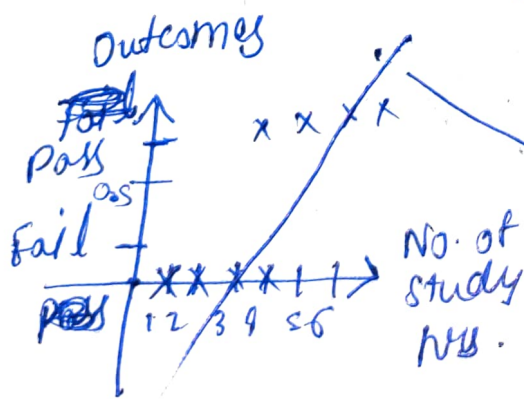


It is use to solve

Logistic Regression (First algorithm for classification) (Binary classification)

No. of study	No. of play	P/I =
—	—	P
—	—	
—	—	F



We can't solve this problem with linear regression

If $h_0(n) < 0.5 \Rightarrow 0 \rightarrow \text{Fail}$
 $h_0(n) \geq 0.5 \Rightarrow 1 \rightarrow \text{Pass}$

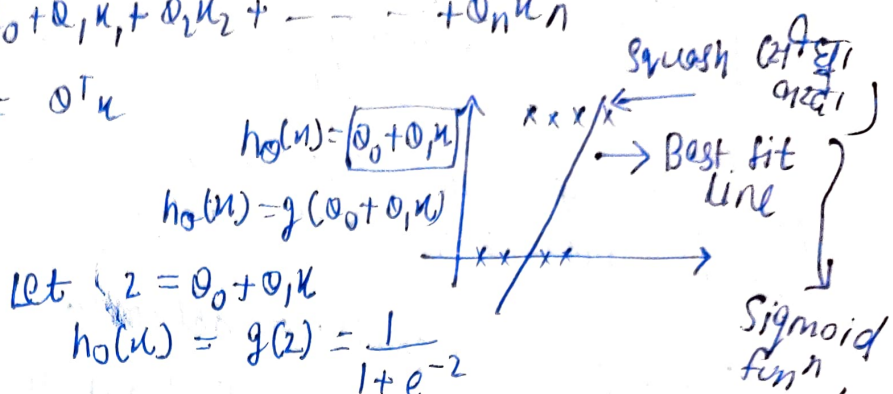
- Why we are not using linear regression?
- When we got outliers, it change the entire line & give wrogy result.
 - Greater than 1 & less than 1 then we will get negative no. which will also be problem.

But we are result = 0 or 1.
 So sigmoid function is the solution for it.

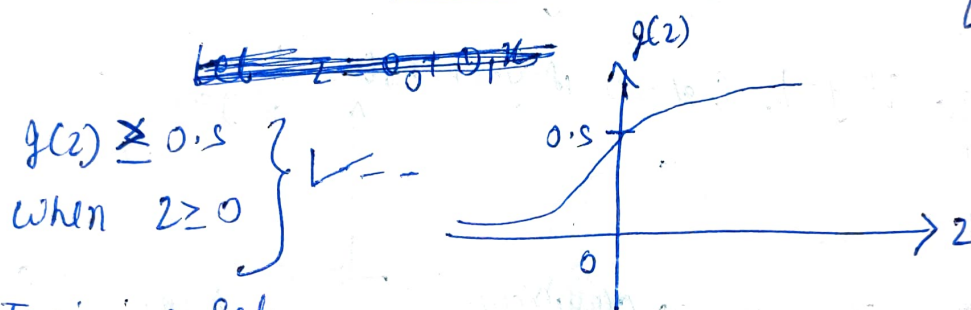
Decision Boundary Logistic Regression

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$h_0(x) = \theta^T x$$



$$h_0(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}} \rightarrow \text{Sigmoid or Logistic fun}^n$$



Training Set

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

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

Change parameter θ_1 ?

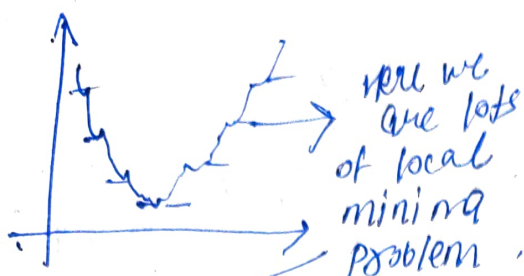
Cost function

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

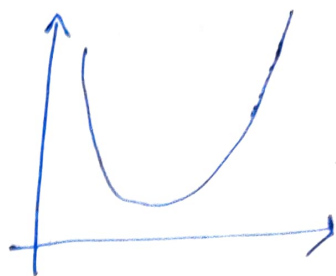
Logistic regression $\rightarrow h_0(x) = \frac{1}{1 + e^{-(\theta_1 x)}}$

Logistic regression cost function $= \frac{1}{2} (h_0(x^{(i)}) - y^{(i)})^2$
 So, $h_0(x) = \frac{1}{1 + e^{-(\theta_1 x)}}$

