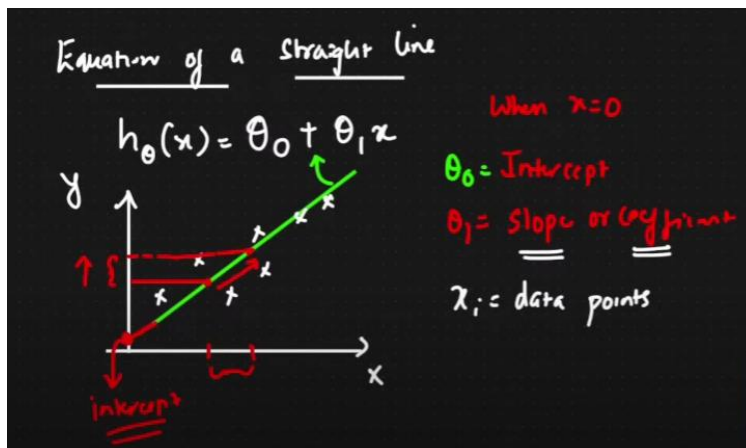
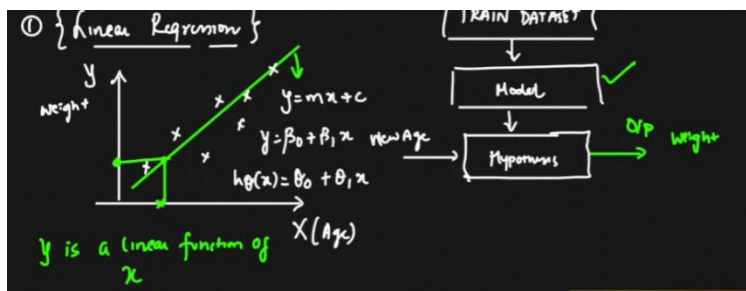
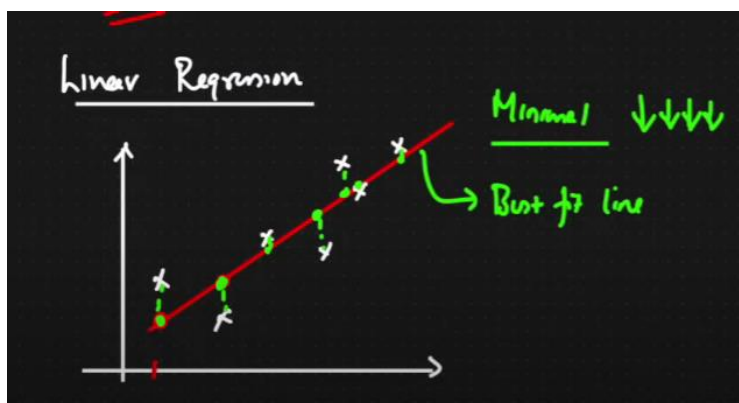


## 1) Linear Regression:



$\theta_0$ : at what point you are meeting the y-axis (intercept)

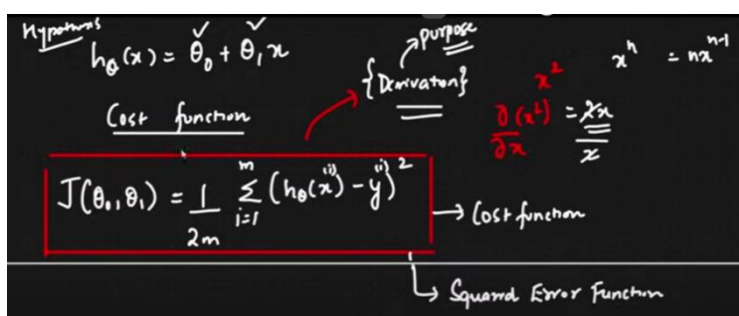
$\theta_1$ : unit movement of x & y axis



Best fit line: the distance between the predicted & the real points should be minimum!

Will get it using, cost function also called as a distance formula and a squared error function!

green dot: real pts. & white circle: predicted pts.



