

# ASSIGNMENT 3

Student ID : 7205329

## Exercise 1.

- Find the gradient  $\nabla_{\underline{\theta}} \mathcal{L}_{A, \Lambda}(\underline{\theta})$

Firstly, we rewrite  $\mathcal{L}_{A, \Lambda}$ ,

$$\begin{aligned}\mathcal{L}_{A, \Lambda}(\underline{\theta}) &= \|y - X\underline{\theta}\|_A^2 + \|\underline{\theta}\|_\Lambda^2 \\ &= \langle y - X\underline{\theta}, y - X\underline{\theta} \rangle_A + \langle \underline{\theta}, \underline{\theta} \rangle_\Lambda \\ &= (y - X\underline{\theta})^T A (y - X\underline{\theta}) + \underline{\theta}^T \Lambda \underline{\theta} \\ &= (y - X\underline{\theta})^T (A_y - A X \underline{\theta}) + \underline{\theta}^T \Lambda \underline{\theta} \\ &= \underbrace{y^T (A_y)}_A - \underbrace{y^T (A X \underline{\theta})}_B - \underbrace{(X \underline{\theta})^T (A_y)}_C + \underbrace{(X \underline{\theta})^T (A X \underline{\theta})}_D + \underline{\theta}^T \Lambda \underline{\theta}\end{aligned}$$

A

B

C

D

E

We will calculate the derivatives individual and sum them up later using the sum rule. The following calculations will use results from section 5.5, lecture "Vector Calculus".

Part A.

$$\frac{\partial}{\partial \underline{\theta}} [y^T (A_y)] = \underline{\theta}^T.$$

Part B.

$$\frac{\partial}{\partial \underline{\theta}} [y^T (A X \underline{\theta})] = \frac{\partial}{\partial \underline{\theta}} \left[ ((y^T A X)^T)^T \underline{\theta} \right] = \underbrace{y^T A X}_\text{using the result } \frac{\partial (\underline{\theta}^T \underline{x})}{\partial \underline{\theta}} = \underline{\theta}^T$$

Part C.

$$\begin{aligned}\frac{\partial}{\partial \theta} [(x\theta)^T (A y)] &= \frac{\partial}{\partial \theta} [(A y)^T (x\theta)] \\&= \frac{\partial}{\partial \theta} [y^T A^T x\theta] \\&= \frac{\partial}{\partial \theta} [y^T A x\theta] \quad \text{since } A \text{ is symmetric} \quad A^T = A \\&= y^T A x \quad (\text{from part b above})\end{aligned}$$

Part D.

$$\begin{aligned}\frac{\partial}{\partial \theta} [(x\theta)^T A x\theta] &= \frac{\partial}{\partial \theta} [\theta^T (x^T A x) \theta] \quad \downarrow \\&\quad \text{using result } \frac{\partial(x^T B x)}{\partial x} = x^T(B + B^T) \\&= \theta^T [x^T A x + x^T A^T x] = 2\theta^T x^T A x.\end{aligned}$$

Part E.

$$\begin{aligned}\frac{\partial}{\partial \theta} [\theta^T \Lambda \theta] &= \theta^T (\Lambda + \Lambda^T) = \underbrace{\theta^T (\Lambda + \Lambda)}_{\text{using result } \frac{\partial(x^T B x)}{\partial x} = x^T(B + B^T), \Lambda \text{ is symmetric}} = 2\theta^T \Lambda.\end{aligned}$$

Finally,

$$\begin{aligned}\nabla_{\theta} L_{A, \Lambda}(\theta) &= 0^T - y^T A x - y^T A x + 2\theta^T x^T A x + 2\theta^T \Lambda \\&= -2y^T A x + 2\theta^T x^T A x + 2\theta^T \Lambda \\&= -2(y^T - \theta^T x^T) A x + 2\theta^T \Lambda \quad \begin{aligned}y^T - \theta^T x^T \\= y^T - (x\theta)^T \\= (y - x\theta)^T\end{aligned} \\&= -2(y - x\theta)^T A x + 2\theta^T \Lambda\end{aligned}$$