Non-convex function



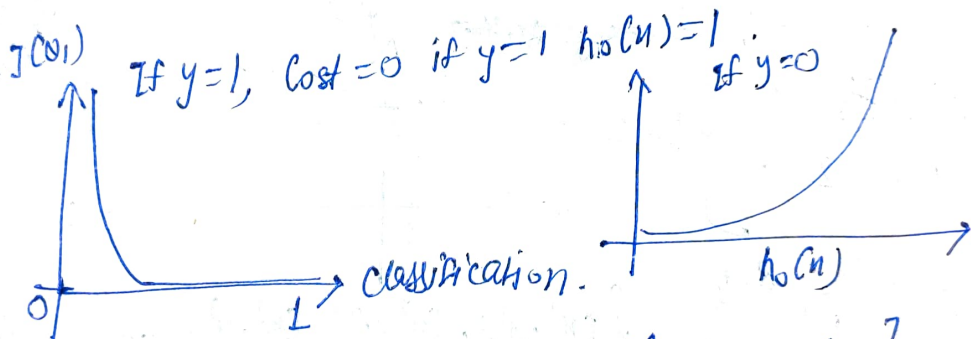
Convex function



In order to solve this we have something called logistic regression cost function -

So, logistic regression actually changes,

$$J(w) = \begin{cases} -\log(h_0(w)) & y=1 \\ -\log(1-h_0(w)) & y=0 \end{cases}$$



$$\text{So, Cost}(h_0(x^i), y) = \begin{cases} -\log(h_0(x^i)) & \text{if } y=1 \\ -\log(1-h_0(x^i)) & \text{if } y=0 \end{cases}$$

$$\boxed{\text{Cost}(h_0(x^i), y) = -y \log(h_0(x^i)) - (1-y) \log(1-h_0(x^i))}$$

Cost for y

If $y=1$

$$\text{Cost}(h_0(x^i), y) = -\log(h_0(x^i))$$

If $y=0$

$$\text{Cost}(h_0(x^i), y) = -\log(1-h_0(x^i))$$

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

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

Repeat until convergence.

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

}

→ We are changing cost funⁿ for non convex function.

Performance metrics { Classification problem }

		Actual - o/p		Pred		Actual	
x_1	x_2	y	\hat{y}			1	0
—	—	0	1	1	3	2	
—	—	0	0	0	1	1	
—	—	1	1				
—	—	0	1				
—	—	1	0				
—	—	0	0				

No. of times

	1	0
1	TP	FP
0	FN	TN

→ Confusion matrix.

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

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

Now, $\begin{matrix} 0 \rightarrow 900 \\ 1 \rightarrow 100 \end{matrix} \left. \vphantom{\begin{matrix} 0 \rightarrow 900 \\ 1 \rightarrow 100 \end{matrix}} \right\} \text{ Imbalanced data}$ $\begin{matrix} 0 \rightarrow 600 \\ 1 \rightarrow 400 \end{matrix} \left. \vphantom{\begin{matrix} 0 \rightarrow 600 \\ 1 \rightarrow 400 \end{matrix}} \right\} \text{ Balanced data}$

1. Precision:

2.1

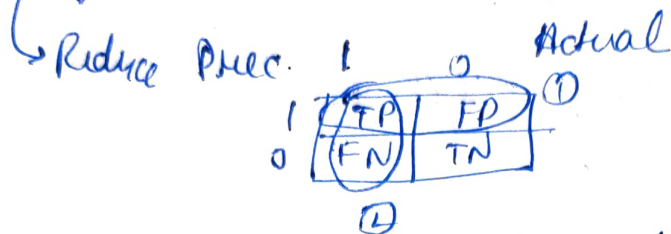
Precision

Recall

F-score

$$\left\{ \frac{TP}{TP + FP} \right\}$$

$$\left\{ \frac{TP}{TP + FN} \right\} \rightarrow \text{Reduce } FN$$



Spam Classification — Precision.

Has Cancer or not — Recall

Tom stock market is going to crash — For people
For company Different for both.

But for both we use F-score.

$$(1 + \beta^2) \frac{\text{Prec} \times \text{Recall}}{\beta^2 \times \text{Prec} + \text{Recall}}$$

$\beta = 1$, $\textcircled{F_1}$, $\frac{2 (\text{Prec} \times \text{Recall})}{\text{Prec} + \text{Recall}} \rightarrow \text{Harmonic mean.}$

$\beta = 0.5$, $\textcircled{F_{0.5}}$, $(1 + 0.5^2) \frac{P \times R}{0.5^2 P + R}$

$\beta = 2$, $\textcircled{F_2}$, $1 + 4 \gg 1 + \beta$

Agenda

- ① Practicals.
- ② Naive Bayes.
- ③ KNN algorithms.

Cross Validation — for doing testing & training of all possible combinations

→ L_1 & L_2 penalty is also present in logistic regression.