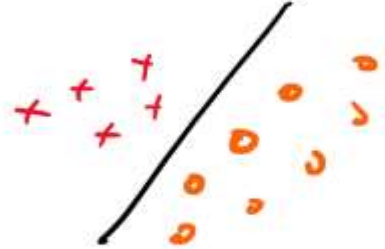B_ Innovotive     <u>Part-2</u>

<u>Agenda</u>

1) Understanding Similarities b/w Logistic regression and perceptron.

2) What is Multi-layered perception.

3) Understanding matematical Notations.

4) Training a Single Neuron Model

5) Training M.L.P.

6) Memoization

# Logistic Regression and perceptron.

A Logistic regression $w$,

Given any point $x_i \longrightarrow \hat{y_i}$

$\hat{y_i}$ is predicted value of $y_i$.

So, in General $\hat{y_i} = \text{Sigmoid} (w^T x_i + b)$

where $x_i \in R^d$, $w \in R^d$, $b \in R$  :)

A dataset $D = \{x_i, y_i\}$

we train L.R $\longrightarrow$ to find $w$ and $b$.

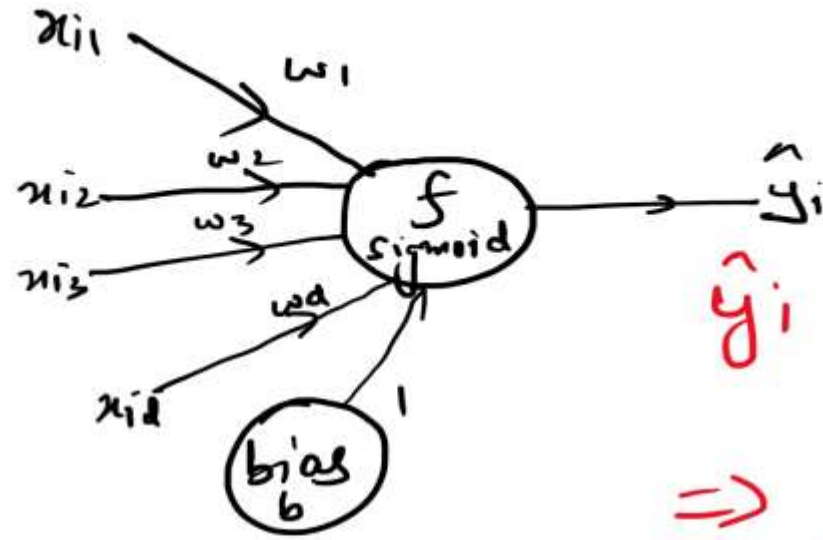$\therefore \quad \hat{y_i} = \text{Sigmoid} (w^T x_i + b).$ } Logistic Regression.

$$\Rightarrow \hat{y_i} = \text{Sigmoid} \left( \sum_{i=1}^{d} w_i x_{ij} + b \right)$$

Let assume a point $x_i$ , $x_j \in R^d$

So we can represent $x_i$ as a vector.

$$x_i = [x_{i1}, x_{i2}, x_{i3} - - - x_{id}]$$

weights $w = [w_1, w_2, w_3 - - - - w_d]$

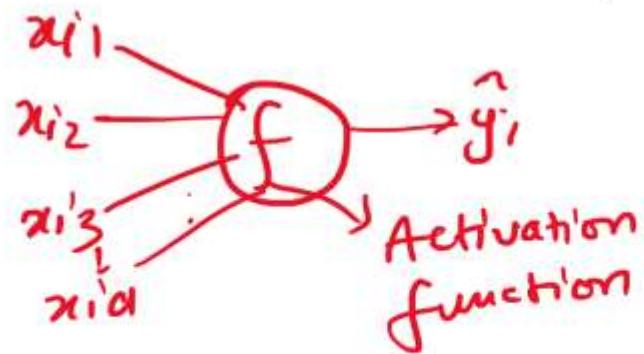Output of Neuron $= f \left( \sum_{i=1}^{d} w_j x_{ij} \right)$

$$\hat{y}_i = f(w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} - \cdots w_d x_{id})$$

$$\Rightarrow \hat{y}_i = f(w^T x_i + b)$$

Training a Neural Network.

