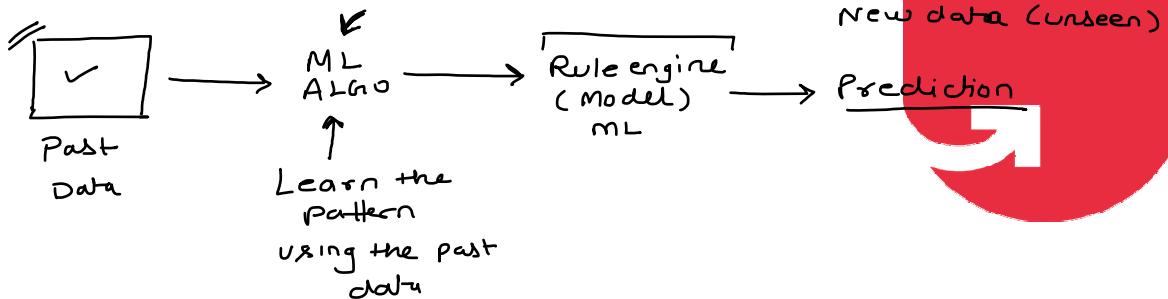
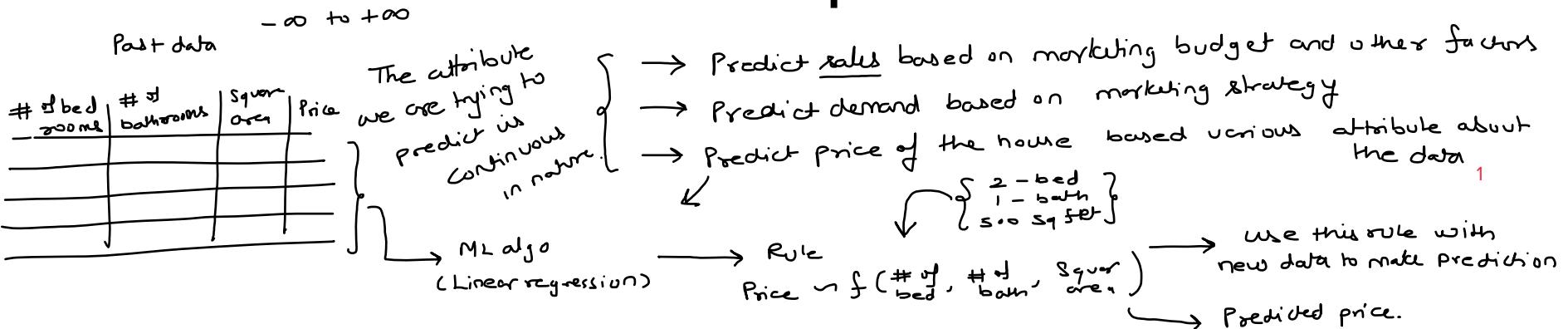


workflow - Predictive machine learning algo.



Data Science: Linear Regression Doubts and Concepts



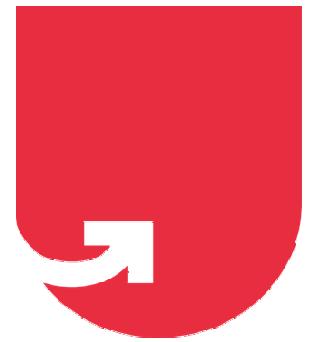
Course : Data Science

Lecture On : Linear Regression

Instructor : Sumit Shukla

Today's Agenda

- 1 Regression Line
- 2 Best Fit Line
- 3 Strength of Simple Linear Regression
- 4 Need of Multiple Linear Regression
- 5 QNA



Introduction to SLR



Machine Learning

	Pop density	avg income	conn c	# com	Sales
s ₁					
s ₂					
s ₃					

upGrad

$$\text{ML Algo (Linear regression)} \rightarrow \text{Sales} = \beta_1 (\text{pop. density}) + \beta_2 (\text{avg income})$$

Predictions.

Introduction to Linear Regression

Choose that location which generate maximum sales

I am the CEO of a hypermarket chain and I want to open new store which should give me the best sales. I am hiring you as a Data Scientist to help me figure out a location where to open the new store

What factor will you be considering to open a new store?

