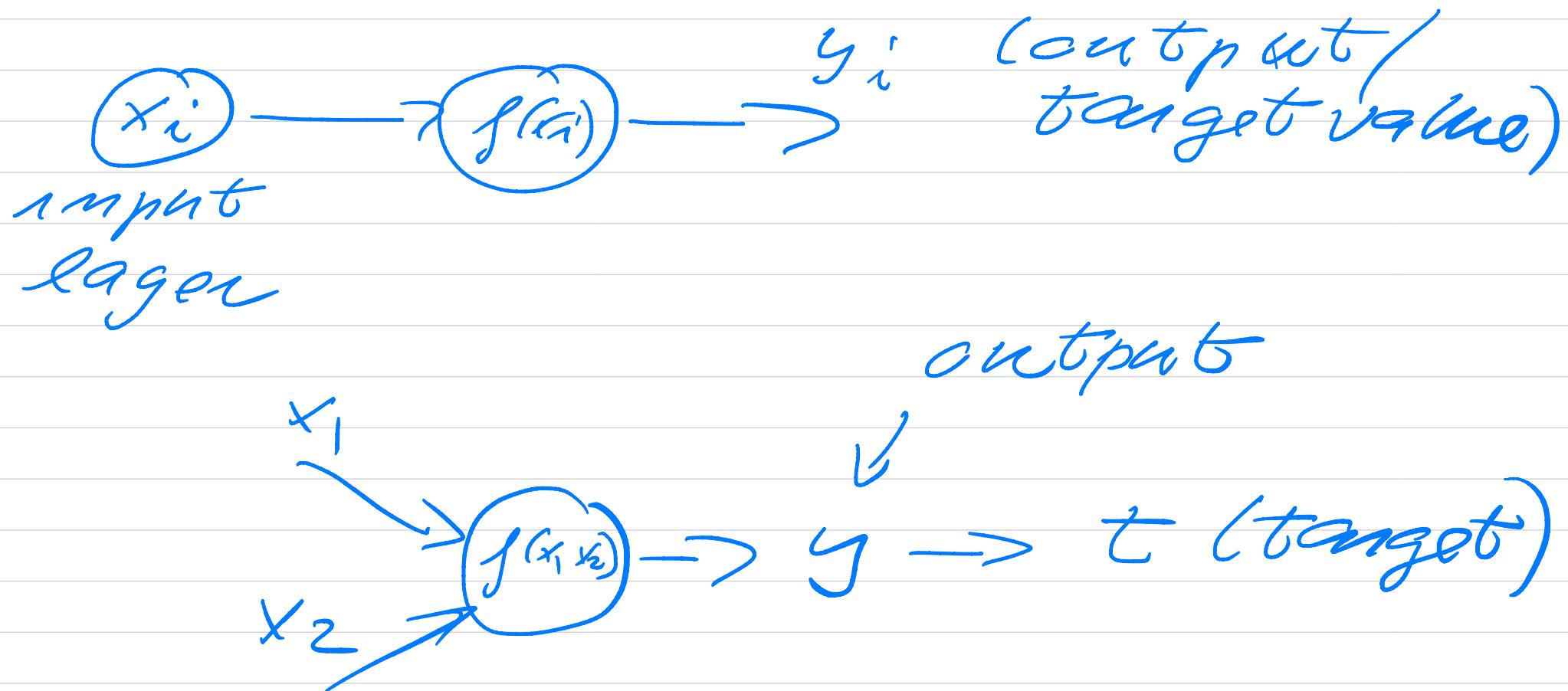


**Erasmus+ course  
on Machine  
Learning, lecture  
November 20, 2023**

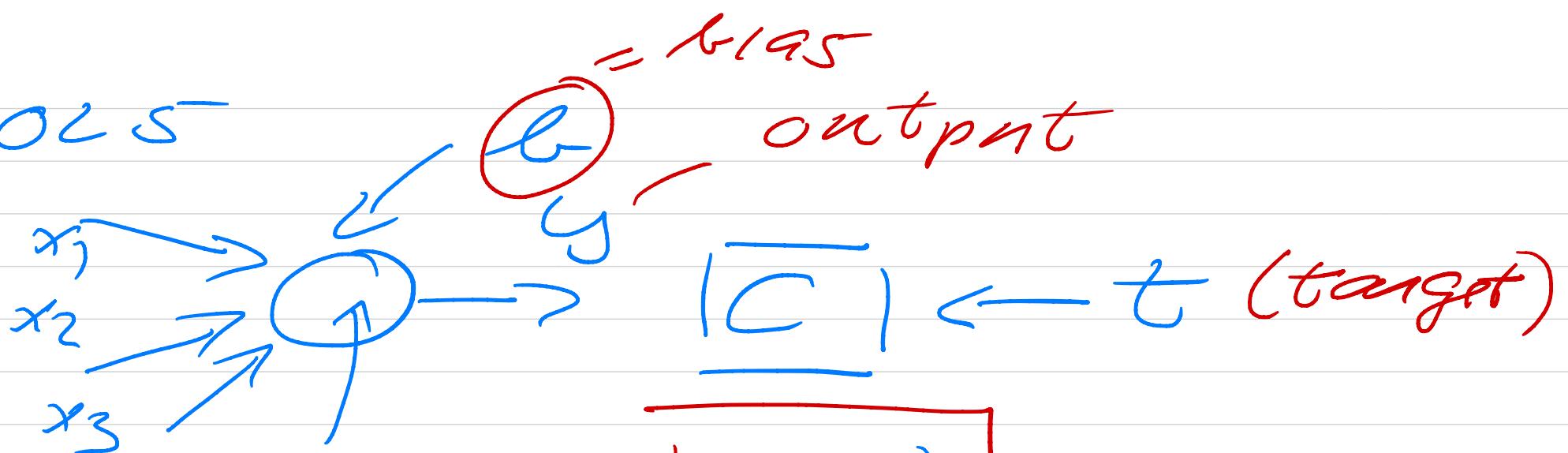
# Neural Networks



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

$$w^T = [w_1 \ w_2]$$
$$x^T = [x_1 \ x_2]$$

OLS



$$f(x_1, x_2, x_3) = \boxed{\sigma(x)}$$

$$\bar{x}^T = [x_1 \ x_2 \ x_3]$$

$$\bar{w}^T = [w_1 \ w_2 \ w_3]$$

weights

activation  
junction

$$C(w, b) = \frac{1}{2} \sum_{i=0}^{n-1} (t_i - y_i)^2$$

$$\frac{\partial C}{\partial w} = 0 \quad \text{and} \quad \frac{\partial C}{\partial b} = 0$$

Classification matrice (Binary)

