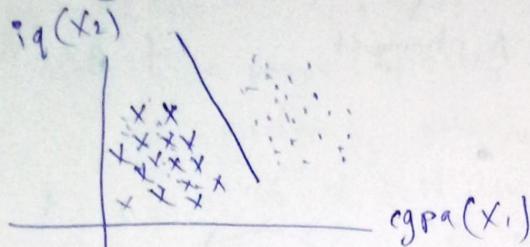
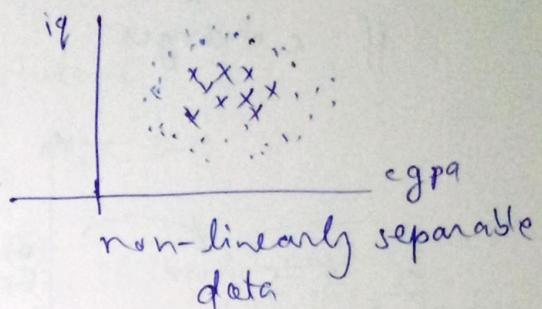


# Logistic Regression

- The prerequisites of applying the logistic regression is that the data is linearly classifiable.



linearly separable data

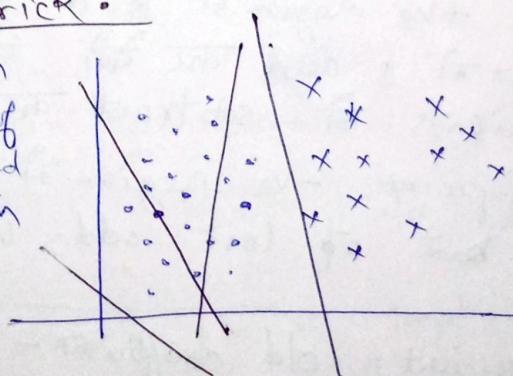


non-linearly separable data

- In logistic regression we use  $AX + BY + C = 0$   
or  $Ax_1 + Bx_2 + C = 0$  for two-D  
 $Ax_1 + Bx_2 + Cx_3 + D = 0$  for 3-D

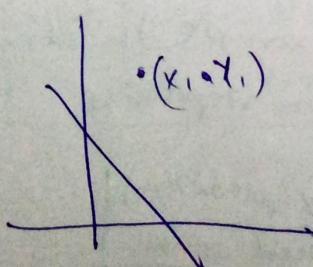
## perception Trick:

loop starts with random values of A, B and C and changes in each loop.



