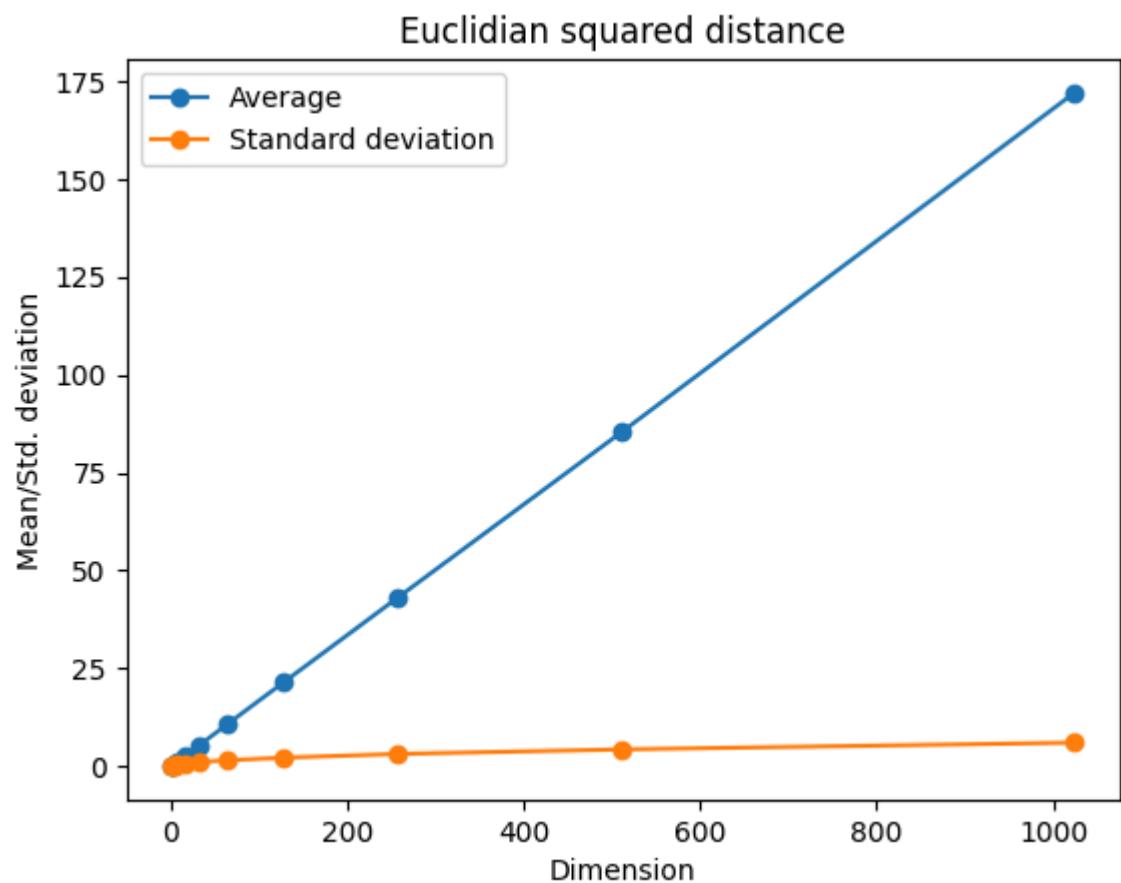


Q 1. (a)





HOMEWORK - 1

Q. 1. (b) For a univariate random variable X with uniform distribution in the interval $[0, 1]$, the probability density function $f_x(x)$ is given by:

$$f_x(x) = \begin{cases} 1, & 0 \leq x \leq 1 \\ 0, & \text{O.W.} \end{cases}$$

$$\therefore E[X] = \int_{-\infty}^{\infty} x f(x) dx = \int_0^1 x dx = \frac{x^2}{2} \Big|_0^1 = \frac{1}{2}$$

$$E[X^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \int_0^1 x^2 dx = \frac{x^3}{3} \Big|_0^1 = \frac{1}{3}$$

