

ML-Assignment.

Ishit Bajpai
2020380

10. To Prove : For Simple linear regression the least square fit passes through the (\bar{x}, \bar{y}) .

Let $y = \vec{w} \cdot \vec{x} + b$ be the best fit line. (w are parameters and b is the bias) & \vec{x} is input vector consisting of features.

Alternatively; $y = w^T \vec{x} + b$

Since let \hat{y}_i be the prediction for \vec{x}_i and error in the prediction be ϵ_i .

$$\hat{y}_i = w^T \vec{x}_i + b$$

if y_i is the true value then

$$y_i = \hat{y}_i + \epsilon_i \quad i=1, \dots, n.$$

$$\Rightarrow y_i = w^T \vec{x}_i + b + \epsilon_i. \quad i=1, \dots, n.$$

Adding for $i=1 \dots n$.

$$\sum_{i=1}^n y_i = w^T \sum_{i=1}^n \vec{x}_i + \sum_{i=1}^n b + \sum_{i=1}^n \epsilon_i.$$

$$\sum_{i=1}^n y_i = w^T \sum_{i=1}^n \vec{x}_i + nb + \sum_{i=1}^n \varepsilon_i \quad \text{--- (1)}$$

On an informal note we can say that errors tend to cancel each other
 $\therefore \sum_{i=1}^n \varepsilon_i = 0$. (as their $\mu = 0$) .

Proving Mathematically.

$$J(w, b) = \frac{1}{2n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \frac{1}{2n} \sum_{i=1}^n (w^T \vec{x}_i + b - y_i)^2$$

We have to minimize the Cost function
 w.r.t to w and b .

$$\therefore \frac{\partial J(w, b)}{\partial b} = \frac{1}{2n} \cdot 2 \sum_{i=1}^n (w^T \vec{x}_i + b - y_i) \times 1$$

$$\frac{\partial J(w, b)}{\partial b} = 0. \quad (\text{Condition for finding critical point, since this is})$$

$$\Rightarrow \sum_{i=1}^n (w^T \vec{x}_i + b - y_i) = 0$$

$$\Rightarrow - \sum_{i=1}^n \varepsilon_i = 0$$

$$\Rightarrow \sum_{i=1}^n \varepsilon_i = 0.$$

Using ①

$$\left(\frac{\sum y_i}{n} \right) = w^T \left(\frac{\sum x_i}{n} \right) + b + 0$$

$$\bar{Y} = w^T \bar{X} + b$$

$\therefore (\bar{Y}, \bar{X})$ passes through the

The least square fit line passes
through (\bar{X}, \bar{Y}) .

b) ~~so~~ let's say two variables X and Y are highly correlated with a third variable Z .

Then it does not necessarily imply that X and Y are also highly correlated with each other.

for pearson's correlation

$$\rho_{XY} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

Using the result, (~~since ρ_{YZ} and ρ_{ZX} are positive~~)

$$\rho_{XY} = \rho_{YZ} \cdot \rho_{ZX} - \sqrt{(1 - \rho_{YZ}^2)(1 - \rho_{ZX}^2)}$$

• if we choose

ex 1 $\rho_{YZ} = 0.75 \quad \rho_{ZX} = 0.8$

We get $\rho_{XY} = 0.2031$.

• if we choose

$$\rho_{YZ} = 0.7 \quad \rho_{ZX} = 0.7$$

We get $\rho_{XY} = -0.02$.

The above two examples shows that even though $r_{xy} \approx r_{xz}$ and $r_{zy} \approx r_{zx}$ are highly correlated, it does not imply $r_{xy} \approx r_{zy}$ will also be highly correlated.

Informal example

Say we have $X = \text{sugar}$, $Y = \text{salt}$ and $Z = \text{water}$.

X, Z are highly correlated in sense that sugar dissolves easily in water. $\therefore Y, Z$ are also highly correlated as salt dissolves with ease in water. However salt and sugar does not dissolve when mixed i.e. that is they are not co-related.

Eg Using Partial Correlation

$$\rho_{xy} = \rho_{xz} \rho_{zy} + \rho_{\cancel{xy}} \left(\sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{zy}^2} \right)$$

~~ρ_{xy}~~
 ~~ρ_{xz}~~
 ~~ρ_{zy}~~

Since $\rho \in [-1, +1]$

$$\Rightarrow \cancel{\rho_{xy}} = \rho_{xz} \rho_{zy} \pm \sqrt{1 - \rho_{xz}^2} \sqrt{1 - \rho_{zy}^2}$$

~~ρ_{xy}~~

if $\rho_{xz} = 0.7$ $\rho_{zy} = 0.8$

then $\rho_{xy} \in [0.131, 0.9884]$

Thus even though correlation bins

ρ_{zy} & ρ_{xz} was high

ρ_{xy} can still ~~not~~ take lone value.

c) To Prove : weak law of large numbers
i.e.

It states that as number of trials approaches ∞ , the observed mean approaches theoretical Expected Mean.