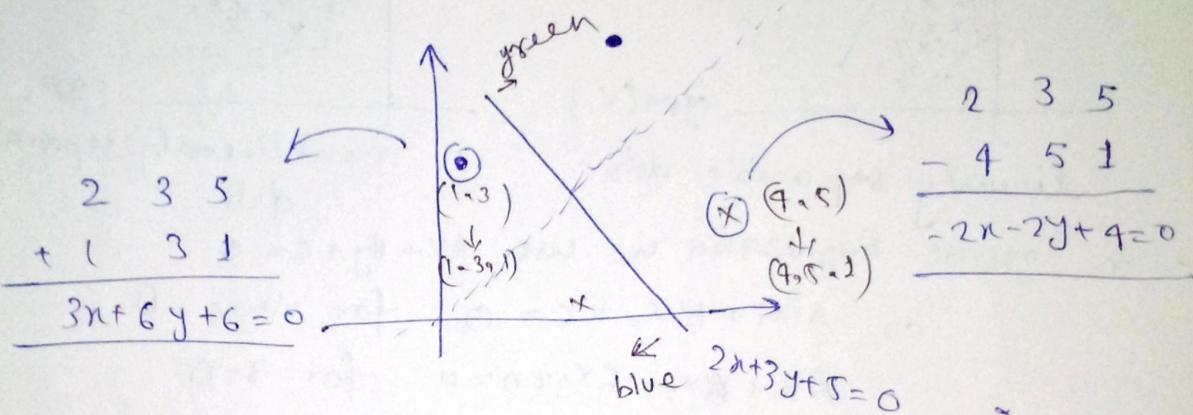
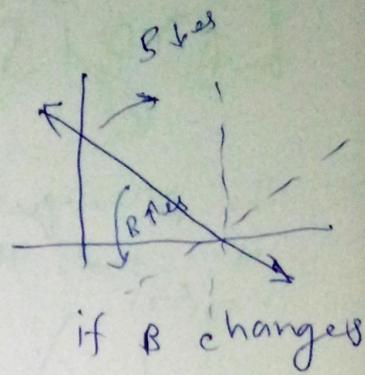
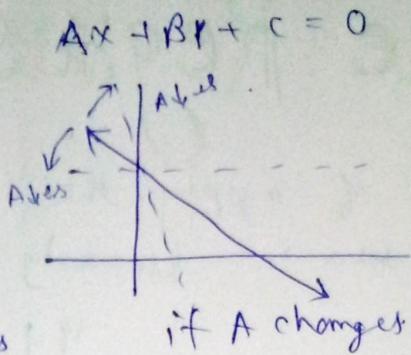
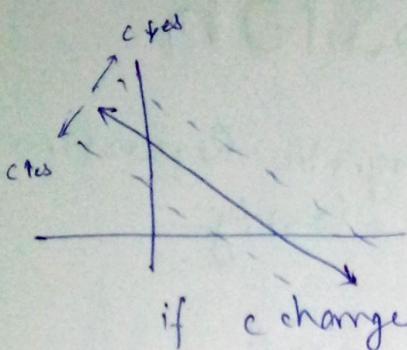
loop continues till convergence (all points are correctly classified).

## How to label region?



equation of line  $AX + BY + C = 0$   
 if  $AX_1 + BY_1 + C > 0$  positive region.  
 if  $AX_1 + BY_1 + C < 0$  negative region  
 if  $AX_1 + BY_1 + C = 0$  on line.

## How to transform line?



Both -ve point +ve region में है तो उसे उसके coordinates में last में 1 add कर देंगे और उस line के coefficients से subtract कर देंगे।  
जब फिर +ve point -ve region में है तो same process as above but at last add. with line's coefficient.

new coefficient = old coefficient -  $\eta$  coordinates.  
learning rate ( $\eta$ ) = 0.01

in not of the long transformation example in above we can't do that, we do this in stepwise

### Algorithm

$x_0$	$x_1$	$x_2$	y (prediction)
1	cgpa	iq	placed or not
1	7.5	81	1
1	8.9	109	1

$$AX + BY + C = 0$$

$$w_0x_0 + w_1x_1 + w_2x_2 = 0 \rightarrow$$

$$w_0 = C, w_1 = A, w_2 = B$$

if  $\sum_{i=0}^n w_i x_i \geq 0$  then 1

else  $\sum_{i=0}^n w_i x_i < 0$  then 0

$$\text{epoch} = 1000 \quad \eta = 0.01$$

for i in range (epochs):

randomly select a student

if  $x_i \in N$  and  $\sum_{i=0}^n w_i x_i \geq 0$

here if model predict 1 but in data it is zero then here we update the W's

→ point → region #  $\frac{\eta}{2}$

$$W_{\text{new}} = W_{\text{old}} - \eta [x_0 \ x_1 \ x_2]$$

$$W_{\text{new}} = W_{\text{old}} - \eta x_i$$

if  $\sum_i \epsilon P$  and  $\sum_{i=1}^n w_i x_i > 0$  then point -ve region #  $\frac{\eta}{2}$

$$W_{\text{new}} = W_{\text{old}} + \eta x_i$$

$x_i \in N$ and $\sum w_i x_i \geq 0$	$W_n = w_0 - \eta x_i$	(ii)
$x_i \in P$ and $\sum w_i x_i < 0$	$W_n = w_0 + \eta x_i$	(iii)

