

day - 52

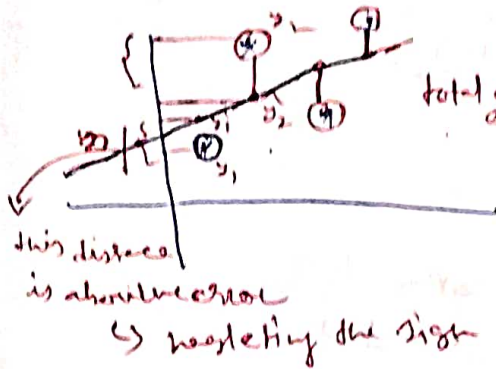
when you apply regression algorithm how how could you know how good is the algorithm.

basically we check the efficiency of regression algorithm, to find this there is many metrics

Regression Metrics

- 1) MAE (mean absolute error)
- 2) MSE (mean squared error)
- 3) RMSE (root mean squared error)
- 4) R2 Score (coefficient of determination)
- 5) adjusted R2 Score

(1) MAE



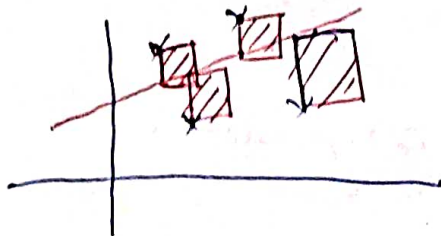
advantage :-

- 1) Same unit (like the error is in LPA, if the data is in LPA)
- 2) Robust to outliers (it can handle the outliers - MAE has also effect on it if outlier present but compared to others it is less)

disadvantage :- not differentiable at 0

As to many optimization technique there we use derivative

(2) MSE



$$mse = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

advantage :-

→ we can use any loss function

→ In previous data let we use mse as our loss function.

Disadvantage :-

→ (unit) \sim

→ for outliers it give more attention (not robust to outliers)

$$R_{mse} = \sqrt{mse}$$

advantage

→ same unit

disadvantage

→ not that robust to outliers

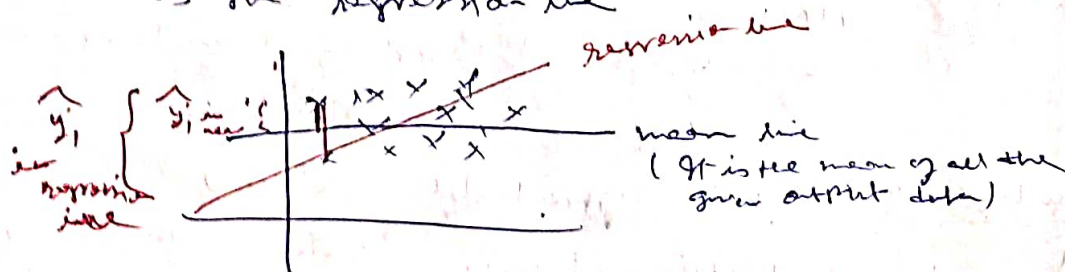
(4) R² score (coefficient of determination)
 ↳ goodness of fit
 cpa | Percentage (cpa)

NOTE:-

If we don't know what is cpa
 then for same body how can you tell
 the package simply by previous package mean

this is not the best but we can
 say something about the package

but when you bring cpa it brings
 us the regression line



So, R² score calculate how good your linear regression
 line better than mean line

$$R^2 = 1 - \frac{SSR}{SSM}$$

SSR = Sum of squared error in regression line

SSM = " " " " " " mean "

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

Now how to interpret R2 score

$$R^2 = 0 \Rightarrow \frac{[]_{reg}}{[]_{tot}} = 1 \Rightarrow []_{reg} = []_{tot}$$

i.e. the mistakes did in regression are same as the mistakes did the mean line

$$R^2 = 1 \Rightarrow \frac{[]}{[]} = 0 \Rightarrow [] = 0$$

⇒ i.e. regression line does not make any mistake

So, after calculate R2 it should goes toward 1

$$R^2 = -ve = \frac{[]}{[]} > 1$$

$$SSR > SST$$

So, SSR line made more mistake than SST

$$\frac{R^2 \rightarrow 0.80}{csp \rightarrow 1 \text{ ipa}} \rightarrow \boxed{csp} \text{ explains } 80\% \text{ of variance in ipa column}$$

i.e. it gives explanation of 80%

but for rest 20% we don't know

for ex:- csp → ipa

80% explanation gives via csp & ipa

rest 20% we don't know it may be random noise or skills?

$(R^2 \times 100)\%$ gives amount of explanation by input column.

ex:- In R^2 score
 cspal lpa
 If we add more input - then R^2 score get increase

this is not the issue
 but - if you give irrelevant - column
 it does not impact on lpa but - same time
 R^2 score gets increase. for ex:- here is temperature column

$$\text{Adjusted } R^2 = 1 - \left[\frac{(1 - R^2)(n - 1)}{(n - 1 - k)} \right]$$

n = no of rows

k = total no of independent - column.

so, if we add temp column $k \uparrow$

$\therefore (n - 1 - k) \downarrow$

$\left[\quad \right] \uparrow$

$1 - \left[\quad \right] \downarrow$

now if we add $\frac{1}{2}$ it is relevant

$k \uparrow$, $(n - 1 - k) \downarrow$, $R^2 \uparrow$, $(1 - R^2) \downarrow$

but $(1 - R^2)$ decrease more faster than $(n - 1 - k)$

$\left[\quad \right] \downarrow$

$1 - \left[\quad \right] \uparrow$

It is very useful when there are multiple columns

ex: multiple regression