

Assignment - 3

Name - Biswajit Karmakar

SR: - 21055

Question: 1

(I) The simple linear regression model is given by

$$\hat{y} = \hat{w}_0 + \sum_{j=1}^M x_j \hat{w}_j \quad \text{--- (1)}$$

For a given set of outputs y in the training set we minimize the cost function $L(\hat{w}) = \|y - x\hat{w}\|^2$ such that we get best w .

where x = feature matrix

\hat{w} = coefficient vector
 $[\hat{w}_0 \ \hat{w}_1 \ \dots \ \hat{w}_M]^T$

$$\begin{aligned} \text{Now } L(\hat{w}) &= \|y - x\hat{w}\|^2 \\ &= (y - x\hat{w})^T (y - x\hat{w}) \\ &= (y^T - \hat{w}^T x^T) (y - x\hat{w}) \\ &= y^T y - y^T x \hat{w} - \hat{w}^T x^T y - \hat{w}^T x^T x \hat{w} \end{aligned}$$

we minimize $L(\hat{w})$ with respect to \hat{w} ,
i.e., we set the derivative w.r.t \hat{w} is 0.

$$\nabla L(\hat{w}) = 0 \Rightarrow 0 - x^T y - x^T y - 2x^T x \hat{w} = 0$$

[$x^T x$ is positive definite, symmetric]
Also $\hat{w}^T x^T y = y^T x \hat{w}$ is a scalar

$$\Rightarrow 2x^T y = 2x^T x \hat{w}$$

$$\Rightarrow \hat{w} = (x^T x)^{-1} x^T y$$

[Here x is full column rank matrix, so $x^T x$ is invertible]

Hence proved.

(II) Now we consider the linear regression in one variable $\hat{y} = \hat{w}_0 + x \hat{w}_1$

Suppose the mean squared error be denoted by S_r i.e.,

$$S_r = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}^i)^2$$

$$= \frac{1}{N} \sum_{i=1}^N (y^i - \hat{w}_0 - x^i \hat{w}_1)^2$$

where (x^i, y^i) are i th data point for $i=1, 2, \dots, N$.

Now Differentiating w.r.t \hat{w}_0 we get

$$\frac{\partial S_r}{\partial \hat{w}_0} = \frac{1}{N} 2 \sum_{i=1}^N (y^i - \hat{w}_0 - x^i \hat{w}_1) (-1) = 0$$

(We minimizing S_r w.r.t \hat{w}_0)

$$\Rightarrow \sum_{i=1}^N (y^i - \hat{w}_0 - x^i \hat{w}_1) = 0$$

$$\Rightarrow \boxed{\sum_{i=1}^N y_i = \sum_{i=1}^N w_0 + \sum_{i=1}^N x^i \hat{w}_1} \quad \text{--- (1)}$$

Differentiating w.r.t \hat{w}_1 ,

$$\frac{\partial S_r}{\partial \hat{w}_1} = \frac{2}{N} \sum_{i=1}^N (y^i - \hat{w}_0 - x^i \hat{w}_1) (-x^i) = 0$$

[minimizing w.r.t \hat{w}_1]

$$\Rightarrow \boxed{\sum_{i=1}^N x^i y^i = \sum_{i=1}^N x^i \hat{w}_0 + \sum_{i=1}^N (x^i)^2 \hat{w}_1} \quad \text{--- (2)}$$

These two equations are the normal equations with two variable.

From equation (1)

$$N \hat{w}_0 = \sum_{i=1}^N y^i - \sum_{i=1}^N x^i \hat{w}_1$$

$$\Rightarrow \hat{w}_0 = \frac{1}{N} \left(\sum_{i=1}^N y^i - \hat{w}_1 \sum_{i=1}^N x^i \right) \quad \text{--- (3)}$$

Substituting ③ in ②, we get

$$\begin{aligned} & \left(\frac{1}{N} \sum_{i=1}^N y_i - \hat{w}_1 \frac{1}{N} \sum_{i=1}^N x_i \right) \sum_{i=1}^N x_i \\ & + \hat{w}_1 \sum_{i=1}^N (x_i)^2 = \sum_{i=1}^N x_i y_i \\ \Rightarrow & \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right) - \hat{w}_1 \left(\sum_{i=1}^N x_i \right)^2 + N \hat{w}_1 \sum_{i=1}^N (x_i)^2 = N \sum_{i=1}^N x_i y_i \\ \Rightarrow & \hat{w}_1 \left[N \sum_{i=1}^N (x_i)^2 - \left(\sum_{i=1}^N x_i \right)^2 \right] = N \sum_{i=1}^N x_i y_i - \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right) \end{aligned}$$

$$\Rightarrow \hat{w}_1 = \frac{N \sum_{i=1}^N x_i y_i - \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right)}{N \sum_{i=1}^N (x_i)^2 - \left(\sum_{i=1}^N x_i \right)^2}$$

Hence derived.