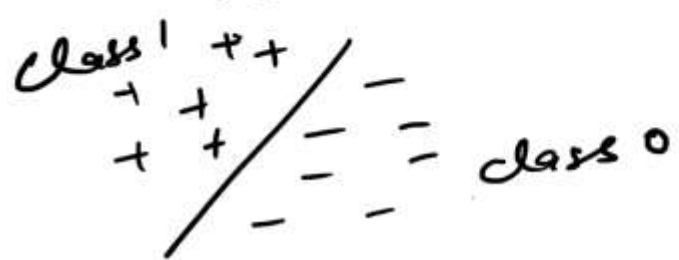→ at states when we have to find weights on edges / vertices.

Let's talk about perceptron.

$x_{i1}$
$x_{i2}$
$x_{i3}$
$x_{id}$

→ $\hat{y_i}$

Activation function

$$f(x) = \begin{cases} 1 & \text{if } w^T x_i + b > 0 \\ 0 & \text{otherwise.} \end{cases}$$

A perceptron is also called a
Linear classifier.

Class 1 + +
+ +
+ + / — —
+ / — — — class 0
/ — —

A perceptron will always tries
to find out the best line
which separates class 1 from
class 0.
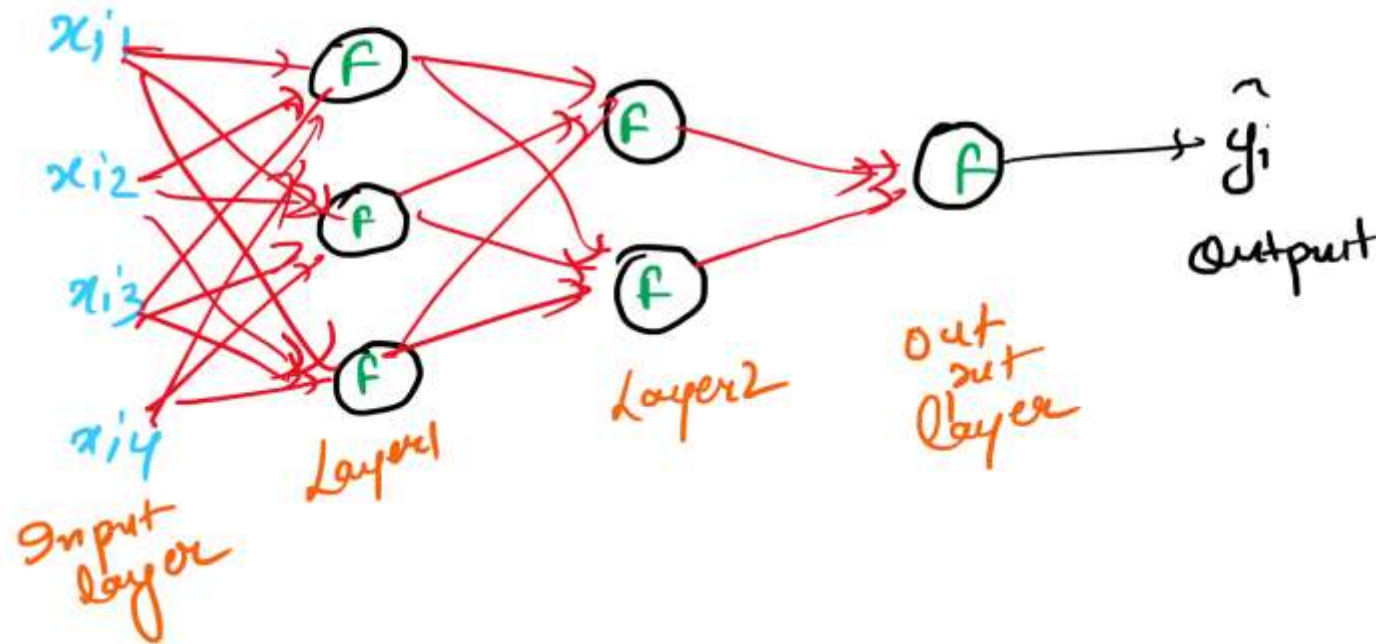
In Logistic Regression → Squashing function
(Sigmoid)

