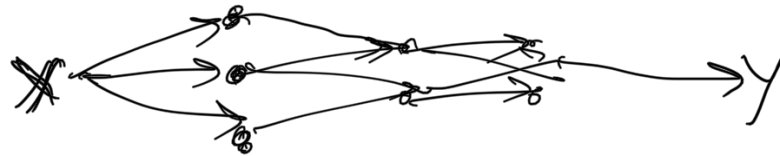


$$y = f(x)$$

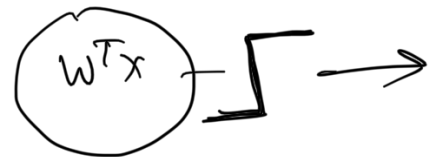
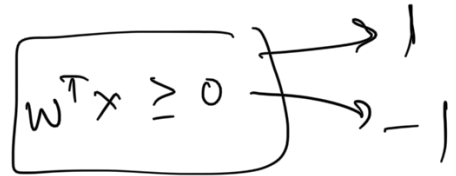
Neural Networks.



$$\underline{\underline{W^T X}}$$

Thresholded
perception

$$\begin{aligned} &w_1 x_1 \\ &+ w_2 x_2 \\ &+ w_3 x_3 \\ &+ w_4 x_4 \\ &+ w_5 x_5 \end{aligned}$$

