

Data Science: Linear Regression Doubts and Concepts

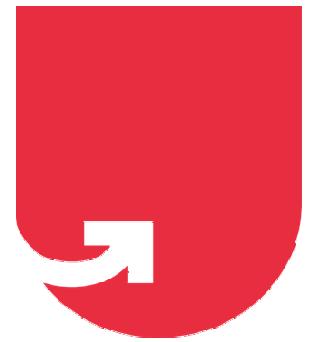
Course : Data Science

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



Introduction to Linear Regression

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

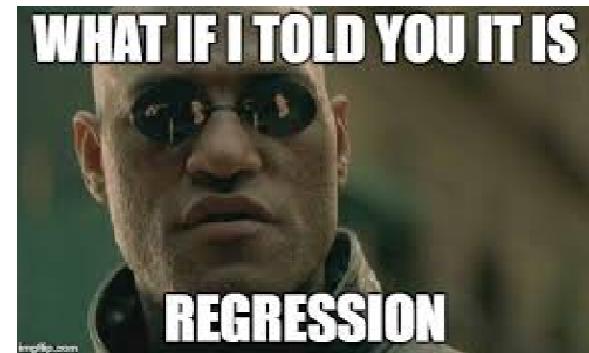
What factor will you be considering to open a new store?

- 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

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 = 10X₁ + 20X₂ + 0.5X₃ + 8X₄ +.....

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

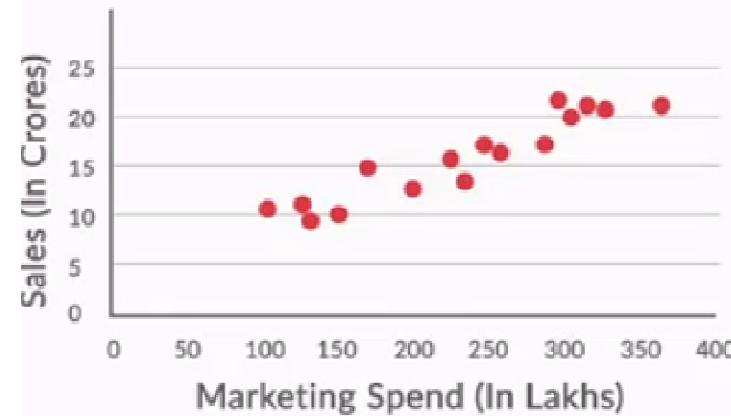
We have historic data of Marketing Spend and Sales, Let's see that

| 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.