$$t = \{0, 1\}$$

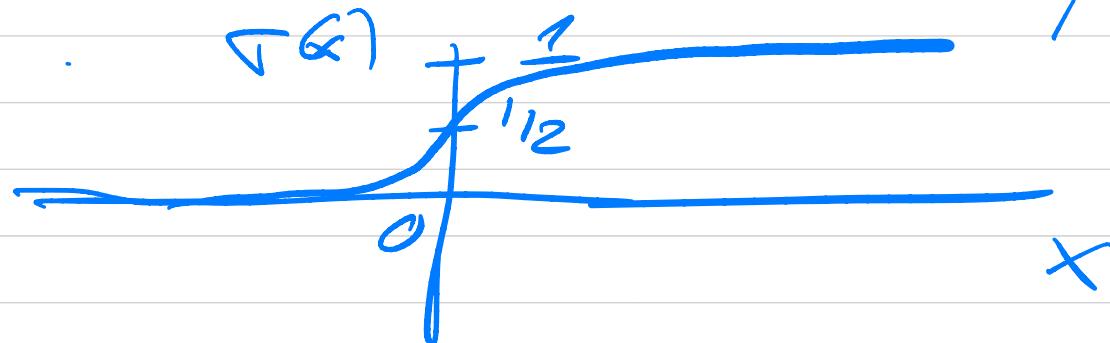
Cross-entropy

$$C(w, b) = - \sum_{i=0}^{n-1} (y_i \log p_i + (1-y_i) \log (1-p_i))$$

$$f(x) \rightarrow p(x)$$

$$\left\{ p_i = p(x_i) \right\}$$

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



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

$$f(x_1, x_2) = y = x_1 w_1 + x_2 w_2 + b$$

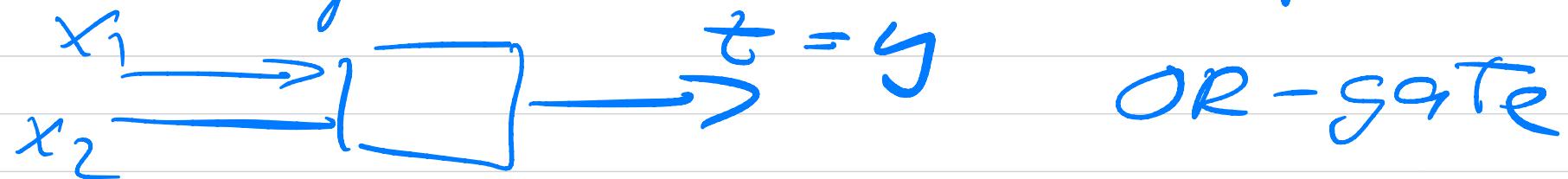
$$f(x_1, x_2) \rightarrow \tau(x_1, w_1, x_2, w_2, b)$$

$$= \tau(x; \Theta)$$

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

New example:

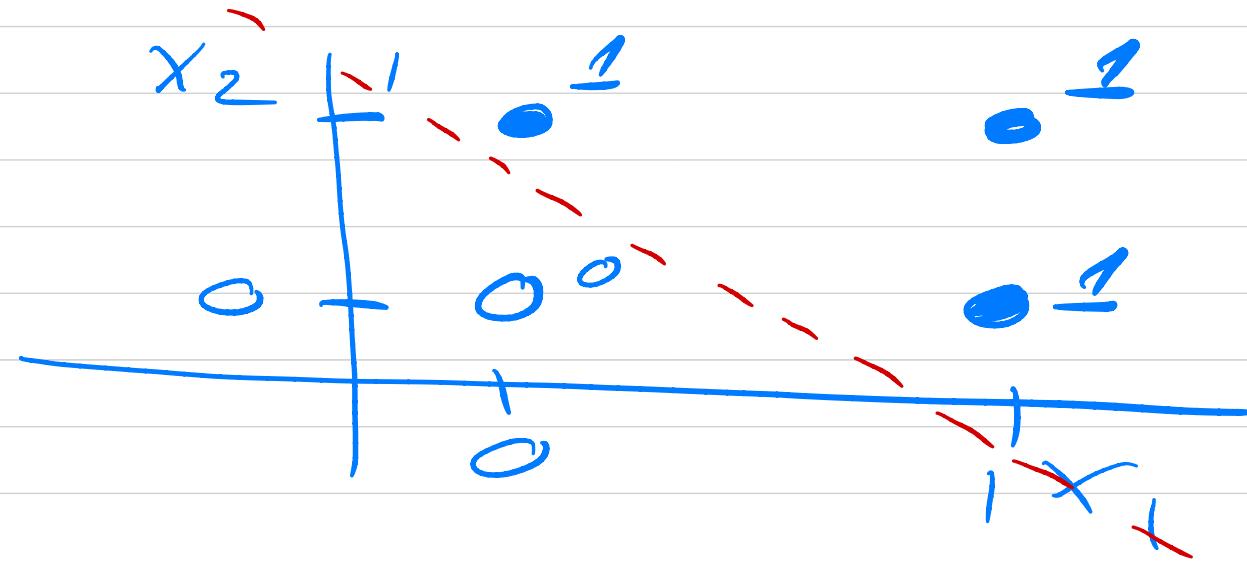
Gates from electronic



$x_1$        $x_2$

		$t$	2 features
		0	$P = 2$
$4$ values $n = 4$	0	1	
	1	0	1
	1	1	1
	1	1	1

$$X^T = \{ [\underline{0}, \underline{0}], [\underline{0}, \underline{1}], [\underline{1}, \underline{0}], [\underline{1}, \underline{1}] \}$$



$$\theta = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} \quad \{ \beta \}$$

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

$$\begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} + b = 0$$

$$\begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix} + b = 1, \text{ etc}$$

