## Neural Networks:-
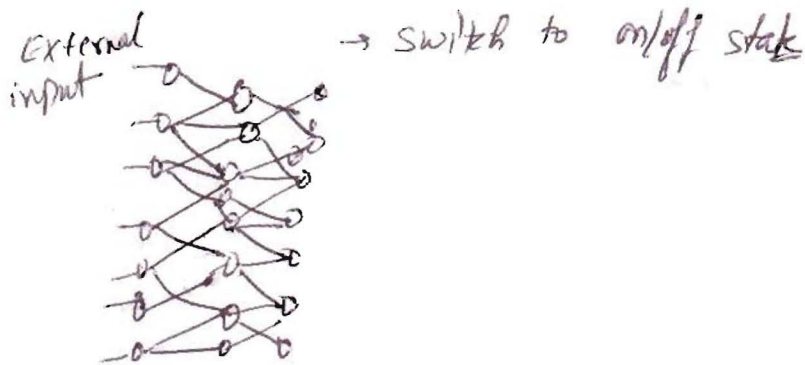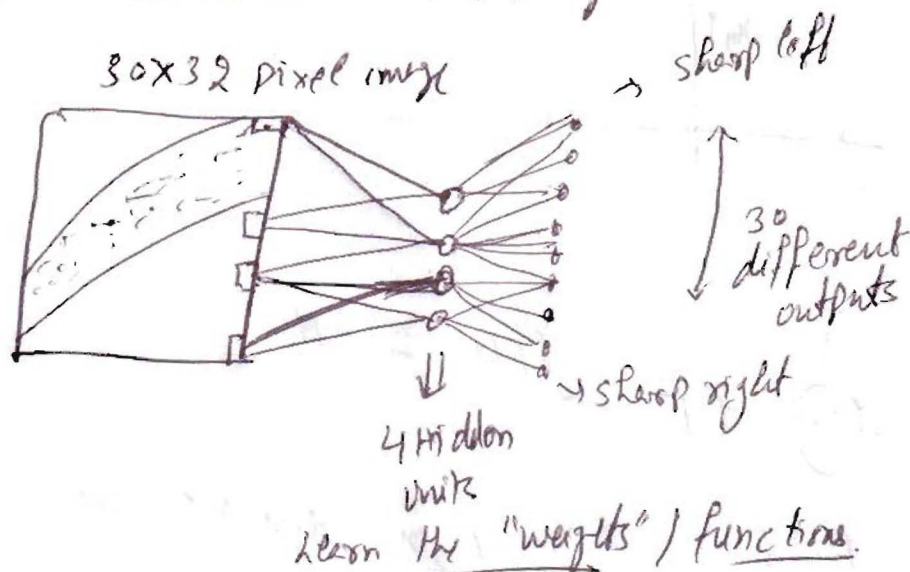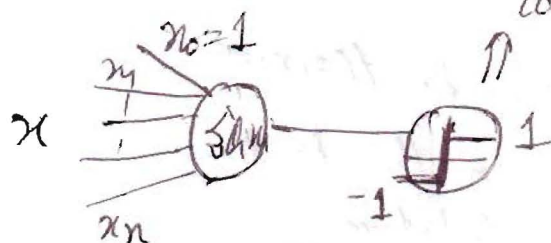
1. One of the first models for learning
2. Motivated by "neurons" in the brain.
3. Massively parallel Architecture

   ~ $10^{10}$ neurons
   ~ $10^{4-5}$ connections per neuron
   ~ switching time $\simeq$ .001 seconds
   ~ scene recognition ~ .1 second

External input      → Switch to on/off state

Example: Autonomous Car Driving

30×32 pixel image      → sharp left

30 different outputs

→ sharp right

4 Hidden units

Learn the "weights" / functions.

