

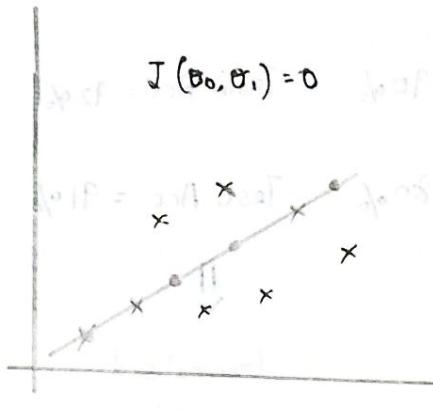
## Lecture - 02:

- Ridge and Lasso Regression .
- Assumption of Linear regression .
- Logistic Regression .
- Confusion Matrix .
- Practicals for linear, Ridge, Lasso and Logistics.

## Ridge and Lasso Regression:

Let's take,

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



Suppose we trained  
the model with help  
of subset of data.  
And also fit the  
line.

Then, new data points are come (x). But the difference between new data points (testing data) and predicted points are huge differences.

It is categorized into two conditions,

### Overfitting:

Model performs well in Training data. [Low Bias]

fails to perform well in Testing data [High Variance]

### Underfitting:

Model Accuracy is bad with training data. [High bias]

Model Accuracy is also bad with Test data. [High Variance]

Bias - The error between average model prediction and the ground truth.

Let's Consider Some Models,

Model 1

Model 2

Model 3

Training Accuracy = 90%

Train Acc = 92%

Train Acc = 70%

Testing Accuracy = 80%

Test Acc = 91%

Test Acc = 65%

Overfitting

Generalised  
model

Underfitting

[Low Bias, High Var] and [Low Bias, Low Variance]

[High bias,  
High Var]

Let Consider the One Scenario,

Assumption,

$$J(\theta_0) = 0$$

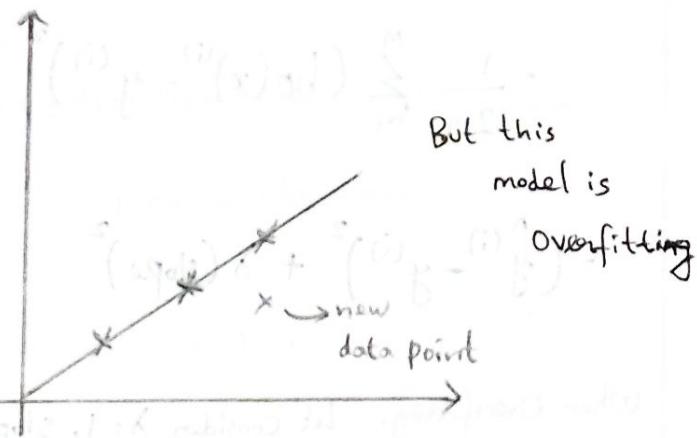
$$\text{Cost function} = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$h_\theta(x) = \hat{y}$$

$$(h_\theta(x) - y)^2 = (\hat{y} - y)^2 = 0$$

The value comes to the zero. Because the  $J(\theta_1) = 0$

With below graph,



Because of overfitting, we conclude this model value does not come to zero. Because, the new data points

and cost function gives different value after applying a best fit line. We need to minimize the cost function.

Then only get generalized model.

So, we need to minimize and prevent overfitting

Conditions. Therefore ridge regression will come

to the picture.

Ridge Regression: [L<sub>2</sub> Regularization]

Ridge Regression Says that the,

$$[(\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2]$$

The Cost function,

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

When,  $J(\theta_1) = 0$

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

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

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

$$= (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2$$

$$\hat{y}^{(i)} = h_{\theta}(x)$$

When Overfitting, Let Consider  $\lambda = 1$ , Slope = 2

$$= 0 + 1(2)^2$$

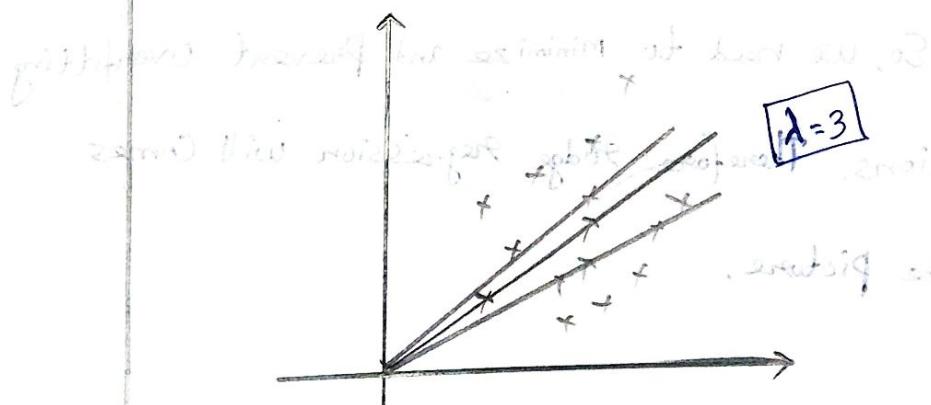
$$= 4 \downarrow$$

Then minimize the Cost function &, In Order to

reduce the overfitting we have to change the value of  $\theta_1$ .