- The hypermarket chain has more than 5000 stores across the world
- It is upstream hypermarket store catering to high end products ↙
- There are more than 100 locations he needs to choose from

(dependent variable)

Target

$$[\text{Sales} = f(\text{pop. density}, \text{avg income}, \text{connect}, \text{avg age}, \dots)]$$

$$\text{Sales} = \beta_1 (\text{pop. density}) + \beta_2 (\text{avg income}) + \beta_3 (\text{connect}) - \beta_4 (\# \text{ of competitor}) + \dots$$

Equation of straight line.

weights / coefficients

The factors that influence your target variable are independent variables

what are the factors that are related to the location and will impact my sales at a particular store?

(independent variables)

}

→ Population density
→ avg income
→ connectivity
→ avg age
→ # of competitor
→ parking
→ education level
→ standalone vs shopping complex 5

Introduction to Linear Regression

What could impact sales ?

- Population Density in the area
- Disposable Income
- Demographics of the region
- Parking size of the location
- No of other grocery stores in around (3km)
- Credit card usage
- Internet penetration/usage
- Average no of cars/household
- Avg family size/household
- No of working people/household



Sales = function (X1, X2, X3, X4,X5,X6.....)

Sales = $10X_1 + 20X_2 + 0.5X_3 + 8X_4 + \dots$

If the function is linear we call it linear regression

Now this was a example of Multiple Linear Regression, Let's stop this topic for now and think in terms of two variables only, Say X and y. This we will called as Simple Linear Regression.

Simple Linear Regression

We have historic data of Marketing Spend and Sales, Let's see that

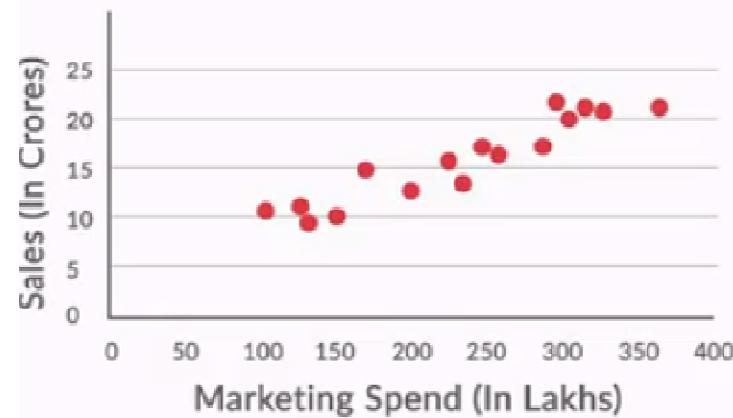
*past
data*

<i>Independent</i>	<i>dependent</i>
Marketing Spend (In Lakhs)	Sales (In Crores)
127.4	10.5
364.4	21.4
150	10
128.7	9.6
285.9	17.4

Now using this data we want to predict what we need to keep the marketing budget to generate X amount of sales.

Simple Linear Regression

Any relationship between Sales and Marketing Budget?



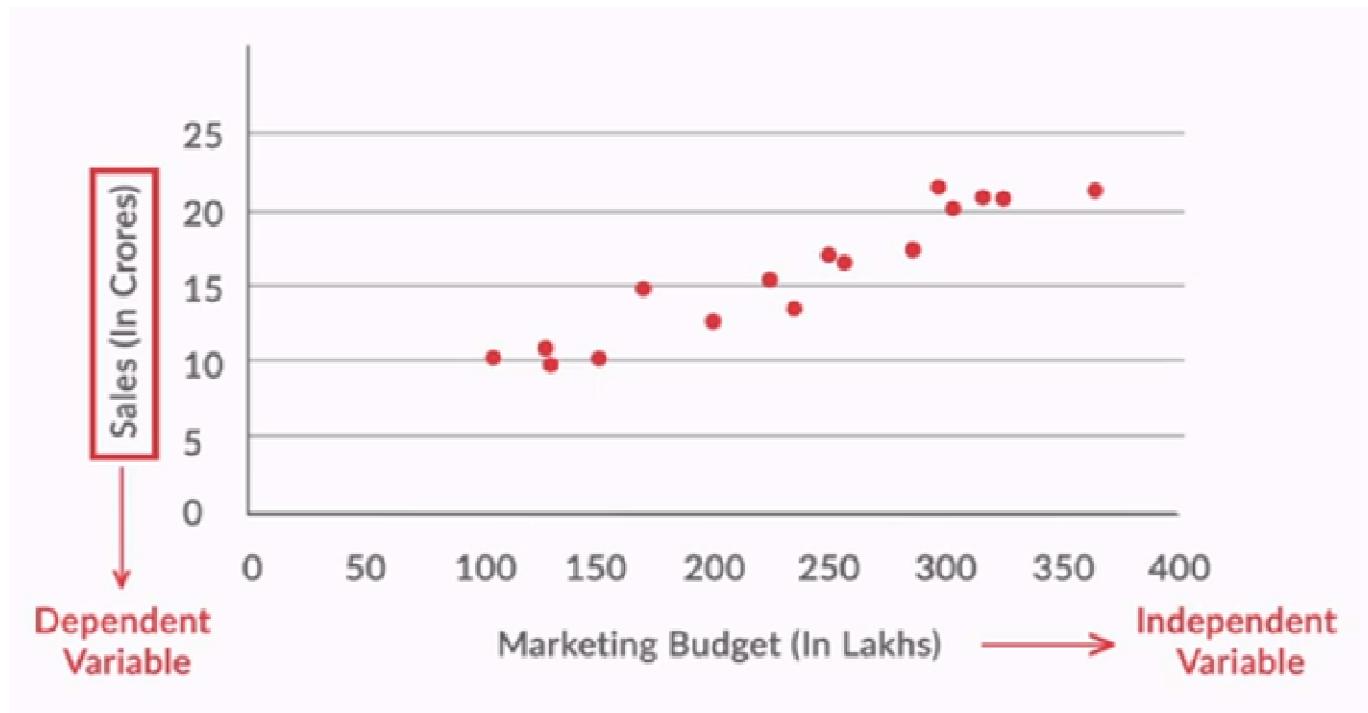
$$\text{Sales} = \beta_1(\text{marketing budget}) + \beta_0$$

