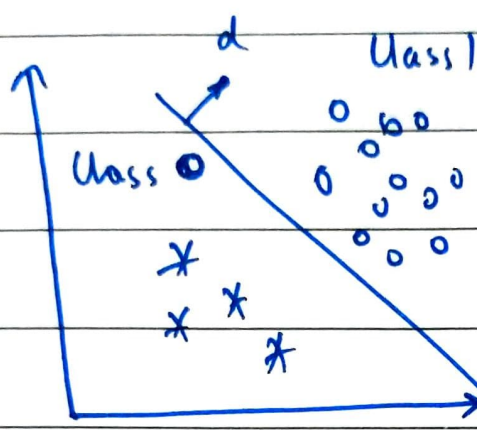


Logistic Regression

* It is a classification technique



This line is
Hypothesis.

Or, hypothesis will look like a plane in higher dimensions.

Eqⁿ of line :

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$$

(This is of form $ax + by + c = 0$)

Now, we are given a point d .

We calculate its distance from hypothesis.

If, $d = \theta_1 x_1 + \theta_2 x_2 + \theta_0 \geq 0$ then point lies Rhs of line & class is 1

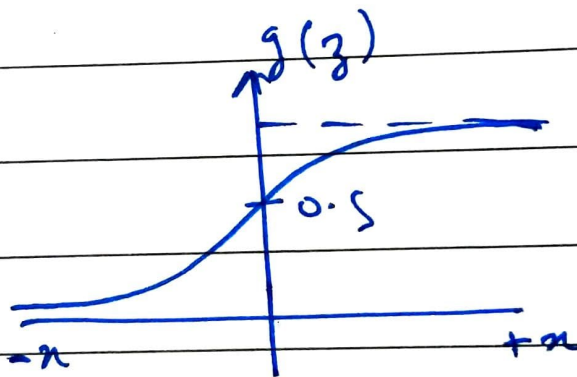
If $d < 0$ the class is 0.

Higher distance means more confidence.
 d can be from $-\infty$ to ∞ , we want to map it in range $[0, 1]$. This will be our probability.

← -ve side of line \rightarrow +ve side of line
 $-\infty$ ∞

Sigmoid function:

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

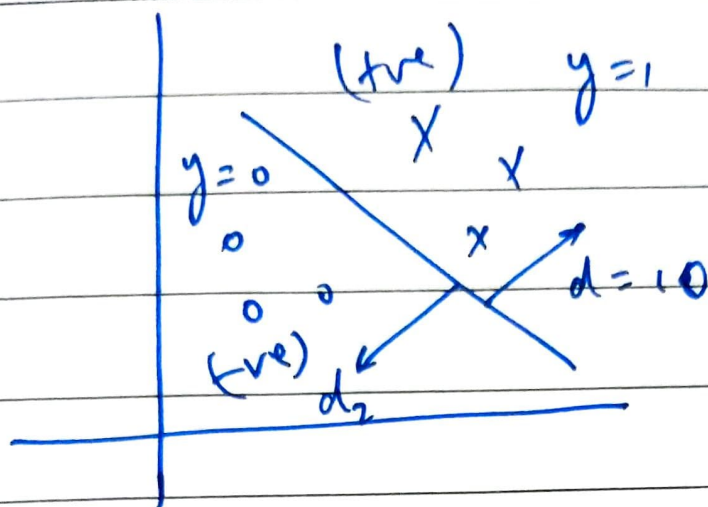


$$g(z) \in (0, 1)$$

Now, hypothesis is going to be of form:

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

where, $d = \theta^T x$ // distance



$$g(z) = \frac{1}{1 + e^{-3}} = 0.8 \text{ (true)}$$

Prob ($Y=1/x$) = 0.8 Belonging to class 1

Prob ($Y=0/x$) = 0.2 Belonging to class 2

* As we move away from the line, the probability that point lies in that class increases.

d_2 will be negative.

$$g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} = 0.2 \text{ (This will be true now)}$$

\Downarrow
-ve

$P(Y=1/x) = 0.2$ // Important

$P(Y=0/x) = 0.8$ // Point belongs here.

