

①

Multi-layer Perceptions

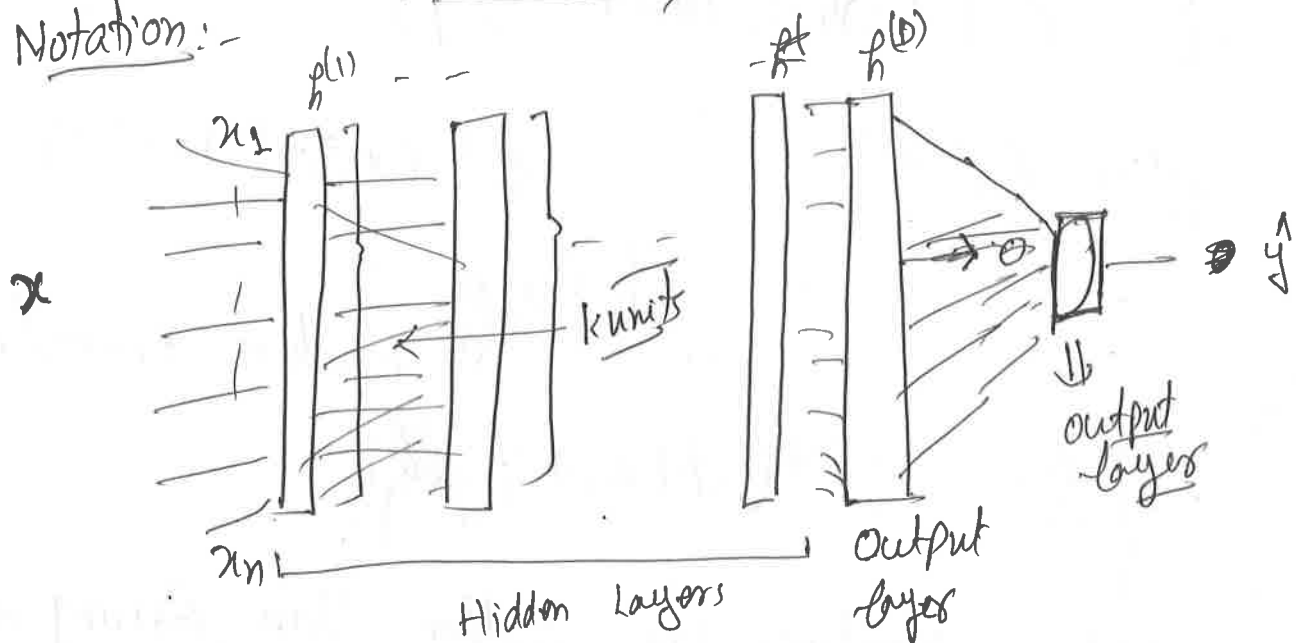
(Deep Learning)

COL 865

July 27, 2017

Forward Connections

Notation:-



Each layer transforms vector x into a new vector x' . At the final layer, output y is produced.

$$y = f(x; \theta) = f^{(0)} \circ f^{(1)} \circ \dots \circ f^{(L)}(x)$$

Approximate

$$y = f^*(x)$$

$$\theta = (\theta, \theta^{(1)})$$

Activation function

Questions-

What is the form of each $f^{(l)}$?

What is the form of $f^{(0)}$?

What do we want to optimize?

i.e. $L(y, y'; \theta)$ Loss (e.g. cross-entropy)

How do we optimize?

How do we optimize $L(y, y')$?

Design questions

Note:-

$$\cancel{h^{(0)}(h^{(D-1)} \dots h^{(1)}(x))} \equiv \underbrace{\phi(x; \theta)}_{\substack{\text{(Non-linear)} \\ \text{Feature transformation}}}$$

$$\cancel{y = f^{(0)}(\phi(x; \theta); \omega^{(0)}, \theta^{(0)})}$$

$$h^{(0)}(h^{(D-1)} \dots h^{(1)}(x; \theta^{(1)}; \theta^{(2)} \dots), \theta^{(0)}) = \phi(x; \theta) \Downarrow \text{Feature transformation}$$

$$y = f^{(0)}(\phi(x; \theta); \omega^{(0)})$$

Kernel Function

$$K(x, x') = \phi(x) \cdot \phi(x')$$

Note:- Rather than having a single (non-linear) feature transformation, (e.g. SVMs), we ϕ is expressed as a ~~combination~~ composition of several non-linear transformations.