simplified formula  
 $W_n = w_0 + \eta (y_i - \hat{y}_i) x_i \quad \text{--- (iii)}$

actual prediction

$$\hat{y}_i \quad y_i \quad y_i - \hat{y}_i$$

case-1

$$1 \quad 1 \quad 0$$

case-2

$$0 \quad 0 \quad 0$$

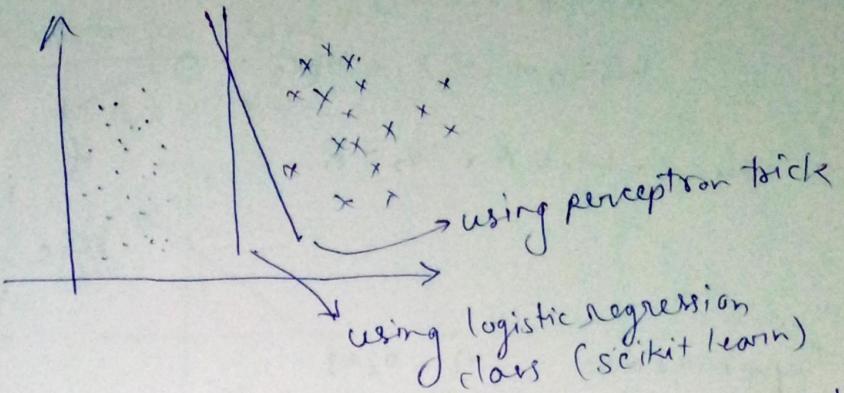
case-3

$$1 \quad 0 \quad 1$$

case-4

$$0 \quad 1 \quad -1$$

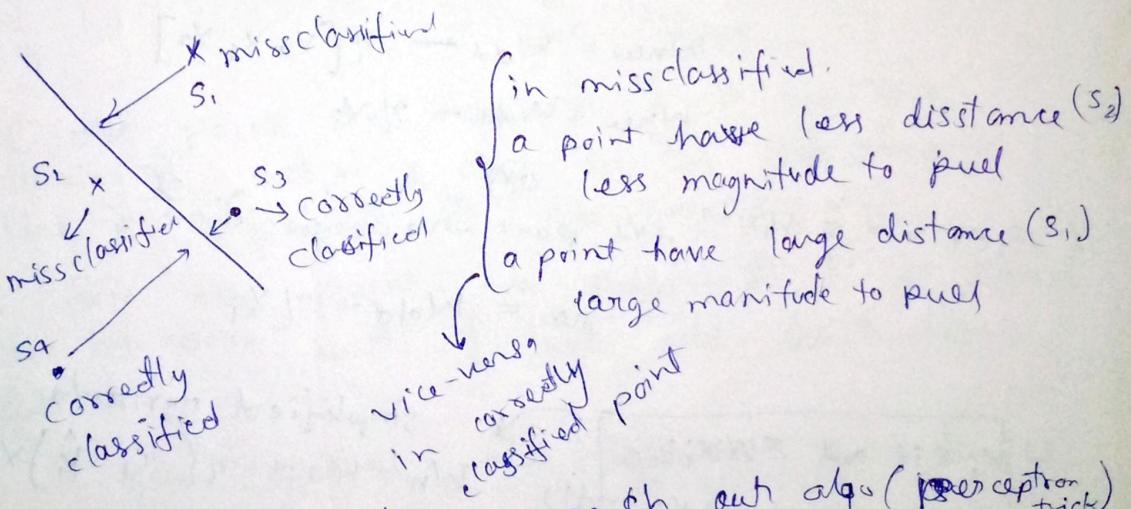
in case 1 and 2 no updates in weight but in case 3 and case -4 weight equation becomes eqn (iii) and (ii)



If we even the perceptron trick model is work as the logistic regression class then here we also considered the correctly classified class point also ~~missclassified~~ correctly classified point line तक पहुँचने से तक पहुँचने तक Till the ~~equilibrium~~ phase.

misclassified — line pull

correctly classified — line push



\* for all of these above approach our algo (perception trick) is work similar as the logistic regression class.

