# History of Neural Networks
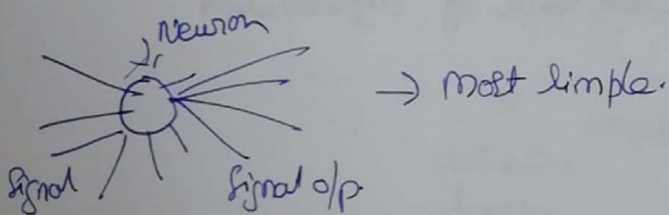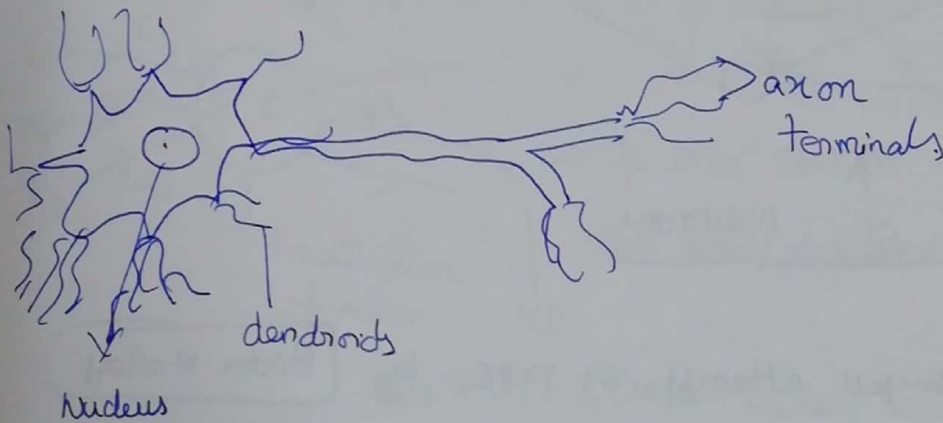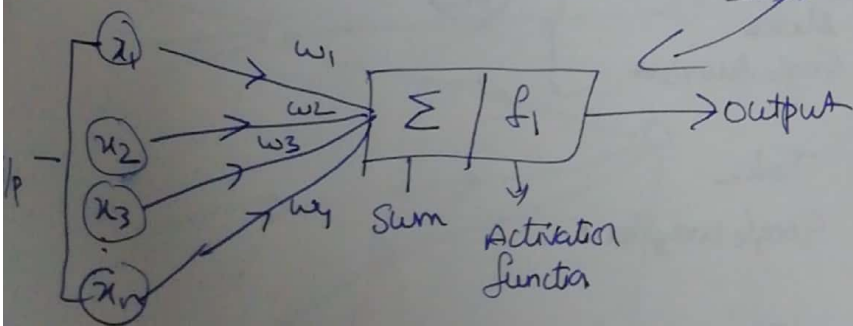
47:1    History of Neural network & Deep Learning

→ Perceptron ——→ 1957 [Rosenblatt]
     ↳ A simple model
     ↳ similar to Logistic Regression



axon
terminals

dendroids

Nucleus



Neuron

→ most simple.

Signal    Signal o/p

Electrical pulses comes in and some process done and sends out the
electrical pulses

Artifical neural n/w



$x_1$   $w_1$
$x_2$   $w_2$
$x_3$   $w_3$
   $w_4$
$x_n$

i/p

$\Sigma$ | $f_1$  → output
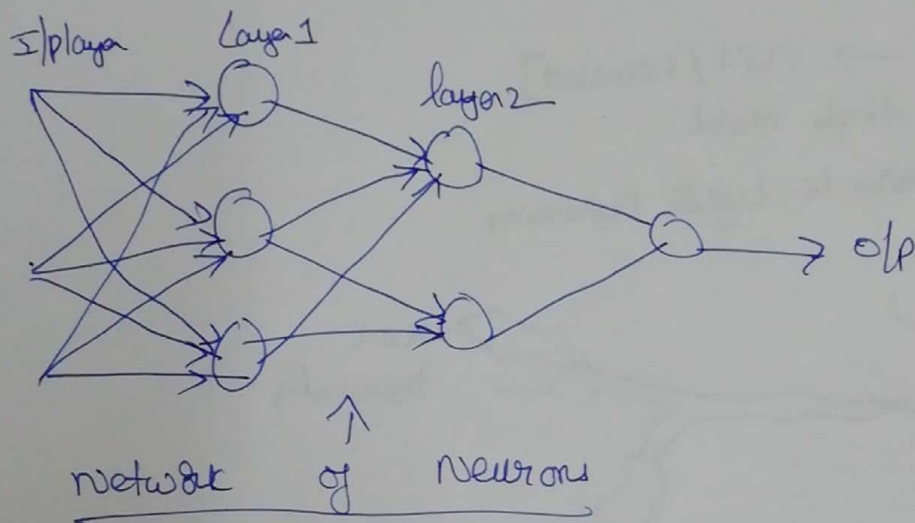
Sum   Activation
     functon

$$f\left(w_1 x_1 + w_2 x_2 + \cdots w_n x_n\right)$$

→ loosely inspired from biology

In biology $\longrightarrow$ Brain

-) There will be multiple neurons not only one

I/p layer    Layer 1

layer2

$\uparrow$
network of Neurons

-) There is successfull attempt in 1986 by [Hinton & ollors]

→ Back propagation Algo → Chain rule of differentiation

→ As depth of

In 2012

4 image net

-) Voice assistant → Siri
                      Cortana
                      Alexa       } DL
                      Google Assistant

Self driving cars → Tesla
                     Google waymo d

Skin Cancer → AI

## Simplified View of Neuron

$x_1$ →
$x_2$ →
$x_3$ ↑

dendroids

axons.

