

Assignment - 1

Q.1) To Prove :- The least square fit line always passes through the point  $(\bar{x}, \bar{y})$ , where  $\bar{x}$  &  $\bar{y}$  represent mean of  $x$  &  $y$ .

To show  $\bar{y}$  &  $\bar{x}$  lie on the line we need to prove  $\bar{y} = w_1 \bar{x} + b$  where  $w_1$  &  $b$  are parameters obtained from regression.

So, let  $y_i$  represent value of  $y$  for  $i^{\text{th}}$  sample & same follows for  $x_i$

Let the linear regression model be  
 $\hat{y} = w_1 x + b$

Now  $y^{(i)} = \hat{y}^{(i)} + \epsilon_i$  (True value = Predicted value + Error).

$$y^{(i)} = w_1 x^{(i)} + \epsilon_i + b$$

$$y^{(1)} = w_1 x^{(1)} + b + \epsilon_1$$

$$y^{(2)} = w_1 x^{(2)} + b + \epsilon_2$$

$$y^{(n)} = w_1 x^{(n)} + b + \epsilon_n$$

Summing above

$$\sum_{i=1}^n y^{(i)} = w_1 \sum_{i=1}^n x^{(i)} + bn + \sum_{i=1}^n \epsilon_i$$

dividing by  $n$

$$\frac{1}{n} \sum_{i=1}^n y^{(i)} = \frac{w_1}{n} \sum_{i=1}^n x^{(i)} + \frac{bn}{n} + \frac{1}{n} \sum_{i=1}^n \epsilon_i$$

So, ~~sum~~ to show that  $\sum_{i=1}^n e_i = 0$ .

$$\text{Sum of mean square errors} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{MSE} = \sum_{i=1}^n (y_i - w_1 x_i - b)^2$$

$$\text{For this to be min } \frac{\partial \text{MSE}}{\partial b} = 0$$

$$\Rightarrow \sum_{i=1}^n 2(y_i - w_1 x_i - b) \cdot (-1) = 0$$

$$\Rightarrow \text{New } e_i = y_i - w_1 x_i - b$$

$$\text{So } -2 \sum_{i=1}^n e_i = 0 \Rightarrow \sum_{i=1}^n e_i = 0$$

$$\text{So, } \bar{y} = w_1 \bar{x} + b + \bar{e}$$

$$\bar{y} = w_1 \bar{x} + b \quad (\text{as } \bar{e} = 0)$$

Hence  $(\bar{y}, \bar{x})$  lie on the least square fit line.

Q1(b)

Here let

 $\rho_{x,y}$  denote correlation coefficient between  $x$  &  $y$ 

Now, given

$$\rho_{x,z} > 0$$

$$\rho_{y,z} > 0$$

Now we claim  $\rho_{x,y}$  need not be greater than 0 always

So, using partial correlation,

$$\rho_{x \cdot y \cdot z} = \frac{\rho_{xy} - \rho_{xz} \cdot \rho_{yz}}{\sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{yz}^2}}$$

$$\Rightarrow \rho_{xy} = \rho_{xz} \cdot \rho_{yz} + \rho_{x \cdot y \cdot z} \sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{yz}^2}$$

Now by condition of partial coeff.  $-1 \leq \rho_{x \cdot y \cdot z} \leq 1$ 

So

$$\rho_{xy} = \rho_{xz} \cdot \rho_{yz} \pm \rho_{x \cdot y \cdot z} \sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{yz}^2}$$

So the upper & lower bound of  $\rho_{xy}$  are

$$\rho_{xz} \cdot \rho_{yz} \pm \sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{yz}^2} \quad \& \quad \rho_{xz} \cdot \rho_{yz} - \sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{yz}^2}$$

respectively.

Now assuming  $\rho_{xz} = 0.75$  &  $\rho_{yz} = 0.8$  which indicates strong correlation b/w means

$$\rho_{xy} \in \left[ (0.75)(0.8) - \sqrt{1-0.75^2} \sqrt{1-0.8^2}, (0.75)(0.8) + \sqrt{1-0.75^2} \sqrt{1-0.8^2} \right]$$

$$\rho_{xy} \in \left[ 0.6 - \sqrt{0.4375} \sqrt{0.36}, 0.6 + \sqrt{0.4375} \sqrt{0.36} \right]$$

$$\rho_{x,y} \in [0.6 - 0.39, 0.6 + 0.39]$$

$$\rho_{x,y} \in [0.21, 0.99]$$

So  $\rho_{x,y}$  can be between 0.21 to 0.99 as well, which indicates a <sup>very</sup> weak correlation b/w  $x, y$ .

Hence a strong correlation is not guaranteed always  
Hence the proof.



