

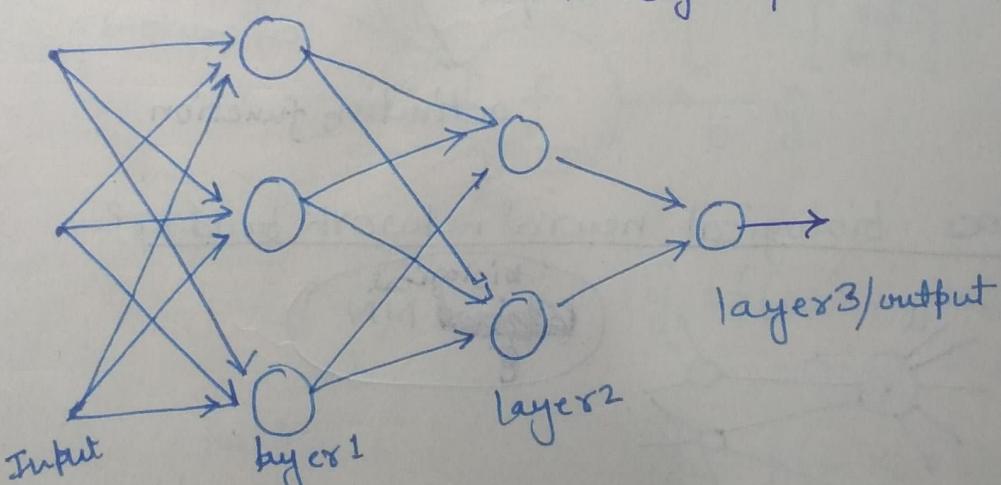
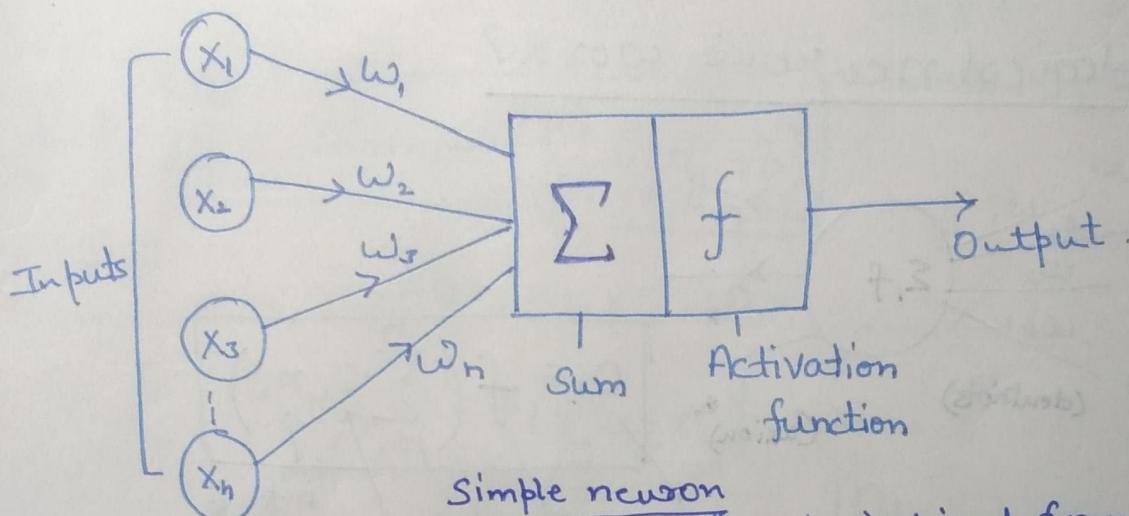
# "DEEP LEARNING"

## Neural Nets.

### ① History of Neural Networks.

→ "Perception" → 1957 [Rosenblatt]

similar to "Logistic regression"



Around 1986

↳ "Back propagation" → chain rule of diff

1990+ AI funding died called "AI winter".

2006 → "Hinton" paper came

2012 → ImageNet (large data) @ stanford.

Image dataset.

Image → object?

DNN model (computation)

"huge margin"

→ this DNN beat every algorithm by huge margin.

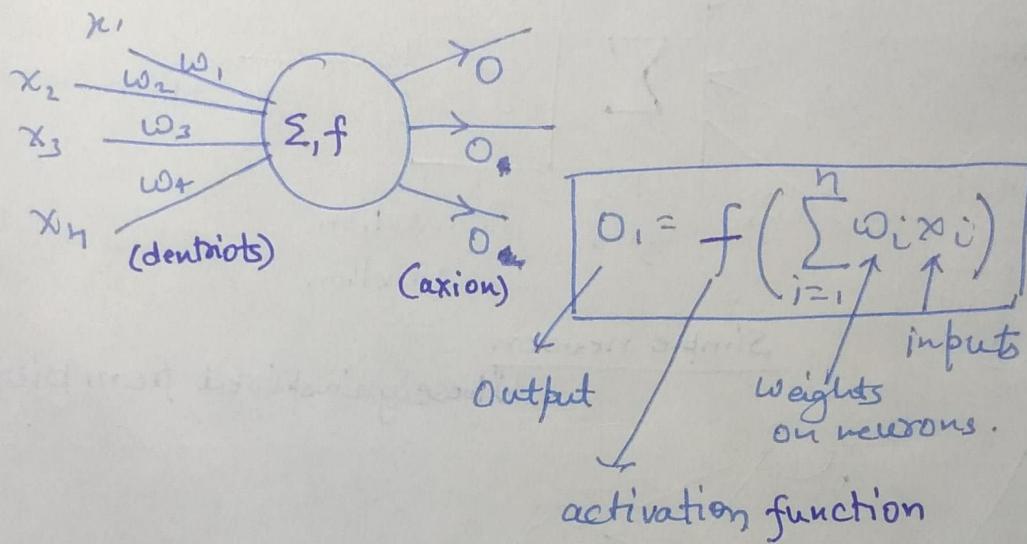
then funding started by Google, Ms, F b, baidu, Amazon and golden days of AI started.

