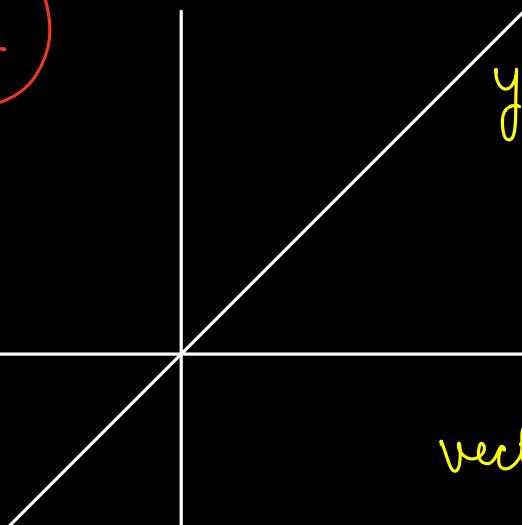


Linear equation in Matrix form

I



$$y = mx + c \Rightarrow ax + by = c$$

$$w_1x + w_2y = w_0$$

OR we can write it in the vector form.

$$X = \begin{bmatrix} 1 \\ x \\ y \end{bmatrix} \text{ or } \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

$$\therefore \boxed{W^T X = 0}$$

What is physical significance of m or Slope or $\frac{dy}{dx}$ or Derivative

$$W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

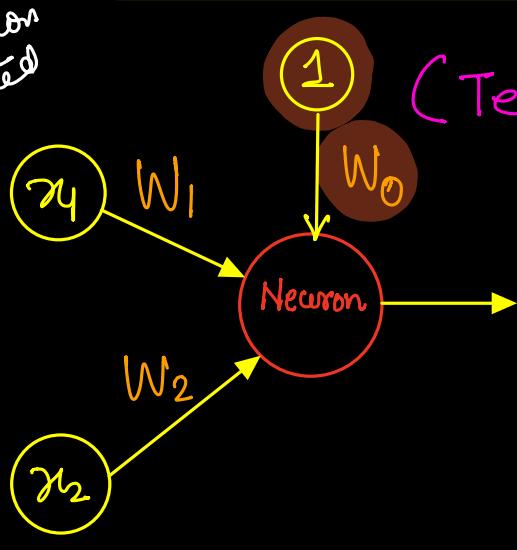
Given a line say $(y = mx + c)$ & a point (x_1, y_1)

How to know its side.

$$y_1 - m x_1 - c \geq 0$$

$$\text{OR } \boxed{y_1 - m x_1 - c < 0}$$

This computation can be abstracted as a neuron.



(Termed as bias)

This neuron just works as an aggregator and do 2 things

- Compute $w_0 + w_1 x_1 + w_2 x_2$
- fires if it is ≥ 0

Derivative

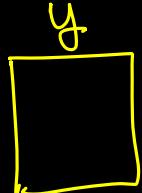
$$m = \frac{dy}{dx}$$

which is rate of change of y wrt x

↓↓

Just thinks x as a knob and we

How much is Δy when $\Delta x = 1$



need to know how much it can effect y .

its sensitivity wrt to the value of y

III

What is a function?

It is just a mapping b/w input & o/p.

$$y = f(x) \quad (1D \text{ function})$$

Plotted in 2D space

Can only be differentiated in one direction

$$\left[\frac{dy}{dx} \right]$$

(Multivariate
Calculus)

$$y = f(x_1, x_2)$$

Plotted in 3D space.

Can be differentiated in 2 directions

$$\left[\frac{\partial y}{\partial x_1} \text{ & } \frac{\partial y}{\partial x_2} \right]$$

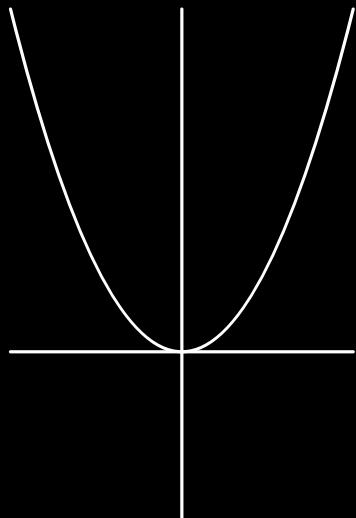
If differentiated wrt vector $x = [x_1 \ x_2]$

But why derivatives are essentially required.

IV

Given any function $f(x)$ how can you maximize or minimize. [Optimization]

Say $y = f(x) = x^2$ (Quadratic fn)



$$\frac{dy}{dx} = 2x \quad \therefore 2x = 0 \\ \Rightarrow x = 0$$

where we can get extremes.

If $y = (x+1)^2$ $\frac{dy}{dx} = 2(x+1) = 0$
 $(x = -1)$

(How to know maxima OR minima)

$$ax + by + cz = d \quad (\text{This is a plane in 3D})$$

$$w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 = 0 \quad (\text{What is this})$$

$$\sum_{i=0}^n w_i x_i = 0 \quad (\text{what is this?}) \quad (x_0 = 1)$$

(Hyper-plane)

Multidimensional
linear function.

Given a multidimensional fn say (F) defined in \mathbb{R}^N

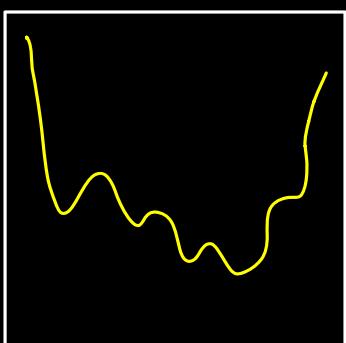
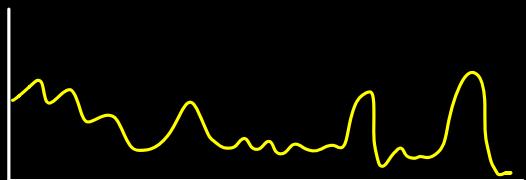
How can we minimize it. This (F) can be a loss function or objective function. (Closed-form)

In some (N) dimensional (say $w_0, w_1, w_2, \dots, w_{N-1}$) parametric space assume that some loss function (F) is defined. We require :

$$w^* = \underset{W}{\operatorname{arg\min}} (F)$$

Just differentiate F wst $w \Rightarrow \frac{dF}{dw}$

This can be done easily if (F) can be represented as some close form formula (as discussed above)



But most of the functions are not well defined as some formula and are represented in a discrete manner.

$$f(x_1), f(x_2), f(x_3), \dots$$

OR $f(x_1, x_2, x_3)$

$$x_1, x_2, x_3$$

Discrete functions

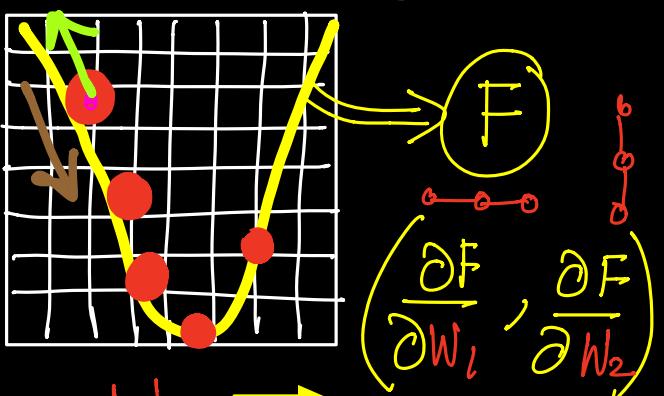
Now in such a scenario how can we optimize a an objective function.

In discrete space we have an iterative solution.

$$W_1^{\text{New}} = W_1^{\text{Old}} - \eta \frac{\partial F}{\partial W_1}$$

$$W_2^{\text{New}} = W_2^{\text{Old}} - \eta \frac{\partial F}{\partial W_2}$$

Gradient decent algorithm for optimization (used as optimizer)



Parametric Search for optimization. [Update till Convergence]

SGD, ADAM, AdaGrad etc.c.

Analytically Computing gradient :

$$\lim_{h \rightarrow 0} \frac{f(x_i + h) - f(x_i)}{h} \quad \boxed{\begin{matrix} A \\ \begin{bmatrix} 0 & 0 & 0 \\ x_{i-1} & x_i & x_{i+1} \end{bmatrix} \end{matrix}}$$

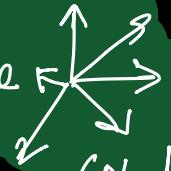
$$\boxed{\begin{matrix} B \\ \begin{bmatrix} \begin{matrix} x_i & x_i & x_i \\ i-1, j-1 & i-1, j & i-1, j+1 \end{matrix} \\ \begin{matrix} x_i & x_i & x_i \\ i, j-1 & i, j & i, j+1 \end{matrix} \\ \begin{matrix} x_i & x_i & x_i \\ i+1, j-1 & i+1, j & i+1, j+1 \end{matrix} \end{bmatrix}} \quad \text{for } (x_i), \text{ gradient in horizontal direction can be : } x_{i+1} - x_{i-1}$$

for $(x_{i,j})$ gradient in horizontal : $x_{i,j+1} - x_{i,j-1}$ OR averaging all 3 Rows

This is gradient w.r.t a neighbourhood.

VII

Hence in that high dimensional space



(N-dim)

- ✳ These are (N) directions to move on.
- ✳ we are in search of optimal set of parameters (W^*) that can minimize the objective fn (F).
- ✳ We need to use Gradient decent iterative algorithm to optimize over (F) in \mathbb{R}^N .

