

# Lecture Fys- Stk3155/4155, October 12, 2023

# NNs (Neural Networks) part 1

- universal approximation theorem
- components of NNs
  - (i) cost / objective function
  - (ii) Architecture of NN  
= model
    - a) number of hidden layers
    - b) number of nodes in hidden layers

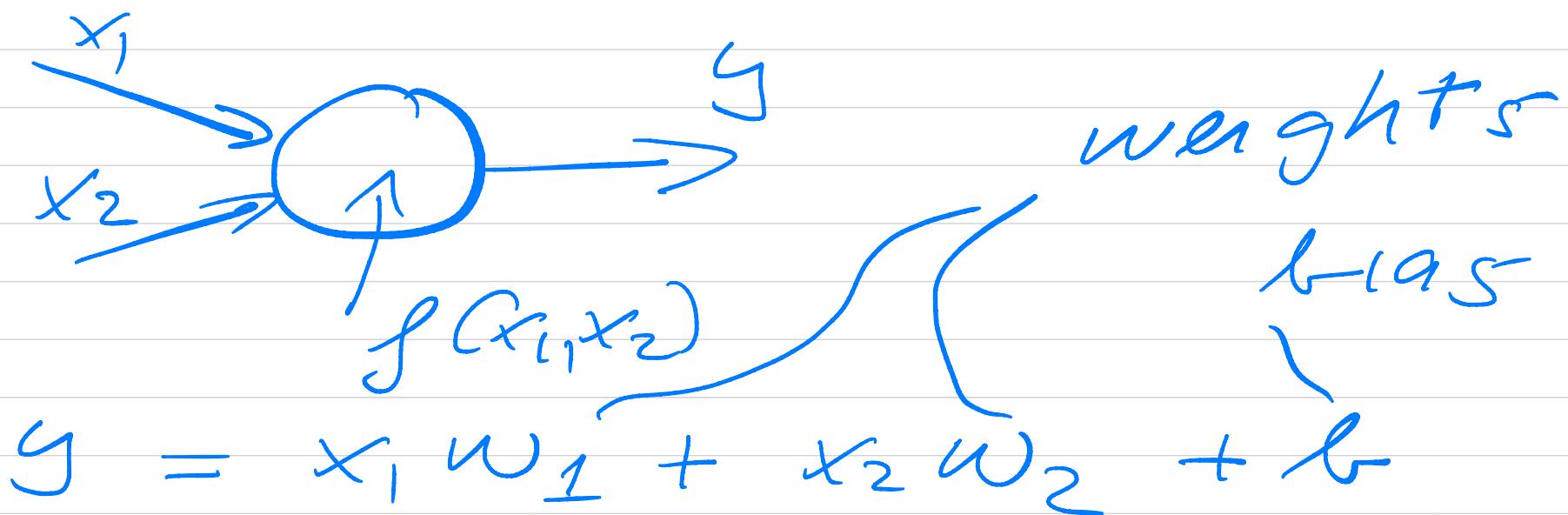
c) type of activation  
functions  $\sigma(z)$

- Back propagation for setting up gradients for training.
- Gradient methods for training (GD, SGD...)
- hyperparameters,  $\lambda$  in  $\lambda \sum \beta_j^2$  etc, learning rates etc etc