DL - CNN  $\rightarrow$  1990's

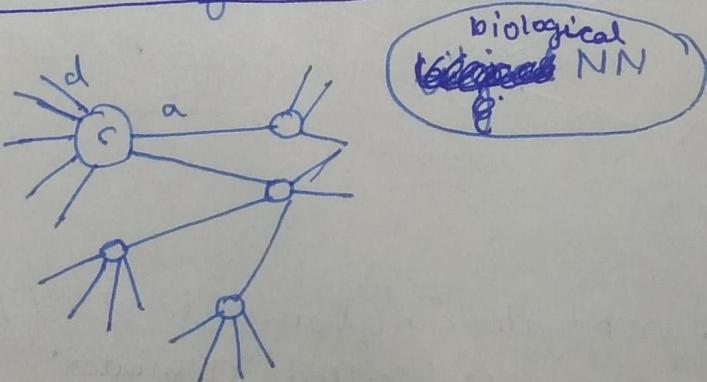
but popular due to  
① Large dataset.  
② Huge Computational powers  
③ New Algorithms.

NN  $\rightarrow$  "classical" 1986  
MLP  $\rightarrow$  deep NN 2006  
CNN  $\rightarrow$   
RNN  $\rightarrow$

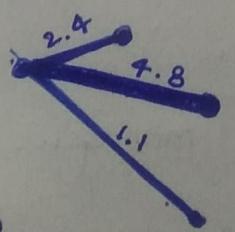
## ② How biological neurons work?



## ③ How does biological neural network build up?

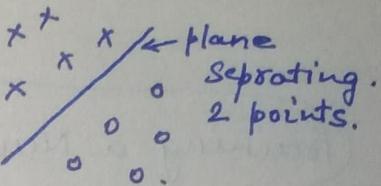


- \* learning  $\Rightarrow$  connecting neurons with edges.
- all biological learning  $\Rightarrow$  weights on neural connections.



- Artificial Neural network is just connection of neural nets and weights on neural connections.

#### 4) Logistic Regression & Perceptron



LR:  $x_i \rightarrow \hat{y}_i \rightarrow$  predicted value of  $y_i$ .

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

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

↓  
Data

Train L.R.  $\Rightarrow \omega, b$  ← weight & bias  
are obtained.  
 $\omega \in \mathbb{R}^d$   
 $b \in \mathbb{R}$  by training.

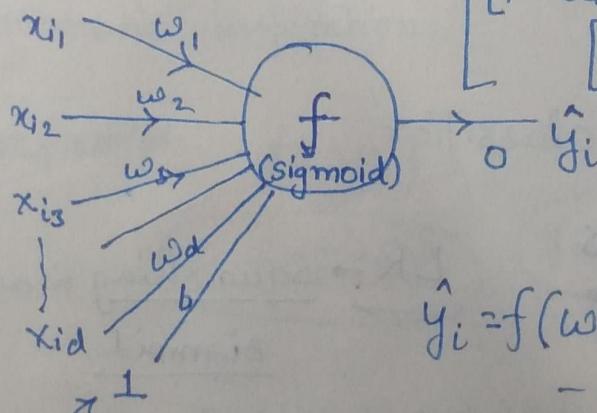
L.R.

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

$$x_{ij} = [x_{i1}, x_{i2}, x_{i3}, x_{i4}, \dots, x_{id}]$$

$$O = f\left(\sum_{j=1}^d w_j x_{ij}\right)$$

Output



[Neuron interpretation]  
[of LR]

$$\hat{y}_i = f(w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} + \dots + w_d x_{id} + 1 \cdot b)$$

for bias

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

if f is sigmoid function then NN is logistic regression.

In L.R.

Task was from ~~D~~  $D = \{x_i, y_i\}$

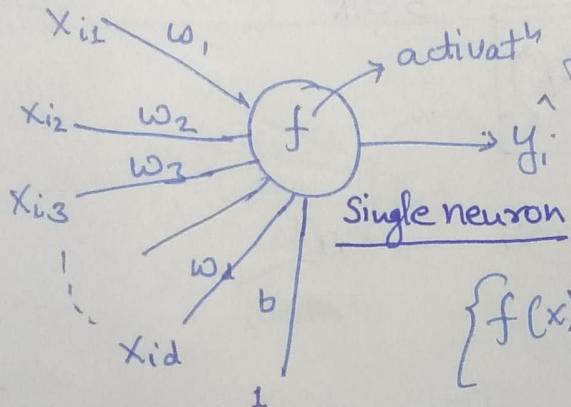
to find Task =  $w_i, b? \forall i = 1 \dots d$

SGD → L.R. regression & optimization

Training a Neural Net  $\rightarrow$  Weights on edges/vertices.