\therefore For $Z = (X-Y)^2 = X^2 + Y^2 - 2XY$,

$$\begin{aligned} E[Z] &= E[X^2 + Y^2 - 2XY] \\ &= E[X^2] + E[Y^2] - 2E[XY] \\ &= E[X^2] + E[Y^2] - 2E[X]E[Y] \quad \{ \because X \text{ & } Y \text{ are indep.} \} \\ &= \frac{1}{3} + \frac{1}{3} - 2 \times \frac{1}{2} \times \frac{1}{2} \\ &= \frac{2}{3} - \frac{1}{2} \\ \therefore E[Z] &= \frac{1}{6} \end{aligned}$$

$$\text{Var}[Z] = E[Z^2] - (E[Z])^2$$

$= E[Z^2]$

$$\begin{aligned} Z^2 &= (X-Y)^4 = X^4 + Y^4 - 4X^3Y + 6X^2Y^2 - 4XY^3 \\ \therefore E[X^4] &= \int_0^1 x^4 dx = \frac{1}{5} \end{aligned}$$

$$E[X^3Y] = E[X^3]E[Y] = \left(\int_0^1 x^3 dx\right) \left(\frac{1}{2}\right) = \frac{1}{8}$$

$$E[X^2Y^2] = E[X^2]E[Y^2] = \frac{8}{9}$$

$$\begin{aligned} \therefore E[Z^2] &= E[X^4] + E[Y^4] - 4E[X^3]E[X] - 4E[Y^3]E[Y] + 6E[X^2]E[Y^2] \\ &= \frac{1}{5} + \frac{1}{5} - 4 \times \frac{1}{8} \times \frac{1}{8} - 4 \times \frac{1}{8} \times \frac{1}{8} + 2 \times \frac{1}{8} \times \frac{1}{8} \\ &= \frac{2}{5} - \frac{1}{3} + \frac{2}{3} \\ &= \frac{1}{15} \end{aligned}$$

