# Logistic Regression (Binary classification)

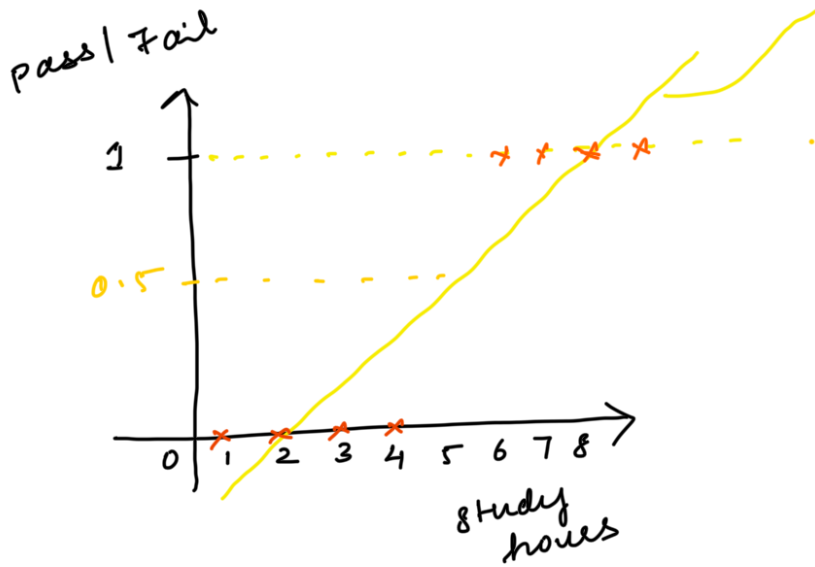## Dataset

| Input Feature | Output Feature | (Binary classification) |
|---|---|---|
| Study hours | Pass / Fail | |
| 2 | Fail | |
| 3 | Fail | |
| 4 | Fail | |
| 5 | Pass | |
| 6 | Pass | |
| 7 | Pass | |

Input Feature → study hours → ML Model → Output Feature → Pass / Fail

Why we cannot use Linear regression for classification?

① outlier → Best fit line changes
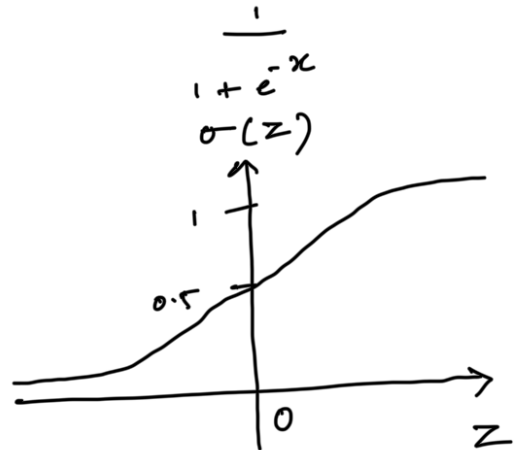
② >1 and <0 → squashing the line is not possible.

How Logistic Regression slower classification

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

pass / fail



study hours

Sigmoid Activation

$$\frac{1}{1 + e^{-x}}$$

$\sigma(z)$



In Logistic regression also we will create a best fit line, and top of the best fit line we will apply the activation function

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

$Z \geqslant 0$

$\sigma(z) \geqslant 0.5$

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

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

$\rightarrow Z = \theta_0 + \theta_1 x_1$

$\rightarrow$ Logistic regression hypothesis

Logistic regression cost function

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

$$J(\theta \cdots) \qquad 2m \qquad \overline{\sum_{i=1}}$$

$$h_\theta(x) = \cfrac{1}{1 + \bar{e}^2} \qquad \Bigg| \qquad Z = \theta_0 + \theta_1 x_1$$

This will not provide Gradient descent

This function provides non convex function



$J(\theta)$ — Local minima

non convex function

Global minima

In the above example, my $\theta$ will 8 tangent at one position.