Now we substitute \hat{w}_1 in ③

$$\begin{aligned} \hat{w}_0 &= \frac{1}{N} \left[\sum_{i=1}^N y_i - \frac{N \sum_{i=1}^N x_i y_i - \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right)}{N \sum_{i=1}^N (x_i)^2 - \left(\sum_{i=1}^N x_i \right)^2} \sum_{i=1}^N x_i \right] \\ &= \frac{1}{N} \left[\sum_{i=1}^N (x_i)^2 \sum_{i=1}^N y_i - \left(\sum_{i=1}^N x_i \right)^2 \sum_{i=1}^N y_i - N \sum_{i=1}^N x_i y_i \sum_{i=1}^N x_i + \left(\sum_{i=1}^N x_i \right)^2 \sum_{i=1}^N y_i \right] \\ &= \frac{\sum_{i=1}^N (x_i)^2 \sum_{i=1}^N y_i - \sum_{i=1}^N x_i y_i \sum_{i=1}^N x_i}{N \sum_{i=1}^N (x_i)^2 - \left(\sum_{i=1}^N x_i \right)^2} \end{aligned}$$

Question 2

②

(a) For ridge regression, the ridge co-efficient minimize a penalised residual sum of squares

$$Z(w_0, w_j) = \sum_{i=1}^N (y^i - w_0 - \sum_{j=1}^M x_j^i w_j)^2 + \lambda \sum_{j=1}^M w_j^2$$

Now with the vector notation, this can be written as

$$Z(w_0, w) = (y - xw - w_0 \mathbf{1})^T (y - xw - w_0 \mathbf{1}) + \lambda w^T w$$

$$\text{where } y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N, \quad w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{pmatrix} \in \mathbb{R}^M \quad \text{--- ①}$$

$$\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^N, \quad x = \begin{pmatrix} x_1^1 & x_2^1 & \dots & x_M^1 \\ \vdots & \vdots & & \vdots \\ x_1^N & x_2^N & \dots & x_M^N \end{pmatrix} \in \mathbb{R}^{N \times M}$$

We will minimize Z with respect to coefficient
so, $\nabla_{w_0} Z(w_0, w) = 0$ and $\nabla_w Z(w_0, w) = 0$

From 1st ~~result~~ condition

$$\nabla_{w_0} Z = 0 \Rightarrow 2(y - xw - w_0 \mathbf{1})^T \cdot \mathbf{1}_{N \times 1} = 0$$

$$\Rightarrow \sum_{i=1}^N y^i - \sum_{i=1}^N x_j^i w_j - w_0 N = 0$$

$$\Rightarrow w_0 = \frac{1}{N} \left(\sum_{i=1}^N y^i - w_j \sum_{i=1}^N x_j^i \right)$$

$$\Rightarrow w_0 = \bar{y} - w_j \bar{x}_j \quad \text{--- ②}$$

The intercept $\hat{w}_0 = \bar{y}$ (Given that $\bar{x} = 0$)

Now From the 2nd ~~result~~ condition

$$\nabla_w Z = 0 \Rightarrow 2(y - xw - w_0 \mathbf{1})^T \cdot (-x) + 2\lambda w^T = 0$$

$$\Rightarrow (y - xw - w_0 \mathbf{1})^T \cdot x = \lambda w^T$$

$$\Rightarrow y^T x - w^T x^T x - w_0 \mathbf{1}^T \cdot x = \lambda w^T \quad \text{--- ③}$$

④

$$= W_0 \cdot (1, 1, \dots, 1)_{1 \times N} \begin{pmatrix} x_1^1 & x_2^1 & \dots & x_M^1 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^N & x_2^N & \dots & x_M^N \end{pmatrix}$$

$$= \omega_0 N \bar{x} = 0 \quad (\because \bar{x} = 0)$$

for the matrices A, B

$$[(AB)^T]^T = B^T A^T$$

$$[(AB)^T]^T = B^T A^T$$

$$A^T + B^T = (A+B)^T]$$

$$A^T + B^T = (A+B)^T]$$

$$(x \in \mathbb{R}^{N \times M}, y \in \mathbb{R}^N)$$

$$(x \in \mathbb{R}^{N \times M}, y \in \mathbb{R}^N)$$

So, $X^T X$ has non zero eigen value.