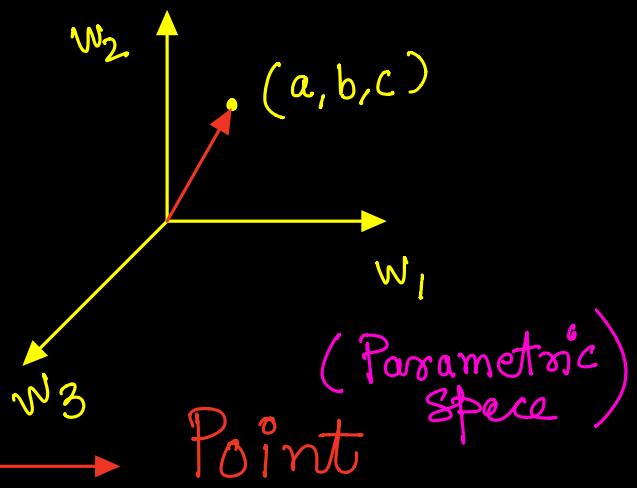
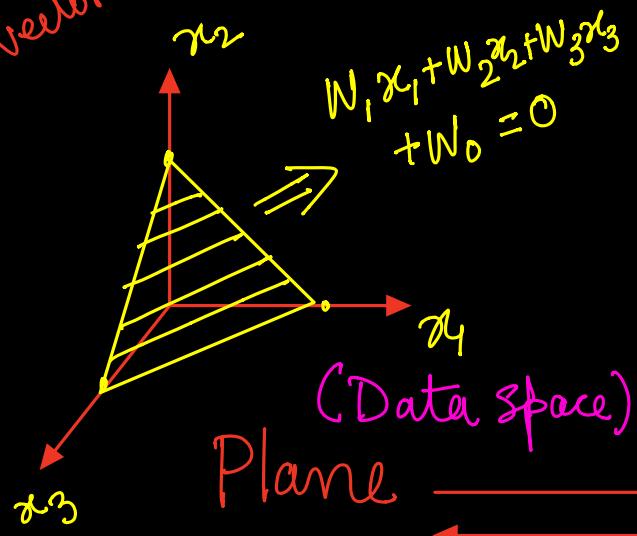
→ Initialize (W) vector $\in \mathbb{R}^N$, randomly.
 → Update $\forall W_i^* = W_i^{\text{old}} + \Delta W_i$

$$\text{where } \left\{ \Delta W_i = -\eta \frac{\partial F}{\partial W_i} \right\}$$

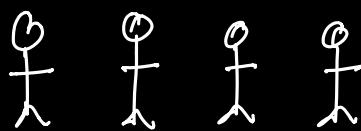
$WX = 0$
 A h-plane in
 vector form is \mathbb{R}^N

Each $[W]$ vector $\in \mathbb{R}^N$, and can be
 seen as a h-plane in the space of (X) .

\mathbb{R}^N



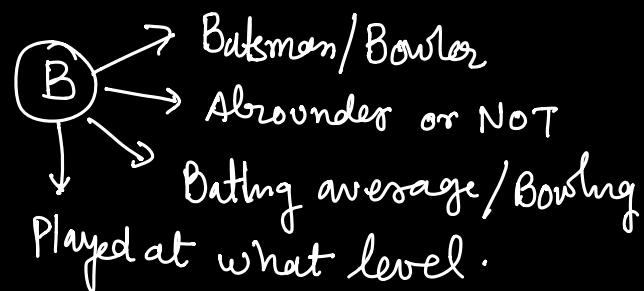
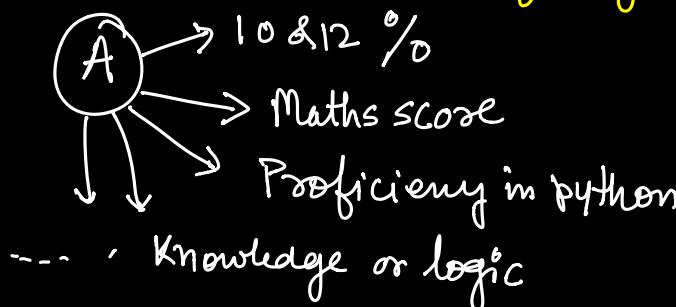
Problem of binary classification



(1000 persons have applied
for
 (a) Join by research group
 (b) Join my cricket team.

Name \otimes (How many attributes to be considered)?
Mobile
Qualification as well as what?
Gender

- Central question
- \otimes During the interview like (father's Name, Village, Uncle, family, friend, 3D gesture, face structure - - - - -)
 what questions we need to ask.
 \otimes What is the weightage of that question.

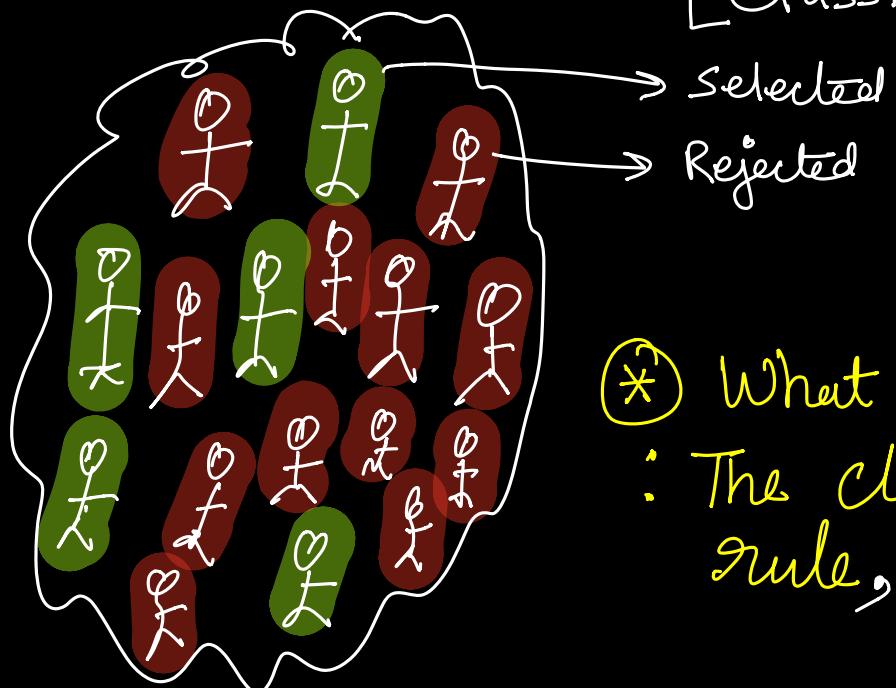


feature selection \otimes One can observe that the questions that we need to ask are depending upon the objective or Task in hand (Problem Statement)

Question and their importance keep on changing with respect to the problem in hand [Feature Engg]

Now as an interviewer, I will ask questions to all the applicants. Get their answers (features) and accept/reject (Based on some Criterion)

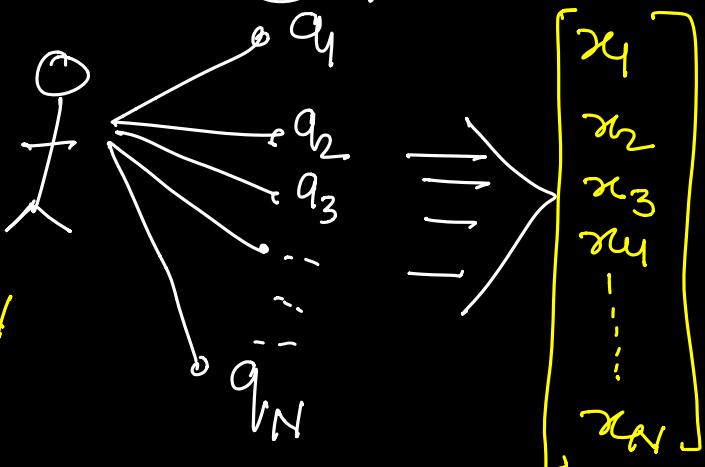
[Classification decision]



(Out of 1000
100 Selected &
900 got rejected)

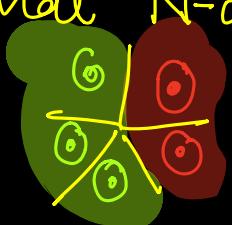
* What we need to learn:
: The classification decision rule, given the questions.

* Assuming we ask N questions to each applicant:



* Basically each data got projected onto \mathbb{R}^N space: data to an N-dim vector

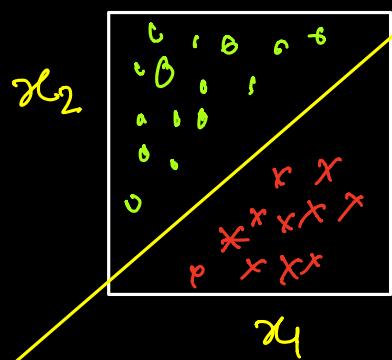
* In that N-d space (\mathbb{R}^N) it will be seen as a point.



Typical (N) values are 128/256 & much bigger.

* We need to ask relevant questions, so as Feature extraction to project data points into (\mathbb{R}^N) in such a way that in that space, accepted & rejected Candidates are well apart.

