# Machine Learning (COL 774)

## Neural Networks: Basics
### Mar 31, 2020

①

Deep Learning Models. COL774

Machine Learning

Mar 31, 2020
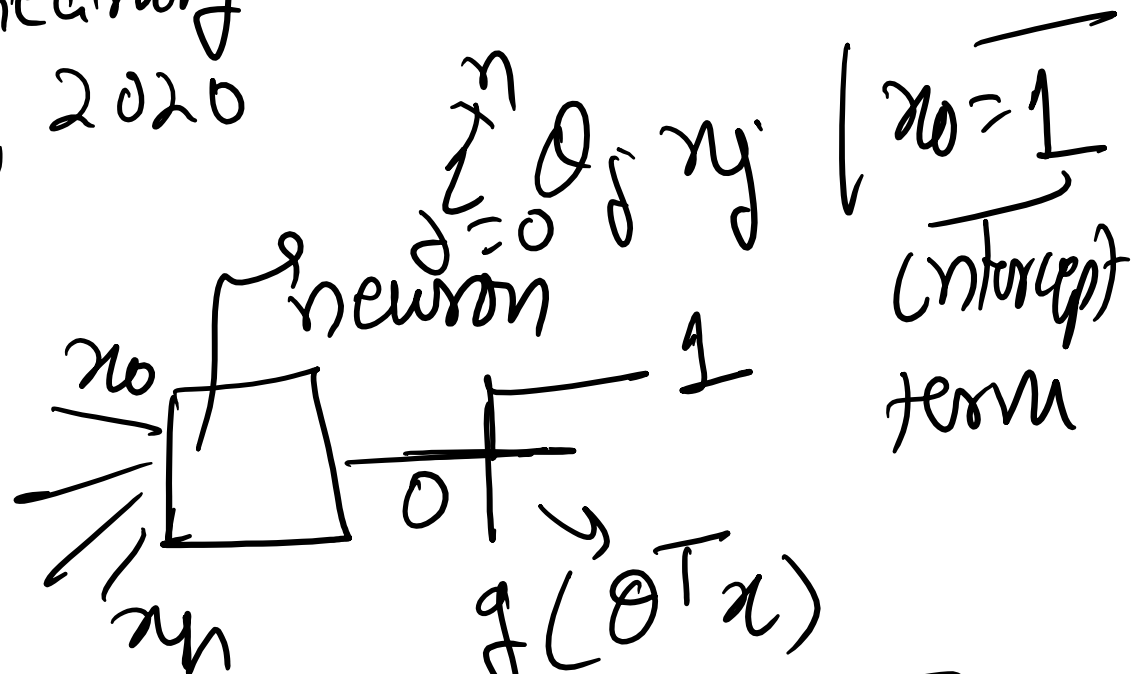
Neural Networks :-
_____

} :- Brain :-

$\sum_{j=0}^{n} \theta_j x_j \Big|_{x_0 = 1}$

neuron

$x_0$

$\theta^T$

1

(intercept term

$g(\theta^T x)$

$\hookrightarrow \begin{cases} 1 & \{\theta^T x \geq 0\} \\ \end{cases}$

$\boxed{x_0 = 1}$

$\sum_{j=1}^{n} \theta_j x_j := -\infty$

→ millions of neurons in human brain

→ patterns of interconnections

$x_n$

if neuron fires

no Layer 1     Layer 2                    Layer (n-1)   Layer n.

