

Logistic Regression is a statistical and machine learning method used for classification. Instead of predicting a continuous value (like in linear regression), it predicts the probability that an observation belongs to a particular class (e.g., Yes/No, 0/1, Pass/Fail).

#### ☐ Key Idea:

We want to predict an outcome that is binary (0 or 1).

If we used linear regression, we could get values outside the range [0, 1], which doesn't make sense for probabilities.

Logistic regression fixes this by applying the logistic (sigmoid) function to map predictions into the range [0, 1].

#### ☐ Real-Life Examples:

Email spam detection → Spam (1) or Not Spam (0)

Credit scoring → Default (1) or Not Default (0)

Medical diagnosis → Has disease (1) or Healthy (0)

Customer churn → Will leave (1) or Stay (0)

#### ☐ Extensions:

Multinomial Logistic Regression → Used for multi-class classification (e.g., predicting if a fruit is apple, orange, or banana).

Regularized Logistic Regression → Ridge (L2), Lasso (L1) to avoid overfitting.

## Different Equations of line

### 1 Slope Intercept Form

$$y = mx + b$$

Usefulness: This form is quickly identifying the slope and intercept of line.

### 2 Point Slope Form

$$y - y_1 = m(x - x_1)$$

Usefulness - This form is helpful when you know a point on the line and its slope.

### 3 Standard Form

$$Ax + By + C = 0$$

Usefulness - This form is useful for comparing different lines and finding intercepts.

### 4 Intercept Form

$$\frac{x}{a} + \frac{y}{b} = 1$$

a is the x intercept  
b is the y intercept

## 5 Vertical Line

$$x = a$$

Vertical line with  $a$  as the  $x$  intercept.

## 6 Horizontal Line

$$y = b$$

Horizontal Line with  $b$  as the  $y$  intercept.

## Logistic Regression.

We are going to use standard line equation

$$ax + by + c = 0$$

Eq. in form of  $\Theta$

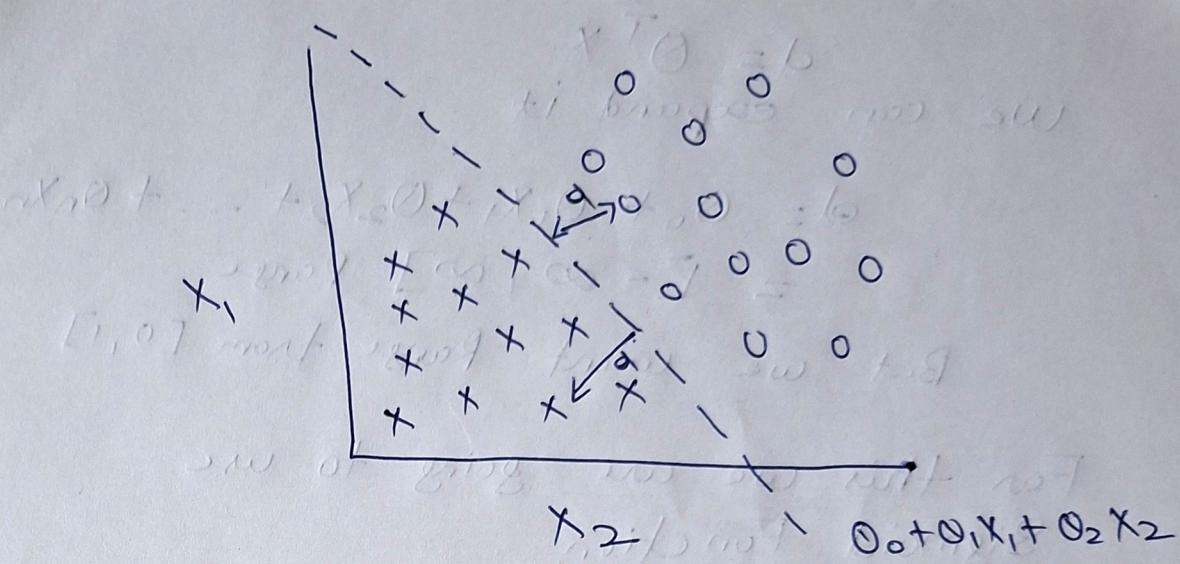
$$\Theta_0 + \Theta_1 x_1 + \Theta_2 x_2 = 0$$

We want a line which is best for separating  $\Theta$  data points.