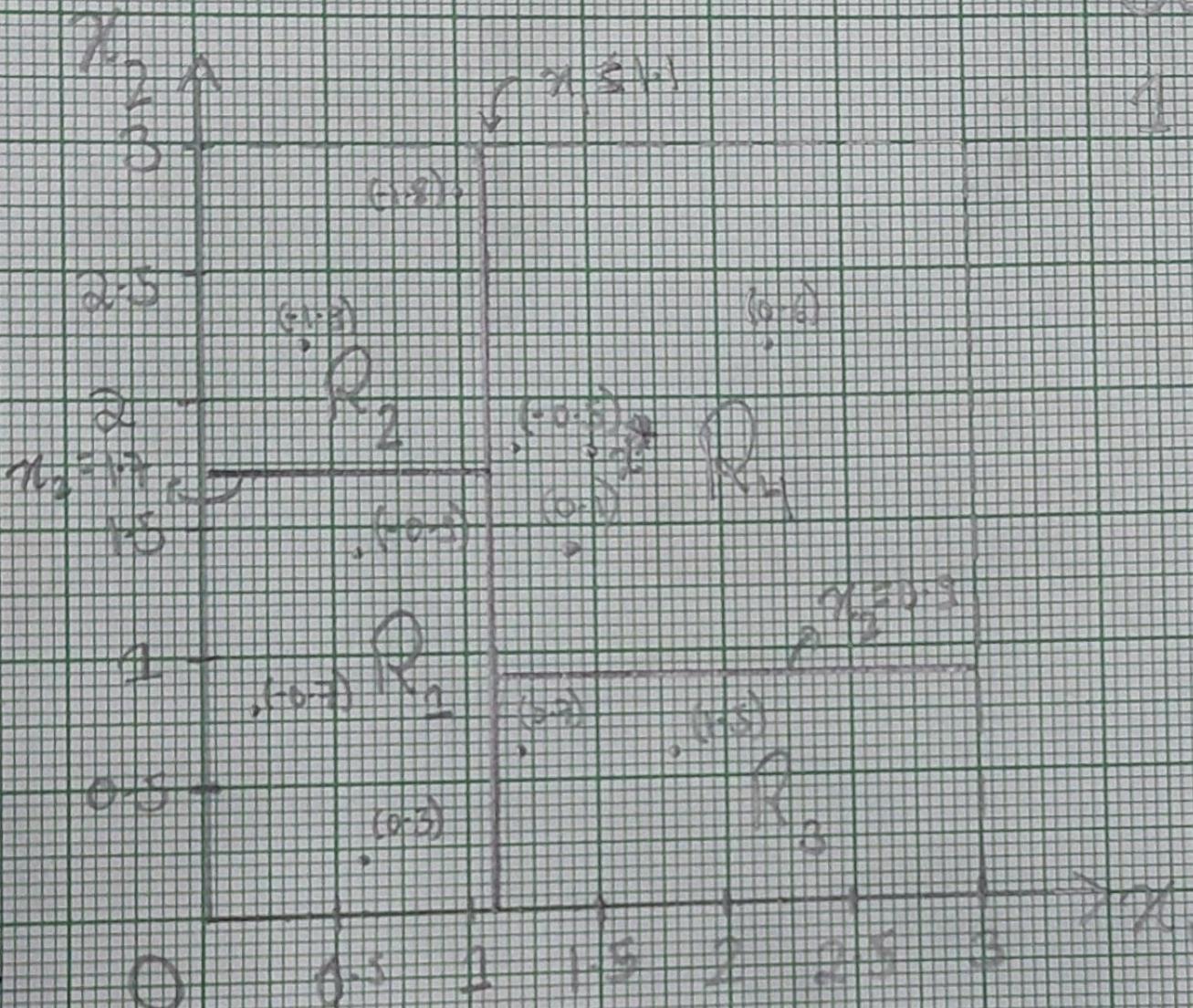
$$\begin{aligned} \therefore \text{Var}[Z] &= E[Z^2] - (E[Z])^2 \\ &= \frac{1}{15} - \frac{1}{36} \\ &= \frac{7}{180} \end{aligned}$$

$$\text{Q. 1. (c) Given, } S = Z_1 + Z_2 + \dots + Z_d$$

$$\begin{aligned} \therefore E[S] &= E[Z_1] + E[Z_2] + \dots + E[Z_d] \\ &= \frac{1}{6} + \frac{1}{6} + \dots + \frac{1}{6} \\ &= \frac{d}{6} \end{aligned}$$

$$\begin{aligned} \therefore \text{Var}[S] &= \text{Var}[Z_1 + Z_2 + \dots + Z_d] \\ &= \text{Var}[Z_1] + \text{Var}[Z_2] + \dots + \text{Var}[Z_d] \quad \{Z_i \text{ s are indep}\} \\ &= \frac{7d}{180} \end{aligned}$$

Q. 3 (a)



Scale:

1 unit = 0.5 (cm in n_2)

g3. (b) From the graph, clearly $x^* = [1.5 \ 1.8]^T$ lies in R_4

$$\therefore y \text{ for } x^* = 0.6 + 0.1 - 0.5 = 0.067$$

$$(c) \text{ For } R_1: \hat{y} = \frac{0.3 - 0.7 - 0.9}{3} = -0.433$$

Split at	\hat{y}_1	$\sum(y_i - \hat{y}_1(j,s))^2$	\hat{y}_2	$\sum(y_i - \hat{y}_2(j,s))^2$	Loss
$x_1 = 0.5$	-0.7	0	-0.3	0.72	0.72
$x_2 = 0.5$	0.3	0	-0.8	0.02	0.02
$x_2 = 1$	-0.2	0.5	-0.9	0	0.5

$\therefore x_2 = 0.5$ is the best split line for R_1 .