Reduce this the  $\theta_1$  will be change.

Let's see how the cost function changes with different values of  $\theta_1$ .



Try multiple Combination of  $\theta_1$  Value. To fit the line and Avoid Overfitting.

$$(optimal \theta_0 + \theta_1 x^{(i)})$$

So, take the next fit line (i.e) change  $\theta_i$  value,

$$\text{Cost of fit} = (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2$$



Obviously the value to be small compare than the previous iteration.

$$= (\text{small value}) + 1(1.36)^2$$

$$\approx 3 \downarrow$$

In this way, The Cost function reduce and fit the best line.

After many iteration, the Value will be small and fit the best line.

$\lambda$  is nothing but Hyperparameter.

iterations {Hyperparameter} ( $\lambda$ )

$\lambda$  decides the how the line steeps makes the

line better with the help of iterations.

$\theta_i$  Value Change is based on the Convergence

Algorithm.

$\lambda \rightarrow$  Hyperparameter

$R^2$ , adjusted  $R^2$  is based on the hyperparameter

Lasso (L1 Regularization):

Lasso Regression

The L1 regularization says that the,

$$= (\hat{y} - y)^2 + \lambda |\text{slope}|$$

↳ It is also used for  
feature selection.

So, the equation will be,

$$h_{\theta}(x) = \hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$= (\hat{y} - y)^2 + \lambda |\theta_0 + \theta_1 + \theta_2 + \theta_3 + \theta_4 + \theta_5 + \dots + \theta_n|$$

In this case, we are not squaring the slope. Instead we are using Mod. So, it is neglecting the features those are not important to perform.

The important of the L1 regularization are the,

1. Prevent Overfitting.  
 2. Feature Selection.

$\rightarrow$  L1 Regularization.

The hyperparameter ( $\lambda$ ) is assigned something

Called as Cross Validation.

In real time, Most of the cases we are performing the L1 regularization and L2 regularization then use the performance Metrics we can use that.

Ridge Regression (L2 Regularization) Overview:

It is also called as L2 Normalization.

$$\text{Cost function} = (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda (\text{slope})^2$$

Purpose: Prevent Overfitting.

Lasso Regression (L1 Regularization) Overview:

It is also called as L1 Normalization.

$$\text{Cost function} = (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda |\text{slope}|$$

Purpose:

$\rightarrow$  Prevent Overfitting.

$\rightarrow$  Feature Selection.

## Assumption of Linear Regression:

1. Normal/Gaussian distribution  $\rightarrow$  Model will get trained well.

2. Standardization {Scaling data}  $\rightarrow$  Applying Z-Score

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

$\rightarrow$  If the distribution is not Gaussian, Make the normal distribution Using mathematical and Statistical methods.

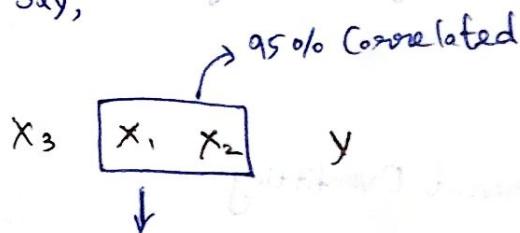
$\rightarrow$  Standardization is Compulsory ?

If the Standardization is not used, The value of gradient descent points to be large. So, It takes to long train time to train the algorithm.

3. Data points Should be linearity.

4. Multi Collinearity.

Let Say,



This feature is Correlated to  $(x)$ . Definitely any One of the feature is dropped.

No need to take both the features for the prediction.

5. Homoscedasticity. [It means two group Variance is

approx Similar]

6. Variance Inflation factor.

Variance Inflation factor:

We need to understand VIF we know,

$$1 - R^2$$

2-Simple Mathematics

$$\left[ \begin{array}{c} x \\ y \end{array} \right]$$

if X Constant, y increases. The value should

be an decrease. Example:  $\frac{2}{6} = 0.333$

if y constant, x increases is just Vice-Versa.