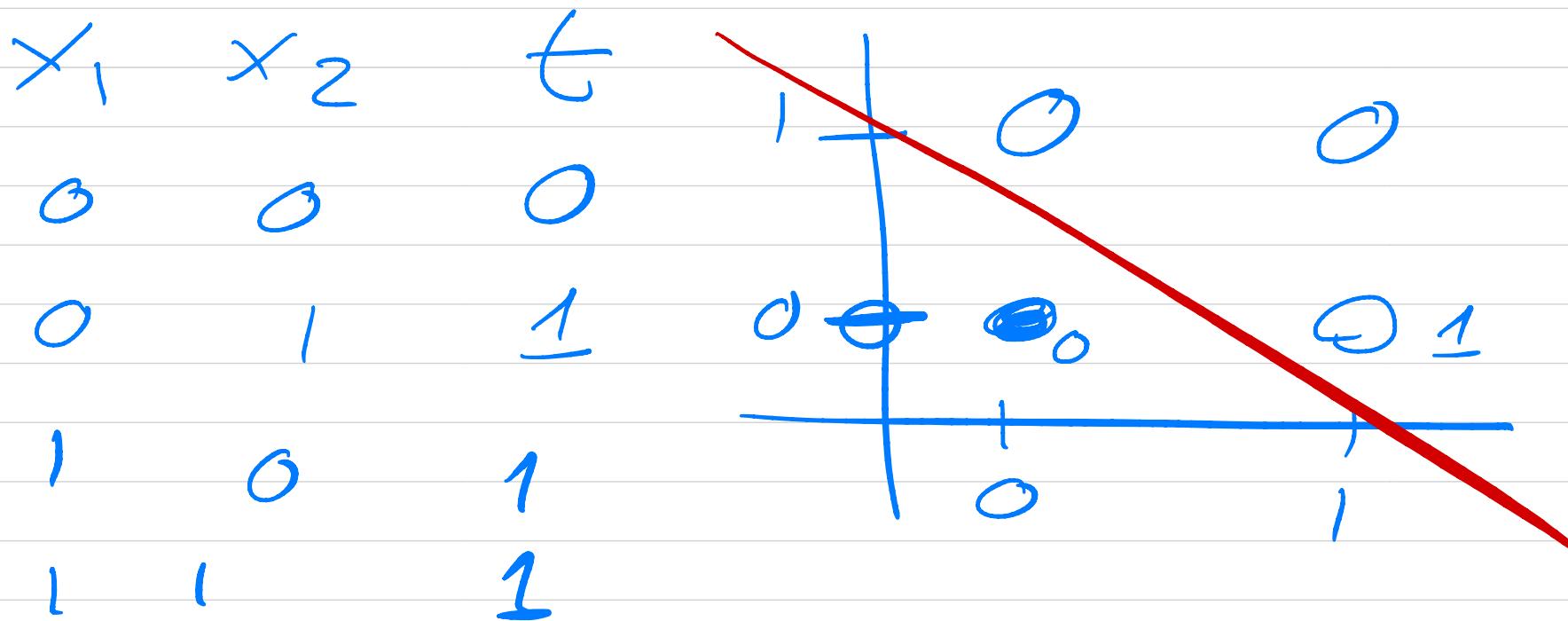
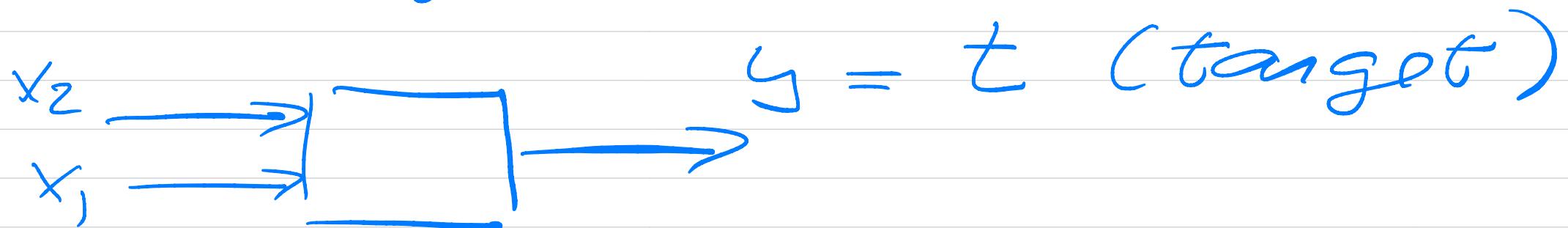
Graphical representation  
Simple network with one node



Example



OR-gate



$$x_1 w_1 + x_2 w_2 + b = y$$

$$w^T = [w_1, w_2]$$

$$x^T = [x_1, x_2]$$

$$y = x^T w + b$$

$$x^T = [0, 0], [0, 1], [1, 0]$$

$$\Theta = \{w, b\} \quad G = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

