

FYS-STK 4155, OCT 6, 2022

Gradient methods

$$\beta^{(n+1)} = \beta^{(n)} - H^{-1}(\beta^{(n)}) \times g(\beta^{(n)})$$

$$\approx \beta^{(n)} - \underset{\substack{\uparrow \\ \text{Learning rate}}}{\gamma^{(n)}} g(\beta^{(n)})$$

- efficient gradient evaluation \rightarrow Stochastic GD + momentum GD
- Learning rate tuning
 - linear $\gamma^{(n)}$ in terms of iterations n
 - exponential
 - other

- Adagrad
- RMS prop
- Adam

steepest descent

$$\gamma^{(n)} = \frac{g^{T(n)} g^{(n)}}{g^{T(n)} H^{(n)} g^{(n)}}$$

Adagrad $\gamma^{(n)} \sim \frac{1}{g^{(n)} g^{T(n)}}$

Neural Networks

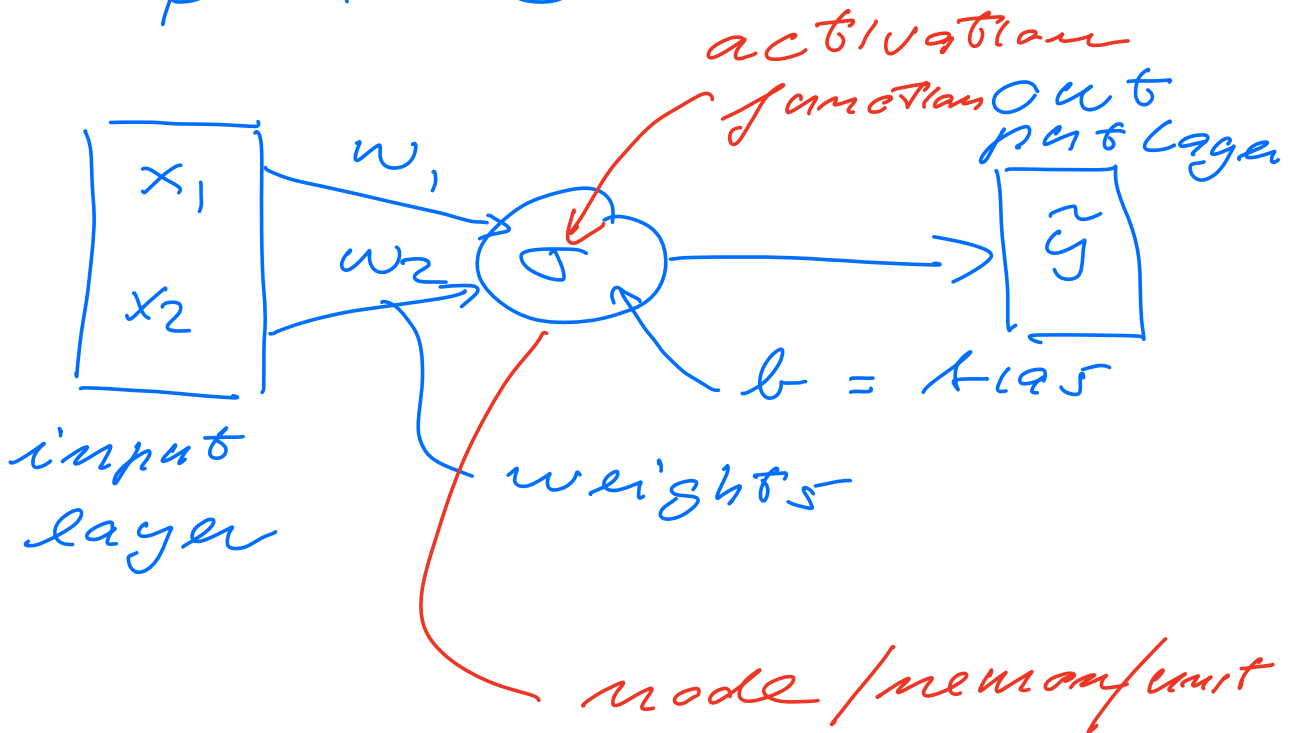
consider a simple case



$y \rightarrow t$ (target)