- 1) predicted pts – real pts:  $h_0(x) - y$
- 2) **squaring**: to get rid of -ve values
- 3)  $\sum$ : summation  $i = 1$  to complete  $m$ ,  $m$ : no. of data pts, distance between predicted & real pts
- 4)  $1/m$ : average,  $1/2m$ : derivation purpose, equation simpler

Later, will minimize the cost function ( $\theta_0, \theta_1$ ) as it's also called Squared Error Function!

What we need to solve

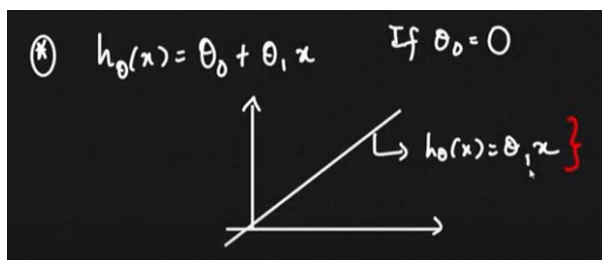
$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

↓

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

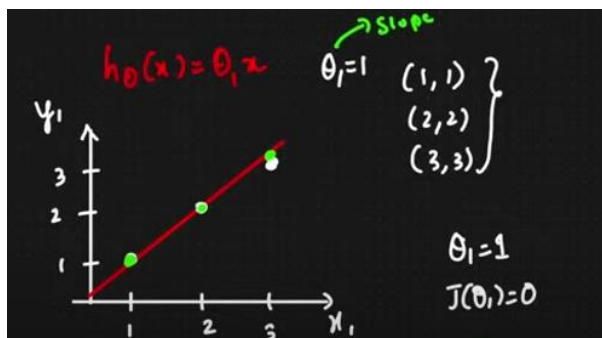
2) Comparing w.r.t hypothesis testing & cost function with an example:

when,  $\theta_0 = 0$ :



i) when,  $\theta_1 = 1$ : ( $\theta_0$  passing through origin) (green circle: real pts.; white circle: predicted pts.)

hypothesis testing:

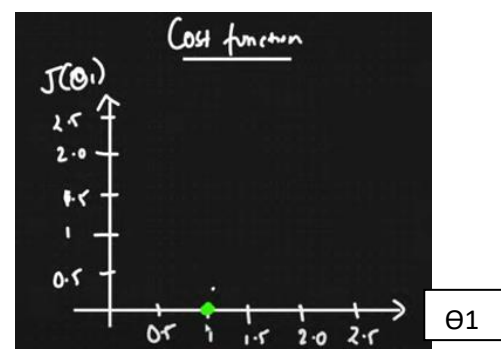


cost function:

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^3 (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

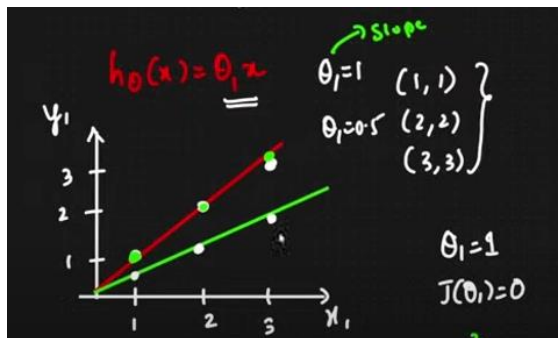
$$= \frac{1}{2m} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$$J(\theta_1) = 0$$



ii) when,  $\Theta_1 = 0.5$ :

hypothesis testing: (green line), (green circle: real pts.; white circle: predicted pts.)

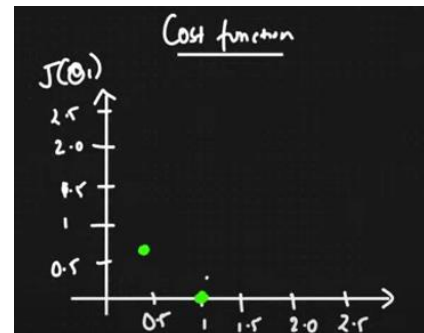


cost function:

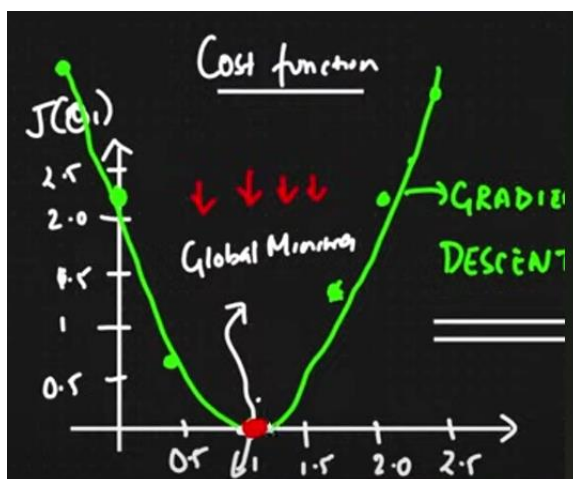
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{2 \times 3} [0.25 + 1 + 2.25] \approx 0.58$$



3) Gradient Descent:



“GRADIENT DESCENT” helps to make sure that we get the right  $\Theta_1$  value!

The most suitable point is the point which has the minimum distance between the real pt and the predicted pt., also known as “GLOBAL MINIMA”

Here, it's 1!

Here, we are assuming  $\Theta_1$  value and moving ahead but this isn't the right way! The right way should be to find the global minima in the least attempts; so, for this will use something called as: “CONVERGENCE ALGORITHM”.

4) Convergence Algorithm: Repeat the equation until it's convergence!

Convergence Algorithm

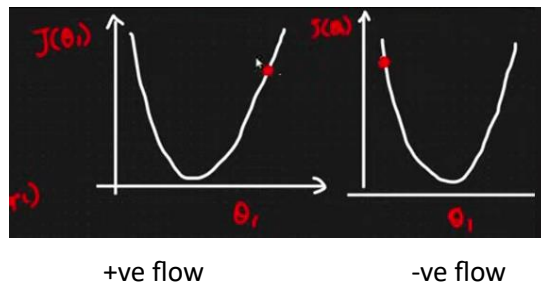
Repeat until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j}$$

*derivative(slope)*

$\alpha$  = learning rate

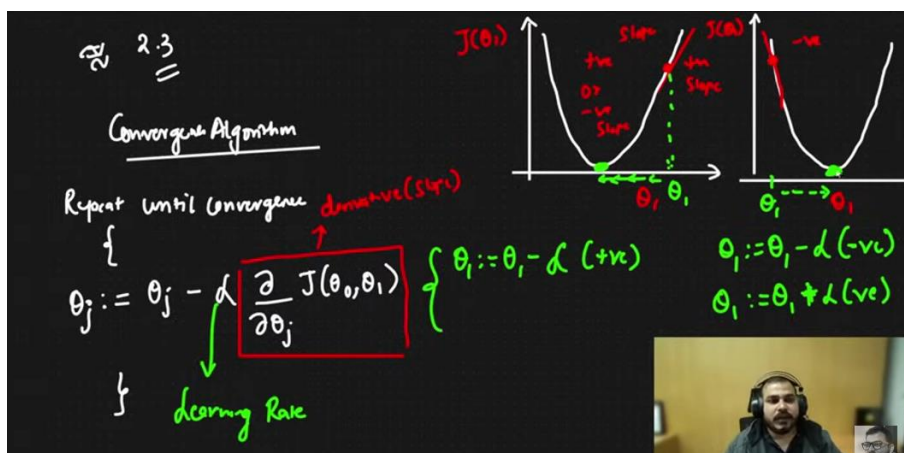
1) find the slope:



i) initial pts marked with red!

ii) apply the derivative on the red pt to find the slope (right: +ve, left: -ve)

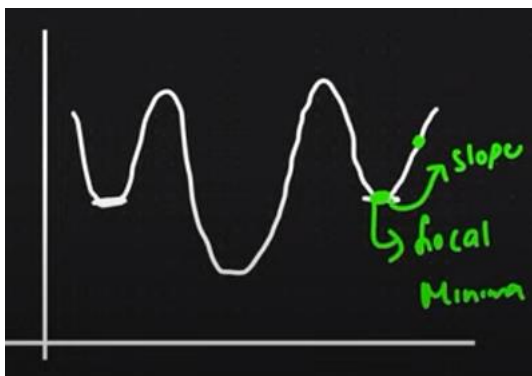
iii) apply convergence algorithm to get the global minima



iv) what is  $\alpha$  a learning rate?