if  $y < \frac{1}{2}$   
 $y = 0$   
 if  $y > 1/2$   
 $y = 1$

$$X^T X = \begin{bmatrix} 4 & 2 & 2 \\ 2 & 2 & 1 \\ 2 & 1 & 2 \end{bmatrix}$$

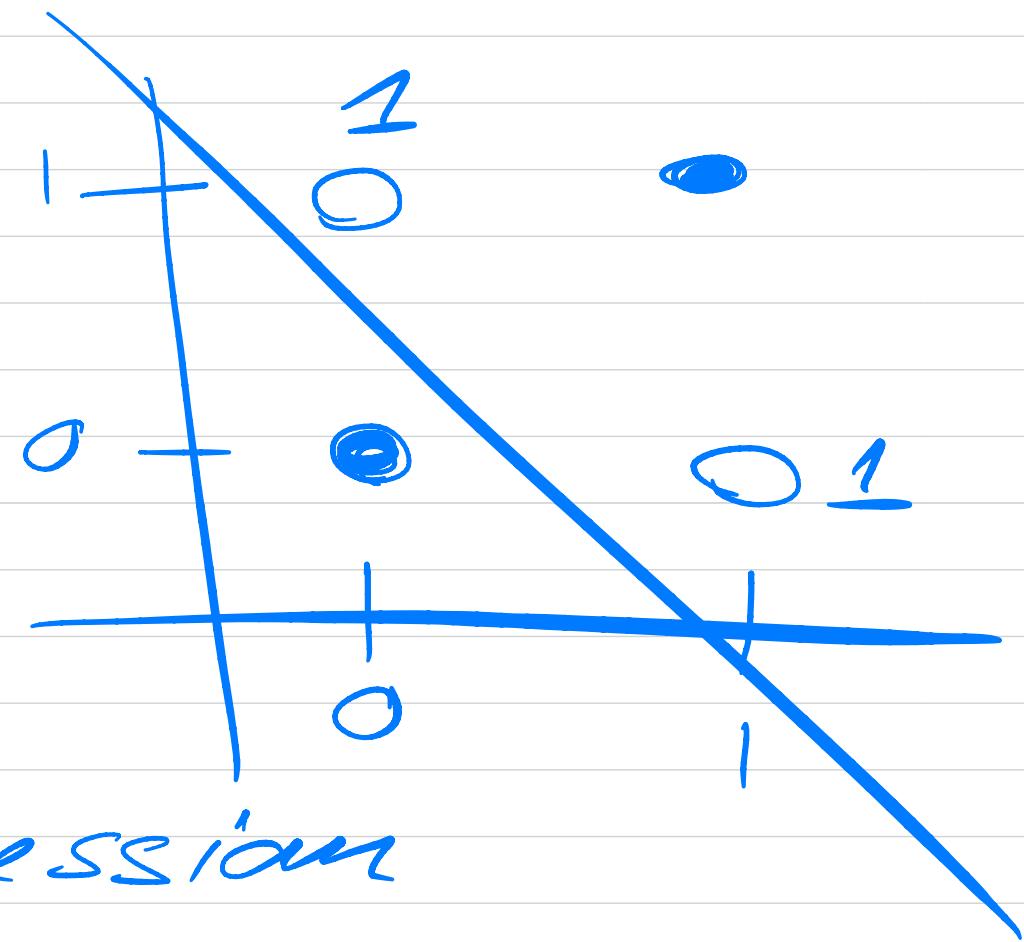
$$\hat{G} = (X^T X)^{-1} X^T y'$$

$t = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$

$$\hat{G} = [1/4, 1/2, 1/2] \quad \tilde{y} = y^T = \left[ \frac{1}{4}, \frac{3}{4}, \frac{3}{4}, \frac{3}{4} \right]$$