**Basic Unit:-**

Perceptron:- (neuron)



computes $g(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$

$y^{(i)} \in \{-1, 1\}$

$h_\theta(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ -1 & \text{if } \theta^T x < 0 \end{cases}$

$\Downarrow$

$g(\theta^T x)$

What kind of function is this?
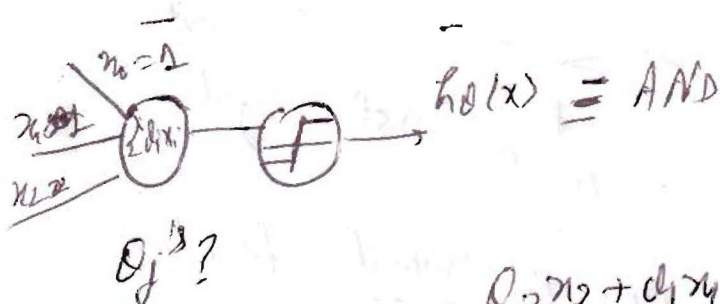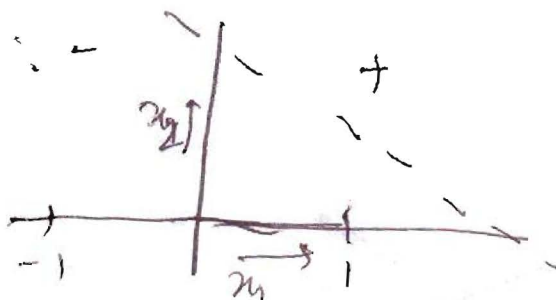(classifier)
$\downarrow$
Linear

~~Can Represent Linearly functions~~

Can classify linearly separable Data

Example:- AND function





$h_\theta(x) \equiv$ AND

$\theta_j$'s?

$\theta_2 x_2 + \theta_1 x_1 + \theta_0 > 0$ if $x_1 = 1$, $x_2 = 1$

$\theta_2 x_2 + \theta_1 x_1 + \theta_0 < 0$ ~~~~ o.w.

$\theta_2 x_2 + \theta_1 x_1 > \theta_0$ } $\begin{cases} \sum_{j=1}^{2} \theta_j x_j > \text{threshold} & \rightarrow +ve \\ \text{else} & \rightarrow -ve \end{cases}$

$\theta_2 = 1, \theta_1 = 1, -\theta_0 = 1.5$

Represents AND

Alternate formulation

③

What about OR function?

$$\theta_2 = 1$$
$$\theta_1 = 1$$
$$\theta_0 = -0.5$$
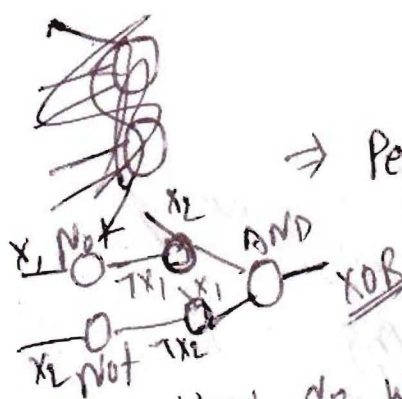
XOR?



There is no line which
separates + & -.

$\Rightarrow \nexists \; \theta_2, \theta_1, \theta_0$
s.t. $\theta_2 x_2 + \theta_1 x_1 + \theta_0 > 0$ for $+ve$
AND $\theta_2 x_2 + \theta_1 x_1 + \theta_0 < 0$ for $-ve$

Example

$\Rightarrow$ Perceptron is a linear classifier.
$\rightarrow$ Need to combine individual perceptrons
to get a non linear classifier. $\boxed{x_1 \oplus x_2 = x_1 \neg x_2 + \neg x_1 x_2}$

How do we learn the weights (parameters)?

Perceptron training rule (delta rule)

$$\Delta \theta_j = \eta \left[ y - h_\theta(x^{(i)}) \right] x_j^{(i)} \quad \text{Similar to gradient descent}$$

where $x^{(i)}$ is a wrongly classified example.

Note $J(\theta) = \sum_{i=1}^{m} \left[ y^{(i)} - h_\theta(x^{(i)}) \right]^2$ is not differentiable

So, not really gradient descent.

Algorithm:  init($\theta$)
do {
  for $i = 1$ to $m$ $\quad x^{(i)} \leftarrow$ wrongly classified example based on
  $\quad \Delta \theta_j = \eta \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$
  $\quad \theta_j \leftarrow \theta_j + \Delta \theta_j$
} until (all examples correctly classified)

**Properties:—** Guaranteed to converge for linearly separable data

(Assuming $\eta$ is sufficiently small).