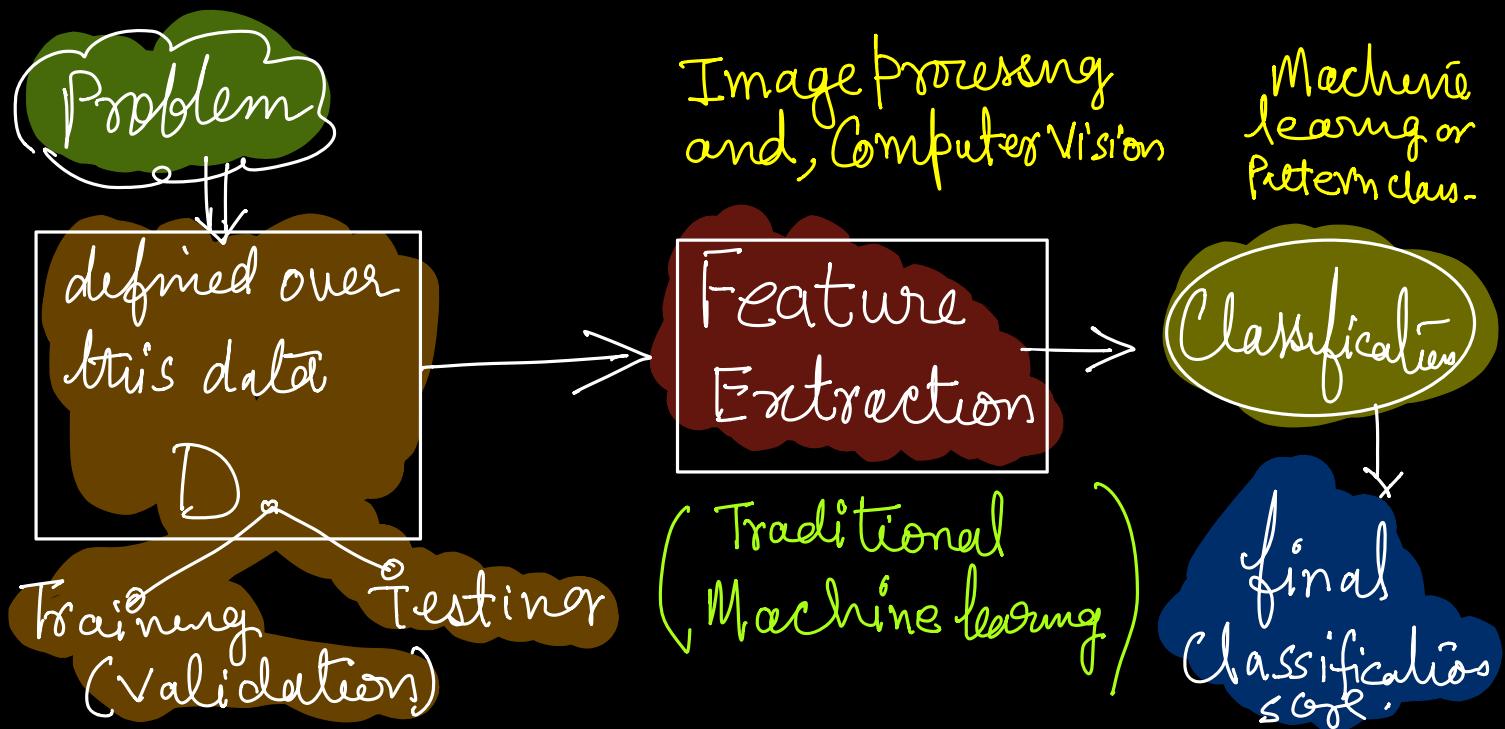
* Assuming we have asked only 2 questions
→ data got projected in 2D space.



Classifier need to learn this Separating line or H-plane Called Decision boundary.

ML to DL

The full flow of a classification system can be





Similar to the above mentioned functional Optimization, we need to search for the best h-plane in the feature space $X = \{x_1, x_2, \dots, x_N\} \in \mathbb{R}^N$

which will be

$$w_0, w_1, w_2, \dots, w_N \left. \begin{array}{l} \text{Coefficients} \\ \text{of the h-plane.} \end{array} \right\}$$

Bias

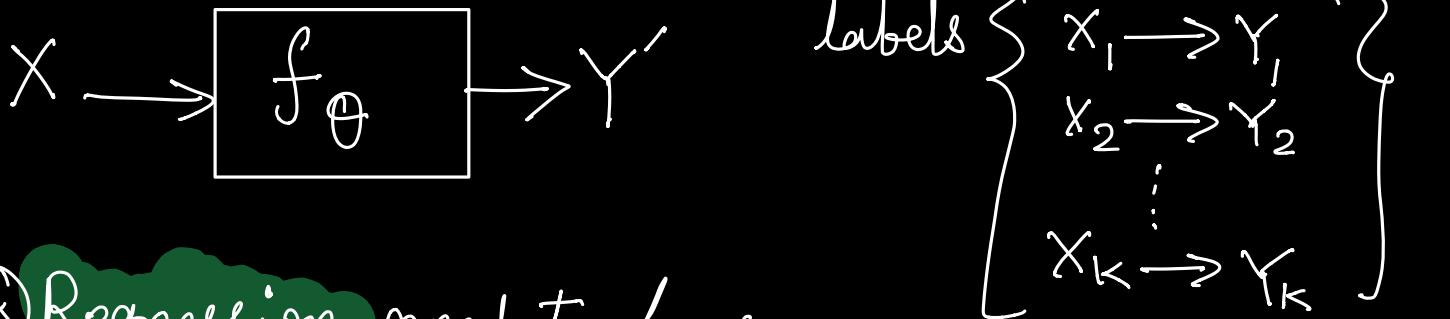
But what is going to be the function that we need to optimize \Rightarrow loss fn

- MSE or MAE (for Regression)
- Cross Entropy (Classification)

Supervised Learning:

These loss fn, for each given i/p data (X) assumes a correct o/p (Y) say ground truth. Using the learned logic, (X) will be converted into the predicted o/p (Y'). Loss functions need to estimate the dissimilarity/deviation b/w Y & Y' $\Rightarrow \{LF(Y, Y')\}$

In case of regression Y & Y' may be a single value
for classification Y & $Y' \Rightarrow$ Probability scores for all the classes.



(A) Regression need to learn

(θ) (Set of parameters) such that

$f_{\theta}(X_i) = Y'_i$ need to be as close as possible to Y_i . Loss fns can be :

$$L_2 \text{ Norm (a)} \quad MSE(Y_i, Y'_i) = \frac{1}{2} (Y_i - Y'_i)^2$$

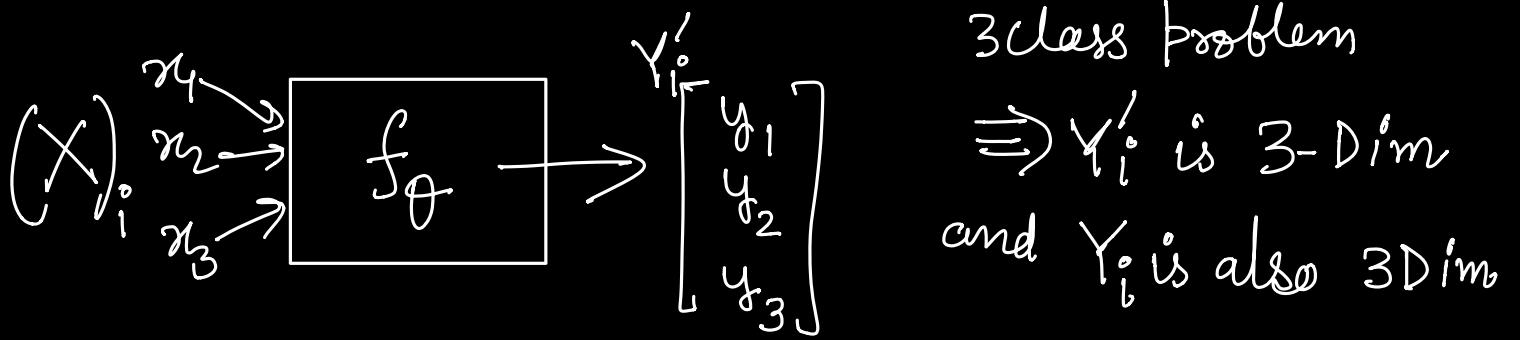
$$L_1 \text{ Norm (b)} \quad MAE(Y_i, Y'_i) = |Y_i - Y'_i|$$

(B) Classification : Y_i 's are represented as 1-Hot encoding.

if 10 class classification then

each $Y_i \Rightarrow$ 10-dim (binary vector)

$$\begin{bmatrix} \text{Cat} \\ \text{dog} \\ \text{car} \\ \vdots \\ \vdots \end{bmatrix} \Rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow \begin{array}{l} \text{implies } Y_i \text{ belongs to} \\ \text{Cat (only one place)} \\ \text{is set} \end{array}$$



$\forall i$ X_i we have their corresponding labels (Y_i) (3D)

Loss fn : Cross entropy (Y_i, Y'_i)

$$\Rightarrow Y_i = \begin{bmatrix} y_1 & y_2 & y_3 \end{bmatrix} \Rightarrow \sum_{i=1}^3 -y_{i\circ} \log(Y'_{i\circ})$$