For R_2 :

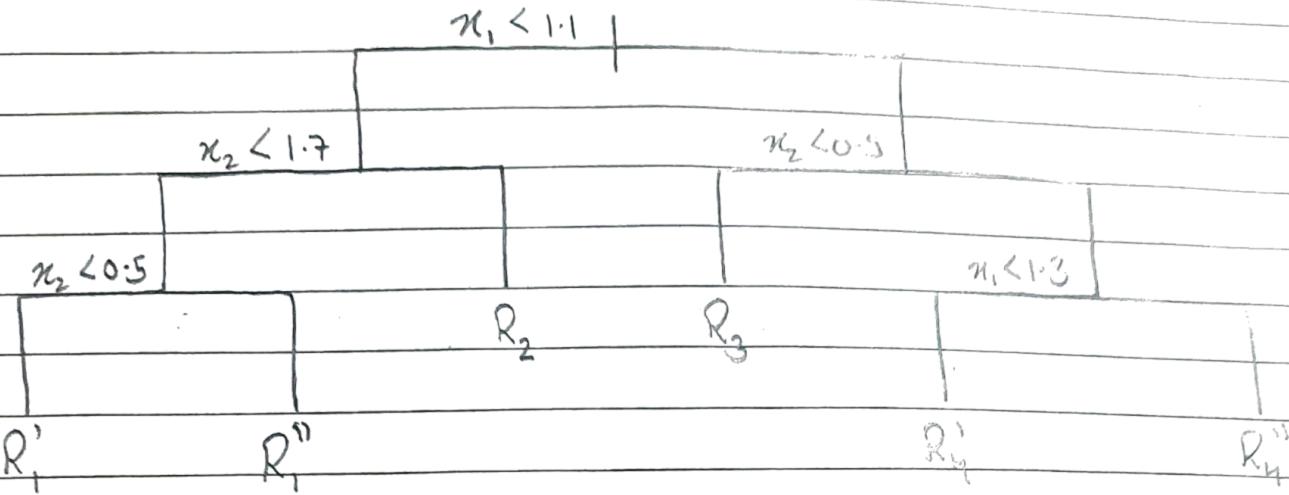
Split at	\hat{y}_1	$\sum(y_i - \hat{y}_1(j,s))^2$	\hat{y}_2	$\sum(y_i - \hat{y}_2(j,s))^2$	Loss
$x_1 = 1.3$	-0.5	0	0.35	0.125	0.125
$x_1 = 2$	-0.2	0.18	0.6	0	0.18
$x_2 = 1.5$	0.1	0	0.05	0.605	0.605
$x_2 = 2$	-0.2	0.18	0.6	0	0.18

$\therefore x_1 = 1.3$ is the best split line for R_2 .

Note that R_2 & R_3 already have only 2 points, thus there is no requirement of splitting them further.

Theoretically, the regions could have been split further such that each region contains a single point to minimize the error, but this would cause overfitting of the model.

The final regression tree looks like:



(d) For $x^* = [1.5 \ 1.8]^T$, the new tree predicts $y = \frac{0.1 + 0.6}{2} = 0.35$

Q4. The training accuracy of the decision tree = 94.16%
The testing accuracy of the decision tree = 92.72%

Q5. (a) The training model for a Linear Regression is given by

$$\hat{y} = x^T \theta + \epsilon$$

where, $x = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}$; $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}$; $\epsilon \rightarrow$ Random error with 0 mean

Now,

We want to find a θ such that the probability of our prediction(\hat{y}) matching the output label(y) is maximum, i.e.

$$\hat{\theta} = \arg \max_{\theta} p(y | x; \theta)$$

We assume that the error ϵ is a random variable and follows a Gaussian distribution with mean ($\mu = 0$) and standard deviation σ_ϵ . i.e. $\epsilon \sim N(\epsilon; 0, \sigma_\epsilon^2)$