- ➔ Due to the speed the pt reaches from the initial to the global minima!
- ➔ It shouldn't be bigger nor smaller!
- ➔ So, it's good to keep  $\alpha = 0.01$

5) What if the cost function has the local minima?



$$\theta_1 := \theta_1 - \alpha(0)$$

$$\theta_1 := \theta_1$$

i) usually, in gradient descent and this type of equations (cost function) we don't find any issue, as we get a:



ii) but in the deep learning when we are learning about gradient descent and ann, we get lots of local minima! Because of that we have various gradient descent algorithms, like: rms prop, adam optimizers which will solve the specific problem!

6) Gradient Descent Algorithm: (MOST IMPORTANT, USES ON DAILY BASIS)

Repeat until Converges

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

Find:

$$\bar{J} = 0 \text{ and } 1$$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

If  $j = 0$ :

$$j=0 \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$1/m$  instead of  $1/2m$  just because when will find derivative of:

$$\frac{1}{2m} x^2 \text{ will get } \rightarrow \frac{x}{m}$$

And so, we have written:  $\frac{1}{m}$

$$j=0 \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x - y^{(i)})^2$$

Here, we have written the hypothesis testing formula, in the place of  $h_{\theta}(x)$ !

Similarly:

$$\left\{ \begin{array}{l} j=0 \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ j=1 \Rightarrow \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \end{array} \right\}$$

$h_{\theta}(x) = \theta_0 + \theta_1 x$   
 $x^2 \rightarrow 2x$

Here,  $j=1$  has  $x^i$  in the end; as we have found the derivative of  $h_{\theta}(x)$  from  $j=1$ , removing the square from  $j=1$ , similarly did for  $j=0$ , in  $j=0$  there's no  $x^i$  as its 0!

Finally, the convergence algorithms are:

$$\begin{array}{l} \text{Repeat until converge} \\ \left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \end{array} \right. \end{array}$$

Will have lots of convex functions when we have multiple features:  $x_1, x_2, x_3, \dots, x_n$

Here, gradient descent will look like a 3D curve:



There will be lots of slopes & our aim will be to get the global minima, down the line!



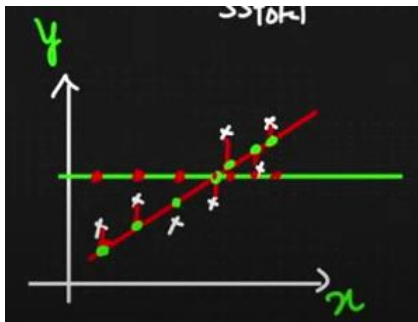
7) Performance Matrix: a)  $R^2$  & b) Adjusted  $R^2$

a)  $R^2$  used to verify, how good our model is w.r.t linear regression!

$R^2$  and Adjusted  $R^2$

$$R^2 = 1 - \frac{SS_{Res}}{SS_{Total}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

$\uparrow$   
 $h_0(x)$



i) green dots: predicted, green line: mean

ii) Which one will be high?

This:  $\sum (y_i - \hat{y}_i)^2$  or this:  $\sum (y_i - \bar{y})^2$

second one; mean of a particular distance will always be higher!

$$R^2 = 1 - \frac{SS_{Res}}{SS_{Total}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = 1 - \left[ \frac{\text{Low}}{\text{High}} \right] = \text{Big number}$$

$\nearrow$  Small number

iii) with the above given steps, will get  $R^2$  as the bigger value, when we subtract from 1

iv) as much bigger the  $R^2$ , better the model!

v) can we get  $R^2$  as a -ve number?

