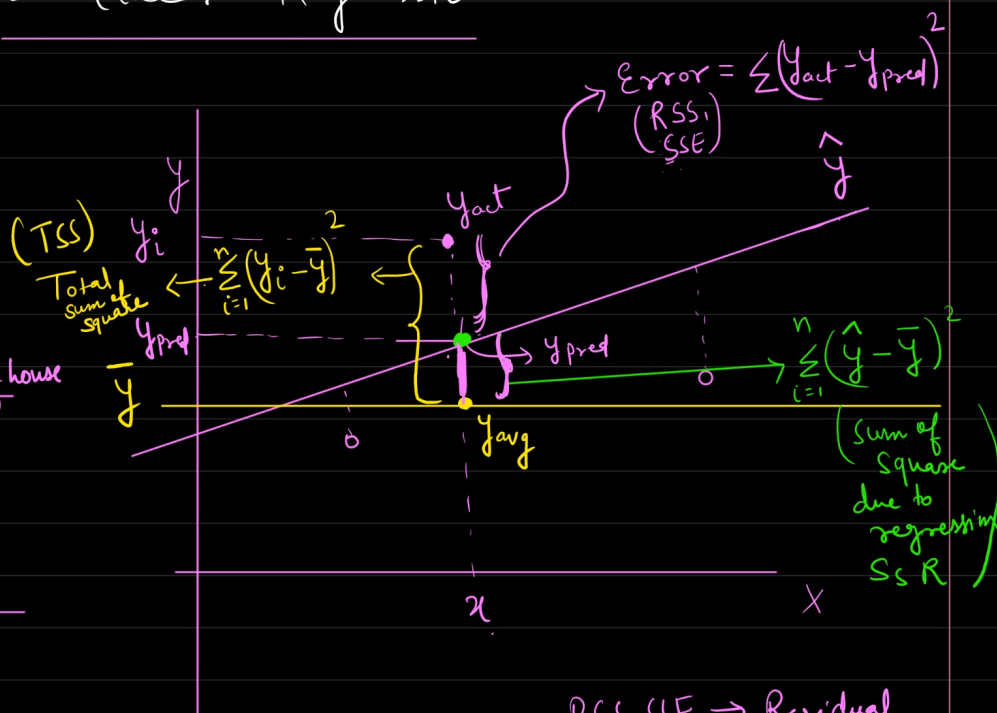


Evaluation metrics for Linear Regression

* Performance

- ① R square
- ② adjusted R square

Area of house	Price of house
1000	50
1100	60
—	—
—	—
—	—
	\bar{y}



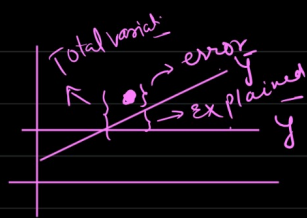
→ \bar{y} becomes your baseline soln

RSS, SSE → Residual Sum of Square,
→ Sum of squared Error.

SSR (sum of square due to regression) ⇒ Explained variation in y by best fit line ⇒ $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$

SSE (Sum of Square Error) ⇒ Unexplained variation ⇒ $\sum_{i=1}^n (y_i - \hat{y}_i)^2$
RSS - (Residual Sum of Square)

TSS = Total Variation in y = Ex Var + Unexplained Var.
(Total sum of Squared) = $\sum_{i=1}^n (y_i - \bar{y})^2$
(SST - sum of square total)



$$R_{\text{square}} = 1 - \frac{RSS}{TSS} \quad \text{or} \quad \frac{SSR}{TSS}$$

$\xrightarrow{\text{unexplained (error) variation}}$ $\xrightarrow{\text{explained variation}}$
 $\xrightarrow{\text{Total variation}}$

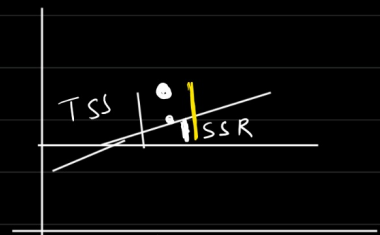
R_{square} = Coefficient of determination \rightarrow Out of (total variation), SSR is variation explained.