By assuming, it will create a model like that
Now, (here we assume that for $y=0$, distance is -ve)

$$P(Y=1 | x, \theta) = h_{\theta}(x)$$

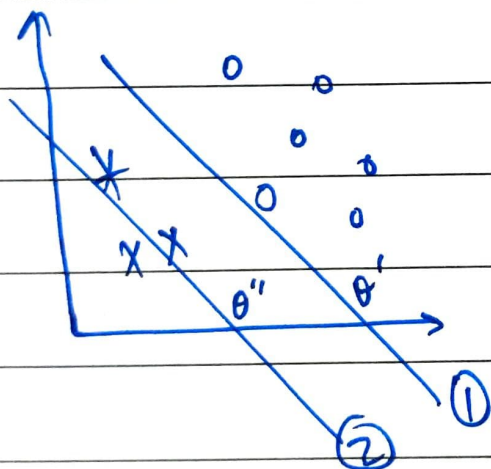
$$P(Y=0 | x, \theta) = 1 - h_{\theta}(x) \rightarrow y_a$$

$$P(Y/x, \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

which is the best line?

[Here y will be 0 or 1 only!]

Concept: Max likely hood estimation.



(Of c line ① has more likely hood)

$$P(Y/x, \theta') > P(Y/x, \theta'')$$

Line ① has higher likely hood than ②

(Basically means all pts. in dataset)

* Each pt. is independent

$P(y_i | x_i) \rightarrow$ For i^{th} item, the label y_i depends on the features of the pt.

Say you have 2 types of fruits apple & mango in a basket. If you pick apple you get juice, if mango then shake.

$P(\text{Juice} | \text{Fruit})$
 $y=1$

$P(\text{Shake} | \text{Fruit})$
 $y=0$

Here, what drink you will get depends on fruit you pick. Independent probabilities.

$P(\text{Both Juice \& Mango}) = \text{Multiply Probabilities}$

Likelihood = $\prod P(y_i | x_i, \theta)$ // Product of all Prob.

We maximize this product.

Because we want product of all probabilities in training set should be maximised.

This is $P(y/x, \theta)$ (on previous page)

* We need to find line such that product of all probabilities (-ve & +ve side both) is maximum.

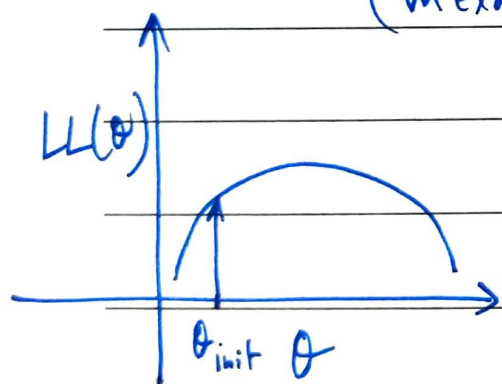
* This is called maximize likelihood

* After that, we take log likelihood.
LL after taking log on both sides of $p(y|x, \theta)$.

$$LL(\theta) = \sum_{i=1}^m y_i \log h_{\theta}(x_i) + (1-y_i) \log (1-h_{\theta}(x_i))$$

(m examples)

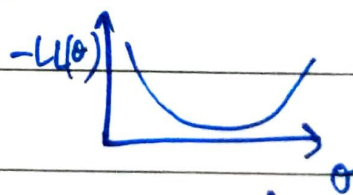
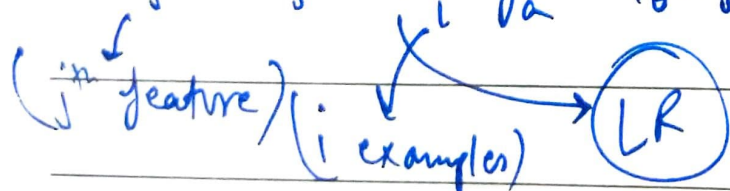
AKA cross entropy