→ In dendroid is thicken, the more weight is associated with it.

$x_1$ → $w_4$
$x_2$ → $w_2$  $\Sigma, f$  o/p↑
$x_3$ → $w_3$

→ activation function.

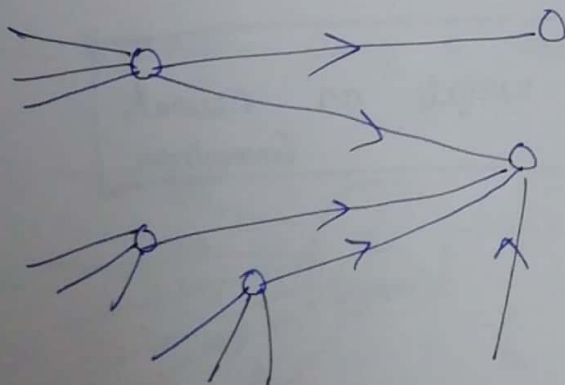$$\Big\{ \begin{array}{l} \text{activated} \\ \text{fired.} \end{array}$$

→ If more Electrical signal in the i/p. the cell fired / Activated.

o/p↑ → $O_1 = f(w_1 x_1 + w_2 x_2 + \cdots + w_n x_n)$

$$O = f\Big( \sum_{i=1}^{n} w_i x_i \Big)$$
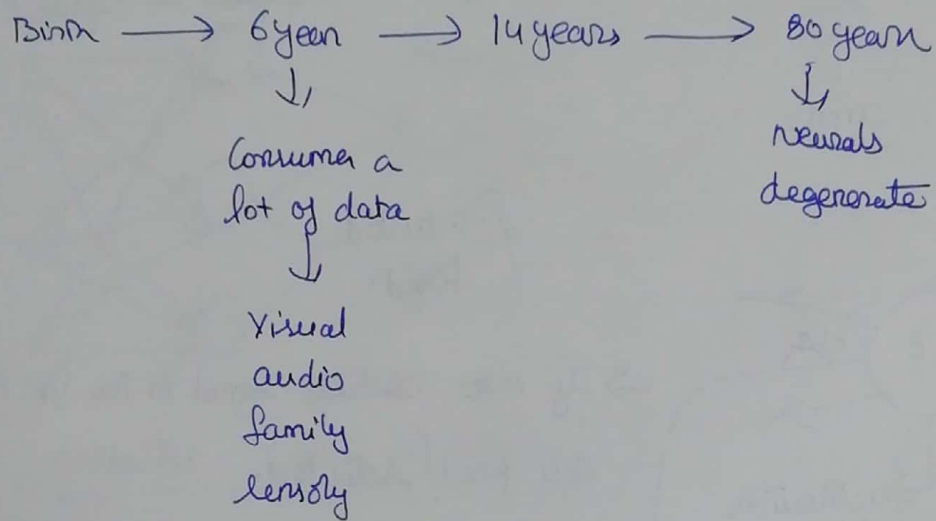
activation function.

weights

i/p's

## 47.3 Growth of biological Neural network

→ Front part of brain    Cerebral Content
$$J_1$$
front part of brain.

learning $\rightarrow$ Connecting neurons with Edges

more - connection $\rightarrow$ more Calories required.

Birth $\longrightarrow$ 6 years $\longrightarrow$ 14 years $\longrightarrow$ 80 years
$\downarrow$ $\downarrow$
Consumer a neurals
lot of data degenerate
$\downarrow$
Visual
audio
family
sensory

$\rightarrow$ Japanese kid VS Indian kid

$\rightarrow$ Connectia are formed based on data.

hot pain

$\circ$

$\circ$ $\circ$

all biological learning $\xrightarrow{+}$ weights on neural Connections

Anu

47.4 : Diagramatic Representation : Logistic Regression & Perceptron.

LR : $x_i \longrightarrow \hat{y}_i \rightarrow$ Predicted Value of $(y_i)$

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

$\mathcal{D} = \{x_i, y_i\}$      Train LR $\rightarrow$ we will find ~~the~~ vector $(w)$ & B