$$\rightarrow \text{inv} (X^T X + \lambda I)^{-1} \text{ exists.}$$

Hence proved

$$\hat{y} = X \hat{w}_{\text{ridge}}$$

$$\hat{y} = X \hat{w}_{\text{ridge}}$$

Where $U \in \mathbb{R}^{N \times D}$ orthogonal matrix with column as u_j ($j=1, \dots, D$)
 $S \in \mathbb{R}^{D \times D}$ Diagonal matrix with $U^T U = I$
 With diagonal entries s_1, s_2, \dots, s_D .

Where $U \in \mathbb{R}^{N \times D}$ orthogonal matrix with column as u_j ($j=1, \dots, D$)
 $S \in \mathbb{R}^{D \times D}$ Diagonal matrix with $U^T U = I$
 With diagonal entries s_1, s_2, \dots, s_D .

$V \in \mathbb{R}^{D \times D}$ orthogonal matrix ($\because V^T V = I$)

(5)

$$\begin{aligned} \text{Thus } \hat{y} &= X (X^T X + \lambda I)^{-1} X^T y \\ &= U S V^T [(U S V^T)^T U S V^T + \lambda I]^{-1} (U S V^T)^T y \\ &= U S V^T [V S^T U^T U S V^T + \lambda I]^{-1} V S U^T y \quad \left[\begin{array}{l} S^T = S \\ \text{as } S \\ \text{diagonal} \end{array} \right] \\ &= U S V^T [V S^2 V^T + \lambda I]^{-1} V S U^T y \quad [U^T U = I_{D \times D}] \\ &= U S V^T [(V S^2 + \lambda (V^T)^{-1}) V^T]^{-1} V S U^T y \\ &= U S V^T (V^T)^{-1} [V S^2 + \lambda (V^T)^{-1}]^{-1} V S U^T y \quad \text{for two matrices } (AB)^{-1} = B^{-1} A^{-1} \\ &= U S [V (S^2 + \lambda V^{-1} (V^T)^{-1})]^{-1} V S U^T y \\ &= U S [S^2 + \lambda (V^T V)^{-1}]^{-1} V^{-1} V S U^T y \\ &= U S [S^2 + \lambda I]^{-1} S U^T y \\ &= \sum_{j=1}^D U_j \left(\frac{1}{S_j^2 + \lambda} \right) S U^T y \\ &= \left(\sum_{j=1}^D U_j S_j \frac{1}{S_j^2 + \lambda} S_j U_j^T \right) y \\ &= \sum_{j=1}^D U_j \frac{S_j^2}{S_j^2 + \lambda} U_j^T \cdot y \end{aligned}$$

Any diagonalizable matrix A has decomposition as

$$A = U \Lambda U^T = \sum_{j=1}^D U_j \lambda_j U_j^T$$

where

- $\lambda_j \in \mathbb{R}$
- $U_j \in \mathbb{R}^D$
- $U_j \in \mathbb{R}^D$

so, $\hat{y} = \sum_{j=1}^D U_j f(S_j, \lambda) U_j^T y$ where $f(S_j, \lambda) = \frac{S_j^2}{S_j^2 + \lambda}$

For the least squares we use

$$W_{ls} = (X^T X)^{-1} X^T Y$$

Then $\hat{y} = X W_{ls}$

$$= X (X^T X)^{-1} X^T Y$$

$$= U S V^T [U S V^T]^T U S V^T (U S V^T)^T Y$$

$$= U S V^T [V S U^T U S V^T]^{-1} V S U^T Y \quad \left[\begin{array}{l} S^T = S \text{ as } S \\ \text{diagonal} \\ U^T U = I_{D \times D} \end{array} \right]$$