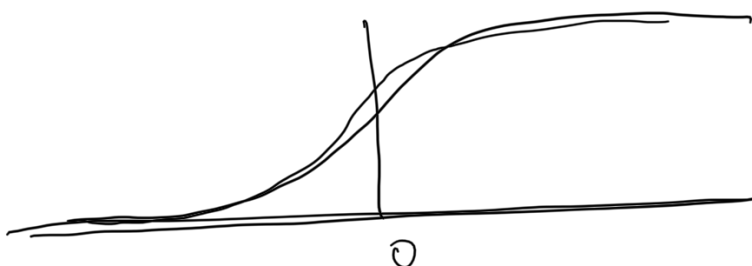


Unit

Layer

Sigmoid unit

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



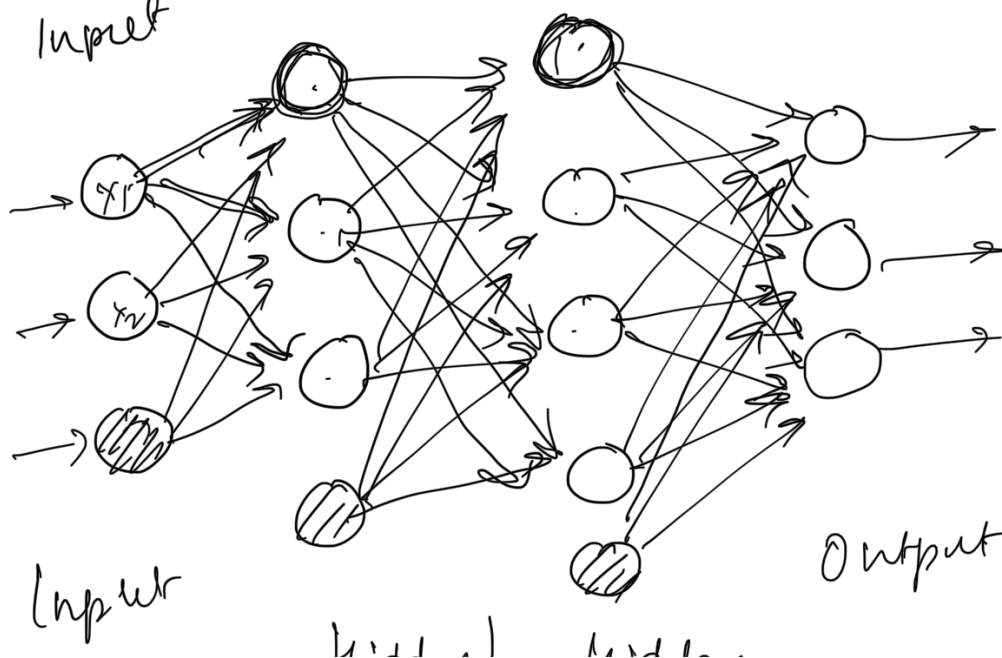
March 12 Friday

- PA2 is out

- Mid-term next Friday

$D=2$

Input



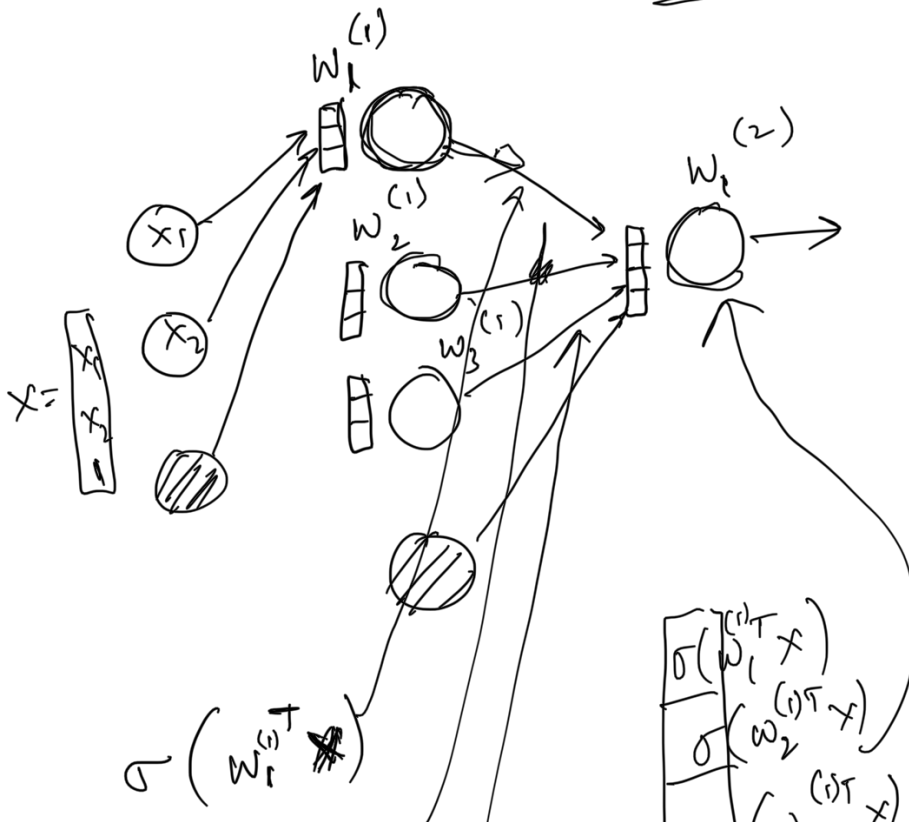
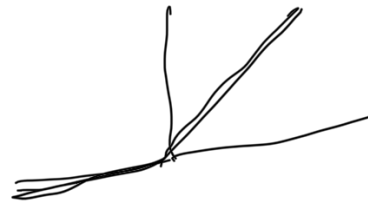
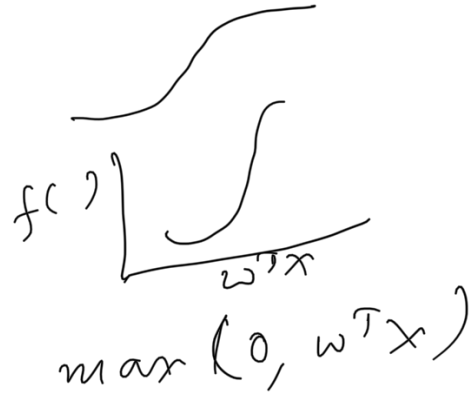
hidden hidden