$$W \in \mathbb{R}^d$$
$$b \in R \text{ (Real number)}$$

we can also write as

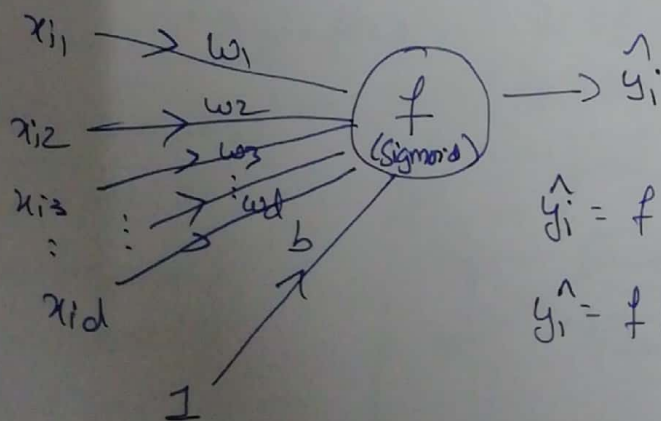$$\hat{y}_i = \text{Sigmoid} \left( \sum_{j=1}^{d} w_j x_{ij} + b \right)$$

$$x_i = [x_{i1}, x_{i2}, x_{i3} \cdots x_{id}]$$
$$w = [w_1, w_2, w_3 \cdots w_d]$$

Part $\sum_{j=1}^{d} (w_j x_{ij})$

o/p of neuron (o) is $f\left( \sum_{j=1}^{d} w_j x_{ij} \right)$

Logistic Regression using a neuron.

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} * \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$
$$\phantom{xx}2 \times 2 \phantom{xx} 2 \times 1$$

$$\begin{bmatrix} 1 \times 2 + 2 \times 3 \\ 3 \times 2 + 4 \times 3 \end{bmatrix} = \begin{bmatrix} 8 \\ 18 \end{bmatrix}$$
$$\phantom{xxxxxxxxxxx}2 \times 1$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} * \begin{bmatrix} 2, & 3 \end{bmatrix}$$
$$\phantom{xx}2 \times 1 \phantom{xxx} 1 \times 2$$

$$\begin{bmatrix} 1 \times 2 & 2 \times 3 \\ 3 \times 2 & 4 \times 3 \end{bmatrix} = \begin{bmatrix} 2 & 6 \\ 6 & 12 \end{bmatrix}$$
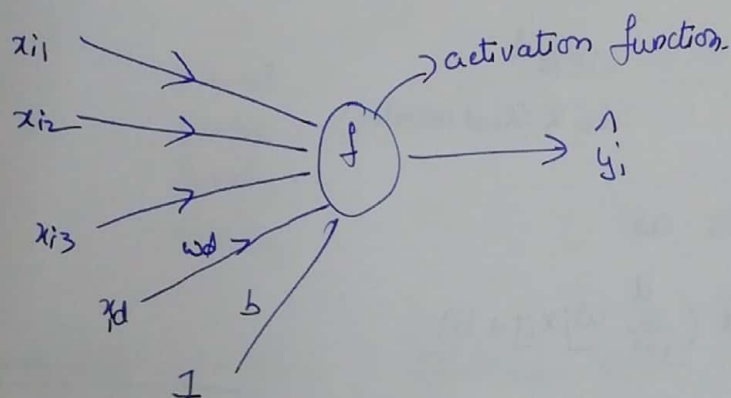$$\phantom{xxxxxxxxx}2 \times 2$$



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

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

If $f$ is sigmoid This is how we will represent LR in the form of neuron.

→ Training a neural network (NN) ⟹ Computing the weights on
Edges/Vertices

→ Represented LR in the form of neuron using a activation
function called Sigmoid

## Perceptron (1957)

$x_{i1}$
$x_{i2}$
$x_{i3}$
$x_d$
$w_d$
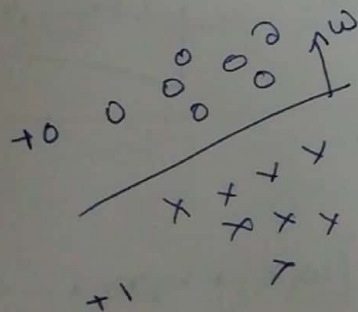$b$
$f$
→ activation function
→ $\hat{y_i}$
1

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

If the activation function fires then $f(x) = 1$
If the activation function doesn't fire $f(x) = 0$

Perceptron is also a br. classifier.

47.5 → multi-layer Perceptron

Perceptron → Single neuron
└→ Logestic Regression

Hidden layers


i/p layer.

o/p layer

Yi
o/P

|← ————— Neural network. ————→|
|← ————— multi layer perceptron ———→|

Ⓠ why should we care about MLP?

(a) biological Inspiration

(b) mathematical

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

└→ $y_i = f(x)$     $x_i \in \mathbb{R}^1$
                    $y_i \in \mathbb{R}$

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

$f_1 \rightarrow$ add ()
$f_2 \rightarrow$ square ()
$f_3 \rightarrow$ sqrt ()
$f_4 \rightarrow$ sin ()
$f_5 \rightarrow$ mul ()

hidden layers.