$$= U S V^T [V S^2 V^T]^{-1} V S U^T Y$$

$$= U S V^T (V^T)^{-1} (S^2)^{-1} V^{-1} V S U^T Y$$

$$= U S I_{D \times D} (S^2)^{-1} I_{D \times D} S U^T Y$$

$$= U S (S^2)^{-1} S U^T Y$$

$$= U S S^{-1} S^{-1} S U^T Y = U I_{D \times D} I_{D \times D} U^T Y$$

For least square prediction

$$\hat{y} = U U^T Y = \sum_{j=1}^D u_j u_j^T Y$$

Here we have that

$$\hat{y}_{\text{ridge}} = \sum_{j=1}^D u_j f(s_j, \lambda) u_j^T Y$$

$$\hat{y}_{ls} = \sum_{j=1}^D u_j u_j^T Y$$

The effective degrees of freedom defined as $df(\lambda) = \sum_{j=1}^D \frac{s_j^2}{s_j^2 + \lambda}$

where s_j is singular value of X ,

(i) Now, when $\lambda = 0$, we have D parameters since there is no penalisation. This corresponds to no shrinkage, gives $df(\lambda) = D$.

This is the case for least square ridge regression.

- (ii) When λ is large, the parameters heavily constrained and the degrees of freedom will effectively be lower, tending ^{to} 0 as $\lambda \rightarrow \infty$.
i.e. all the weights w_j are shrink to zero.

Ridge regression shrinks the co-ordinates with respect to the orthonormal basis formed by the principal components, ~~with smaller~~ co-ordinates with respect to principal components u_j with smaller variance are shrunk more.
The larger λ is, the more the projection is shrunk in the direction of u_j .

Question 3

(8)

Given that $z_1(w)$

$$= (y - xw)^T (y - xw) + \lambda_2 \|w\|^2 + \lambda_1 \|w\|_1$$

$$z_2(w) = (\hat{y} - \tilde{x}w)^T (\hat{y} - \tilde{x}w) + c \lambda_1 \|w\|_1$$

where $c = \frac{1}{\sqrt{1+\lambda_2}}$, $\tilde{x} = c \begin{pmatrix} x \\ \sqrt{\lambda_2} I_d \end{pmatrix}$, $\hat{y} = \begin{pmatrix} y \\ 0_{d \times 1} \end{pmatrix}$

We want to show

$$z_1(cw) = z_2(w)$$

Now $x = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_d^1 \\ x_1^2 & x_2^2 & \dots & x_d^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^n & x_2^n & \dots & x_d^n \end{bmatrix}_{n \times d}$

$y = \begin{pmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{pmatrix}_{n \times 1} \in \mathbb{R}^n$ $w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}_{d \times 1} \in \mathbb{R}^d$

where (x^i, y^i) $i = 1, 2, \dots, n$, $x^i \in \mathbb{R}^d$
are data points

Now the modified data

$$\tilde{x} = c \begin{pmatrix} x \\ \sqrt{\lambda_2} I_d \end{pmatrix} = c \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_d^1 \\ x_1^2 & x_2^2 & \dots & x_d^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^n & x_2^n & \dots & x_d^n \\ \sqrt{\lambda_2} & 0 & \dots & 0 \\ 0 & 0 & \dots & \sqrt{\lambda_2} \end{bmatrix}_{(n+d) \times d}$$

$\hat{y} = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{(n+d) \times 1} \in \mathbb{R}^{n+d}$

Hence

$$Z_1(w) = (y^T - w^T x^T)(y - xw) + \lambda_2 w^T w + \lambda_1 |w|_1$$

$$= y^T y - y^T x w - w^T x^T y + w^T x^T x w + \lambda_2 w^T w + \lambda_1 |w|_1$$

$$Z_1(cw) = y^T y - y^T x(cw) - (cw)^T x^T y + (cw)^T x^T x(cw)$$

$$+ \lambda_2 (cw)^T cw + \lambda_1 |cw|_1$$

$$= y^T y - c y^T x w - c w^T x^T y + c^2 w^T x^T x w + c^2 \lambda_2 w^T w + c \lambda_1 |w|_1$$

$$Z_2(w) = (\tilde{y} - \tilde{x}w)^T (\tilde{y} - \tilde{x}w) + c \lambda_1 |w|_1$$

$$= \tilde{y}^T \tilde{y} - \tilde{y}^T \tilde{x}w - w^T \tilde{x}^T \tilde{y} + w^T \tilde{x}^T \tilde{x}w + c \lambda_1 |w|_1$$

calculating $\tilde{y}^T \tilde{y}$ in terms of y

$$\tilde{y}^T \tilde{y} = \sum_{i=1}^{n+d} \tilde{y}^i \tilde{y}^i = \sum_{i=1}^n y^i y^i + \sum_{i=n+1}^{n+d} 0 \cdot 0 \quad \left[\begin{array}{l} y^i = y^i \text{ for } i=1, \dots, n \\ = 0 \text{ for } i=n+1, \dots, n+d \end{array} \right]$$