If we have two features then

$$\Theta = [\Theta_0, \Theta_1, \Theta_2]$$

i.e  $(n+1)$  if we have  $n$  features



Distance of a point  
from a line  $d_1 = \frac{\theta_0 + \theta_1 x_1 + \theta_2 x_2}{\sqrt{\theta_1^2 + \theta_2^2}}$

bisecting to signed distance  $d_2 = \frac{\theta_0 + \theta_1 x'_1 + \theta_2 x'_2}{\sqrt{\theta_1^2 + \theta_2^2}}$

In both  $d_1$  and  $d_2$  denominator remain same only  $x_1$  and  $x_2$  changes

So

$$d \propto \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$d$  is +ve for ~~over~~ one side of line

$d$  is -ve for other side of line.

$$d = \theta^T X$$

We can expand it

$$d = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

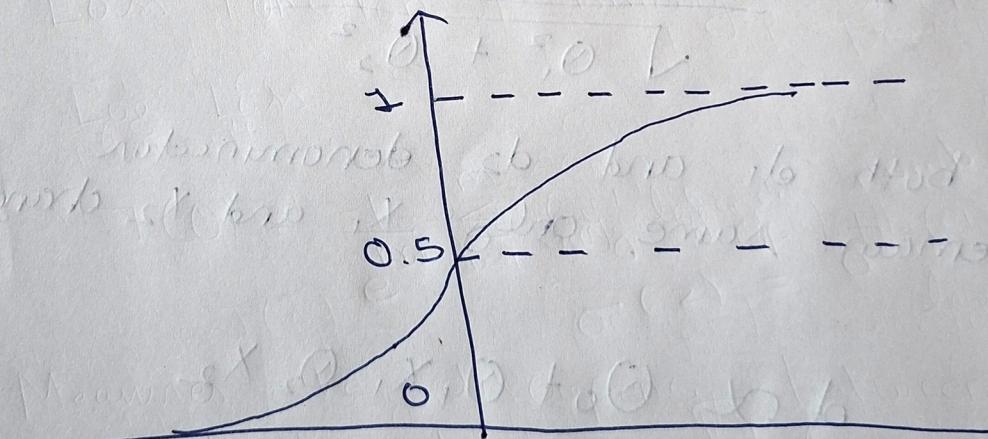
$$= [-\infty, \infty] \text{ Range}$$

But we want range from  $[0, 1]$

For this we are going to use  
Sigmoid Function.

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

Let's see the graph of sigmoid with different functions.



$$\text{when } z = \infty \quad \sigma(z) = \frac{1}{1+e^{-\infty}} = 1$$

$e^{-\infty}$  is very small approaching to zero

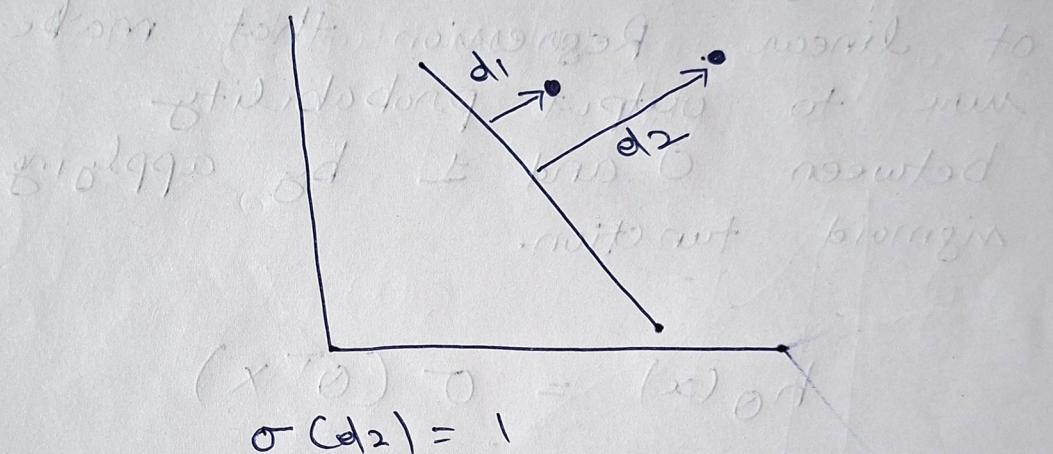
$$\text{so } \sigma(z) = 1 \text{ when } z = \infty$$