VIF definition,

VIF is a term or measurement through

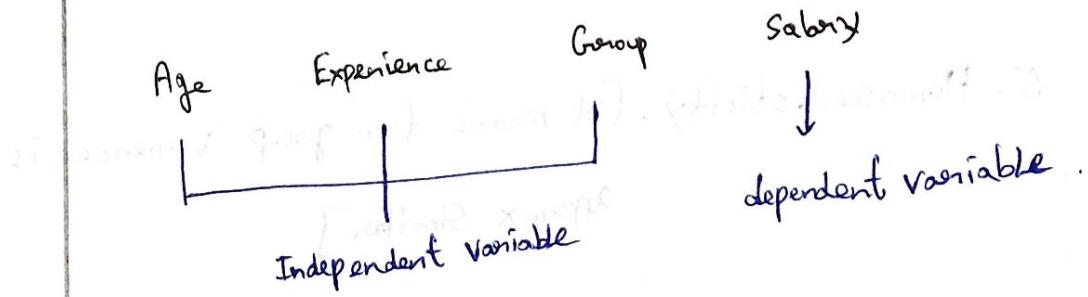
which we can know which variables in the data

highly correlated with other variables.

Example:

Consider the employee dataset.

not much correlation between independent variables  
Features,



Suppose we know One-one Correlation Coefficient, we use  
normally matrix.

like,

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}$$
$$\begin{bmatrix} x_1 & 0.0 & 0.1 & 0.2 \\ x_2 & - & - & - \\ x_3 & - & - & - \end{bmatrix}$$

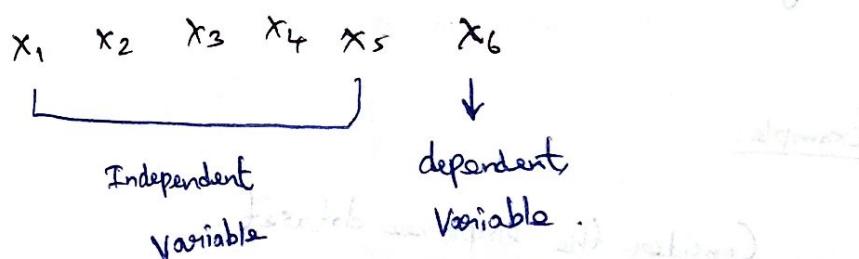
We know only two features Correlation. But the  
VIF is Used to one feature alone and Compare

all the independent features in the dataset and give

Single Coefficient Value.

and all in collision. So it's good not to do this.

Suppose,



We need to compare  $x_1$  feature to all other independent features.

How?

Take the feature as independent Variables,

$$\boxed{x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5}$$

Take  $x_1$  feature and Create Separate model of regression.

After fit linear regression,

$$\boxed{x_1 = \theta_0 + \theta_1 x_2 + \theta_2 x_3 + \theta_3 x_4 + \theta_4 x_5}$$

After fit this model, the  $R^2$  value will be give as

Something.

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

$$\text{If } R^2 \text{ Something } 1.237, \text{ The VIF} = \frac{1}{1 - 1.237} = \frac{1}{-0.237} \\ = -4.2194$$

$1.237 R^2$  is away from the magnitude of  $R^2$ . So, in this situated at very rare cases.

$$\text{If } R^2 \text{ Something } 0.9, \text{ Then } VIF = \frac{1}{1-0.9} = \frac{1}{0.1} = 10$$

This way the VIF is calculated, The final output of all VIF of each feature would be,

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
6	10	2-1	5-6	5-8

The feature  $x_3$  is low VIF, So the feature  $x_3$  will be considered for analysis.

The high VIF will be dropped based on the features.

The  $x_4$  and  $x_5$  have similar same VIF,

Suppose the  $x_5$  is 5-8. Definitely, drop any one of the feature. It is something called as the Multicollinearity.

The VIF is play a vital role in the case of multicollinearity.

Logistic Regression: [Classification Algorithm]

Logistic Regression Very well works in Binary Classification.

It also used to solve multiclassification problems.

Example,

No of Study	No of play	Pass/Fail
-	-	Pass
-	-	Fail

Suppose if a student did not fail after doing 3 hours of study then

Why linear regression is not used in Classification

Problems?

Suppose used linear regression, it would be,

(In dataset says less than 3 hours study leads to fail)