of units in a layer \equiv width of a layer

Sigmoid

tanh

Relu

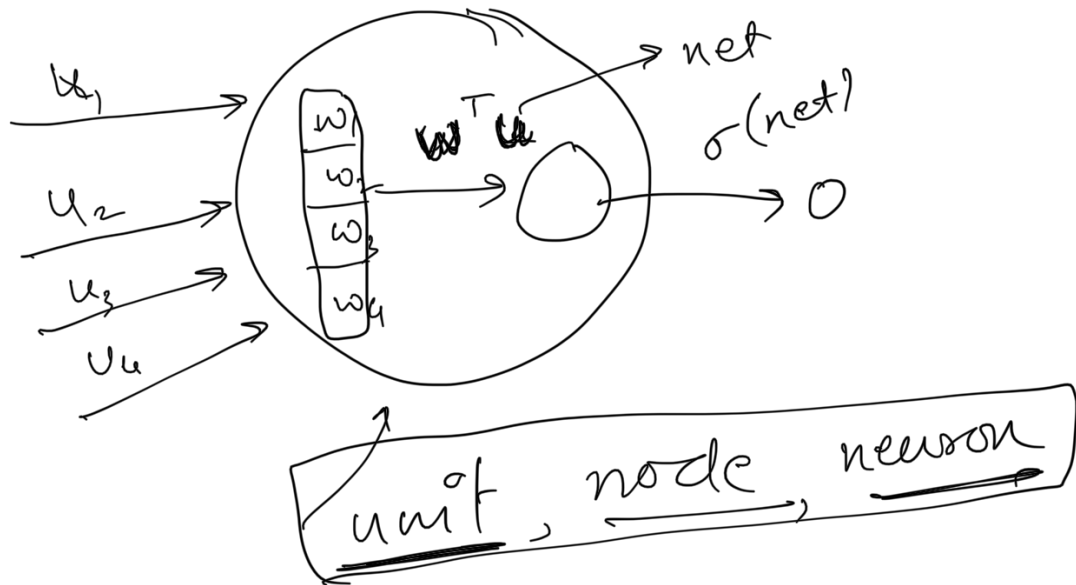


$$\sigma(w_2^{(1)T} x)$$

$$\sigma(w_3^{(1)T} x)$$

$$\sigma(w_3)$$

$$1$$



Sigmoid

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

Tanh

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

Simple example

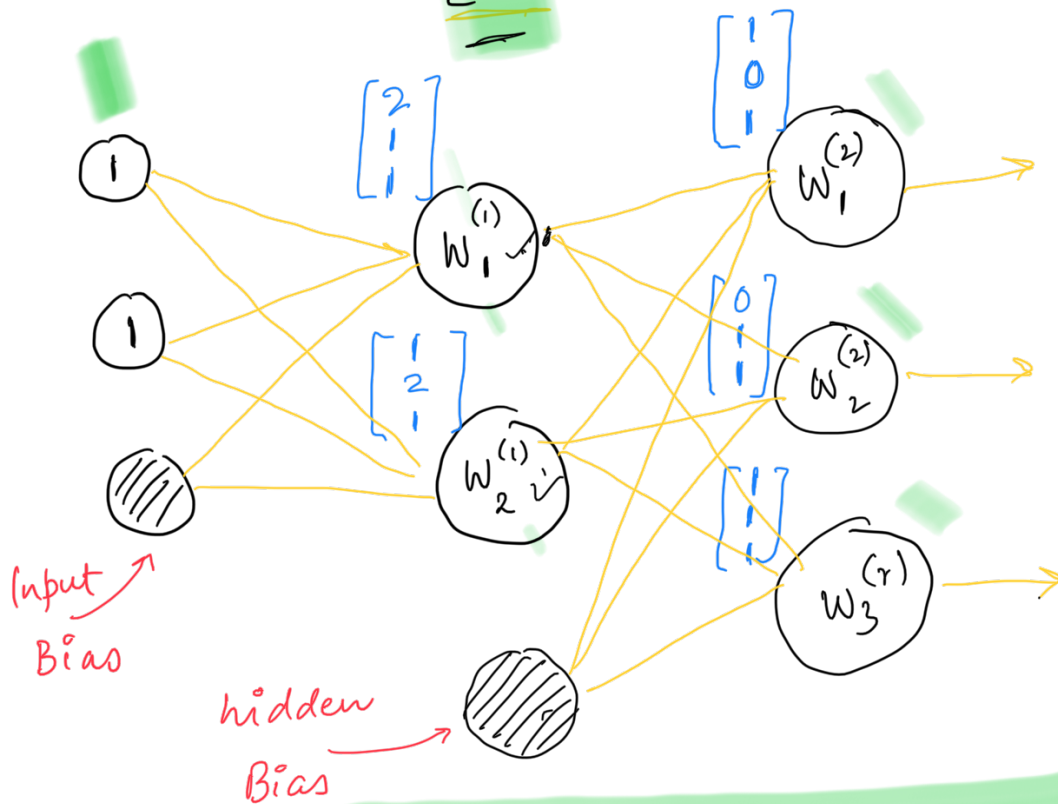
$D = 2$ (2-d data)