What is total percentage of variation explained by model?

$$r_{\text{square}} = \frac{SSR}{TSS} \times 100$$

\rightarrow The percentage variation in y explained by x is called as R -Square.

$$\underline{r_{\text{square}}} \Rightarrow r^2 \Rightarrow R^2 = \frac{SSR}{TSS}$$



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

$$\frac{TSS}{TSS} \Rightarrow 1$$

The maximum value of r_{square} is 1

$$SSR \approx 0$$

$$\frac{SSR}{TSS} = 0$$

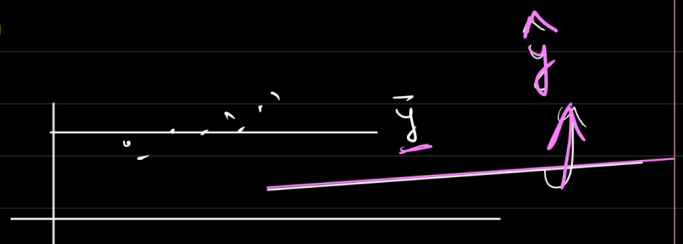
r_{square} min value = 0



r_{square} (0,1)

Can r_{square} be negative?

\rightarrow yes if the best fit



line is very very far from \bar{y} (baseline model) which is not possible given the nature of algorithm.

$m_1 \rightarrow R_{\text{square}} \underline{0.80} \checkmark$

$m_2 \rightarrow R_{\text{square}} 0.60$

② adjusted R-square

$r^2 \rightarrow$ % age explained variance in y due to x .

x_1
(Area of house)

x_2
(No of rooms)

x_3
Parking space

y Price of house

$x_1 - y \rightarrow \underline{80\%}$ r_{square}

$x_1, x_2 - y \rightarrow 85\%$

$x_1, x_2, x_3 - y \rightarrow \underline{88\%}$

* As we add more features r_{square} will improve or remain as it is (constant) $x_1, x_2, \dots, x_{10} - y \rightarrow \underline{89\%}$

x_1
(Area of house)

x_2
(No of rooms)

x_3
(Parking space)

x_4, x_5, \dots, x_{10}
(gender)

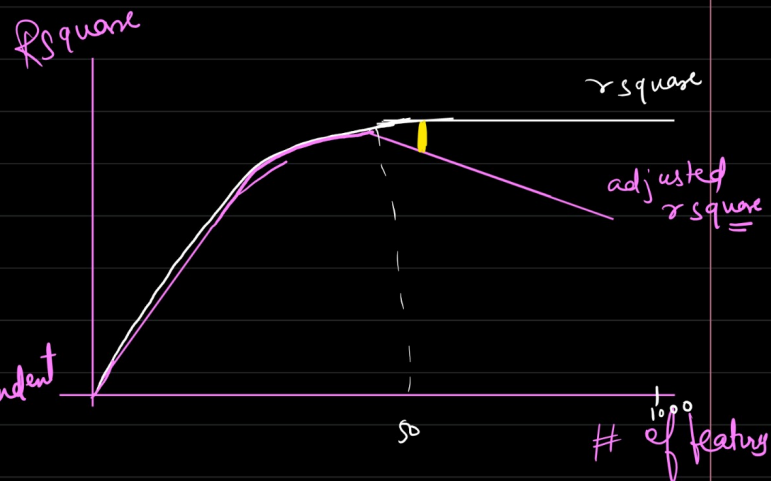
y Price of house

Adjusted R-square

$$= 1 - \frac{(1 - R^2)(N - 1)}{N - p - 1}$$

$N \rightarrow$ No of dp's

$p =$ No of independent features.



Scen-1 $R^2 = 80\%$, $N = 11$, $p = 2$

$$\text{adj } R^2_{\text{square}} = 1 - \frac{(1 - 0.8)(11 - 1)}{11 - 2 - 1}$$

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

points

→ Adj r square $<$ R square

Scen-2 $R^2 = 80\%$, $N = 11$, $p = 8$

$$\text{adj } r^2_{\text{square}} = 1 - \frac{(0.2)(11 - 1)}{11 - 8 - 1}$$

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

* Only add features in the model if the difference between R square & adjusted R square is not more than 3-5%.