$$\text{When } z = 0 \quad \sigma(z) = \frac{1}{1+e^0} = \frac{1}{2} = 0.5$$

$$\text{When } z = -\infty \quad \sigma(z) = \frac{1}{1+e^{-\infty}} = 0$$

$e^\infty$  is very large number so

$$\sigma(z) = 0$$



It means farther is the point from the line we get more confidence in predicting  $\hat{y}$ .

Smaller the distance less is the confidence.

### Hypothesis for Logistic Regression

$$y_p = h_\theta(x) = \sigma(\theta^T x)$$

where  $x_0 = 1$

$$p = \frac{1}{1 + e^{-\theta^T x}}$$

where  $\theta^T x = \sum_{i=0}^n \theta_i x_i$

## Logit Model

The Logit Model is modification of linear regression that make use to output probability between 0 and 1 by applying sigmoid function.

$$h_\theta(x) = \sigma(\theta^T x)$$

## Loss Function - Binary Cross Entropy

### Log Loss

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log p^{(i)}) + (1-y^{(i)}) \log (1-p^{(i)})$$

Measures the discrepancy between predictions and true labels.

Case I +ve Samples

If  $y^{(i)} = 1$

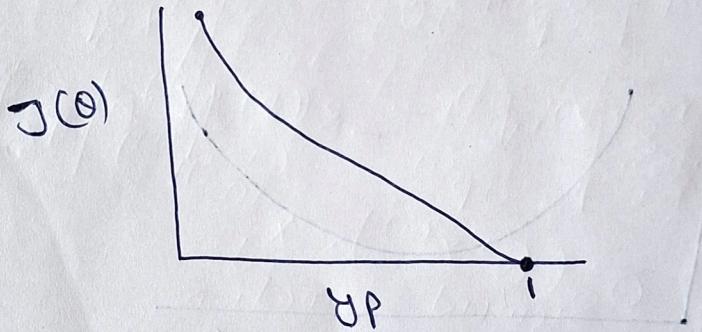
second term becomes zero

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log y_P^{(i)}$$

If  $y_P = 1 \log 1 = 0$

but  $\log 0 = -\infty$

If  $y_P = 0 \log 0 = +\infty$



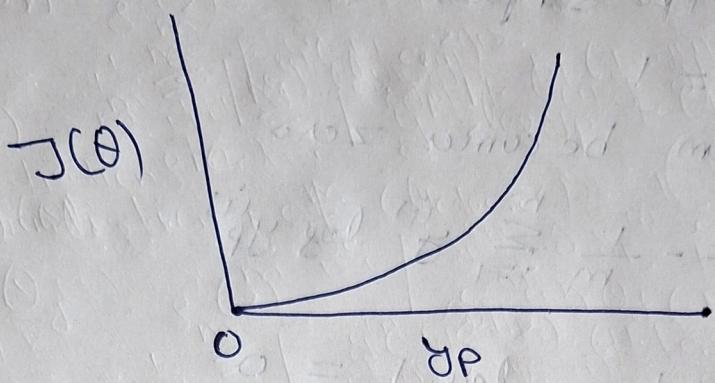
Case II -ve samples

$y_i = 0$  first term becomes zero

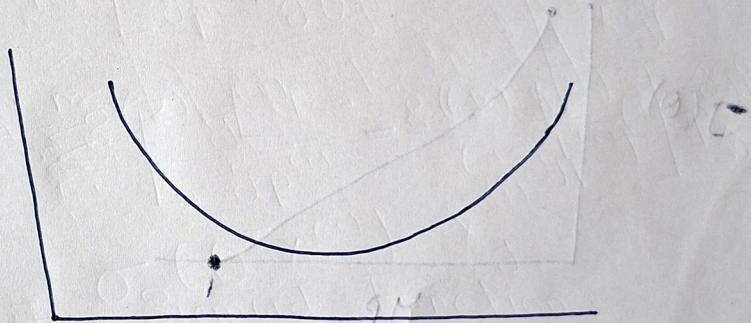
$$J(\theta) = -\frac{1}{m} (1-y^{(i)}) \log (1-y_P^{(i)})$$

If  $y_P = 0 1 \cdot \log 1 = 0$

If  $y_P = 1 1 \cdot \log 0 = +\infty$



If we combine two functions and their graph.



We can use Gradient Descent and it is a convex function (function in which we have one global minima and maxima).

## Differentiation of Log Loss Function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log y_p^{(i)} + (1-y^{(i)}) \log (1-y_p^{(i)}) \right]$$

Differentiation of  $\log x = \frac{1}{x}$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m \left[ \frac{y^{(i)}}{y_p^{(i)}} \frac{\partial