so for LR  $\rightarrow$  Neuron was made  
& Sigmoid was Activation function

Perception  $\rightarrow$  1957

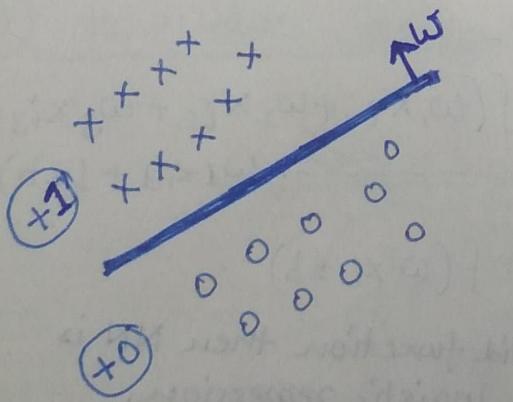


$$\omega^T x = [w_1, w_2, w_3, \dots, w_d] \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ \vdots \\ x_{id} \end{bmatrix}$$

$$= \sum_{j=1}^d w_j x_{ij} + b$$

$$\left\{ \begin{array}{ll} f(x) = 1 & \text{if } \omega^T x + b > 0 \\ = 0 & \text{otherwise} \end{array} \right.$$

Perception  $\rightarrow$  Lr. classifier



LR  $\rightarrow$  Squashing was done ~~done~~ by Sigmoid.

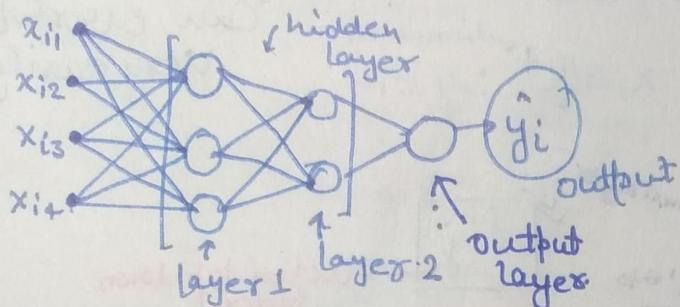
Perception  $\rightarrow$  Activation function.

(5)

## MULTI-LAYERED PROGRAM PERCEPTRON (MLP)

Perception  $\rightarrow$  single neuron

$\hookrightarrow$  similar to logistic regression



(Q) why should we care about (MLP)

(a) biological inspiration

Neuroscience  $\rightarrow$  human, rats, monkeys.

(b) Mathematic inspiration

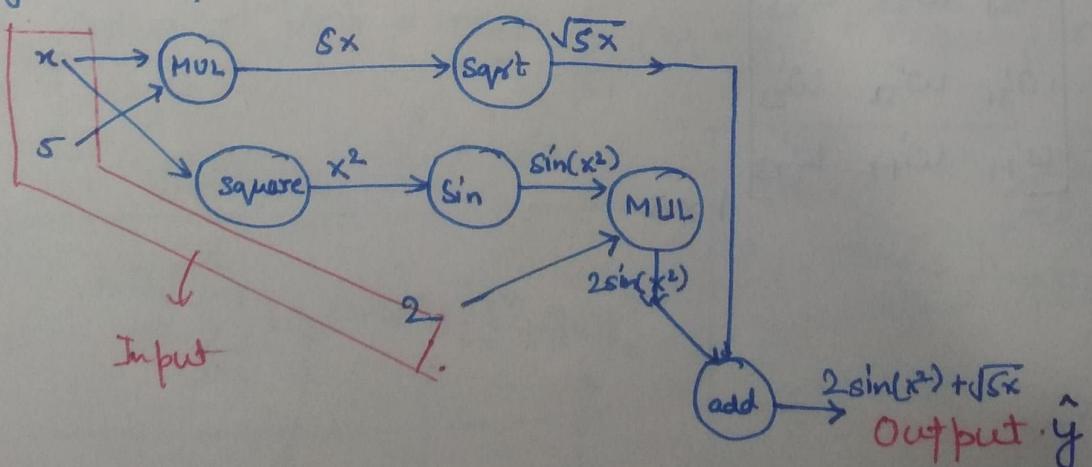
of regression  $\{x_i, y_i\} = D$

$$y_i = F(x_i)$$

$$y_i \in \mathbb{R}$$

blocks

Q. Say a complex function  $\Rightarrow 2 \sin(x^2) + \sqrt{5x}$



So complex problem can be solved easily.

This is MLP-like structure.