$x_n$

Recognize
grand mother

Activation
function

$g(z) = g(\theta^T x) = 1 \leq \theta^T x \geq 0$
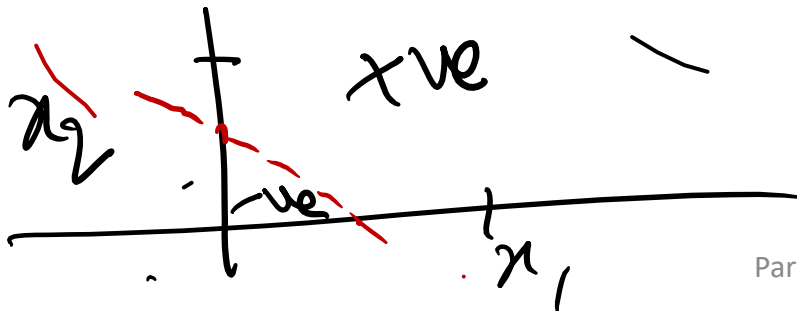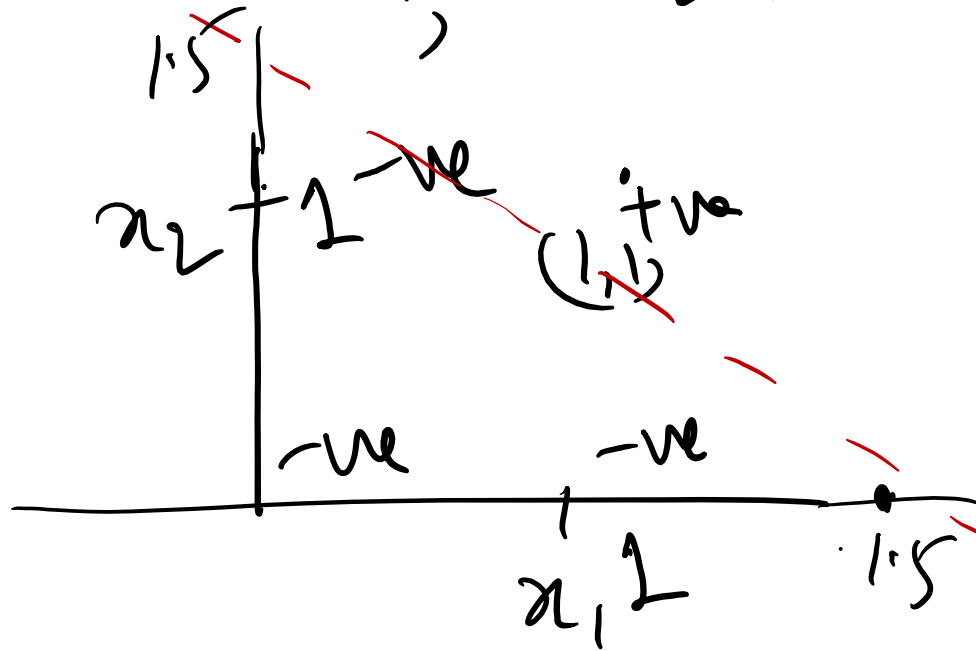
$h_\theta(x) = g(\theta^T x)$  ① Logistic Reg

② SUMs

# Representing Boolean Functions :—

AND:—  $f(x) = x_1 \wedge x_2$

$$x_1, x_2 \in \{0,1\}$$



$$h_\theta(x) = 1\{\theta_2 x_2 + \theta_1 x_1 + \theta_0 \geq 0\}$$

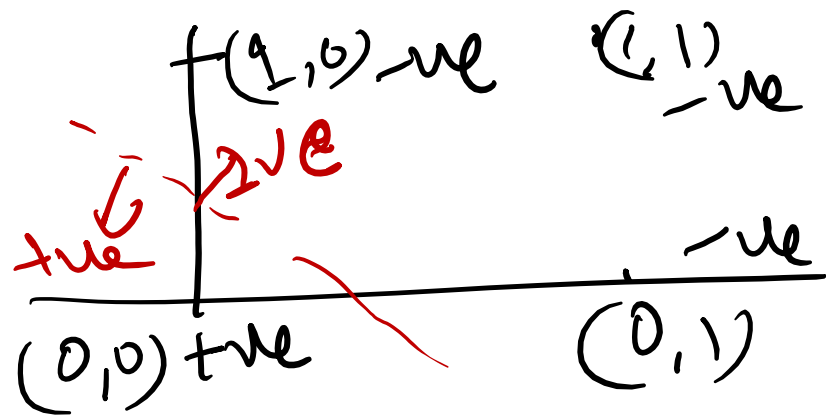$$= 1 \quad \text{if } x_2 = 1 \wedge x_1 = 1$$

AND :

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

OR :

$$h_\theta(x) = 1\{\theta^T x \geq 0\}$$

$$\theta_2 = 1, \quad \theta_1 = 1, \quad \theta_0 = -0.5$$

NAND :- $\neg(x_1 \wedge x_2)$

$f(1,0) \,\, -ve$  $(1,1)\,\, -ve$

true $\nearrow$ +ve

$-ve$

$(0,0)$ +ve   $(0,1)$

$h_\theta(x) = \mathbb{1}\{\theta^T x \geq 0\}$  NAND

$\theta_2 = -1, \quad \theta_2 = -1, \quad \theta_0 = 0.5$

HW:-   $f(x) = \neg x_1$   $h_\theta(x) = \mathbb{1}\{\theta_1 x_1 + \theta_0 \geq 0\}$
$\hookrightarrow$ NOT

XOR:-   $f(x) = x_1 \oplus x_2$   $\rightsquigarrow$ 1 iff
Exactly one of
$x_1$ and $x_2$ = 1

Parag Singla @ IIT Delhi

$x_2$ | tve $\quad\quad\quad\quad$ $\overset{-ve}{(1,1)}$ $\quad$ $f(x) = x_1 \oplus x_2$

can not be represented

by a single perception.