~~Its~~ won't work non-linearly separable data

**Issues:—**
1. Sort of ad-hoc?
2. Non-differentiable error function
3. Can not handle non linear data

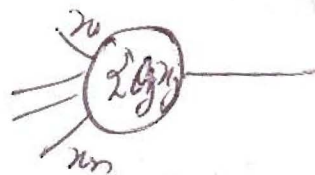**Context:—** Much before then some modern algorithms came about.

**why studying:—**
1. Historical context
2. Works well for certain class of problems (digit recognition)
3. Similarity to working on brain
4. Deep Networks:— Very powerful/effective.
   Belief
   On going research

**Rectify some of the problems:—**

$$h_\theta(x) = \theta^T x.$$

Directly optimize the unthresholded output.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{m} \left[ y^{(i)} - h_\theta(x^{(i)}) \right]^2$$



$$\Downarrow$$
$$\theta^T x^{(i)}$$

Similar to linear regression—
Difference:— $y^{(i)} \in \{-1, 1\}$

⑤

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{2} \sum_{i=1}^{m} 2(y^{(i)} - h_\theta(x^{(i)})) \cdot \frac{\partial}{\partial \theta_j} (\cdot) x_j^{(i)}$$

Aside:-

$\left[ \frac{\partial}{\partial \theta_j} \theta^T x^{(i)} = x_j^{(i)} \right]$   We want to minimize $J(\theta)$

$$\Rightarrow \theta_{new} = \theta_j - \eta \frac{\partial}{\partial \theta_j} J(\theta)$$

$\Rightarrow$ update rule:-

$$\theta_j = \theta_j + \Delta\theta_j$$

$$\Delta\theta_j = \eta \sum_{i=1}^{m} [y^{(i)} - h_\theta(x^{(i)})] x_j^{(i)}$$

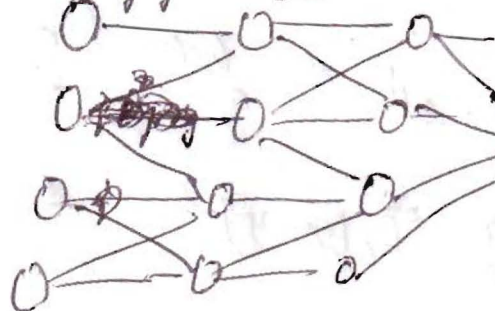Batch Gradient / Stochastic Gradient
Descent / Descent

Issues:-  ① $y^{(i)} \in \{-1, 1\}$

$h_\theta(x^{(i)}) \in (-\infty, \infty)$

$\Rightarrow$ Bigger issue:- Can not handle non
linear data!