$$f(x_1, x_2) \rightarrow f(x) \rightarrow \sigma$$

$$\beta \rightarrow \Theta$$



$$\hat{y} = \sigma(\underbrace{x_1, x_2}_{\text{inputs}}, \underbrace{w_1, w_2, b}_{\Theta})$$

x

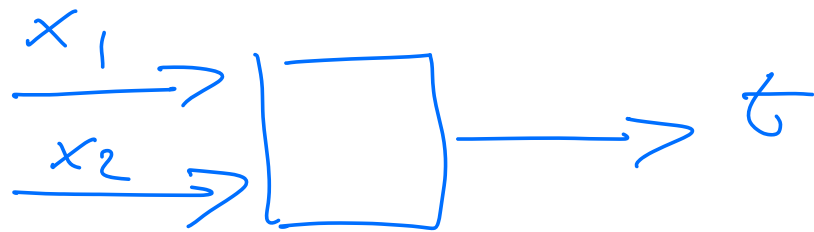
$$x = [x_1, x_2]^T \quad w = [w_1, w_2]^T$$

$$\hat{y} = x^T \underbrace{w + b}_{\Theta}$$

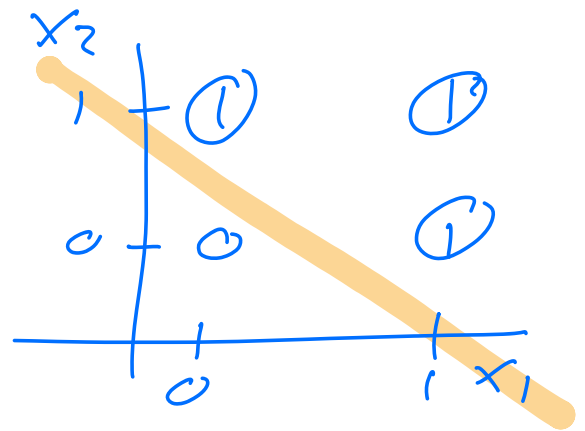
parameters

OR, XOR & AND gate

OR gate



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



$$\tilde{y}_i = x_1(i) w_1 + x_2(i) w_2 + b$$

$$X = \left\{ \begin{matrix} [0, 0]^T \\ [0, 1]^T \\ [1, 0]^T \\ [1, 1]^T \end{matrix} \right\}$$

$\begin{matrix} x_1 & x_2 \end{matrix}$

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

$$t=1 : \begin{bmatrix} c & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + b$$

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

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

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

$$\Theta = [b, w_1, w_2]$$

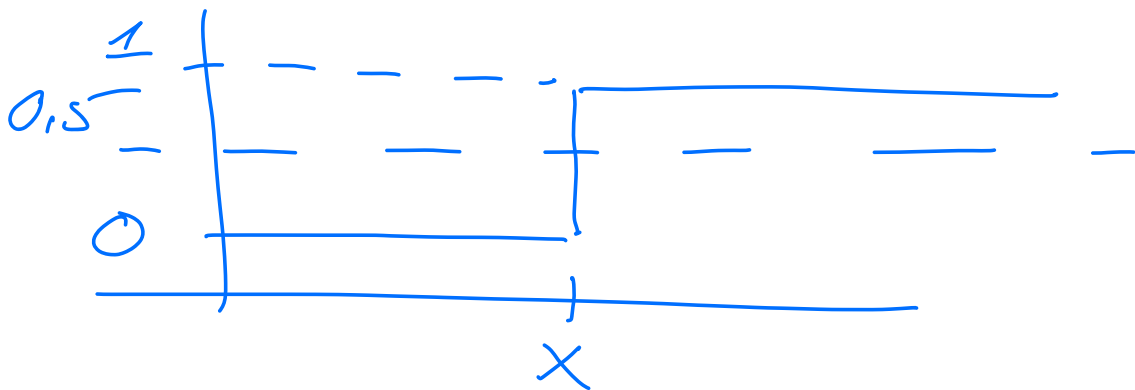
$$(X^T X)^{-1} = \begin{bmatrix} 3/4 & -1/2 & -1/2 \\ -1/2 & 1 & 0 \\ -1/2 & 0 & 1 \end{bmatrix}$$