I/P



multi layered perceptron LIKE structure

→By using multi layered structure we can come up with complex multi-layered functions

→ multi layered structure → Enormous power to the model.

$$\begin{cases} \text{linear model} \longrightarrow \text{Simplest model.} \\ \text{non-linear} \rightarrow \text{RBF-SVM, RF, GBT} \end{cases}$$

## Function Composition

$$f \cdot g(x) = g \cdot f(x)$$

$F(x) = 2\sin x^y + \sqrt{5x}$

Let's pick   $2 * \sin x^y$       ,   $\sqrt{5x}$

$$f_5 (2, f_4(f_2(x)))$$       $$f_3(f_5(x,5))$$

$$F(x) = f_1 (f_5(2, f_4(f_2(x))), f_3(f_5(x,5))) \rightarrow (MLP)$$

$$f(g(x)) = f \cdot g(x)$$

$$g(f(x)) = g \cdot f(x)$$
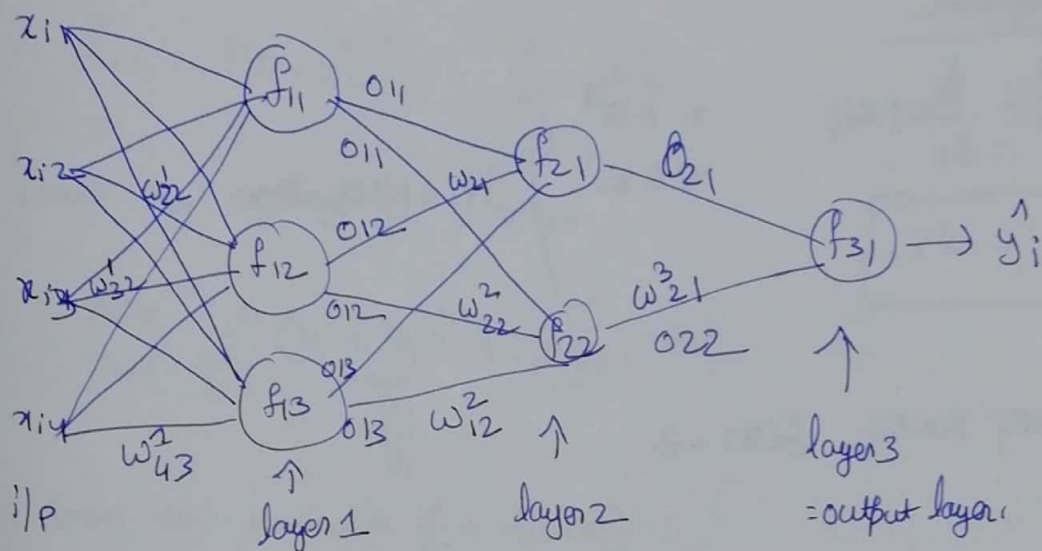
In MLP : graphical way of representing fog, gof

→ Having a multilayered structure result in powerful models we can overfit very Easily

## 47.6 Notation

$D: \{x_i, y_i\}$ ; $x_i \in \mathbb{R}^4$ ; $y \in \mathbb{R}$ (Regression)

$x_{ij}$
↗ ↑
point feature



$x_{i1}$
$x_{i2}$
$x_{i3}$
$x_{i4}$

$f_{11}$   $O_{11}$
$O_{11}$
$f_{12}$   $O_{12}$
$O_{12}$
$f_{13}$   $O_{13}$
$O_{13}$

$\omega_{21}^1$   $\omega_{32}^1$   $\omega_{43}^1$

$\omega_{24}$   $f_{21}$   $O_{21}$
$\omega_{22}^2$   $f_{22}$   $O_{22}$
$\omega_{12}^2$

$\omega_{21}^3$
$f_{31}$ → $\hat{y}_i$

i/p          layer 1          layer 2          layer 3 = output layer

$f_{ij}$
↗ ↖
layer   index

$O_{ij}$
↗ ↑
$i^{th}$ layer   neuron

$\omega_{ij}^{K}$ ← next layer
↗ ↑
from   TO

$$W^1 = \begin{bmatrix} \omega_{11}^1 & \omega_{12}^1 & \omega_{13}^1 \\ \omega_{21}^1 & \omega_{22}^1 & \omega_{23}^1 \\ \omega_{31}^1 & \omega_{32}^1 & \omega_{33}^1 \\ \omega_{41}^1 & \omega_{42}^1 & \omega_{43}^1 \end{bmatrix}_{4 \times 3}$$

$$W^2_{3 \times 2} = \begin{bmatrix} & & \\ & & \end{bmatrix}_{3 \times 2} \qquad W^3_{2 \times 1} = \begin{bmatrix} & \\ & \end{bmatrix}_{2 \times 1}$$

47.7   (Training) a Single-neuron model

find the best Edge-weights using $D_{Tr}$

Perceptron & LR $\rightarrow$ Single neuron models for classification

(Early) $\rightarrow$ Lr. Reg $\rightarrow$      "      "    for regression
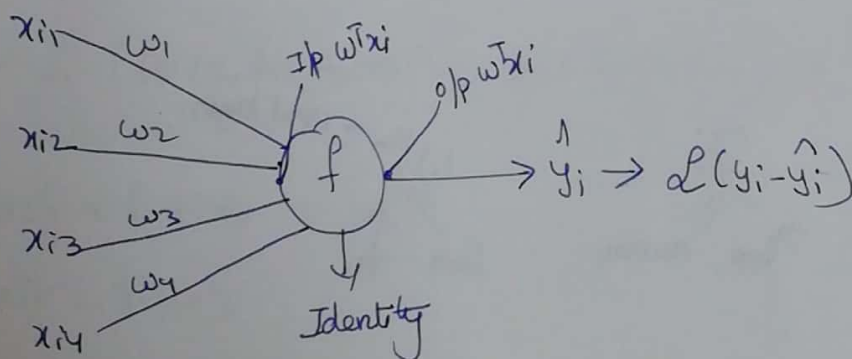
## Lr. Regression