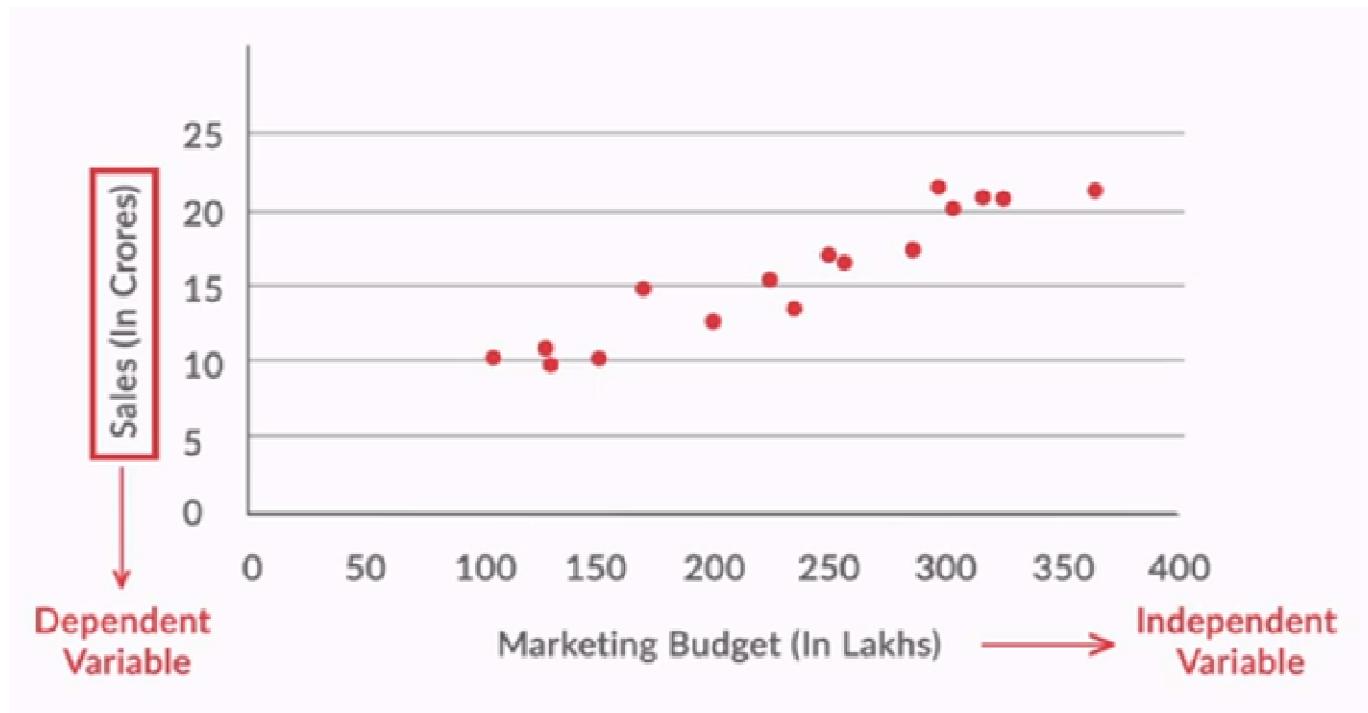
Simple Linear Regression

Any relationship between Sales and Marketing Budget?



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.

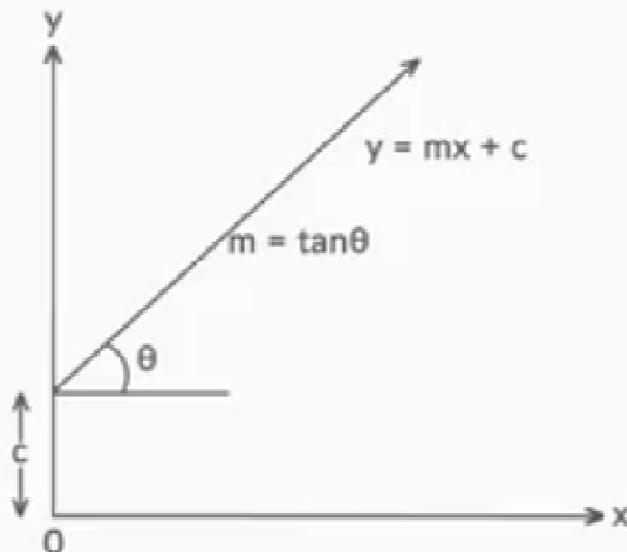
Simple Linear Regression



How to Make Prediction?