$$(\beta) = \Theta = (x^T x)^{-1} x^T y$$

$$y = [0 \ 1 \ 1 \ 1]^T$$

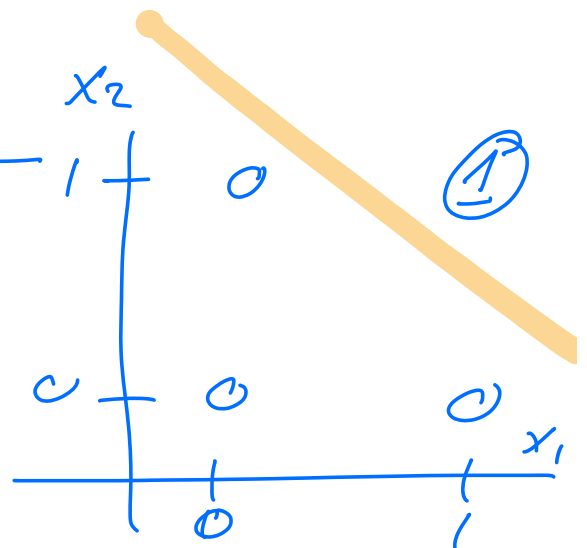
$$\Theta = \left[\frac{1}{4}, \frac{1}{2}, \frac{1}{2} \right]^T$$

$$\hat{y} = X\Theta = \left[\frac{1}{4}, \frac{3}{4}, \frac{3}{4}, \frac{5}{4} \right]$$



AND Gate

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1



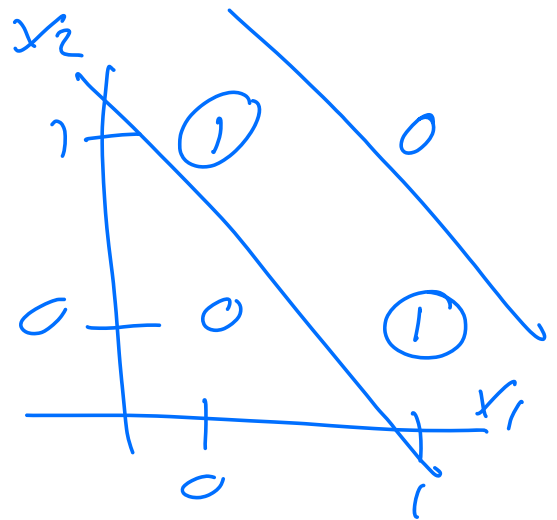
$$G = \begin{bmatrix} -1/4 & 1/2 & 1/2 \end{bmatrix}^T$$