In order to prevent this what we can do is that to bring the convex function

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

we know that

$$h_\theta(x)^i = \cfrac{1}{1 + e^{-Z}} \qquad \Bigg| \qquad Z = \theta_0 + \theta_1 x$$

Let's denote $\left( h_\theta(x)^i - y^i \right)^2$ as

$$cost\left( h_\theta(x)^i, \; y^i \right)$$

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

$\Downarrow$

we basically call this as

Log Loss

this Log Loss help us to come up with a convex function. — so that we can acheive a global minima.

therefore,

$$cost\left(h_\theta(x)^i, y^i\right) = -y\log\left(h_\theta(x)\right) - (1-y)\log\left(1-h_\theta(x)\right)$$

$\Downarrow$            $\Downarrow$

$y=1$            $y=0$

therefore, the cost function can be denoted as.

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

⇓

this cost function will give me the concave function, Probably we can call it as global minima.

Minimize cost function $J(\theta_0, \theta_1)$ by changing $\theta_0$ and $\theta_1$

convergence Algorithm

Repeat

{

$$\theta_j \approx \theta_j - \alpha \frac{d}{d\theta_j} J(\theta_0, \theta_1)$$
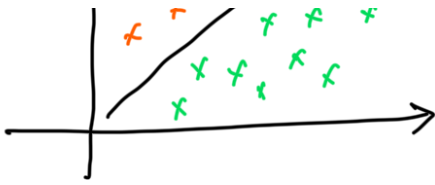
| $j = 0$ and $1$

}

Performance Metrics, Accuracy, Precision, Recall and F - Beta Score



Data set

$\hat{\Omega}$

| $x_1$ | $x_2$ | $y$ | $\hat{y}$ |
|-------|-------|-----|-----------|
| — | — | 0 | 1 |
| — | — | 1 | 1 |
| — | — | 0 | 0 |
| — | — | 1 | 1 |
| — | — | 1 | 1 |
| — | — | 0 | 1 |
| — | — | 1 | 0 |

$x_1$ & $x_2$ → I/P
     feature

$y$ → Output
     feature

$\hat{y}$ → predicted
     feature.

derived from
above data set.

① confusion Matrix

Basically its 2×2

Actual values

|  | 1 | 0 |
|---|---|---|
| 1 | ③ | ②|
| 0 | ① | ① |

Predicted values

Actual values

• wrong prediction
• correct prediction

Actual values

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

Predicted values.

here, TP and TN

TP - True positive,
when actual is 1 and
predicted is 1
FP - False positive,
when actual is 0 and
Predicted is 1
FN - False negative,
when actual is 1, and
predicted is 0
TN - True negative, when

are the correct results, and FP and FN are the wrong results.
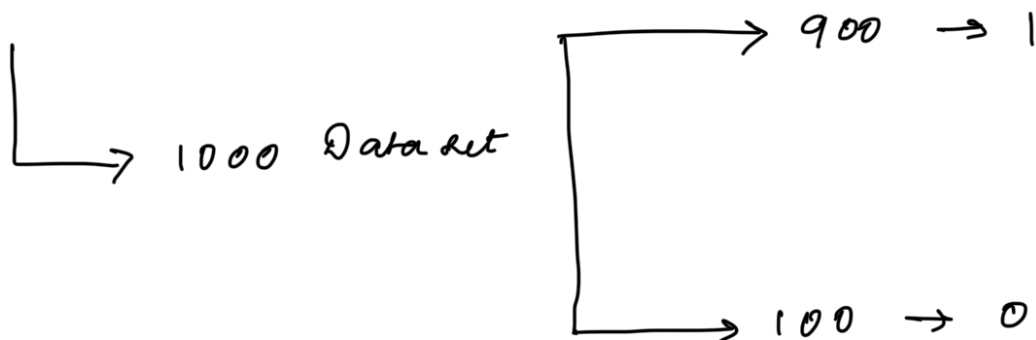
actual is 0, and predicted is 0

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

Lets compute the accuracy for our dataset

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

$$= \frac{4}{7}$$

$$= 0.5714$$

Data set

Binary classification

```
      ┌─────────→ 900 → 1
──┐   │
  └──→ 1000 Data set │
      │
      └─────────→ 100 → 0
```

The above classification is imbalanced data set, and we cannot

directly use accuracy.

⇓

In order, to prevent this we can use precision and recall.

## Precision

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

⇓

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

## Recall

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

⇓

out of all the predicted value, how many are correctly predicted.

1. Then should we use precision and

when should we use recall. ?