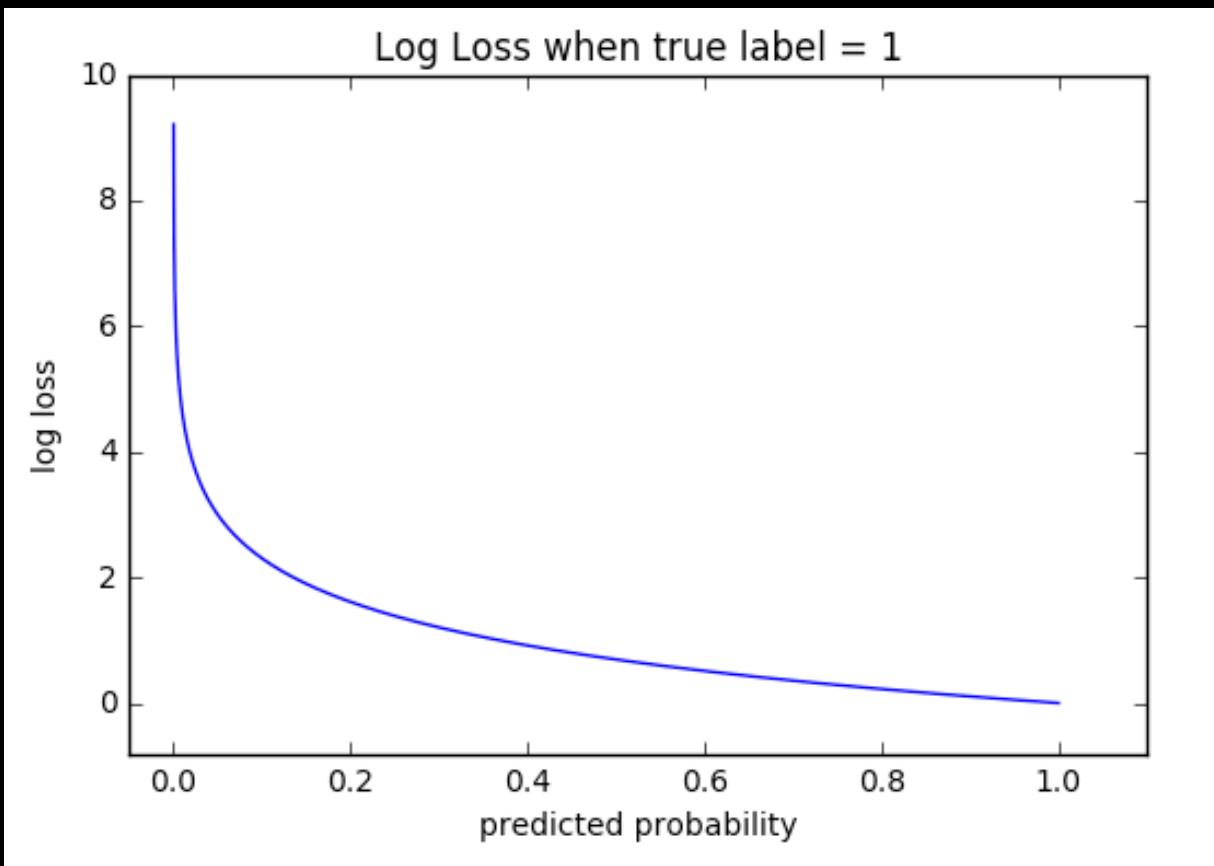
y P

 $1-y$ $1-P$

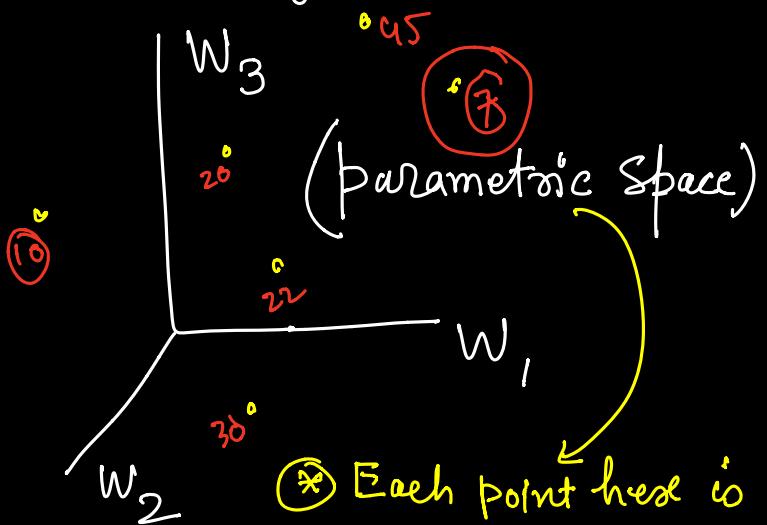
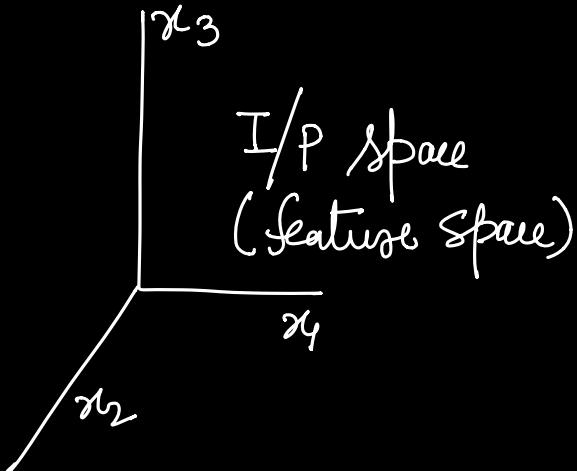
\Rightarrow Binary cross entropy

$$-(y \log P + (1-y) \log(1-P))$$

This is $\Rightarrow - (y_1 \log(Y'_1) + y_2 \log(Y'_2) + y_3 \log(Y'_3))$



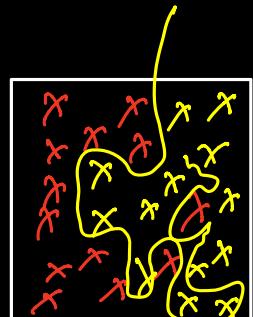
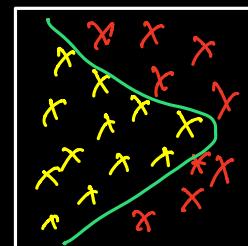
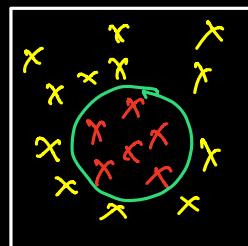
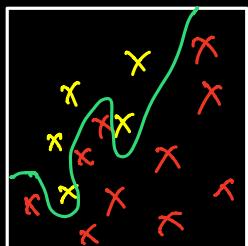
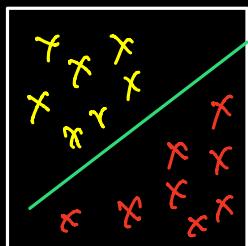
* So, now we are having 2 spaces



* we need to explore the parametric space in order to optimize the Cost fn efficiently

(This is done by applying GD or SGD)

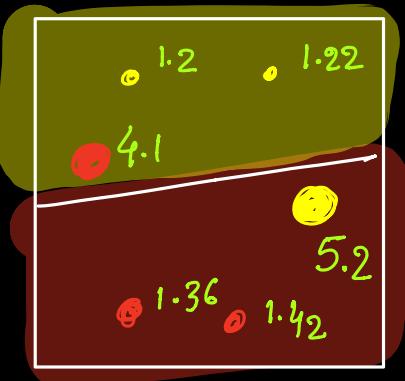
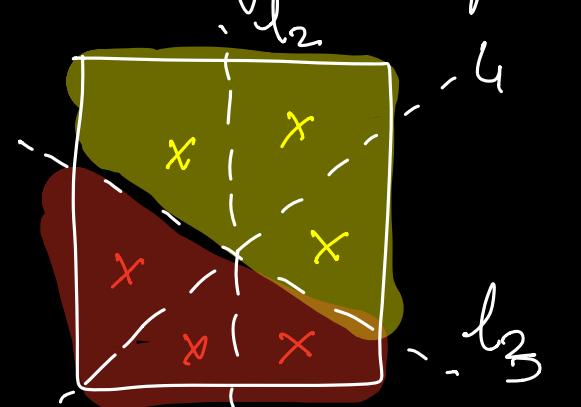
* One neuron can be used to learn one such h-plane if our data is linearly separable. But mostly data is not linearly separable. Then how to learn non-linear boundaries.



