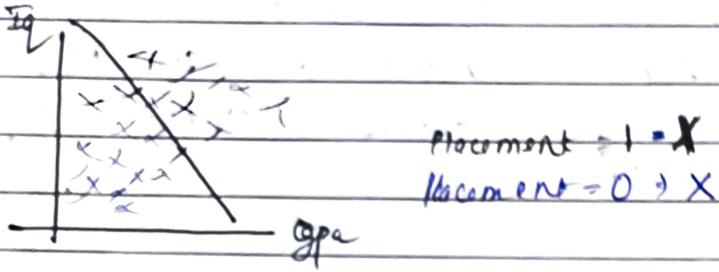


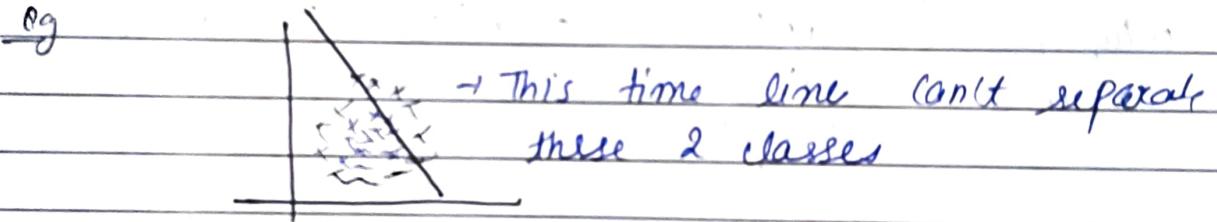
LOGISTIC REGRESSION

→ refer to Notepad for kypets

→ works well on linear / sort of linear data
eg cgpa iq placement



→ works not well on non linear data



→ multiple interpretations are there for LR

SOME BASIC GEOMETRY

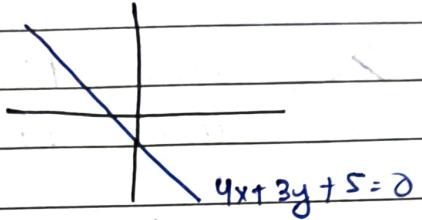
- In linear regression we used $y = mx + b$ as a eqn of line
- In logistic we will use $\Rightarrow Ax + By + c = 0$

Concepts

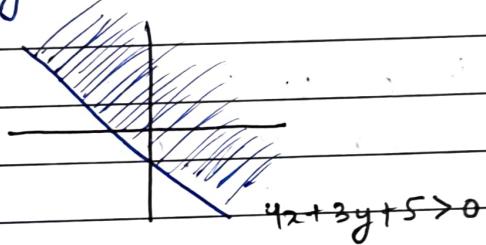
M	T	W	T	F	S	S
Page No.:						
Date:					YOUVA	

① Every line has a +ve side and -ve side

e.g. \rightarrow line having line $\Rightarrow 4x + 3y + 5 = 0$

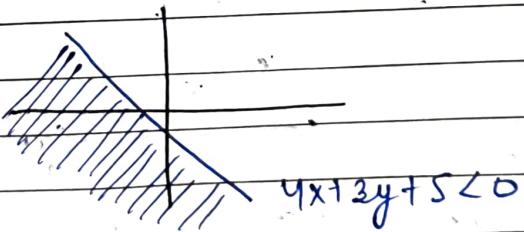


- positive side of line
- $4x + 3y + 5 > 0$



- Negative side / region of line

- $4x + 3y + 5 < 0$



② How to find out if a pt lies on line

- $4x + 3y + 5 = 0 \rightarrow$ line

- $(5, 2) \rightarrow$ pt

- Find if pt lies on that line

- $4x + 3y + 5 = 0$

- $4(5) + 3(2) + 5$

$$20 + 6 + 5 = 31, 31 \neq 0, \therefore \text{pt doesn't lie on line}$$

③ How to find if given pt is on positive side of line or the negative side

- if $Ax_1 + bx_2 + c \geq 0$, then pt is in the region
- if $Ax_1 + bx_2 + c < 0$, then pt is in -ve region

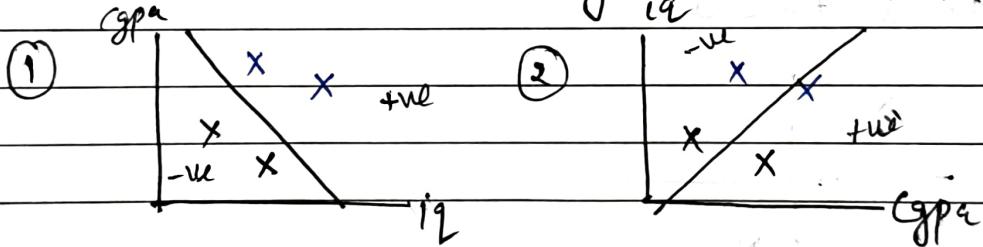
$$\text{eg} \rightarrow 4x + 3y + 5, \text{ pt} \Rightarrow (2, 3)$$

$$= 4(2) + 3(3) + 5$$

$$= 8 + 9 + 5$$

= 22, since it is > 0 ; then pt is in the region of line

- Ques We have 2 models for classification which model you will select.



