

Ayush Madhavi Sohini

2019156

M2-Assignment 1



Q1  $Y = X\theta + \epsilon \rightarrow \text{loss} = \epsilon = Y - X\theta$

$Y \rightarrow$  Vector of predicted values

$X \rightarrow$  Vector of i/p features

$\epsilon \rightarrow$  error value

The loss function needed to be minimized  $\rightarrow$

$$L(\theta) = \frac{1}{2N} \sum_{i=1}^N (y(x^i) - x^i \theta)^2 \rightarrow \text{Mean Squared Loss}$$

$x_i \in \mathbb{R}^d \rightarrow i^{\text{th}}$  sample from data set of size  $N$

In vector form  $\rightarrow$

$$L(\theta) = \frac{1}{2N} (Y - X\theta)^2$$

$\downarrow$   
predicted                       $\hookrightarrow$  label vector

$$= \frac{1}{2N} (Y - X\theta)^T (Y - X\theta)$$

$$= \frac{1}{2N} (Y^T - (X\theta)^T) (Y - X\theta)$$



$$= \frac{1}{2N} (Y^T - \theta^T X^T) (Y - X\theta)$$

$$= \frac{1}{2N} (Y^T Y - Y^T X\theta - \theta^T X^T Y + \theta^T X^T X \theta)$$

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{1}{2N} [-X^T Y - X^T Y + 2X^T X \theta]$$

$$= \frac{1}{2N} [2X^T X \theta - 2X^T Y]$$

To move in the direction of the optimal solution,  
re equate the derivative to zero

$$\frac{1}{2N} [2X^T X \theta - 2X^T Y] = 0$$

$$\therefore X^T X \theta = X^T Y$$

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



Q2

Whenever the inverse of  $X^T X$  exists, the closed form solution exists. If  $X^T X$  is a **singular** matrix, the closed form solution won't exist.

Q3

The closed form solution is a better option **ONLY** when the size of the i/p matrix  $X$  is small or  $X$  is sparse. When  $X$  is a very large (suppose  $A$  has  $10^5$  entries),  $X^T X$  would be a  $10^5 \times 10^5$  matrix i.e. it has  $10^{10}$  entries, which would be very **difficult to store**. Also performing  $(X^T X)^{-1}$  is also **computationally inefficient** on such a large matrix.

Also if  $X^T X$  is **singular**, the inverse doesn't exist anyways. In such cases, the iterative methods like gradient descent are a better choice.

Q4

$$y(x^i) = \theta_0 + x_1^i \theta_1 + x_2^i \theta_2 \dots x_n^i \theta_n + \epsilon^i$$

$$y(x^1) = \theta_0 + x_1^1 \theta_1 + x_2^1 \theta_2 \dots x_n^1 \theta_n + \epsilon^1$$

$\vdots$

$$y(x^m) = \theta_0 + x_1^m \theta_1 \dots x_n^m \theta_n + \epsilon^m$$

---

$m$



$$y(x^1) + \dots + y(x^m) = m\theta_0 + \theta_1 \sum_{i=1}^m x_1^i + \theta_2 \sum_{i=1}^m x_2^i + \dots + \theta_n \sum_{i=1}^m x_n^i + \sum_{i=1}^m \epsilon^i$$

Dividing both sides by  $m$

$$\frac{y(x^1) + \dots + y(x^m)}{m} = \theta_0 + \frac{\theta_1 \sum_{i=1}^m x_1^i}{m} + \dots + \frac{\theta_n \sum_{i=1}^m x_n^i}{m} + \frac{\sum_{i=1}^m \epsilon^i}{m}$$

$= 0$  as  $\epsilon$  follows a zero mean gaussian

$$\therefore \bar{y} = \theta \bar{x} + \theta_0 \rightarrow \text{eqn of a line}$$

$$\text{where } \bar{x} = \begin{bmatrix} \sum_{i=1}^m x_1^i \\ \vdots \\ \sum_{i=1}^m x_n^i \end{bmatrix} \times \frac{1}{m}$$



$$\Theta = [\Theta_1 \ \Theta_2 \ \dots \ \Theta_n]$$

$$\bar{Y} = \begin{bmatrix} y \ (x^1) \\ \vdots \\ y \ (x^m) \end{bmatrix} \times \frac{1}{m}$$

Q5

Linear Regression Model throws a continuous output. So if we can set a particular threshold i.e. if the o/p of the linear regressor is above a threshold then it belongs to class A else class B, it can work as a binary classifier. The threshold values can be determined using tools like the ROC-AUC curve. However the data is rarely distributed as a gaussian, so this is not a good classifier as in linear regression the error in the data is assumed to be a gaussian.