outcome

value of passing always at least 0.5, and we use Linear Regression

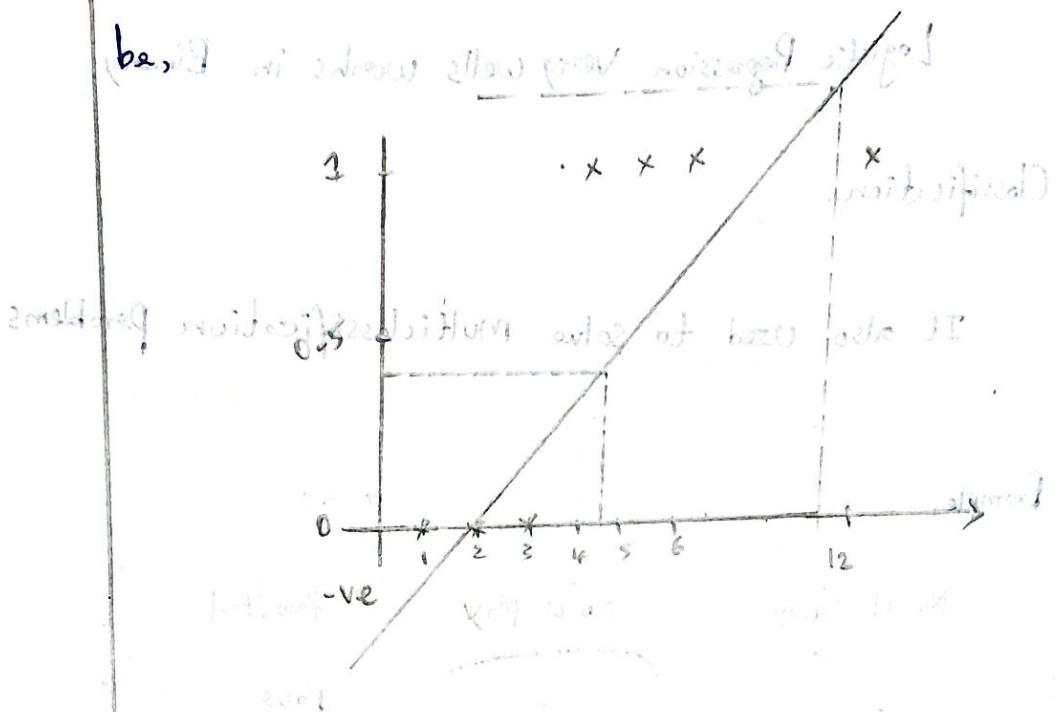
so outcome will be  $h_0(x) \leq 0.5 \Rightarrow 0 \rightarrow \text{fail}$

value of passing always at least 0.5, and we use  $h_0(x) \geq 0.5 \Rightarrow 1 \rightarrow \text{Pass}$



Outcomes based on No. of Study hours

[Outliers are there, The graph should be,]



The linear line change because of outliers. So,

the student study 4.5 hours also fail in this case. But it is not correct. In some times, The Outlier or any related values predict the output as high magnitude and above 1 in the graph and same like sometimes below 0 in the graph. It leads to causes the problem.

To overcome this, we need to squash the function

With the help of Sigmoid function. These two reasons

to avoid linear regression in Classification Problem

for one,

→ Outliers.

→ More over magnitude. [upto 1 and below 0]

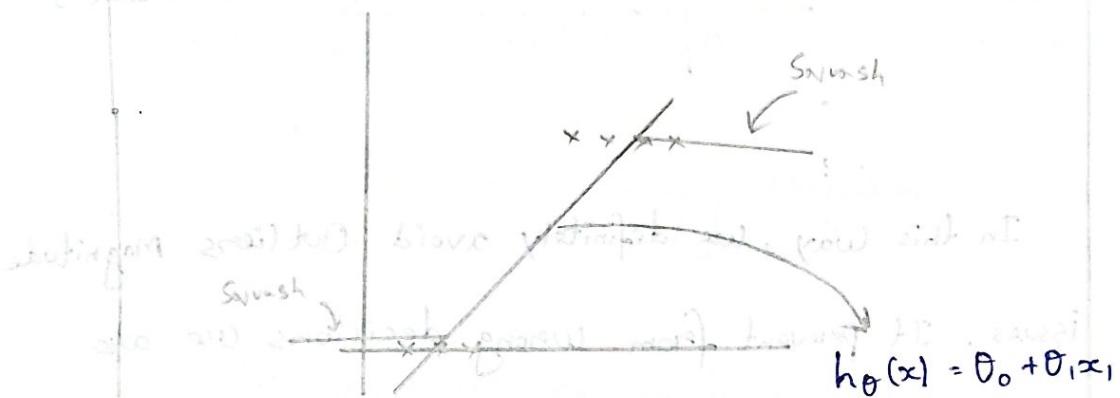
## Decision Boundary Logistic function:

Let,

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n.$$

Written as,

$$h_{\theta}(x) = \theta^T x$$



We need to Squash, So applied

Some function. So,

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1)$$

$$\text{Let } Z = \theta_0 + \theta_1 x_1$$

$$h_{\theta}(x) = g(Z)$$

Sigmoid or logistic function.

$$h_{\theta}(x) = \frac{1}{1 + e^{-Z}}$$

[This function helps to squash]