\* if we apply the above approach then our plan have to make changes in the equation (iii)  
or we can also say that the case 3 and case 4

$$w_n = w_0 + \gamma(Y_i - \hat{Y}_i)x_i$$

↓  
come zero due to this

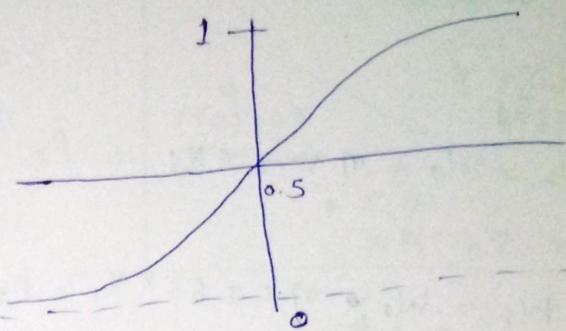
$y_i - \hat{y}$  is come zero due to  $\sum w_i x_i$  this function that is step function that output always the 0 or 1 hence now we use the sigmoid function.

sigmoid function

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

if  $-\infty < z < \infty$

then  $0 < y < 1$



$$z = \sum w_i x_i$$

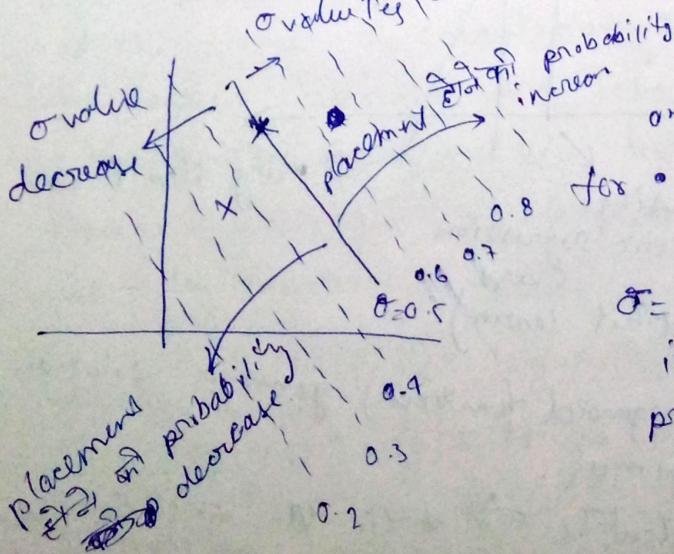
if  $z$  +ve. then  $\sigma(z) > 0.5$

if  $z$  -ve. then  $\sigma(z) < 0.5$

steps:- first calculate  $z$

Step. 2 - put this  $z$  in equation (iv)  
and calculate

Step. 3 if  ~~$\sigma(z) = 0.5$~~   $\rightarrow$  output  
 $z$  +ve then  $\sigma(z) \geq 0.5 \rightarrow \boxed{1}$   
 $z$  -ve then  $\sigma(z) < 0.5 \rightarrow \boxed{0}$



online  $\sum w_i x_i = 0 \quad \sigma = 0.5$

for  $\sum w_i x_i > 0 \quad \sigma > 0.5$

$\sigma = 0.5$  means  $P(0.5) = ?$   
it means placement  $\stackrel{\text{at}}{\text{at}}$  probability is 0.5

