Data → X, y → LR Algo → LR Model
                                  ↙  ↓  ↘
                                  c   m   eq.

$X_{train}$ → LR model → $\hat{y}$

$X_{test}$

$\hat{y} = mx + c$

Form

$y_{train}$ | $\hat{y}_{train}$
$y_{test}$ | $\hat{y}_{test}$

$MSE = \frac{1}{N} \sum (y - \hat{y})^2$

$MAE = \frac{1}{N} \sum |y - \hat{y}|$

$RSME = \frac{1}{N} \sqrt{\sum (y - \hat{y})^2}$

$\hat{y} = mx + c \pm RSME$

SD → RMS

15.95 - 16.00

6

15.99

4

14

$\hat{y} = x + 9.99$

13.99

0.01

15.99

16

± 0.01

LR Algo → Gradient Descent ✓

↳ OLS → Ordinary least Sq

In life

hit & trial

(you learn by making mistakes)

OLS

SGD

you saw from others

Reduce the chance of mistakes

G.D.

0 - 100

→ OLS

θ

②

①

Not that much Imp

least Sq. = $\frac{Cov(x,y)}{Var(x)}$

⇒ $\frac{\sum (x-\bar{x})(y-\bar{y})}{\frac{\sum (x-\bar{x})^2}{n-1}}$ = $\frac{\sum (x-\bar{x})(y-\bar{y})}{\sum (x-\bar{x})^2}$

I

II

Strenght

r, ρ → low

III

IV

LR → Strenght → Coeff. of Regression | $(R^2)$ } Coeff - of Determination

$-1 \leq \gamma \leq 1$

$0 \leq R^2 \leq 1$

$\frac{SS}{N} = \sigma^2$

worst    Best

$R^2 \; \boxed{\geq 0.7}$   → Then only, the model is of good quality model

$\bar{y}$ → Avg.
→ LR Model

$R^2 = \dfrac{\Sigma \, \text{Explained Var}}{\text{Total variation}}$

$\Rightarrow \dfrac{\text{Total var} - \text{UnExplained var}}{\text{Total variation}} \Rightarrow 1 - \dfrac{\text{Unexplained Var}}{\text{Total Variance}}$

$\Rightarrow \boxed{1 - \dfrac{\Sigma (y - \hat{y})^2}{\Sigma (y - \bar{y})^2} = R^2}$

$\bar{X}$
$A, \; \bar{H}$

$\hat{y}$   $\sigma^2$

Total var

$\hat{y}$    Error → UnExplained   $\hat{y} = mx + c$

$\bar{y}$   Explan

Explan

$x$   X

$\rightarrow$ $R^2 \rightarrow$ quality of model $(\hat{y} = a + bx)$

$\rightarrow$ strncy of Rel^ b/y $(x, y)$   $10\mu$

$\boxed{10hnz} \rightarrow \boxed{24hn}$

$c.0. \rightarrow$ $\frac{}{24}$

$\rightarrow$ $\underline{Adj\ R^2}$

$R^2 = 1 - \dfrac{SSE/\underline{N}}{SST/\underline{N}}$

$Adj\ R^2 = 1 - \dfrac{SSE/d.f.\varepsilon}{SST/d.f.T}$

$\bar{R}^2 \Rightarrow$ +ve

$\boxed{d.f.T = N-1}$

$\boxed{d.f.\varepsilon \Rightarrow N_{exp} - 1}$

$\text{Penalty system}$

$d.f.\varepsilon \leq d.f.T$

$\downarrow \Uparrow$

$1$

$d_{x_n} \rightarrow$   slope $(+)$ $(\leftarrow ve)$

data + LR$_\rightarrow$   $\boxed{9a}$ $(0)$

$\not\Rightarrow$



$(R^2 = 0.2)$  $0.1$

data thr  $100$  $\downarrow$  $91 \rightarrow$ error



$\boxed{R^2 = 0.7}$  $0.69$

→ Use $R^2$ & $Adj\ R^2$ for feature Engg.

$$SSE = 10$$
$$SST = 10$$

|  | $R^2$ | $Adj\ R^2$ |
|---|---|---|
| $y \to \quad x_1$ | 0.8 | 0.74 |

→ good features

$N = 5$

⑤
↓
③ → Predictor

$y \to$

| | $R^2$ | $Adj\ R^2$ |
|---|---|---|
| $+ \to x_2$ | 0.82 | 0.8 |

$$1 - \frac{SSE\ /\ df_{error}}{SSF\ /\ df_{T}}$$

$d.f.\ error = 2$

$d.f.\ Total = 4$

| | $R^2$ | $Adj\ R^2$ |
|---|---|---|
| $+ \to x_3$ | 0.67 | 0.6 |

↓ bad features

$$1 - \frac{10\ /\ 5}{10\ /\ 5}$$

$$1 \quad \sim 0$$

$$1 - \frac{\boxed{10\ /\ ②}}{10\ /\ 4}$$

$$1 - \frac{⑤}{②.⑤} = 1 - 2 = \boxed{-1}$$

$(m_1, c_1)$

NCERT

$X_1$

$m_2, c_2$

$x_2$

$m_3, c_3$

$x_3$ $x'$

$x_4$

$(m_4, c_4)$

Dublicate learning-s

R.D. sharma

$X_4$

→ Multi colinearity.

→ The model can learn from any $\underline{1}$ variable array N highly correlated var

$\downarrow$

Feature & Feature → No High corr

→ Target & Feature → High corr is req.

10 - n.c.v.

$\downarrow$

$\textcircled{\underline{1}}$

$\rightarrow$ Corr Matrix $\Longrightarrow$

$R^2 = 0.70$
$Adj R^2 = 0.68$
$> 0.7$

$\bigcirc$ 10 Features. 1 Target.

$F_1$ , $F_2$ , $F_p$ $\rightarrow$ $r/\rho \uparrow$

$\downarrow$

7 Feat. + 1 g

8 Features

$\rightarrow$ VIF ( Variance Inflation Factor )

$VIF = \dfrac{1}{1-R^2}$ $\rightarrow$ $[1, \infty]$

$0 \le R^2 \le 1$

If $VIF \le 5$ $\rightarrow$ M.C. features exist.

Multicolinearity will Exist

$\times$ $F_2$ , $F_6$ , $F_{10}$ $\rightarrow$ Low/None $\rightarrow$ Trash

If $VIF > 5$ $\rightarrow$ $\checkmark$

$(g+g)$

$\downarrow$

5 features

$\rightarrow$ $R^2$, Adj $R^2$ $\rightarrow$ 3 Feat