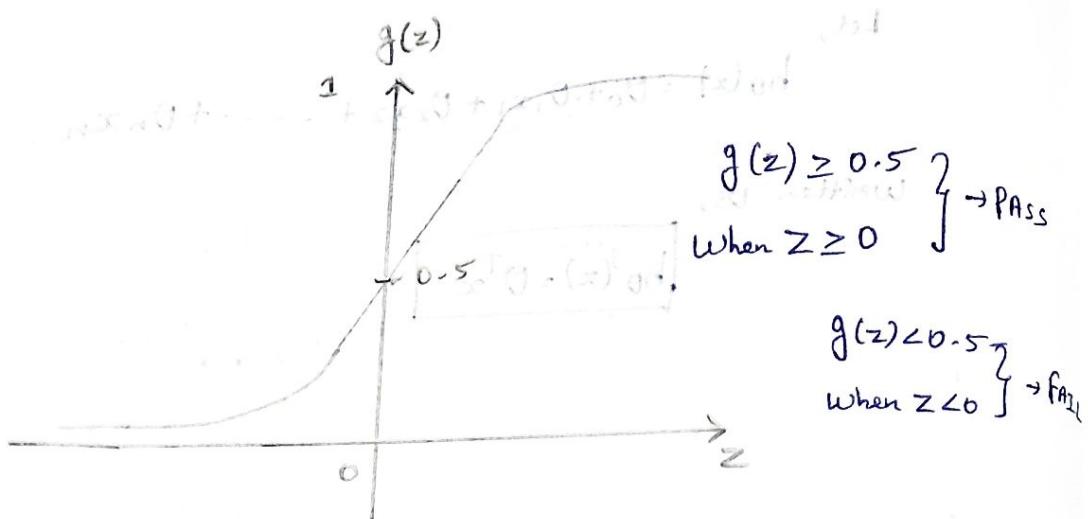
Finally, the equation will be

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

This is known as the logistic function (or)

Regression.

The Sigmoid function look like,



In this way, we definitely avoid Outliers Magnitude issues. It prevent from wrong decisions we are made.

Solve problem statement:

Training set,

$$\{(x^1, y^1), (x^2, y^2), (x^3, y^3), \dots, (x^n, y^n)\}$$

$(x^i, y^i) \rightarrow$  datapoints.

and  $y \in \{0, 1\} \rightarrow$  2/o/p

[Binary Classification]

Let's take,

$$(x_0) \text{ with } h_\theta(z) = \frac{1}{1 + e^{-z}} \quad \text{where } z = \theta_0 + \theta_1 x_1 \quad (\because \theta_0 = 0)$$

$$z = \theta_1 x_1$$

Change the parameter  $\theta_1$ ?

Yes, Because we need to fit the line.

Cost function,

Linear Regression

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

↳ Cost function  
of linear.

Logistic Regression equation,

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_1 x)}} \rightarrow \text{Equation}$$

We just apply logistic eq to linear Cost function,

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

where  $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_1 x)}}$

This is the cost function of logistic regression.

But we cannot use this cost function for logistic

regression. Because  $h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_1 x)}}$  this equation

is a non-convex function. } Non-Convex function.

Convex and Non-Convex functions are related to the

Gradient descent.

The Convex function mostly occurs in the linear regression,  $h_\theta(x) = \theta_0 + \theta_1 x$

$$h_\theta(x) = \theta_0 + \theta_1 x$$

So, the function looks like,

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

Convex  
function

Cost function.

Global Minima

In linear regression Cost function easily to reach the Global minima. So, it is called as Convex function.

But Logistic Regression,

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

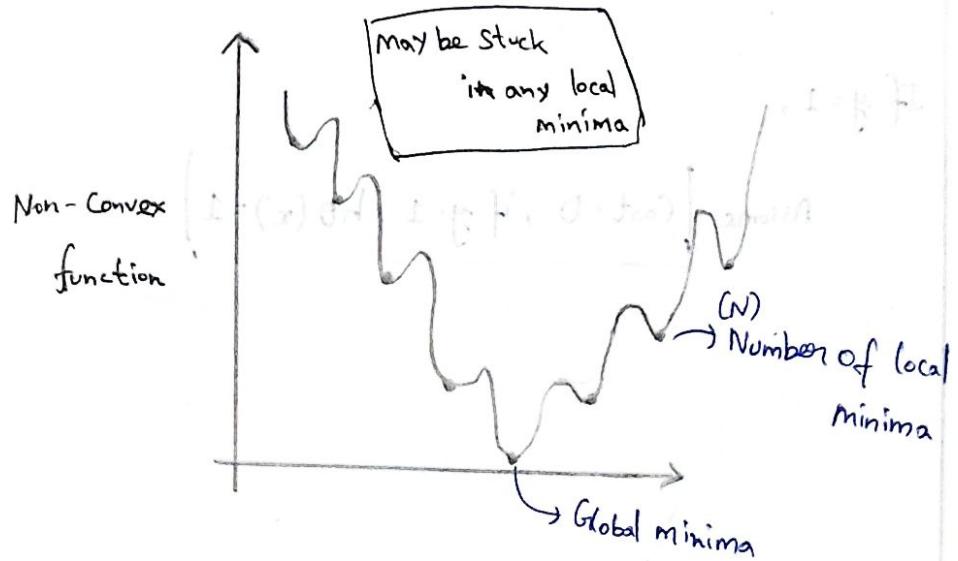
The non-convex function mostly occurs in the logistic (or) Sigmoid.