|||

$$\& y = m x + c$$

Let's first check if there is any relationship between Marketing Budget and Sales? If there is no relationship we would not put any money in our Marketing budget as this is not contributing to sales.

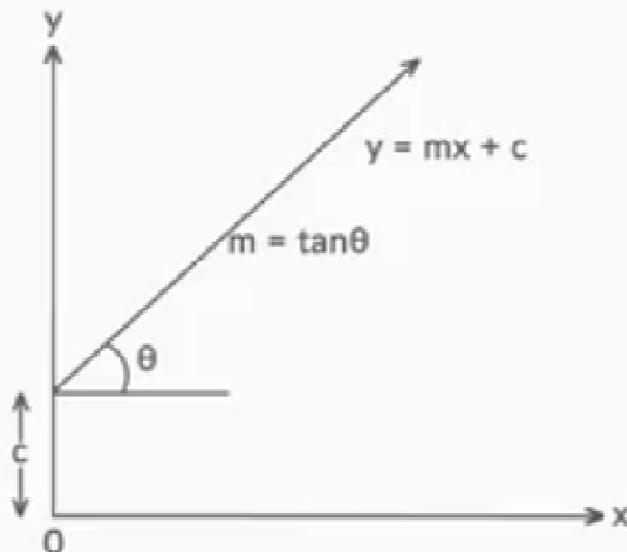
Simple Linear Regression



How to Make Prediction?



EQUATION OF STRAIGHT LINE

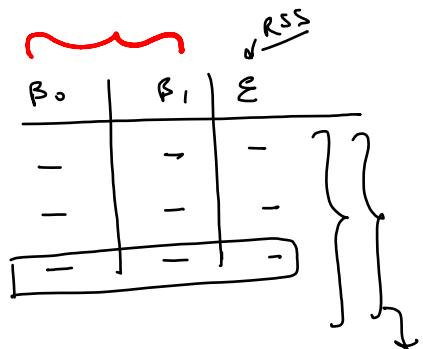


If there is zero marketing budget still we will observe C amount of sales.

SLR



upGrad



OLS = ordinary least squares method.



$$RSS = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x))^2$$

$$\begin{aligned} & y = \beta_0 + \beta_1 x \\ & \textcircled{1} \quad y = \beta_0' + \beta_1' x \\ & \textcircled{2} \quad y = \beta_0'' + \beta_1'' x \end{aligned}$$

[OLS out of all the possible lines (coefficients/weights), we will choose that particular line for which the error is minimum \equiv best fit line.]

