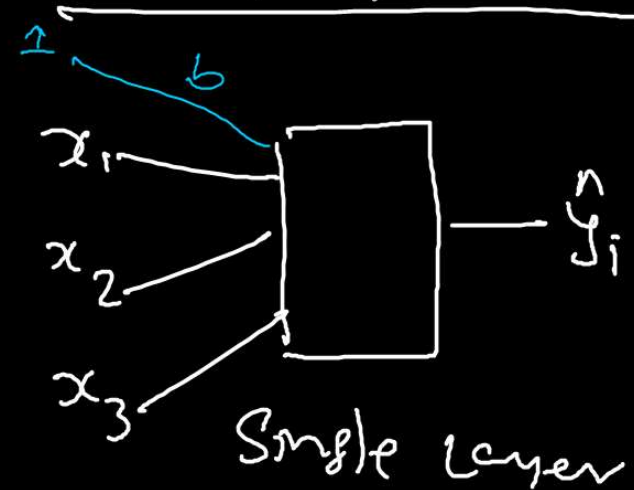# Deep Neural Net~~

—— * —— * ——

Agenda :-
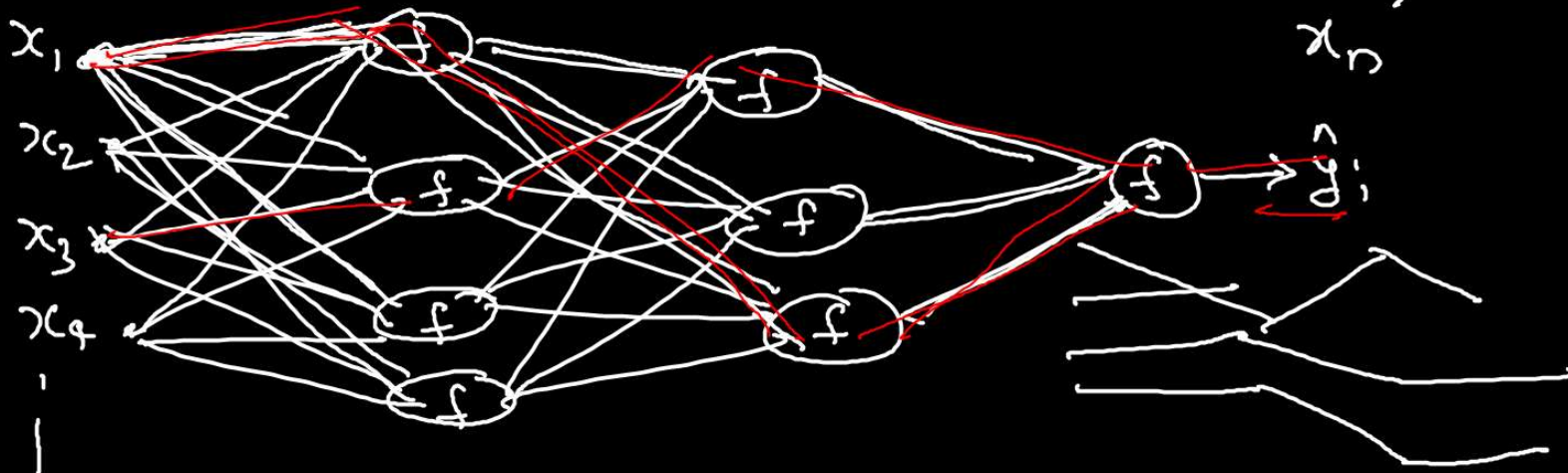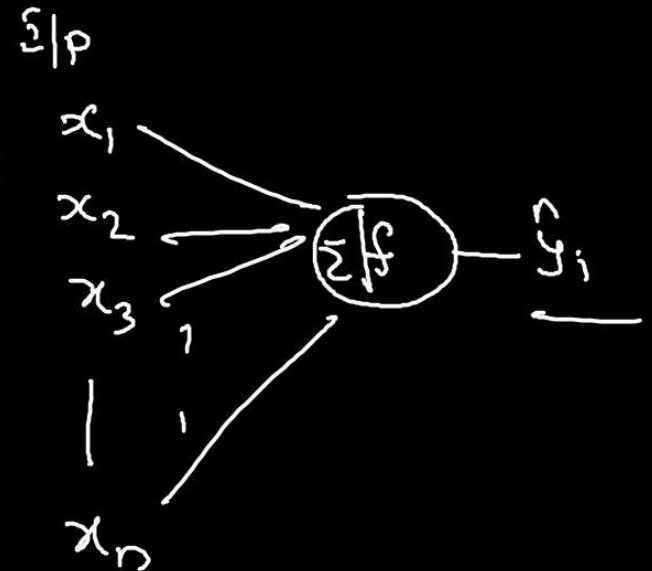
① Notation

