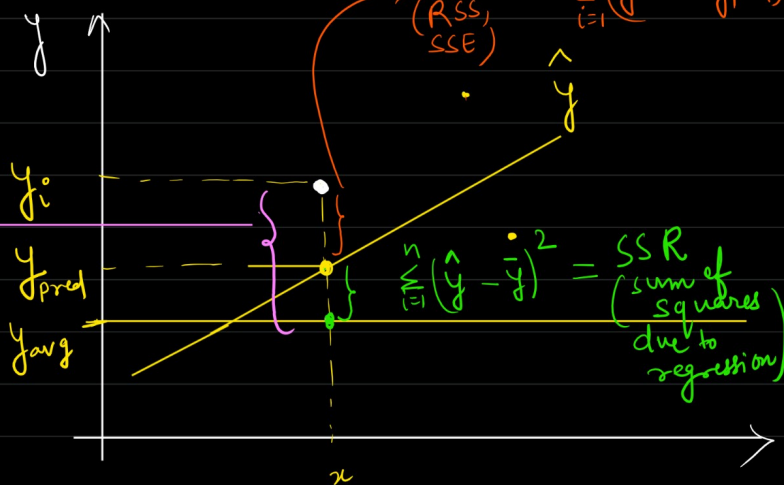# Evaluation metrics for Linear Regression

* **Performance**

① Rsquare

② adjusted R square

$$\sum_{i=1}^{n}\left(y_i - \bar{y}\right)^2$$

TSS
Total sum of squares

$$Error = \sum_{i=1}^{n}\left(y_{act} - y_{pred}\right)^2$$
(RSS, SSE)

$\hat{y}$

$y_i$

$y_{pred}$

$y_{avg}$

$y$

$x$

$$\sum_{i=1}^{n}\left(\hat{y} - \bar{y}\right)^2 = SSR$$
(sum of squares due to regression)

| Area of house | Price of house |
|---|---|
| 1000 | 50 |
| 1100 | 60 |
| ⋮ | ⋮ |
| $\bar{y}$ ← | avg of y. |

$\to \bar{y} \to$ if you don't build any model, the average value is the baseline solution.

* Error/RSS/SSE   if $1/n$ divided then it becomes MSE

Residual Sum of Square

Sum of Squared Error

$$SSR \left(\begin{array}{c}\text{Sum of}\\ \text{Squares due}\\ \text{to regression}\end{array}\right) = \begin{array}{c}\text{Explained}\\ \text{Variation}\\ \text{in y by}\\ \text{best fit line}\end{array} = \sum_{i=1}^{n}\left(\hat{y} - \bar{y}\right)^2$$

$$\begin{array}{c}SSE\\ \left(\begin{array}{c}\text{Sum of}\\ \text{Square of}\\ \text{Error}\end{array}\right)\\ RSS\left(\begin{array}{c}\text{Residual}\\ \text{Sum}\\ \text{of Square}\end{array}\right)\end{array} \Rightarrow \begin{array}{c}\text{Unexplained}\\ \text{Variation}\end{array} \Rightarrow \sum_{i=1}^{n}\left(y_i - \hat{y}\right)^2$$

$$\begin{array}{c}TSS\\ \left(\begin{array}{c}\text{Total Sum of}\\ \text{Squared}\\ \text{Error}\end{array}\right)\end{array} \Rightarrow \begin{array}{c}\text{Total Variation}\\ \text{in y /}\\ \text{difference in}\\ \text{y wrt to}\\ \text{baseline} (\bar{y})\end{array} = \sum_{i=1}^{n}\left(y_i - \bar{y}\right)^2$$

$$TSS \begin{pmatrix} Total \\ variation \\ in\ y \end{pmatrix} = \text{Explained Variation} + \text{Unexplained Variation}$$
$$(SSR) \qquad\qquad (RSS, SSE)$$

Total var → Unexplained

$y_{act}$
$y_{pre}$
$y_{av}$

Explained var

$$\underset{\substack{\uparrow \\ Unexplained \\ Variation\ (Error)}}{} \qquad \underset{\substack{\uparrow \\ Explained \\ Variation}}{}$$

$$R\ squared = 1 - \frac{RSS}{TSS} \quad or \quad \frac{SSR}{TSS}$$

↓ Total Variation.    ↳ Total variation

R square = Coefficient of determination ⟹ Out of total variation,

SSR is Variation explained

Total percentage of variation explained by model.