SLR

upGrad

x	y	\hat{y}	$\epsilon = (y_i - \hat{y})$	ϵ^2
1	3	3.5	-0.5	0.25
2	7	4	3	9
4	6	5	1	1
6	4	6	-2	4
8	8	7	1	1
				<u>15.25</u>
			<u>Total</u>	
				<u>= 5.6</u>

$$e_1^2 + e_2^2 + e_3^2 + e_4^2 + e_5^2$$

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

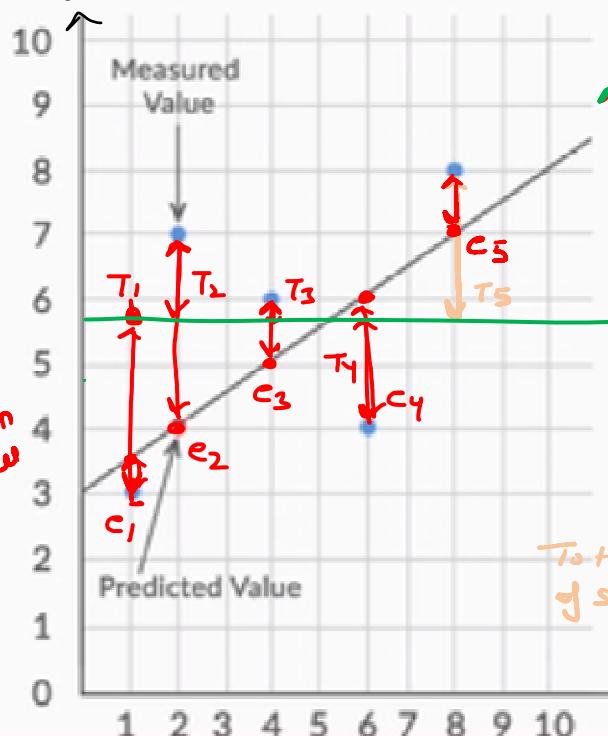
error by regression model

$$RSS = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

Residual sum of squares

(Dependent)

RESIDUALS



$$\bar{y} = \frac{\sum y}{n}$$

$$y = \beta_0 + \beta_1 x \quad] \text{ Regression model}$$

model create using linear regression

↳ aug-model. ($y_i - \bar{y}$)

$$TSS = T_1^2 + T_2^2 + T_3^2 + T_4^2 + T_5^2$$

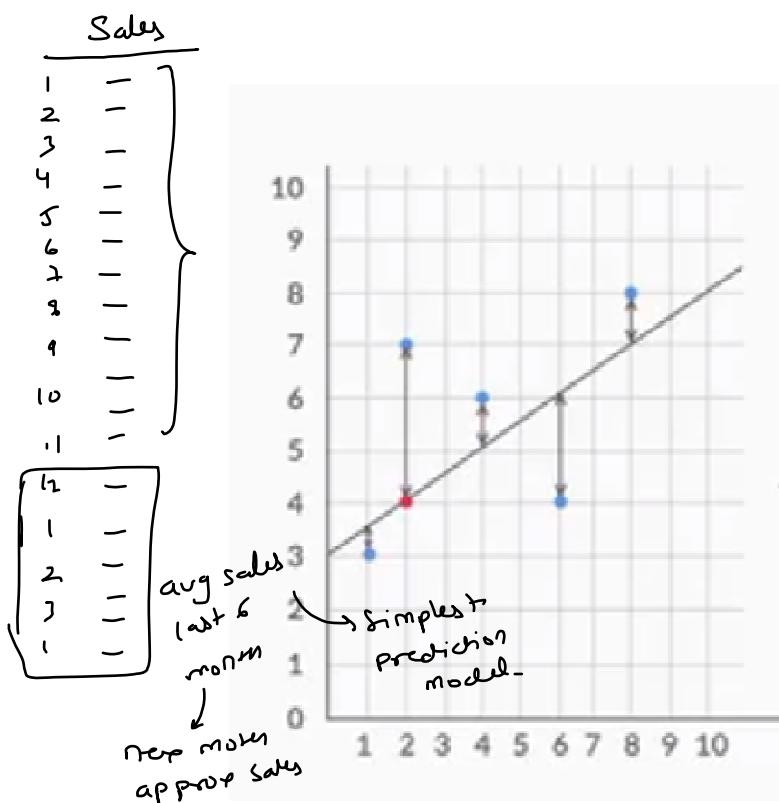
$$\text{Total sum of squares} - TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

↳ Error by the aug model

SLR

Metric that can help us to understand if the regression model is good or bad.

upGrad



Best Fit Line

$$Y = \beta_0 + \beta_1 X$$

↓ ↓

Intercept Slope

$$e_i = Y_i - Y_{\text{pred}}$$

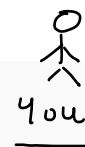
Ordinary Least Squares Method:

↓

$$e_1^2 + e_2^2 + \dots + e_n^2 = \text{RSS} \text{ (Residual Sum Of Squares)}$$

$$\text{RSS} = (Y_1 - \beta_0 - \beta_1 X_1)^2 + (Y_2 - \beta_0 - \beta_1 X_2)^2 + \dots + (Y_n - \beta_0 - \beta_1 X_n)^2$$

$$\text{RSS} = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$



shornaji
ka ladka.
reference

avg model

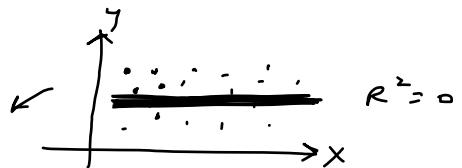
SLR



$$RSS = TSS \quad R^2 = 0$$

RSS < TSS ; improvement

upGrad



RSS and TSS

Coefficient of determination

$$R^2 = 1 - \frac{RSS}{TSS} \rightarrow \begin{array}{l} \text{Error by regression model} \\ \text{corr by avg-model.} \end{array}$$

80 %

$$TSS = (Y_1 - \bar{Y})^2 + \dots + (Y_n - \bar{Y})^2$$

$$\text{Or } \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$R^2 = 1 - \frac{RSS}{TSS}$$

Where

RSS - Residual sum of squares

TSS - Total sum of squares

Demo using: <https://da-upgrad.shinyapps.io/rsquared/>

[Your regression model is 80%.
better than the avg-model.]

TSS represent error of the very basic model that
we can build without having any independent variable.
So any model that we build should be better from this basic model.

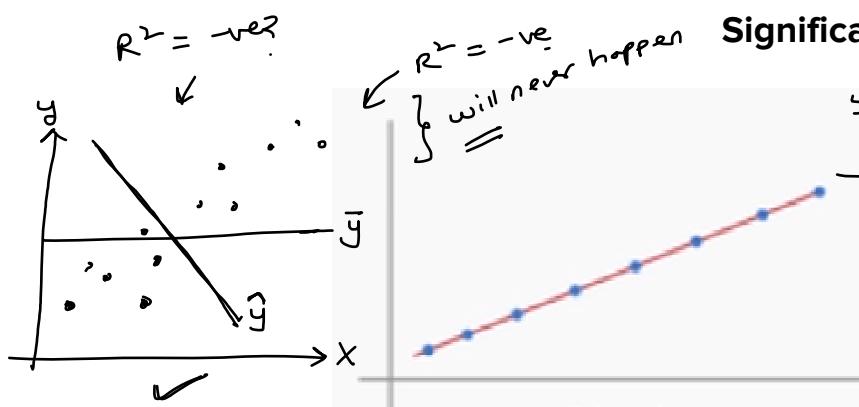
SLR

$$\overline{TSS} < \overline{RSS}$$

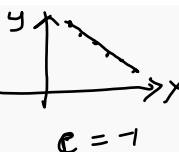
$$R^2 = 1 - \frac{RSS}{TSS}$$

upGrad

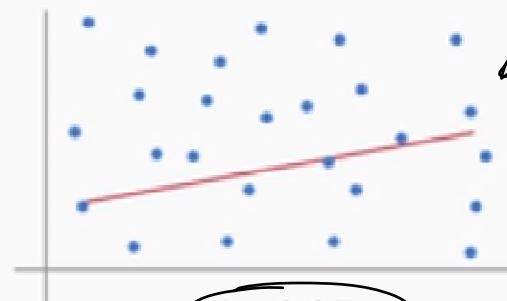
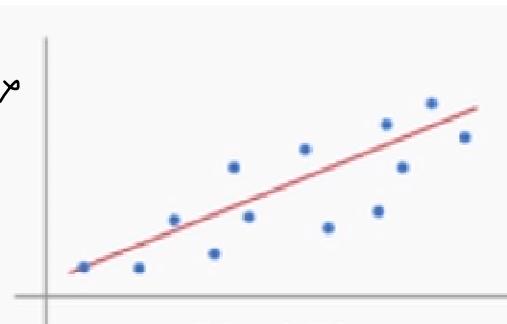
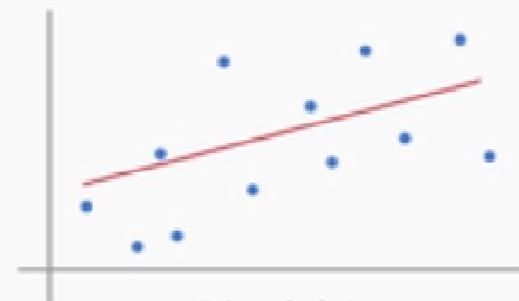
$$\downarrow R^2 \propto \frac{1}{TSS} \uparrow$$