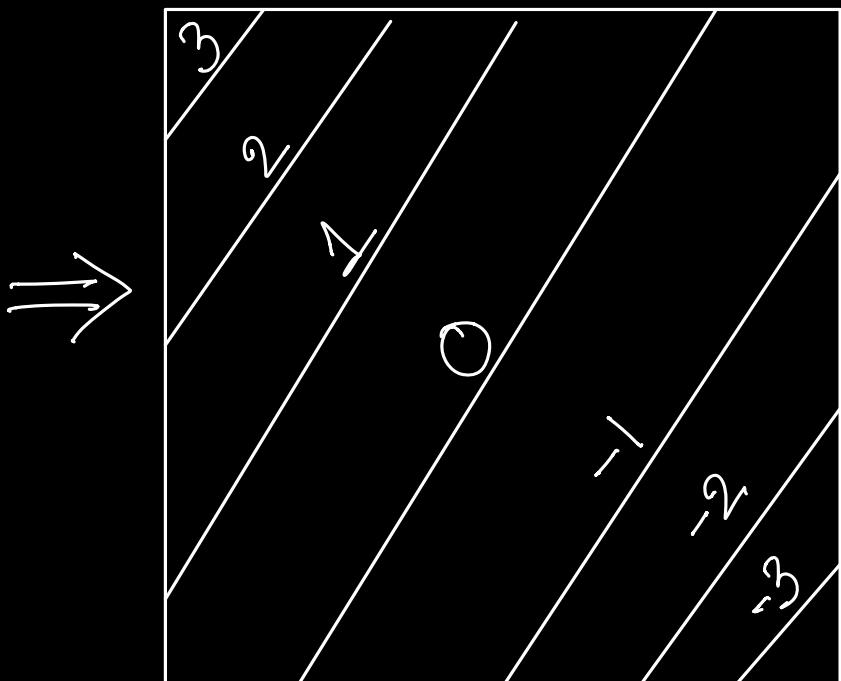
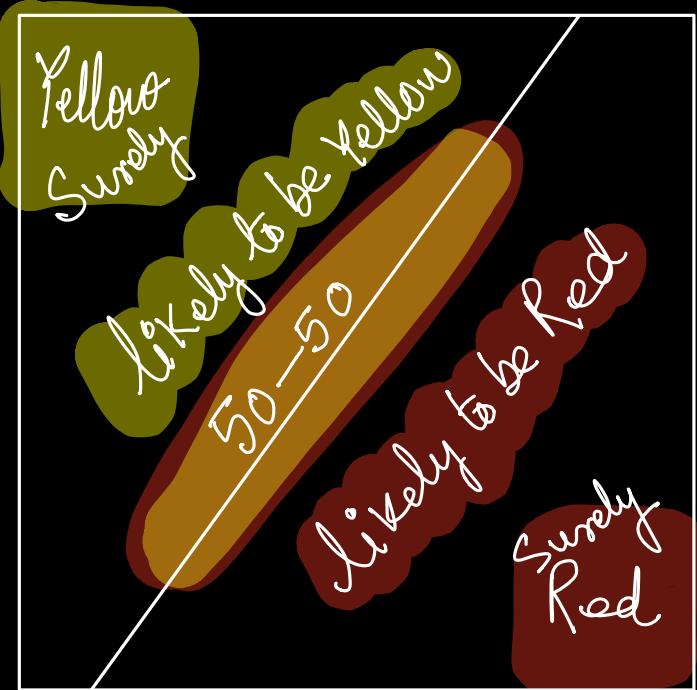
We require more neurons, may be in multiple layers.

Intuition of Neural Networks

- ✳) Neural networks can be giants and may have millions of neurons & thousands of layers.
- ✳) Essentially they are partitioning the feature space with an objective of minimizing misclassification.
Basically Splitting the data
 - one can start from some very bad h-plane.
 - Need to update it parameter so as to minimize the objective fn (mis-classification)
 - We can utilize SGD & GD iteratively.
But there is an issue \Rightarrow Loss fn need to be differentiable
(Here it is discrete)
- ✳) We can use logistic loss fn ?
 - Each and every data is Penalized. Correctly classified then less penalty & mis-classified then more penalty
 - Error $\Rightarrow 1.2 + 1.22 + 4.1 + 5.2 + 1.36 + 1.42$ [continuous]

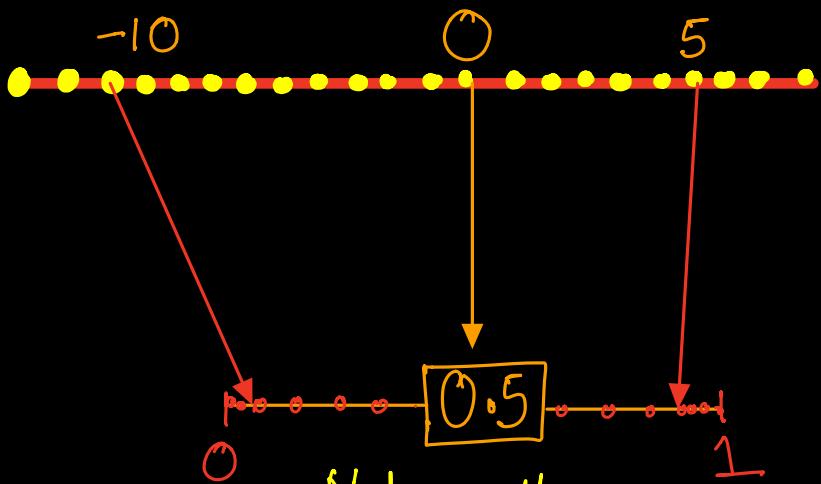


Now the question is how to get such loss fn?



It can be done using
Sigmoid fn :

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

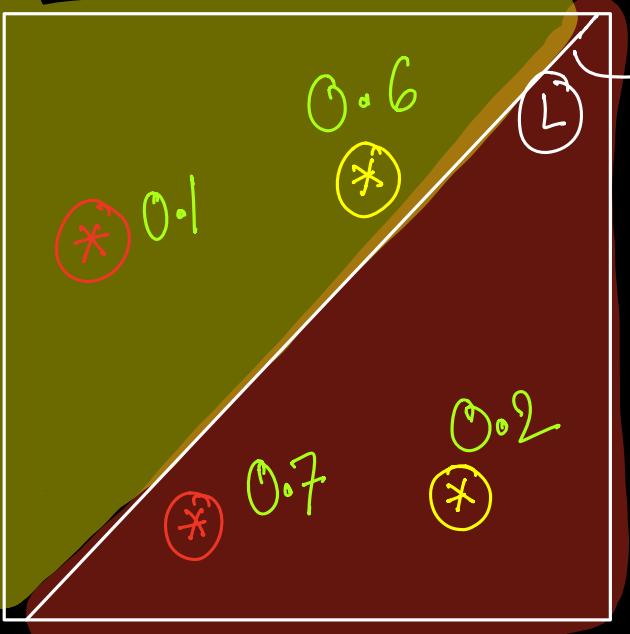


Not equally spaced \Rightarrow Basically a non-linear mapping



They are called Activation function, (Nonlinear)
like, Sigmoid, Tanh, ReLU, LReLU, PReLU - - - .

Let us consider this Example.

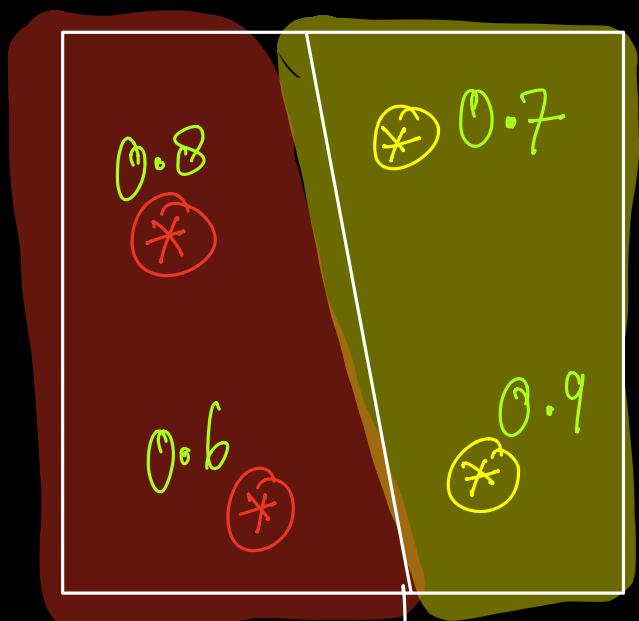


→ Decision boundary (L) separating Yellow & Red Region.

- 2 Classes Yellow / Red
- Shown values are the probability that they belong to the correct class wrt given decision boundary (L).

∴ Joint probability that all points got correctly classified by (L) : $0.6 \times 0.2 \times 0.1 \times 0.7 \Rightarrow 0.0084$
(This can be termed as likelihood.)

∴ Likelihood can be used as a gain fn & need to be maximized.



for above shown boundary

$$0.6 \times 0.2 \times 0.1 \times 0.7$$

$$= 0.0084$$

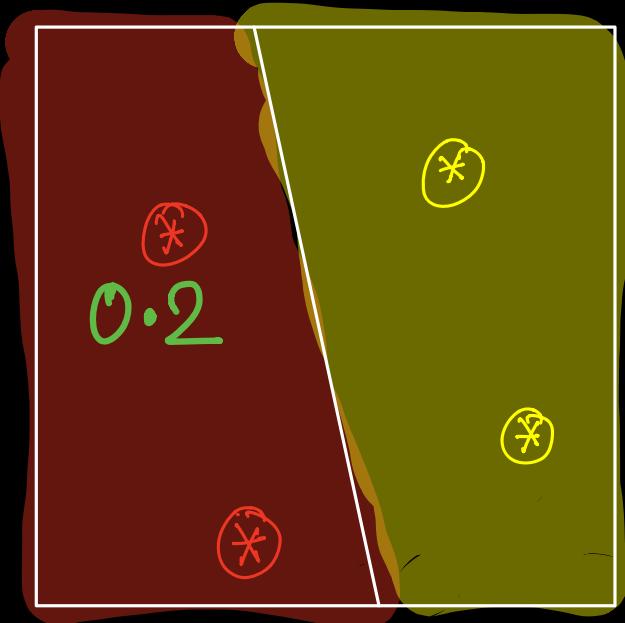
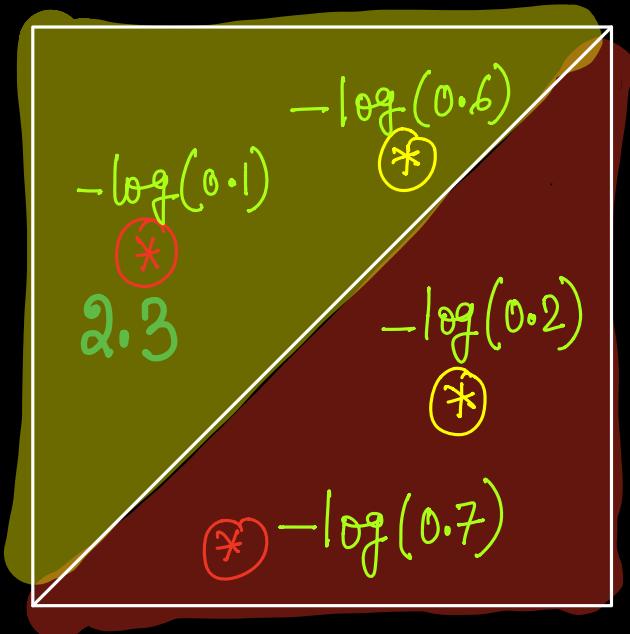
for L'

(maximized)

$$0.7 \times 0.9 \times 0.8 \times 0.6 = 0.3024$$

- These issue is it is multiplication
- The values are very small.