include intercept

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

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

$$\hat{\theta} = (X^T X)^{-1} X^T t$$

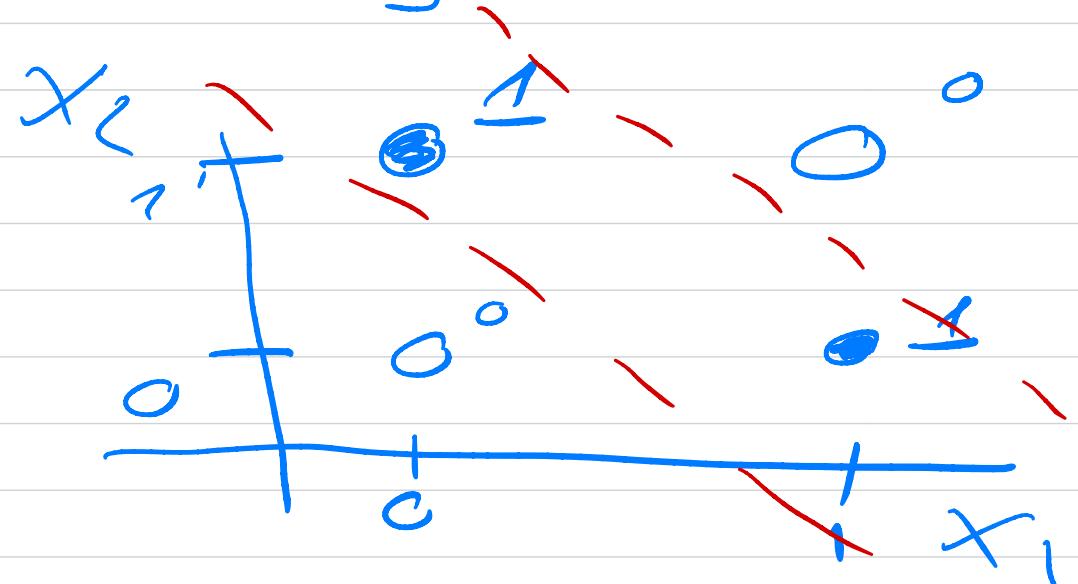
$$y = X \hat{\theta} = \begin{bmatrix} 1/4, 3/4, 3/4, 3/4 \\ 0, 1, 1, 1, 1 \end{bmatrix}$$