$$r\ square = \frac{SSR}{TSS} \times 100$$

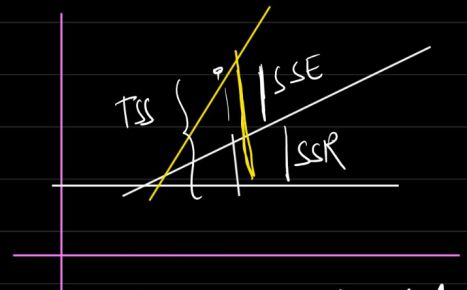→ The percentage variation in y explained by X is called R square.

$$r^2 = R^2 = \frac{SSR}{TSS}$$

TSS {    } SSE
         SSR

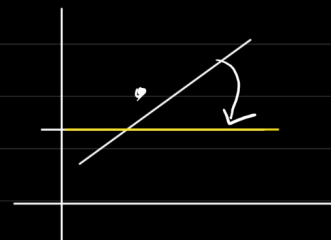$SSR \approx TSS$, it explains the variation completely

⟹ A perfect model.

$$\frac{TSS}{RSS + TSS} = 1$$

→ The maximum value of r square is 1.

→ The minimum value of rsquare is 0

$$\frac{SSR}{TSS}$$

$$SSR \approx 0$$
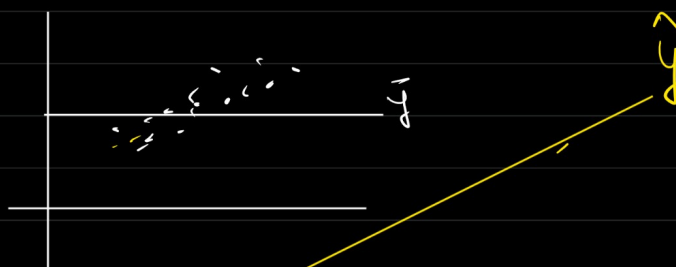
Range of r-square = 0 to 1

Closer to 1 better model will be.

* <u>Can r square be negative ?</u>

Yes it can be negative.
if the best fit line
is very very far from $\bar{y}$ (baseline model) which
is not possible given the nature of algorithm.

$m_1 →$ Rsquare 0.80 ✓
$m_2 →$ ,, 0.60

② <u>adjuster r-square</u>

$r^2 →$ %age explained variance in
y due to x

$x_1$ (Area of house)　$x_2$ (No of room)　$x_3$ Parking space　　　　$y$ price of house

* As we add more features r square will improve! Or remain as it is (constant)

r square

$x_1 - y \rightarrow$ 80%.

$x_4 x_2 - y \rightarrow$ 85%.

$x_1 x_2 x_3 - y \rightarrow$ 88%.

$x_1 x_2 x_3 \underline{x_4} - y \rightarrow$ 88.1 %

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad - \quad - \quad x_n \qquad y$

Area ⎣ #of room ⎦ parking space　　gender

To Understand if the added feature | IV is Contributing to the performance or not, we have adjusted r-square.

Adjusted r-square

$$1 - \frac{(1 - R\text{square})(N-1)}{N - P - 1}$$

where $n$ is no. of dl's
$p$ is no. of IV.

Rsquare



5%. r square

adjusted r square

# of feature

Scen-1  Rsq = 80%., N=11, p=2

adj rsq = $1 - \frac{(1 - 0.8)(11-1)}{11-2-1}$

= $1 - \frac{0.2 \times 10}{8}$ = 0.75

Scen-2

Rsq = 80, N=11, p= 8

adj rsq = $1 - \frac{(0.2)(11-1)}{11-8-1}$

= $1 - \frac{2}{2}$ = 0

① R square will always be lesser than adj rsquare. (Adj rsq < R square)

② The difference between Rsq and adj R square should not be more the 5%.

<u>Best practice</u> → Only add features in the model if the difference b/w r square & adjusted r square is not more than 3-5%.