$$\Rightarrow \boxed{\tilde{y}^T \tilde{y} = y^T y} \quad \text{--- (5)}$$

$$\tilde{y}^T \tilde{x} = (y^1, y^2, \dots, y^n, 0, \dots, 0) e \begin{pmatrix} x_1^1 & x_1^2 & \dots & x_1^d \\ x_2^1 & x_2^2 & \dots & x_2^d \\ \vdots & \vdots & \ddots & \vdots \\ x_n^1 & x_n^2 & \dots & x_n^d \\ \sqrt{\lambda_2} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda_2} \end{pmatrix}$$

$$= c \left(\sum_{i=1}^n y^i x_i^1, \sum_{i=1}^n y^i x_i^2, \dots, \sum_{i=1}^n y^i x_i^d \right) \begin{matrix} 0 & 0 & \dots & \sqrt{\lambda_2} \\ 1 \times d \end{matrix}$$

$$\Rightarrow \boxed{\tilde{y}^T \tilde{x} = c y^T x} \quad \text{--- (6)}$$

Transposing this $\boxed{\tilde{x}^T y = c x^T y} \quad \text{--- (7)}$

Now $\tilde{X}^T \tilde{X}$

$$= C \begin{bmatrix} x_1^1 & \dots & x_1^{d+n} \\ x_2^1 & \dots & x_2^{n+d} \\ \vdots & & \vdots \\ x_d^1 & \dots & x_d^{n+d} \end{bmatrix}_{d \times (n+d)} C \begin{bmatrix} x_1^1 & \dots & x_d^1 \\ x_1^2 & \dots & x_d^2 \\ \vdots & & \vdots \\ x_1^{n+d} & \dots & x_d^{n+d} \end{bmatrix}_{(n+d) \times d}$$

$$= C^2 \begin{bmatrix} \sum_{i=1}^{n+d} x_1^i x_1^i & \dots & \sum_{i=1}^{n+d} x_1^i x_d^i \\ \vdots & & \vdots \\ \sum_{i=1}^{n+d} x_d^i x_1^i & \dots & \sum_{i=1}^{n+d} x_d^i x_d^i \end{bmatrix}_{d \times d}$$

Where $x_1^{n+1}, x_1^{n+2}, \dots, x_1^{n+d}$ are all zero for $i=1, \dots, d$ except $x_i^{n+i} = 1$

for diagonal component

$$\sum_{i=1}^{n+d} x_j^i x_j^i = \sum_{i=1}^n x_j^i x_j^i + \sum_{i=n+1}^{n+d} x_j^i x_j^i$$

$$= \sum_{i=1}^n x_j^i x_j^i + \lambda_2$$

For non-diagonal entries.

$$\sum_{i=1}^{n+d} x_j^i x_k^i \quad k, j = 1, \dots, d \text{ and } k \neq j$$

$$= \sum_{i=1}^n x_j^i x_k^i + \sum_{i=n+1}^{n+d} 0 \cdot 0$$

$$= \sum_{i=1}^n x_j^i x_k^i$$

$x_k^i = 0$ for $k = 1, 2, \dots, d$
 $i = n+1, \dots, n+d$

$$\text{So, } \tilde{X}^T \tilde{X} = c^2 \begin{bmatrix} \sum_{i=1}^n x_i^1 x_i^1 + \lambda_2 & \dots & \sum_{i=1}^n x_i^1 x_d^1 \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^n x_d^1 x_i^1 & \dots & \sum_{i=1}^n x_d^1 x_d^1 + \lambda_2 \end{bmatrix} \quad (11)$$

$$= c^2 (X^T X + \lambda_2 I) \quad \text{--- (8)}$$

Now we substitute result (5), (6), (7), (8) in (4).

we get

$$\begin{aligned} Z_2(W) &= y^T y - c y^T X W - c W^T X^T y + W^T c^2 (X^T X + \lambda_2 I) W \\ &\quad + c \lambda_1 \|W\|_1 \\ &= y^T y - c y^T X W - c W^T X^T y + c^2 W^T X^T X W + c^2 \lambda_2 W^T W \\ &\quad + c \lambda_1 \|W\|_1 \\ &\geq Z_1(W) \end{aligned}$$

Hence proved that the elastic problem can be reduced to a lasso problem on modified data.