

Erasmus+ course on
machine learning,
lecture November
27, 2023

Basic elements of NNs

Define model (architecture)

- # of hidden layers, - and nodes
- activation functions, -

Define hyperparameters (shallow case)

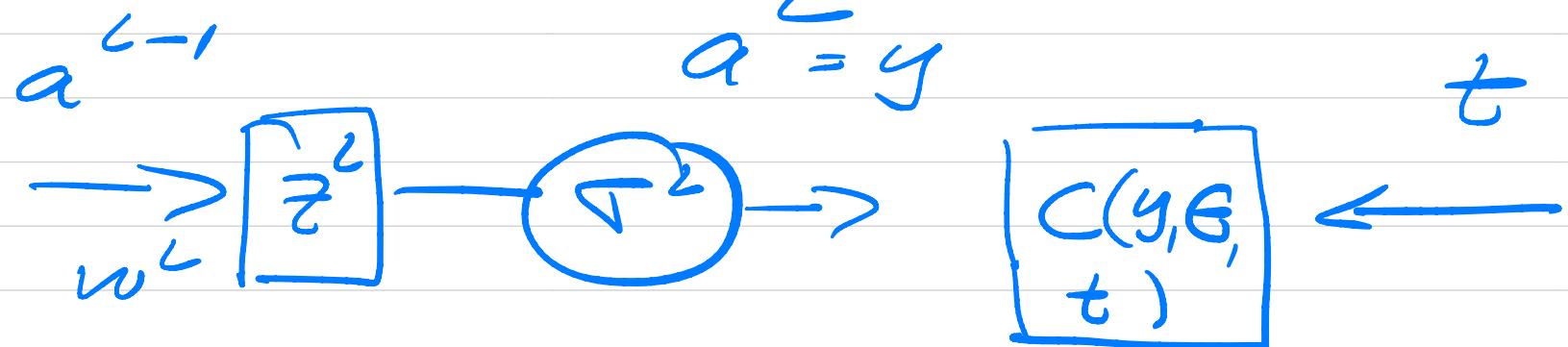
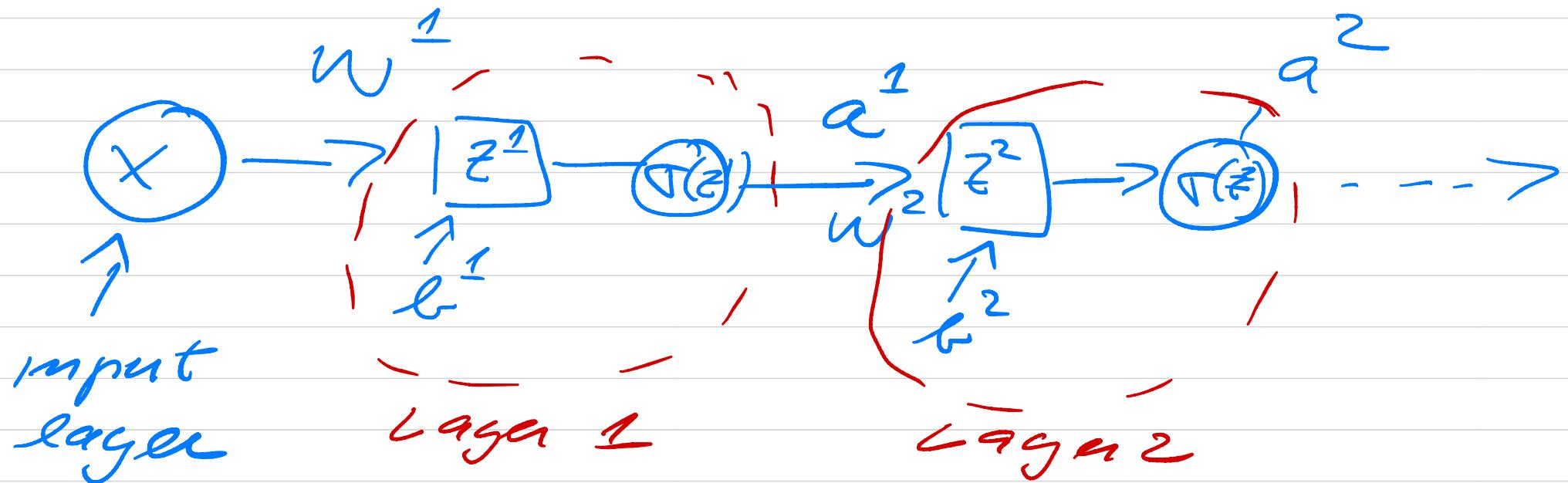
$$\lambda \|\mathbf{w}\|_2^2$$

Gradient descent (GD) algorithm

- stochastic GD
- momentum
- update learning {
ADAM
RMSprop
Adagrad}

Define cost function $C(\Theta)$

$$C = \left\{ \left(\frac{\mathbf{w}^1, b^1}{\theta^1} \right), \left(\mathbf{w}^2, b^2 \right), \dots, \left(\mathbf{w}^L, b^L \right) \right\}$$