~~What is form of each~~

Intuition:- Mimicks the computation in Brains of transfer learning

Two parts: - (a) Feature transformation } Hidden layer
(b) Output processing } Output layer

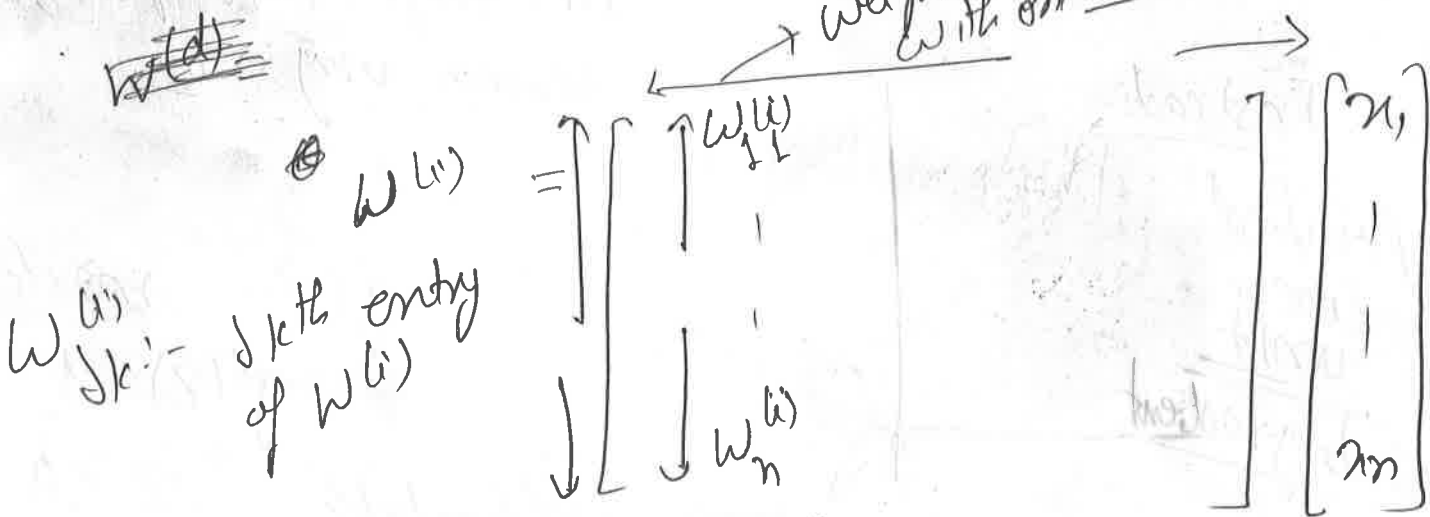
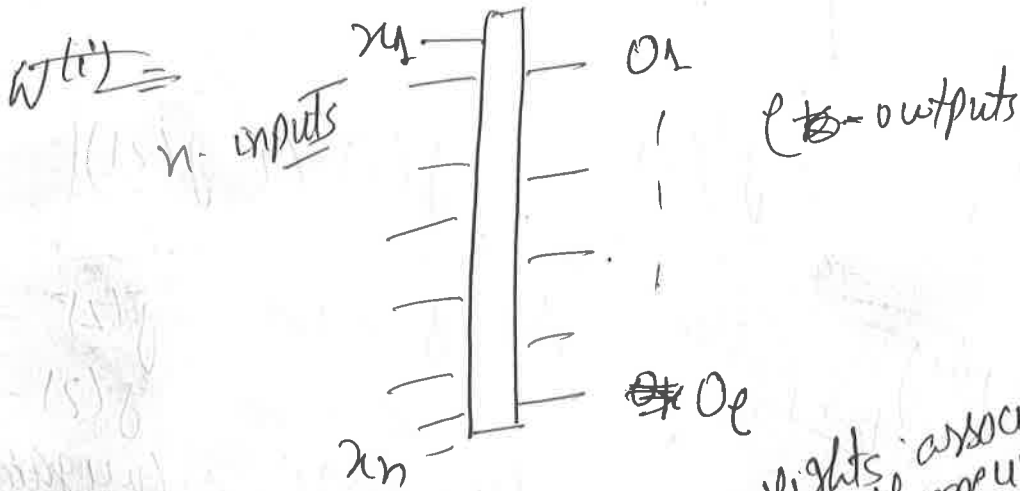
~~$f^{(i)}$~~ Each $f^{(i)}$ implements a non-linear function.