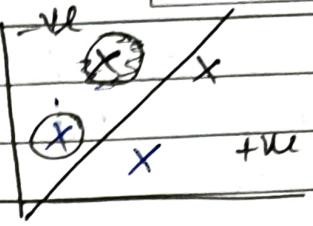
Ans • Model 1 performs better, as it was able to separate 2 blue pts and 2 black pts.

• Model 2 is doing misclassification.

• Model 1 is great model as it is correctly classified all the pts whereas model 2 is not good as it is doing certain misclassif.

- In classification, we have to make sure the +ve pts land in +ve region and -ve pts lands in -ve region
- How did we come to conclusion that model 1 is better?
- Because of the no. of missclassification pts ie in model 1 there were no missclassified pts, but in model 2 there were 2 missclassified pts
- How to check through code if the pt is missclassified or not?
 - • loop on all pts.
 - considering ^{blue} green pts in data are positive and black pts in data to be negative
 - we will loop through all pts
 - and ask what is the color of the pt eg: → we have blue, and now will ask our eqⁿ of the is $Ax + By + c = 0$ that the blue pt is true or not. If $Ax + By + c > 0$ that means according to model the green pt is in +ve region and in actually also it was a +ve pt
 - It is correct classification

- Whereas :- if model is \Rightarrow



- let's say we picked up blue circle pt. Acc to data blue pt is +ve pt.

- let's ask the model line eqn $\Rightarrow Ax_1 + By_1 + C = 0$
- say line says $Ax_1 + By_1 + C < 0$, which means acc to line eqn it is -ve but our data is saying blue pt is +ve
- This is example of missclassification

II logic would be to count no of missclassifications pts for particular model.

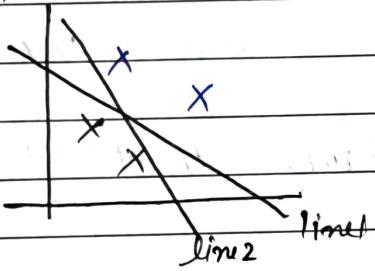
count = 0

```
if ( pt → blue & Ax1 + By1 + C < 0 ) {
    missclassification
    count++
}
```

```
else if ( pt → black & Ax1 + By1 + C > 0 )
    missclassification
    count++
```

- We will apply this logic on both models.
- In whichever model the count is more, it is not good model.

- This is not logistic Regression, it is very simplified version for classifications.
- And there is also a problem eg we have 2 models, where in both cases the model has zero missclassified pts or you can say equal no of missclassified pts.
- In this case, which model would you choose. Eg in diagram



- Both lines were able to separate the 2 classes, now which line is the best we don't know.
- This is the problem with this approach
- Now, we will study, how Logistic regression solves this problem.
- Classification should not be so black and white i.e. means if we should not be so straight forward
- a better approach would instead of directly calling pts position above the line, we should assign some

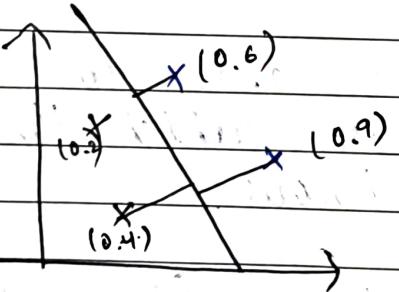
some no and the pts that are more far away from line, the number should be that larger representing that this pt is more positive.

And if the pt is more closer to line the number should be smaller representing that the pt is less far away from line.

Same goes for -ve region, the pts that are closer to line we should assign less no whereas pts that are far away from line should be assigned larger no.

This is for just to show that distance matter.

Eg



Analogy would be :- our basic approach was if our percentage is less than 50, we will not get job and if more than 50, we will get that job.

But our new approach is saying that if our % is more than 50 placement to ho gayegi but agar aapke percentage 95, 96 around hai

- But our new approach is saying 50% jitna zada upar hai in 96, 98 apki job or achi lagni chye. and same goes for jitna zada niche from so wohi bhar job.

