

# VE472 Lecture 6

Jing Liu

UM-SJTU Joint Institute

Summer

- Although a simple training-test split is done randomly, only one such split is actually used in the above approach. For a small  $n$ , the resulting  $\lambda_{\text{opt}}$  based on the above approach might vary a lot from one random split to another.
- Leave-one-out cross-validation (LOOCV) can remove this layer of variability.  
split:  $n-1$  as training set, 1 as test set

---

### Algorithm 1: Determining $\lambda_{\text{opt}}$ using LOOCV

---

**Input** : Data matrix  $\mathbf{X}$ , data vector  $\mathbf{y}$ , set of parameters  $\mathcal{S}_\lambda$

**Output** : Optimal shrinkage parameter  $\lambda_{\text{opt}}$

```

1 Function LamLOOCV( $\mathbf{X}, \mathbf{y}, \mathcal{S}_\lambda$ ):
2   for  $i \leftarrow 1$  to  $n$  do
3      $\mathbf{X}_{-i} \leftarrow \mathbf{X}[-i,];$  /* Create  $n$  Training sets, the  $i$ th set */
4      $\mathbf{y}_{-i} \leftarrow \mathbf{y}[-i];$  /* consists all cases except the  $i$ th case */
5   end for
6    $\lambda_{\text{opt}} \leftarrow \arg \min_{\lambda \in \mathcal{S}_\lambda} \left\{ \sum_{i=1}^n \left( \mathbf{y}[i] - \mathbf{X}[i,] (\mathbf{X}_{-i}^T \mathbf{X}_{-i} + \lambda \mathbf{I})^{-1} \mathbf{X}_{-i}^T \mathbf{y}_{-i} \right)^2 \right\};$ 
7   return  $\lambda_{\text{opt}};$ 
8 end

```

每一行作为test set的时候的error

without  $i$ th row

- Notice the last algorithm is only realistic for relatively small  $n$  since

$$(\mathbf{X}_{-i}^T \mathbf{X}_{-i} + \lambda \mathbf{I})^{-1}$$

n大了算起来太慢  
就算用LU, etc也太慢了

need to be done  $n$  times for every  $\lambda$ , that is, one for every  $i = 1, 2, \dots, n$ .

Q: Why does the following give a more efficient way to implement Algorithm 1.

### Theorem 0.1

Let  $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T$  and  $\mathbf{D}$  be a diagonal matrix of  $n \times n$  with

$$\mathbf{D}_{ii} = \frac{1}{1 - \eta_{ii}}, \quad \text{where} \quad \eta_{ii} = \mathbf{X}_i (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}_i^T,$$

and  $\mathbf{A}_{-i} = \mathbf{X}_{-i}^T \mathbf{X}_{-i} + \lambda \mathbf{I}$ , then the following holds

只要算一遍inverse

$$\sum_{i=1}^n (y_i - \mathbf{X}_i \mathbf{A}_{-i}^{-1} \mathbf{X}_{-i}^T \mathbf{y}_{-i})^2 = \|\mathbf{D} (\mathbf{I} - \mathbf{H}) \mathbf{y}\|^2$$

前一页的sum of squares

where  $\|\cdot\|$  denotes the usual Euclidean norm.