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

$$\left. \begin{array}{l} x_i \in \mathbb{R}^d \\ y_i \in \mathbb{R} \end{array} \right\} \text{lr. optimization.}$$

$$\boxed{\hat{y_i} = w^T x_i}$$

Identity function $f(z) = z$



Input $w^T x_i$, output $w^T x_i$, $\hat{y_i} \rightarrow \mathcal{L}(y_i - \hat{y_i})$

Identity

## Lr-reg

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

# pt in tr data

---

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

$\mathcal{D} = \{x_i, y_i\}_{i=1}^{n}$

$x_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R}$

① Define loss-function

$$\mathcal{L} = \sum_{i=1}^{n} (y_i - \hat{y_i})^2 + \boxed{reg} \rightarrow \text{keep away for now}$$

$$\mathcal{L}_i = (y_i - \hat{y_i})^2 \qquad \hat{y_i} = w^T x_i$$

② write the optimation problem.

$$\min_{w_i} \sum_{i=1}^{n} \underbrace{(y - w^T x_i)^2}_{\hat{y_i}} + \boxed{reg} \rightarrow y$$

from NN perspective $\hat{y_i} = f(w^T x_i)$

$\hookrightarrow$ Lr. reg $f$ is Identity (I)

$\hookrightarrow$ Log. reg $f$ is sigmoid & logistic.

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

$\hookrightarrow$ I : lr. reg

$\hookrightarrow$ Sigmoid : log reg

$$w^{+} = \underset{w}{argmin} \sum_{i=1}^{n} (y_i - f(w^T x_i))^2 + reg$$

③ Solve the optimization problem

ⓐ Initialization of $w_i$'s $\rightarrow$ random

ⓑ $\nabla_w \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial w_2} \\ \frac{\partial \mathcal{L}}{\partial w_3} \\ \frac{\partial \mathcal{L}}{\partial w_4} \end{bmatrix}$

$\rightarrow$ vector represent

$w \in \mathbb{R}^d$

$x_i \in \mathbb{R}^d$

© $w_{new} = w_{old} - \eta \left[ \nabla_w L \right]_{w_{old}}$

$(w_i)_{new} = (w_i)_{old} - \eta \left[ \dfrac{\partial L}{\partial w_i} \right]_{(w_i)}$

for iter = 1 to k

$\{$

$\nabla_w L = \left[ \dfrac{\partial L}{\partial w_1}, \dfrac{\partial L}{\partial w_2}, \dfrac{\partial L}{\partial w_3}, \dfrac{\partial L}{\partial w_4} \right]^T$

$\dfrac{\partial L}{\partial w_1} = \dfrac{\partial L}{\partial f} \cdot \dfrac{\partial f}{\partial w_1} \rightarrow$ chain rule of diff



$L = \sum\limits_{i=1}^{n} (y_i - \hat{y}_i)^2 + reg \lambda$

$= \sum\limits_{i=1}^{n} (y_i - \underbrace{f(\vec{w} x_i)}_{\hat{y}})^2$

$\dfrac{\partial L}{\partial f} = \sum\limits_{i=1}^{n} 2(y_i - f(\vec{w}^T x_i))$

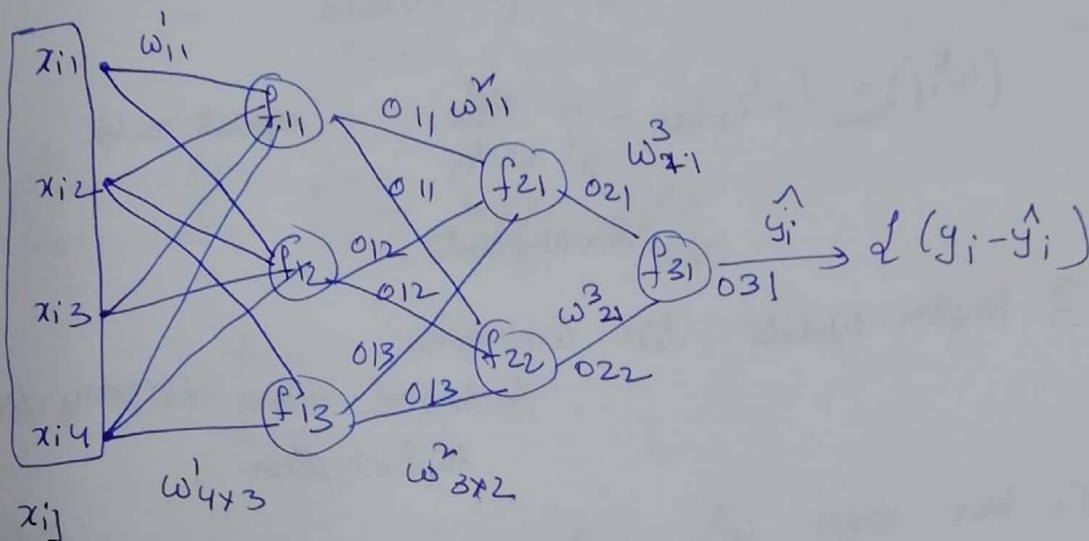$\dfrac{\partial f}{\partial w_1} = \dfrac{\partial f(w^T x_i)}{\partial w_i} = x_i$