Multi-layered  $\Rightarrow$  enormous power

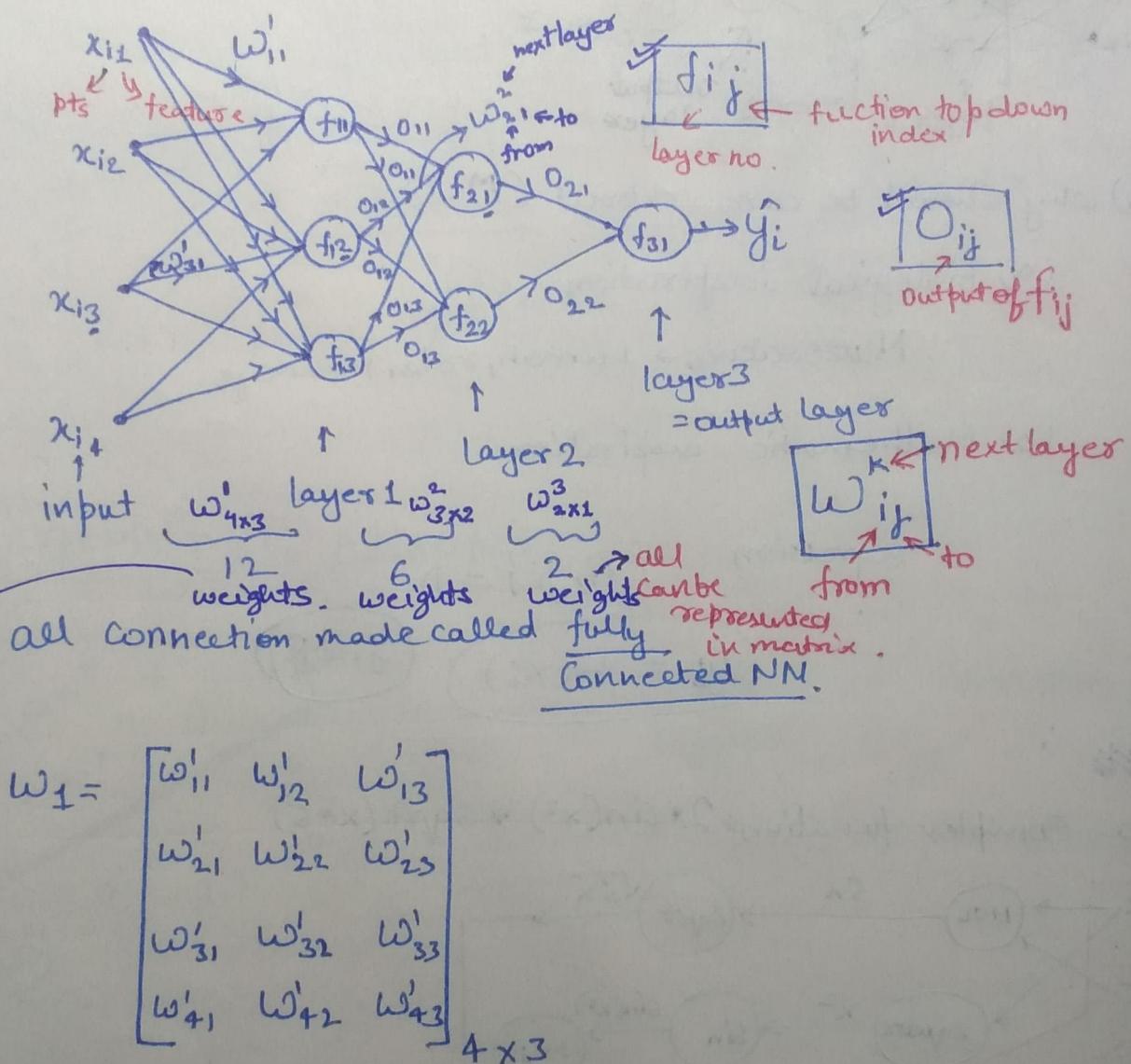
Similar concept of function composition:  $g(f(x))$  or  $f(g(x))$

M LS (Multilayer structure)  $\rightarrow$  easily overfit.

$\hookrightarrow$  powerful model and powerful models  
can overfit  
Very easily.

## ⑥ Notation

$$\mathcal{D} = \{x_i, y_i\} : x_i \in \mathbb{R}^4 \xleftarrow{\text{dimension}} ; y_i \in \mathbb{R}$$

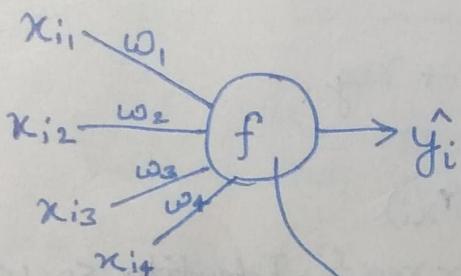


## 7) Train a single Neuron model

find the best edge weights using  $D_{Tr.} \leftarrow$  training Data.

perception & LR  $\rightarrow$  single neuron models for classification  
 we use regression L.R. Regression  $\rightarrow$  " " " for regression  
 because maths is easy.

(En Reg) & (optim.)  $\rightarrow$  SGD  
 linear



L.R. Reg.

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

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

$y_i \in \mathbb{R}$

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

in this case

$$\begin{cases} \text{Input} = w^T x_i \\ \text{Output} = w^T x_i \end{cases} \rightarrow$$

f is identity function.

for this neuron to  
be working as L.R regression.

In case of Log. Reg.

f = logistic function.

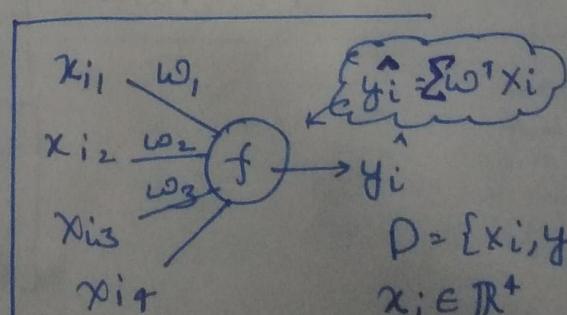
## L.R. Regression

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

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

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

steps,  
are given  
in next page.