(\Rightarrow Gradient Ascent)

$$\theta = \theta + \eta LL'(\theta)$$

$$\theta_j = \theta_j + \eta \sum (y_i - h_{\theta}(x_i)) x_i^j$$



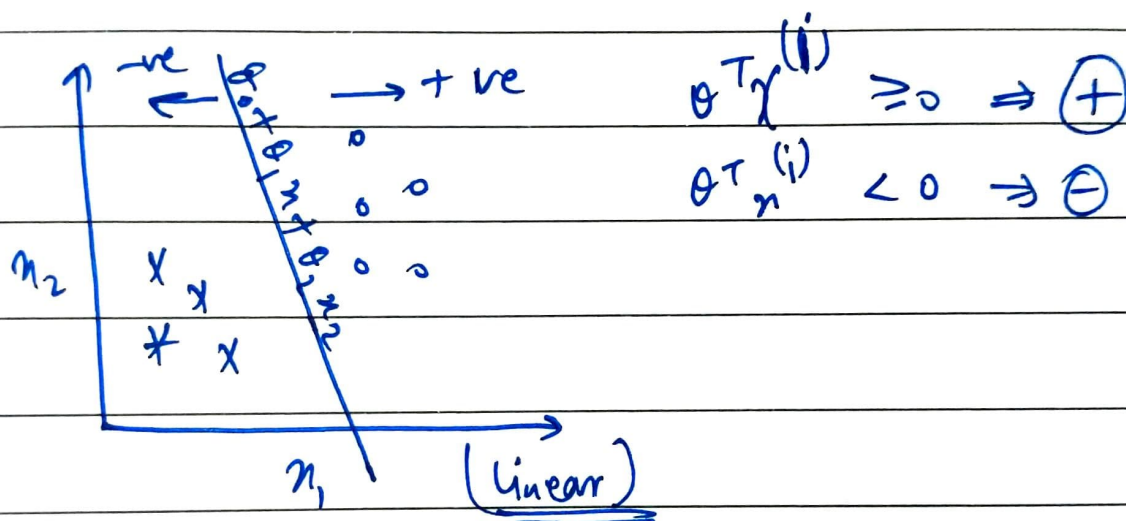
(OR, loss = $-LL(\theta)$ & minimize this.)

Prediction:

$$h_{\theta}(n) = g(\theta^T n) \quad \therefore \text{Prob } * \text{ belongs to class 1}$$

$$g(\theta^T n) \geq 0.5 \Rightarrow \text{Class 1}$$

$$< 0.5 \Rightarrow \text{Class 0}$$

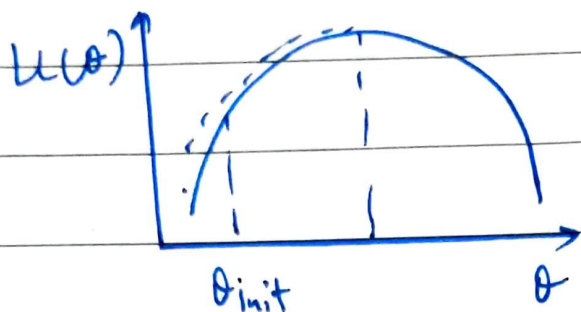
Find, optimal [best set of parameters] θ

$g(\theta^T n) \rightarrow (0, 1)$, it represents prob.

update rule:

$$\theta = \theta + \eta \nabla_{\theta} LL(\theta)$$

learning rate



At local maxima, gradient will become 0.

$$\nabla_{\theta} L(\theta) = \frac{\partial (L(\theta))}{\partial \theta}$$

one property, $g(z) = \frac{1}{1+e^{-z}}$

$$g'(z) = g(z)(1-g(z))$$

$$\frac{\partial (L(\theta))}{\partial \theta_j} = (y - h_{\theta}(x)) x_j \quad // \text{Identical to LR}$$

$$\theta_j = \theta_j + \eta \sum_{i=1}^m (y_i - h_{\theta}(x^{(i)})) x_j^{(i)}$$

// This is our update rule to reach local maxima.

(j^{th} feature of i^{th} example)

So, we are summing up gradients over all examples for all the features.

In the end, our $x\theta$ is made in such a way that for a point $\in y=0$, distance is -ve & hypothesis $h(x)$ gives $P(y=1/x)$.
