

# CS 335/337

## Assignment. 1

1.1 
$$mse(w, b) = \sum_{i=1}^N \frac{1}{2N} (wx_i + b - y_i)^2$$

$$\nabla mse(w, b) = \begin{pmatrix} \frac{\partial mse}{\partial w} \\ \frac{\partial mse}{\partial b} \end{pmatrix}$$

$$\begin{aligned} \therefore \frac{\partial mse}{\partial w} &= \sum_{i=1}^N \frac{1}{2N} \cdot 2 (wx_i + b - y_i) \cdot (x_i) \\ &= \sum_{i=1}^N \frac{x_i (wx_i + b - y_i)}{N} \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial mse}{\partial b} &= \sum_{i=1}^N \frac{1}{2N} \cdot 2 (wx_i + b - y_i) \cdot (1) \\ &= \sum_{i=1}^N \frac{(wx_i + b - y_i)}{N} \end{aligned}$$

$$\nabla mse(w, b) = \sum_{i=1}^N \begin{bmatrix} \frac{x_i (wx_i + b - y_i)}{N} \\ \frac{(wx_i + b - y_i)}{N} \end{bmatrix}$$

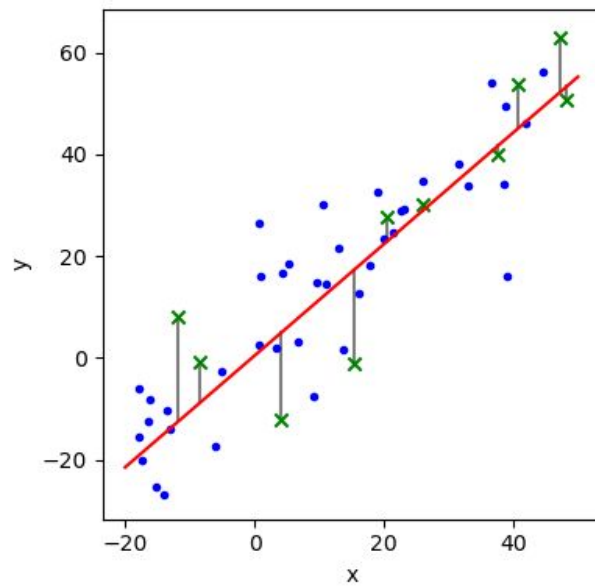
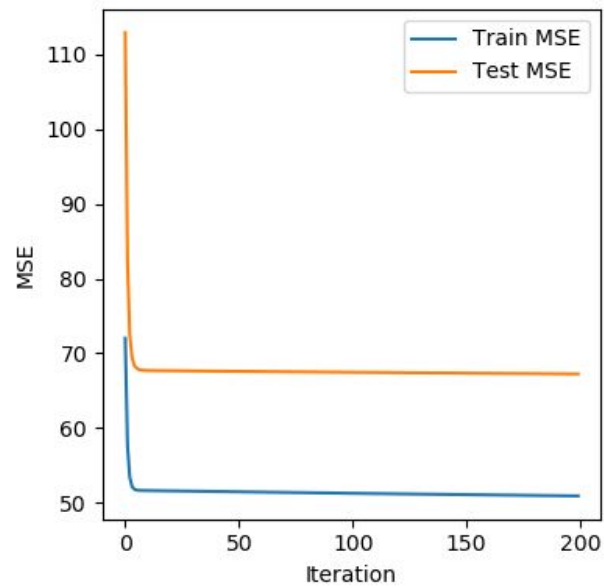
$$= \sum_{i=1}^N \left( \frac{wx_i + b - y_i}{N} \right) \begin{bmatrix} x_i \\ 1 \end{bmatrix}$$

Vectorised:-

(Note:- All summations are  $\sum_{i=0}^{n_{\text{samples}}}$ ).

## Assignment 1 Graphs: Niraj Mahajan: 180050069

Q1.2)d)



Q 1.2) (d) The desired image is on the previous page.

The plot on the right has  $x, y$  coordinates of the data scattered in the 2d plane. The red line is the least square fit (line that minimises sum of squared distance from the line) of the data (train).

The green points are the test datapoints and the subsequent green lines are the deviation of the test data from their prediction.

So, assuming the <sup>train</sup> data follows a linear relation, the red line is the best possible function that estimates/predicts the <sup>train</sup> data most appropriately by minimising the total MSE error.

2.1] Assuming :- Bias term is included in  $W$  while preprocessing  $X$ .

For a given sample  $X_i = [x_{i1} \ x_{i2} \ \dots \ x_{id}]^T$   
and for given  $w = [w_1 \ w_2 \ \dots \ w_d]^T$

$$\text{prediction} = (w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id}).$$

Hence, ~~off~~ for  $X_{N \times D}$ ,  $W_{D \times 1}$

$$\hat{Y}_{N \times 1} = X_{N \times D} \cdot W_{D \times 1}$$

$$\left\{ \text{For each individual } \hat{y}_i = x_{i \times D}^T \cdot W_{D \times 1} \right\}$$

$$b) \text{mse} = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

$$= \frac{1}{2N} \sum_{i=1}^N (w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} - y_i)^2$$

$$\frac{\partial \text{mse}}{\partial \vec{w}} = \left[ \frac{\partial \text{mse}}{\partial w_1} \quad \frac{\partial \text{mse}}{\partial w_2} \quad \dots \quad \frac{\partial \text{mse}}{\partial w_d} \right]^T$$

For a general  $j$ , st  $d \geq j \geq 1$ ,

$$\frac{\partial \text{mse}}{\partial w_j} = \frac{1}{2N} \sum_{i=1}^N 2 \cdot (w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} - y_i) \times x_{ij}$$

$$\frac{\partial \text{mse}}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i) x_{ij}$$

Vectorized Formula.

$$\frac{\partial \text{mse}}{\partial w_j} = \frac{1}{N} \text{sum} \left( (X_{N \times d} \cdot w_{d \times 1} - Y_{N \times 1}) .* X[:, j] \right).$$

$$\frac{\partial \text{mse}}{\partial \vec{w}_{d \times 1}} = \frac{1}{N} \left[ (X_{N \times d} \cdot w_{d \times 1} - Y_{N \times 1})^T \cdot X_{N \times d} \right]^T$$

c) In case of Ridge Regression:-

Assuming  
type in problem  
statement here

$$mse = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2 + \lambda \|W\|^2.$$

Let

$mse_1$

$mse_2$

Define  $mse_1$  and  $mse_2$  as above.