② Backpropogation method

③ Data Normalization

④ early stopping

⑤ case study

⑥ Dropout

## Multilayer perceptron



Single Layer

$\hat{y}_i$

# Multilayer Perceptron

## Perceptron — Single Neuron →

$i/p$

$x_1$

$x_2$

$x_3$

$1$

$1$

$x_n$

$\Sigma | f$ — $\hat{y}_i$

$x_1$

$x_2$

$x_3$

$x_4$

$1$

$f$

$f$

$f$

$f$

$f$

$f$

$f$

$\hat{y}_i$

Q:- Why should we care about MLP ? -Inter-

① Biological inspiration etc

    Neuroscience → human , rets , ants , monkeys etc .
    -Neural Structure Network -

② Mathematical Arguments :-

    regression :- $\{x_i , y_i\}$ , $\boxed{x_i \in \mathbb{R}^n , y \in \mathbb{R}}$ → Real num

    eg :- $\boxed{2 * \sin(x^2) + \text{sqrt}(x * 5)}$

    LR :-
    $y = m_1 x_1 + c + m_2 x_2$ thus

Plot $\left( 2 * \sin(x^2) + \sqrt{5x} \right)$

$y = f(x)$

| $x_1$ | $y$ |
|-------|-----|
| 1 | 1 |
| 2 | 4 |
| 3 | 9 |

square(x)

$y = f(x)$

NN

Let,

$f_1 = add()$

$f_2 = square()$

$f_3 = sqrt()$

$f_4 = sin()$

$f_5 = multi()$

MLP - structure

I/p

x

5

2

multi → $5x$ → sqrt → $\sqrt{5x}$

square → $x^2$ → sin → $\sin(x^2)$

multi → $2 * \sin(x^2)$ → add

$2 * \sin(x^2) + \sqrt{5x}$

$\boxed{\text{high - School}}$ :- Functional composition ✓

$$f(g(x)) = f \circ g(x) \qquad g \circ f(x) \to g(f(x))$$

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

$f1 = add$
$f2 = square$
$\underline{f3 = \text{sqrt}}$
$f4 = \sin$
$\quad \dot{f5} = \text{multi}$

$$x * 5 = f5(x, 5)$$

$$\underline{\text{sqrt}(x * 5)} = f3\left(f5(x, 5)\right) \qquad \underline{2 * \sin(x^2) + \text{sqrt}(x * 5)}$$

$$x^2 = f2(x)$$
$$\sin(x^2) = f4(f2(x)) \qquad = f1\left(f5(2 f4(f2(x))), f3(f5(x, 5))\right)$$

$$2 * \sin(x^2) = f5\left(2, f4(f2(x))\right)$$

$\underline{\text{MLP}}$ - Graphical way of representation $(f \circ g, g \circ f)$ $\underline{\text{Functi}}$
$\boxed{\text{Compu}}$

MLP is very powerful Structure which can handle linear or non-linear or both by using DNN/MLP connect

* <u>Notation</u>

$\mathbb{R}$ - Real Num

$$D = \{x_i, y_i\} \quad , x_i \in \mathbb{R}^q \quad , y_i \in \mathbb{R} \quad - \text{Regression}$$