How to solve using  
neural nets?  
in short how to find.

different weights?

### ① Define Loss function

loss  $\rightarrow L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg.}$   
 across all  
 the data

loss  $\leftarrow l_i = (y_i - \hat{y}_i)^2 \Rightarrow \hat{y}_i = \omega^T x_i$   
 across just  
 training point &  
 actually  $L = \sum_{i=1}^n l_i$   
 sum loss of  
 across all points.

### ② Optimization problem

$$L_r \underset{\text{reg.}}{\rightarrow} \min_{\omega_i} \sum_{i=1}^n (y_i - \underbrace{\hat{y}_i}_{\hat{y}_i = f(\omega^T x_i)})^2 + \text{reg.}$$

$\hookrightarrow L_r$  reg is Identity, so we write  $\omega^T x_i$  as it is.

Neuron  $\rightarrow \min_{\omega_i} \sum_{i=1}^n [y_i - f(\omega^T x_i)]^2 + \text{reg}$

$\hookrightarrow L_r \rightarrow$  reg

$\hookrightarrow$  sigmoid  $\rightarrow$  log-regression

$w^* = \arg \min_w \sum_{i=1}^n (y_i - f(\omega^T x_i))^2 + \text{reg}$  general optimization problem of neuron

### ③ Solve the optimization problem.

a) Initialize of  $\omega_i$ 's  $\rightarrow$  random

b) Compute  $\nabla_{\omega} L =$

$$\begin{bmatrix} \frac{\partial L}{\partial \omega_1} \\ \frac{\partial L}{\partial \omega_2} \\ \vdots \\ \frac{\partial L}{\partial \omega_n} \end{bmatrix}$$

Watch  
optimization  
Chapter for  
understanding  
properly.

c)  $\omega_{\text{new}} = \omega_{\text{old}} - \eta [\nabla_{\omega} L]$

$\omega_{\text{old}}$  learning rate

$(\omega_i)_{\text{new}} = (\omega_i)_{\text{old}} - \eta \left[ \frac{\partial L}{\partial \omega_i} \right] (\omega_i)_{\text{old}}$  we derive next

update weight with respect to  $(\omega_i)_{\text{old}}$

G for iter 1 — K till converge.

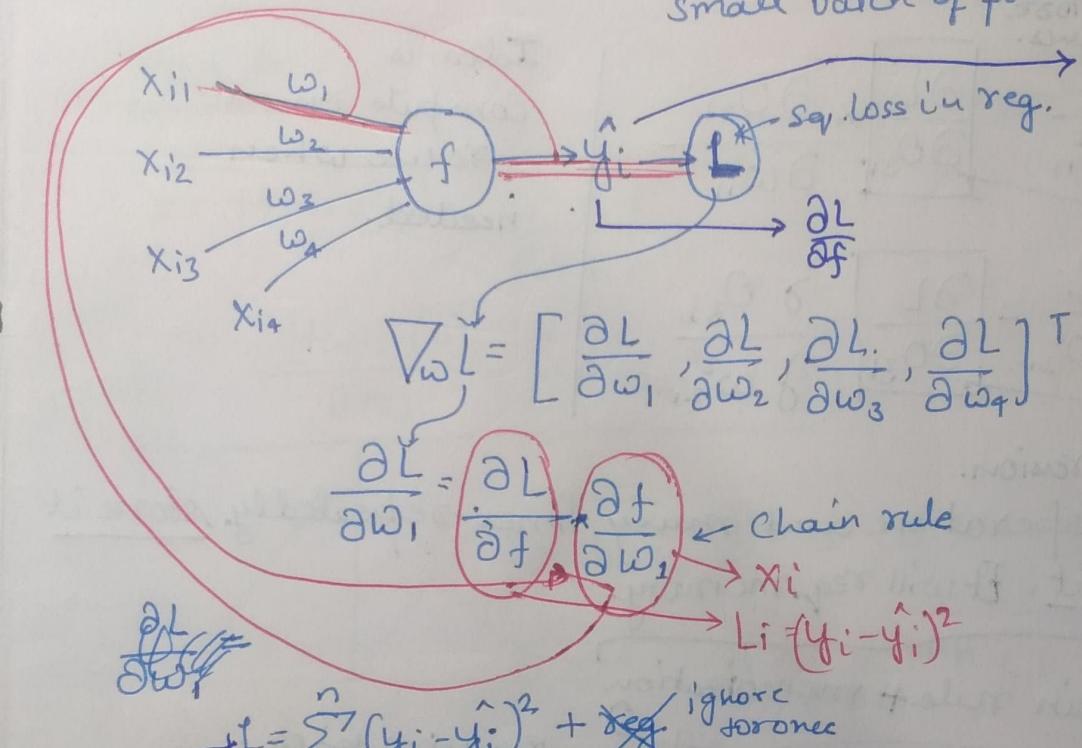
GD :-

Gradient descent:  $\nabla_{\omega} L \rightarrow x_i's \text{ & } y_i's$  In SGD we do this for small batch of points or individual min-batch point one by one.

SGD:  $\nabla_{\omega} L \approx \text{one pt } \{x_i, y_i\}$

small batch of points  $\xleftarrow{\text{batch}} S_{\text{GD}}$

Now, this  $y$  predicted & actual may have some error based on that we define Sq loss.