② $f^{(l)}(x; \theta^{(l)}) = g[\underbrace{w^{(l)T}x + b^{(l)}}]$

↓
Vector
valued
function

↓
Non-linearity
(activation
function)
↳ can be shared across layers

↓
Linear function



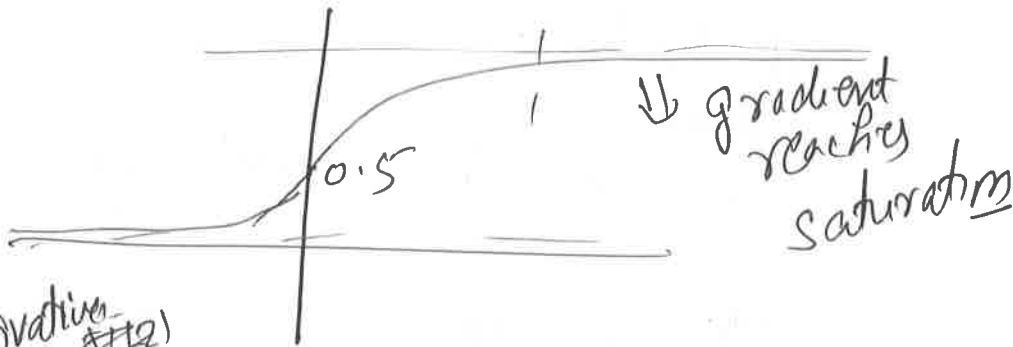
$$b^{(l)} = \begin{bmatrix} b_1^{(l)} \\ \vdots \\ b_c^{(l)} \end{bmatrix}$$

$$x' = g^{(l)}(\underbrace{w^{(l)T}x + b^{(l)}}_{\text{element wise application}})$$

$$\begin{bmatrix} x'_1 \\ \vdots \\ x'_c \end{bmatrix} = g \left(\begin{bmatrix} w_{11}^{(l)} & \dots & w_{1n}^{(l)} \\ \vdots & \ddots & \vdots \\ w_{c1}^{(l)} & \dots & w_{cn}^{(l)} \end{bmatrix} x + \begin{bmatrix} b_1^{(l)} \\ \vdots \\ b_c^{(l)} \end{bmatrix} \right)$$

g : activation unit (non-linearly)

choices: - $g(z) = \frac{1}{1 + e^{-z}}$: - sigmoid function



Left Derivative
 $\lim_{z \rightarrow 0^-} g'(z)$
 Right derivative
 $\lim_{z \rightarrow 0^+} g'(z)$

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

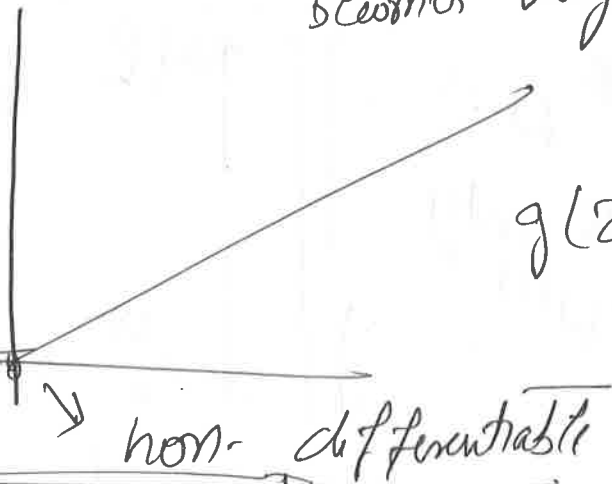
Why is this not a problem?

