligned}$$

The asymptotic normality in this case becomes

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_{MD} - \beta_n) &= \overbrace{\sqrt{n}(\hat{\beta} - \beta_n)}^{\xrightarrow{d}\mathcal{N}(0, V_\beta)} - \underbrace{\hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}R'}_{\equiv \delta^*} \overbrace{\sqrt{n}(\hat{\beta} - \beta_n) - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}R'\delta\sqrt{n}}^{\xrightarrow{d}\mathcal{N}(0, V_\beta)} \\ &\xrightarrow{d} \mathcal{N}(0, V_\beta(W)) - \hat{W}^{-1}R(R'\hat{W}^{-1}R)^{-1}R'\delta \\ &= \mathcal{N}(\delta^*, V_\beta(W)). \end{aligned}$$