$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \cancel{\text{ignore for once}}$$

$$\frac{dL}{d\omega_1} \text{ then } \frac{dL}{d\omega_1} = \sum_{i=1}^n (y_i - f(\omega^T x_i))^2$$

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

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

$$\frac{\partial L_i}{\partial \omega_1} = -2x_i(y_i - \hat{y}_i)$$

$$\frac{\partial L}{\partial \omega_1} = \sum_{i=1}^n -2x_i(y_i - \hat{y}_i)$$

We are given with data.

We know actual  $y_i$  with respect to old weights.

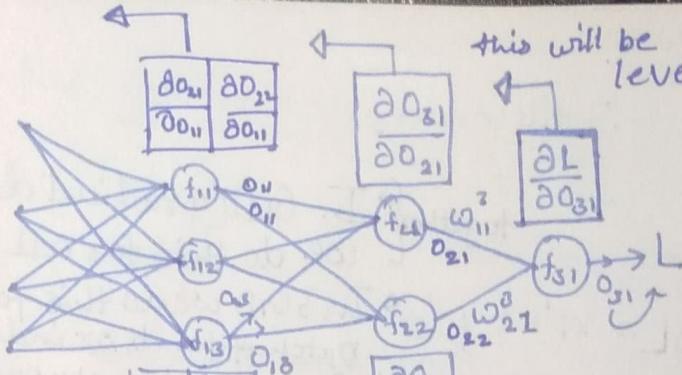
We get this with respect to old weights.

8) Train an MLP chain rule.  
9) Memoization :-  $\xrightarrow{\text{go to last.}}$

Algorithm  $\rightarrow$  Dynamic Programming

Not memorisation

Storing the same expression value rather calculating derivatives again & again.



Here.  
2 outputs →  
are contrib. update  $w_{11}$  &  $w_{21}$   
to this value so we store  
2 values.

$$w^3 = \frac{\partial L}{\partial w^3} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial w^3}$$

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

Same expression.

Idea is  
compute once &  
reuse when  
needed.

if any operation used many times repeatedly store it and use it. It will req. memory.

chain rule + memoization

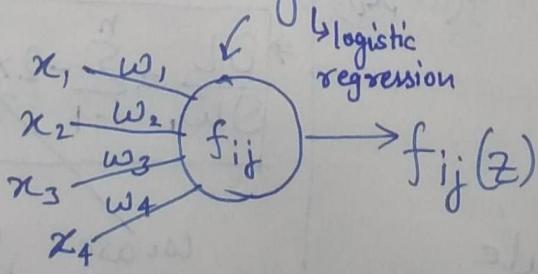
⑩ Backpropagation → go to last

⑪ Activation functions ( $f_{ij}$ )

Backpropagation  
Algorithm.

1980 & 1990's → Sigmoid & tanh. ← 2 of the oldest activation function.

easily differentiable.



$$\sum_i w_i x_i = w^T x$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

Sigmoid.

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

easily differentiable.

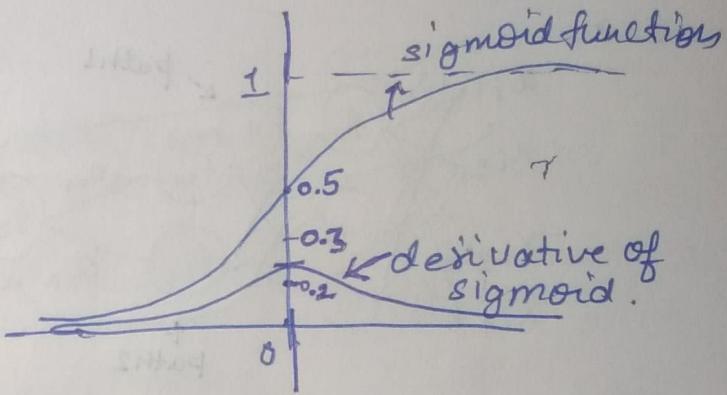
\* this seems difficult to differentiate but  
it is just  $\sigma(z)[1-\sigma(z)]$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

## Req. of Activation

- 1) differentiable
- 2) easy to differentiate

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



Sigmoid diff come in term of sigmoid only.

tanh :-

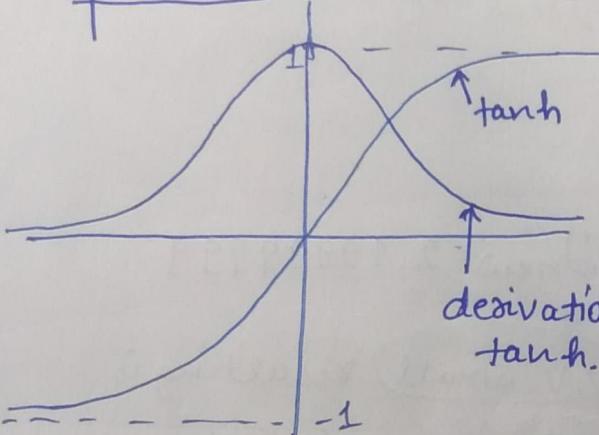
for proof do it yourself.