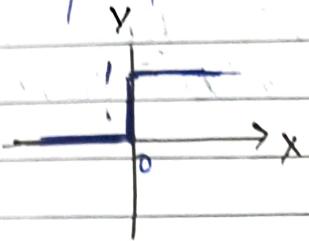
- How Normal approach is doing prediction?

data \Rightarrow	iq	cgpa	placement	
	60	6	0	\Rightarrow
	40	4	0	x
	80	8	1	x
Query \Rightarrow	90	9	1	x

- Classification model $\Rightarrow Ax + By + C = 0$; Assume $= 3x + 5y + 6 = 0$
- Query $\Rightarrow \{90, 9\}$
- $3(90) + 5(9) + 6 = 150$
- $150 > 0$, Hence output $\Rightarrow 1$ (that means blue)
- if $150 < 0$, then output would be 0 (that means black)
- In machine learning ex
- classification model $\Rightarrow Ax + By + C \leq 0 \Rightarrow$
 - $\beta_0 + \beta_1 x_1 + \beta_2 x_2$
 - $\beta_0, \beta_1, \beta_2$ are coeff, x_0, x_1 are pts x_1 and x_2
- Technically we are finding $\Rightarrow \beta_0 + \beta_1 x_1 + \beta_2 x_2$
- Let's call it $Z = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
- Then we are checking on Z , i.e. if $Z > 0 \rightarrow 1$, and if $Z < 0 \rightarrow 0$

$Z < 0$	0
$Z > 0$	1

- if you notice this is a mathematical $f^n \rightarrow$ step function



- Technically, we are calculating $z = B_0 + B_1 x_1 + B_2 y_2$, and we are putting it in a Step fn \Rightarrow Step(z) which gives 0, or 1 which is used for prediction
- But in new method, we don't want to use fn that outputs 0 or 1, we want a continuous output for each pt, that represents if a no assigned is more, the ~~pt~~ is more far away and if no is less, then pt is less far away from line.
- Since we don't want behaviour of 0 or 1 therefore we will stop using Step fn
- We want something that if $z = B_0 + B_1 x_1 + B_2 y_2$ is more, then we want a large no in decimal

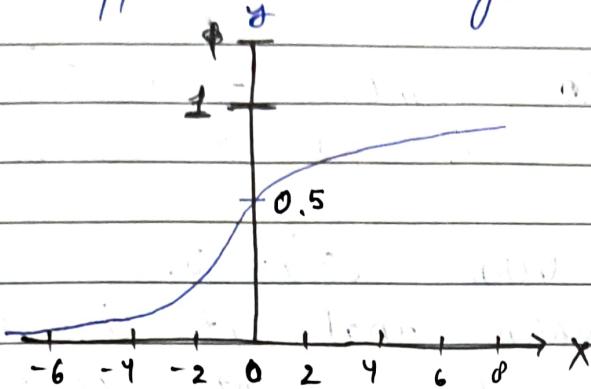
$z \uparrow = 0.9$ (z too much big, decimal also big)

$z \uparrow = 0.6$ (z big, decimal big)

$z \downarrow = 0.1$ (z too much small, decimal too much small)

$z \downarrow = 0.2$ (z small, decimal small)

- This is where a new function appears \rightarrow Sigmoid fn

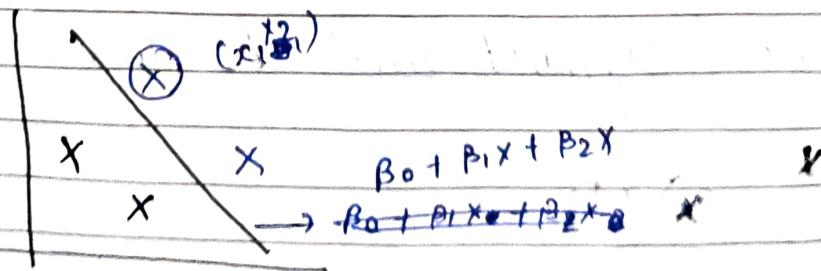


$$y = \frac{1}{1 + e^{-x}}$$

- If x value is large, then it tries to touch 1 (At infinite, y will be 1)
- If x value is too much small ($-ve$) then it tries to touch 0 (At infinite small value of x , y value will be 0)
- If x value is 0, then $y = 0.5$

- $x \rightarrow \text{infinity} (\infty), y = 1$
- $x \rightarrow 0, y = 0.5$
- $x \rightarrow -ve \text{ infinity } (-\infty), y = 0$

Let's exchange it with step



- Let's calculate z for $(\cancel{x_1}, \cancel{x_2}) (x_1, x_2)$

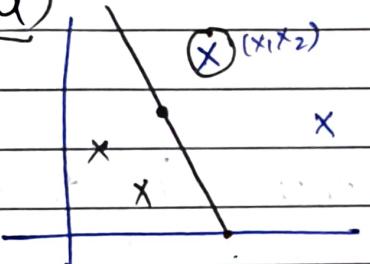
$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2, \quad (x_1, x_2) \Rightarrow (z, y)$$