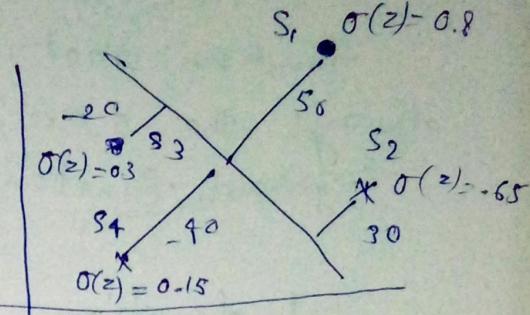
$$z = \sum w_i x_i$$

Impact on Sigmoid

$$w_n = w_0 + \eta(\hat{y}_i - y_i) x_i$$

$$y_i = \sigma(z)$$

$$\text{where } z = \sum w_i x_i$$



for  $S_1$

$$w_n = w_0 + \eta \times 0.2 \times x_i = p_1$$

for  $S_2$

$$w_n = w_0 + \eta \times 0.65 \times x_i = p_3$$

$y_i$	$\hat{y}_i$	$y_i - \hat{y}_i$
1	0.8	0.2
0	0.65	-0.65
1	0.3	0.7
0	0.15	-0.15

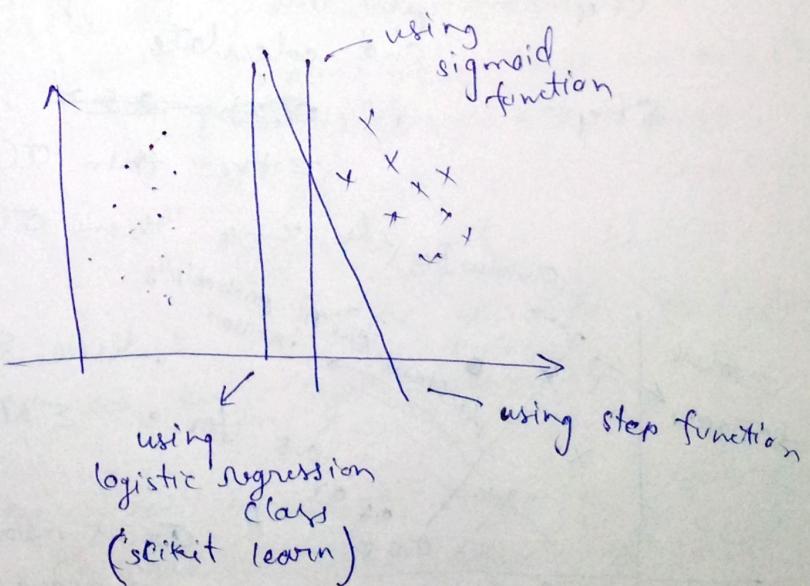
for  $S_3$

$$w_n = w_0 + \eta \times 0.7 \times x_i = p_2$$

for  $S_4$

$$w_n = w_0 + \eta \times 0.15 \times x_i = p_4$$

$$p_2 > p_1 \text{ etc.}$$



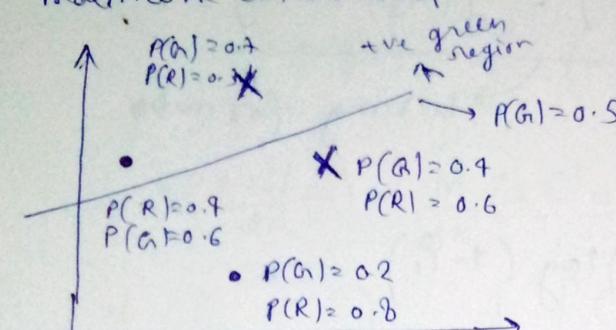
- \* regression line (using sigmoid function) gives a solution but it is not a optimal.

- \* इसीले एही loss function calculate करे और उसके बराबर करने के लिए

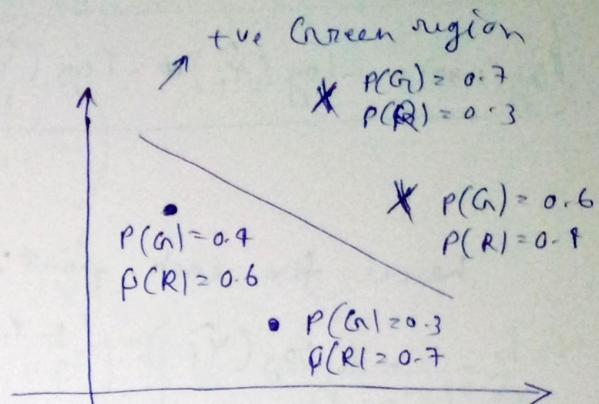
Here we now for deriving the loss function we use the maximum likelihood.