$\dfrac{\partial L_j}{\partial w_1} = -2(y_i - \hat{y}_i)^2 x_i$

$\dfrac{\partial L}{\partial w_i} = \sum\limits_{i=1}^{n} -2x_i (y_i - \hat{y}_i)$

## 47.8   Training a MLP : Chain Rule

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

$$\left.\begin{array}{l} x_i \in \mathbb{R}^4 \\ y_i \in \mathbb{R} \end{array}\right\} \text{regression problem}$$



$$x_{ij}$$

<u>step 1</u>   $\mathcal{D} = \{x_i, y_i\}$

determine   $\underset{12}{w^1_{4\times3}}, \quad \underset{6}{w^2_{3\times2}} \quad \underset{2}{w^3_{2\times1}} \quad = 20$

① $\mathcal{L} = \underset{i=1}{\overset{n}{\sum}} (y_i - \hat{y}_i)^2 + \underbrace{reg.}$

$\qquad \underbrace{\qquad\qquad}_{sq. loss.}$

$\qquad\qquad\qquad\qquad \underset{i,j,x}{\sum} (w^k_{ij})^2 \rightarrow L2$

$\qquad\qquad\qquad\qquad \underset{i,j,k}{\sum} |w^k_{ij}| \rightarrow L1$

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

$\mathcal{L} = \underset{i=1}{\overset{n}{\sum}} \mathcal{L}_i + reg$

<u>optimization</u> : min the sq.loss $\mathcal{L} \overset{sq.loss}{\underset{reg}{\longleftarrow}}$

$\qquad\qquad\quad w^1, w^2, w^3$

$$\boxed{\underset{w^k_{i,j}}{min} \; \mathcal{L}}$$

② SGD ठी GD

$$\frac{\partial L}{\partial w_{ij}^k} \quad (\text{Generic})$$

ⓐ Initialize $w_{ij}^k$ randomly $\}$ → lots of technique

ⓑ Suppose we have a weight $(w_{ij}^k)_{new}$

$$(w_{ij}^k)_{new} = (w_{ij}^k)_{old} - \eta \underbrace{\frac{\partial L}{\partial w_{ij}^k}}_{\nearrow} \qquad (\text{update rule})$$

learning rate

ⓒ Perform update till Convergence

↳) old & new value are very close
to Each other.

Let's pick weight $w_{11}^3$ ($w_{11}^3$ impacts $O_{31}$, which impact $\hat{y}_i$) 
loss fnctn
↑
$y_i$

(look before page)

$$\frac{\partial L}{\partial w_{11}^3} = \underbrace{\frac{\partial L}{\partial O_{31}}}_{1} \cdot \underbrace{\frac{\partial O_{31}}{\partial w_{11}^3}}_{2} \quad \leftarrow \text{chain rule} \left. \begin{matrix} \\ \\ \\ \end{matrix} \right\} w^3$$

$$\frac{\partial L}{w_{21}^3} = \frac{\partial L}{\partial x_{31}} \cdot \frac{\partial O_{31}}{\partial w_{21}^3}$$

$$\frac{\partial L}{w_{11}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{11}^2} \left. \begin{matrix} \\ \\ \\ \\ \\ \end{matrix} \right\} w^2$$
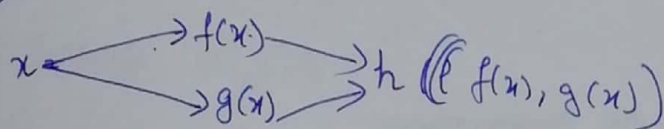
$$\frac{\partial L}{w_{21}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{21}^2}$$

$$\frac{\partial L}{w_{31}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{31}^2}$$

$$\frac{\partial h}{\partial w} =$$



$$x \underset{\to g(x)}{\overset{\to f(x)}{\rightleftarrows}} \to h \,(\!(\; f(x), g(x)\,)$$
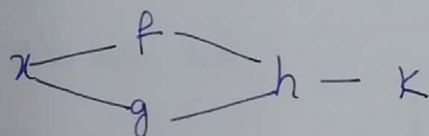
$$\frac{\partial h}{\partial x} = \boxed{\frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x}} \; \boxed{+} \left( \frac{\partial h}{\partial g} + \frac{\partial g}{\partial x} \right) \rightarrow \text{chain rule}$$

$$\downarrow \text{sum.}$$

$$\frac{\partial f}{\partial w_{11}} = \left( \frac{\partial f}{\partial o_{31}} \quad \frac{\partial o_{31}}{\partial o_{21}} \quad \frac{\partial o_{21}}{\partial o_{11}} \quad \frac{\partial o_{11}}{\partial w_{11}} \right)$$

$$x \underset{g}{\overset{f}{\rightleftarrows}} \to h - k$$

$$\frac{\partial k}{\partial x} = \frac{\partial k}{\partial h} \cdot \frac{\partial h}{\partial x} \longrightarrow \boxed{\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial o_{31}} \cdot \boxed{\frac{\partial o_{31}}{\partial w_{11}}}$$

$$\frac{\partial o_{31}}{\partial w_{11}} = \frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{ij}} \; \boxed{+}$$

$$\frac{\partial o_{31}}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}}$$