Multiplication → Addition & value scaling.



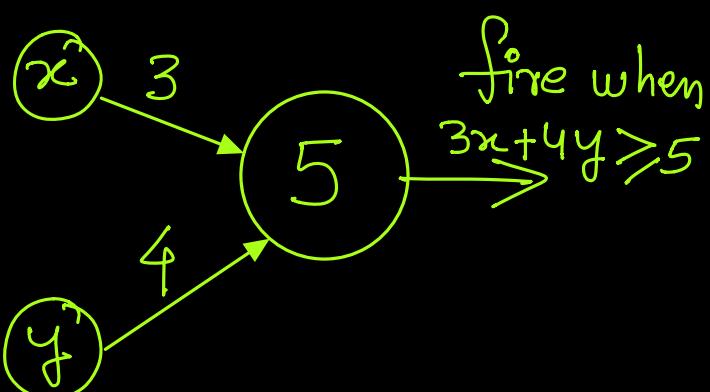
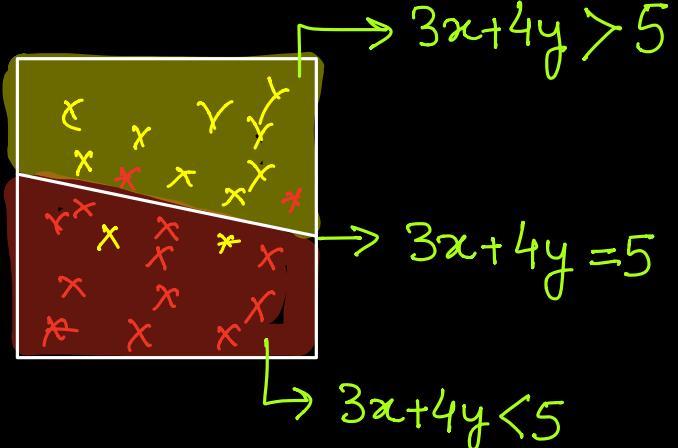
$$0.6 * 0.2 * 0.1 * 0.7 \\ = (0.0084)$$

$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024 \\ \text{(Maximization of likelihood)}$$

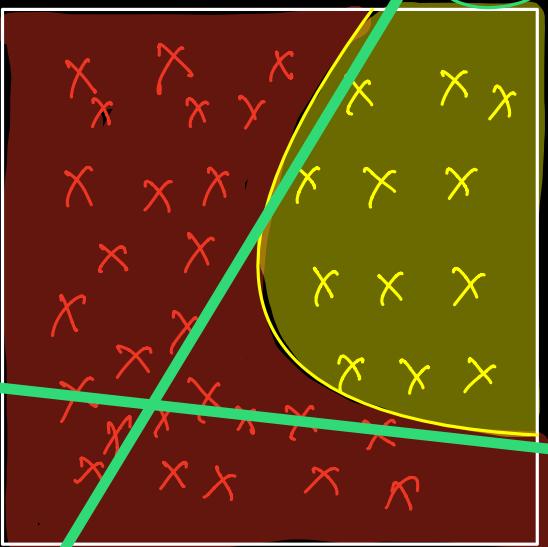
$$-\log(0.6) - \log(0.2) - \log(0.1) \\ - \log(0.7) = 4.8 \\ \downarrow \\ (2.3)$$

$$-\log(0.7) - \log(0.9) \\ - \log(0.8) - \log(0.6) \\ \downarrow \\ (0.2)$$

Log likelihood minimization. \Rightarrow Optimize & get the linear boundary.



Single neuron can learn decision boundary if data is linearly separable.

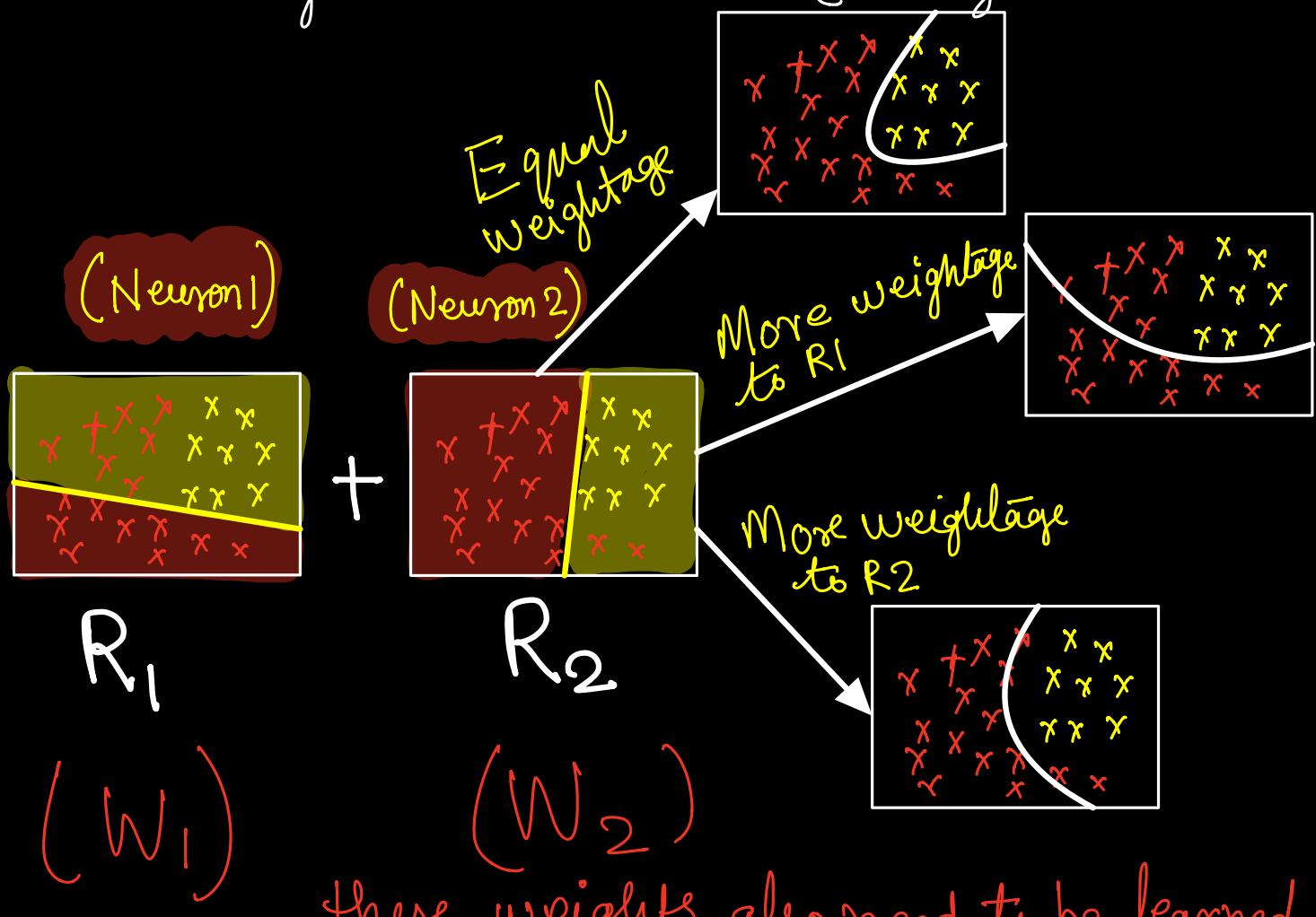


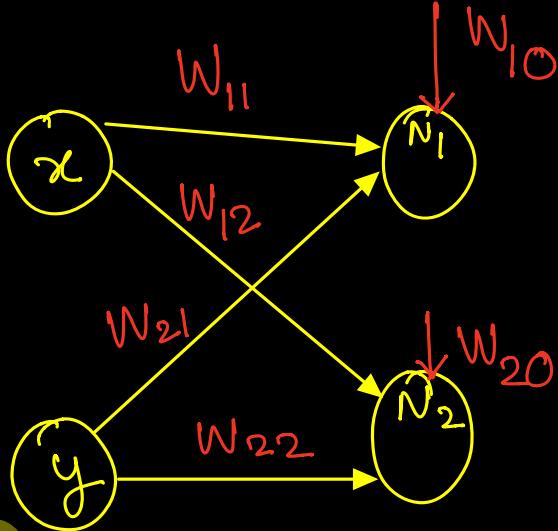
(L₁)

How to learn the non linearity?