if output  $y \geq 1/2$  then  $y = 1$   
 else  $y \leq 1/2$  then  $y = 0$

$$\vec{e}^T = [1/q, 1/z, 1/z]$$

### XOR

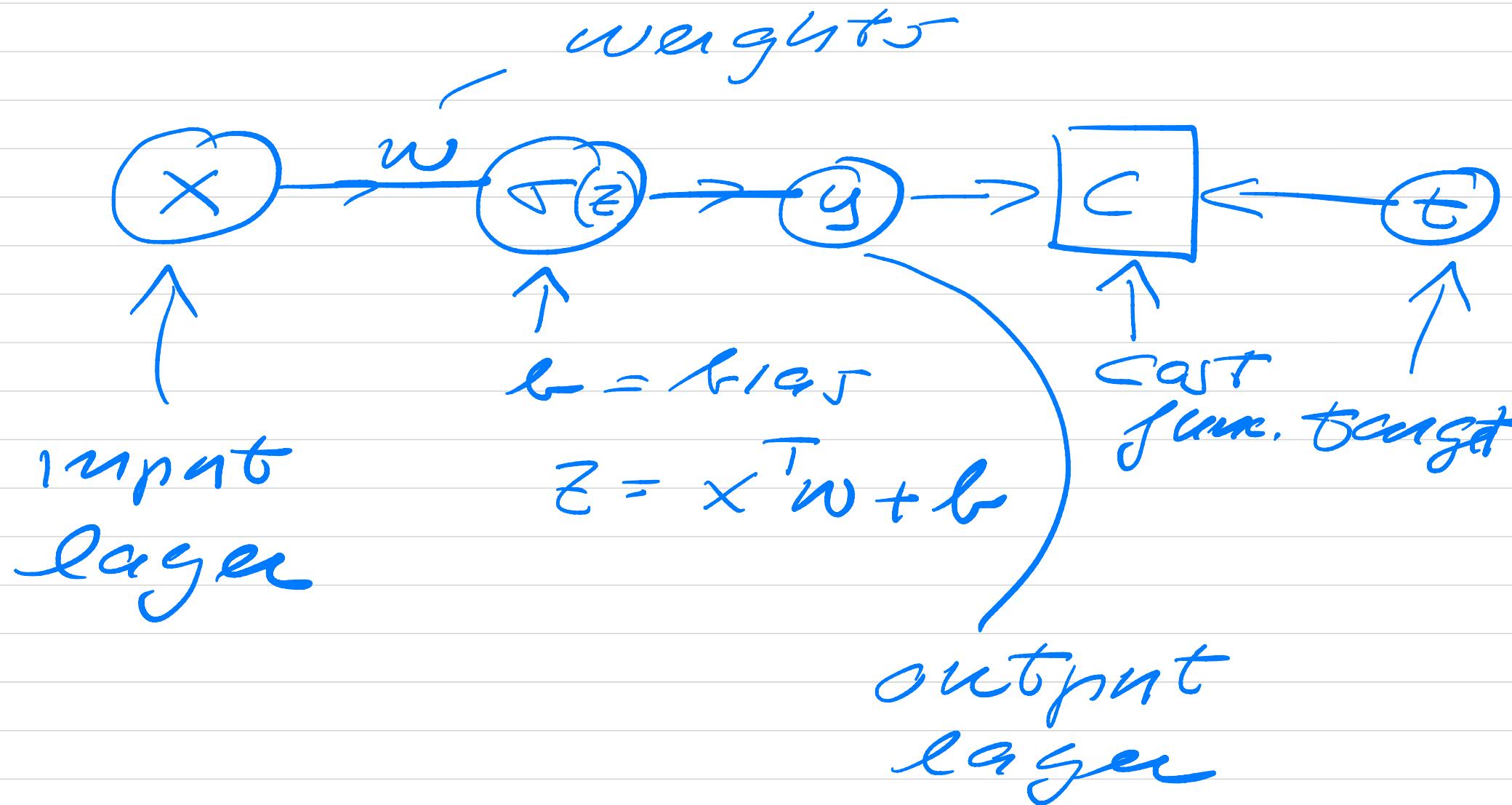
$x_1$	$x_2$	$t$
0	0	0
0	1	1
1	0	1



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

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

with one input and one output layer, we cannot reproduce the XOR gate



The universal approximation theorem states that given a function  $F \in [0, 1]^d$  and error  $\epsilon > 0$ , there exists a one-layer (one hidden) neural network  $f(x; \theta)$  with  $\theta = \{w, b\}$

$$w \in \mathbb{R}^{m \times n} \quad \text{and} \quad b \in \mathbb{R}^m$$

for which

$$|f(x; \theta) - F(x)| < \epsilon$$

$$\text{for all } x \in [0, 1]^n$$

any continuous function  
 $y = Fx$ ) supported on unit  
interval can be approximated  
with a given error  $\epsilon$  with  
one hidden layer.

## Model (architecture of network)

- a given number of hidden layers
- a given number of nodes in a hidden layer
- a specific activation function  $\Gamma(z)$

- Sigmoid  $\sigma(z) = \frac{1}{1+e^{-z}}$
- tanh  $\sigma(z) = \tanh(z)$
- ELU  $\sigma(z) \sim \exp(z)$
- ReLU  $\sigma(z) = \max(0, z)$

⋮

Bounded, continuous / with  
continuous derivatives.

## Optimization

- gradient descent method -

- simple GD with constant learning rate  $\eta$
- Stochastic GD with # epochs and mini-batches
- ADAM
- RMSprop
- Adagrad
- ⋮
- Cost function  $c(\theta) = c(w, b)$
- Hyperparameters  $\lambda \|w\|_2^2$   
 $\chi \|w\|_1$

$$\frac{\partial c}{\partial w} = 0 \quad \text{and} \quad \frac{\partial c}{\partial b} = 0$$

Back propagation algo:

$$E = \{W, b\}$$

General structure

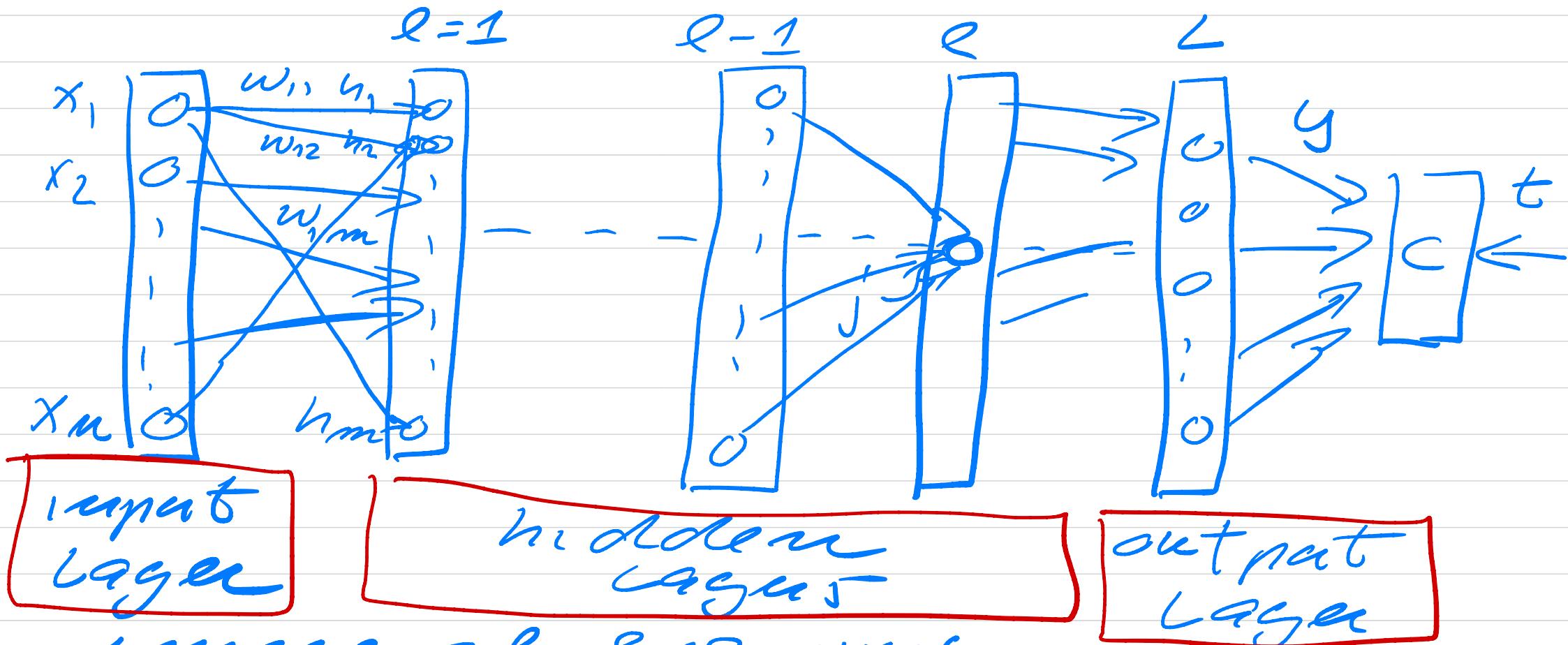
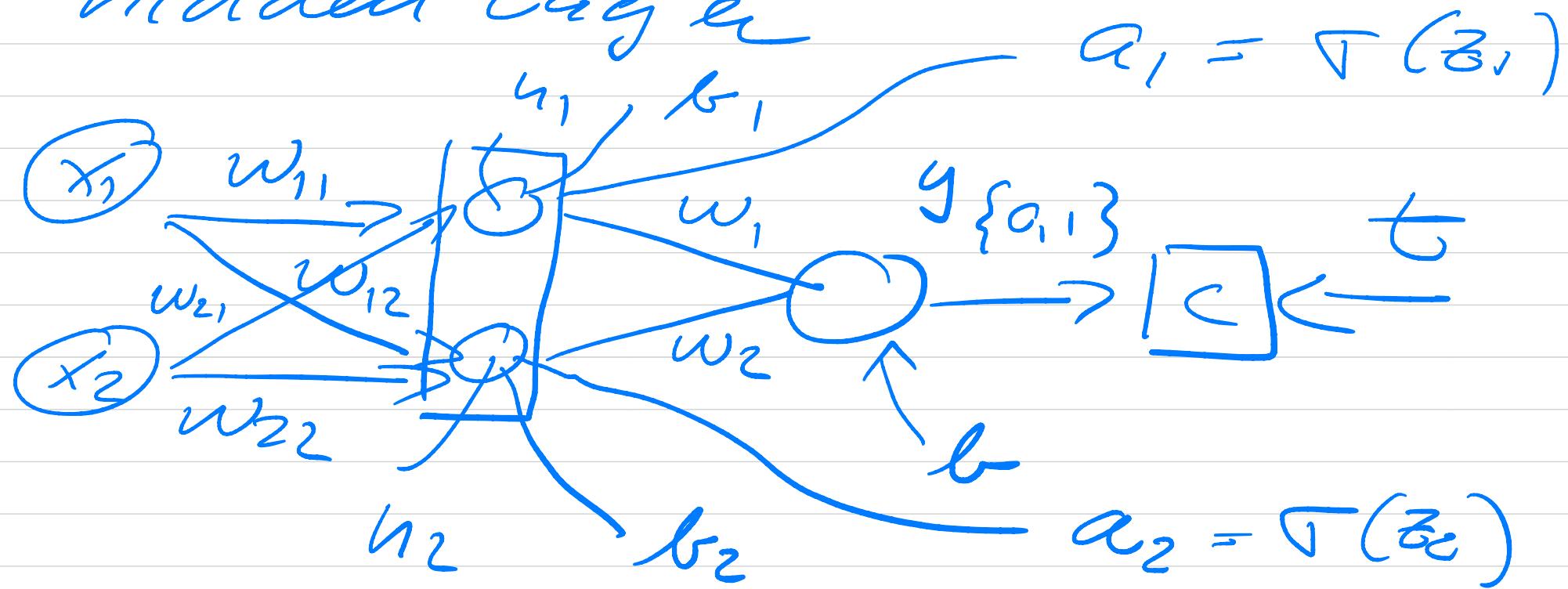