$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

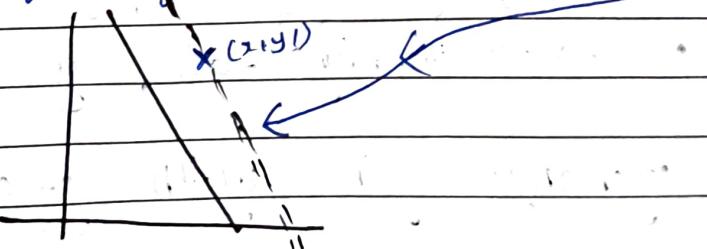
- give this z to sigmoid $\Rightarrow \sigma(z)$
- From this you will get a no between 0 and 1
- Note if $z = 0$, then $\sigma(z) = 0.5$, and z is 0 when that pt lies on line, that means if any pt lies on line, then its sigmoid = 0.5

(A) $\rightarrow (+w \text{ side})$

I.

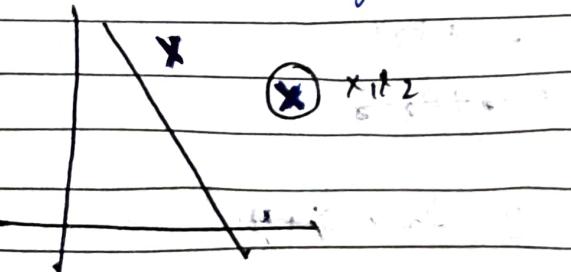


- Circled pt $\neq x_1, x_2$, for this calculate z , obviously its $z > 0$, not too much greater than 0, \therefore sigmoid will be slight more than $\sigma = 0.5$, let's say $\Rightarrow \sigma(z) = 0.6$, that means if any pt comes on,

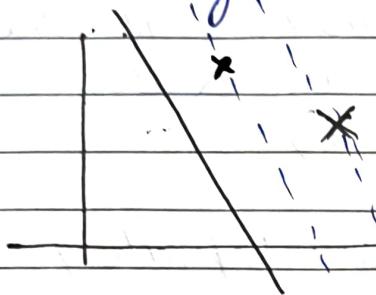


then it's sigmoid will be 0.6

II.

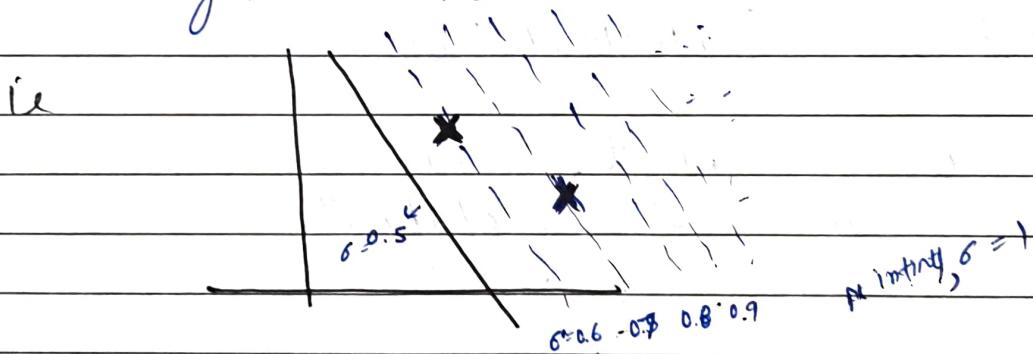


Now the circled pt z will move than the other, therefore its sigmoid will be slight higher than before let's say $\sigma(z) = 0.7$



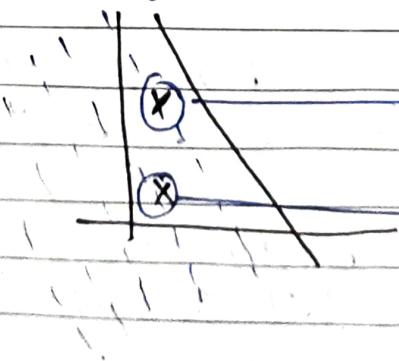
→ that means any pt on this line sigmoid will be 0.7

- III. Similarly, thoda aur age gayeng sigmoid will be more, and if infinite move then sigmoid will be 1



(B) (-ve side)