---

**Algorithm 2:** Determining  $\lambda_{\text{opt}}$  using LOOCV with Woodbury identity

---

**Input** : Data matrix  $\mathbf{X}$ , data vector  $\mathbf{y}$ , set of parameters  $\mathcal{S}_\lambda$

**Output** : Optimal shrinkage parameter  $\lambda_{\text{opt}}$

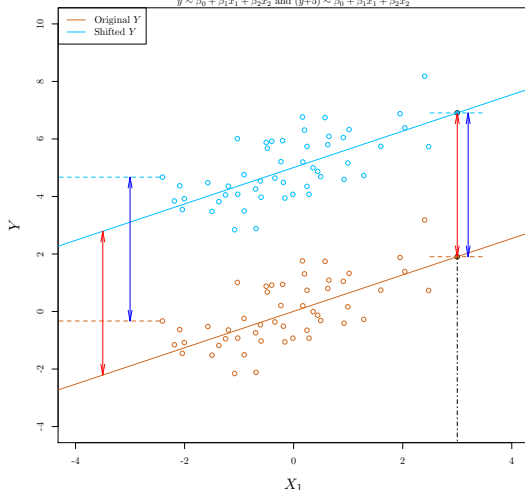
```
1 Function LamLOOCVWoodbury( $\mathbf{X}, \mathbf{y}, \mathcal{S}_\lambda$ ):  
    /* Create functions of  $\lambda$  needed for the optimisation */  
2     $\mathbf{A}^{-1}(\lambda) \leftarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}$ ; /* matrix-valued function of  $\lambda$  */  
3     $\mathbf{H}(\lambda) \leftarrow \mathbf{X} \mathbf{A}^{-1} \mathbf{X}^T$ ; /* matrix-valued function of  $\lambda$  */  
4    for  $i \leftarrow 1$  to  $n$  do  
5         $\eta_{ii}(\lambda) \leftarrow \mathbf{X}[i,] \mathbf{A}^{-1} \mathbf{X}[i,]^T$ ; /* real-valued functions of  $\lambda$  */  
6    end for  
7     $\mathbf{D}(\lambda) \leftarrow \text{diag} \left( \frac{1}{1 - \eta_{11}(\lambda)}, \frac{1}{1 - \eta_{22}(\lambda)}, \dots, \frac{1}{1 - \eta_{nn}(\lambda)} \right)$ ;  
    /* matrix-valued function of  $\lambda$  */  
8     $\lambda_{\text{opt}} \leftarrow \arg \min_{\lambda \in \mathcal{S}_\lambda} \left\{ \left\| \mathbf{D}(\lambda) (\mathbf{I} - \mathbf{H}(\lambda)) \mathbf{y} \right\|^2 \right\}$ ;  
9    return  $\lambda_{\text{opt}}$ ;  
10 end
```

---

- Notice the following desirable property that the least squares estimator has.

### Usual least squares estimate

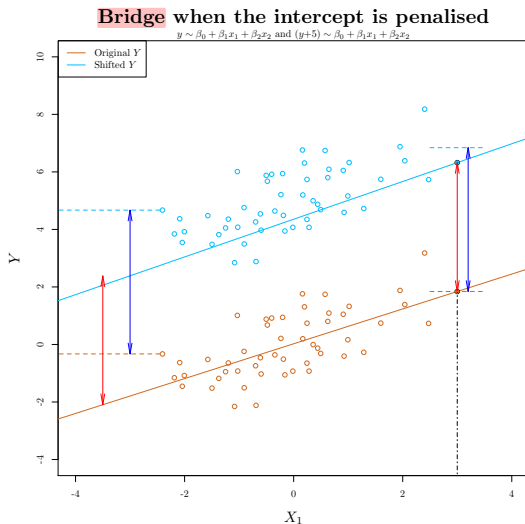
$$y \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2 \text{ and } (y+5) \sim \beta_0 + \beta_1 x_1 + \beta_2 x_2$$



所有的y加了5

slope is the  
same,  
intercept  
changed

- However, it isn't true in our current formulation of  $\hat{\beta}_{\text{ridge}}$ , i.e. penalising  $\hat{\beta}_0$



shift  $y \rightarrow$   
 slope intercept  
 都变了

- The above simulation illustrates a desirable property of  $\hat{\beta}_{\text{lse}}$  that our current formulation of  $\hat{\beta}_{\text{ridge}}$  does not have, namely, altering the centre of  $\mathbf{y}$ ,

目标: slope不change

$$\bar{y} = \frac{1}{n} \mathbf{1}^T \mathbf{y}, \quad \text{where } \mathbf{1} \text{ denotes the vector of ones,}$$

y bar: mean of all components

only alters the intercept component of  $\hat{\beta}_{\text{lse}}$ , while it alters  $\hat{\beta}_{\text{ridge}}$  completely.