image of  $8 \times 8$  pixel  
64 input nodes

Think of XOR gate with one hidden layer



Design an input matrix

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 2}$$

$$w = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \in \mathbb{R}^{2 \times 2}$$

$$xw \in \mathbb{R}^{4 \times 2}$$

$$xw + b$$

$$z_1 = xw + b_1$$

$$a_1 = \sigma(z_1)$$

$$a_2 = \sigma(z_2)$$

in more detail for layer

$l$  :  $\boxed{z_j^l} = \sum_{i=1}^{N_{l-1}} w_{ij} a_i^{l-1} + b_j^l$

output from node - j - in  
layer - l -

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

$$\sigma = \frac{1}{1 + e^{-z}}$$

intermediate quantities -

$$z_j^l = \sum_i w_{ij}^l a_i^{l-1} + b_j^l$$

$$a_i^{l-1} = \sigma(z_i^{l-1})$$

$$\frac{\partial z_j^l}{\partial w_{ij}^l} = a_i^{l-1} \quad \frac{\partial z_j^l}{\partial a_i^{l-1}} = w_{ij}^l$$

$$\frac{\partial a_j^e}{\partial z_j^e} = \left| \begin{array}{l} a_j^e = \tau(z_j^e) \\ = \frac{1}{1 + e^{-z_j^e}} \end{array} \right.$$

$$= \tau(z_j^e)(1 - \tau(z_j^e))$$

$$= a_j^e(1 - a_j^e)$$

$$c = c(x; \epsilon)$$

$$= c(x; \epsilon^1, \epsilon^2, \dots, \epsilon^L)$$

$$\frac{\partial C(\epsilon)}{\partial w_{jk}^l} = ? \quad \frac{\partial C}{\partial z_j^l} = ?$$

assume that  $C = \frac{1}{2} \sum_i (t_i - y_i)^2$

note implicit dependence  
on  $\epsilon$  in  $y_j^l = a_j^l$

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

$$\delta_j^l = \frac{g_j^l(1-g_j^l)(a_j^{l-1} - t_j^l)}{\tau^l(z_j^e) \frac{\partial C}{\partial g_j^l}}$$

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

$$\delta_j^l = \frac{\partial C}{\partial h_j^l}$$

collect all results at output  
layer  $l = L$

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

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

$$\delta_j^l = \frac{\partial C}{\partial b_j^l}$$

$l \rightarrow e$

$$\delta_j^e = ?$$

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

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

$$\bar{z}_j^{e+1} = \sum_i w_{ij}^{e+1} \underbrace{a_i^e}_{+ b_j^{e+1}}$$

$$\frac{\partial z_k^{e+1}}{\partial z_j^e} = \begin{matrix} \\ \uparrow \\ \tau(z_j^e) \end{matrix}$$

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

Set up gradient:

$$w_{jk}^e \leftarrow w_{jk}^e - \eta \delta_j^{e-1} q_k^{e-1}$$

$$b_j^e \leftarrow b_j^e - \eta \frac{\partial c}{\partial b_j^e}$$

$$= b_j^e - \eta \delta_j^e$$

$$e = 1, 2, 3, \dots L-1, L$$

Next week: set up algorithm  
+ code for Neural network

Link Back propagation with  
automatic differentiation.