

logistic Regression

* for classification problem

Example

* Problem to predict Pass or Fail based on study hours & play hours

IIT, JEE

| Study/hor | Play hours | o/p (Pass/Fail) |
|-----------|------------|-----------------|
| 1 | 8 | Fail |
| 2 | 7 | Fail |
| 3 | 7 | Fail |
| 6 | 3 | Pass |
| 2 | 4 | Pass |

→ out here

* we are predicted Pass or Fail

* we can make model to predict

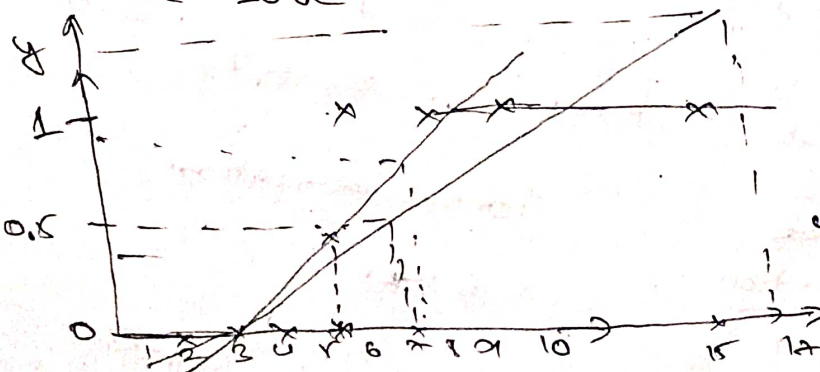
logistic regression \Rightarrow Buy or Not

Dataset (UPSC)

| Study hours | o/p (Pass/Fail) |
|-------------|-----------------|
| 2 | Fail |
| 3 | Fail |
| 4 | Fail |
| 6 | Pass |
| 7 | Pass |
| 8 | Pass |

1 \Rightarrow Pass 0 \Rightarrow Fail

① Can we solve the problem using regression?



* Threshold is 0.5 \Rightarrow 0 and 1

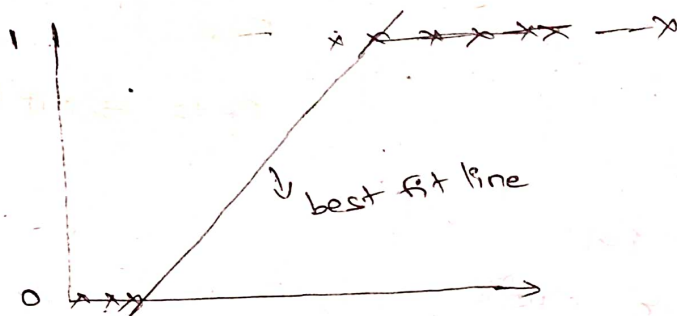
* In regression

if outliers are there then the complete best fit line will be change if we see the 7 hours study student will fail it will give wrong

* we can't remove outliers (Important)

* we squash the line in between 0 and 1

* find best fit then squash by sigmoid Activation



using
Sigmoid Activation

\Downarrow

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

\Downarrow

o/p \Rightarrow 0 to 1

steps logistic

① $z = h(\theta, x) = \theta_0 + \theta_1 x$ \rightarrow fit line $z = \theta_0 + \theta_1 x$

② sigmoid fn $= \frac{1}{1 + e^{-z}}$ \Rightarrow get in between 0 and 1

① create a best fit line

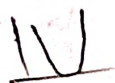
② squashing \rightarrow sigmoid function

Linear regression Cost function

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

\Downarrow MSE
 $h(\theta, x) = \theta_0 + \theta_1 x$
 \Downarrow

Global minimum \Rightarrow Convex function



Logistic Regre Cost Function

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

$z = \theta_0 + \theta_1 x$

$$h(\theta, x) = \sigma(\theta_0 + \theta_1 x)$$

\Downarrow
sigmoid activation

$$= \sigma(z)$$

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

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

\Downarrow
0 to 1

- * Big issue with equation
- * when we apply in logistic it give non convex function
- * more
- * Cost function not convex or non convex function