In perceptron → No Squashing function

Simple Neuron model, only difference is activation function.

# Multi layered perceptron.

↳ Bunch of Connected Neurons.



$x_{i1}$

$x_{i2}$

$x_{i3}$

$x_{i4}$

Input layer

Layer1

Layer2

out put layer

$\hat{y_i}$

Output

Multi layered perceptron or Neural Network.

Why should we care about Multi layered perceptron?

→ Since perception has been terine of Biological Neural Network, on Neuro science → Humans, rats, monkeys, ants, found collection of Neurons interconnected in smartest way.

How do you Connect it and make it a Network?

→ Lot of mathematical arguments happened.

Let's take Regression problem, $D = \{x_i, y_i\}$

Determine $y_i = f(x_i)$, $y_i \in R$.

Let's assume, $y$

$$\mathcal{D} = \{x_i, y_i\}$$

$x_i \in R^1$ } (one
$y_i \in R$ } dimensional)

Objective $\rightarrow$ To find $f$ in , $y_i = f(x_i)$

Case-1    Let say    $f(x_i) = y_i$

$\therefore x_i = y_i$

$$f(x) = 2 \sin(x^2) + \text{sqrt} \ (5x).$$

$\hookrightarrow$ Composition of function. } Basically
Complex function
gives enormous
power to ry solute
problem.

# Composition of function.

→ $fog(x)$ or $gof(x)$

→ MLP is a Graphical way of representing simple function Composition.

→ It results in powerful model.

→ we can also overfit easily.

