

Lecture - 49

Simple Linear Regression - (SLR)

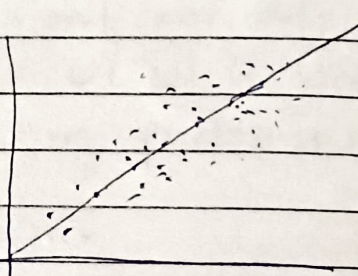
(LR) \rightarrow Supervised (Regression)

\rightarrow SLR \Rightarrow 1 input column & 1 o/p col

\rightarrow MLR \rightarrow multiple " " " " " "

\rightarrow PLR

\rightarrow SLR :-

CQIA	Package	i) Plot the graph
6.6	3.01	
7.1	4	
6.2	4.2	
1	1	

- This graph is sort of linear (Not exactly linear).

- Draw a line which passes closely to every points present. This is called best fit line

- Eqⁿ of line is: $y = mx + b$

$$\Rightarrow \text{package} = m \times \text{CQIA} + b$$

m tells about ~~pa~~ weightage

i) there are two ways to find m, c

i) Closed form Solⁿ -

\hookrightarrow we can derive mathematical formula

\hookrightarrow doesn't have diff and integrals

\hookrightarrow we use OLS (Ordinary Least

Squares)

Square)

ii) Non-closed form

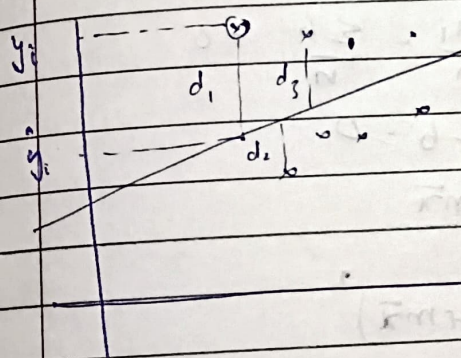
↳ we use approximation to find m, c

↳ called technique is Gradient descent

↳ used when dealing with higher descent

$$i) \quad b = \bar{y} - m\bar{x} \quad m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Derivation -



$$\Rightarrow \bar{E} = d_1 + d_2 + d_3 + \dots + d_n$$

↳ since d can be below line so add square and we don't do mode because its

not differentiable at $x=0$

$$\text{So, } \bar{E} = d_1^2 + d_2^2 + \dots + d_n^2$$

$$\bar{E} = \sum_{i=1}^n d_i^2$$

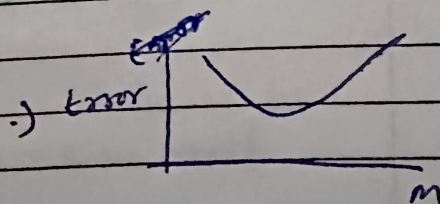
$$\text{as } d_i = (y_i - \hat{y}_i) \quad \text{So, } \bar{E} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

we want that m, b which minimizes error

$$\Rightarrow E(m, b) = \sum_{i=1}^n (y_i - mx_i - b)^2$$

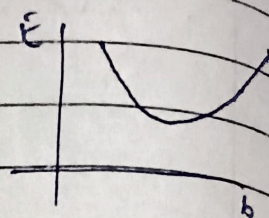
i) keeping $b=0$.

$$\bar{E}(m) = \sum_{i=1}^n (y_i - mx_i)^2$$



i) Keeping m constant

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



$$\rightarrow E = \sum_{i=1}^n (y_i - mx_i - b)^2$$

$$\frac{\partial E}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \frac{\partial}{\partial b} (y_i - mx_i - b)^2 = 0$$

$$= \frac{\sum_{i=1}^n -2(y_i - mx_i - b)}{n} = 0$$

$$= \frac{\sum y_i}{n} - \frac{\sum mx_i}{n} - \frac{\sum b}{n} = 0$$

$$\bar{y} - m\bar{x} - b = 0$$

$$\Rightarrow b = \bar{y} - m\bar{x}$$

$$\rightarrow E = \sum (y_i - mx_i - \bar{y} + m\bar{x})^2$$

$$\frac{\partial E}{\partial m} = 0$$

$$\Rightarrow \sum 2(y_i - mx_i - \bar{y} + m\bar{x})(-x_i + \bar{x}) = 0$$

$$= \sum (y_i - mx_i - \bar{y} + m\bar{x})(x_i - \bar{x}) = 0$$

$$\sum ((y_i - \bar{y}) - m(x_i - \bar{x}))(x_i - \bar{x}) = 0$$

$$\therefore \sum (y_i - \bar{y})(x_i - \bar{x}) = \sum m(x_i - \bar{x})^2$$

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

Regression Metrics

↳ It tells how good a regression model predictions are

1) MAE (Mean Absolute Error)

$$\frac{|y_1 - \hat{y}_1| + |y_2 - \hat{y}_2| + \dots + |y_n - \hat{y}_n|}{n}$$

$$\boxed{MAE = \frac{\sum |y_i - \hat{y}_i|}{n}}$$

↳ Advantage -

- ↳ unit of MAE = unit of o/p
- ↳ Robust to outlier

↳ Disadvantage -

- ↳ graph (MAE) is not differentiable at 0

2) MSE (Mean Square Error)

$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

↳ Advantage

- ↳ can use at loss fn (as differentiable)

Disad \Rightarrow MSE unit = (o/p)²

- ↳ Not Robust to outlier

3) RMSE $= \sqrt{MSE}$

↳ output has same unit

Q. R^2 Score (coeff. of determination) Date

↳ tells how well regression model explains

↳ also called goodness of fit

$$R^2 = 1 - \frac{SS_p}{SS_m}$$

$$R^2 = 1 - \frac{\left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right]_{\text{res}}}{\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right]_{\text{mean}}}$$

$$\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right]_{\text{mean}}$$

$R^2 = 1$ \rightarrow It means Regression line passes through all pts

$R^2 = 0$ \rightarrow we are just using avg

So R^2 should be close to 1

$R^2 < 0$ means our regression is doing even more mistakes

• CGPA | lac

$$R^2 = 0.8$$

means 80% of variation in lac is explained by CGPA

Remaining 20% depends on other factors

• Adjusted R^2 Score

Suppose by CGPA | lac we get R^2 score as 0.8

Now we add irrelevant columns like temp.

So R^2 score should decrease but it either remains same or increases

So

$$R^2_{\text{adj}} = 1 - \frac{(1 - R^2)(n-1)}{(n-1-k)}$$

$n \rightarrow$ # rows, $k =$ indep. columns