## 47.9 Training an MLP : Memoization

In Computer Science

        $\hookrightarrow$ Algorithms $\rightarrow$ Dynamic Programming

### Memoization

If there is any operation that is used many times. repeatedly
{It's a good idea to compute it once $\rightarrow$ save it $\rightarrow$ reused it

    $\Downarrow$          $\Downarrow$

 Speed up     Take some memory

    (Math)        (Computer science trick)
  Chain rule  + memoization

             $\downarrow$

     Back Propagation

## 47.10 Back Propagation

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

1) Initialize $w^k_{i,j}$'s

2) For Each $x_i$ in $\mathcal{D}$:

    a) Pass $x_i$ forward through the n/w $\rightarrow$ Forward Propagation.

    b) Compute the loss $\mathcal{L}(\hat{y_i}, y_i)$

    c) Compute all the derivative using chain rule & memoization

    d) update weights from End of the n/w to the start

        for Ex+

           Let's take $w^3_{11}$

$$\frac{\partial L}{\partial O31} \cdot \frac{\partial O31}{\partial w^3_{11}}$$

$$(w^3_{11})_{new} = (w^3_{11})_{old} - \eta \left(\frac{\partial L}{\partial w^3_{11}}\right)$$

        Backword propagation.

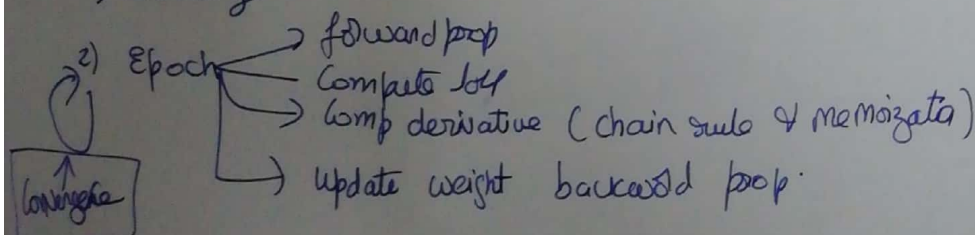3) Repeat step2 till Convergence.

Epoch: Input all the points in the dataset once it is called an Epoch
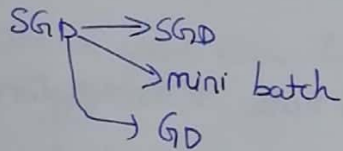
    If we 5 times then It Is Called 5 Epoch

In real time we will run multiple Epoch's

Back propation

  1) Initialize

  2) Epoch
        $\rightarrow$ forward prop
        $\rightarrow$ Compute loss
        $\rightarrow$ Comp derivative (chain rule & memoizata)
        $\rightarrow$ update weight backward prop.

Convergence

<u>V.V. Imp</u> → Back propagation will only work if our activation functions are differentiable

i/p

$$x_1 \rightarrow \boxed{1} \rightarrow \boxed{2} \rightarrow \text{(output)} \rightarrow y_1$$
$$\quad w^1 \qquad w^2 \qquad w^3$$

$$SGD \rightarrow SGD$$
$$\searrow \text{mini batch}$$
$$\searrow GD$$

{ Keeping all the data points in RAM & Computing "$\theta$" using D is extremely time Consuming

mini-batch based back propagation is widely used technique

① 10k points in D

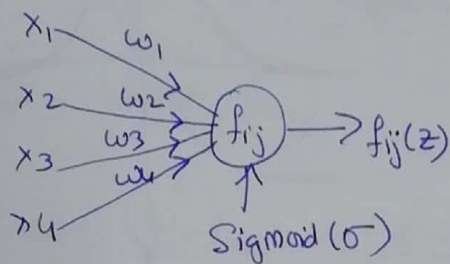    ↳ mini batch = 100 → $\boxed{64, 128, 256, 32}$ → RAM

         Epoch : $\dfrac{10,000}{100} = 100$

② For each batch of size 100

       Forward prop, $\mathcal{L}$, $\dfrac{\partial \mathcal{L}}{\partial w_{ij}}$, update, back prop

47.1❶ Activation functions

In 1980 & 90's only two activation function sigmoid & tanh.



$$z = \sum w_i x_i = w^T x$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

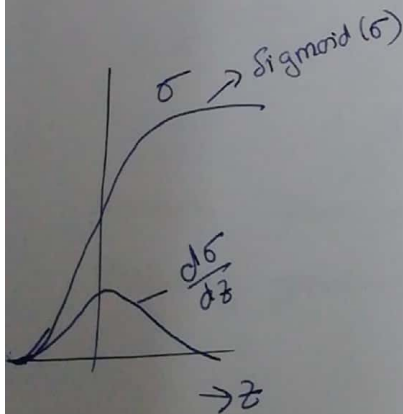$$\sigma(z) = \frac{e^z}{1+e^z} \quad \& \quad \frac{1}{1+e^{-z}} \quad (\text{This is we will use in } LR)$$

Activation function
  ↳ should be differentiable
  ↳ should be easy & fast.