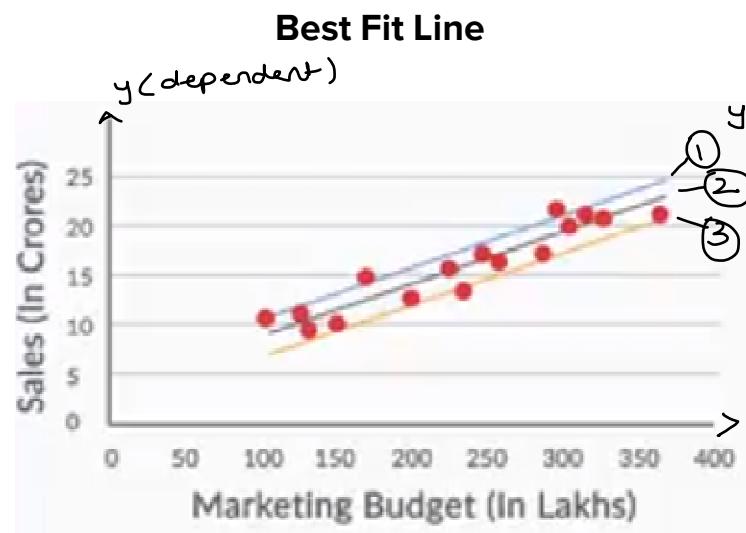
EQUATION OF STRAIGHT LINE



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

SLR

upGrad
Easy to handle while processing
Smaller terms



$$y = \beta_0 + \beta_1 x$$

$$y = \beta'_0 + \beta'_1 x$$

$$y = \beta''_0 + \beta''_1 x$$

$$RSS_1/n$$

$$RSS_2/n$$

out of all the possible options of β_0 & β_1 , we

RSS_3/n will choose that particular line for which your RSS is minimum.

SLR



upGrad

RESIDUALS

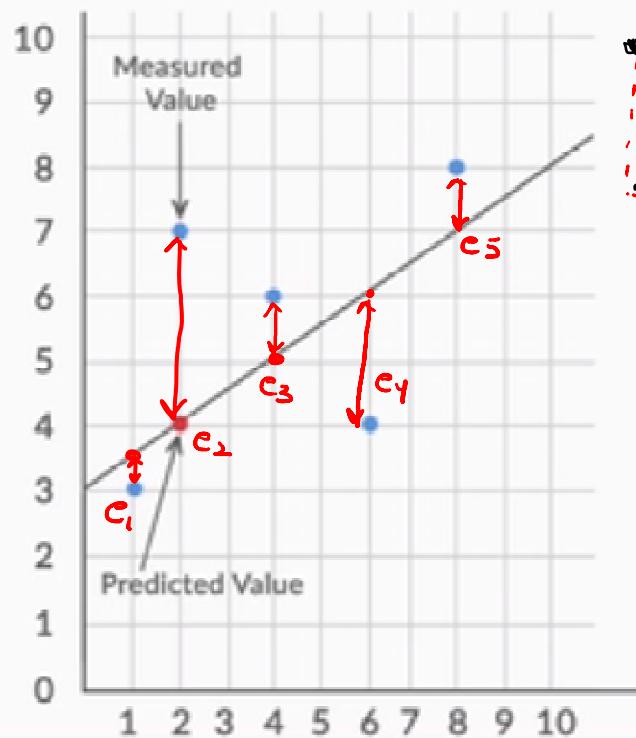
$$RSS = e_1^2 + e_2^2 + e_3^2 + e_4^2 + e_5^2$$

$$e = (y_i - \hat{y})$$

$$\hat{y} = \beta_0 + \beta_1 x$$

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

find out that value of β_0, β_1 for which the RSS term is minimum

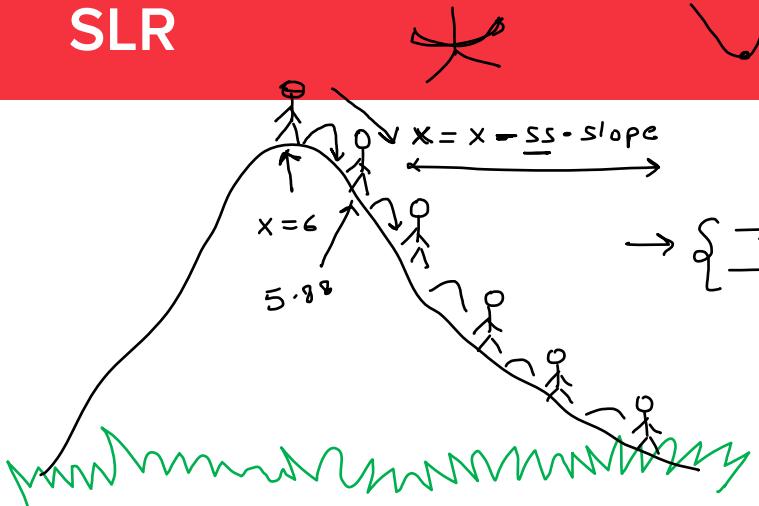


$n = \# \text{ of observation/rows in my data}$

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

Cost function for linear regression.