- To understand why, let  $\mathbf{y}_c = \mathbf{y} - \bar{y}\mathbf{1}$ , then for any centre  $\bar{y}$  we can rewrite

yc centered at 0

$$\begin{aligned} \mathbf{b}_{\text{lse}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y}_c + \bar{y}\mathbf{1}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_c + \bar{y} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{1} \end{aligned}$$

first column of  $\mathbf{X}^T \mathbf{X}$  (X第一列都是1, 因为有intercept)

- Note  $\mathbf{X}^T \mathbf{1}$  is the 1st column of  $\mathbf{X}^T \mathbf{X}$ , so  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{1} = [1 \quad \mathbf{0}_{1 \times k}]^T$  and

$$\arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 \right\} = \arg \min_{\mathbf{b}_c \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y}_c - \mathbf{X}\mathbf{b}_c\|^2 \right\} + \begin{bmatrix} \bar{y} & \mathbf{0}_{1 \times k} \end{bmatrix}^T$$

intercept = 0的部分

changing according to changing of intercept

- However, in terms of our current formulation of  $\hat{\beta}_{\text{ridge}}$ , we have

$$\begin{aligned}\mathbf{b}_{\text{ridge}} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{y}_c + \bar{y} \mathbf{1}) \\ &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}_c + \bar{y} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{1}\end{aligned}$$

- Note the vector  $\mathbf{X}^T \mathbf{1}$  is not the 1st column of  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ , thus in general,

$$\underline{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{1} \neq \begin{bmatrix} 1 & \mathbf{0}_{1 \times k} \end{bmatrix}^T}$$

hence the estimate will change more than just the first component,

$$\begin{aligned}\mathbf{b}_{\text{ridge}} &= \arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\} \\ &\neq \arg \min_{\mathbf{b}_c \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y}_c - \mathbf{X}\mathbf{b}_c\|^2 + \lambda \|\mathbf{b}_c\|^2 \right\} + \begin{bmatrix} \bar{y} & \mathbf{0}^T \end{bmatrix}^T\end{aligned}$$

- A small change in our formulation of  $\hat{\beta}_{\text{ridge}}$  will allow us to rectify it.



- Suppose  $\mathbf{X} = [\mathbf{1} \quad \mathbf{X}_{-0}]$  and  $\mathbf{b}^T = [b_0 \quad \mathbf{b}_{-0}^T]$ , then we have the following

$$\begin{aligned}\|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 &= \|\mathbf{y}_c + \bar{y}\mathbf{1} - b_0\mathbf{1} - \mathbf{X}_{-0}\mathbf{b}_{-0}\|^2 \\ &= \|\mathbf{y}_c - (b_0 - \bar{y})\mathbf{1} - \mathbf{X}_{-0}\mathbf{b}_{-0}\|^2 = \|\mathbf{y}_c - \mathbf{X}\mathbf{b}_c\|^2\end{aligned}$$

from which we can conclude  $\mathbf{b}_c^T = [b_0 - \bar{y} \quad \mathbf{b}_{-0}^T]$  for any centre  $\bar{y}$ .

- Note  $\|\mathbf{b}_c\|$  varies as  $\bar{y}$  varies due to **the 1st component**, so does the penalty

$$\lambda\|\mathbf{b}_c\|^2$$

for the same  $\mathbf{b}$  and  $\lambda$ , which suggests an approach of rectifying the difference

$$\begin{aligned}\mathbf{b}_{\text{ridge}} &= \arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \lambda\|\mathbf{b}\|^2 \right\} \\ &\neq \arg \min_{\mathbf{b}_c \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y}_c - \mathbf{X}\mathbf{b}_c\|^2 + \lambda\|\mathbf{b}_c\|^2 \right\} + [\bar{y} \quad \mathbf{0}^T]^T\end{aligned}$$

Q: Can you guess what this approach is?

- Consider the following alternative formulation  $\hat{\beta}_{\text{ridge}}$

$$\mathbf{b}_{\text{ridge}} = \arg \min_{b_0 \in \mathbb{R}, \mathbf{b}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - b_0 \mathbf{1} - \mathbf{X}_{-0} \mathbf{b}_{-0}\|^2 + \lambda \|\mathbf{b}_{-0}\|^2 \right\}$$

which differs from our original formulation by exempting  $b_0$  from penalty.

- Differentiating the above objective function with respect to  $b_0$ , we have

$$-2\mathbf{1}^T \mathbf{y} + 2nb_0 + 2\mathbf{1}^T \mathbf{X}_{-0} \mathbf{b}_{-0}$$

setting which to zero, we have

$$\text{intercept } b_0 = \frac{1}{n} (\mathbf{1}^T \mathbf{y} - \mathbf{1}^T \mathbf{X}_{-0} \mathbf{b}_{-0}) = \overset{\text{center of y}}{\bar{y}} - \underbrace{\frac{1}{n} (\mathbf{1}^T \mathbf{X}_{-0}) \mathbf{b}_{-0}}_{\text{?}}$$

row vector of sample mean of each variable

which means if we apply a mean-centring shift to  $\mathbf{X}_{-0}$ , then the value of  $b_0$  that attains the minimum is given by  $b_0 = \bar{y}$  independent of  $\mathbf{b}_{-0}$  and  $\lambda$ .