2. let  $\nabla_{\underline{\theta}} \mathcal{L}_{A, \Lambda}(\underline{\theta}) = 0$ , we have

$$-\cancel{2} (\underline{y} - \underline{x}\underline{\theta})^T A X + \cancel{2}\underline{\theta}^T \Lambda = 0$$

$$\Leftrightarrow -\cancel{2} \underline{y}^T A X + \cancel{2}\underline{\theta}^T X^T A X + \cancel{2}\underline{\theta}^T \Lambda = 0$$

$$\Leftrightarrow \cancel{2}\underline{\theta}^T (X^T A X + \Lambda) = \cancel{2} \underline{y}^T A X$$

$$\Leftrightarrow \underline{\theta}^T = \underline{y}^T A X (X^T A X + \Lambda)^{-1}$$

$$\Leftrightarrow \underline{\theta} = [(X^T A X + \Lambda)^T]^{-1} (\underline{y}^T A X)^T$$

$$\Leftrightarrow \underline{\theta} = (X^T A X + \Lambda)^{-1} X^T A \underline{y}$$

We have  $X^T A X = \langle X, X \rangle_A$  is asymmetric and positive definite  
(property of inner product).

Since the sum of a symmetric and asymmetric matrix is asymmetric,  
and the sum of two positive definite matrix is positive definite.

$\Rightarrow X^T A X + \Lambda$  is an asymmetric positive definite  $D \times D$  matrix and  
thus it is invertible, making our  $\underline{\theta}$  defined.

3. Substituting  $A = I$  and  $\Lambda = \lambda I$ , we have

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

4. For standard regularisation term  $\lambda \|\underline{\theta}\|_2^2$ , we are applying the same weight of  $\lambda$  to all  $\theta_i$ 's.

$$\begin{aligned}\lambda \|\underline{\theta}\|_2^2 &= \lambda (\theta_1^2 + \theta_2^2 + \dots + \theta_D^2) \\ &= \lambda \theta_1^2 + \lambda \theta_2^2 + \dots + \lambda \theta_D^2\end{aligned}$$

On the other hand, using  $\|\underline{\theta}\|_\Lambda^2$  let us apply different weights to different parameters  $\theta_i$ . This way, we can specify the sensitivity of  $\hat{y}$  to each  $x_i$ .

For example, let  $\Lambda = \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0.1 & 0 \\ 0 & 0 & 0.3 \end{bmatrix}$  and  $\underline{\theta} = [\theta_1 \ \theta_2 \ \theta_3]^T$

$$\Rightarrow \underline{\theta}^T \Lambda \underline{\theta} = [\theta_1 \ \theta_2 \ \theta_3] \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0.1 & 0 \\ 0 & 0 & 0.3 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$= [0.5 \theta_1 \quad 0.1 \theta_2 \quad 0.3 \theta_3] \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$= 0.5 \theta_1^2 + 0.1 \theta_2^2 + 0.3 \theta_3^2.$$

## Exercise 2.

1. We have

$$L_2(\mu, D) = \sum_{i=1}^n (\mu - x_i)^2$$

$$= (\mu - x_1)^2 + (\mu - x_2)^2 + \dots + (\mu - x_n)^2$$

$$= \mu^2 - 2\mu x_1 + x_1^2 + \mu^2 - 2\mu x_2 + \dots + \mu^2 - 2\mu x_n + x_n^2$$

$$= n\mu^2 - 2\mu \sum x_i + \sum_{i=1}^n x_i^2$$

$$\Rightarrow \frac{dL_2}{d\mu} = \frac{d(n\mu^2 - 2\mu \sum x_i + \sum x_i^2)}{d\mu}$$
$$= 2n\mu - 2\sum x_i$$

$$\begin{aligned} \text{Setting } \frac{dL_2}{d\mu} &= 0 \quad \text{we have} \quad 2n\mu - 2\sum x_i &= 0 \\ &\quad 2n\mu &= 2\sum x_i \\ &\quad \mu &= \frac{1}{n} \sum x_i \end{aligned}$$

Therefore, the optimal choice for  $\mu$  is the mean.

d. We have  $L_1 = \sum_{i=1}^N |\mu - x_i|$ , we consider 2 cases of dataset D.

i) D has an odd number of elements.

Suppose we have  $\mu$  that minimises  $L_1(\mu, D)$  (★), with  $n_l = \text{number of points } x_i < \mu$  and  $n_r = \text{number of points } x_i > \mu$ .