Significance of R-Square

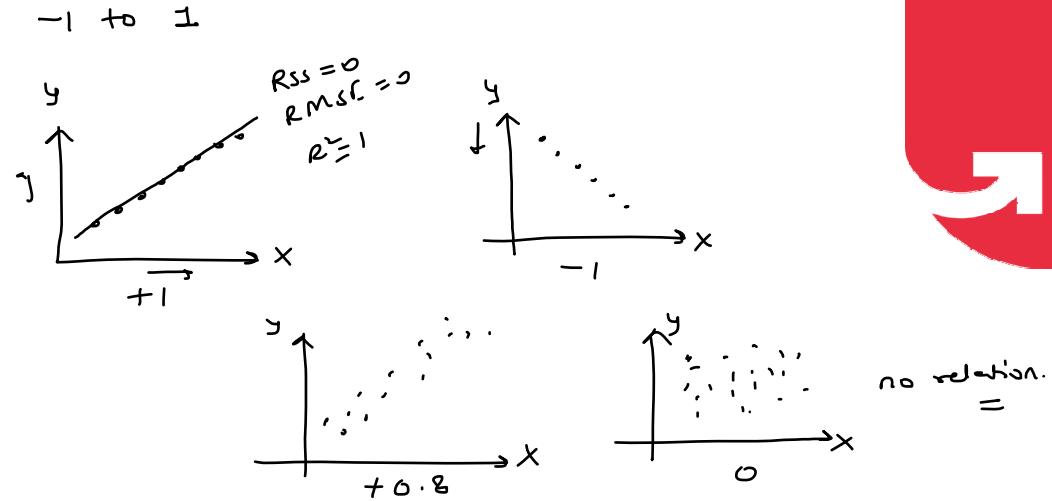
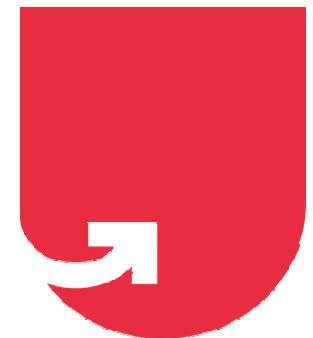


$$R^2 = 1$$



$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

The smaller the means squared error, the closer you are to finding the line of best fit.



Let's Test our Knowledge

Let's Test

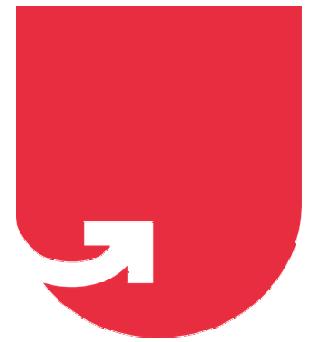
$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{2.75}{6.8} = 0.59 \approx 60\%$$

upGrad

$$\frac{\sum y_i}{n}$$

\Rightarrow my regression model is 60% better than the avg model.

X	\bar{Y}_i	$\hat{y} = 1.5 + 1.5x$	$(y_i - \hat{y})$	$(y_i - \hat{y})^2$	Distance between Y values and their mean $(y_i - \bar{y})$	$(y_i - \bar{y})^2$
		Predicted \hat{Y} Value	Error	Error Squared	Mean distances squared	
1	3	$1.5 + 1.5 \times 1 = 3$	$3 - 3 = 0$	0	$3 - 4.8 = -1.8$	3.24
2	4	$1.5 + 1.5 \times 2 = 4.5$	$4 - 4.5 = -0.5$	0.25	$4 - 4.8 = -0.8$	0.64
2	5	$1.5 + 1.5 \times 2 = 4.5$	$5 - 4.5 = 0.5$	0.25	$5 - 4.8 = 0.2$	0.04
3	6	$1.5 + 1.5 \times 3 = 6$	$6 - 6 = 0$	0	$6 - 4.8 = 1.2$	1.44
4	6	$1.5 + 1.5 \times 4 = 7.5$	$6 - 7.5 = -1.5$	2.25	$6 - 4.8 = 1.2$	1.44
Mean:	$\bar{y} = 4.8$	Sum: $\sum_{RSS} 2.75$		Sum: $\sum_{TSS} 6.8$		



Thank You!

21

23/05/19