- Similarly will happen to -ve side



$$\therefore z < 0, \sigma(z) = 0.4 \\ (\text{not too much})$$

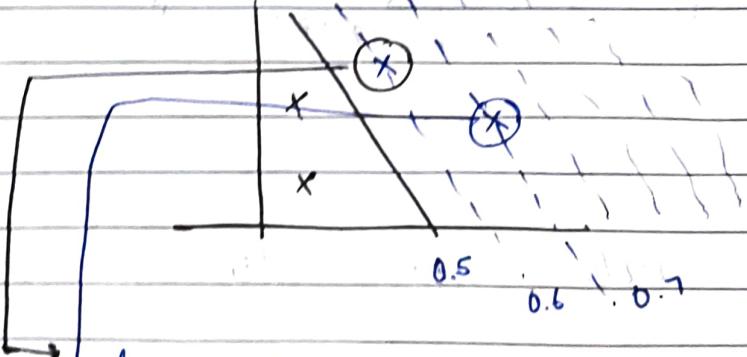
$$\therefore z < 0, \text{ smaller}, \sigma(z) = 0.3$$

- Similarly, when z is - large $\sigma(z) = 0$

M	T	W	T	F	S	S
Page No.						
Date						YODVA

- That means, our classification region is converted into gradient
- If we move infinitely in +ve, then sigmoid = 1, that means if there is pt at infinity we are 100% sure it is positive (+)
- Similarly, if we move infinitely in -ve and there is pt at -infinity then we are 100% sure that it is negative (-)
- In a way, notice there is nice probabilistic interpretation.
- Now for a point, we can say in terms of probability

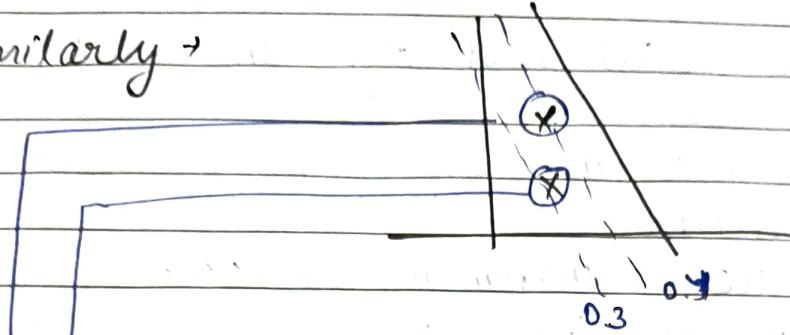
Eg



for this pt, the probability of this pt to be +ve is 0.6 & -ve is 0.4

for this pt, probability of this pt to be +ve is 0.7 & -ve is 0.3

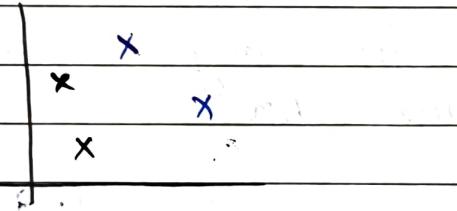
Similarly →



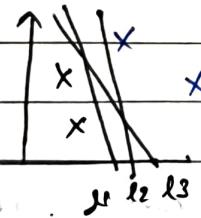
- Prob of this pt being +ve is 0.4,
Hence p of this pt -ve is 0.6
- Prob of this pt +ve is 0.3, ∴ p of
this pt -ve is 0.7

Maximum Likelihood

• Data



- Ultimately we want a line that does best classification
- There could be multiple lines

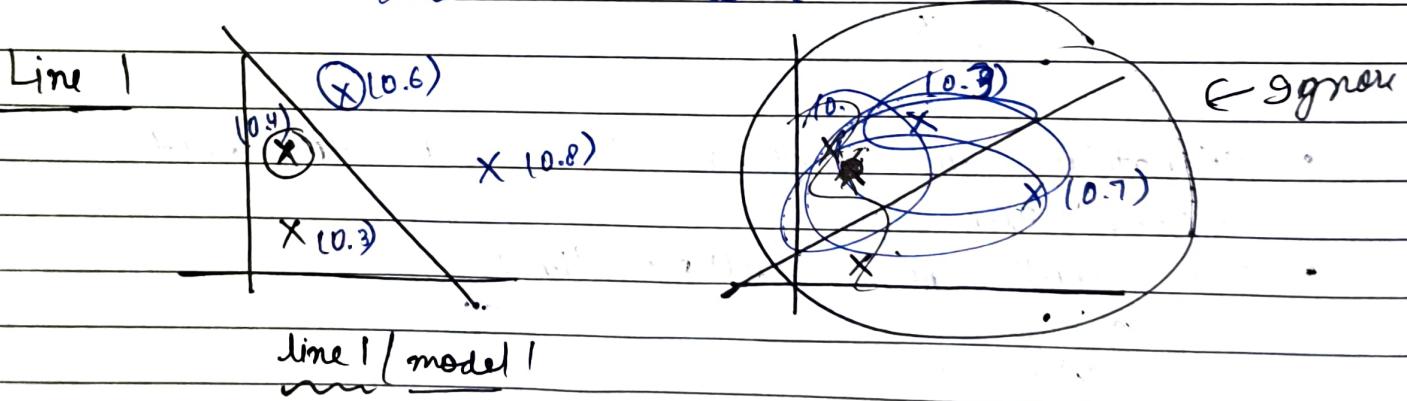


- These all lines are doing classification properly
- But we want a best line.
- To find best line in Machine Learning, we use loss function
- we minimize the loss function using for given parameters and for which parameters the loss fn is minimum is the best line