SLR

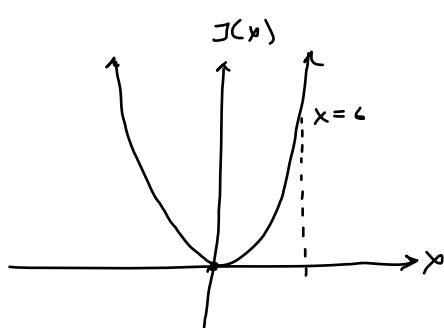


$$x^4 = 4x^3 \frac{d}{dx} x^3$$

upGrad

Gradient Descent → Small step in direction of slope.

optimisation algo that helps you to find the minimum value of a given function.



Slope of any given function is the derivative of that function

$$\begin{aligned}\frac{d}{dx} J(x) &= \frac{d}{dx} x^2 \\ &= 2x^{2-1} \\ &= 2x\end{aligned}$$

$$\boxed{\frac{d}{dx} x^n = n x^{n-1}}$$

$J(x) = x^2$
I need that value of x for which $J(x)$ is minimum

$$\rightarrow x = 6$$

$$\rightarrow \text{next } x = \boxed{x - \underbrace{SS \times \text{slope}}_{\text{current}}}$$

$$\boxed{SS = 0.01} \quad \boxed{\text{slope} = 2x}$$

$$= 6 - 0.01 \times (2 \times 6)$$

$$= 6 - 0.12$$

$$= 5.88$$

Iteration
1000

$$5.88 - SS \times \text{slope}$$

$$\boxed{x = 0}$$

$$\begin{aligned} J(x) &= x^2 \\ &= 2x^{n-1} \\ &= 2x \end{aligned}$$

$$\begin{aligned} J(x) &= (x+2x)^2 \\ &= 2(x+2x)^{2-1} \frac{d}{dx}(x+2x) \\ &= 2(x+2x) \left\{ \frac{d}{dx}x + \frac{d}{dx}2x \right\} \\ &= 2(x+2x)(1+2) \\ \boxed{\frac{d}{dx}J(x)} &= 6(x+2x) \end{aligned}$$

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

find out that value of
 β_0 and β_1 for which
 $J(\beta_0, \beta_1)$ is minimum

Slope wrt β_0

$$\begin{aligned} \frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1) &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta_0} (y_i - (\beta_0 + \beta_1 x_i))^2 \\ &= \frac{1}{n} \sum_{i=1}^n 2(y_i - (\beta_0 + \beta_1 x_i)) \frac{\partial}{\partial \beta_0} (y_i - (\beta_0 + \beta_1 x_i)) \\ &= \frac{1}{n} \sum_{i=1}^n 2(y_i - (\beta_0 + \beta_1 x_i)) \left\{ \frac{\partial}{\partial \beta_0} y_i - \frac{\partial}{\partial \beta_0} \beta_0 - \frac{\partial}{\partial \beta_0} \beta_1 x_i \right\} \\ &= \frac{1}{n} \sum_{i=1}^n 2(y_i - (\beta_0 + \beta_1 x_i)) \left\{ 0 - 1 - 0 \right\} \\ &= -\frac{2}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) = \boxed{-\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y})} \end{aligned}$$

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

$$\begin{aligned}
 \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1) &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta_1} (y_i - (\beta_0 + \beta_1 x_i))^2 \\
 &= \frac{1}{n} \sum_{i=1}^n 2(y_i - (\beta_0 + \beta_1 x_i)) \frac{\partial}{\partial \beta_1} (y_i - (\beta_0 + \beta_1 x_i)) \\
 &= \frac{2}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) \left\{ \frac{\partial}{\partial \beta_1} y_i - \frac{\partial}{\partial \beta_1} \beta_0 - \frac{\partial}{\partial \beta_1} \beta_1 x_i \right\} \\
 &= \frac{2}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) \left\{ 0 - 0 - x_i \right\} \\
 &= -\frac{2}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i)) (x_i) \\
 &\boxed{= -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}) (x_i)}
 \end{aligned}$$

$$\begin{aligned} \text{corr}(x_2, x_3) & \\ \text{corr}(x_2, x_4) & \end{aligned} \left. \begin{array}{l} \text{Limited to only} \\ \text{two series.} \end{array} \right\} \quad [x_2 \sim x_3, x_4, x_1]$$

