

STOR 767, Advanced Machine Learning

Homework 1, Theoretical Part

Zhenghan Fang

May 22, 2019

1. Answer to question 1:

(i)

$$Y = \mathbf{X}\boldsymbol{\beta} + \epsilon \quad (1)$$

$$\text{where } Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_n^T \end{bmatrix}, \text{ and } \epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}.$$

(ii) The least square estimate (LSE) of $\boldsymbol{\beta}$ is

$$\hat{\boldsymbol{\beta}}_{LS} = \arg \min_{\boldsymbol{\beta}} (Y - \mathbf{X}\boldsymbol{\beta})^T (Y - \mathbf{X}\boldsymbol{\beta}) \quad (2)$$

Set the derivative of $(Y - \mathbf{X}\boldsymbol{\beta})^T (Y - \mathbf{X}\boldsymbol{\beta})$ w.r.t $\boldsymbol{\beta}$ to zero, we get

$$\mathbf{X}^T (Y - \mathbf{X}\hat{\boldsymbol{\beta}}_{LS}) = 0 \quad (3)$$

$$\Rightarrow \mathbf{X}^T Y - \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}}_{LS} = 0 \quad (4)$$

$$\Rightarrow \hat{\boldsymbol{\beta}}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y \quad (5)$$

The LSE of σ^2 is

$$\hat{\sigma}_{LS}^2 = \frac{1}{n} (Y - \hat{Y})^T (Y - \hat{Y}) \quad (6)$$

$$= \frac{1}{n} (Y - \mathbf{X}\hat{\boldsymbol{\beta}}_{LS})^T (Y - \mathbf{X}\hat{\boldsymbol{\beta}}_{LS}) \quad (7)$$

$$= \frac{1}{n} (Y - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y)^T (Y - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y) \quad (8)$$

(iii)

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (9)$$

\mathbf{H} represents the linear transform from $\{Y_i\}_{i=1}^n$ to the estimations of $\{Y_i\}_{i=1}^n$ by the linear regression model. \mathbf{H} implies that this transform is determined only by $\{\mathbf{X}_i\}_{i=1}^n$.

(iv)

$$\mathbb{E}(\hat{\boldsymbol{\beta}}_{LS}) = \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y) \quad (10)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(Y) \quad (11)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \quad (12)$$

$$= \boldsymbol{\beta} \quad (13)$$

$$\mathbb{E}(\hat{\sigma}_{LS}^2) = \mathbb{E}\left(\frac{1}{n}(Y - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y)^T (Y - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y)\right) \quad (14)$$

$$= \mathbb{E}\left(\frac{1}{n}(Y - \mathbf{H}Y)^T (Y - \mathbf{H}Y)\right) \quad (15)$$

$$= \frac{1}{n} ((\mathbf{I} - \mathbf{H})\mathbf{X}\boldsymbol{\beta})^T ((\mathbf{I} - \mathbf{H})\mathbf{X}\boldsymbol{\beta}) + \frac{\sigma^2}{n} \text{Tr}((\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H})^T) \quad (16)$$

$$= \frac{\sigma^2}{n} \text{Tr}(\mathbf{I} - \mathbf{H}) \quad (17)$$

$$= \sigma^2 \left(1 - \frac{1}{n} \text{Tr}(\mathbf{H})\right) \quad (18)$$

$$\text{Cov}(\hat{\boldsymbol{\beta}}_{LS}) = \text{Cov}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y) \quad (19)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Cov}(Y) ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T)^T \quad (20)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Cov}(Y) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad (21)$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Cov}(\epsilon) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad (22)$$

$$= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \quad (23)$$

(v)

$$\hat{\boldsymbol{\beta}}_{ML}, \hat{\sigma}_{ML}^2 = \arg \max_{\boldsymbol{\beta}, \sigma^2} \Pr(\mathbf{X}, Y | \boldsymbol{\beta}, \sigma^2) \quad (24)$$

$$= \arg \max_{\boldsymbol{\beta}, \sigma^2} \Pr(Y | \mathbf{X}, \boldsymbol{\beta}, \sigma^2) \quad (25)$$

$$= \arg \max_{\boldsymbol{\beta}, \sigma^2} \prod_{i=1}^n \Pr(Y_i | \mathbf{X}_i, \boldsymbol{\beta}, \sigma^2) \quad (26)$$

$$= \arg \max_{\boldsymbol{\beta}, \sigma^2} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2}{2\sigma^2} \right] \quad (27)$$

$$= \arg \max_{\boldsymbol{\beta}, \sigma^2} \sum_{i=1}^n \left\{ -\frac{1}{2} \log(\sigma^2) - \frac{(Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2}{2\sigma^2} \right\} \quad (28)$$

$$= \arg \max_{\boldsymbol{\beta}, \sigma^2} -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2 \quad (29)$$