# Notations:



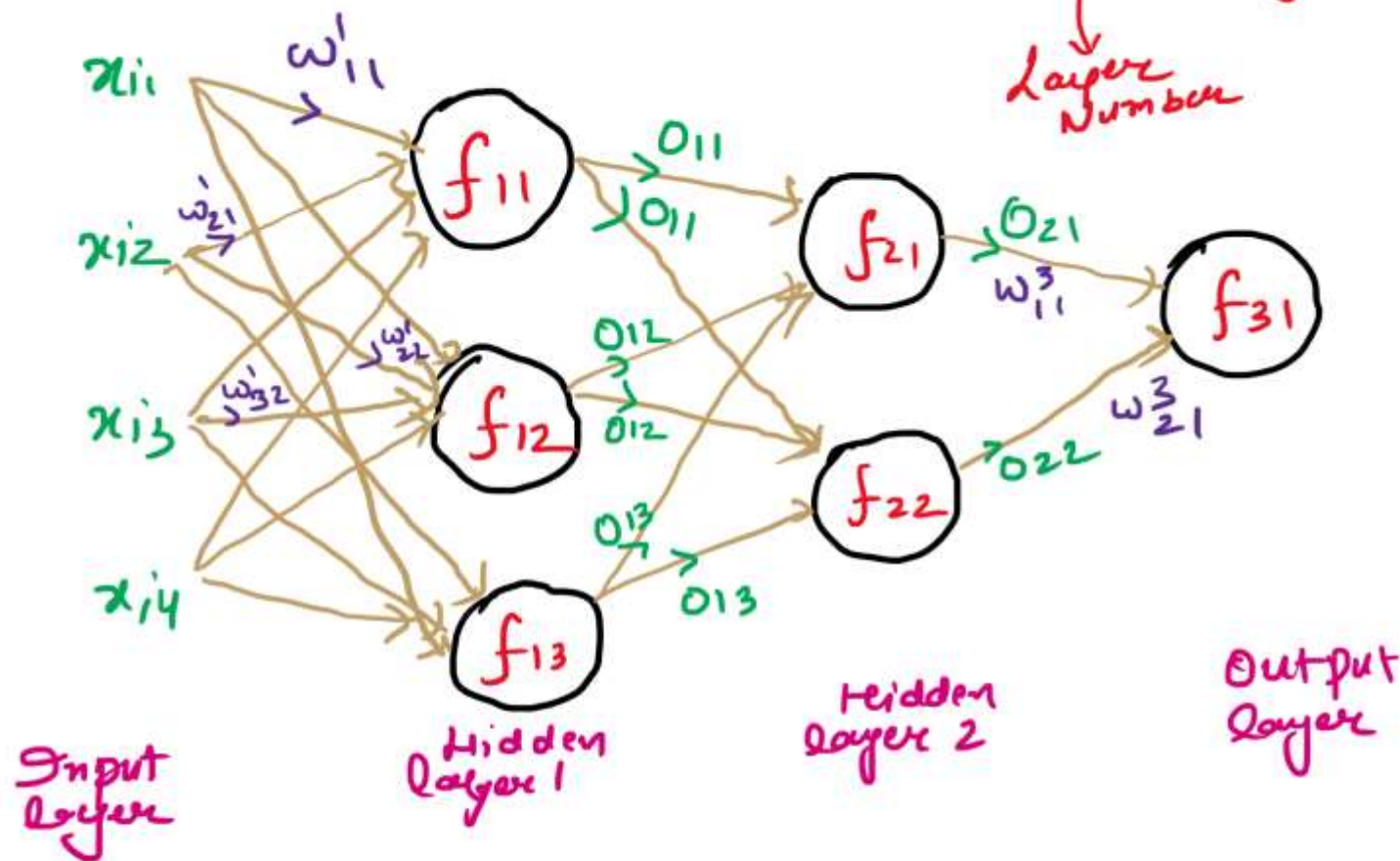$f_{ij}$ → function number

Layer Number

$f_{11}$ → Layer 1, function 1
$f_{22}$ → Layer 2, function 2

$O_{11}$ → output of layer 1 from function 1

$O_{12}$ → output of layer 1 from function 2

$O_{13}$ → output of layer 1 from function 3.

$w_{ij}^{k}$   $k$ → what is next layer number.
$i$ → from which neuron
$j$ → to which Neuron.