$f_{ij}$ → Feature
Index
Layer

$O_{ij}$ → output
Index
Layer

$\overset{u}{\underset{From}{\longrightarrow}}$ To
$W$ = weight

Input Layer

Hidden Layer I

Hidden Layer II

Output Layer

$X_{ij}$ → Input
Point → Feature

$X_{i1}$

$W^1_{11}$  $W^1_{12}$  $W^1_{21}$  $W^1_{22}$  $W^1_{31}$  $W^1_{32}$  $W^1_{33}$  $W^1_{42}$  $W^1_{43}$

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

$W^2_{11}$  $W^2_{12}$

$X_{i2}$

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

$W^2_{21}$  $W^2_{22}$  $W^2_{31}$  $W^2_{32}$

$f_{2i}$  $O_{21}$

$W^3_{11}$

$f_{22}$  $O_{22}$

$W^3_{21}$

$f_{31}$  $O_{31}$  $\hat{y_i}$

$X_{i3}$

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

$X_{i4}$

basic flow in Neural network

③  ②  ①

④  ③

$W^1 = 4 \times 3 = 12$

$W^2 = 3 \times 2 = 6$

$W^3 = 2 \times 1 = 2$

Total $W = 12 + 6 + 2$
$= 20$

Bias $= \dfrac{6}{26}$

$$W_1 = \begin{bmatrix} W_{11}^1 & W_{12}^1 & W_{13}^1 \\ W_{21}^1 & W_{22}^1 & W_{23}^1 \\ W_{31}^1 & W_{32}^1 & W_{33}^1 \\ W_{41}^1 & W_{42}^1 & W_{43}^1 \end{bmatrix}$$

$4 \times 3$

how weight calculation —

Weight

calculation —

copper

Xevier/Glorut

$1/e \cdots 1$

$$W_3 = \begin{bmatrix} W_{11}^3 \\ W_{21}^3 \end{bmatrix}$$

$2 \times 1$

$$W_2 = \begin{bmatrix} W_{11}^2 & W_{12}^2 \\ W_{21}^2 & W_{22}^2 \\ W_{31}^2 & W_{32}^2 \end{bmatrix}$$

$3 \times 2$

Mute | Start Video | Security | Participants 16 | Chat | New Share | Pause Share | Annotate | Remote Control | Apps | More

You are screen sharing | Stop Share

$$\dfrac{10}{SC_1}$$

$$\boxed{100}$$

$$f_{11} \dfrac{10}{0}$$

$$\boxed{110}$$

$$\dfrac{11}{0}$$

3

$$\boxed{110}$$

$$\dfrac{10}{0}$$

$$x_2$$

$$x_3$$

$$x_4 \quad W_{53}^1$$

$$f_{12} \ 0$$

$$f_{13} \ 0$$

$$W_{41}$$

$$f_{23} \ \boxed{0}$$

6

10

100

110

110

$$\dfrac{10}{}$$

$$x_5$$

$$0$$

$$0$$

$$0$$

$$6$$

$$6$$

1

0

$$TW = 330$$

$$x_6$$

$$f_{16} \ 0$$

$$W_{67}^2$$

$$0$$

$$0$$

$$f_{35} \boxed{0}$$

$$B = 32$$

$$x_7$$

$$0$$

$$0$$

$$0$$

$$6$$

$$362$$

$$x_8$$

$$f_{19} \ 0$$

$$0$$

$$6$$

$$x_9$$

$$f_{110} \ 0$$

$$0$$

$$0$$

$$x_{10}$$

$$0$$

$$0$$

Feature

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $\hat{y}_i$ |
|---|---|---|---|---|
| 1 | 2 | 3 | 4 | Yes |
| 1 | 2 | 5 | 6 | NO |
| 5 | 3 | 2 | 1 | NO |
| 10 | 20 | 30 | 40 | Yes |