y_p^{(i)}}{\partial \theta_j} + \frac{(1-y^{(i)})}{(1-y_p^{(i)})} \frac{\partial y_p^{(i)}}{\partial \theta_j} \right]$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m \left[ \frac{y^{(i)}}{y_p^{(i)}} \frac{\partial y_p^{(i)}}{\partial \theta_j} + \frac{(1-y^{(i)})}{(1-y_p^{(i)})} \frac{\partial y_p^{(i)}}{\partial \theta_j} \right]$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m \left[ \frac{y^{(i)}}{y_p^{(i)}} - \frac{(1-y^{(i)})}{(1-y_p^{(i)})} \right] \frac{\partial y_p^{(i)}}{\partial \theta_j} \quad (1)$$

## Hypothesis Sigmoid Function

$$\begin{aligned} y_p^{(i)} &= h_\theta(x^{(i)}) \\ &= \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n) \\ &= \sigma(\theta^\top x^{(i)}) \end{aligned}$$

## Differentiation of Sigmoid Function

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

$$\sigma'(z) = \frac{\partial \sigma(z)}{\partial z}$$

Differentiation of  $\frac{1}{x} = \frac{1}{x^2}$  and  $e^x = 1 \cdot e^x$

$$g'(z) = \frac{e^{-z}}{(1+e^{-z})^2}$$

$$g'(z) = \frac{e^{-z}}{(1+e^{-z})^2}$$

We can also write it as

$$g'(z) = \frac{1}{1+e^{-z}} - \frac{1}{(1+e^{-z})^2}$$

$$g'(z) = g(z) [1 - g(z)] \quad \text{--- (2)}$$

$$\frac{\partial y_p^{(i)}}{\partial \theta_j} = \sigma(\theta^T x^{(i)}) \cdot x_j^{(i)}$$

$$y_p^{(i)} = \theta^T x^{(i)}$$

$$\frac{\partial y_p^{(i)}}{\partial \theta_j} = y_p^{(i)} (1 - y_p^{(i)}) x_j \quad \text{--- (3)}$$

Using Eq (1)

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m \left[ \frac{y^{(i)} - y_p^{(i)} - \delta p^{(i)} + \delta^{(i)} y_p^{(i)}}{\delta p^{(i)} (1 - \delta p^{(i)})} \right] \frac{\partial \delta p^{(i)}}{\partial \theta_j}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[ \frac{y^{(i)} - y_p^{(i)}}{\delta p^{(i)} (1 - \delta p^{(i)})} \right] y_p^{(i)} (1 - \delta p^{(i)}) x_j$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] x_j$$

$$\boxed{\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - y_p^{(i)}] x_j}$$