Since all data points (x_i^T, y_i) are independent, we may write
 $p(y) = p(y_1, y_2, y_3, \dots, y_N) = p(y_1)p(y_2)\dots p(y_N)$

$$\therefore p(y|X; \theta) = \prod_{i=1}^N p(y_i|x_i^T; \theta)$$

$$\text{or, } \ln p(y|X; \theta) = \sum_{i=1}^N \ln p(y_i|x_i^T; \theta)$$

$$\text{Also, } y_i = x_i^T \theta + \epsilon \quad \text{---} \quad \textcircled{1}$$

$$\therefore p(y_i|x_i^T; \theta) = \frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp\left(-\frac{(\epsilon - \mu)^2}{2\sigma_\epsilon^2}\right)$$

Using \textcircled{1} & $\mu = 0$:

$$p(y_i|x_i^T; \theta) = \frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} \exp\left(-\frac{(y_i - x_i^T \theta)^2}{2\sigma_\epsilon^2}\right)$$

$$\therefore \ln p(y_i|x_i^T; \theta) = -\ln \sqrt{2\pi\sigma_\epsilon^2} - \frac{(y_i - x_i^T \theta)^2}{2\sigma_\epsilon^2}$$

$$\therefore \ln p(y|X; \theta) = -\underbrace{\sum_{i=1}^N \ln \sqrt{2\pi\sigma_\epsilon^2}}_{\text{constant}} - \sum_{i=1}^N \frac{(y_i - x_i^T \theta)^2}{2\sigma_\epsilon^2}$$

As since log is an increasing function, maximising $\ln p(y|X; \theta)$
 \Rightarrow $\underset{\text{maximising}}{p(y|X; \theta)}$

$$\therefore \hat{\theta} = \underset{\theta}{\operatorname{arg\max}} \ln p(y|X; \theta) = \underset{\theta}{\operatorname{arg\max}} \left(-\sum_{i=1}^N \frac{(y_i - x_i^T \theta)^2}{2\sigma_\epsilon^2} \right)$$

$$= \underset{\theta}{\operatorname{arg\min}} \left(\sum_{i=1}^N (y_i - x_i^T \theta)^2 \right)$$

Now, Define $J(\theta)$ as:

$$J(\theta) = \sum_{i=1}^N (y_i - x_i^T \theta)^2 = \|y - X\theta\|_2^2 = (y - X\theta)^T (y - X\theta)$$

$$= (y^T - (X\theta)^T) (y - X\theta)$$

$$= y^T y - (X\theta)^T y - y^T X\theta + (X\theta)^T (X\theta)$$

$$= y^T y - 2(\theta^T X^T) y + (X\theta)^T (X\theta)$$

To minimize $J(\theta)$ w.r.t. θ ,

$$\frac{\partial J}{\partial \theta} = 0$$

$$\text{or, } \frac{\partial}{\partial \theta} E(y|y) = -2 \frac{\partial}{\partial \theta} ((x^T \theta)^T y) + \frac{\partial}{\partial \theta} ((x^T \theta)^T (x^T \theta)) = 0$$

y is const.

$$\text{or, } -2 \frac{\partial}{\partial \theta} (\theta^T x^T y) + \frac{\partial}{\partial \theta} (\theta^T x^T x \theta) = 0$$

$$\text{or, } -2 (x^T y)^T + \theta^T (x^T x + x x^T) = 0$$

$$\text{or, } -2 x^T y + 2 \theta^T x^T x = 0$$

$$\therefore (\theta^T x^T x)^T = (x^T y)^T$$

$$\therefore x^T x \hat{\theta} = x^T y$$

(b) When $x^T x$ is not invertible, the closed form equation derived in part (a) does not have a unique solution. Rather it possesses infinite number of solutions where each solution $\hat{\theta}$ achieves the same minimal square error. This case typically arises when the number of samples of the training data (N) is less than the number of features (or the dimension) of the input data (p), i.e., $N < p$. This can be fixed by reducing dimensions by eliminating the redundant features or transforming the features in a new way.