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

$$\frac{\partial [\tanh(z)]}{\partial z} = [1 - \tanh^2(z)]$$

$$\frac{\partial a}{\partial z} = 1 - a^2; a \rightarrow \tanh \text{ function.}$$

easily differentiated.



$$\textcircled{1} * -1 \leq z \leq 1 \leftarrow \tanh$$

$$\textcircled{2} 0 \leq \frac{\partial \tanh}{\partial z} \leq 1$$

\* One has to note shape of Sigmoid & tanh is same it is just the range.  
Sigmoid  $\rightarrow (0, 1)$

Modern techniques like ReLU: tanh  $\rightarrow [-1, 1]$   
are very popular.

② Leaky ReLU.

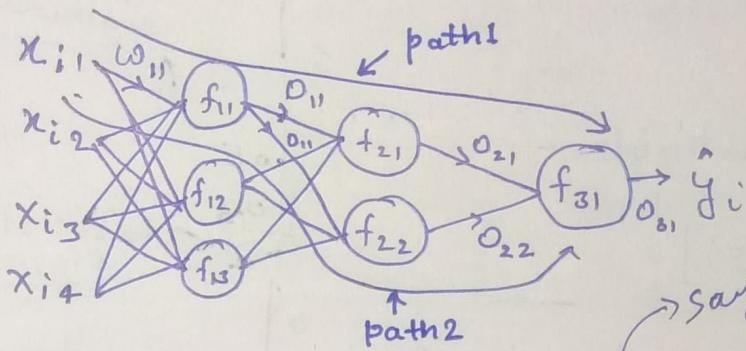
12

Vanishing gradients

80, 90, 00's  $\rightarrow$  NN problem.  
 $f_{ij} \rightarrow$  Sigmoid.

Most important reason why people lost interest in deep-learning because of problem of vanishing gradient.

Ex-



$$(\omega'')_{\text{new}} = (\omega'')_{\text{old}} - n \frac{\partial L}{\partial \omega''}$$

vv small

$$\frac{\partial L}{\partial \omega''} = \frac{\partial L}{\partial O_{31}} \left[ \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial \omega''} + \frac{\partial O_{31}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial \omega''} \right]$$

$\Rightarrow O_{31} \approx$   
 Output of sigmoid  
 function & these  
 diff is b/w 0.25.  
 $< 10^{-3}$

If  $(\omega'')_{\text{old}} = 2.5$   
 $(\omega'')_{\text{new}} = 2.499$

Say if deep layers.

then

$$(\omega'')_{\text{old}} = 2.5 \quad (\omega'')_{\text{new}} = 2.49999999$$

① Vanishing gradient  $\frac{\partial L}{\partial \omega_{ij}^k} = \text{v.v. small}$  resulting in  
 solution Relu.

because derivative of Sigmoid/tanh both less than 1.

Updating take many many iteration.

so people were only training 2 hidden layers.

to build deep NN  $\begin{cases} 10 \\ 20 \\ 100 \end{cases} \rightarrow$  then  $\frac{\partial L}{\partial \omega_{ij}^k}$  would be of  $10^{-10}$  or  $10^{-20}$  or  $10^{-100}$ .

② Exploding gradient this is the reason we cannot add deep neural network or large no. of layers.

\* although we don't have it  $\frac{\partial L}{\partial \omega''}$  is greater or v.v. large

In this case  $(\omega'')_{\text{new}} = (\omega'')_{\text{old}} - n \frac{\partial L}{\partial \omega''}$  then  $(\omega'')_{\text{new}}$  will never converge and go crazy.

\* even after 100 epoch we will never converge.

large derivative is even bigger problem than small derivation.

### (13) Bias-Variance Trade off

① # layers  $\uparrow \Rightarrow$  more weights/params  
 generally we use ↓  
 more layers and have Chances of overfitting  $\leftarrow$  why?  
 problem of high variance ↓  
 (overfitting) high variance

ex: 1 layer = Log. Reg.  $\rightarrow$  + + → higher chance  
 with sigmoid activation function. of underfit.

To solve overfit/underfit add regularizer:- underfit

MLP  $\xrightarrow{?}$  multiple layers  $\rightarrow$  overfitting/high var

optimization  $\rightarrow$  add regularization

$$L = \sum_{i=1}^n \text{loss}(y_i, \hat{y}_i) + \sum_{i,j,k} (w_{i,j}^k)^2$$

L<sup>2</sup> regularizer

this regularization helps to avoid overfitting. (high variance).

$\Rightarrow$  L<sub>2</sub> regularization or L<sub>1</sub> regularization

$$L = [\text{loss}] + \lambda L_2$$

reg on weights  $\rightarrow$  we just add  $\lambda$  to give more weightage to regularizer.

problem with L<sub>1</sub> reg. larger  $\lambda \Rightarrow$  lesser overfit  
 L<sub>1</sub> reg  $\Rightarrow$  sparsity  $\Rightarrow$  some weights becomes 0.  $\downarrow$  visit old lectures for clarity.  
 fewer edges become non-zero.

Then in L<sub>2</sub> reg.  $\lambda$  become hyperparameter.