- Let us denote the mean-centring shift in  $\mathbf{X}_{-0}$  by

$$z_{ij} = x_{ij} - \bar{x}_j \quad \text{for } j = 1, 2, \dots, k \quad \text{and} \quad i = 1, 2, \dots, n$$

where  $\bar{x}_j$  is the mean of the  $j$ th column of  $\mathbf{X}_{-0}$ , in matrix notation, we have

$$\mathbf{Z}_{-0} = \mathbf{X}_{-0} - \mathbf{1}\bar{\mathbf{x}}^T \quad \text{where} \quad \bar{\mathbf{x}}^T = \frac{1}{n}\mathbf{1}^T\mathbf{X}_{-0}$$

$$\begin{aligned} \text{so } \mathbf{b}_{\text{ridge}} &= \arg \min_{b_0 \in \mathbb{R}, \mathbf{b}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - b_0\mathbf{1} - \mathbf{X}_{-0}\mathbf{b}_{-0}\|^2 + \lambda\|\mathbf{b}_{-0}\|^2 \right\} \\ &= \arg \min_{b_0 \in \mathbb{R}, \mathbf{b}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - b_0\mathbf{1} - (\mathbf{Z}_{-0} + \mathbf{1}\bar{\mathbf{x}}^T)\mathbf{b}_{-0}\|^2 + \lambda\|\mathbf{b}_{-0}\|^2 \right\} \\ &= \arg \min_{b_0 \in \mathbb{R}, \mathbf{b}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - (b_0 + \bar{\mathbf{x}}^T\mathbf{b}_{-0})\mathbf{1} - \mathbf{Z}_{-0}\mathbf{b}_{-0}\|^2 + \lambda\|\mathbf{b}_{-0}\|^2 \right\} \\ &= \arg \min_{b_0 \in \mathbb{R}, \mathbf{b}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y}_c - \mathbf{Z}_{-0}\mathbf{b}_{-0}\|^2 + \lambda\|\mathbf{b}_{-0}\|^2 \right\} \end{aligned}$$

where the minimiser  $b_0 = \bar{y} - \frac{1}{n}(\mathbf{1}^T\mathbf{X}_{-0})\mathbf{b}_{-0}$ , and  $\mathbf{y} = \mathbf{y}_c - \bar{y}\mathbf{1}$  are used.

- Therefore, we can conclude  $\mathbf{b}_{-0}$  is invariant under mean-centring

$$\begin{aligned}\mathbf{b}_{\text{ridge}} &= \arg \min_{\mathbf{b}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - b_0 \mathbf{1} - \mathbf{X}_{-0} \mathbf{b}_{-0}\|^2 + \lambda \|\mathbf{b}_{-0}\|^2 \right\} \\ &= \arg \min_{\mathbf{v}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - v_0 \mathbf{1} - \mathbf{Z}_{-0} \mathbf{v}_{-0}\|^2 + \lambda \|\mathbf{v}_{-0}\|^2 \right\}\end{aligned}$$

where  $v_0 = \bar{y}$  and  $\mathbf{b}_{-0} = \mathbf{v}_{-0}$ , and prediction is also invariant

$$\begin{aligned}\hat{y}_i &= b_0 + \mathbf{x}_{i,}^T \mathbf{b}_{-0} = (v_0 - \bar{\mathbf{x}}^T \mathbf{b}_{-0}) + \mathbf{x}_{i,}^T \mathbf{b}_{-0} \\ &= v_0 + (\mathbf{x}_{i,}^T - \bar{\mathbf{x}}^T) \mathbf{b}_{-0} \\ &= v_0 + \mathbf{z}_{i,}^T \mathbf{v}_{-0}\end{aligned}$$

since the last expression is the prediction for  $Y$  in terms of the centred data

$$\mathbf{z}_{i,} = \mathbf{x}_{i,} - \bar{\mathbf{x}}$$

- Note this invariant property is actually true for any shift  $\mathbf{Z}_{-0} = \mathbf{X}_{-0} - \mathbf{1}\alpha^T$ .