XOR

$x_1$	$x_2$	$t$
0	0	0
0	1	1
1	0	1
1	1	0



Linear regression

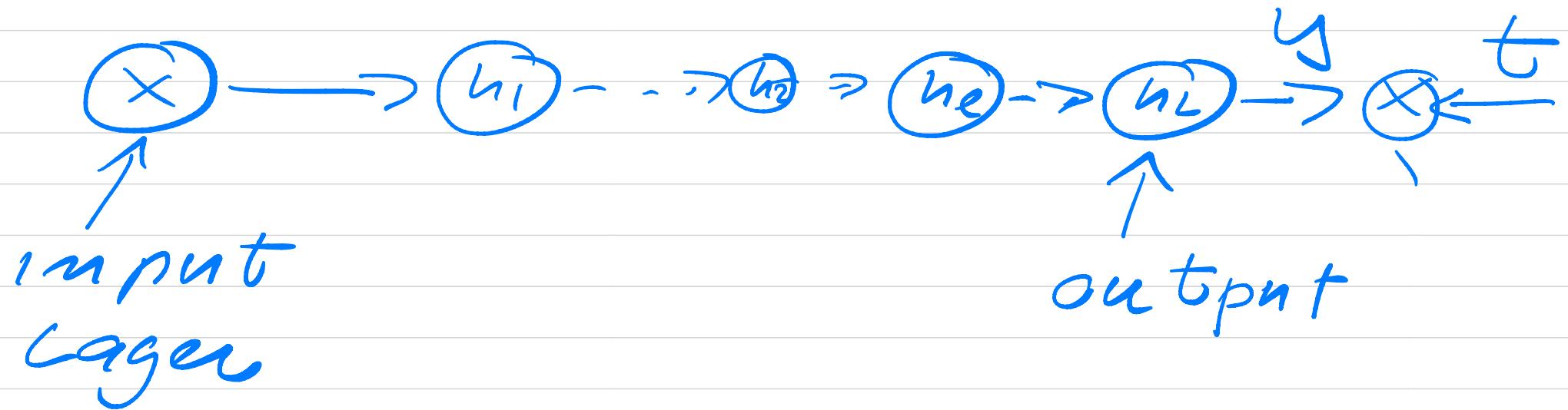
$$y = \left[ \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2} \right]$$

$$t = [0, 1, 1, 0]$$

Universal approximation theorem  
(Gybenko 1989, Hornik 1991)

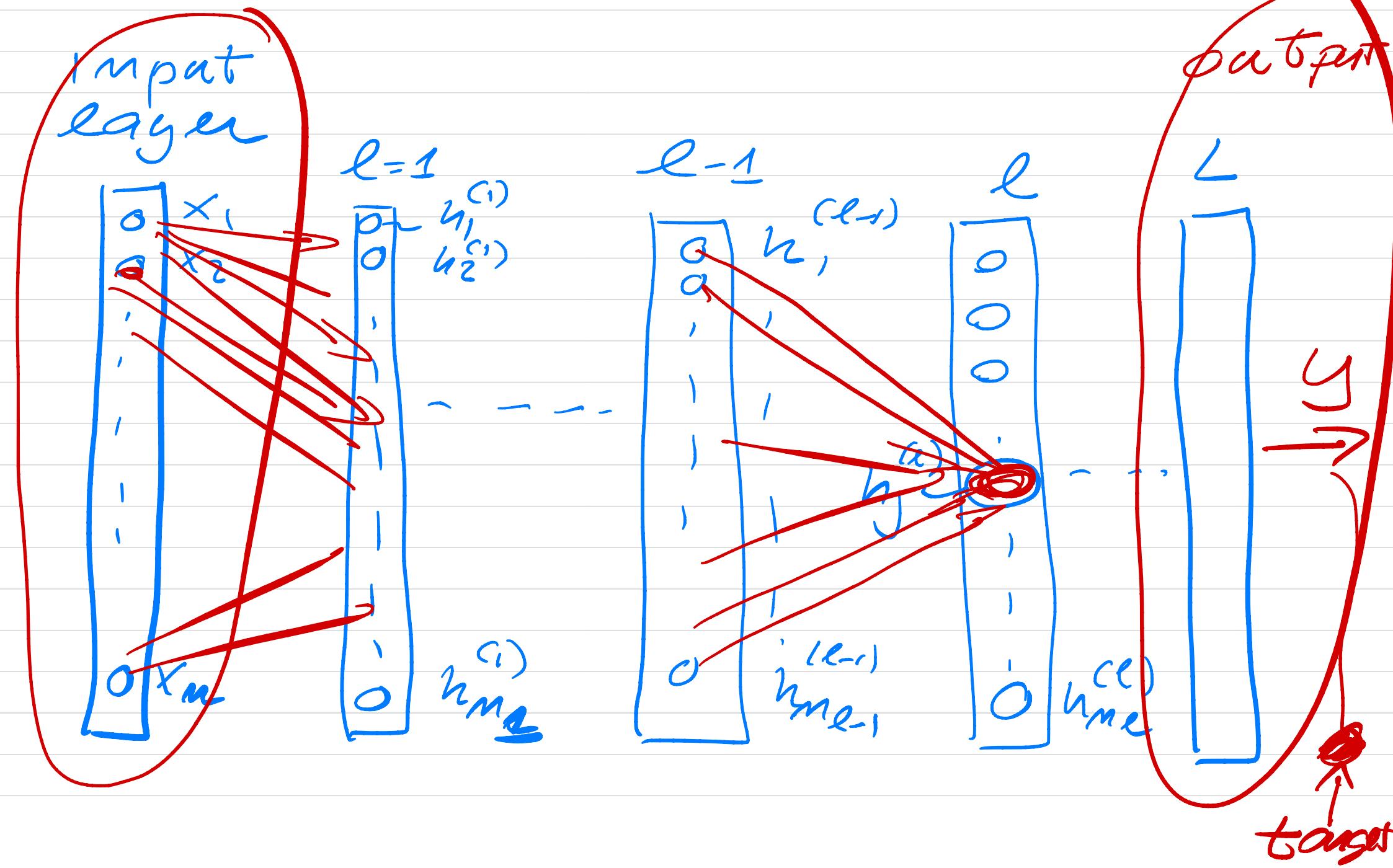
Given function  $F \in [0,1]^d$   
and  $\varepsilon > 0$ , there exists/ $\exists$   
a one hidden-layer NN  
 $f(x; \Theta)$  of the form with  
 $\Theta = \{w, b\}$   $w \in \mathbb{R}^{m \times n}$   
and  $b \in \mathbb{R}^m$  for  
where  $|f(x; \Theta) - F(x)| < \varepsilon$   
for all  $x \in [0,1]^d$