• If  $n_r < n_l$ ,

For example, let  $D = \{0, 1, 2, 3, 5\}$  and  $\mu = 3 \Rightarrow n_l = 4, n_r = 1$ .

$$\begin{aligned} \Rightarrow L_1 &= |3-0| + |3-1| + |3-2| + |3-3| + |3-5| \\ &= \underbrace{3 + 2 + 1}_{S_l} + 0 + \underbrace{2}_{S_r} \end{aligned}$$

now if we move  $\mu$  to be  $\mu' = 2 \Rightarrow \Delta = |\mu' - \mu| = |2-3| = 1$ .

$$\begin{aligned} \Rightarrow L'_1 &= |2-0| + |2-1| + |2-2| + |2-3| + |2-5| \\ &= \underbrace{2 + 1 + 0}_{\text{each decreases by } \Delta} + \underbrace{1 + 3}_{\text{each increases by } \Delta} \end{aligned}$$

$$S'_l = S_l - n_l \cdot \Delta$$

$$S'_r = S_r + n_r \cdot \Delta$$

$$\begin{aligned} \Rightarrow L'_1 &= S'_l + S'_r = S_l - n_l \Delta + S_r + n_r \Delta \\ &= (S_l + S_r) + \Delta(n_r - n_l) \\ &= L_1 + \Delta(n_r - n_l) \end{aligned}$$

Since  $n_r < n_l$ , the net change is negative and we can further minimise  $L_1$ , which contradicts our assumption (★).

- If  $n_\ell < n_r$ : Let  $D = \{0, 1, 2, 3, 5\}$  and  $\mu = 1$  to be the optimal choice  $\Rightarrow n_\ell = 1$  and  $n_r = 3$ .

$$\text{With } \mu = 1, L_1 = |1 - 0| + |2 - 1| + |1 - 2| + |1 - 3| + |1 - 5|$$

$$= \underbrace{|1|}_S_\ell \quad 0 \quad \underbrace{|1|}_S_r \quad 2 \quad 4$$

if we move  $\mu$  to be  $\mu' = 2$ ,  $\Delta = |\mu' - \mu| = |2 - 1| = 1$

$$\Rightarrow L_1' = |2 - 0| + |2 - 1| + |2 - 2| + |2 - 3| + |2 - 5|$$

$$= \underbrace{|2|}_S_\ell + \underbrace{|1|}_S_r + \underbrace{|0|}_S_r + \underbrace{|1|}_S_r + \underbrace{|3|}_S_r$$

increases by  $\Delta$

$$S_\ell' = S_\ell + n_\ell \cdot \Delta$$

each decreases by  $\Delta$

$$S_r' = S_r - n_r \cdot \Delta$$

$$L_1' = S_\ell' + S_r' = S_\ell + S_r + \Delta(n_\ell - n_r)$$

$$= L_1 + \Delta(n_\ell - n_r)$$

since  $n_\ell < n_r$ , the net change is also negative and that we can further minimise  $L_1$ .

Thus, as  $n_\ell \neq n_r$ , if we move  $\mu$  to the left (decrease  $\mu$  by a value of  $\Delta$ ) or to the right (increase by  $\Delta$ ), then we will get a different  $\mu$  that can further minimises  $L_1$ .

It follows the minimum  $L_1$  can be reached with a  $\mu$  value such that  $n_\ell = n_r$ . For  $D$  with an odd number of elements, it is easy to point out that the optimal choice for  $\mu$  is the median of  $D$  (by definition).

- 2) For  $D$  with an even number of elements, we start off with any number between the two middle numbers.

For example,  $D = \{1, 2, 4, 6\}$

If we choose  $\mu$  to be any number between 2 and 4 (e.g.  $\mu = 3$ ) then  $n_r = n_l$ , and thus will minimise  $L_1$ .

However, we also find that the endpoints of this interval also minimise  $L_1$ .

Indeed,  $x_{n/2}$  and  $x_{n/2+1}$  divide the set  $D$  into two smaller groups, each with  $n/2$  elements.

$$D = \{x_1, \dots, x_{n/2}, x_{n/2+1}, \dots, x_n\}$$

$$\text{Let } \mu = x_{n/2}$$

$$\Rightarrow L_1 = |x_{n/2} - x_1| + \dots + |x_{n/2} - x_{n/2}| + |x_{n/2} - x_{n/2+1}| + \dots + |x_{n/2} - x_n|$$

$$\text{if we move } \mu \text{ to be } \mu' = x_{n/2+1}, \Delta = |x_{n/2+1} - x_{n/2}|$$

$$\Rightarrow L'_1 = |x_{n/2+1} - x_1| + \dots + \underbrace{|x_{n/2+1} - x_{n/2}|}_{\text{each increases by } \Delta} + |x_{n/2+1} - x_{n/2+1}| + \dots + \underbrace{|x_{n/2+1} - x_n|}_{\text{each decreases by } \Delta}$$

$$\Rightarrow s'_l = s_l + \Delta \cdot (n/2)$$

$$\Rightarrow s'_r = s_r - \Delta \cdot (n/2)$$

$$\Rightarrow L'_1 = s'_l + s'_r = s_l + \Delta \cdot (n/2) + s_r - \Delta \cdot (n/2) = s_l + s_r = L_1.$$

Finally, we found that  $L_i$  is minimised when

- $\mu$  is the median (when  $D$  has an odd number of elements), and
- $\mu$  is any number between  $x_{n/2}$  and  $x_{n/2+1}$  (when  $D$  has an even number of elements).