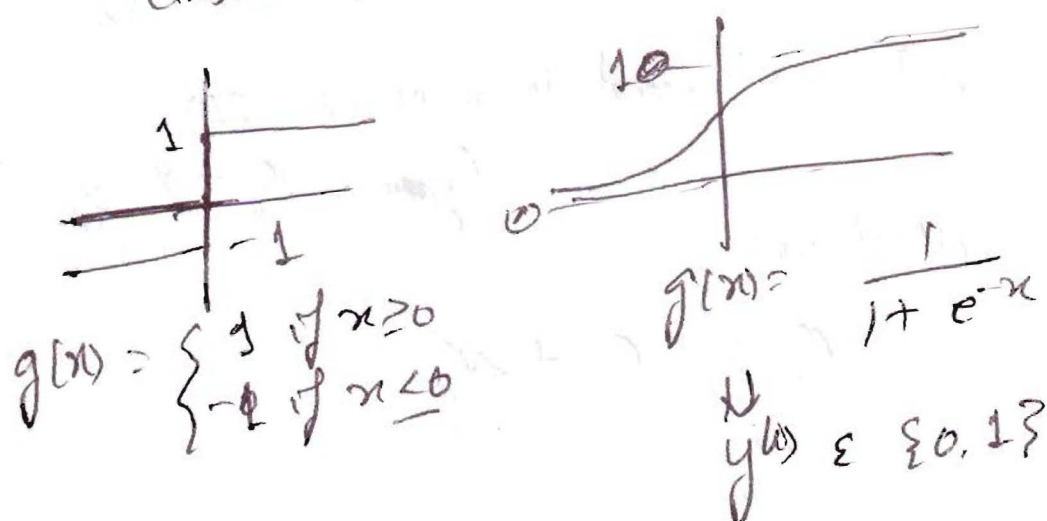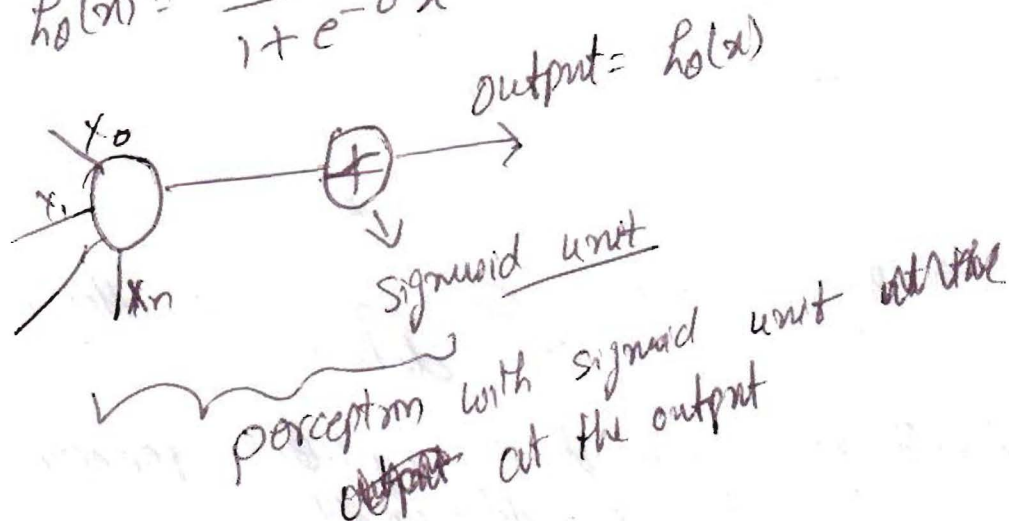[*Could use combination thresholded perceptrons
→ but non differentiable]

$\{\theta_j^T x_j \atop \}$   Each unit computes a linear fn



output is a linear
∏  function of inputs

How do we solve this?

-Idea:- Use a sigmoid unit
instead of thresho(ding)

$g(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x \leq 0 \end{cases}$

$g(x) = \dfrac{1}{1 + e^{-x}}$

$y^{(i)} \in \{0, 1\}$

$\Rightarrow \quad h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$

output $= h_\theta(x)$

$x_0$
$x_1$
$x_n$

Sigmoid unit

perceptron with sigmoid unit at the output

§ Properties:-
1. Differentiable
2. can represent Non-linear functions.

$J(\theta) = \dfrac{1}{2} \sum\limits_{i=1}^{m} [y^{(i)} - h_\theta(x^{(i)})]^2$

$\dfrac{\partial J(\theta)}{\partial \theta_j} = \dfrac{1}{2} \sum\limits_{i=1}^{m} [2 (y^{(i)} - h_\theta(x^{(i)}))] \cdot (-1) \dfrac{\partial}{\partial \theta_j} h_\theta(x^{(i)}))$

$\dfrac{\partial h_\theta(x^{(i)})}{\theta_j} = \dfrac{\partial}{\partial [\theta^T x^{(i)}]} \underbrace{h_\theta(x^{(i)})}_{\downarrow \text{ sigmoid}} \cdot \dfrac{\partial [\theta^T x^{(i)}]}{\partial \theta_j}$