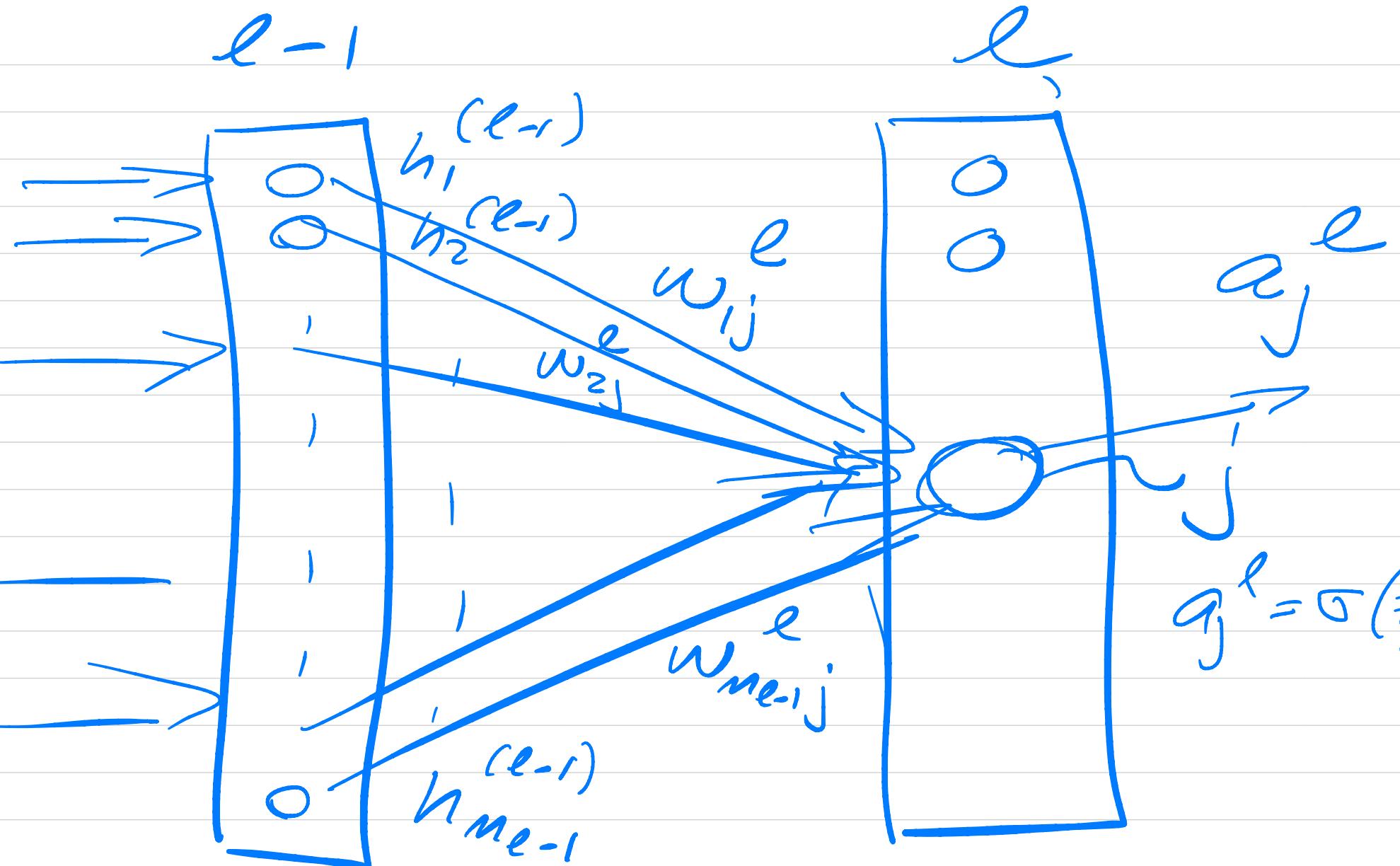
# Back propagation Algo



- Feed Forward stage
- Back propagation stage

$$G = \{ w, b \}$$



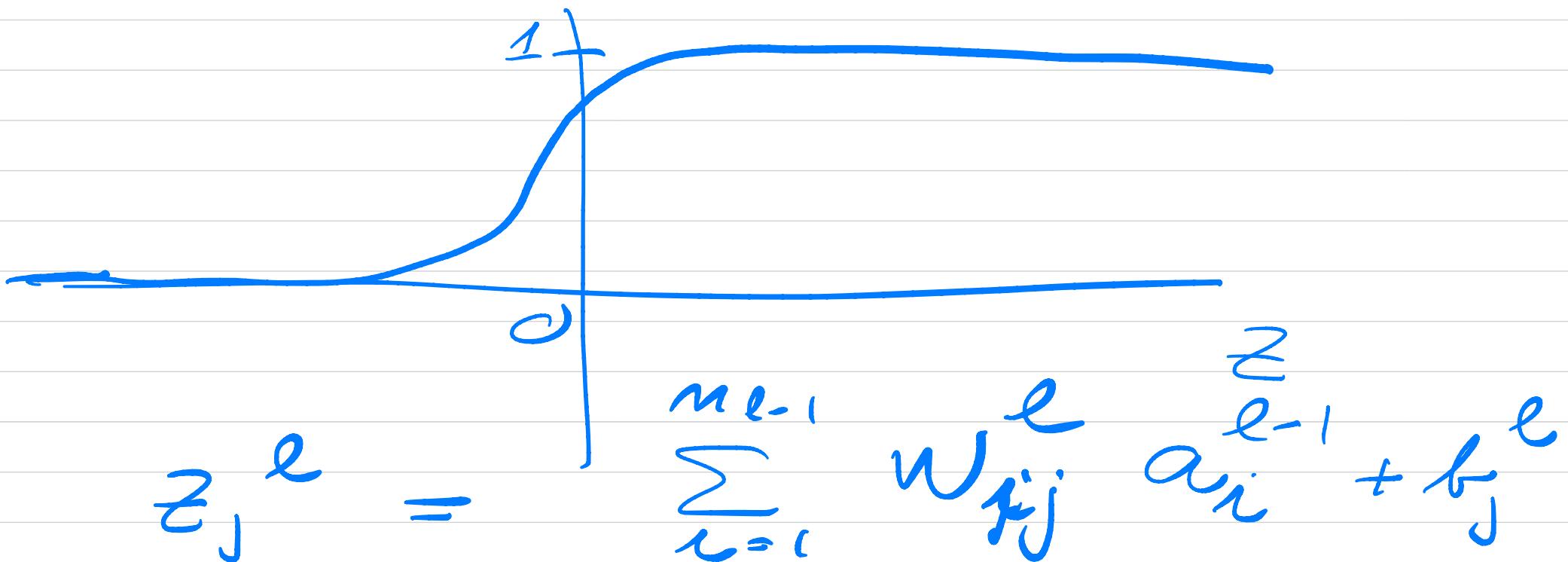


input to  $j$  in layer  $l$   
 $z_j^l$

$$a_j^e = \sigma(z_j^e)$$

Sigmoid is much used

$$\sigma(z_j^e) = \frac{1}{1 + e^{-z_j^e}}$$



$$z^e = [W^{(e)}]^T a^{(e-1)} + b^{(e)}$$

$W$  is a matrix containing the unknown parameters to train.

