# Characteristics of NN

1)  Mapping capability      input pattern to their associal output pattern.

2)  learn by example
         train then    used to infer or predict new object

3)  can generalized : therefore can predict new outcomes. from past trends.

4)  Robust & fault tolerant : can recall full pattern from partial or noisy pattern.

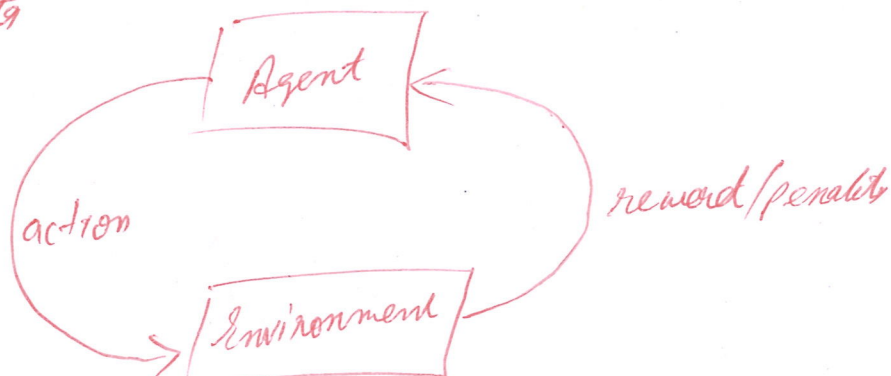5)  can process informat in parallel in a distributed manner.

# LEARNING METHODS

Supervised

| input featur | output labels |
|---|---|
| | |

     labelled data

Unsupervised
     un labelled data

Reinforce.
in absence of data



Agent — action — Environment — reward/penalty

Robotic / Drone stability / etc

Hebbian learning :   Fire together - wire togeth

Based on  correlative weight adjustment

oldest learning mechanism inspired by biology.


Gradient descent

base on minimizat of error E defin in terms of weight & the

activation function of the network.

→ activate fun should be differencia  as weight update depends

on gradient of error E

weigh updat $\Delta W_{ij} = \eta \dfrac{\partial E}{\partial W_{ij}}$
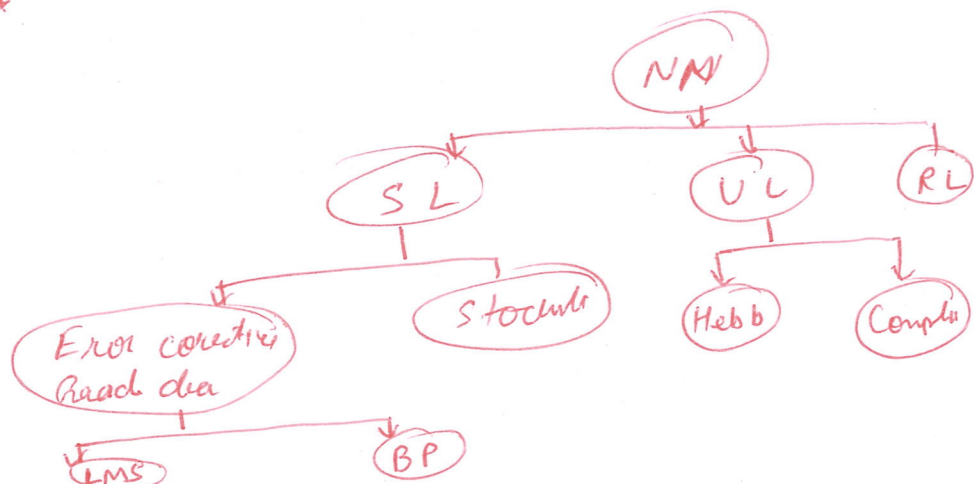

$\eta$ is learning rate   &   $\dfrac{\partial E}{\partial W_{ij}}$  is error gradient

wrt weight
$W_{ij}$


Competitive learning   winner - take - all

those neuron  which responds strongly to input stimuli

have their weight updated

when input is presnt all neura compete & the winner neuron.

undergos weight updat

Stochastic Learning

# LOSS FUNCTON

Euclidan distan.

~~Hamilton~~

Manhatan distanc

MSE

$$C_x(w,b) = \left\| \hat{q}(x) - q(x|w,b) \right\|^2$$

desird        Netu oupu
Targ

For enti Tan D.

$$C(w,b) = \frac{1}{n} \sum_n C_x(w,b)$$

No of imes.

Lows Cos: vous

Inti Wei & Bicn of Neti Rendu Vo

2  Comp  Actual fo eac Tranin Iron

3)  Compad Cosd fe Snti Trin Dale

4)  Wpde Weig & Bias con  Gracch Dcn

5)  Repti Step  2 - 4  un Cos is Rou to an acu Len

$x_1$

$b_1^2$

$b_2^2$

$b_3^3$

$x_2$

$b_1^3$   Signore Neur

$l^2$
Input layer

MNIST (hand write Digits)      60,000   image

28 X 28

$2$ $\Rightarrow$ [ N N ] $\frac{0}{0}$ — NO

Net by   Nielsen 2015

784 neur          30 neur          10 output ourer.
$x_i$

$x_{784}$

$2$

Torh
0        0   0
1        0   0
2        1   0
3        0   0
         0   0
         0   2
         0   0
         0   0
         0   6
9        0   6

Activah.

| Inu | 1 | 2 | 3 |
|-----|---|---|---|
| Labr | 1 | 0 | 3 |
| Activah | | | |
| $\hat{q}(x)$ | | | |

Desire Activah.