(7)

$$\frac{d}{dx} g(x) = g(x)(1-g(x))$$

$$g(x) = \frac{1}{1+e^{-x}} \qquad h_\theta(x) = g(\theta^T x)$$

$$\Rightarrow \frac{d}{\theta_j} h_\theta(x^{(i)}) = h_\theta(x^{(i)})(1-h_\theta(x^{(i)})) \, x_j^{(i)}$$

$$\Rightarrow \frac{d}{d\theta_j} J(\theta) = \sum_{i=1}^{m} \left[ y^{(i)} - h_\theta(x^{(i)}) \right] (-1) \, h_\theta(x^{(i)}) (1-h_\theta(x^{(i)})) \, x_j^{(i)}$$

$$\theta_j \leftarrow \theta_j + \eta \sum_{i=1}^{m} \left[ y^{(i)} - h_\theta(x^{(i)}) \right] \left[ h_\theta(x^{(i)}) \right] \left[ 1 - h_\theta(x^{(i)}) \right] x_j^{(i)}$$

① Different from logistic regression → what is different? why?
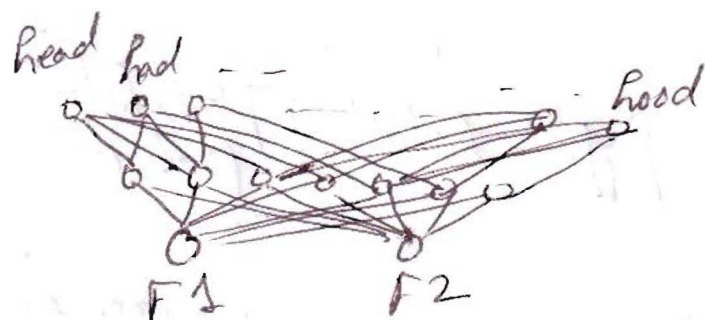
② Non-convex optimization ⇒ Local minima ≠ Global minima.

2. Representing Non-Linear Functions :-



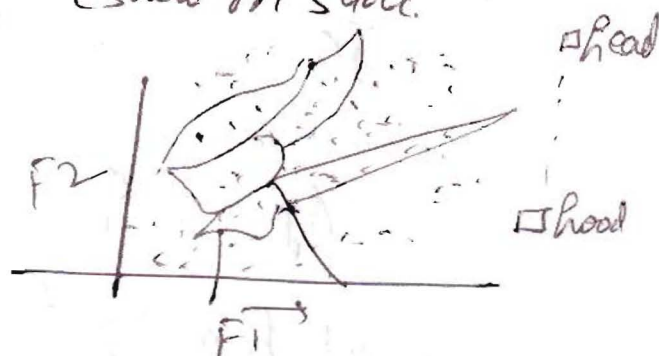Final output

$\pi$

Layer 2

one layer 1

Layer n

# Examples:-

1. Non-linear boundaries

   Recognition of different sounds given two different frequency inputs (characterizing sound)



   Decision surface (show on slide.

2. Video:- Andrew Ng Lecture 6
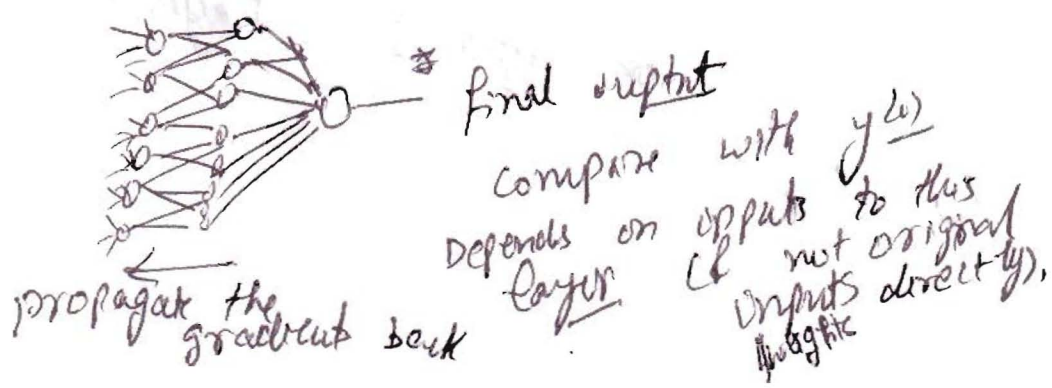
   1. 37:03 - 39:45       (Digit Recognition)
   2. 41:30 - 42:50       (Producing Speech) from Text