\downarrow w_1 w_2

$$\tilde{y} = XG = \begin{bmatrix} -1/4 & 1/4 & 1/4 & 3/4 \end{bmatrix}^T$$

XOR

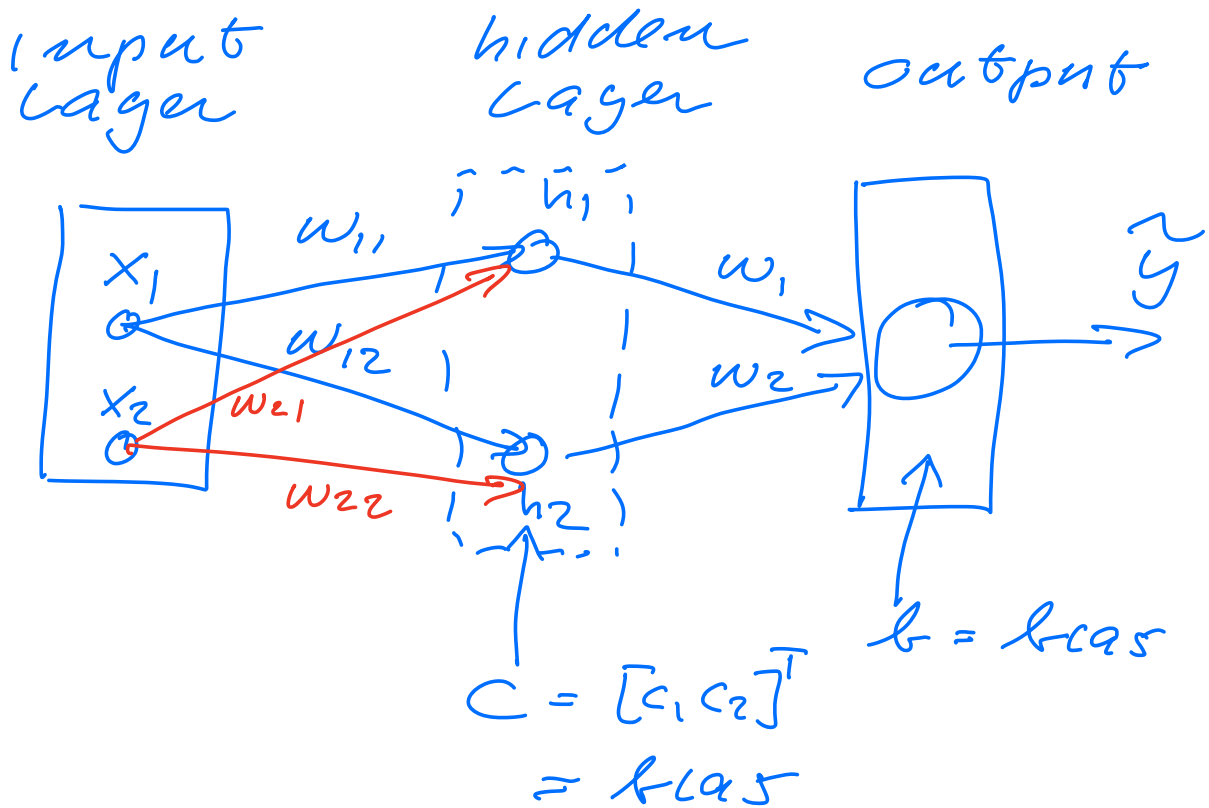
x_1	x_2	t
0	0	0
0	1	1
1	0	1
1	1	0



$$G = \begin{bmatrix} \frac{1}{2} & 0 & 0 \end{bmatrix}^T$$

$$\tilde{y} = XG = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \end{bmatrix}^T$$

First NN encounter,



- Model : number of hidden layers
of nodes
- activation functions; outputs from h_i and outputs from output layer.

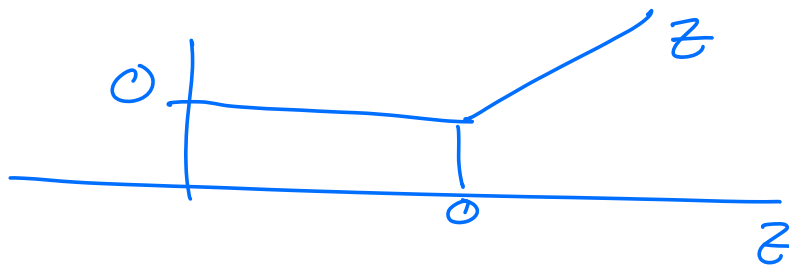
- Examples are

- step function t
- sigmoid $\sim \frac{e}{1+e^t}$
- tanh

ReLU

- ELU ...

$\rightarrow \sigma(z) = \max(0, z)$



- From input to hidden

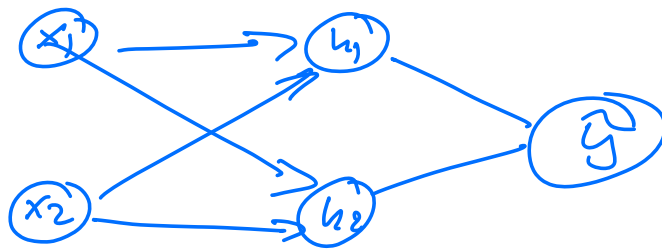
$$h = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + c$$

$$\sigma^{(1)}(h) = \sigma^{(1)}(x; W, c)$$

$$\tilde{y} = \sigma^{(2)}(h; w, b)$$

$$= \sigma^{(2)}(\sigma^{(1)}(x; W_c); w_k)$$

Simple model



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

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

$$w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

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

XOR - gate

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

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

Add bias vector

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

ReLU function application

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

multiply with w

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad ?$$

xor gate