$C(y, t ; \epsilon)$ cost/cost/error

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

$$z^1 = w^1 x + b^1 \quad \sigma^1(z') = a^1$$

$$a^L(x; \Theta)$$

Algorithm : Define network

Feed Forward through all

$l=1, 2, \dots L$ layers

$$a^1, a^2, \dots \boxed{a^L}$$

initialize w and b

have $c(a^L; t; \Theta)$

- Backpropagation part
update of w_{jk}^e and b_j^e

$$l = L, L-1, L-2, \dots, 1$$

one Feed Forward pass +
one Backprop pass = 1 iteration

no explicit dependence
on θ

$$\frac{\partial C}{\partial \theta^L} = ?$$

$$C(\theta) = \frac{1}{2} \| t - a^L(x; \theta) \|_2^2 + \lambda \| w \|^2$$

$$\frac{\partial C}{\partial w^l}$$

$$, \quad \frac{\partial C}{\partial b^l}$$

$$C(a^l; t; \theta) = \frac{1}{2} \| t - a^l \|_2^2$$

$$\frac{\partial C}{\partial w^l} = \frac{\partial C}{\partial a^l} \frac{\partial a^l}{\partial z^l} \frac{\partial z^l}{\partial w^l}$$

$$z^l = w^l \cdot a^{l-1} + b^l$$

$$a^l = \sigma^l(z^l)$$

$$\boxed{\delta^l = (\sigma^l(z^l))' \frac{\partial C}{\partial a^l}} \Rightarrow$$

$$\boxed{\frac{\partial C}{\partial w^l} = \delta^l a^{l-1}}$$

$$\frac{\partial C}{\partial b^L} = \delta^L$$

Back propagation algo starts

$$\delta^L \rightarrow \delta_j^L$$

$$l = L, L-1, L-2, \dots, 2, 1$$

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

BASIC gradient update

$$w_{jk}^L \leftarrow w_{jk}^L - \eta \delta_j^L a_k^{L-1}$$

$$b_j^L \leftarrow b_j^L - \eta \delta_j^L$$

Example

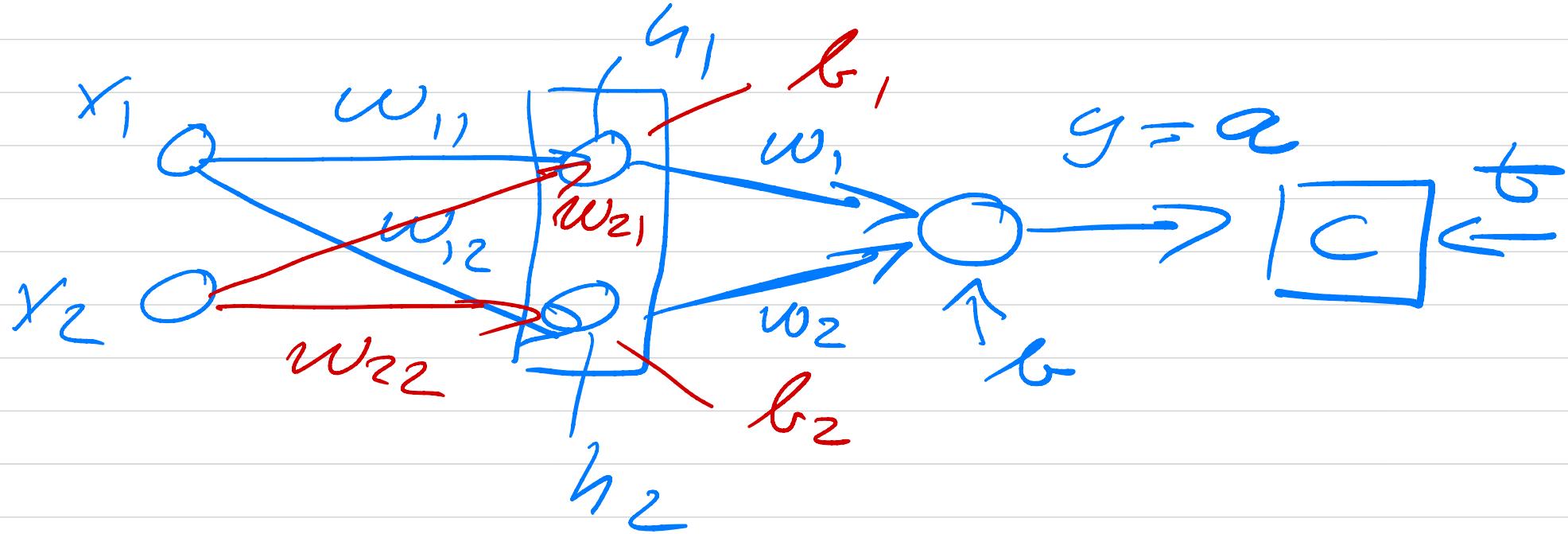
OR - GATE or XOR - GATE

Two features (x_1, x_2)

OR-GATE

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

$$X = \begin{bmatrix} x_1 & x_2 \\ 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$



$$z_1 = w_{11}x_1 + \boxed{w_{21}}x_2 + b_1$$

$$z_2 = \boxed{w_{12}}x_1 + w_{22}x_2 + \boxed{b_2}$$

$$a_1 = \tau(z_1) \quad 1 \quad a_2 = \tau(z_2)$$

$$z = \boxed{w_1}a_1 + a_2 \boxed{w_2} + \boxed{b}$$

$$a = \tau(z)$$

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

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

$$Xw \in \mathbb{R}^{4 \times 2} \Rightarrow$$

$$a_{1,2} \in \mathbb{R}^{4 \times 2}$$

$$w^T = [w_1 \ w_2] \in \mathbb{R}^{1 \times 2}$$

$$w \in \mathbb{R}^{2 \times 1} \Rightarrow aw \in \mathbb{R}^{4 \times 1}$$

Comparison with linear regression

$$\tilde{y}_0 = \beta_0 + \beta_1 x_0 + \beta_2 x_0^2 + \dots + \beta_{p-1} x_0^{p-1}$$

$$\tilde{y}_1 = \beta_0 + \beta_1 x_1 + \dots$$

$$\tilde{y}_2 = \beta_0 + \beta_1 x_2 - \dots$$

, |

$$\tilde{y}_{n-1}$$

$$\tilde{y} = X\beta$$