Q7 (a) For a logistic regression model, the logistic function $h(x)$ is given by:

$$h(x) = \frac{e^x}{1 + e^x}$$

$$\begin{aligned} \therefore \frac{dh(x)}{dx} &= \frac{e^x(1+e^x) - (e^x)^2}{(1+e^x)^2} = \frac{e^x + (e^x)^2 - (e^x)^2}{(1+e^x)^2} = \frac{e^x}{(1+e^x)^2} \\ &= \underbrace{\frac{e^x}{1+e^x}}_{h(x)} \times \underbrace{\left(1 - \frac{e^x}{1+e^x}\right)}_{1-h(x)} \end{aligned}$$

$$\therefore \frac{dh(x)}{dx} = h(x)(1-h(x))$$

(b) For a logistic regression model, we define

$$g(\underline{x}; \underline{\theta}) = h(\underline{x}^T \underline{\theta}) = \frac{e^{\underline{x}^T \underline{\theta}}}{1 + e^{\underline{x}^T \underline{\theta}}}$$

Now,

$$p(y=1 | \underline{x}; \underline{\theta}) = g(\underline{x}; \underline{\theta}) = \frac{e^{\underline{x}^T \underline{\theta}}}{1 + e^{\underline{x}^T \underline{\theta}}}$$

$$p(y=0 | \underline{x}; \underline{\theta}) = 1 - g(\underline{x}; \underline{\theta}) = \frac{e^{-\underline{x}^T \underline{\theta}}}{1 + e^{-\underline{x}^T \underline{\theta}}}$$

$$\hat{\underline{\theta}} = \arg \max_{\underline{\theta}} p(y | \underline{x}; \underline{\theta})$$

Maximising $p(y | \underline{x}; \underline{\theta})$ is as good as maximising $\log(p(y | \underline{x}; \underline{\theta}))$

$$\therefore \hat{\underline{\theta}} = \arg \max_{\underline{\theta}} \ln(p(y | \underline{x}; \underline{\theta}))$$

$$\begin{aligned} p(y | \underline{x}; \underline{\theta}) &= \prod_{i=1}^N p(y^{(i)} | \underline{x}^{(i)}; \underline{\theta}) \\ \ln(p(y | \underline{x}; \underline{\theta})) &= \sum_{i=1}^N \ln(p(y^{(i)} | \underline{x}^{(i)}; \underline{\theta})) \end{aligned}$$

$$\begin{aligned} \therefore \hat{\underline{\theta}} &= \arg \max_{\underline{\theta}} \sum_{i=1}^N \ln(p(y^{(i)} | \underline{x}^{(i)}; \underline{\theta})) \\ &= \arg \min_{\underline{\theta}} - \sum_{i=1}^N \ln(p(y^{(i)} | \underline{x}^{(i)}; \underline{\theta})) \end{aligned}$$

$$\therefore \hat{\theta} = \arg \min_{\theta} \sum_{i=1}^N \left[y^{(i)} \ln(g(\underline{x}^{(i)}; \underline{\theta})) - (1-y^{(i)}) \ln(1-g(\underline{x}^{(i)}; \underline{\theta})) \right]$$

$$\therefore J(\underline{\theta}) = - \sum_{i=1}^N \left[y^{(i)} \ln(g(\underline{x}^{(i)}; \underline{\theta})) + (1-y^{(i)}) \ln(1-g(\underline{x}^{(i)}; \underline{\theta})) \right]$$