So, the function look like,

Applied Logistic Cost  
function (It is not  
correct one)

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

$$\text{where } h_\theta(x^{(i)}) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$



In this Case, The Global minima never reached. Because, So many number of local minima. So, we cannot get the Global Minima.

In order to avoid the above Scenario,

The actual Cost function of logistic regression is

$$J(\theta_0) = \begin{cases} -\log(h_{\theta}(x^{(i)})) & ; \text{if } y=1 \\ -\log(1-h_{\theta}(x^{(i)})) & ; \text{if } y=0 \end{cases}$$

$$\text{where } h_{\theta}(x) = \frac{1}{1+e^{-(\theta_0 x)}}$$

This is the required Cost function of Logistic

regression.

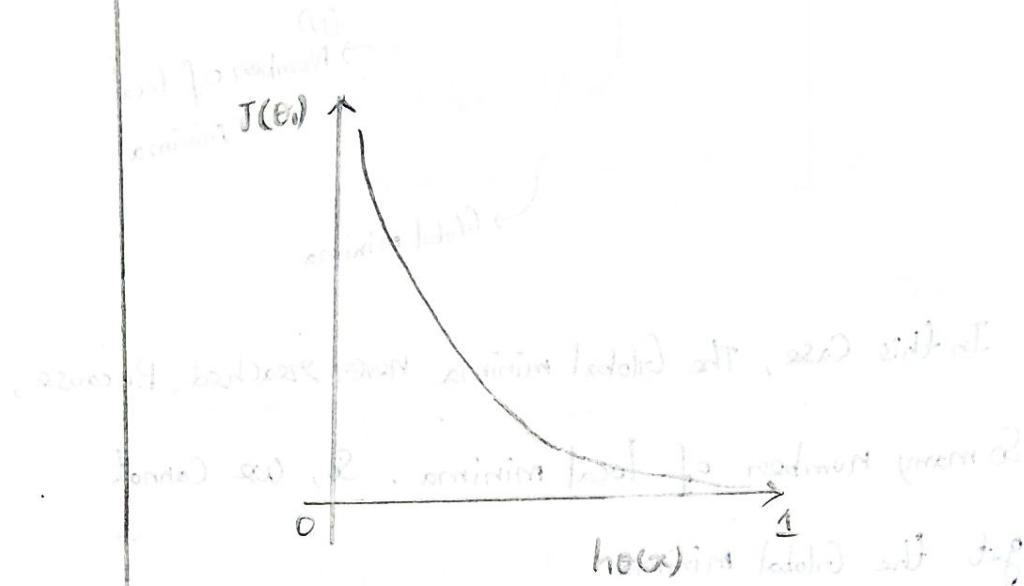
Training of the learning rate and model

The above Cost function basically mean,

If  $y=1$ ,

Assume

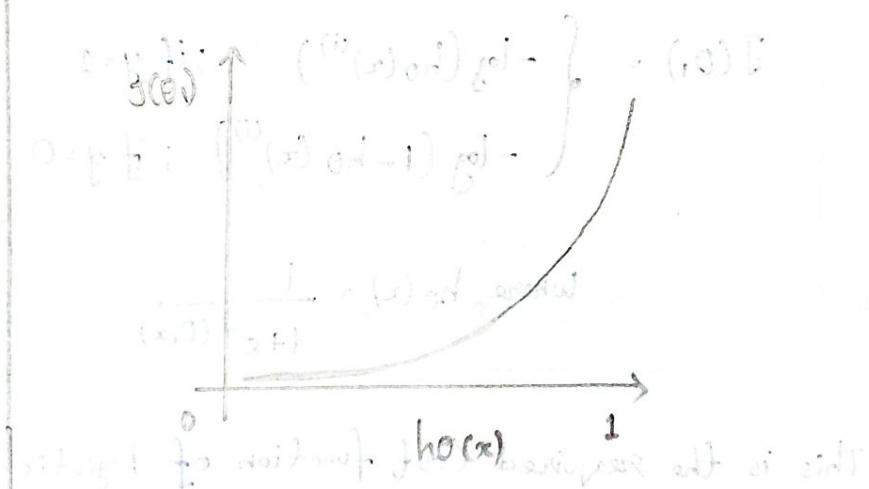
$$\text{Cost} = 0, \text{ if } y=1, h\theta(x)=1$$



If  $y=0$ ,

Assume

$$\text{Cost} = 0, \text{ if } y=0, h\theta(x)=0$$



When Combine above two graph, we get the gradient

descent. Then, find the Global Minima.