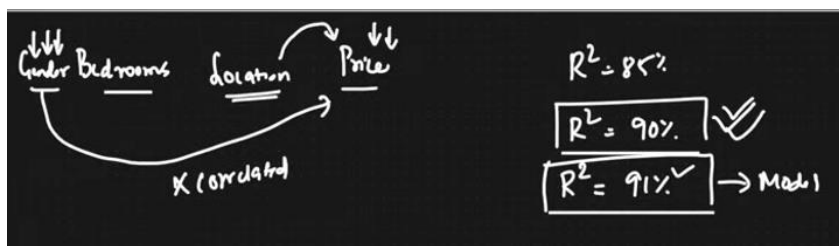
→ it's only possible when the best fit line is above the mean, due to which we get:

$$1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\text{High}}{\text{Low}}$$

usually, it doesn't happen; at least will try to fit a line which is good, having low distance between predicted and real pts then the mean.

vi) to get a good model and  $R^2$  value, sometimes we take a column which isn't necessary / not correlated to the target variable!

Example:



To get rid of this situation, will use Adjusted  $R^2$

b) Adjusted  $R^2$ :

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

Where

$R^2$  Sample R-Squared

$N$  Total Sample Size

$p$  Number of independent variable

Handwritten notes on a blackboard:

- For  $p=2$ ,  $R^2 = 90\%$ ,  $R^2_{\text{adjusted}} = 86\%$
- For  $p=3$ ,  $R^2 = 91\%$ ,  $R^2_{\text{adjusted}} = 82\%$

Arrows indicate that as  $p$  increases,  $R^2$  increases (upward arrow) but  $R^2_{\text{adjusted}}$  decreases (downward arrow).

Why have the Adjusted  $R^2$  decrease?

When the independent variables are strongly related to the dependent variable, both  $R^2$  and Adjusted  $R^2$  will go up.

However, if an independent variable doesn't really help explain the dependent variable,  $R^2$  might still be high, even though the model isn't that useful.

In this case, Adjusted  $R^2$  will be lower, giving a more accurate picture of how well the model works.

So, Adjusted  $R^2$  helps fix the problem of  $R^2$  being misleading, especially when there are many independent variables, some of which aren't helpful. It's a better way to choose the best model in these cases.

8) Overfitting & Underfitting:

Goal: Using cost function, get the minimum distance and find the best fit lines!

Handwritten formula for the cost function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

when,  $\theta_0 = 0$ ,  $J(\theta_0, \theta_1) = 0$ , we got the best fit line by training the sample data!

But this is not the case in the real-world scenario, we either get good result on the training dataset or on testing dataset; it's impossible to be exactly accurate on the training and testing dataset!

It happens due to overfitting & the underfitting, explained below:



## Overfitting vs Underfitting:

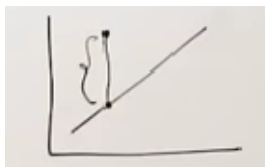
Aspect	Underfitting	Overfitting
What it means	Model is too simple to capture important patterns.	Model is too complex and captures unnecessary details.
Training Performance	Poor (doesn't learn the patterns well).	Excellent (learns the training data too well).
Test Performance	Poor (doesn't generalize well to new data).	Poor (fails to generalize to new, unseen data).
Bias	High (makes strong assumptions and misses patterns).	Low (fits the training data very closely, even the noise).
Variance	Low (model doesn't change much across different data).	High (model's performance is sensitive to the training data).
Example	Child says "all fruit is the same size" and can't tell apples from oranges.	Child memorizes every small detail (like blemishes) and gets confused with new fruit.
Key Issue	Too simple, ignores important details.	Too complex, memorizes irrelevant details.
Ideal Solution	Make the model more complex, include important factors.	Simplify the model, avoid overcomplicating with unnecessary details.

To handle the overfitting and underfitting, will work on regularization (ridge & lasso)!

## 9) Regularization:

- i) used to reduce the complexity of the model (mainly the training phase)
- ii) shrink the magnitude value of a regression coefficient, if high!
- iii) make it computational (should be cheap & take no time)

### a) Ridge Regression (L2 regularization):



$$\text{Ridge } R = \text{Loss} + \alpha \|W\|^2 \quad (\text{penalty})$$

\*loss = predicted value – actual value!