(c)

$$\frac{\partial J}{\partial \theta_j} = - \sum_{i=1}^N \frac{\partial}{\partial \theta_j} (y^{(i)} \ln(h(\underline{x}^{(i)} \underline{\theta}))) + (1-y^{(i)}) \ln(1-h(\underline{x}^{(i)} \underline{\theta}))$$

$$= - \left[\sum_{i=1}^N \frac{y^{(i)} h(\underline{x}^{(i)} \underline{\theta}) (1-h(\underline{x}^{(i)} \underline{\theta})) \underline{x}_j^{(i)}}{h(\underline{x}^{(i)} \underline{\theta})} + \frac{(1-y^{(i)}) (-h(\underline{x}^{(i)} \underline{\theta}) (1-h(\underline{x}^{(i)} \underline{\theta}))) \underline{x}_j^{(i)}}{1-h(\underline{x}^{(i)} \underline{\theta})} \right]$$

$$= - \left[\sum_{i=1}^N \left\{ \frac{y^{(i)}}{h(\underline{x}^{(i)} \underline{\theta})} - \frac{(1-y^{(i)})}{1-h(\underline{x}^{(i)} \underline{\theta})} \right\} h(\underline{x}^{(i)} \underline{\theta}) (1-h(\underline{x}^{(i)} \underline{\theta})) \underline{x}_j^{(i)} \right]$$

$$= - \left[\sum_{i=1}^N \left\{ y^{(i)} (1-h(\underline{x}^{(i)} \underline{\theta})) - (1-y^{(i)}) (h(\underline{x}^{(i)} \underline{\theta})) \right\} \underline{x}_j^{(i)} \right]$$

$$= - \left[\sum_{i=1}^N \{ y^{(i)} - h(\underline{x}^{(i)} \underline{\theta}) \} \underline{x}_j^{(i)} \right]$$

In matrix form, we can write it as:

$$\frac{\partial J}{\partial \underline{\theta}} = - \underline{X}^T (\underline{y} - \underline{h}(\underline{X} \underline{\theta}))$$

where, $\underline{X} = \begin{bmatrix} \underline{x}_1^{(1)} & \underline{x}_2^{(1)} & \dots & \underline{x}_p^{(1)} \\ \underline{x}_1^{(2)} & \underline{x}_2^{(2)} & \dots & \underline{x}_p^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ \underline{x}_1^{(N)} & \underline{x}_2^{(N)} & \dots & \underline{x}_p^{(N)} \end{bmatrix} ; \underline{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix} ; \underline{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}$

$$\therefore \frac{\partial J(\underline{\theta})}{\partial \underline{\theta}} = \underline{X}^T (\underline{h}(\underline{X} \underline{\theta}) - \underline{y})$$

(d) Differentiating $\frac{\partial J}{\partial \theta_j}$ w.r.t. θ_k gives:

$$\frac{\partial^2 J}{\partial \theta_j \partial \theta_k} = - \sum_{i=1}^N \left[-h(\underline{x}^{(i)} \underline{\theta}) (1-h(\underline{x}^{(i)} \underline{\theta})) \underline{x}_j^{(i)} \underline{x}_k^{(i)} \right]$$

$$= \sum_{i=1}^N h(\underline{x}^{(i)} \underline{\theta}) (1-h(\underline{x}^{(i)} \underline{\theta})) \underline{x}_j^{(i)} \underline{x}_k^{(i)}$$

∴ The Hessian matrix is given by:

$$H = \begin{bmatrix} \frac{\partial^2 J}{\partial \theta_0 \partial \theta_0} & \frac{\partial^2 J}{\partial \theta_1 \partial \theta_0} & \dots & \frac{\partial^2 J}{\partial \theta_n \partial \theta_0} \\ \vdots & & & \\ \frac{\partial^2 J}{\partial \theta_0 \partial \theta_1} & \dots & \dots & \frac{\partial^2 J}{\partial \theta_n \partial \theta_1} \\ \vdots & & & \\ \frac{\partial^2 J}{\partial \theta_0 \partial \theta_n} & \dots & \dots & \frac{\partial^2 J}{\partial \theta_n \partial \theta_n} \end{bmatrix}$$

where, $\frac{\partial^2 J}{\partial \theta_j \partial \theta_k} = \sum_{i=1}^N h(x^{(i)T} \theta) (1 - h(x^{(i)T} \theta)) x_j^{(i)} x_k^{(i)}$