⇒ We need to learn the non-linear boundaries in a modular fashion & hierarchically, (using a layered N/w)

First we need to understand Region Combination (weighted)

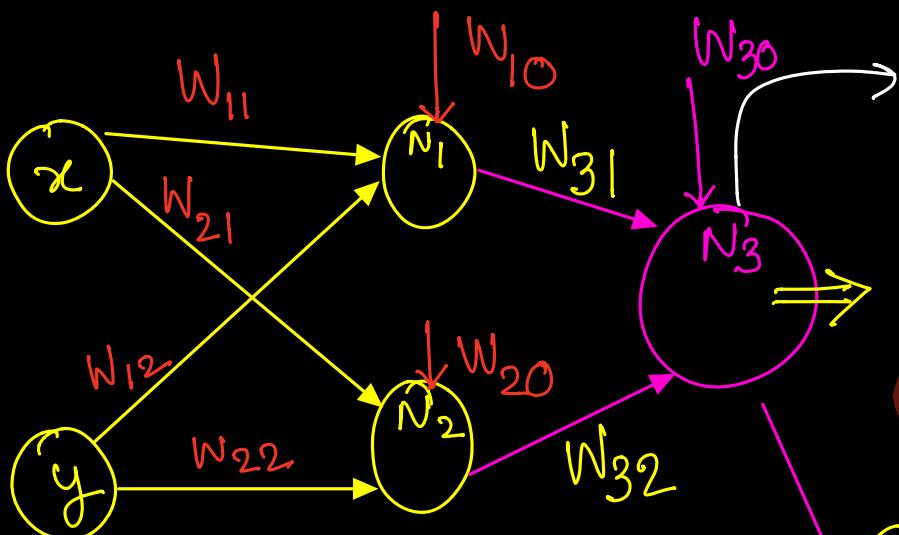




→ This neuron is learning one h-plane (learning Region1)

→ This neuron is learning another h-plane. (learning Region2)

$w_{ij} \Rightarrow i^{\text{th}}$ neuron, j^{th} input



This neuron is learning how to merge 2 regions, (Basically 2 decision merging)

④ (N_3) not seeing the actual IP.

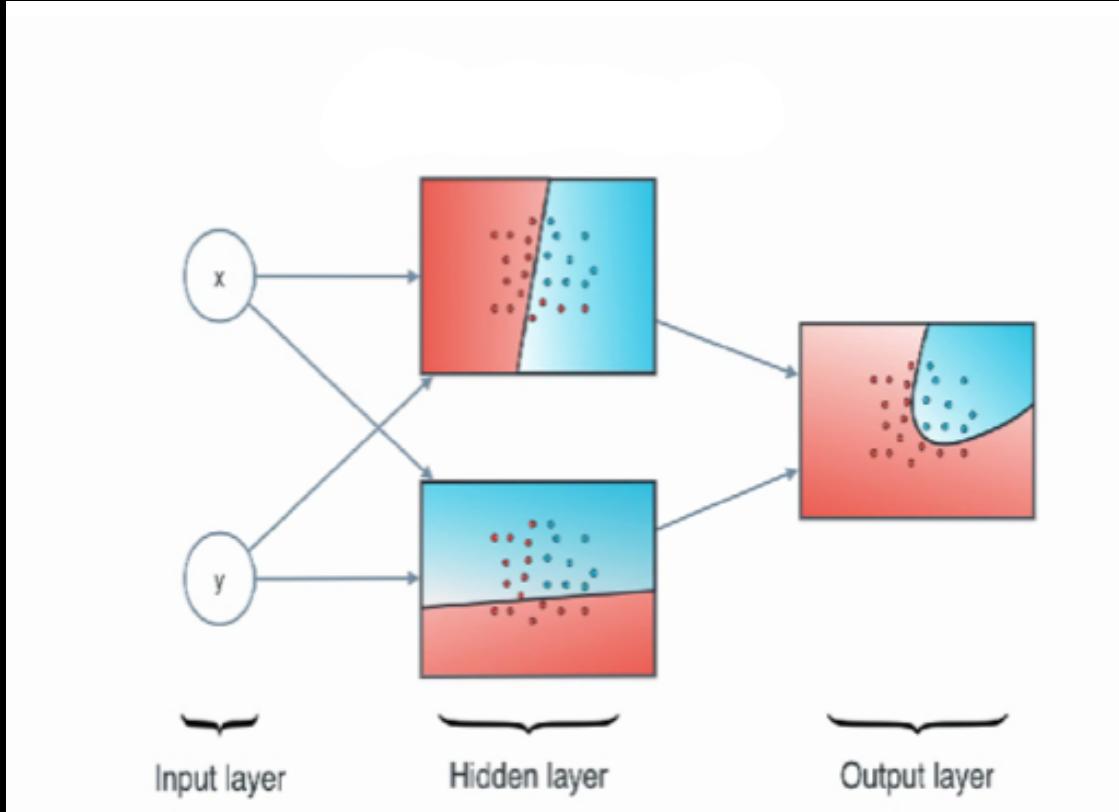
* Objectives of N_1 & N_2 (?)

$$\begin{matrix} \text{neuron-1} \\ \text{neuron-2} \end{matrix} \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} w_{10} \\ w_{20} \end{bmatrix}$$

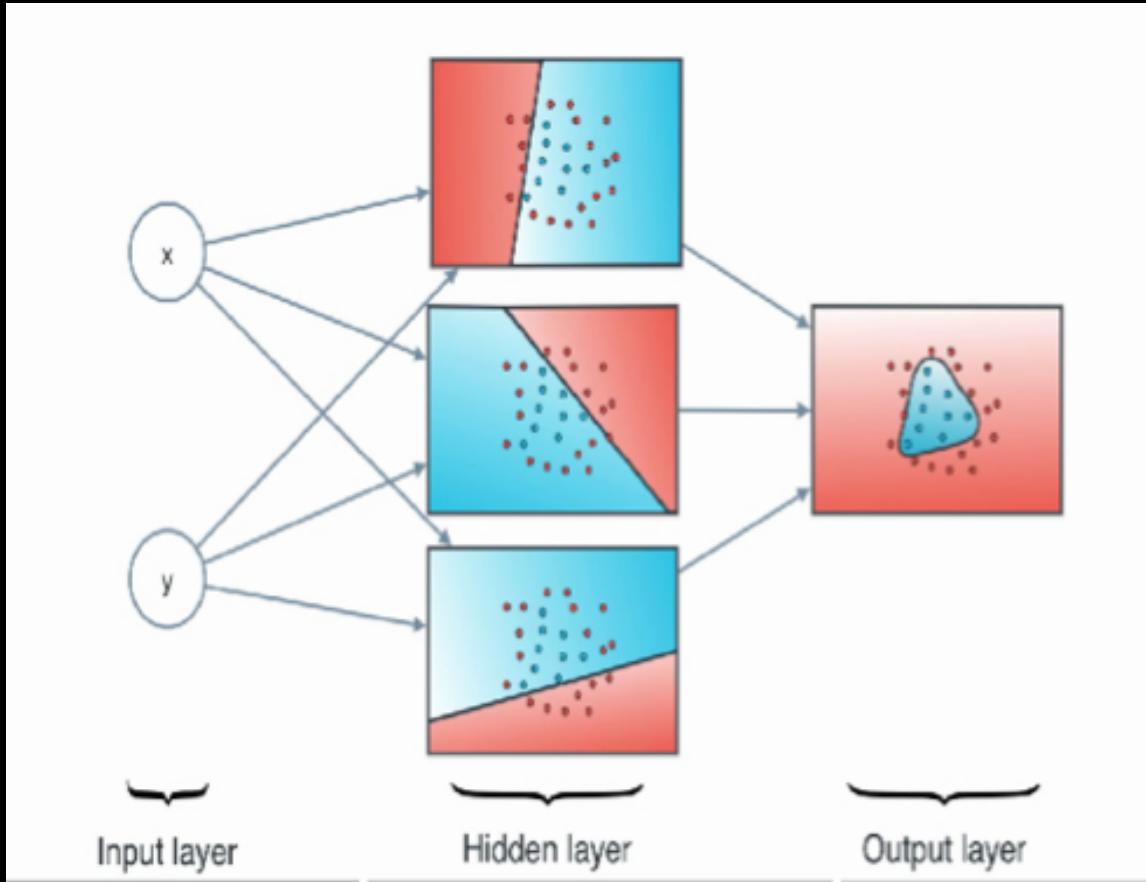
$$\Rightarrow \begin{bmatrix} w_{11}x + w_{12}y + w_{10} \\ w_{21}x + w_{22}y + w_{20} \end{bmatrix} = \begin{bmatrix} I_1 \\ I_2 \end{bmatrix}$$

$$\begin{bmatrix} w_{31} & w_{32} \end{bmatrix} \begin{bmatrix} I_1 \\ I_2 \end{bmatrix}$$