As $z \rightarrow \infty$ $g(z) \rightarrow 1$
 $z \rightarrow -\infty$ $g(z) \rightarrow 0$

For large z (weight), learning becomes very slow.

Instead:-

Rectified linear unit
 Sub gradient



$$g(z) = \max(0, z)$$

$$g'(z) = 1$$

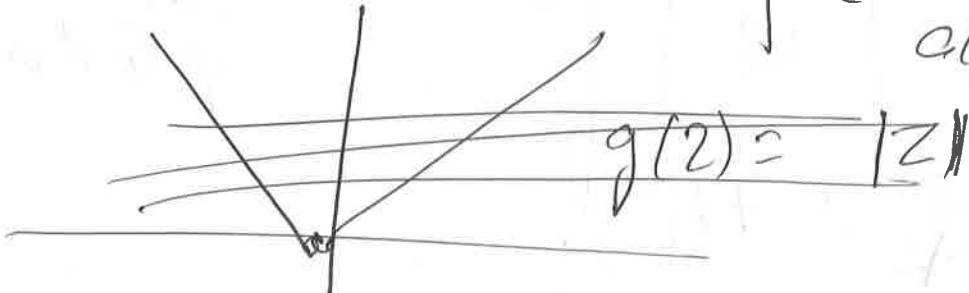
$z > 0$

gradient does not saturate

(when unit is active)

$$g(z) = \max(0, z) + \alpha \min(0, z)$$

↓ Variations

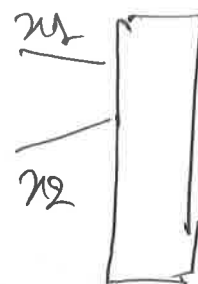


③

Example -

XOR

Not separable



Hidden Layer
(2 units)

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$b = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

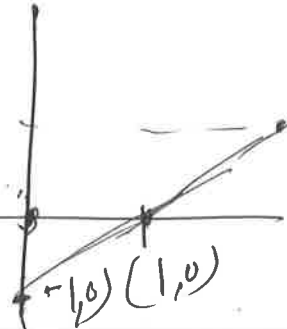
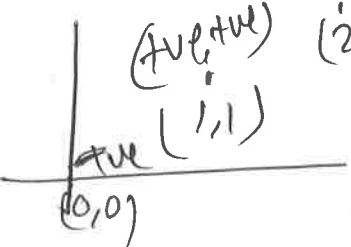
$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & x^{(4)} \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$W^T X + b$$

$$= \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 2 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

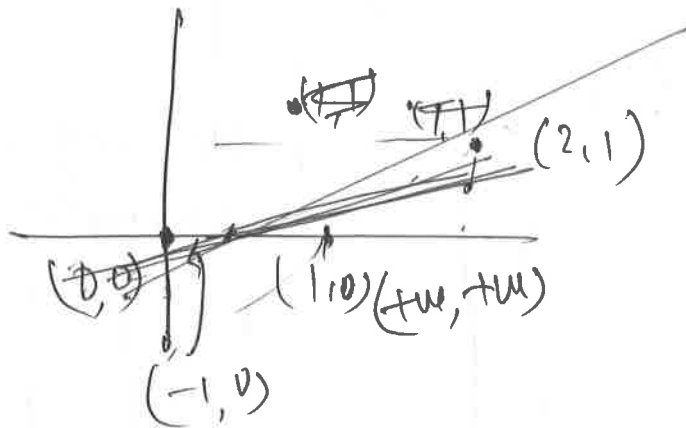
Not linearly separable

$$= \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & x^{(4)} \\ 0 & 1 & 0 & 1 \\ -1 & 0 & 0 & 1 \end{bmatrix}$$



⇒ Now apply $g(z) \equiv \max(0, z)$

$$\Rightarrow g(w^T x + b) = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$