- Rock band
- 2 singers $\begin{cases} \text{male} \\ \text{female} \end{cases}$
 - 1 Drum
 - 1 keyboard
 - 2 guitars $\begin{cases} \text{Same speed} \\ \text{Same pitch} \\ \text{Some tone} \end{cases}$
 - \downarrow
 - delivering the same information

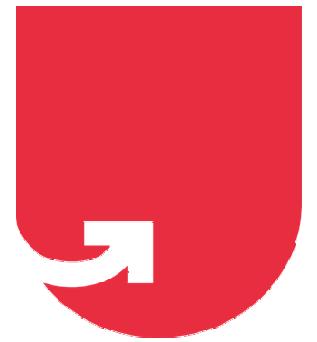
Multicollinearity : The presence of correlated independent features.

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 \quad \left. \begin{array}{l} R^2 = 89\% \\ \text{using only independent variable} \end{array} \right\}$$

| | R^2 | VIF |
|---|---|-----|
| 1 | $x_1 = \alpha_0 + \alpha_1 x_2 + \alpha_2 x_3 + \alpha_3 x_4$ | 10% |
| 2 | $x_2 = \alpha'_0 + \alpha'_1 x_1 + \alpha'_2 x_3 + \alpha'_3 x_4$ | 30% |
| 3 | $x_3 = \alpha''_0 + \alpha''_1 x_1 + \alpha''_2 x_2 + \alpha''_3 x_4$ | 50% |
| 4 | $x_4 = \alpha'''_0 + \alpha'''_1 x_1 + \alpha'''_2 x_2 + \alpha'''_3 x_3$ | 80% |

$[x_1, x_2 \text{ and } x_3 \text{ together are able to explain } 80\% \text{ of } x_4]$

$\text{VIF} \geq 5$] Higher multicollinear



Thank You!

20

23/05/19



What is Multicollinearity?

Imagine you went to watch a rock band's concert. There are 2 singers, a drummer, a keyboard player, and 2 guitarists. You can easily differentiate between the voice of singers as one is male and other is female but you seem to have trouble telling who is playing better guitar. Both guitarists are playing on the same tone, same pitch and at the same speed. If you could remove one of them then it wouldn't be a problem since both are almost same. The benefit of removing one guitarist is cost-cutting and fewer members in the team. In machine learning, it is fewer features for training which leads to a less complex model.

Consider the simplest case where Y is regressed against X and Z and where X and Z are highly positively correlated. Then the effect of X on Y is hard to distinguish from the effect of Z on Y because any increase in X tends to be associated with an increase in Z.

Another way to look at this is to consider the equation. If we write $Y=b_0+b_1X+b_2Z+e$, then the coefficient b_1 is the increase in Y for every unit increase in X while holding Z constant. But in practice, it is often impossible to hold Z constant and the positive correlation between X and Z mean that a unit increase in X is usually accompanied by some increase in Z at the same time.



What is VIF?

A variance inflation factor(VIF) detects multicollinearity in regression analysis. Multicollinearity is when there's correlation between predictors (i.e. independent variables) in a model; it's presence can adversely affect your regression results. The VIF estimates how much the variance of a regression coefficient is inflated due to multicollinearity in the model.

$$\text{VIF} = \frac{1}{1 - R_i^2}$$

Variance inflation factors range from 1 upwards.

- 1 = not correlated.
- Between 1 and 5 = moderately correlated.
- Greater than 5 = highly correlated.