The final cost function,

$$\text{Cost}(h_{\theta}(x^{(i)}), y) = \begin{cases} -\log(h_{\theta}(x)) & ; \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & ; \text{if } y=0 \end{cases}$$

We combine the above cost function we get,

$$\boxed{\text{Cost}(h_{\theta}(x^{(i)}), y) = -y \log(h_{\theta}(x^{(i)})) - (1-y) \log(1-h_{\theta}(x))}$$

↳ Cost function

if we apply  $y=1$ ,

$$\text{cost}(h_{\theta}(x^{(i)}), y) = -\log(h_{\theta}(x^{(i)}))$$

if  $y=0$ ,

$$\text{cost}(h_{\theta}(x^{(i)}), y) = -\log(1-h_{\theta}(x^{(i)}))$$

↳ Cost  
function

So finally,

why  $y^i$ ?  
Because, It is binary classification  
there are two outputs.

$$J(\theta_1) = -\frac{1}{2m} \sum_{i=1}^m [y_i \log(h_{\theta}(x^{(i)})) + (1-y_i) \log(1-h_{\theta}(x^{(i)}))]$$

$$\text{where } h_{\theta}(x^{(i)}) = \frac{1}{1+e^{-\theta_1 x}}$$

↳ This is Required final  
Cost function for logistic  
regression

Repeat until Convergence

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

↳ This is the  
Convergence algorithm of  
Logistic regression

## Performance Metrics [Classification problem]

The Consider One Problem for binary Classification  
Let say,

$x_1, x_2$  are  $x$  and  $y$  [Actual O/P] and  $\hat{y}$  [Predicted O/P]

$(x_1, x_2) \text{ got } \hat{y} = 0$   $(x_1, x_2) \text{ got } \hat{y} = 1$

$(x_1, x_2) \text{ got } \hat{y} = 0$   $(x_1, x_2) \text{ got } \hat{y} = 1$

$(x_1, x_2) \text{ got } \hat{y} = 1$   $(x_1, x_2) \text{ got } \hat{y} = 0$

$(x_1, x_2) \text{ got } \hat{y} = 1$   $(x_1, x_2) \text{ got } \hat{y} = 0$

Based on the  $y$  and  $\hat{y}$ , we calculate the how good

the model is with the help of Confusion Matrix.

The first  $y$  and  $\hat{y}$  output are different. In this

case our first prediction was wrong.

The accuracy of the Model is Calculated with the help of Confusion Matrix.

Actual o/p ( $y$ )

	1	0
1	1	1
0	0	1

$y$        $\hat{y}$       Put

0	1	1
1	1	1
0	0	1

Actual o/p ( $y$ )

	1	0
1	2+1	2
0	1	1

$y$        $\hat{y}$       Put

1	1	Already $1+1=2$
1	1	Already $2+1=3$
0	1	Already $1+1=2$

The final Confusion matrix,

	1	0
1	3	2
0	1	1

The binary Classification Confusion matrix look like,

$\text{out} \leftarrow 0$

$\text{out} \leftarrow 1$

$\text{out} \leftarrow 0$

$\text{out} \leftarrow 1$

		Actual o/p	
		1	0
Predicted o/p	1	TP	FP
	0	FN	TN

→ This is called as

Confusion matrix.

TP and TN are the right because the actual input equal to the predicted input. But FN and FP are just vice-versa. So, The better accuracy we need to decrease FN and FP.

So, the accuracy formula will be,

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$\left[ \because \text{Acc} = \frac{\text{No. of Predict}^{\text{Corct}}}{\text{Total no of Pred}} \right]$

$$= \frac{3+1}{3+2+1+1} = \frac{4}{7}$$
$$= 0.57 \Rightarrow 57\%$$

We reduce FP and FN, we get better accuracy.

Suppose, the binary classification problem dataset look like

this,

In a dataset,

$$0 \rightarrow 900$$

$$1 \rightarrow 100$$

In this Scenario, We called it  
imbalanced dataset (or) biased  
dataset.

Suppose,

If  $0 \rightarrow 600$  and  $1 \rightarrow 400$ , In this case, We say 'Balanced  
dataset' (or) Normal dataset?

So, let's take,

$0 \rightarrow 900$ ,  $1 \rightarrow 100$  Biased dataset

The above case, The model was predicted mostly as 0. Because the model trained more 0 because of high count. But the Model accuracy is good and high.

$$\left\{ \text{model} \rightarrow 0 = \frac{900}{1000} = 90\% \right\}$$

↑  
mostly  
predicted

$\therefore \text{Accuracy} = \frac{\text{Correct}}{\text{Total no of data pred}}$

The accuracy is good and fine. But the prediction is mostly zero.

If the data is balanced, Then we go Accuracy performance metrics. If the data is Unbalanced, we go to the following one,

→ Recall.

→ Precision.

→ F-Score.

Recall, precision and F-Score will be based on

What type of the problem actually Given.