\*to compensate the loss, we have added the penalty; due to this we

reduced the loss by some amount & shrink the coefficient magnitude value

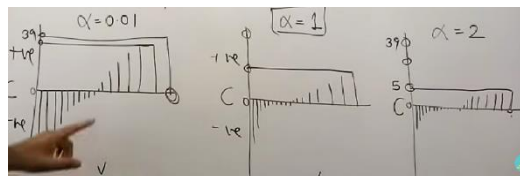
\* $\alpha$  = constant (can be any value); \* $\|w\|^2$ : (vector of the coefficient) =  $w_1^2 + w_2^2 + w_3^2 + \dots + w_n^2$

### Example:

$$y = 0.9 + 1.2x_1 + 20x_2 + 39x_3$$

$$y = 0.9 + 0.7x_1 + 2x_2 + 5x_3$$

→



as,  $\alpha$  value increases the coefficient of the magnitude decreases/shrinks/scale downs almost to 0;

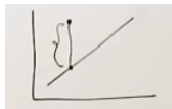
$x_3 = 39$  decreases to 5!

Achieved the regularization, by reducing the complexity & by reducing the expenses of the computational. Prevent overfitting!

b) **Lasso Regression** (L1 regularization):

$$\text{Lasso } R = \text{Loss} + \alpha \|W\|$$

(penalty)



\*loss = predicted value – actual value!

\*to compensate the loss, we have added the penalty; due to this we reduced the loss by some amount & shrink the coefficient magnitude value

\* $\alpha$  = constant (can be any value)

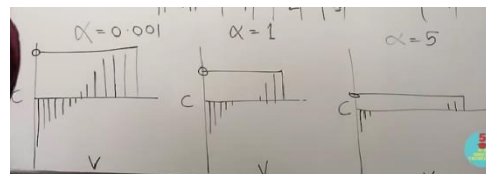
$$\|W\| = |w_1| + |w_2| + |w_3| + \dots + |w_n|$$

Example:

$$y = 0.9 + 1.2x_1 + 20x_2 + 39x_3$$

$$y = 0.9 + 0.7x_1 + 2x_2 + 5x_3$$

→



as,  $\alpha$  value increases the coefficient of the magnitude decreases/shrinks/scale downs to exactly 0;  
 $x_2 = 20$  decreases to 0!

Achieved the regularization, by reducing the complexity & by reducing the expenses of the computational.

Also, those who gets 0, those variables/features are not important and can be directly eliminated w.r.t variable  $y$  (a dependent variable), known as **feature selection**.

Prevent Overfitting & achieves Feature Selection!

- Try both the regularization and the one with good performance matrix, select/use that!

**Keep checking all the below given assumption in the start, mid or in the end; whenever required!**

10) Assumption of Linear Regression:

i) if our features follow: Normal / Gaussian Distribution → model will get trained well!

if it doesn't follow: Normal / Gaussian Distribution, will try to do feature transformation to convert the data into Normal / Gaussian Distribution, if possible.

iii) Standardization: Scaling the data → z-score,

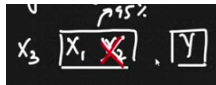
$$\mu = 0, \sigma = 1$$

use whenever there's Gradient Descent and find the global minima value, as soon as possible!

Helps to optimize the model & use to increase the training time of the model.

iv) Linear regression works w.r.t linearity.

v) Check multicollinearity:



The diagram illustrates multicollinearity. It shows a set of features  $x_1$  and  $x_2$  enclosed in a box, with a red 'X' between them, indicating a problem. Above the box, it is handwritten that they are 95% correlated. To the left of the box is  $x_3$ , and to the right is the target variable  $y$ .

over here, if we try to see the collinearity of  $x_1$  &  $x_2$ , and we get to know it's 95% correlated, will use just one among the 2 features & do the prediction!

**Variation Inflation Factor (VIF)** is also used to solve the multicollinearity.

There's also **Homoscedasticity**!

Reference:

- 1) [Krish Naik - Linear Regression \(from 33 minutes onwards\)](#)
- 2) [5 Minutes Engineering - Overfitting & Underfitting](#)
- 3) [5 Minutes Engineering - Ridge Regression Explanation](#)
- 4) [5 Minutes Engineering - Lasso Regression Explanation](#)
- 5) [Krish Naik - Assumption of Linear Regression \(from 41:36 minutes onwards\)](#)