Inur — Truh → [ N N ] Neh Activ → | Cost Comp | → Cost

784      30      10      $y$     $t$

$w_i$ $b_j$

$\underbrace{\phantom{xxxxx}}$ Cost fnct

$C \equiv Cost$

$b_j \equiv biases$

$w_i = weights$

init cost

$func = f(x, y)$      Taylor Series

& to fnd $f(x + \Delta x, x + \Delta y)$ with line approx...

$$\cong f(x, y) + \left\{ \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial y} \Delta y \right\}$$

$$= f(x, y) + \Delta f(x, y)$$

when,

$$\Delta f(x, y) = \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial y} \Delta y$$

Similar

$$\Delta C = \Delta C_{(w, b_j)} = \frac{\partial C}{\partial w_i} \Delta w_i + \frac{\partial C}{\partial b_j} \Delta b_j$$

In vector form:

$$\Delta C = \begin{bmatrix} \frac{\partial C}{\partial w_i} & \frac{\partial C}{\partial b_j} \end{bmatrix} \cdot \begin{bmatrix} \Delta w_i \\ \Delta b_j \end{bmatrix}$$

$$= \underset{GRADIENT}{} \cdot \underset{change\ in\ weight\ \&\ bias}{}$$

$$= \nabla C \cdot \Delta \sigma$$

$$\Delta C = \nabla C \cdot \Delta v$$

(change in wees & bias)

Let $\Delta v = -\nabla C$



we need to move opposite to gradient

$\text{or} -\nabla C$

$\therefore$ we use a fraction of this

$\text{wn} -\eta \nabla C$

$\downarrow$

learn rate

$$\therefore \Delta C = \nabla C \cdot (-\eta \nabla C)$$
$$= -\eta (\nabla C \cdot \nabla C)$$
$$= -\eta \| \nabla C \|^2$$

$\searrow$ Magnitude

$\therefore$

$$\boxed{\begin{array}{l} w_i \rightarrow w_i' = w_i - \eta \dfrac{\partial C}{\partial w_i} \\[3mm] b_j \rightarrow b_j' = b_j - \eta \dfrac{\partial C}{\partial b_j} \end{array}}$$

Update rule
$\rightarrow$ need to find deriv of
Cost wrt all
parameter

$\eta \rightarrow$ Small      more computate      slow convergus

$\eta \rightarrow$ large      less compute      faster conv $\left\{ \text{but may over shoot} \right.$

## From Taylor Series

$$f(x+\Delta x, y+\Delta y) \simeq f(x,y) + \frac{\partial f}{\partial x}\Delta x + \frac{\partial f}{\partial y}\Delta y$$

$$\overbrace{\Delta f(x,y)}$$

$$\Delta f(x,y) \equiv \frac{\partial f}{\partial x}\Delta x + \frac{\partial f}{\partial y}\Delta y$$

Applying to Cost func (Loss fun)

$$\Delta C \simeq \frac{\partial C}{\partial w_i}\Delta w_i + \frac{\partial C}{\partial b_j}\Delta b_j$$

$$= \begin{bmatrix} \frac{\partial C}{\partial w_i} & \frac{\partial C}{\partial b_j} \end{bmatrix} \begin{bmatrix} \Delta w_i \\ \Delta b_j \end{bmatrix}$$

Change in Cost:

$$\Delta C = \nabla C \cdot \Delta v$$

↓ Gradient

↘ Chang in Weigh & bias



$\nabla C$

we need opposid direcfo $\nabla C$

Let $\Delta v = -\eta \nabla C$

↘ Learn rate

$$\therefore \Delta C = \nabla C \cdot -\nabla C = -\|\nabla C\|^2$$

$$\boxed{\begin{array}{l} w_i \to w_i' = w_i - \eta \frac{\partial C}{\partial w_i} \\[2mm] b_j \to b_j' = b_j - \eta \frac{\partial C}{\partial b_j} \end{array}}$$

Wpade rule
→ need $\#$ find derivat of cost wrt all the paremeters

$\eta$ small $\rightarrow$ more compulat (slow convergen)

$\eta$ - large $\rightarrow$ fasters (fast convergen) but may over shoot.

# GRADIENT DECENT for Entire networs

$$\Delta v = \begin{bmatrix} \Delta w_{jk}^{(\ell)} \\ \vdots \\ \Delta b_j^{(\ell)} \\ \vdots \end{bmatrix} \quad ; \quad \nabla C = \begin{bmatrix} \frac{\partial C}{\partial w_{jk}}^{(\ell)} \\ \vdots \\ \frac{\partial C}{\partial b_j}^{(\ell)} \\ \vdots \end{bmatrix} \quad ; \quad \begin{array}{l} j : \text{neuron} \\ k : \text{inputs} \\ \ell : \text{layers.} \end{array}$$

$$\boxed{\Delta v = -\eta \nabla C} \quad \text{is still valid}$$

$$w_{jk}^{(\ell)} \longrightarrow w_{jk}^{(\ell)\prime} = w_{jk}^{(\ell)} - \eta \frac{\partial C}{\partial w_{jk}}^{(\ell)}$$

$$b_j^{(\ell)} \longrightarrow b_j^{(\ell)\prime} = b_j^{(\ell)} - \eta \frac{\partial C}{\partial b_j}^{(\ell)}$$