# Training Multi-layered Networks :-
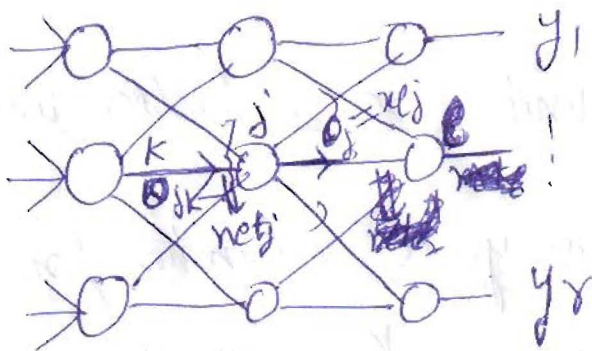
## Backpropagation Algorithm

Main Issue:-   $J(\theta)$ not directly a function of original inputs/weights



final output

compare with $y^{(i)}$

Depends on inputs to this layer (& not original inputs directly).

propagate the gradient back

Deriving backpropagation :-



$$net_j = \sum_K \theta_{jk} \cdot x_{jk}$$

We wibd will do it for one example

Let the output of the units in last(output)
layer be:    $O_\lambda$ ___    or (implicitly a function of $x$).

$$J(\theta) = \frac{1}{2} \sum_{\ell=1}^{r} (y_\ell - O_\ell)^2$$

Note: $\theta_{jk}$ can influence $J(\theta)$ only
through $net_j$

Now,

$$\boxed{\frac{\partial}{\partial \theta_{jk}} J(\theta) = \frac{\partial}{\partial net_j} J(\theta) \cdot \frac{\partial net_j}{\partial \theta_{jk}} = \frac{\partial J(\theta)}{\partial net_j} \cdot x_{jk}}$$

$$= \frac{\partial J(\theta)}{\partial net_j} \cdot \cancel{x_{jk}}$$

① Output units :-

Note: $net_j$ can
influence $J(\theta)$
only through
$O_j$ ~~$O_{jk}$~~

~~[since we are
considering change
in $w_{jk}$ & $w_{jk}$ are
unique for each unit]~~

$$\frac{\partial J(\theta)}{\partial net_j} = \frac{1}{2} \frac{\partial}{\partial net_j} \sum_{\ell=1}^{r} (y_\ell - O_\ell)^2$$

$$= \frac{1}{2} \frac{\partial}{\partial net_j} (y_j - O_j)^2$$

$$= \frac{1}{2} \times 2 (y_j - O_j)(-1) \frac{\partial O_j}{\partial net_j} = -(y_j - O_j) \cdot (O_j)(1 - O_j)$$

Let

$$\frac{\partial J(\theta)}{\partial net_j} = S_j$$

$$\boxed{S_j = (y_j - O_j)\, O_j(1 - O_j)}$$
$\partial \varepsilon$ output layer

Not, hidden layers:-

Let unit $j$ be a hidden unit

Aside:-

$f(y_i - y_k)$ suppose $y_j$ is a function of $x$

Then $\dfrac{\partial f(y_i - y_k)}{\partial x} = \sum\limits_{j=1}^{k} \dfrac{\partial f(y_i - y_k)}{\partial y_j} \cdot \dfrac{\partial y_k}{\partial x}$

Using this:-

$$\frac{\partial J(\theta)}{\partial net_j} = \sum_{e \in \text{downstream}(j)} \frac{\partial J(\theta)}{\partial net_e} \cdot \frac{\partial net_e}{\partial net_j}$$

$$\frac{\partial net_e}{\partial net_j} = \frac{\partial(\theta_{ej}\, x_{ej})}{\partial O_j} \qquad \frac{\partial O_j}{\partial net_j}$$

$$= \theta_{ej} \cdot O_j(1 - O_j)$$

$$\Rightarrow \frac{\partial J(\theta)}{\partial net_j} = \sum_{e \in \text{downstream}(j)} - S_e \cdot \theta_{ej} \cdot O_j(1 - O_j)$$

$$\boxed{\Rightarrow S_j = O_j(1 - O_j) \sum_{e \in \text{downstream}(j)} S_e \cdot \theta_{ej}}$$ $\partial \varepsilon$ hidden unit