Differentiate the object function w.r.t. $\boldsymbol{\beta}$ and σ^2 , and set the derivative to zero, we get

$$\hat{\boldsymbol{\beta}}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y \quad (30)$$

$$\frac{n}{2\hat{\sigma}_{ML}^2} - \frac{\sum_{i=1}^n (Y_i - \mathbf{X}_i^T \hat{\boldsymbol{\beta}}_{ML})^2}{2(\hat{\sigma}_{ML}^2)^2} = 0 \quad (31)$$

$$\Rightarrow \hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \hat{\boldsymbol{\beta}}_{ML})^2 \quad (32)$$

$$\Rightarrow \hat{\sigma}_{ML}^2 = \frac{1}{n} (Y - \mathbf{H}Y)^T (Y - \mathbf{H}Y) \quad (33)$$

Proof of $\hat{\boldsymbol{\beta}}_{ML} \perp \hat{\sigma}_{ML}^2$:

$$\mathbb{E}(\hat{\boldsymbol{\beta}}_{ML}) = \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y) \quad (34)$$

$$= \boldsymbol{\beta} \quad (35)$$

$$\mathbb{E} [\hat{\boldsymbol{\beta}}_{ML} \hat{\sigma}_{ML}^2] = \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y \hat{\sigma}_{ML}^2] \quad (36)$$

$$= \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \boldsymbol{\beta} + \epsilon) \hat{\sigma}_{ML}^2] \quad (37)$$

$$= \mathbb{E} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \hat{\sigma}_{ML}^2 + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \hat{\sigma}_{ML}^2] \quad (38)$$

$$= \mathbb{E} [\boldsymbol{\beta} \hat{\sigma}_{ML}^2 + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \hat{\sigma}_{ML}^2] \quad (39)$$

$$= \boldsymbol{\beta} \mathbb{E} [\hat{\sigma}_{ML}^2] + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E} [\epsilon \hat{\sigma}_{ML}^2] \quad (40)$$

where

$$\mathbb{E} [\epsilon \hat{\sigma}_{ML}^2] = \mathbb{E} [\epsilon \mathbf{X} \boldsymbol{\beta}^T (\mathbf{I} - \mathbf{H}) \mathbf{X} \boldsymbol{\beta} + \epsilon \epsilon^T (\mathbf{I} - \mathbf{H}) \mathbf{X} \boldsymbol{\beta} \quad (41)$$

$$+ \epsilon \mathbf{X} \boldsymbol{\beta}^T (\mathbf{I} - \mathbf{H}) \epsilon + \epsilon \epsilon^T (\mathbf{I} - \mathbf{H}) \epsilon] \quad (42)$$

$$= 0 + \sigma^2 (\mathbf{I} - \mathbf{H}) \mathbf{X} \boldsymbol{\beta} + \sigma^2 (\mathbf{I} - \mathbf{H}) \mathbf{X} \boldsymbol{\beta} + 0 \quad (43)$$

$$= 0 \quad (44)$$

Therefore,

$$\mathbb{E} [\hat{\boldsymbol{\beta}}_{ML} \hat{\sigma}_{ML}^2] = \boldsymbol{\beta} \mathbb{E} [\hat{\sigma}_{ML}^2] \quad (45)$$

$$= \mathbb{E} [\hat{\boldsymbol{\beta}}_{ML}] \mathbb{E} [\hat{\sigma}_{ML}^2] \quad (46)$$

Therefore, $\hat{\boldsymbol{\beta}}_{ML} \perp \hat{\sigma}_{ML}^2$.

(vi) From the calculation above, we get

$$\hat{\boldsymbol{\beta}}_{ML} = \hat{\boldsymbol{\beta}}_{LS} \quad (47)$$

$$\hat{\sigma}_{ML}^2 = \hat{\sigma}_{LS}^2 \quad (48)$$

This suggests that the MLE of $\boldsymbol{\beta}$ and σ^2 under the assumption that error ϵ follows Gaussian distribution yields the LSE of $\boldsymbol{\beta}$ and σ^2 .