$M = 2$

$K = 3$ (3- outputs / classes)

Assume sigmoid activation

Data: $\mathbf{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$



At hidden unit 1

$$\text{net}_1^{(1)} = \mathbf{w}_1^{(1)\top} \mathbf{x}$$

$$= [2 \ 1 \ 1] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 4$$

$$z_1 = \sigma(\text{net}_1^{(1)}) = \sigma(4)$$

$$= 0.98$$

At hidden unit 2

$$\text{net}_2^{(1)} = w_2^{(1)} x$$

$$= [1 \ 2 \ 1] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 4$$

$$z_2 = \sigma(\text{net}_2^{(1)}) = \sigma(4)$$

$$= 0.98$$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ 1 \end{bmatrix}$$

At output unit 1

$$\text{net}_1^{(2)} = w_1^{(2)T} \mathbf{z}$$

$$= [1 \ 0 \ 1] \begin{bmatrix} 0.98 \\ 0.98 \\ 1 \end{bmatrix} = 1.98$$

$$o_1 = \sigma(1.98) = \frac{1}{1 + \exp(-1.98)}$$

$$= 0.88$$

At output unit 2

$$\text{net}_2^{(2)} = w_2^{(2)T} \mathbf{z}$$

$$net_2 = \begin{bmatrix} 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.98 \\ 0.98 \\ 1 \end{bmatrix} = 1.98$$