Second hyperparameter is no. of layers.

now, we have 2 hyperparameters ~~more~~

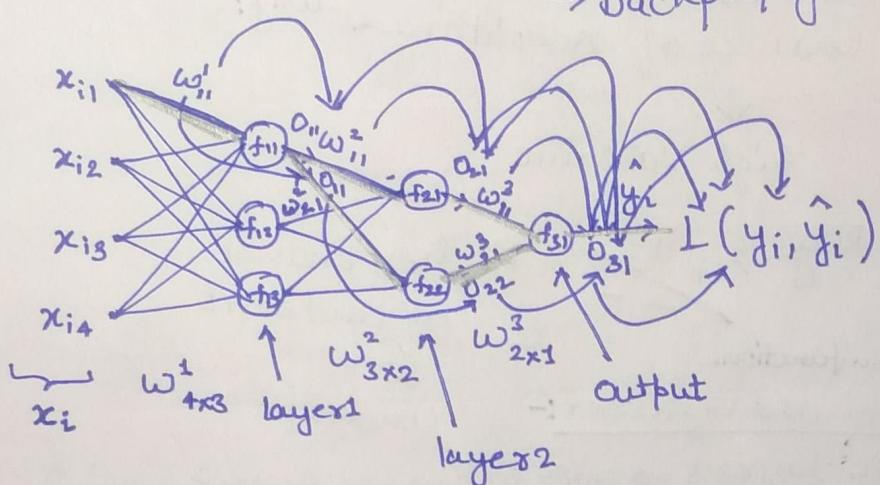
①  $\lambda \rightarrow$  weight of regularizer.

② No. of layers. (we don't want less or more layers).

⑭ Playground - fusesflow.org

⑮ Training an MLP chain Rule.

↳ Backpropagation Algorithm



Now we have find  $4 \times 3 + 3 \times 2 + 2 \times 1 = 20$  weights.

$$\textcircled{1} \quad L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg}$$

↓  
Define loss?

$\sum_{i,j,k} (w_{ijk})^2$  this is something like L2 regularizer.

$\sum_{i,d,k} |w_{ijk}|$  this is L1 regularizer

② Optimization problem:- Sq. Loss + reg.

$$\min_{w^1, w^2, w^3} (L)$$

② SGD or GD (solving optimization problem)

a) initialize  $w_{ij}^k$  } smart initialisation  
 randomly } lot of technique.  
 learning rate

$$\textcircled{b} \quad (w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{ij}^k} \right]_{(w_{ij}^k)_{\text{old}}}$$

Converge  $w_{ij}^k$  until  $(w_{ij}^k)_{\text{new}} \approx (w_{ij}^k)_{\text{old}}$ .

$$\text{ex} \rightarrow \frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3} \Rightarrow \frac{\partial L}{\partial w_{11}^3}$$

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

$w^3$ 's derivatives w.r.t. L are found.  
② weights.

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

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

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

$$\frac{\partial L}{\partial w_{12}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial w_{12}^2}$$

$$\frac{\partial L}{\partial w_{22}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial w_{22}^2}$$

$$\frac{\partial L}{\partial w_{32}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial w_{32}^2}$$

$w^2$ 's derivatives w.r.t L are found.  
③ weights.

Now we have 2 path rather than just 1

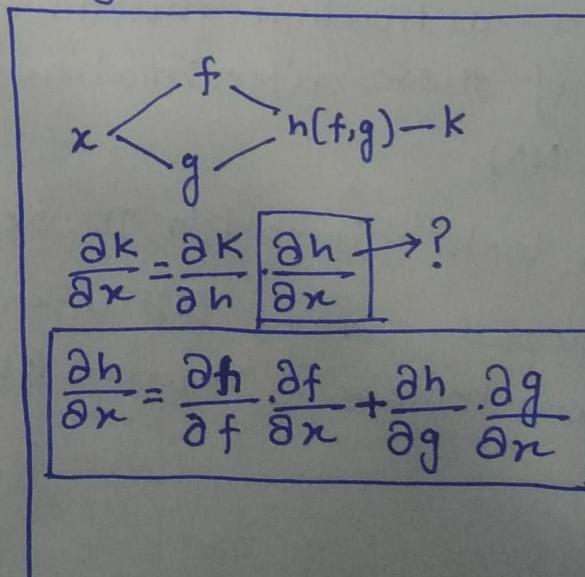
Now this way we can get other  $w^1$ 's derivatives.

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^1} - ??$$

$$\frac{\partial O_{31}}{\partial w_{11}^1} = \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}^1}$$

+

$$\frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}^1}$$



## ⑩ Back propagation.

$$\mathcal{D} = \{x_i, y_i\} \leftarrow \text{Dataset}$$

- ① initialize  $w_{ij}^k$ 's
  - ② for each  $x_i$  in  $\mathcal{D}$ 
    - a) pass  $x_i$  forward through the network
    - b) Compute  $L(\hat{y}_i, y_i)$
    - c) Compute all the derivative chain-rule & memorization.
    - d) Update weights for end of the network to the start.
- $\mathcal{D}$  GDS  
(one by one point)
- $$(w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}} - \eta \left[ \frac{dL}{dw_{ij}^k} \right]$$
- back propagation.

- ③ repeat step 2 till convergence

$$\text{if } (w_{ij}^k)_{\text{new}} \approx (w_{ij}^k)_{\text{old}},$$

when you pass all the datapoints for once we call it 1 epoch.  
we typically run for multiple epochs & these  $x_i$  are chosen randomly in each epoch rather than one after other.

vv important Backpropagation algorithm only work when activation function is differentiable.  
and if these differentiation is fast & easy we can get fast NN.

# we will look mini-batch to optimise the result.

\* we generally avoid GD when all the data backpropagation is used because it require lot of memory.