(vii) We denote the new independent sample to be predicted as (\mathbf{X}_m, Y_m) ,

because we have defined $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_n^T \end{bmatrix}$ and $Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$.

$$\text{MSE}(\hat{Y}_{LS}) = \mathbb{E} \left(Y_m - \hat{Y}_{LS} \right)^2 \quad (49)$$

$$= \mathbb{E} \left(\mathbf{X}_m^T \boldsymbol{\beta} + \epsilon - \hat{Y}_{LS} \right)^2 \quad (50)$$

$$= \mathbb{E} \left(\mathbf{X}_m^T \boldsymbol{\beta} - \hat{Y}_{LS} \right)^2 + \mathbb{E}(\epsilon^2) \quad (51)$$

$$= \left[\mathbb{E} \left(\mathbf{X}_m^T \boldsymbol{\beta} - \hat{Y}_{LS} \right) \right]^2 + \text{Cov}(\mathbf{X}_m^T \boldsymbol{\beta} - \hat{Y}_{LS}) + \sigma^2 \quad (52)$$

$$= \left[\mathbb{E} \left(\mathbf{X}_m^T \boldsymbol{\beta} - \hat{Y}_{LS} \right) \right]^2 + \text{Cov}(\hat{Y}_{LS}) + \sigma^2 \quad (53)$$

where

$$\mathbb{E} \left(\mathbf{X}_m^T \boldsymbol{\beta} - \hat{Y}_{LS} \right) = \mathbf{X}_m^T \boldsymbol{\beta} - \mathbb{E} \left(\hat{Y}_{LS} \right) \quad (54)$$

$$= \mathbf{X}_m^T \boldsymbol{\beta} - \mathbb{E} \left(\mathbf{X}_m^T \hat{\boldsymbol{\beta}}_{LS} \right) \quad (55)$$

$$= 0 \quad (56)$$

and

$$\text{Cov}(\hat{Y}_{LS}) = \mathbf{X}_m^T \text{Cov}(\hat{\boldsymbol{\beta}}_{LS}) \mathbf{X}_m \quad (57)$$

$$= \sigma^2 \mathbf{X}_m^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_m \quad (58)$$

Therefore,

$$\text{MSE}(\hat{Y}_{LS}) = \sigma^2 \left(\mathbf{X}_m^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_m + 1 \right) \quad (59)$$

2. Answer to question 2:

Let

$$\mathbf{X}_1 = \begin{bmatrix} \mathbf{X}_{11}^T \\ \vdots \\ \mathbf{X}_{n1}^T \end{bmatrix} \quad (60)$$

and split the new independent sample \mathbf{X}_m as

$$\mathbf{X}_m = \begin{bmatrix} \mathbf{X}_{m1} \\ \mathbf{X}_{m2} \end{bmatrix} \quad (61)$$

where $\mathbf{X}_{m1} \in \mathbb{R}^{p_1}$, $\mathbf{X}_{m2} \in \mathbb{R}^{p_2}$.

(i)

$$\mathbf{MSE}(\hat{Y}_{LS}^{(1)}) = \sigma^2 (\mathbf{X}_{m1}^T (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_{m1} + 1) \quad (62)$$

$$\mathbf{MSE}(\hat{Y}_{LS}^{(1,2)}) = \sigma^2 (\mathbf{X}_m^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_m + 1) \quad (63)$$

(ii) For $\mathbf{MSE}(\hat{Y}_{LS}^{(1)})$, the bias of $\hat{Y}_{LS}^{(1)}$ is not zero:

$$\mathbb{E}(\mathbf{X}_m^T \boldsymbol{\beta} - \hat{Y}_{LS}^{(1)}) = \mathbf{X}_m^T \boldsymbol{\beta} - \mathbb{E}(\hat{Y}_{LS}^{(1)}) \quad (64)$$

$$= \mathbf{X}_m^T \boldsymbol{\beta} - \mathbf{X}_{m1}^T \mathbb{E}(\hat{\boldsymbol{\beta}}_{LS}^{(1)}) \quad (65)$$

where

$$\mathbb{E}(\hat{\boldsymbol{\beta}}_{LS}^{(1)}) = \mathbb{E}((\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T Y) \quad (66)$$

$$= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X} \boldsymbol{\beta} \quad (67)$$

Therefore,

$$\mathbf{MSE}(\hat{Y}_{LS}^{(1)}) = (\mathbf{X}_m^T \boldsymbol{\beta} - \mathbf{X}_{m1}^T (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X} \boldsymbol{\beta})^2 \quad (68)$$

$$+ \sigma^2 (\mathbf{X}_{m1}^T (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_{m1} + 1) \quad (69)$$

Besides,

$$\mathbf{MSE}(\hat{Y}_{LS}^{(1,2)}) = \sigma^2 (\mathbf{X}_m^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_m + 1) \quad (70)$$

Therefore, the condition under which

$$\mathbf{MSE}(\hat{Y}_{LS}^{(1,2)}) \leq \mathbf{MSE}(\hat{Y}_{LS}^{(1)}) \quad (71)$$

is

$$\sigma^2 \mathbf{X}_m^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_m \leq \quad (72)$$

$$(\mathbf{X}_m^T \boldsymbol{\beta} - \mathbf{X}_{m1}^T (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X} \boldsymbol{\beta})^2 + \sigma^2 \mathbf{X}_{m1}^T (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_{m1} \quad (73)$$

3. Answer to question 3:

(i) Let

$$\tilde{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b} \quad (74)$$

Because $\mathbf{\Sigma}$ is real symmetric, it can be decomposed as

$$\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T \quad (75)$$

where \mathbf{Q} is an orthogonal matrix whose columns are the eigenvectors of $\mathbf{\Sigma}$, and $\mathbf{\Lambda}$ is a diagonal matrix whose entries are the eigenvalues of $\mathbf{\Sigma}$.