- Therefore we will first try to find out least loss function.
- Loss fn is just the logic which tells which line is better
- We will use Maximum Likelihood, which is the loss fn for Logistic regression

INTUITION OF WHICH LINE IS BETTER

- We will use a likelihood function: first we will assign (for now without calculation) each pt the probability of pt of 1 (in blue) w.r.t to our line.
- pts that are far away, \therefore prob will lie more and vice versa



- We will calculate likelihood: which basically means multiplying probabilities of each model for actual class.
- ie means P of $pt \rightarrow X_i$.
ie means eg the circled pt is actually blue in data, & we also calculate prob of that pt blue $\rightarrow 0.6$

- eg2: The second circle pt is actually black in data, but we finds its prob of blue ie 0.4, but likelihood says p of actual class it's actual class is black $\therefore P \Rightarrow 1 - 0.4 = 0.6$

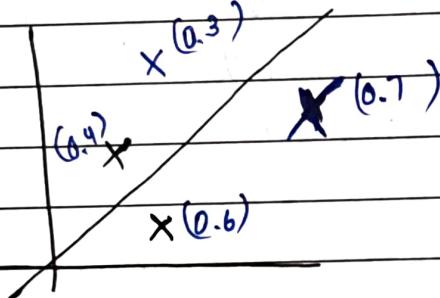
$$\therefore \text{Likelihood} = 0.6 \times 0.8 \times (1 - 0.4) \times (1 - 0.3)$$

$$= 0.6 \times 0.8 \times 0.6 \times 0.7$$

$$\boxed{\text{Likelihood 1} = 0.20}$$

Similarly we have line 2

Line 2:



$$= 0.3 \times 0.7 \times (1 - 0.4) \times (1 - 0.6)$$

$$= 0.3 \times 0.7 \times 0.6 \times 0.4$$

$$\boxed{\text{Likelihood 2} = 0.05}$$

- whosver likelihood is more, is better model
- It's actually logical, as if this number is big that means, har rang ka usi gong hone ka prob' is high. ie if prediction (prob) aligns with color, that means our prediction is right
- and if this no is less, that means aap blue ko black and black ko blue bula raha ho

- Maximum likelihood is a good parameter for judging because human has got no probability among grade hai go no hai hai

- Hence we have to find that like whose likelihood is maximum.

- Problem with using this approach

- In real data, there are lakhs of pts, which means lakhs of probability being multiplied which introduces the problem of underflow.

Solution

- We will use log of probabilities
- i.e. \Rightarrow we had $\Rightarrow 0.8 \times 0.6 \times 0.6 \times 0.7$
- Applying log $\Rightarrow \log(0.8 \times 0.6 \times 0.6 \times 0.7)$
 $= \log(0.8) + \log(0.6) + \log(0.6) + \log(0.7)$

- We solved the underflow problem
- But there is also a problem i.e. log of (0 to 1) is -ve no.
- To prevent from this, we will take -ve of log.

$$= -\log 0.8 - \log 0.6 - \log 0.6 - \log 0.7$$

by which we solved the -ve problem

- Also the product in likelihood we wanted was maximum

- but for logs, we want whose log is minimum. Ques is why?
 - because log of small no is more -ve and we reverse it using -ve of log.
- Eg. $\log(0.1) = -1$ $-\log(0.1) = +1$
 $\log(0.5) = -0.3$ $-\log(0.5) = +0.3$

- till now, we have:

$$-\log 0.8 - \log 0.8 - \log 0.6 - \log 0.7$$

- since there were 4 pts -

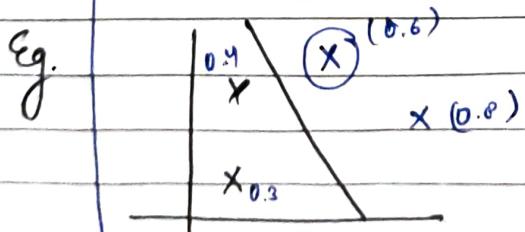
$$\Rightarrow -\log(\hat{y}_1) - \log(\hat{y}_2) - \log(\hat{y}_3) - \log(\hat{y}_4)$$

$$\Rightarrow \text{where } \hat{y}_i = \sigma(z), \text{ & } z = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