Precision: Out of all TP and FP, how many are predicted correctly. (TP).

$$\boxed{\text{Precision} = \frac{TP}{TP+FP}}$$

Recall: Out of all TP and FN, how many are predicted correctly. (TP)

$$\boxed{\text{Recall} = \frac{TP}{TP+FN}}$$

The Confusion matrix look like,

		Actual value	
		TP	FP
Value	0	FN	TN
	1	TP	FP

Real time example of Confusion Matrix, [Recall]

Classification evaluation metrics.

Based on primary building blocks - TP, TN, FP, FN.

Our Classification problem about whether the person is Covid or not.

Covid - POSITIVE CLASS

Not covid - NEGATIVE CLASS.

[Unauthorized] will (not) help if (not) -

TP <u>Actual</u> = Covid <u>Model_pred</u> = Covid <u>Outcome</u> : Patient Admitted.	FP <u>Actual</u> : Not covid. <u>Model</u> : Covid. <u>Outcome</u> : Patient Admitted
FN <u>Actual</u> : Covid <u>Model</u> : Not Covid <u>Outcome</u> : Patient is not admitted risk of infecting Others.	TN <u>Actual</u> : Not covid. <u>Model</u> : Not Covid <u>Outcome</u> : Patient is not admitted.

Out of above four Cases, FN has need to reduce.

Because, the model predict 'Not covid' but actual is Covid. It may be infecting Others. So, we need to focus of FN.

So, we use the Performance Metrics

Called as Recall.

Another Example,

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

↓  
Recall

Consider another one: [Precision]

Our classification problem is about the whether the mail is Spam (or) Ham. [Spam classification]

		Actual	
		Spam	Ham
Predict	Spam	TP	FP
	Ham	FN	TN

Let's breakdown,

[TP] Actual  $\rightarrow$  spam but predict  $\rightarrow$  spam [ok] ✓

[FP] Actual  $\rightarrow$  Ham but Predict  $\rightarrow$  spam [P] (x)

↳ It causes danger because the spam

Message also consider in the inbox.

[maybe missing some original msg]

[FN] Actual  $\rightarrow$  spam but predict  $\rightarrow$  Ham

↳ If also in the inbox and

mixing in original message .

[TN] Actual  $\rightarrow$  Ham but predict  $\rightarrow$  Ham [ok] ✓

The FP and FN has more danger. Compare both  
FP is more important.

Suppose message is Spotted Ham, but not consider  
We lose something. [FP]

And other side machine does not make 100%  
correct prediction. [Actual  $\rightarrow$  Spam, pred  $\rightarrow$  Ham]. we also  
need some Safetyness. [FN]

So, therefore in this case, we use Precision  
formula to evaluate.

$$\boxed{\text{Precision} = \frac{TP}{TP+FP}}$$

Last Example: [F - Beta]

Our Classification Problem is about Tomorrow Share

Market going to Crash.

Confusion matrix, Actual

		Crash	Not Crash
Predict	Crash	TP	FP
	Not Crash	FN	TN

The above Confusion Matrix,

TP ( $\rightarrow$ ) Says the bank Crash

FP  $\rightarrow$  Actual the market  $\rightarrow$  not Crash

Predict  $\rightarrow$  Crash

So, in this Case the People are sells the

Stock because of Crash and the people does not

buy any stocks.

FN  $\rightarrow$  Actual the market  $\rightarrow$  Crash

Predict  $\rightarrow$  Not Crash.

In this case, the market going to be Crash

but prediction is not Crash. So, the market Owners

are not ready to recover.

In this both Cases, we need to focus both

FP and FN. Then we are go to the F-Beta.

$$F\text{-Beta} = \frac{(1+\beta^2) \cdot \text{Precision} \times \text{Recall}}{\beta^2 \times \text{precision} + \text{Recall}}$$

Let's Consider,

$$\beta = 1 \quad = \frac{(1+1) \text{ Precision} \times \text{Recall}}{(1) \text{Precision} + \text{Recall}}$$
$$= \frac{2 (\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}} \quad \begin{array}{l} \text{Harmonic mean} \\ \text{Precision} + \text{Recall} \end{array}$$

It is something called

$$\text{Harmonic mean} = \frac{2xy}{x+y}$$

Suppose,  $\beta = 0.5$ ,

$$= (1+(0.5)^2) \frac{\text{P} \times \text{R}}{(0.25) \text{P} + \text{R}}$$

Same  $\beta = 2$ , . . . .

Suppose Use  $\beta = 2$  is called  $F_2$ -Score.

Suppose Use  $\beta = 0.5$  is called  $F_{0.5}$ -Score.

Suppose Use  $\beta = 1$  is called  $F_1$ -Score.

What is and why Beta Value Used?

Beta is deciding parameter of F-Score.

If  $\beta = 1$ , The Same importance for both FP and FN.

If  $\beta < 1$  like  $\beta = 0.5$ , we give more important to FP than FN.

If  $\beta > 1$  like  $\beta = 2$ , we give more important to FN than FP.