Let

$$\mathbf{A} = (\mathbf{Q}\mathbf{\Lambda}^{-\frac{1}{2}})^T \quad (76)$$

$$\mathbf{b} = -\mathbf{A}\boldsymbol{\mu} \quad (77)$$

Then

$$\text{cov}(\tilde{\mathbf{x}}) = \mathbf{A}\text{cov}(\mathbf{x})\mathbf{A}^T \quad (78)$$

$$= \mathbf{A}\mathbf{\Sigma}\mathbf{A}^T \quad (79)$$

$$= (\mathbf{Q}\mathbf{\Lambda}^{-\frac{1}{2}})^T \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T (\mathbf{Q}\mathbf{\Lambda}^{-\frac{1}{2}}) \quad (80)$$

$$= \mathbf{I} \quad (\mathbf{Q}^T \mathbf{Q} = \mathbf{I}) \quad (81)$$

$$E(\tilde{\mathbf{x}}) = \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \quad (82)$$

$$= \mathbf{0} \quad (83)$$

(ii) a) If $\mathbf{y}^T \mathbf{x} \neq 0$,
 $\mathbf{x}\mathbf{y}^T$ has one nonzero eigenvalue

$$\lambda_1 = \mathbf{y}^T \mathbf{x} \quad (84)$$

with corresponding eigenvector

$$\mathbf{v}_1 = \mathbf{x} \quad (85)$$

and $(n - 1)$ zero eigenvalues with corresponding eigenvectors

$$\{\mathbf{v} \mid \mathbf{v}^T \mathbf{y} = 0\} \quad (86)$$

- b) If $\mathbf{y}^T \mathbf{x} = 0$ and $\mathbf{y} \neq \mathbf{0}$ and $\mathbf{x} \neq \mathbf{0}$,
 $\mathbf{x}\mathbf{y}^T$ has $(n-1)$ zero eigenvalues with corresponding eigenvectors

$$\{\mathbf{v} \mid \mathbf{v}^T \mathbf{y} = 0\} \quad (88)$$

- c) If $\mathbf{y} = \mathbf{0}$ or $\mathbf{x} = \mathbf{0}$,

$$\mathbf{x}\mathbf{y}^T = \mathbf{0} \quad (89)$$

- (iii) Because \mathbf{A} is real symmetric,

$$\sum_{i=1}^p \sum_{j=1}^p a_{ij}^2 = \text{Tr}(\mathbf{A}\mathbf{A}) \quad (90)$$

and the eigendecomposition of \mathbf{A} yields

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T \quad (91)$$

where \mathbf{Q} is orthogonal and $\mathbf{\Lambda}$ is a diagonal matrix whose entries are eigenvalues of \mathbf{A} , i.e. $\lambda_1, \dots, \lambda_p$. Then,

$$\mathbf{A}\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T \quad (92)$$

$$= \mathbf{Q}\mathbf{\Lambda}\mathbf{\Lambda}\mathbf{Q}^T \quad (93)$$

$$= [\mathbf{q}_1 \quad \dots \quad \mathbf{q}_p] \mathbf{\Lambda}\mathbf{\Lambda} \begin{bmatrix} \mathbf{q}_1^T \\ \vdots \\ \mathbf{q}_p^T \end{bmatrix} \quad (94)$$

$$= \sum_{i=1}^p \lambda_i^2 \mathbf{q}_i \mathbf{q}_i^T \quad (95)$$

where $\mathbf{q}_1, \dots, \mathbf{q}_p$ are the columns of \mathbf{Q} . Therefore,

$$\sum_{i=1}^p \sum_{j=1}^p a_{ij}^2 = \text{Tr}(\mathbf{A}\mathbf{A}) = \text{Tr} \left(\sum_{i=1}^p \lambda_i^2 \mathbf{q}_i \mathbf{q}_i^T \right) \quad (96)$$

$$= \sum_{i=1}^p \lambda_i^2 \text{Tr}(\mathbf{q}_i \mathbf{q}_i^T) \quad (97)$$

$$= \sum_{i=1}^p \lambda_i^2 \quad (\mathbf{q}_i \text{'s are orthonormal vectors}) \quad (98)$$