But, we directly can't write it as this, as you can see $\hat{y} = \sigma(z)$, which means it is the probability of pt being blue(1) (it only considered P(blue), which means it is directly putting of blue class actual class excluding)

- Therefore it. can be transformed as:-

$$= y_1 \log \hat{y}_1 - (1 - y_1) \log (1 - \hat{y}_1)$$



- $\Rightarrow \text{actual pt} = \text{blue}(1) = y_1$

- $1 \times (\log(0.6)) - (1 - 1) \log(1 - \hat{y}_1)$

$$= \frac{\log 0.6 - 0}{\log 0.6}$$

Basically this eqn: $-y_1 \log \hat{y}_1 - (1-y_1) \log (1-\hat{y}_1)$

it is the logic that is finding prob of actual class, in this we know this pt is blue; which means p of blue should be assigned

- similarly let's say actual data pt is black (0), i.e. $y_2 = 0$
 - It prevents (eq) from assigning directly the p of blue, but assigns the actual probability of class
- $$= -y_2 \log \hat{y}_2 - (1-y_2) \log (1-\hat{y}_2)$$
- $$= 0 \times (\log \hat{y}_2) - (1-0) \log (1-0.4)$$
- $$= 0 - \log(0.6)$$
- $$= -\log 0.6$$

That's how it prevents from assigning all prob of blue without knowing actual class

\therefore We had new eqn

$$= -y \log$$

∴ we transformed it into :-

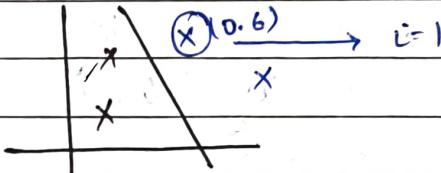
The output we wanted was :-

$$-\log 0.8 - \log 0.6 - \log 0.6 - \log 0.7$$

which we can achieve it by using :-

$$[-y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)]$$

for $i=1 \rightarrow$



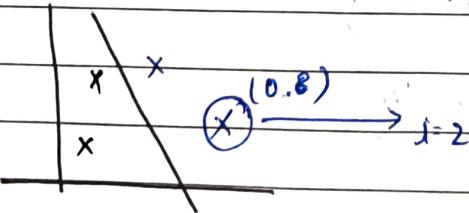
$$y_i = \text{Actual color} = \text{blue} = 1$$

$$\hat{y}_i = \sigma(z_i) = 0.6$$

put these in eqn

$$\begin{aligned}
 &= -1 \times \log 0.6 - (1-1) \log (1-0.6) \\
 &= -1 \times \log 0.6 = -\log 0.6 \\
 &= -\log 0.6
 \end{aligned}$$

for $i=2$



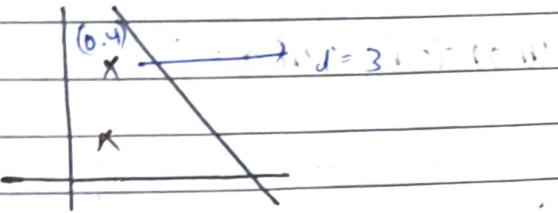
$$y_i = 1 \text{ (Actual } \Rightarrow \text{ blue)}$$

$$\hat{y}_i = \sigma(z_i) = 0.8$$

putting in eqn

$$\begin{aligned}
 &= -1 \log 0.8 - (1-1) \log (1-0.8) \\
 &= -\log 0.8
 \end{aligned}$$

for ($i=3$)



$$y_3 = 0 \text{ (Actual color = black)}$$

$$\hat{y}_3 = \sigma(z_3) = 0.4 \text{ (Prob. of blue)}$$

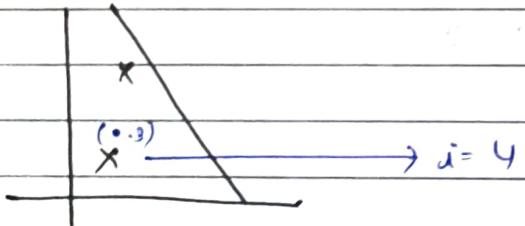
put in eqn

$$= -0 \times \log 0.4 - (1-0) \log (1-0.4)$$

$$= 0 - 1 \times \log (0.6)$$

$$= -\log 0.6$$

for $i=4$



$$y_4 = 0 \text{ (Actual color \(\rightarrow\) black)}$$

$$\hat{y}_4 = \sigma(z_4) = 0.3 \text{ (Prob. of this pt being blue)}$$

$$= -0 \times \log 0.3 - (1-0) \log (1-0.3)$$

$$= -\log 0.7$$

- That's how this eqn $[-y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)]$ works for all points

- Since we have n pts

$$= - \sum_{i=1}^n y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$$

If you want to find aug-i divide by 3

$$L \rightarrow -\frac{1}{n} \sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$

- This is the loss function of Logistic regression also called as log loss / Binary cross entropy.

