

## Week 36 exercise 1

(a) Show that the optimal parameters

$$\hat{\beta}_{\text{ridge}} = (\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \underline{y}$$

with  $\underline{I} \in \mathbb{R}^{P \times P}$  being the identity matrix.

We have to solve the Ridge regression minimization problem

$$\min_{\beta \in \mathbb{R}^P} \frac{1}{n} \|\underline{y} - \underline{X}\underline{\beta}\|_2^2 + \lambda \|\underline{\beta}\|_2^2 \quad (1)$$

The minimization problem is a question of finding the optimal vector  $\underline{\beta}$ .

Since this is a sum we split it into two parts. I start with the latter and utilize the identity  $\frac{\partial}{\partial \underline{x}} \|\underline{x}\|_2 = \underline{x}$

$$\frac{\partial}{\partial \underline{\beta}} \left[ \lambda \|\underline{\beta}\|_2^2 \right] = 2 \lambda \underline{\beta} \quad (2)$$

Secondly we have the first term

$$\frac{1}{n} \frac{\partial}{\partial \beta} \left[ \| y - \underline{x}\beta \|^2_2 \right]$$

$$= \frac{1}{n} \frac{\partial}{\partial u} \left[ \| u \|^2_2 \right] \frac{\partial}{\partial \beta} \left[ y - \underline{x}\beta \right]$$

$$= \frac{1}{n} \left( 2(y - \underline{x}\beta) \right) (-\underline{x})$$

$$= -\frac{2}{n} \left( y - \underline{x}\beta \right) \underline{x}$$

$$= \frac{2}{n} \underline{x}^T (\underline{x}\beta - y) \quad (3)$$

To do the actual optimization we now find the extremes by  $(3)+(2)=0$

$$\frac{2}{n} \underline{x}^T (\underline{x}\beta - y) + 2\lambda \beta = 0$$

$$\frac{2}{n} \underline{x}^T \underline{x}\beta - \frac{2}{n} \underline{x}^T y + 2\lambda \beta = 0$$

By multiplying by  $\frac{n}{2}$  we get

$$\underline{x}^T \underline{x}\beta + n\lambda \beta = \underline{x}^T y$$

Extracting  $\underline{\beta}$  from LH we have to multiply  $\underline{I}$  into the equation for it to be valid. Further since  $n$  is a const we can absorb it into  $\lambda$ ,  $\lambda = \lambda' = n\lambda$

$$(\underline{X}^T \underline{X} + \lambda \underline{I}) \underline{\beta} = \underline{X}^T \underline{y}$$

Assuming invertability

$$\Rightarrow \hat{\underline{\beta}}_{\text{ridge}} = (\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \underline{y} \blacksquare$$

$$1b) \hat{Y}_{OLS} = \underline{X} \underline{\beta} \quad (1)$$

For  $\underline{\beta}$  we have the sol.  
from (1a)

$$\underline{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y} \quad (2)$$

and  $\underline{X}$  can be expressed  
by the SVD as

$$\underline{X} = \underline{U} \underline{\Sigma} \underline{V}^T \quad (3)$$

where :

$\underline{U} \in \mathbb{R}^{n \times n}$ ,  $\underline{\Sigma} \in \mathbb{R}^{n \times p}$ ,  
 $\underline{V} \in \mathbb{R}^{p \times p}$  matrix on the form

$$\underline{\Sigma} = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix}, D = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \end{bmatrix}$$

where  $\sigma$ 's are the first

$r = \text{Rank}(\underline{X})$  singular values  
 $> 0$ .

Combining (1) and (2):

$$\hat{y}_{OLS} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$$

To replace the  $\underline{X}$  matrices we have

$$\hookrightarrow \underline{X} : (3)$$

$$\hookrightarrow (\underline{X}^T \underline{X})^{-1} : (\underline{U}^T \underline{\Sigma}^T \underline{V} \underline{U} \underline{\Sigma} \underline{V}^T)^{-1}$$

$$= (\underline{V} \underline{\Sigma}^T \underline{\Sigma} \underline{V}^T)^{-1}$$

utilising the orthogonality of  $\underline{V}$  we get

$$(\underline{X}^T \underline{X})^{-1} = \underline{V}^T (\underline{\Sigma}^T \underline{\Sigma})^{-1} \underline{V}$$

$$\hookrightarrow \underline{X}^T = \underline{U}^T \underline{\Sigma}^T \underline{V}$$