i=1  1 2 3 4 = Yes

i=2  1 2 5 6 = NO

$x_{ij} \rightarrow 3$

1

$x_{i1}$
$x_{i2}$
$x_{i3}$
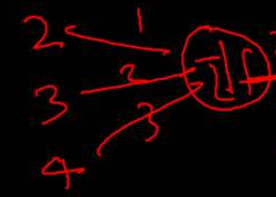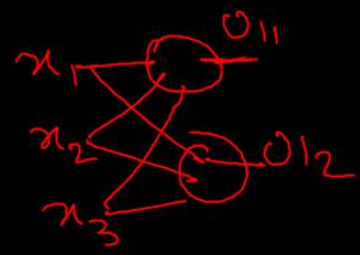$i \leq 4 \quad x_{i4}$

$x_{11}$

$x_i - i = point$

$f = i$

$x_1$ — $O_{11}$
$x_2$ — $O_{12}$
$x_3$

(1) — $y_i^h$

$2x_1 + 3x_2 + 4x_3$

$= 2 + 6 + 12$
$= 20$

2 1
3 2     20 = 20
4

$z = 20$

$f(20) = 20$

# Foreward Propogation method



$$10 \times 2 +$$
$$20 \times 1 +$$
$$30 \times 2$$
$$= 20 + 20$$
$$+ 60$$
$$100$$

$$\begin{array}{|c|c|c|c|} \hline x_1 & x_2 & x_3 & y_i \\ \hline 10 & 20 & 30 & 100 \\ \hline \end{array}$$

Step Act

$$Z = \sum x_i w_i$$

Step $$= (10 * 1 + 20 * 2 + 30 \times 3)$$

$$w_1 = 1 \quad \boxed{Z = 140}$$

$$w_2 = 2 \quad \underset{\text{Activate}}{f(Z)} = f(140) = 140 = \hat{y}_i$$

$$w_3 = 3 \quad \downarrow$$
$$\text{linear}$$

$$y_i = 140 \quad -40 = Loss = Zero$$

## Error / Loss / Cost

$$y_i - \hat{y}_i = -40$$

Interview



— 2006 — Prof Geoffrey Hinton 20yrs

Backpropogation Algorithm ← Chain Rule &
Memoization &

Forward Prop

$\dfrac{\partial L}{\partial w_{31}} \ne \dfrac{\partial O_{31}}{\partial w_{11}^3}$

1.5

I/P

$x_{i1}$  $W_{11}^1$  1  $W_{11}^2$  $f_{11}$  $O_{11}$  $W_{21}^2$  $f_{21}$  $O_{21}$  $W_{11}^3$  O/P

$x_{i2}$  $f_{12}$  $O_{12}$  $f_{31}$  $O_{31}$  $\hat{y}_i$  = Loss

$x_{i3}$  $f_{22}$  $O_{22}$  $W_{21}^3$

$x_{i4}$  $f_{13}$  $O_{13}$  $W_{32}^2$

$W_{43}$

Epochs = 5000

1 Epochs = 1 BW + 1 FB

↓
forward iteration

# Chain Rule

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \cdot \frac{\partial u}{\partial v} \cdot \frac{\partial v}{\partial x}$$

$$\frac{\partial L}{\partial w_{11}^3} = \boxed{\frac{\partial L}{\partial o_{31}}} * \frac{\partial o_{31}}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{21}^3} = \boxed{\frac{\partial L}{\partial o_{31}}} * \frac{\partial o_{31}}{\partial w_{21}^3}$$

$$\frac{\partial L}{\partial w_{11}^2} = \boxed{\frac{\partial L}{\partial o_{31}}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{32}^2} = \boxed{\frac{\partial L}{\partial o_{31}}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial w_{32}^2}$$

$$\frac{\partial L}{\partial w_{11}'} = \boxed{\frac{\partial L}{\partial o_{31}}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \boxed{\frac{\partial o_{11}}{\partial w_{11}'}} +$$

$$\boxed{\frac{\partial L}{\partial o_{31}}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{11}} * \boxed{\frac{\partial o_{11}}{\partial w_{11}'}}$$

$$\frac{\partial L}{\partial w_{11}'} = \boxed{\frac{\partial L}{\partial o_{31}}} * \frac{\partial o_{11}}{\partial w_{11}'} \left\{ \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} + \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{11}} \right.$$

$$x \xrightarrow{f(x)} h \left($$
$$x \xrightarrow{g(x)} h \right.$$

$$\frac{\partial L}{\partial x} = \frac{\partial h}{\partial f} * \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} * \frac{\partial g}{\partial x}$$