### Loss Function Calculation:

maximum likelihood



model-1



model-2

\* in the above  $\hat{y}_i$  is calculated as a probability term

model-1 (maximum likelihood)

$$0.7 \times 0.4 \times 0.4 \times 0.8 = 0.089$$

model-2 (maximum likelihood)

$$0.7 \times 0.6 \times 0.6 \times 0.7 = 0.176$$

\* here the model 2 is better than the model-1

\* in the above it is easy to calculate due to less number what about large dataset?

at the time of calculation of maximum likelihood the values becomes very less and then this comparison (compare b/w two model) is useless hence we again re-formulate the formula and use log() function.

$$\log(a \cdot b) = (\log a + \log b)$$

$$\log(\text{mod}) = \log(0.7) + \log(0.4) + \log(0.4) + \log(0.8)$$

$$\log(\text{mod}) = -\log(0.7) - \log(0.4) - \log(0.4) - \log(0.8) \rightarrow \text{cross entropy}$$

Cross entropy में हम values को minimize करते हैं, ताकि यह model Best होती जिसका cross entropy minimum होता।

$$\log(\max) = -\log(\hat{Y}_1) + -\log(\hat{Y}_2) + \dots + \log(\hat{Y}_n)$$

→ wrong formula

Hence for each point:

$$-\gamma_i \log(\hat{Y}_i) - (1-\gamma_i) \log(1-\hat{Y}_i)$$

$$\text{Loss function}(l) = \sum_{i=1}^m -\gamma_i \log(\hat{Y}_i) - (1-\gamma_i) \log(1-\hat{Y}_i)$$

for average

$$L = \frac{1}{n} \sum_{i=1}^n \gamma_i \log(\hat{Y}_i) + (1-\gamma_i) \log(1-\hat{Y}_i)$$

The above equation is called log loss error or Binary cross entropy

minimize it using and find  $w_1, w_2, \dots, w_n$  such that it have minimum value

\* for calculating closed form of the above loss function is using gradient descent

$$\text{Loss function}(L(w_1, w_2, b)) = \frac{1}{n} \sum_{i=1}^n L(Y_i, f(x_i)) + \alpha R(w_1, w_2)$$

$$E(w, b) = \frac{1}{n} \sum_{i=1}^n L(Y_i, f(x_i)) + \alpha R(w)$$

Regularization

Here  $L$  = Loss function and  $R$  = Regularization

for perceptron -  $L(Y_i, f(x_i)) = \max(0, -Y_i f(x_i))$

Here  $f(x_i) = w_1 x_1 + w_2 x_2 + b$ ;  $n = \text{rows}$

# Derivative of Sigmoid function:

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

$$\sigma'(x) = \frac{d}{dx} \left( \frac{1}{1+e^{-x}} \right)$$

$$\therefore \frac{d}{dx} \left( \frac{1}{x} \right) = \frac{d}{dx} (x)^{-1} = -\frac{1}{x^2}$$

$$\frac{d}{dx} \left[ \frac{1}{1+e^{-x}} \right] = \frac{d}{dx} \left[ (1+e^{-x})^{-1} \right] = -\frac{1}{(1+e^{-x})^2} \frac{d}{dx} (1+e^{-x})$$