This yields as

$$\begin{aligned}
 \tilde{\mathbf{x}}_{OLS} &= \underline{\mathbf{U}} \underline{\Sigma} \underline{\mathbf{V}^T} \underline{\mathbf{V}}^T (\underline{\Sigma}^T \underline{\Sigma})^{-1} \underline{\mathbf{V}} \underline{\mathbf{U}^T} \underline{\Sigma}^T \underline{\mathbf{V}} \\
 &= \underline{\mathbf{U}} \underline{\Sigma} \underline{\mathbf{V}^T} \underline{\mathbf{V}} (\underline{\Sigma}^T \underline{\Sigma})^{-1} \underline{\mathbf{V}^T} \underline{\mathbf{V}} \underline{\Sigma}^T \underline{\mathbf{U}^T} \\
 &= \underline{\mathbf{U}} \underline{\sum} \underbrace{\frac{\mathbf{I}_P}{\underbrace{\mathbf{P} \times \mathbf{P}}_{n \times n \cdot n \times p}}} \underbrace{\frac{1}{\underline{\Sigma}^T \underline{\Sigma}}} \underbrace{\frac{\mathbf{I}_P^T}{\underbrace{\mathbf{P} \times \mathbf{P}}_{\mathbf{P} \times \mathbf{n} \cdot \mathbf{n} \times \mathbf{p}}}} \underline{\mathbf{U}^T} \quad (4) \\
 &\quad \text{~~~~~} \qquad \text{~~~~~} \qquad \text{~~~~~} \\
 &= \mathbf{n} \times \mathbf{p} \qquad = \mathbf{P} \times \mathbf{p} \qquad = \mathbf{P} \times \mathbf{n}
 \end{aligned}$$

From the dimensions we see that the products add up. Further we can observe what each product will behave as:

$\mathbf{U} \Sigma$ : The rows/columns = 0 in  $\Sigma$  will go to 0 for the product thus we will have a  $n \times p$  matrix containing the prod. of the s.v. and the eigenvectors of the design matrix e.g.  $\sigma_i \underline{\mathbf{u}}_i$ , where  $\underline{\mathbf{U}} = [\underline{\mathbf{u}}_1, \underline{\mathbf{u}}_2 \dots]$

$(\Sigma^T \Sigma)^{-1}$ : will yield a diagonal matrix with  $\frac{1}{\sigma_i^2}$  in the entries.

$\underline{\Sigma}^T \underline{U}^T$ : again the o of  $\Sigma$

makes the matrix product where  
the singular values is one col  
and the transposed  $U$  vectors  
filling the remaining cols

However we can easily see that  
we can cancel matrices in (4):

$$\tilde{y}_{OLS} = \underline{U} \underline{\Sigma}$$

$$\frac{1}{\underline{\Sigma}^T \underline{\Sigma}}$$

$$\underline{\Sigma}^T \underline{U}^T \cancel{Y}^{(4)}$$

$$= \underline{U} \underline{U}^T \cancel{Y}$$

For  $i \neq j$  the inner product = 0  
therefor

$$\begin{aligned}\tilde{y}_{OLS} &= (\underline{u}_0 \underline{u}_0^T + \dots + \underline{u}_p \underline{u}_p^T) \cancel{Y} \\ &= \left[ \sum_{j=1}^p \underline{u}_j \underline{u}_j^T \right] \cancel{Y}\end{aligned}$$



$$\hat{y}_{\text{ridge}} = \underline{\underline{X}} \underline{\underline{\beta}}_{\text{ridge}}$$

$$= \underline{\underline{X}} (\underline{\underline{X}}^T \underline{\underline{X}} + \lambda \underline{\underline{I}})^{-1} \underline{\underline{X}}^T \underline{\underline{y}}$$

$$= \underline{\underline{U}} \underline{\Sigma} \underline{\underline{V}}^T \left( \underline{\underline{V}}^T \underline{\underline{U}}^T \underline{\underline{U}} \underline{\Sigma} \underline{\underline{V}}^T + \lambda \underline{\underline{I}} \right)^{-1}$$

$$(\underline{\underline{U}} \underline{\Sigma} \underline{\underline{V}}^T)^T \underline{\underline{y}}$$

Using  $\underline{\underline{U}}^T \underline{\underline{U}} = \underline{\underline{I}}$  and  $\underline{\underline{V}}^T \underline{\underline{V}} = \underline{\underline{I}}$

$$\hat{y}_{\text{ridge}} = \underline{\underline{U}} \underline{\Sigma} (\underline{\Sigma}^2 + \lambda \underline{\underline{I}})^{-1} \underline{\Sigma} \underline{\underline{U}}^T \underline{\underline{y}}$$

Since  $\underline{\Sigma} \underline{\Sigma}^{-1} = \underline{\underline{I}}$  we get

$$= \underline{\underline{U}} (\underline{\Sigma}^2 + \lambda \underline{\underline{I}})^{-1} \underline{\underline{U}}^T \underline{\underline{y}}$$

Similarly as for OLS  $j \neq i$  yields an inner prod  $= 0$  and we therefore only get one summation

Further  $(\Sigma^2 + \lambda I)^{-1}$  will yield

$$\frac{\sigma_j^2}{\sigma_j^2 + \lambda}$$

giving us

$$\tilde{y}_{\text{ridge}} = \sum_{j=1}^p u_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} u_j^T y$$

The interpretation of this is that the  $\tilde{y}$  is the weighted sum of the  $u_j$ 's with the weight being dictated by the regularization parameter.

A  $\lambda = 0$  will yield  $\tilde{y}_{\text{OLS}}$  since  $\frac{\sigma_j^2}{\sigma_j^2 + 0} = 1$ . However a large  $\lambda$  will decrease the value of the sum since

$$\lim_{\lambda \rightarrow \infty} \frac{\sigma_j^2}{\sigma_j^2 + \lambda} = 0$$

This will stabilize the model.