- Recall what we have defined

$$\begin{aligned}\mathbf{X} &= [\mathbf{1} \quad \mathbf{X}_{-0}]; & \mathbf{b}^T &= [b_0 \quad \mathbf{b}_{-0}^T]; & \mathbf{y} &= \mathbf{y}_c + \bar{y}\mathbf{1} \\ \mathbf{Z}_{-0} &= \mathbf{X}_{-0} - \mathbf{1}\bar{\mathbf{x}}^T; & \mathbf{v}^T &= [v_0 \quad \mathbf{v}_{-0}^T]; & \mathbf{v}_{-0} &= \mathbf{b}_{-0}\end{aligned}$$

Q: Have we rectified our first formulation of  $\hat{\beta}_{\text{ridge}}$  with our second formulation

$$\begin{aligned}\mathbf{b}_{\text{ridge}} &= \arg \min_{b_0 \in \mathbb{R}, \mathbf{b}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - b_0\mathbf{1} - \mathbf{X}_{-0}\mathbf{b}_{-0}\|^2 + \lambda \|\mathbf{b}_{-0}\|^2 \right\} \\ &= \arg \min_{v_0 \in \mathbb{R}, \mathbf{v}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y} - v_0\mathbf{1} - \mathbf{Z}_{-0}\mathbf{v}_{-0}\|^2 + \lambda \|\mathbf{v}_{-0}\|^2 \right\} \\ &= \arg \min_{v_0 = \bar{y}, \mathbf{v}_{-0} \in \mathbb{R}^k} \left\{ \|\mathbf{y}_c - \mathbf{Z}_{-0}\mathbf{v}_{-0}\|^2 + \lambda \|\mathbf{v}_{-0}\|^2 \right\} \\ &= \begin{bmatrix} \bar{y} & \left( (\mathbf{Z}_{-0}^T \mathbf{Z}_{-0} + \lambda \mathbf{I})^{-1} \mathbf{y}_c \right)^T \end{bmatrix}^T\end{aligned}$$

Q: Is the above estimator invariant under scaling of  $\mathbf{y}$  or columns of  $\mathbf{Z}_{-0}$  ?

no

- Let  $\mathbf{M}$  be an invertible diagonal matrix of  $k \times k$ , and  $\mathbf{U} = \mathbf{Z}_{-0}\mathbf{M}^{-1}$ , then

$$\begin{aligned} \underline{(\mathbf{Z}_{-0}^T \mathbf{Z}_{-0} + \lambda \mathbf{I})^{-1} \mathbf{y}_c} &= \left( (\mathbf{U}\mathbf{M})^T \mathbf{U}\mathbf{M} + \lambda \mathbf{I} \right)^{-1} \mathbf{y}_c \\ &= (\mathbf{M}^T \mathbf{U}^T \mathbf{U} \mathbf{M} + \lambda \mathbf{I})^{-1} \mathbf{y}_c \neq \underline{(\mathbf{U}^T \mathbf{U} + \lambda \mathbf{I})^{-1} \mathbf{y}_c} \end{aligned}$$

since  $\mathbf{U}^T \mathbf{U} \mathbf{M} \neq \mathbf{M} \mathbf{U}^T \mathbf{U}$  and  $\mathbf{M}$  is not orthogonal in general, which means our second formulation is not invariant under column scaling of  $\mathbf{Z}_{-0}$ .

Q: Should we scale the the data matrix  $\mathbf{Z}_{-0}$ ? If so, what should we use for  $\mathbf{M}$ ?

- Unless the independent variables have the same units, to be “fair” we scale the columns of the data matrix so that the sample variances become 1,

$$\mathbf{U} = \mathbf{Z}_{-0} \text{diag}(s_1, s_2, \dots, s_k)^{-1}, \quad \text{where} \quad s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_{ij} - \bar{x}_j)^2$$

which is the usual formulation of  $\hat{\beta}_{\text{ridge}}$  that is widely implemented.