⇒ The best representative  $\mu$  for any data set is the median of that dataset.

3. For the  $L_\infty$  loss function, we consider 3 cases of  $\mu$ :

- $\mu \leq x_1$ , then  $L_\infty(\mu, D) = \max_{1 \leq i \leq n} |\mu - x_i| = |\mu - x_n|$  since  $x_n$  has the furthest distance from  $\mu$ .

$$\Rightarrow \frac{dL_\infty}{d\mu} = \frac{d(|\mu - x_n|)}{d\mu} = \frac{(\mu - x_n)\mu}{|\mu - x_n|}$$

Setting  $\frac{dL_\infty}{d\mu} = 0$ , we have

$$\Leftrightarrow \frac{(\mu - x_n)\mu}{|\mu - x_n|} = 0$$

$$\Leftrightarrow (\mu - x_n)\mu = 0$$

$$\Leftrightarrow \mu^2 - \mu x_n = 0 \quad \Leftrightarrow \mu = x_n.$$

$\Rightarrow$  In this case,  $\mu = x_n$  is the best representative for  $D$ .

- $\mu \geq x_n$ , then  $L_\infty(\mu, D) = \max_{1 \leq i \leq n} |\mu - x_i| = |\mu - x_1|$  since  $x_1$  has the furthest distance from  $\mu$ .

$$\Rightarrow \frac{dL_\infty}{d\mu} = \frac{d(|\mu - x_1|)}{d\mu} = \frac{(\mu - x_1)\mu}{|\mu - x_1|}$$

Setting  $\frac{dL_\infty}{d\mu} = 0$ , we have

$$\Leftrightarrow \frac{(\mu - x_1)\mu}{|\mu - x_1|} = 0$$

$$\Leftrightarrow (\mu - x_1)\mu = 0$$

$$\Leftrightarrow \mu^2 - \mu x_1 = 0 \quad \Leftrightarrow \mu = x_1.$$

Using the same argument as above, we found in the case  $\mu > x_n$ ,  $\mu = x_1$  is the optimal choice for minimising  $L_\infty$ .

.  $x_1 < \mu < x_n$ , then  $l_{\infty}(\mu, D) = |\mu - x_1|$  or  $|\mu - x_n|$  depends on  $\mu$  is closer to  $x_n$  or to  $x_1$ , respectively.

→ Using the same argument, the best representative for this case can be either  $\mu = x_1$  or  $\mu = x_n$ .

Finally, for all 3 cases, we found that  $\mu$  can only be either  $\mu = x_1$  or  $\mu = x_n$ .  
Thus, the best representative  $\mu$  for any data set is the average of  $x_1$  and  $x_n$ .

$$\Rightarrow \mu = \frac{x_1 + x_n}{2}.$$

4.

$\mu$	$L_1$	$L_2$	$L_\infty$
$D_1$	2	2.2	2.5
$D_2$	2	21.2	50
$D_3$	30	35	50

- For  $L_1(\mu, D) = \sum_{i=1}^N |\mu - x_i|$ , we see that this loss function is robust and not too sensitive to outliers (comparing to  $L_2$ ).

Indeed,

$$D_2 : L_1 = |2 - 0| + |2 - 1| + |2 - 2| + |2 - 3| + |2 - 100| \\ = 2 + 1 + 0 + 1 + 98 = 102$$

- For  $L_2(\mu, D) = \sum (\mu - x_i)^2$ , we see that this loss function is very sensitive to outliers as it can easily be skewed towards larger values (large squared distances).

$$D_2 : L_2 = (21.2 - 0)^2 + (21.2 - 1)^2 + (21.2 - 2)^2 + (21.2 - 3)^2 + (21.2 - 100)^2 \\ = 449.44 + 408.04 + 368.64 + 331.24 + 6209.44 \\ = 7766.8$$

- For  $L_\infty = \max_{1 \leq i \leq n} |\mu - x_i|$ , we see that this loss function is also not too sensitive to outliers.

$$D_2. L_\infty = \max(|2.5 - 0|, |2.5 - 1|, |2.5 - 2|, |2.5 - 3|, |2.5 - 100|) \\ = \max(2.5, 1.5, 0.5, 0.5, 97.5) \\ = 97.5.$$