$$\Rightarrow w_{31}I_1 + w_{32}I_2 + w_{30}$$

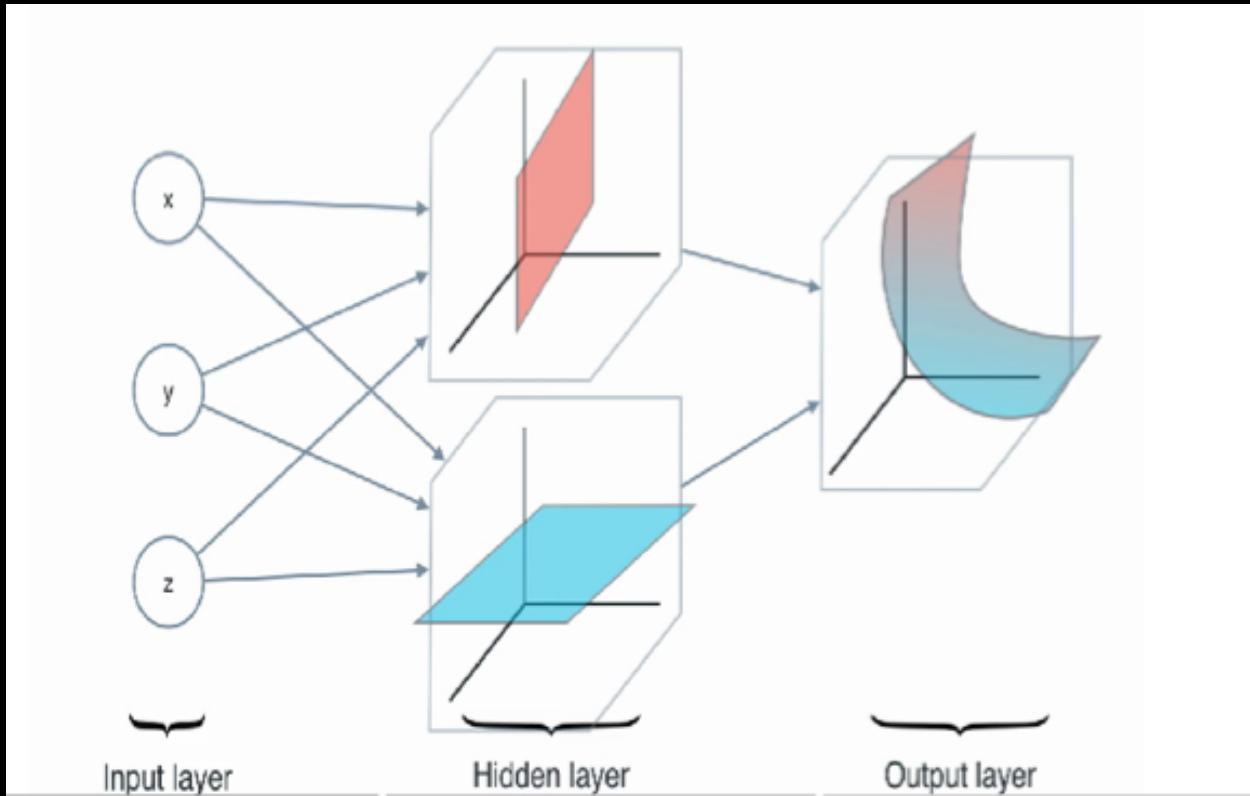


Region Merging

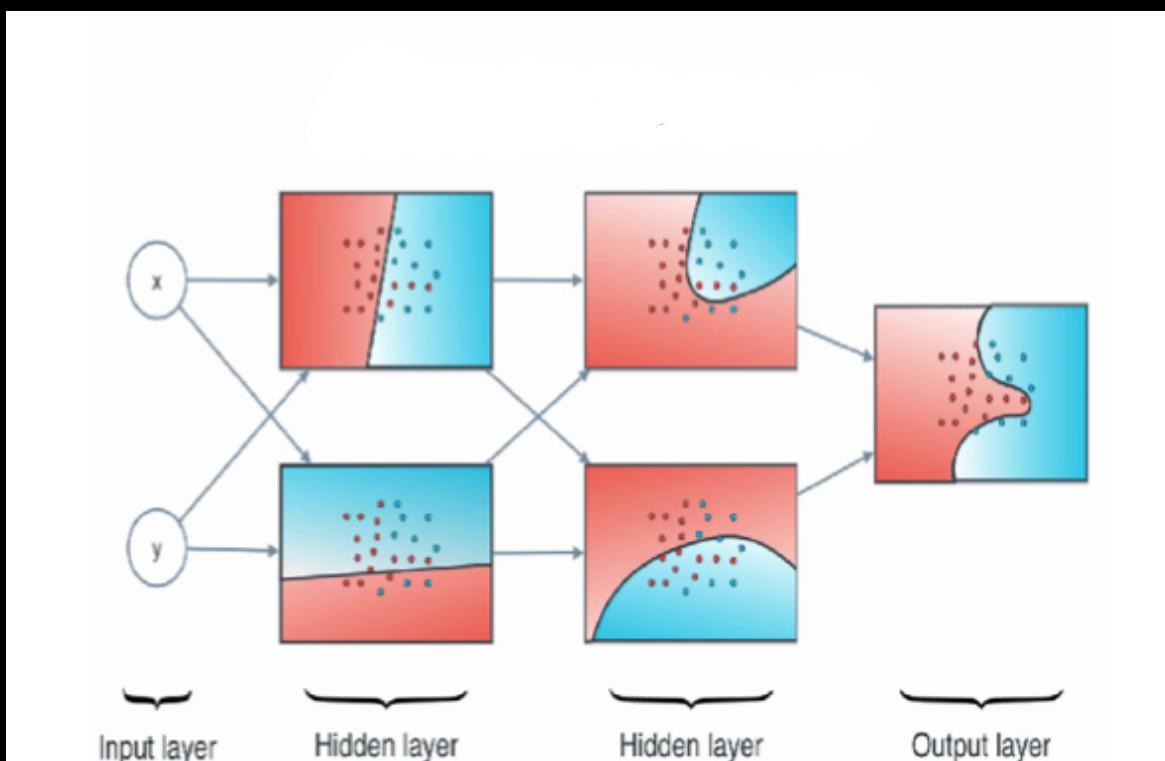


Adding
More
neurons
in a
layer

Scaling Vertically

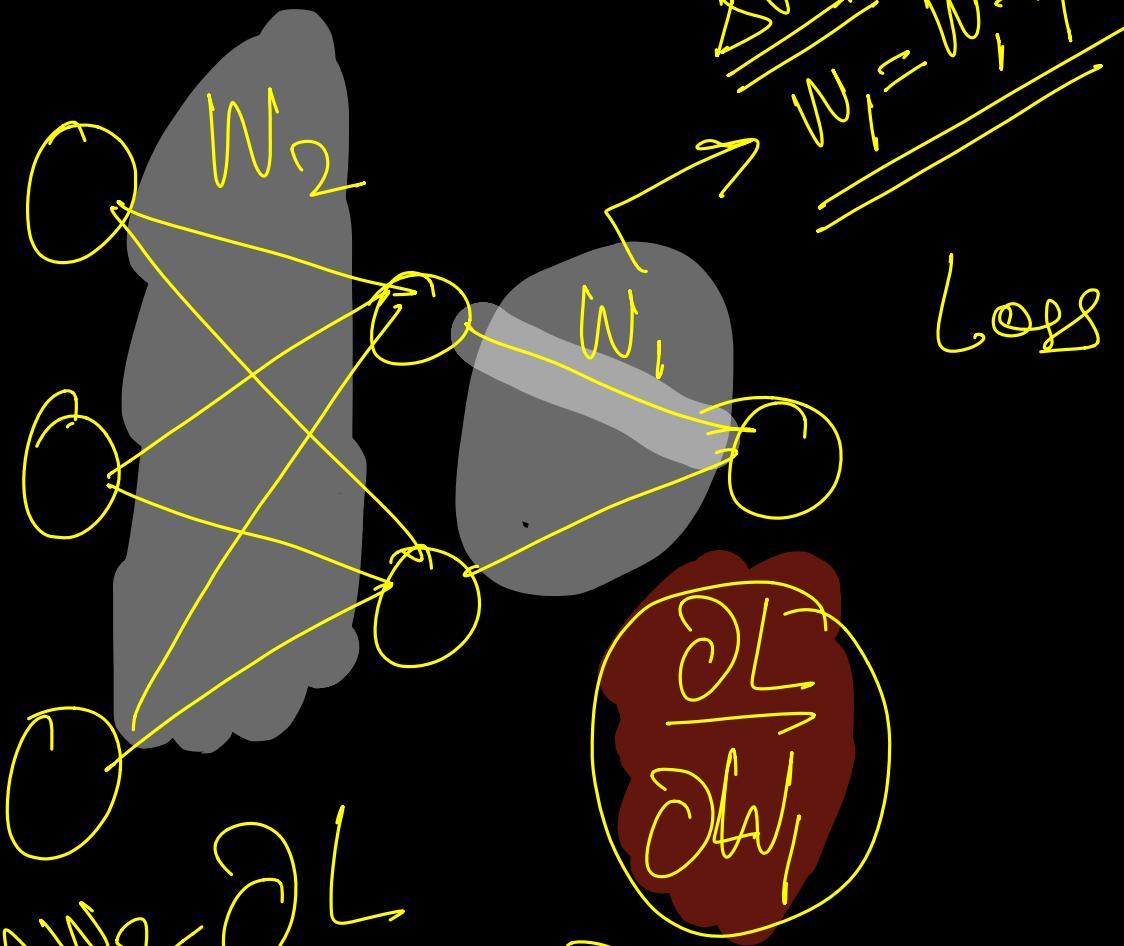


Input dimension > 2



Scaling horizontally.

... N + ... - NDN₁



$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial W_2} = \frac{\partial L}{\partial y} \cdot W_2$$

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial W_1} = \frac{\partial L}{\partial y} \cdot W_1$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial b_1} = \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial b_2} = \frac{\partial L}{\partial y}$$