$x_{i1}$

$w_{11}^{1}$

$w_{21}^{1}$

$x_{i2}$

$w_{22}^{1}$

$w_{32}^{1}$

$x_{i3}$

$x_{i4}$

$f_{11}$   $O_{11}$   $O_{11}$

$f_{12}$   $O_{12}$   $O_{12}$

$f_{13}$   $O_{13}$   $O_{13}$

$f_{21}$   $O_{21}$   $w_{11}^{3}$

$f_{22}$   $O_{22}$   $w_{21}^{3}$

$f_{31}$

Input Layer

Hidden Layer 1

Hidden Layer 2

Output Layer

$x_{i1}$

$w'_{11}$

$f_{11}$

$w'_{12}$

$w'_{21}$

$x_{i2}$

$w'_{22}$

$w'_{31}$

$f_{12}$

$w'_{23}$

$x_{i3}$

$w'_{32}$

$w'_{13}$

$w'_{41}$

$x_{i4}$

$w'_{42}$

Input layer
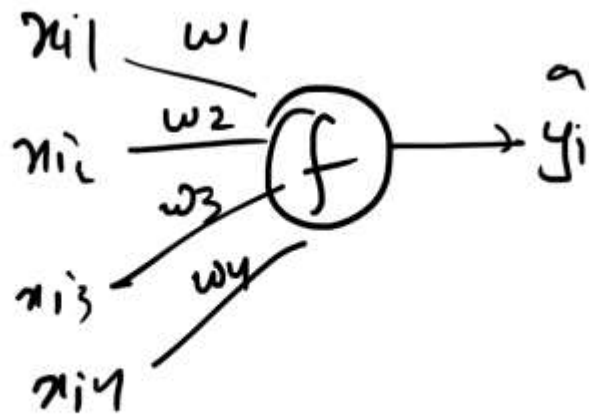
$w'_{43}$

$f_{13}$

Layer 1

How many weights.

4 inputs and 3 units (Neurons)

$4 \times 3 \Rightarrow 12$ weights.

$$\begin{array}{c} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{array} \begin{bmatrix} w'_{11} & w'_{12} & w'_{13} \\ w'_{21} & w'_{22} & w'_{23} \\ w'_{31} & w'_{32} & w'_{33} \\ w'_{41} & w'_{42} & w'_{43} \end{bmatrix} {\scriptstyle 4\times 3}$$

Step by Step process of Training a
Single Neuron Model.

→ Training refers to find the best edge weights
in a network model using training data.

→ perceptron and Logistic regression → **Single Neuron model for classification.**

→ for Linear Regression → **Single Neuron model for regression.**

$$x_{i1} \xrightarrow{w_1}$$
$$x_{i2} \xrightarrow{w_2} \;\; \boxed{f} \xrightarrow{a} \hat{y}_i$$
$$x_{i3} \xrightarrow{w_3}$$
$$x_{i4} \xrightarrow{w_4}$$

**Linear regression**

$$\hat{y}_i = \sum_{j=1}^{d} w_j x_{ij}$$

$$\hat{y}_i = w^T x_i$$

$$x_i \in R^d$$
$$y_i \in R$$

Activation function in Linear Regression cu

$$f(z) = z \quad (\text{called Identity function})$$

In Linear regression

$$\min_{w_i} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + reg.$$

where $\hat{y}_i = w^T x_i$

$n \to$ No. of points in training data.

$$\Rightarrow \min_{w_i} \sum_{i=1}^{n} (y_i - w^T x_i)^2 + \|w\|_2^2 \Bigg\}$$

Defining a loss function.

$$\sum_{i=1}^{n} \mathcal{L}_i \Rightarrow \mathcal{L} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + reg.$$

becomes a optimization problem

Let's take
on one
training point.

$$\mathcal{L}_i = (y_i - \hat{y}_i)^2$$

$$\therefore \quad \hat{y}_i = w^T x_i$$

Computing loss function on top of $y_i$

$$\mathcal{L}(y_i, \hat{y}_i)$$

② Write the optimization problem.

$$\min_{w_i} \sum_{i=1}^{n} \left(y_i - \underbrace{w^T x_i}_{\downarrow \ \hat{y}_i}\right)^2$$

$$\hat{y}_i = f(w^T x_i)$$
$$\hookrightarrow \text{linear function,}$$

$$\therefore \quad \min_{w_i} \sum_{i=1}^{n} \left(y_i - f(w^T x_i)\right)^2 + \text{reg.}$$

for perceptron, Activation function
↓
Threshold function.

∴ updated weights:

$$w^* = \arg\min_{w} \sum_{i=1}^{n} (y_i - f(w^T x_i))^2 + \text{reg}$$

③ Solve the optimization problem.

a) Initialization of weights.
  ↳ random initialization. $w_i$'s

b) Computing the derivative of $L$ w.r.t $\omega$.

$$\nabla_\omega L = \begin{bmatrix} \dfrac{\partial L}{\partial \omega_1} \\ \dfrac{\partial L}{\partial \omega_2} \\ \vdots \\ \dfrac{\partial L}{\partial \omega_4} \end{bmatrix} \qquad -: x_i \in R^4.$$

c) $W_{new} = W_{old} - \eta \left( \nabla_w L \right)_{W_{old}}$

$\downarrow$ updated weight.