$$= -\frac{1}{(1+e^{-x})^2} \frac{d}{dx} (e^{-x})$$

$$= -\frac{e^{-x}}{(1+e^{-x})^2} \frac{d}{dx} (-x)$$

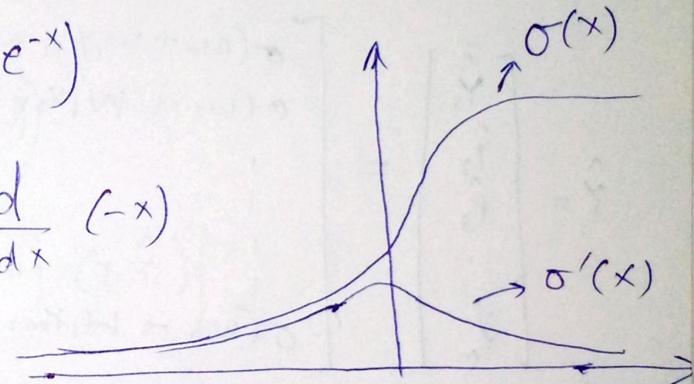
$$= \frac{e^{-x}}{(1+e^{-x})^2}$$

$$= \frac{1 \cdot e^{-x}}{(1+e^{-x})(1+e^{-x})} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}}$$

$$= \sigma(x) \left[ \frac{e^{-x}}{1+e^{-x}} \right] = \sigma(x) \left[ \frac{1+e^{-x}-1}{1+e^{-x}} \right]$$

$$= \sigma(x) \left[ \frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right]$$

$$\sigma'(x) = \sigma(x) [1 - \sigma(x)] =$$



# Gradient Descent on logistic Regression's Loss Function:

Dataset: Rows = m columns = n

$$\begin{array}{ccccccc}
 L & 2 & 3 & - & - & - & n & Y \\
 X_{11} & X_{12} & X_{13} & - & - & - & X_{1n} & Y_1 \\
 X_{21} & X_{22} & X_{23} & - & - & - & X_{2n} & Y_2 \\
 & & & & & & & \\
 & & & & & & & \\
 & & & & & & & \\
 X_{m1} & X_{m2} & X_{m3} & - & - & - & X_{mn} & Y_m
 \end{array}$$

$$\hat{Y} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \hat{Y}_3 \\ \vdots \\ \hat{Y}_m \end{bmatrix} = \begin{bmatrix} \sigma(w_0 + w_1 X_{11} + w_2 X_{12} + \dots + w_n X_{1n}) \\ \sigma(w_0 + w_1 X_{21} + w_2 X_{22} + \dots + w_n X_{2n}) \\ \vdots \\ \sigma(w_0 + w_1 X_{m1} + w_2 X_{m2} + \dots + w_n X_{mn}) \end{bmatrix}$$

$$\hat{Y} = \sigma(XW)$$

$$\hat{Y} = \sigma \left( \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1n} \\ 1 & X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{m1} & X_{m2} & \dots & X_{mn} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \right)$$

$$L = -\frac{1}{m} \sum_{i=1}^m Y_i \log(\hat{Y}_i) + (1-Y_i) \log(1-\hat{Y}_i)$$

or

$$L = -\frac{1}{m} \left[ \sum_{i=1}^m Y_i \log(\hat{Y}_i) + \sum_{i=1}^m (1-Y_i) \log(1-\hat{Y}_i) \right]$$

for 1st term

$$\sum_{i=1}^m \gamma_i \log(\hat{\gamma}_i) = \gamma_1 \log(\hat{\gamma}_1) + \gamma_2 \log(\hat{\gamma}_2) + \gamma_3 \log(\hat{\gamma}_3) + \dots + \gamma_m \log(\hat{\gamma}_m)$$

$$= [\gamma_1 \ \gamma_2 \ \gamma_3 \ \dots \ \gamma_m] \cdot \begin{bmatrix} \log \hat{\gamma}_1 \\ \log \hat{\gamma}_2 \\ \vdots \\ \log \hat{\gamma}_m \end{bmatrix}$$

$$= [\gamma_1 \ \gamma_2 \ \dots \ \gamma_m] \cdot \log \left( \begin{bmatrix} \hat{\gamma}_1 \\ \hat{\gamma}_2 \\ \vdots \\ \hat{\gamma}_m \end{bmatrix} \right)$$