Q1(c) Given:  $x_1, x_2, x_3, \dots, x_n$  a sequence of iid random variables.  
 $E(x_1) = E(x_2) = \dots = E(x_n) = \mu$

To Prove: - For  $n \rightarrow \infty$  it  $\frac{x_1 + x_2 + \dots + x_n}{n} \rightarrow \mu$

ie

~~$P\left(\left|\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right| \geq \epsilon\right)$~~

ie let  $S_n = \frac{x_1 + x_2 + \dots + x_n}{n}$

It  $\lim_{n \rightarrow \infty} P\left(\left|\frac{S_n}{n} - \mu\right| \geq \epsilon\right) \rightarrow 0$  for any  $\epsilon > 0$

Proof

$$S_n = x_1 + x_2 + \dots + x_n$$

$$\bar{S}_n = \frac{S_n}{n}$$

$$\bar{S}_n = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\bar{S}_n = \frac{\mu n}{n}$$

$$\bar{S}_n = \mu$$

Now, By Chebyshev's inequality,

$$P[|\bar{S}_n - \mu| \geq \epsilon] \leq \frac{\text{Var}(\bar{S}_n)}{\epsilon^2}$$

Now

$$\text{Var}(\bar{S}_n) = \text{Var}\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right)$$

$$= \frac{1}{n^2} \text{Var}(X_1 + X_2 + \dots + X_n)$$

$$= \frac{1}{n^2} [\text{Var}(X_1) + \dots + \text{Var}(X_n)] \quad (\text{as } X_i \text{ are iid})$$

$$= \frac{n \sigma^2}{n^2}$$

(as  $\text{Var}(X_1) = \text{Var}(X_2) = \dots = \text{Var}(X_n) = \sigma^2$  as they are iid).

$$= \frac{\sigma^2}{n}$$

$$\text{so, } P[|\bar{S}_n - \mu| \geq \epsilon] \leq \frac{\sigma^2}{n \epsilon^2}$$

$$\text{so as } n \rightarrow \infty, P[|\bar{S}_n - \mu| \geq \epsilon] \leq \lim_{n \rightarrow \infty} \frac{\sigma^2}{n \epsilon^2}$$

$$< 0$$

$$(\text{as } \lim_{n \rightarrow \infty} \frac{1}{n} = 0)$$

Hence Proved.

Q② To Derive: Maximum A Posteriori (MAP) Sol<sup>n</sup> for linear regression

Now,  $y = w^T x + e$

here

$$e \sim \text{Normal}(0, \sigma^2)$$

Hence  $y|x \sim \text{Normal}(w^T x, \sigma^2)$

$$\Rightarrow P(y|x_i) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x_i^T w - y_i)^2}{2\sigma^2}}$$

now,

$$w = \underset{w}{\operatorname{argmax}} P(w | y_1, x_1, y_2, x_2, \dots, y_n, x_n)$$

$$w = \underset{w}{\operatorname{argmax}} \frac{P(y_1, x_1, \dots, y_n, x_n | w) P(w)}{P(y_1, x_1, \dots, y_n, x_n)} \quad (\text{Bayes' Rule})$$

$$= \underset{w}{\operatorname{argmax}} P(y_1, x_1, \dots, y_n, x_n | w) P(w)$$

Now we assume  $w$  to be  $w \sim N(\vec{0}, \tau^2 I)$

$$-\frac{1}{2} \left( \frac{w-0}{\tau} \right)^2$$

$$\text{So, } P(w) = \frac{1}{\sqrt{2\pi}\sigma^2} \frac{1}{\sqrt{2\pi}\tau^2} e$$

$$P(w) = \frac{1}{\sqrt{2\pi}\tau^2} e^{-\frac{w^2}{2\tau^2}}$$

$$\text{Now } P(y_1, x_1, \dots, y_n, x_n | w) = \prod_{i=1}^n P(y_i, x_i | w) \quad \checkmark$$

$$\text{Now } P(x_i, y_i | z) = P \cdot$$

(as  $y_i$  depends

$$P(y, x | z) = P(y | x, z) P(x | z)$$

on  $x_i$  & are  
iid)

$$\text{So, } P(y_1, x_1, \dots, y_n, x_n | w) = \prod_{i=1}^n [P(y_i | x_i, w) P(x_i | w)]$$

$$w = \underset{w}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | x_i, w) P(x_i | w) \cdot P(w)$$

$$= \underset{w}{\operatorname{argmax}} \left[ \prod_{i=1}^n P(y_i | x_i, w) \right] P(x_i) \cdot P(w).$$

(as  $x_i$  independent  
of  $w$ ).

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n [\log P(y_i | x_i, w) + \log P(w)] \rightarrow (\text{as } P(x_i) \text{ is constant})$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \left[ \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log \left( e^{-\frac{(x_i^T w - y_i)^2}{2\sigma^2}} \right) \right] + \left[ \log \left( \frac{1}{\sqrt{2\pi\tau^2}} \right) + \log \left( e^{-\frac{w^T w}{2\tau^2}} \right) \right]$$

$$= \underset{w}{\operatorname{argmin}}$$

$$\underset{w}{\operatorname{argmin}} \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i^T w - y_i)^2 + \frac{1}{2\tau^2} w^T w \quad \checkmark$$

as  $\log \left( \frac{1}{\sqrt{2\pi\tau^2}} \right)$  is  
constant so ignored.

$$w = \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (x_i^T w - y_i)^2 + \lambda (w^T w) \quad \left( \text{where } \lambda = \frac{\sigma^2}{n\tau^2} \right)$$