$\downarrow$ old weight

$\searrow$ Learning rate.

$$(W_i)_{new} = (W_i)_{old} - \eta \left( \frac{\partial L}{\partial W_i} \right)_{(W_i)_{old}}$$

for iteration    1 to k

Computing    Gradient    Descent.

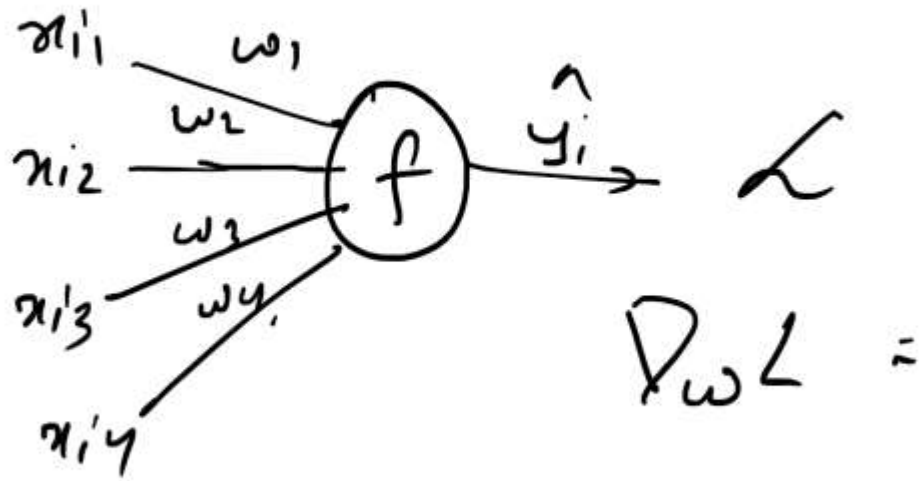$$\nabla_w L \longrightarrow x_i's \text{ and } y_i's.$$

Computing    Stochastic    Gradient    Descent.

$$\nabla_w L \sim \text{one point}$$

Computing  mini  batch  Stochastic  Gradient
descent
$$\{x_i, y_i\}.$$

How to Compute these derivatives.

$x_{i1}$ $w_1$

$x_{i2}$ $w_2$

$x_{i3}$ $w_3$

$x_{i4}$ $w_4$

$f$ $\hat{y_i}$ $\rightarrow$ $L$

Squared loss in regression problem.

$$D_w L = \begin{bmatrix} \dfrac{\partial L}{\partial w_1} \\ \dfrac{\partial L}{\partial w_2} \\ \vdots \\ \dfrac{\partial L}{\partial w_4} \end{bmatrix}^T$$

Applying Chain Rule.

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_1}$$

Always Check Loss function and the Weight (say $w_1$), what comes in between and start differentiating term.

$$w_1 \xrightarrow[\text{path.}]{\textcircled{f}} L$$

$$\mathcal{L} = \sum_{\hat{i}=1}^{n} (y_i - \hat{y}_i)^2 + \text{regularisiere.}$$

$$\Rightarrow \sum_{i=1}^{n} (y_i - f(\omega^T x_i))^2$$

$$\Rightarrow \sum_{i=1}^{n} (y_i - f)^2$$

$$\frac{\partial \mathcal{L}}{\partial f} = (-) \sum_{i=1}^{n} \cdot 2 (y_i - f(\omega^T x_i))$$

$$\frac{\partial f}{\partial \omega_1} = x_i$$

$$\Rightarrow \quad \frac{\partial L}{\partial w_1} = -2 \cdot x_i \left( y_i - \hat{y_i} \right)$$

$$\Rightarrow \quad \frac{\partial L}{\partial w_1} = \sum_{i=1}^{n} (-2) x_i \left( \hat{y_i} - \hat{y_i} \right)$$