- for this, we have to find $\beta_0, \beta_1, \beta_2$ (coef for line) for which our loss function is minimize.

- \Rightarrow where are $\beta_0, \beta_1, \beta_2 = ?$

Ans. $\hat{y}_i = \sigma(z_i)$

$$z_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$$

- We are at pt, when we have data, there can be multiple lines and each line have its set of coef ($\beta_0, \beta_1, \beta_2$)

- We want such coef (line) for which our loss function is minimum

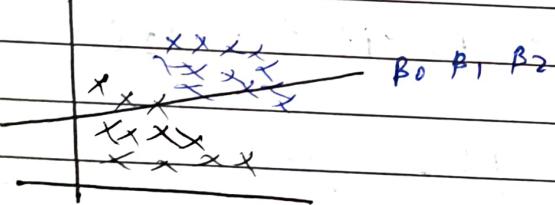
Now we will apply Gradient Descent

GRADIENT DESCENT

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

- Problem with this loss function
- There is no closed form solution
- We can't directly use some formula
- We can only approximate, which can be done using Gradient Descent
- How? Let's see

- Data →



- We start with random line, it will have $\beta_0, \beta_1, \beta_2$
- We will update these values

$$\beta_0 = \beta_0 - \eta \left[\frac{\partial L}{\partial \beta_0} \right] \rightarrow \begin{array}{l} \text{Differentiation of LF} \\ \text{w.r.t. } \beta_0 \end{array}$$

$$\beta_1 = \beta_1 - \eta \left[\frac{\partial L}{\partial \beta_1} \right]$$

$$\beta_2 = \beta_2 - \eta \left[\frac{\partial L}{\partial \beta_2} \right]$$

- finding $\frac{\partial L}{\partial \beta_0}, \frac{\partial L}{\partial \beta_1}, \frac{\partial L}{\partial \beta_2}$

$$L.f = -\frac{1}{n} \sum_{i=1}^n y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$$

- let's assume there is only 1 pt

$$\therefore = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

- we have to find $\frac{\partial L}{\partial \beta_1}$

$$\frac{\partial L}{\partial \beta_1} = \frac{\partial}{\partial \beta_1} [-y \log \hat{y} - (1-y) \log (1-\hat{y})]$$

$$\text{that pt is } \rightarrow \begin{bmatrix} x_1 & x_2 & y & \hat{y}_i \\ 28 & 13 & 0 & 0.63 \end{bmatrix}$$

$$= -y \frac{\partial \hat{y}}{\partial \beta_1} - (1-y) \frac{\partial (1-\hat{y})}{\partial \beta_1} \quad (\because \frac{\partial \log x}{\partial x} = \frac{1}{x})$$

$$= -y \left(\sigma(z) [1 - \sigma(z)] \right) \frac{\partial z}{\partial \beta_1} \quad (\text{Solving T term})$$

$$= -y \hat{y} (1-\hat{y}) x_1$$

1) $\hat{y} = \sigma(z)$
 $z = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

$$\frac{\partial z}{\partial \beta_1} = x_1$$

2) $\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1-\sigma(x))$

Solving 2nd term

$$\begin{aligned}
 &= -\frac{(1-y)}{(1-\hat{y})} \frac{\partial}{\partial \beta_1} (\hat{y} - y) \\
 &= -\frac{(1-y)}{(1-\hat{y})} \frac{\partial}{\partial \beta_1} ((1-\sigma(z))) \\
 &= -(1- \\
 &= \frac{(1-y)}{(1-\hat{y})} \hat{y}(1-\hat{y}) x_1 \\
 &= (1-y)\hat{y} x_1
 \end{aligned}$$

Combining both underline of

$$\frac{\partial L}{\partial \beta_1} = -y(1-\hat{y})x_1 + (1-y)\hat{y}x_1$$

$$\frac{\partial L}{\partial \beta_1} = (\hat{y} - y)x_1$$

for n pts

$$\boxed{\frac{\partial L}{\partial \beta_1} = (\hat{y}_i - y_i)x_{1,i}} \quad \textcircled{1}$$

Similarly

$$\boxed{\frac{\partial L}{\partial \beta_2} = (\hat{y}_i - y_i)x_{2,i}} \quad \textcircled{2}$$

$$\boxed{\frac{\partial L}{\partial \beta_0} = \hat{y} - y} \quad \textcircled{3}$$

Apply in code

Summary

Step 1 - Sigmoid

Step 2 - maximum likelihood

Step 3 - log \rightarrow Binary Cross Entropy

Step 4: Gradient Descent $\rightarrow \beta_0, \beta_1, \beta_2$
which minimizes L.