Output Unit:-

$$f^{(0)}(\phi(x; \theta); \theta^{(0)}) \quad g^{(0)}$$

What is form of $f^{(0)}$?

$$\textcircled{1} \quad \theta^{(0)} = (\omega^{(0)}, b^{(0)})$$

$$f^{(0)}(\phi(x; \theta); \theta^{(0)})$$

$$= \omega^{(0)T} \phi(x; \theta) + b^{(0)} = \hat{y}$$

Linear Transformation

~~Not a~~

$$L(y, \hat{y}; \theta) = \|y - \hat{y}\|^2$$

$$= \|y - f(x; \theta)\|^2$$

Total Loss:-

$$J(\theta) = \sum_{x \in \mathcal{X}} \|y - f(x; \theta)\|^2$$

Not a very
good metric.
Especially for
classification

All possible
instances

when
 $y \in \{0, 1\}$

Instead Use $f(\theta)$ to define a distribution over y (given x & θ) as a function of θ .

Assume:
 $y \in \{0, 1\}$

$$f(\theta) [\phi(x; \theta); \theta(\theta)]$$

$$= g(\theta) / \omega^{\omega T} \phi(x; \theta) + b(\theta)$$

\Downarrow
Sigmoid
function.

can define Bernoulli
rand distribution

$$p(y=1) = \frac{g(\theta) / (2)}{1 + e^{-2}}$$

$$p(y=1) = g(\theta) / (2)$$

$$p(y=0) = 1 - g(\theta) / (2) = g(\theta) / (-2)$$

General
on nature

$$p(y) =$$

$$g(\theta) [(2y-1)z]$$

Why isn't
good thing
do? $y=1$?
we log
likelihood

Finally:

$$L(y, \hat{y}(\theta)) =$$

$$L(y, \hat{y}(\theta)) = -\log [p(y)] \Rightarrow$$

$$J(\theta) = E_{(x, y \sim \text{Data})} - \log [p_{\theta}(y)]$$

Approximate
by training
loss

Extra

Deep Feedforward Networks

Also known as:- feedforward neural networks
multilayer perceptrons

Goal:- $y = f^*(x)$

Define a mapping $y = f(x; \theta)$ parameters of the model

f should be as close to f^* as possible

Information flows forward

(Feed forward)

Non-linear transformation $\phi(x; \theta)$

Inputs

Hidden layer

Layer

several outputs

output layer

Behavior known only at the output layer

output layer

Linear function

$f(x; \theta, w)$
 $\phi(x; \theta, w)$

Feedback

As opposed to RNNs (recurrent neural networks)

Foundational building blocks for CNNs, RNNs, many other deep learning models.

Network:- network of connections (going forward)
depth of the model

$y = f^{(n)}(f^{(n-1)}(\dots(f^{(1)}(x))\dots))$ n layered network.

compositional structure
network structure

implements a compositional structure.

Neural:- Resembles the functioning of neuron in the brain. Each neuron in a layer acts in parallel.

Each layer:- vector-valued function.

Theme in Deep Learning: - Figure out the transformation

$$\phi(x; \theta)$$

↳ hidden layer

Feature transformation

Design choices:-

Output layer:-

(1) Cost Function $L(y, \phi(x^T W))$

(2) ~~Deep~~ Optimizer

Cost Function $L(y, \phi(x^T W))$
Form of output layer:- softmax

Hidden Layer:-

(Introduce non-linearity)

↳ Activation Function (in each layer)

Finally:- Network architecture

For training:- Backpropagation
(gradient based) & generalizations

Hidden Layer:-

$$h^{(1)} = f^{(1)}(x; W, C)$$

$$\downarrow$$

$$\text{Vector valued Function} = \text{Activation Function} \downarrow \text{Linear function}$$

$$h^{(1)} = g(W^T x + C)$$

(element-wise)

$$W \in \mathbb{R}^{n \times k}$$

$$W^T = \begin{bmatrix} w^T(1) \\ \vdots \\ w^T(k) \end{bmatrix}$$

$$+ \begin{bmatrix} 1 \\ c \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} x \\ 1 \end{bmatrix}$$

k such units