$$O_2 = \sigma(1.98) = 0.88$$

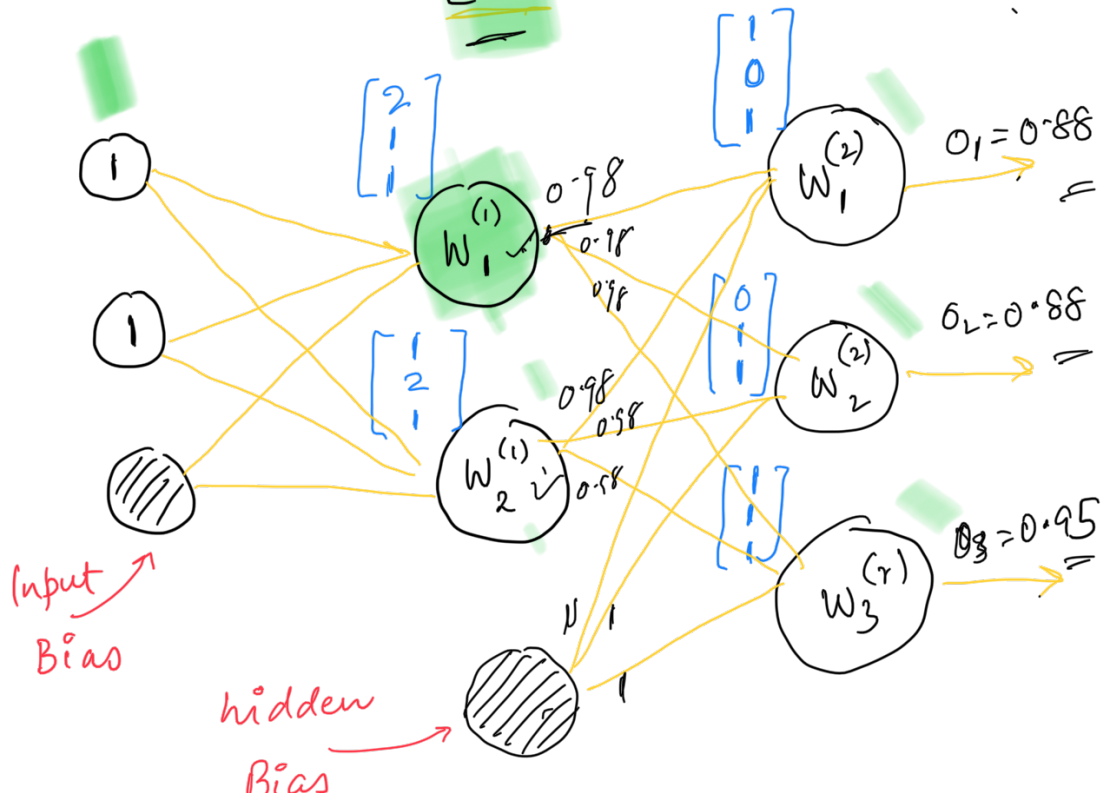
At output unit 3

$$net_3^{(2)} = W_3^{(2)T} Z = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0.98 \\ 0.98 \\ 1 \end{bmatrix}$$

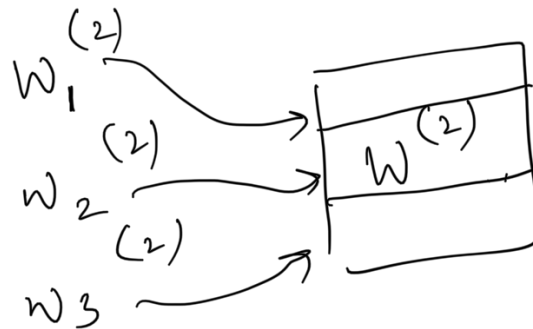
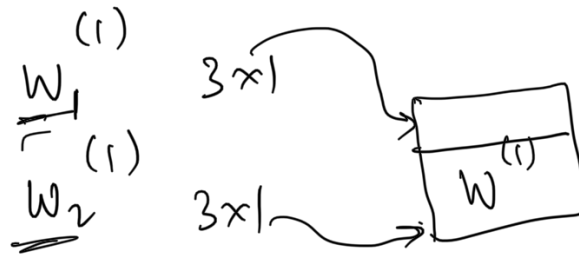
$$= 2.96$$

$$O_3 = \sigma(2.96) = 0.95$$

Data: $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$



softmax output : $\left(\frac{o_1}{o_1 + o_2 + o_3}, \frac{o_2}{o_1 + o_2 + o_3}, \frac{o_3}{o_1 + o_2 + o_3} \right)$



$\sigma \left(\underline{w}^{(1)} \times \underline{x} \right)$

$2 \times 3 \quad 3 \times 1$

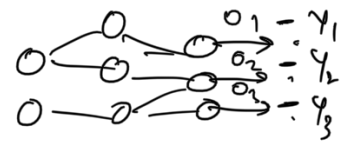
$= \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

$\underline{z} = \begin{bmatrix} 2 \\ z_2 \\ 1 \end{bmatrix}$

$\sigma \left(\underline{w}^{(2)} \underline{z} \right) = \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix}$

$J(w^{(1)}, \dots, w^{(2)}, \dots)$

$= \sum_{i=1}^N J_i$



$$J_i = \frac{1}{2} \sum_{l=1}^k (y_{il} - o_{il})^2$$

$\underline{o_{il}}$ → output for the i^{th} training example at output unit l

$$J = \frac{1}{2} \sum_{i=1}^N \sum_{l=1}^k (y_{il} - o_{il})^2$$

y_{il} → true output for the i^{th} example at l^{th} output unit.

E.g. 3 class classifier.

x	D = 2	y
3.7, 4.8		2
3.4, 1.2		1

One-of-k encoding

Dummy encoding

let $K = 3$

2 →

0	1	0
---	---	---

1 →

1	0	0
---	---	---

3 →

0	0	1
---	---	---

Gradient Descent

gradient descent

$J \rightarrow$ is a function of all
the weights.

∂J