OR

Sample Mean \rightarrow True Mean
 $n \rightarrow \infty$

where n is the # of IID Trials

Let X_i be ~~IID~~ ^{IID} be ϵ IID's.
($i=1 \dots n$)

with mean μ and variance $= \sigma^2$

$$\text{True Mean} = E(X_i) = \mu \quad (i=1 \dots n)$$

$$\text{Var}(X_i) = \sigma^2 \quad i=1 \dots n$$

∴ we have to prove

It $\underset{n \rightarrow \infty}{\lim} \left(\frac{x_1 + x_2 + \dots + x_n}{n} \right) = \mu$

OR

$$P \left(\left| \frac{x_1 + x_2 + \dots + x_n}{n} - \mu \right| \geq \epsilon \right) = 0$$

$$\text{Var} \left(\frac{x_1 + x_2 + \dots + x_n}{n} \right)$$

$$= \text{Var} \left(\frac{\sum_{i=1}^n x_i}{n} \right)$$

$$= \frac{1}{n^2} \left(\sum_{i=1}^n \text{Var}(x_i) + 2 \sum_{i \neq j} \text{Cov}(x_i, x_j) \right)$$

$$= \frac{1}{n^2} (n \times \sigma^2 + 0) \quad (\text{as Covariance of IID variables is equal to 0})$$

$$= \left(\frac{\sigma^2}{n} \right)$$

$$\begin{aligned} \text{Cov}(X, Y) &= E[XY] - E[X]E[Y] \\ &= E[X]E[Y] - E[X]E[Y] \\ &= 0 \end{aligned}$$

$$\therefore \text{Var} \left(\frac{x_1 + x_2 + \dots + x_n}{n} \right) = \frac{\sigma^2}{n}$$

Using Chebychev inequality

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

$$\frac{1}{\epsilon^2} \text{Var} \left(\frac{x_1 + x_2 + \dots + x_n}{n} \right)$$

$$\Rightarrow P \left(\left| \frac{x_1 + x_2 + \dots + x_n}{n} - \mu \right| \geq \epsilon \right) \leq \frac{\sigma^2}{n\epsilon^2}$$

$$\text{as } n \rightarrow \infty, P \left(\left| \frac{x_1 + \dots + x_n}{n} - \mu \right| \geq \epsilon \right) \leq 0.$$

$$\Rightarrow \lim_{n \rightarrow \infty} P\left(\left|\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right| \geq \epsilon\right) = 0.$$

$$\Rightarrow \lim_{n \rightarrow \infty} \frac{x_1 + x_2 + \dots + x_n}{n} = \mu.$$

Hence Weak LLN is proved.

Pseudocode.

Let's say x_i is $\sim \text{Bernoulli}(\frac{1}{2})$ i.e

~~probability of a~~ ~~tail~~

$x_i = 1$ if we get heads $P(x_i=1) = 1/2$
 $x_i = 0$ if we get tails $P(x_i=0) = 1/2$

~~Sum~~ Expectation Sum = 0
for $i = 1 \dots n$.

Sum + = ~~Sum~~ x_i

~~point~~ Expectation = Sum / i.

As $n \rightarrow \infty$ Expectation $\rightarrow 1/2$

(True Expectation)

(d) We have to derive MAP solution for linear regression.

Assumption: gaussian priori distribution of weights

Let $y = w^T X$ be the best fit line
~~whose~~

Assuming ϵ to be error while making a prediction, we can write true value as

$$y_i = w^T x_i + \epsilon.$$

Since $\epsilon \sim N(0, \sigma^2)$.

$$\Rightarrow y_i \sim N(w^T x_i, \sigma^2)$$

$$\Rightarrow P(y_i | \vec{x}_i, w) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(w^T \vec{x}_i - y_i)^2}{\sigma^2}}$$

Since we wish to derive MAP formulae for MAP,

$$P(w|D) = \frac{P(D|w) P(w)}{P(D)}$$

We want to maximise $P(w|D)$ i.e we want to choose w that is most likely given the Data D.

$$\bullet P(w|D) = \underset{w}{\operatorname{argmax}} \frac{P(D|w) P(w)}{P(D)}$$

(ignoring $P(D)$ since it remains constant)

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

~~P(w)~~

Since (y_i, x_i) are i.i.d.s.

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

(ignoring $P(D)$ since it remains constant)

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

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

$$\{ P(A, B | C) = \frac{P(A, B \cap C) \cancel{\times} \cancel{P(C \cap B)}}{P(C)} \}$$

$$= \frac{P(ABC)}{P(BC)} \times \frac{P(B \cap C)}{P(C)} = P(A|BC) P(B|C)$$

Using above result

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

But x_i is independent of w

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

Since $P(x_i)$ is independent of w , $\therefore P(x_i)$
~~is cons~~ is constant for given Data.

Converting ~~Products~~ to ~~the~~ log-likelihood,

$$\Rightarrow \underset{w}{\operatorname{argmax}} \sum_{i=1}^n (\log P(y_i|x_i, w) + \log P(w)) \quad -②$$

Since it is given we have to assume gaussian prior distribution of weights.

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

Pluggin Values in ②

$$\underset{w}{\operatorname{argmax}} \left(\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(w^T x_i - y_i)^2}{2\sigma^2}} \right) + \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(w^T w)/2\sigma^2} \right) \right)$$

$$= \underset{w}{\operatorname{argmax}} \left(\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2} \frac{(w^T x_i - y_i)^2}{\sigma^2} + \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{w^T w}{2\sigma^2} \right)$$

ignoring Constant terms

$$\underset{w}{\operatorname{argmax}} \left(-\frac{1}{2\sigma^2} (w^T x_i - y_i)^2 - \frac{1}{2\sigma^2} w^T w \right)$$

→ ignoring σ^2 as constant

6 DATE _____
PAGE _____

$$= \underset{w}{\operatorname{argmin}} \left(\frac{1}{n} \sum_i (w^T x_i - y_i)^2 + \frac{\lambda}{n} (w^T w) \right)$$

Multiplying by σ^2 & dividing by n .

$$= \underset{w}{\operatorname{argmin}} \left(\frac{1}{n} \sum_i (w^T x_i - y_i)^2 + \frac{\sigma^2}{n} (w^T w) \right)$$

$$w = \underset{w}{\operatorname{argmin}} \left(\frac{1}{n} \sum_i (w^T x_i - y_i)^2 + \lambda (w^T w) \right)$$

$$\lambda = \sigma^2 / n \cdot \zeta^2$$

This is the MAP sol'n for Linear
Regression.