$$= \gamma \cdot \log(\hat{\gamma})$$

$$= \gamma \cdot \log(\sigma(xw))$$

for 2nd term

$$\sum_{i=1}^m (1-\gamma_i) \log(1-\hat{\gamma}_i)$$

$$= (1-\gamma) \log(\sigma(1-xw))$$

then Now whole formula for loss function

$$L = -\frac{1}{m} \left[ \gamma \log \hat{\gamma} + (1-\gamma) \log(1-\hat{\gamma}) \right]$$

when  $\hat{\gamma} = \sigma(xw)$

Here the dot product of  $x$  and  $w$

loss function in matrix form:

$$\text{Loss function} = \frac{1}{m} \left[ \underbrace{\gamma \log(\sigma(wx))}_{\text{for } y=1} + \underbrace{(1-\gamma) \log(1-\sigma(wx))}_{\text{for } y=0} \right]$$

यहाँ पर हमें सारे coefficients ( $w$ ) का एक समान से समान value त्रिकोणित है जो पर loss function is minimum. To solve this problem we use gradient descent.

approach is

- start with random value of  $w = [$

- for  $i = 1$  in epochs:

$$w = w - \eta \frac{\partial L}{\partial w} \quad \text{--- (i)}$$

$$\left( \frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n} \right)$$

$$L = \frac{1}{m} \left[ \underbrace{\gamma \log \hat{y}}_{\text{for } y=1} + \underbrace{(1-\gamma) \log(1-\hat{y})}_{\text{for } y=0} \right]$$

$$\frac{\partial L}{\partial w} = ?$$

$$\text{1st form} \quad \frac{d}{dw} \gamma \log \hat{y} = \gamma \frac{d}{dx} \log \hat{y}$$

$$= \frac{\gamma}{\hat{y}} \frac{d}{dx} (\hat{y})$$

$$= \frac{\gamma}{\hat{y}} \frac{d}{dx} (\sigma(wx))$$

$$= \frac{\gamma}{\hat{Y}} \cdot \sigma(wx) [1 - \sigma(wx)] \frac{d}{dw} (\omega x)$$

$$= \frac{\gamma}{\hat{Y}} \cdot \hat{Y} (1 - \hat{Y}) x$$

$$= \gamma (1 - \hat{Y}) x$$

2nd term

$$\frac{d}{dw} (1 - \gamma) \log(1 - \gamma) = (1 - \gamma) \frac{d}{dw} \log(1 - \gamma)$$

$$= \frac{(1 - \gamma)}{(1 - \hat{Y})} \frac{d}{dw} [1 - \hat{Y}]$$

$$= - \frac{(1 - \gamma)}{(1 - \hat{Y})} \frac{d}{dw} \sigma(wx)$$

$$= - \frac{(1 - \gamma)}{(1 - \hat{Y})} [\sigma(wx) \cdot [1 - \sigma(wx)]] \frac{d}{dw} (\omega x)$$

$$= - \frac{(1 - \gamma)}{(1 - \hat{Y})} \hat{Y} (1 - \hat{Y}) x$$

$$= - \hat{Y} (1 - \gamma) x$$

Now the whole term is

$$\frac{dL}{dw} = -\frac{1}{m} \left[ \gamma (1 - \hat{Y}) x - \hat{Y} (1 - \gamma) x \right]$$

$$= -\frac{1}{m} \left[ \gamma (1 - \hat{Y}) - \hat{Y} (1 - \gamma) \right] x$$

$$= -\frac{1}{m} \left[ \gamma - \gamma \cancel{Y} - \hat{Y} + \hat{Y} \cancel{Y} \right] x$$

$$\boxed{\frac{dL}{dW} = -\frac{1}{m} (\hat{Y} - Y) X}$$

replace this equation in (i)

$$W = W + \eta \frac{1}{m} (\hat{Y} - Y) X$$

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}_{(m, n+1)} \quad Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{bmatrix}_{(m, 1)} \quad \hat{Y} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_m \end{bmatrix}_{(m, 1)}$$