### Algorithm 3: Ridge Estimation

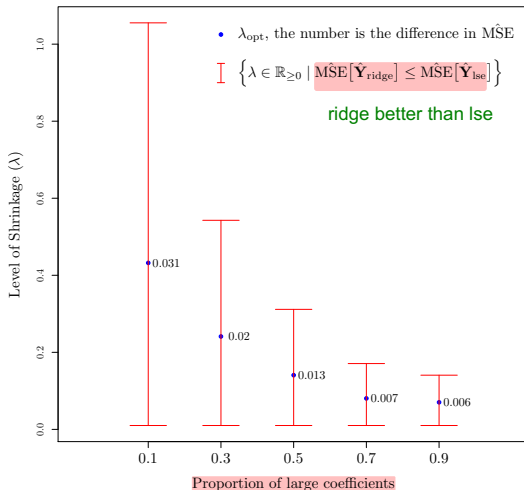
**Input** : Data matrix  $\mathbf{X}_{n \times (k+1)}$ , data vector  $\mathbf{y}_{n \times 1}$ , set of parameters  $\mathcal{S}_\lambda$

**Output** : Vector of estimated parameters  $\mathbf{b}_{\text{ridge}}$

```
1 Function RIDGE( $\mathbf{X}, \mathbf{y}, \mathcal{S}_\lambda$ ):
2    $\mathbf{X}_{-0} \leftarrow \mathbf{X}[:, -1]$  ; /* remove the column of ones */
3    $\mathbf{Z}_{-0} \leftarrow \mathbf{X}_{-0} - \frac{1}{n} \mathbf{1}^T \mathbf{X}_{-0}$  ; /* mean-centre the columns of  $\mathbf{X}_{-0}$  */
4   for  $j \leftarrow 1$  to  $k$  do
5      $s_j = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (\mathbf{Z}_{-0}[i, j])^2}$  ; /* column standard deviation
6     */
7   end for
8    $\mathbf{U} \leftarrow \mathbf{Z}_{-0} \text{diag}(s_1, s_2, \dots, s_k)^{-1}$  ; /* scale the columns of  $\mathbf{Z}_{-0}$  */
9    $\mathbf{b}_{\text{ridge}}[1] \leftarrow \frac{1}{n} \mathbf{1}^T \mathbf{y}$  ; /* set the mean of  $\mathbf{y}$  as  $b_0$  */
10   $\mathbf{y}_c \leftarrow \mathbf{y} - \mathbf{b}_{\text{ridge}}[1] \cdot \mathbf{1}$  ; /* mean-centre  $\mathbf{y}$  */
11   $\lambda_{\text{opt}} \leftarrow \text{LamLOOCVWoodbury}(\mathbf{U}, \mathbf{y}_c, \mathcal{S}_\lambda)$  ; /* find the optimal  $\lambda$  */
12   $\mathbf{b}_{\text{ridge}}[-1] \leftarrow (\mathbf{U}^T \mathbf{U} + \lambda_{\text{opt}} \mathbf{I})^{-1} \mathbf{y}_c$  ; /* Slope estimates */
13 return  $\mathbf{b}_{\text{ridge}}$  ;
14 end
```

- $\hat{\beta}_{\text{ridge}}$  only slightly outperform  $\hat{\beta}_{\text{lse}}$  for a narrow range of  $\lambda$  if  $\beta$  is mostly big.

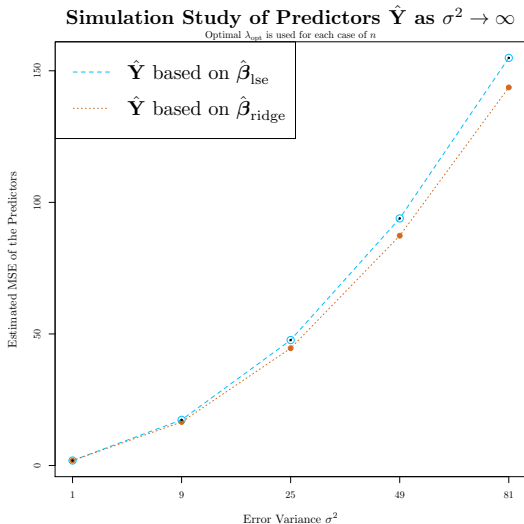
### The effect of true coefficients



small proportion of large coefficients  $\rightarrow$  a large number of lambda that we can choose to outperform  $\hat{\beta}_{\text{lse}}$