## Use case ①

Spam classification

$$Mail \rightarrow spam$$
$$Model \rightarrow spam \quad \Big\} \quad \text{Good}$$

$$Mail \rightarrow Not \; spam$$
$$Model \rightarrow spam \quad \Big\} \quad Blender$$

|  | 1 | 0 | Actual |
|---|---|---|---|
| Predicted 1 | TP | FP → should be reduced |  |
| 0 | FN | TN |  |

1 - Spam

0 - Not spam

For this case we can use

Precision

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

## use case 2

To predict whether person has

diabities or not

Truth $\longrightarrow$ Diabities
Model $\longrightarrow$ Dees'nt Diabities $\Bigg\}$ $\rightarrow$ Blender.

Truth $\longrightarrow$ Diabities
Model $\longrightarrow$ Diabities $\Bigg\}$ hood

Truth $\longrightarrow$ Not diabities
model $\longrightarrow$ Diabities $\Bigg\}$ ee train

Actual
Diabities    No diabities

|  | Actual Diabities | No diabities |
|---|---|---|
| Predicted diabities | TP | FP |
| No diabities | FN | TN |

should be reduced in this use case.

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

F - Beta score

Precission $\times$ Recall

$$F-\text{Beta Score} = (1 + \beta^2) \frac{}{\text{Precission} + \text{Recall}}$$

① If FP and FN are both impartant

$$\beta = 1$$

therefore,

$$F1 \text{ Score} = (1 + 1^2) \frac{\text{Precission} \times \text{Recall}}{\text{Precission} + \text{Recall}}$$

This is reffered as Harmonic Mean.

② If FP is more impartant than FN

$$\beta = 0.5$$

therefore,

$$F0.5 \text{ Score} = (1 + 0.5^2) \frac{\text{Precission} \times \text{Recall}}{\text{Precission} + \text{Recall}}$$
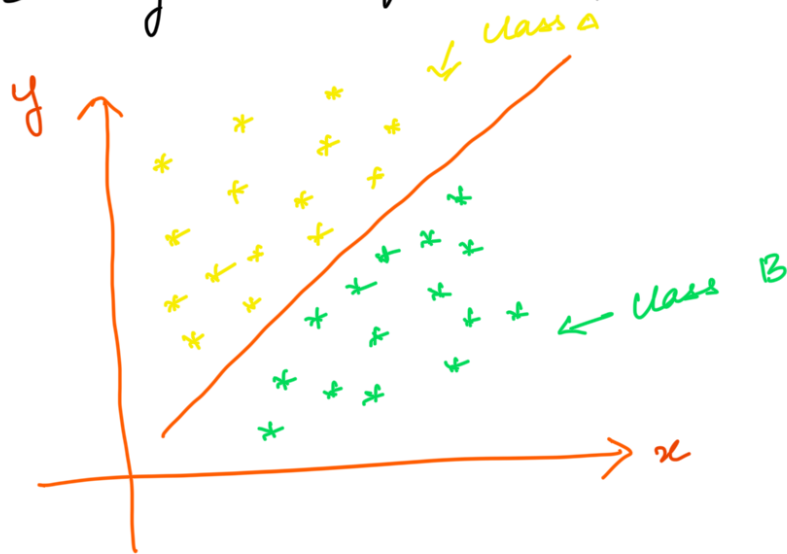
③ If FN is more impartant than FP