-ve $\quad$ tve
$\quad\quad\quad$ 1

———————————

$x_1$

$\hookrightarrow \quad (\neg x_1 \wedge x_2) \vee (x_1 \wedge \neg x_2)$

$\underbrace{\quad\quad\quad\quad\quad}_{AND} \quad \underbrace{\quad\quad\quad\quad\quad}_{AND}$

$\underbrace{\quad\quad}_{OR}$

$x_0$

$x_1$ $\quad$ $\boxed{\neg x_1 \wedge x_2}$ $\quad y_1$ $\quad\quad$ $\boxed{y_1 \vee y_2}$

$x_2$ $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ OR

$x_0$

$x_1$ $\quad$ $\boxed{x_1 \wedge \neg x_2}$ $\quad y_2$

$\neg x_1 \wedge x_2 \quad :- \qquad h_\theta(x) = 1\{\theta^T x \ge 0\}$

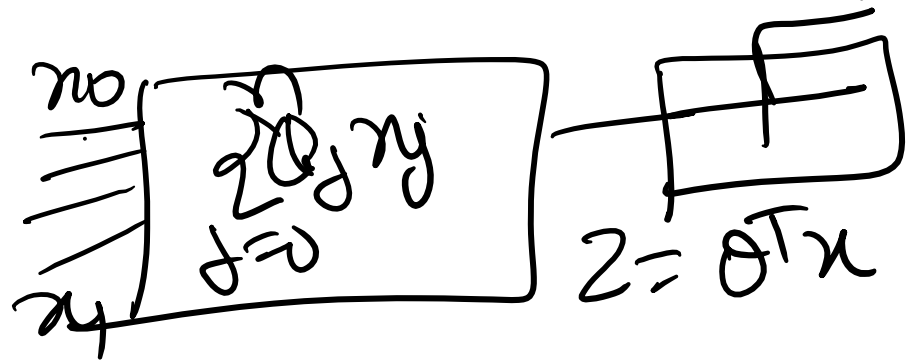$\hookrightarrow \quad \theta_2 = 1, \quad \theta_1 = -1, \quad \theta_0 = 0.5$

Similarly $\quad x_1 \wedge \neg x_2 :-$

$$XOR \equiv y_1 \vee y_2 \quad \text{where } y_1 = \neg x_1 \wedge x_2$$
$$y_2 = x_1 \wedge \neg x_2$$

$\hookrightarrow$ Min # of perceptrons :- 3

Learning the parameters of a perceptron :-

Delta Rule :-

$$g(z) = \{1 \text{ iff } \theta^T x \geq 0\}$$

$$\sum_{j=0}^{n} \theta_j x_j \qquad \rightsquigarrow \text{Activation Units.}$$

$$z = \theta^T x \qquad \hookrightarrow \text{can not pass gradient}$$

① Linearly seprable Data

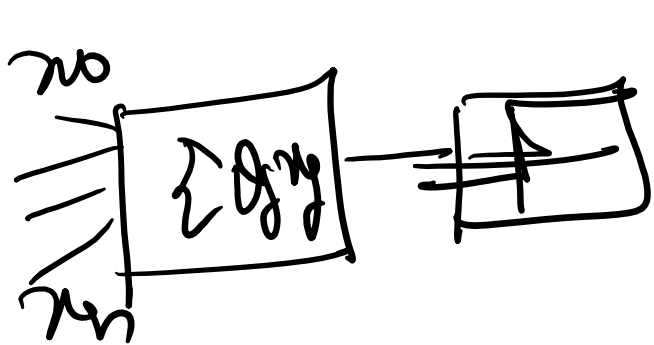② $\eta :-$ (Learning rate) sufficiently
small.