$b$  is so-called bias

Need a cost function.

Here  $C(e) = \frac{1}{2} \sum_{i=1}^{m_L} (y_i - t_i)^2$

$$z_j^e = \sum_{l=1}^{m_{e-1}} w_{ij}^e a_l^{e-1} + b_j^e$$

$$a_l^{e-1} = \tau(z_l^{e-1})$$

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

| in examples  
new

Some intermediate steps

$$\frac{\partial z_j^e}{\partial w_{ij}^e} = a_l^{e-1} \cdot \frac{\partial z_l^e}{\partial a_l^{e-1}} = w_{ij}^e$$

$$\frac{\partial a_l^e}{\partial z_j^e} = \tau(z_j^e)(1 - \tau(z_j^e))$$

$$a_l^e (1 - a_l^e)$$

we start at  $\ell = L$  (output)

$$C(e) = \frac{1}{2} \sum_i (a_i^L - t_i)^2$$

$$\frac{\partial C}{\partial w_{jk}^L} = (a_j^L - t_j) \frac{\partial a_j^L}{\partial w_{jk}^L}$$

$$\frac{\partial a_j^L}{\partial w_{jk}^L} = \frac{\partial a_j^L}{\partial z_j^L} \frac{\partial z_j^L}{\partial w_{jk}^L} = a_j^L(1 - a_j^L) a_k^{L-1}$$

$$\frac{\partial C}{\partial w_{jk}} = (q_j^{(l)} - t_j) q_j^{(l)} (1 - q_j^{(l)}) \underbrace{q_k^{(l-1)}}_{}$$

from  
specific  
activation  
function

$$\frac{d \sigma(z)}{d z}$$

$$\sigma_j^{(l)} = \underbrace{q_j^{(l)} (1 - q_j^{(l)})}_{\sigma'(z)} \underbrace{(q_j^{(l)} - t_j)}_{\frac{\partial C}{\partial q_j^{(l)}}}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \delta_j^l a_k^{l-1}$$

$$\delta_j^l = \frac{\partial C}{\partial r_j^l} \frac{\partial r_j^l}{\partial z_j^l} = \frac{\partial C}{\partial r_j^l}$$

$$\delta_j^l = \tau'(z_j^l) \frac{\partial C}{\partial a_j^l}$$

expressions for layer  $\underline{l}$   
 what about layer  $-l-$ ?

$$l \Rightarrow e$$

$$\delta_j^e = \frac{\partial C}{\partial z_j^e} = \sum_k \frac{\partial C}{\partial z_k^{e+1}} \frac{\partial z_k^{e+1}}{\partial z_j^e}$$

$$= \sum_k \delta_k^{e+1} \frac{\partial z_k^{e+1}}{\partial z_j^e}$$

$$z_j^{e+1} = \sum_{i=1}^{N_e} w_{ij}^{e+1} a_i^e + b_j^{e+1}$$

$$\Rightarrow \frac{\partial z_k^{e+1}}{\partial z_j^e} = w_{kj}^{e+1} \sigma'(z_j^e)$$

$$\delta_j^e = \sum_k \delta_k^{e+1} w_{kj}^{e+1} \tau'(z_j^e)$$

updates -

$$w_{jk}^e \leftarrow w_{jk}^e - \gamma \delta_j^{e+1} a_k^{e-1}$$

$$\begin{aligned} b_j^e &\leftarrow \tau_j^e - \gamma \frac{\partial C}{\partial \tau_j^e} \\ &= \tau_j^e - \gamma \delta_j^e \end{aligned}$$

# Algorithm

- Define architecture
- Initialize  $\{W, b\} = \emptyset$
- Define nodes, hidden layers + activation functions
- Hyperparameters
  - learning rate + gradient approx
- Define cost function

- set up  $X$
- Perform first feed forward
- continue to layer  $L$

$$a^L = \sigma^L(\sigma^{L-1}(\dots \sigma^1(z^1))\dots)$$

- compute  $\delta_j^L$
- Then Backpropagate

$$l = L-1, L-2, \dots, \underline{1}$$

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} \sigma'(z_j^l)$$

Training (gradient part)

$$w_{jk}^e \leftarrow w_{jk}^e - \gamma \delta_j^{e-1} a_k^{e-1}$$

$$b_j^e \leftarrow b_j^e - \gamma \delta_j^e$$

Continue till convergence