$$\beta = 2$$

therefore,

$$F2 \text{ Score} = (1 + 2^2) \frac{\text{Precission} \times \text{Recall}}{\text{Precission} + \text{Recall}}$$
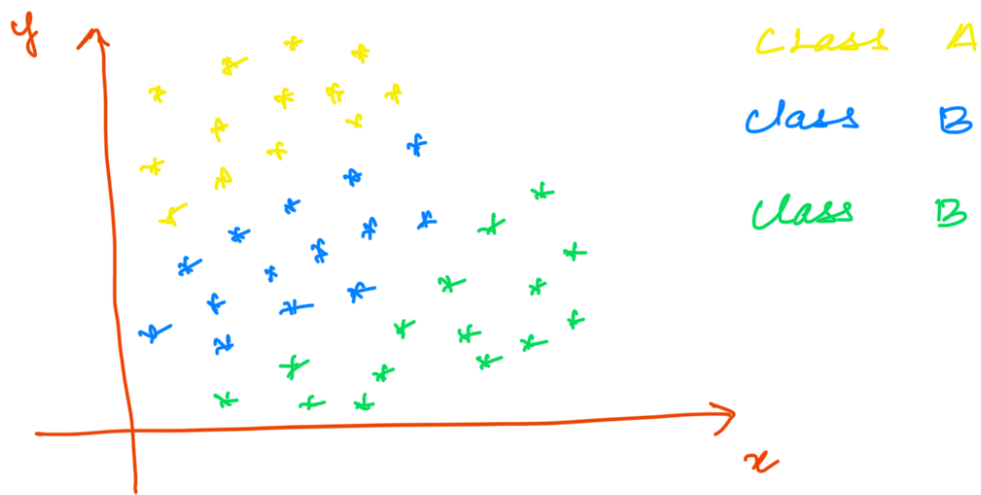
# Logistic Regression (one versus rest)

In Binary classification,



we can deal with only two output classification, what if i have more than two output classification.



Class A

Class B

Class B

The above representation is the multiclass classification problem

Multiclass classification problem can be solved with the help of one versus
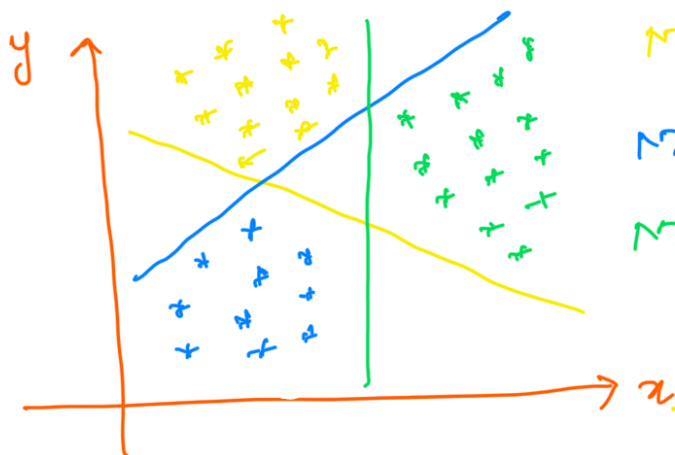
Rest (OVR)

Dataset      we have 3
_____     classification $O_1 O_2 O_3$
                    ↓

| $f_1$ | $f_2$ | $f_3$ | o/p |
|-------|-------|-------|-----|
| -     | -     | -     | $O_1$ |
| -     | -     | -     | $O_2$ |
| -     | -     | -     | $O_3$ |
| -     | -     | -     | $O_1$ |
| -     | -     | -     | $O_3$ |
| -     | -     | -     | $O_2$ |

with the help
of one versus
Rest what we do is
that, we try to
create internally
multiple model, and
each model act as
a binary classification

here we have 3 classes, hence OVR
create 3 internal model



$M_1 \rightarrow$ Binary classification

$M_2 \rightarrow$ Binary classification

$M_3 \rightarrow$ Binary classification

Basically,

$M_1 \rightarrow$ class A + class B and class c

$M_2 \rightarrow$ Class B + Class A and Class C

$M_3 \rightarrow$ Class C + Class A and Class C

$\Downarrow$

Each model, $M_1$, $M_2$, $M_3$ will perform binary classification

We form the results using one hot encoding.

| O/P | $O_1$ | $O_2$ | $O_3$ |
|-----|-------|-------|-------|
| $O_1$ | 1 | 0 | 0 |
| $O_2$ | 0 | 1 | 0 |
| $O_3$ | 0 | 0 | 1 |
| $O_1$ | 1 | 0 . | 0 |
| $O_3$ | 0 | 0 | 1 |
| $O_2$ | 0 | 1 | 0 |

## Model

$M_1 \rightarrow$ Input $\rightarrow O_1$

$M_2 \rightarrow$ Input $\rightarrow O_2$

$M_3 \rightarrow$ Input $\rightarrow O_3$

New Test Data

$M_1 \rightarrow 0.25 \quad O_1$

$M_2 \rightarrow 0.20 \quad O_2$

$M_3 \rightarrow 0.55 \quad O_3$