$\theta^{(0)} = \text{init}();$
$t = 0;$
while $(\exists i : y^{(i)} \neq h_{\theta^{(t)}}(x^{(i)}))$

---

① Not very principled

② Does not work when date
is NOT Linearly Separable

$$\begin{cases} \theta_j^{(t+1)} \leftarrow \theta_j^{(t)} \\ + \eta [y^{(i)} - h_{\theta^{(t)}} x^{(i)}] \\ \} t \leftarrow t+1; \quad x_j^{(i)} \end{cases}$$
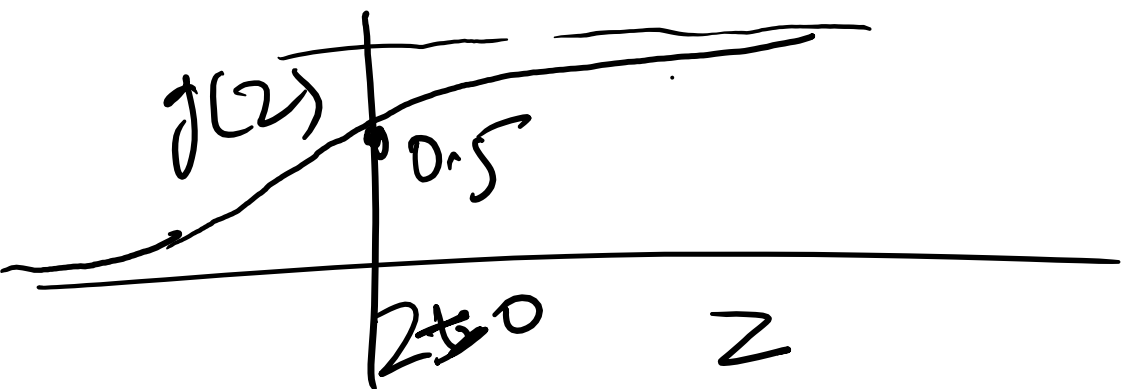
$$g(z) = \begin{cases} 1 & \theta^T x \geq 0 \\ 0 & \end{cases}$$

a step fn

$$z = \theta^T x$$

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

$$= \frac{1}{1 + e^{-\theta^T x}}$$

very similar

logistic



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

No probalistic interpretation

$z \Longrightarrow$ For $z > 0$

$g(z) \to 1$

For $z < 0$

$g(z) \to 0$

Learn Using gradient Descent

Logistic:- $-LL(\theta)$

$\underbrace{}_{\text{Error Function}}$

$\underset{\theta}{argmin} \ -LL(\theta)$

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

Contrast with objective fn. for logistic ~~Regression~~

$$\nabla_\theta J(\theta)$$

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

$\rightarrow$ derivative of Sigmoid fn.

$$h_\theta(x^{(i)}) = g(\theta^T x^{(i)})$$

$$\nabla_z g(z)$$

$$= g(z)(1 - g(z))$$

$$\nabla_\theta g(\theta^T x^{(i)}) = g(\theta^T x^{(i)})(1 - g(\theta^T x^{(i)}))$$

$$\cdot \boxed{\nabla_\theta (\theta^T x^{(i)})} \; \textcolor{red}{x^{(i)}}$$

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

$$\underbrace{\qquad\qquad\qquad\qquad} \rightarrow \text{highly non convex}$$

$$h_\theta(x^{(i)}) = g(\theta^T x^{(i)})$$

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

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

$\hookrightarrow$ extra term compared to Logistic Reg.

→ Initialization

{ Gradient Descent Update:-

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \cdot \nabla_\theta J(\theta) \big|_{\theta(t)}$$

{ $t \leftarrow t+1$

↓ Learning rate

while ( ! Converged)

Completes Learning
for a Perceptron

---

Multi Layered   Neural Network → Hidden Layers

$x_0$

$x_n$

$x_0$

$x_{n}^{(1)}$

$+ \square - O_1$

$\square - O_r$

Output
Layer

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \sum_{e=1}^{r} \left( y_e^{(i)} - O_e \right)^2 \quad [\text{r output units}]$$

$J(\theta) \downarrow$
entire
set of parameters

Assume:- $\underline{m = 1}$ (Only to simplify the notation)

$\nabla_\theta J(\theta)$ $\Big\{$ parameters in the output layer

parameters in the hidden layers

$$J(\theta) = \sum_{e=1}^{r} \left( y_e - O_e \right)^2 \quad [m = 1]$$