→ To find  $\frac{\partial mse_1}{\partial \vec{W}}$ , refer to part (b).

~~reason  $\frac{\partial mse_2}{\partial \vec{W}}$~~

→ To find  $\frac{\partial mse_2}{\partial \vec{W}}$ .

$$mse_2 = \lambda \cdot W^T \cdot W.$$

$$= \lambda (w_1^2 + w_2^2 + \dots + w_D^2).$$

$$\frac{\partial mse_2}{\partial \vec{W}} = \left[ \frac{\partial mse_2}{\partial w_1} \quad \frac{\partial mse_2}{\partial w_2} \quad \dots \quad \frac{\partial mse_2}{\partial w_D} \right]^T$$

For a general  $j \in [0, D]$

$$\frac{\partial mse_2}{\partial w_j} = 2\lambda w_j$$

$$\therefore \boxed{\frac{\partial mse_2}{\partial \vec{W}} = 2\lambda \vec{W}} \quad \text{--- ①}$$



Hence,  $\frac{\partial \text{mse}}{\partial \vec{w}} = \frac{\partial \text{mse}_1}{\partial \vec{w}} + \frac{\partial \text{mse}_2}{\partial \vec{w}}$

From part (b), eqn ①

$$\therefore \frac{\partial \text{mse}}{\partial \vec{w}_{0 \times 1}} = \frac{1}{N} \left[ \left( X_{N \times D} \cdot W_{0 \times 1} - Y_{N \times 1} \right)^T \cdot X_{N \times D} \right]^T + 2\lambda \vec{w}_{0 \times 1}$$

Q-3.1

$$E(w) = \frac{1}{2n} \sum_{i=1}^n \delta_i^2 (y_i - w^T x_i)^2.$$

$$= \frac{1}{2n} \sum_{i=1}^n (y_i \delta_i - w^T x_i \delta_i)^2.$$

each  $x_i, y_i$  are being transformed by  $\delta_i$

Consider a diagonal matrix  $R = \text{diag}(\delta_1, \delta_2, \dots, \delta_n)$

Hence, in our linear regression, we have the following substitution,

$$X_{N \times D} \rightarrow R_{N \times N} X_{N \times D}$$

$$Y_{N \times 1} \rightarrow R_{N \times N} \cdot Y_{N \times 1}$$

→ For a linear regression,  $X_{N \times D} W_{0 \times 1} = Y_{N \times 1}$ , the closed form can be obtained as

$$X^T X \cdot W = X^T Y$$

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

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

Hence, in our case, with the transformation of  $R$ , closed form solution:

$$W = ((RX)^T(RX))^{-1} (RX)^T(RY)$$

$$W = (X^T R^2 X)^{-1} X^T R^2 Y$$

Q. 4.1 The program throughs error when  $X^T X$  is singular, as it will be non invertible.

If we look at the dataset we observe that  $X_0$  and  $X_2$  are related by the following relation:-

$$X_2 = 3X_0.$$

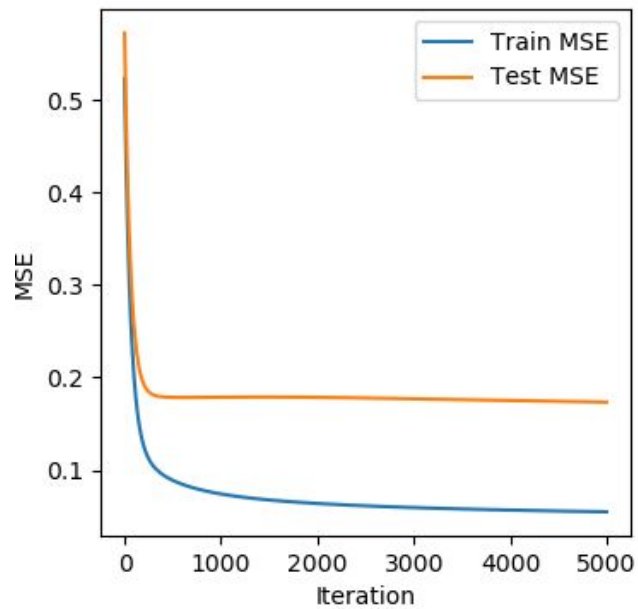
This is the reason for the singularity of  $X^T X$ . Hence, we can simply remove either of  $X_0, X_2$  as one of them is redundant and not adding anything to our training process. Also, doing so will make  $X^T X$  non singular.

Q. 4.2 There exists no solution for the closed form of OLS when the matrix  $X$  is rank deficient, i.e., one or more columns of  $X$  are correlated ~~by~~ <sup>by</sup> some ~~non~~ <sup>linear</sup> relation. In this case  $X^T X$  will be singular and no closed form solution will exist for the linear regression.

→ On the other hand, Gradient Descent will still converge to some solution as the loss function (MSE) is convex and the algorithm will try to reach the global optima.

## Assignment 1 Graphs: Niraj Mahajan: 180050069

Q2) OLS



Q2) Ridge

