# Logistic Regression (classification problem)

## IIT JEE Example

| Study | Play hrs | o/p (Pass/Fail) |
|-------|----------|-----------------|
| 1 | 8 | Fail |
| 2 | 7 | Fail |
| 3 | 7 | Fail |
| 6 | 3 | Pass |
| 7 | 2 | Pass |
| 6 | 4 | Pass |
| 5 | 3 | Pass |

☐ → outliers.

## Data Set

| Study hours | o/p (Pass/Fail) |
|-------------|-----------------|
| 2 | Fail |
| 3 | Fail |
| 4 | Fail |
| 5 | Pass |
| 6 | Pass |
| 7 | Pass |
| 8 | Pass |
| 9 | Pass |
|   | Fail |

UPSC

$1 \Rightarrow$ Pass
$0 \Rightarrow$ Fail

① Can we slove this program using regression.

# Regression

$O.5 \Rightarrow$ <mark>Threshold</mark>

$$y \leq 0.5 = 0$$

$$y > 0.5 = 1$$

<mark>Regression</mark>

$y \leq 0.5 = 0$

$y > 0.5 = 1$

<mark>$>1 \& <0$</mark>



$\rightarrow$ best fit line

$O$   2   3   4   5   6   7   8   9   10   12   15   18

7 $\rightarrow$ fail

---

<mark>Sigmoid Activation</mark>



$\rightarrow$ best fit line

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

$\Downarrow$

<mark>Sigmoid Activation</mark>

$\Rightarrow o/p = 0$ to

---

① $Z = h_\theta(x) = \theta_0 + \theta_1 x$

sigmoid fn

$$\frac{1}{1+e^{-z}} \Rightarrow 0 \text{ to } 1$$

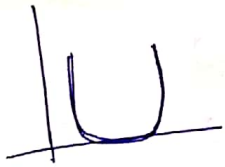$Z = \theta_0 + \theta_1 x$

## Linear Regression Cost function

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

MSE

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

Convex function

1 global minima



① Create a best fix line

② Squashing → Sigmoid function

$$\sigma = \frac{1}{1+e^{-z}}$$

$$Z = \theta_0 + \theta_1 x$$

## Logistic Regression Cost function

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

→ best fit line.

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

↓

Sigmoid Activation

$$= \sigma(z)$$

$$= \sigma(z)$$

$$= \frac{1}{1+\bar{e}^z}$$

$$z = \theta_0 + \theta_1 x$$

$$h_\theta(x) = \frac{1}{1+\bar{e}^z}$$

$$\boxed{h_\theta(x) = \frac{1}{1+\bar{e}^{(\theta_0 + \theta_1 x)}}} \Rightarrow 0 \text{ to } 1$$

$\Rightarrow$ out will be
b/w (0 to 1)

$$\leq 0.5 \Rightarrow 0 \Rightarrow \text{fail}$$
$$> 0.5 \Rightarrow 1 \Rightarrow \text{Pass}$$

$0 \Rightarrow$ fail

$1 \Rightarrow$ Pass

$$\boxed{\text{Threshold} - 0.5}$$

let
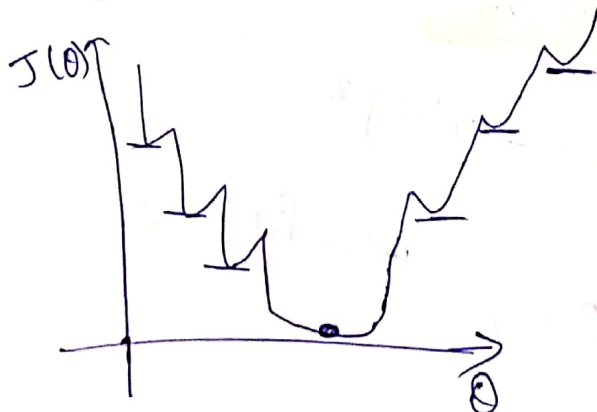$\Rightarrow 0.35 \Rightarrow 0$

$0.25 \Rightarrow 0$

$0.95 \Rightarrow 1$

$1 \Rightarrow 0.7$

$0.5 \Rightarrow 0$

$0.54 \Rightarrow 1$

Non - convex function



Convex function

+ log loss Cost function

$$\text{Cost}\left(h_\theta(x_0)^{(i)}, y^{(i)}\right) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1-h_\theta(x)) & \text{if } y = 0 \end{cases}$$

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

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

↳ Convex function

↳ Never get local minima

minimize Cost function $J(\theta_0, \theta_1)$ by changing $\theta_0, \theta_1$

Converging algorithm

Repeat converging

$$j = 0 \text{ and } 1$$

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

Threshold = 0.5

# Performance metrics

① Confusion matrix

② Accuracy

③ Precision

④ Recall

⑤ F- Beta Score

**Data set**

| $f_1$ | $f_2$ | o/p |
|-------|-------|-----|
| — | — | 0 |
| — | — | 1 |
| — | — | 0 |
| — | — | 1 |
| — | — | 0 |
| — | — | 1 |
| — | — | 0 |

$\hat{y} \Rightarrow y_{cal}$

→ model out

## Confusion matrix:

| | 1 | 0 ⟹ y |
|---|---|---|
| 1 | 2 | 3 |
| 0 | 1 | 1 |

↓ Predicted

↳ Actual value

↳ Confusion matrix

⟹ Actual

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

⟹ Predict

TP ⟹ True prediction (Positive)

TN ⟹ True negative

FP ⟹ False prediction (Positive)

FN ⟹ False Negative

Accuracy =

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

$$= \frac{2+1}{2+3+1+1} = \frac{1}{7}$$

**Dataset:** Binary Classification

→ 1000 data points
- → 900 → 1
- → 100 → 0

} Imbalanced Dataset

Dum model → 1 ⇒ 90% Accuracy ⇒ → x sufficient

⊕ **Precision** =

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

Actual

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

‖ Predicted

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

**Problem statement:**

**Mail → Spam or Ham**

i am getting a mail which is spam (1) and model predicted as spam(1). then it is TP (true Positive)

while ⇒ pay i am getting an important (1) and model is predicting as spam (1) then it is FP ⇒ (false Positive)

So to reduce the false Possitive

So we have to focus on false Possitive//
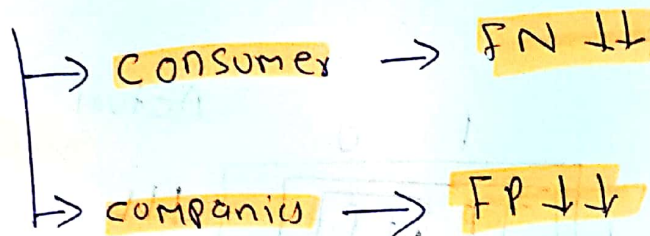
**Model → diabetes or not Diabetes**

Recall = $\boxed{\dfrac{TP}{TP+FN}}$ ⇒ out of all the predicted values how many are correctly predicted.

Ex:-

Tomorrow the stock market is going to crash

⇩

→ consumer → FN ↓↓

→ companies → FP ↓↓

⇒ Actual

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

⇩
predicted

then we use

⊕ F - Beta Score

$$\frac{(1+\beta^2)\; Precison * Recall}{(\beta^2)\; Precison + Recall}$$

① If FP and FN are both important

$$\beta = 1$$

$$F1\ score = \frac{2\; P*R}{P+R}$$

② if FP is more important than FN

$$\beta = 0.5$$

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

③ If FN >> FP

$$F_2 \text{ score} = \frac{(1+4)\ P*R}{(4*P+R)}$$