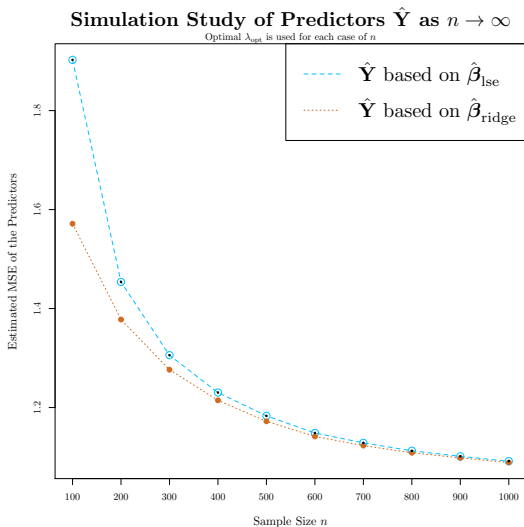


- As  $\sigma^2$  grows, the tradeoff between variance and bias becomes more notable.



noise多, 用ridge好

- However, as  $n$  increases, the performance of  $\hat{\beta}_{\text{lse}}$  catches up with  $\hat{\beta}_{\text{ridge}}$ .



$n$ 很大, both  
converge to real  
beta

- Our motivation for  $\hat{\beta}_{\text{ridge}}$ , i.e. bypass variable selection and achieve a smaller

$$\text{MSE}(\hat{Y}_i) = \text{MSE}(\mathbf{x}_i^T \hat{\beta})$$

ridge其实适用的是medium size data  
n或者k太大不行

becomes weak if  $n$  as well as  $k$  is relatively large.

- Although it is not exclusive to big/complex datasets, rank deficiency issue,

$$\det(\mathbf{X}^T \mathbf{X}) \approx 0$$

is often more prominent in more complex datasets.

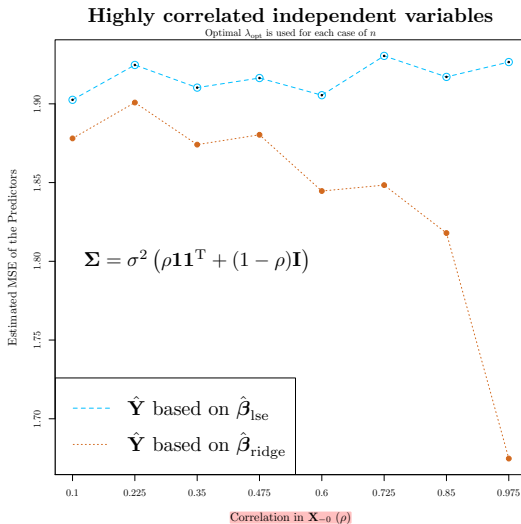
- We will show later on the motivation for using the ridge estimate for large  $n$

$$\mathbf{b}_{\text{ridge}} = \begin{bmatrix} \bar{y} \\ (\mathbf{U}^T \mathbf{U} + \lambda \mathbf{I})^{-1} \mathbf{U}^T \mathbf{y}_c \end{bmatrix}$$

is the fact that the matrix  $\mathbf{U}^T \mathbf{U} + \lambda \mathbf{I}$  is always invertible for  $0 < \lambda < \infty$ .

- Highly Correlated

有rank deficiency:  
用ridge更好



## Recap: Regression

- Linear least squares regression has zero bias but suffers from high variance.

$$\text{MSE}(\hat{Y}_i) = \mathbb{E} \left[ \left( \hat{Y}_i - Y_i \right)^2 \right] = \text{Var}[\hat{Y}_i] + \underbrace{\left( \mathbb{E}[\hat{Y}_i] - \mathbb{E}[Y_i] \right)^2}_{\text{Bias}} + \sigma^2$$

- Ridge regression reduces mean square error without doing variable selection.
- Various things, the size of  $k$  and  $n$ , the relative size of the true coefficients, and the relationship between independent variables, decides when we prefer ridge over linear least square regression when building a predictive model.
- It is used for building predictive models, where  $k$  is relatively large while  $n$  is relatively small. It is particularly good if the independent variables are highly correlated, and there is a subset of true coefficients that are relatively small.

Q: What should we do in terms of regression if  $n$  is very large? Very large  $k$ ?