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

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

σ → Sigmoid (σ)



$\frac{d\sigma}{dz}$

→ z

$0 \le \frac{d}{dz} \le 1$

lie b/w 0 & 1 (σ)

$$\frac{da}{dz} = \frac{1+e^{-z}-1}{(1+e^{-z})^2}$$

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

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

$$= a - a^2$$

$$\boxed{= a(1-a)}$$

$$\boxed{\frac{da}{dz} = a(1-a)}$$

$$a = (1+e^{-z})^{-1}$$

$\underbrace{\phantom{(1+e^{-z})^{-1}}}_{g(z)}$

$\underbrace{\phantom{(1+e^{-z})}}_{a(g)}$

$$\frac{da}{dz} = \frac{da(g)}{g} \cdot \frac{dg(z)}{dz}$$

$$= \frac{d}{dg} g^{-1} \cdot \frac{d}{dz}(1+e^{-z})$$

$$= -1 g^{-2} \cdot \frac{d}{dz} 1 + \frac{d}{dz} e^{-z} \qquad h = -z$$

$$= \frac{-1}{g^2} \cdot 0 + \frac{d}{dh} e^h \cdot \frac{d}{dz}(-z)$$

$e^h \cdot -1$

$-e^h$

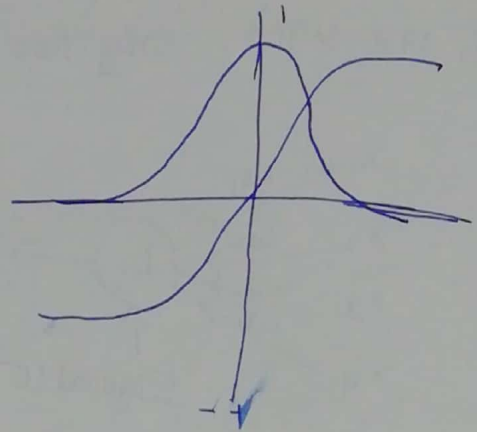$\boxed{-(e^{-z})}$

$$= \frac{-1}{g^2}$$

$$= -(e^{-z})$$

$$\cdot (0 + -e^{-z})$$

$$= e^{-z}/g^2 = \frac{e^{-z}}{(1+e^{-z})^2}$$

# Tanh function

$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

☞ $\dfrac{d\tanh}{dz} = 1 - \tanh^2(z)$



→ tanh lies to $-1$ to $1$

→ $\quad 0 \leq \dfrac{\tanh}{dz} \leq 1$

47.12    Vanishing Gradient Descent ⟵

→ { 80's, 90's, 00's } → nvn faced a problem called

→ Typically Ppl used Sigmoid.

To update a weight $\left(\textcircled{$\omega^1_{11}$}\right)$

$$\left(\omega^1_{11}\right)_{new} = \left(\omega^1_{11}\right)_{old} - \eta \left(\frac{\partial L}{\partial \omega^1_{11}}\right)$$

$$\frac{\partial L}{\partial \omega^1_{11}} = \frac{\partial L}{\partial O31} \left[ \frac{\partial O31}{\partial O21} \cdot \frac{\partial O21}{\partial O11} \cdot \frac{\partial O11}{\partial \omega^1_{11}} + \cancel{\partial O31} \right.$$

$$\left. \frac{\partial O31}{\partial O22} \cdot \frac{\partial O22}{\partial O11} \cdot \frac{\partial O11}{\partial \omega^1_{11}} \right]$$

$$\frac{\partial O31}{\partial O21} = \frac{\partial f_{31}}{\partial O21}$$

$$0 \leq \overbrace{\frac{\partial (f_{31})}{\partial O2}}^{\sigma} < 1$$

Vanishing grad. → V.V. small.

Exploding grad → VV large.

47-13    Bias - Variance Tradeoff

① # layers ↑ ⟹ more weights/params
                    ⇓
            higher chance of overfitting
                    ⇓
            high Variance.

② 1-layer = Log.Res + $\frac{t}{1}$ → higher chance of underfit
                                         ⇓
                                    high bias.

→ Typically there is a higher chance of overfitting
   ↳ we can restrict by adding a regularization parameter

$$\mathcal{L} = \sum_{i=1}^{n} loss(y_i, \hat{y}_i) + \sum_{i,j,k} (\omega_{ij}^k)^v$$

we can either L1 & L2 regularization.

$$\mathcal{L} = \sum_{i=1}^{n} loss_i + \lambda \; reg \; on \; weight$$

        larger $\lambda$ ⟹ lesser overfit

L1 reg ⟹ sparsity  ⟹ some $\omega_{ij}^k = 0$
                                ⇓
                        MLP sparse

→ here '$\lambda$' is a hyper parameter  } hyper parameters
→ # of layers ↑ ⟹ var ↑

47.14    Decision surface: Play Ground

playground. tensorflow.org / # activation = tanh & batchsize = 10 - - - —————