## Memoization :

it is not memorization 'it is called

memoization $\longrightarrow$ Dynamic Programming

computer Science

$$\frac{\partial L}{\partial w^3_{11}} = \boxed{\frac{\partial L}{\partial O_{31}}} * \frac{\partial O_{31}}{\partial w^3_{11}}$$

$$\frac{\partial L}{\partial w^3_{21}} = \boxed{\frac{\partial L}{\partial O_{31}}} * \frac{\partial O_{31}}{\partial w^3_{21}}$$

Compute once & reuse it.

* The core idea of memoization is this

$\longrightarrow$ if there is any operation that is used many times, it's good idea to compute once & then store it for reuse Purpose

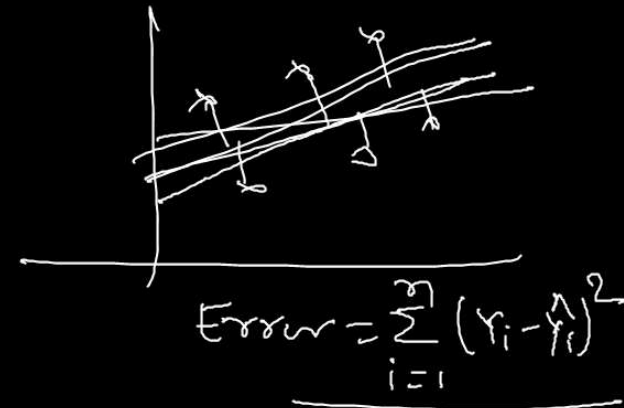$\longrightarrow$ This concept is known as memoization

① Define Loss Function (example - Regression)

$$Y = m_1 x_1 + m_2 x_2 + c$$

$$L = \sum_{i=1}^{n} \left( y_i - \hat{y}_i \right)^2 + regularization$$

$$\underbrace{\qquad\qquad\qquad}_{}$$

Sq-loss fun:



$$Error = \sum_{i=1}^{n} \left( y_i - \hat{y}_i \right)^2$$

② Optimization

Loss

(minimize)

↓

Weight & bias

Mute | Start Video | Security | Participants 17 | Chat | New Share | Pause Share | Annotate | Remote Control | Apps | More

You are screen sharing | Stop Share

③ Gradient Descent ( Batch GD, Mini-batch GD, SGD )

(i) initialize $W_{ij}^{k}$ randomly — lots of techniqu → we will cover later

(ii) update weight — Rule :-

Formula:

$$\left(W_{ij}^{k}\right)_{new} = \left(W_{ij}^{k}\right)_{old} - \eta * \frac{\partial L}{\partial W_{old}}$$

$\frac{\partial L}{\partial W_{old}}$ → V.V.Imp BACKPROP.

→ eta/ learning Rate (0 to 10)

$W_{old} = 50$

$W_{new} = 50 - 1 * 5$

$W_{new} = 45$

(ii) Perform updates till the <u>convergence</u>.

<u>NOTE</u> :- convergence means the new weight value and the old weight value are very close to each other.

$$W_{old} = 50$$
$$W_{new} = 49.9995$$
$$W_{new_2} = 49.9998$$
$$W_{new_3} = 49.9997$$

2.45

# Back-Prop

① Initialize weight - randomly

② epochs $\longrightarrow$ Forward Prop

(1BP+1FB) $\longrightarrow$ compute loss

$\longrightarrow$ Compute derivatives (Chain Rule + memoized)

$\longrightarrow$ update weight in a backward prop

③ repeat Step 2 till the convergence.