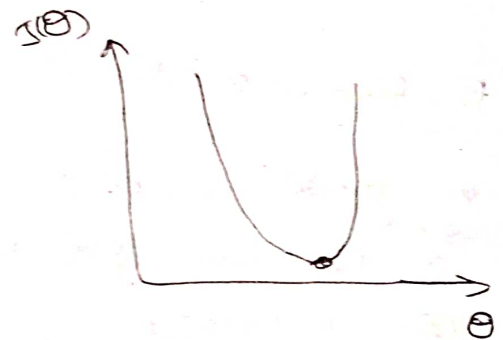
Non Convex function



- * lots of local minima, coefficient stuck
- * global minima

* to fix non convex function we use different fun cost
like log loss cost function

Convex function



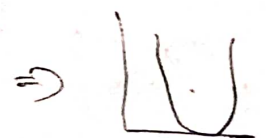
- * linear
- * 1 global minima

* log loss cost function

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

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

$\Downarrow \Downarrow \Downarrow$



gk

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

* through the cost function we will not get local minima * we create convex function
 # increase value reduced by 40000
 minimize cost function $J(\theta_0, \theta_1)$ by changing θ_0, θ_1
 by Convergence Algorithm

Repeat Convergence

$J = 0$ and

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

?

Performance Metrics:

* model is perform well or not using

- ① Confusion matrix
- ② Accuracy
- ② Precision
- ④ Recall
- ⑤ F-Beta score

③ Confusion matrix %
 binary classification

| | | | | |
|----|-----------|---|---|----------|
| | | 1 | 0 | y actual |
| 1 | 3 | 2 | | |
| 0 | 1 | 1 | | |
| 87 | Predicted | | | |

| | | | |
|---|----|----|---|
| | | 1 | 0 |
| 1 | TP | FP | |
| 0 | FN | TN | |

1 \Rightarrow Posit
 0 \Rightarrow Negat

| | | | | | |
|--|--|----------------|----------------|---|------------------|
| | | Dataset | | | |
| | | f ₁ | f ₂ | y | model prediction |
| | | | | 0 | 1 |
| | | | | 1 | 1 |
| | | | | 0 | 0 |
| | | | | 1 | 1 |
| | | | | 0 | 1 |
| | | | | 1 | 0 |

(Accuracy)

$$\textcircled{2} \quad \text{Acc} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$\Rightarrow \frac{3+1}{2+2+1+1} = \frac{4}{4} \approx 50\%$$

* Dataset \rightarrow binary classification

\hookrightarrow 1000 data points $\begin{cases} 900 \rightarrow 1 \\ 100 \rightarrow 0 \end{cases}$ } imbalance data set

Dumb model $\rightarrow 1 \Rightarrow 90\%$ Accuracy \rightarrow X sufficient not

* solve this problems

$\textcircled{3}$ * Precision %

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

| Actual | | | | | |
|---------|---|----|----|----|----|
| 1 | 0 | | | | |
| Predict | <table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table> | TP | FP | FN | TN |
| TP | FP | | | | |
| FN | TN | | | | |

\Rightarrow out of all the actual values how many are correctly predicted.

* main aim reduce False Positive

Problem statement

FP and FN error

mail \rightarrow Spam or Ham

* focus on FP and FN reduce

Problem statement

Diabetes or not Diabetes.

* focus on FN $\downarrow \downarrow$

$\textcircled{4}$ Recall

Recall

\div

$$\frac{TP}{TP + FN}$$

\Rightarrow out of all the predicted values how many are correctly predicted

Tomorrow stock market going to crash

→ Consumer → FN ↓ ↓
 → Companies → FP ↓ ↓

| | | |
|---|----|----|
| | 1 | 0 |
| 1 | TP | FP |
| 0 | FN | TN |

use

* F-Beta Score %

$$\frac{(1+\beta)^2 \cdot \text{Precision} \cdot \text{Recall}}{(\beta^2 \cdot \text{Precision} + \text{Recall})}$$

① If FP and FN are both important

$$\beta = 1$$

$$F1 \text{ score} = 2 \frac{P \cdot R}{P + R}$$

② If FP is more important than FN

$$\beta = 0.5$$

$$F_{-0.5} \text{ score} = \frac{(1+0.25) P \cdot R}{(0.25 \cdot P + R)}$$

③ If FN >> FP

$$F2 \text{ score} = \frac{(1+4) P \cdot R}{(4 \cdot P + R)}$$

* If imbalanced data set is there then we focus

FP ⇒ precision ↓ ↓

FN ⇒ recall ↓ ↓