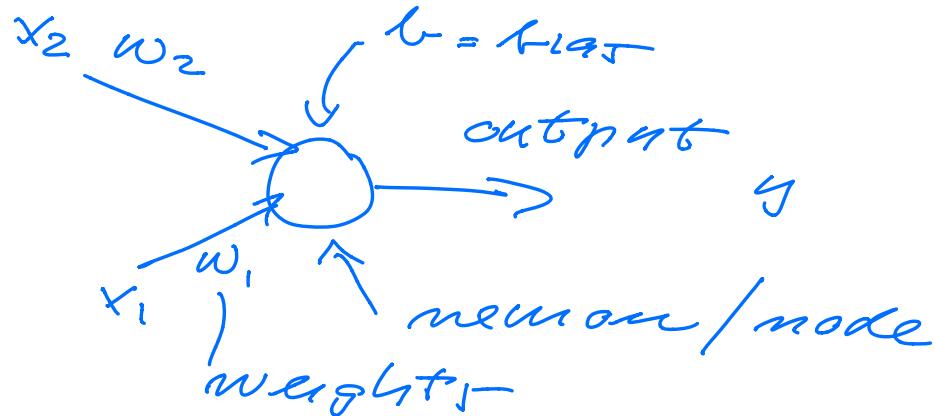


# Lecture January 25

## neural networks

- single perceptron model



$$y = f(u) \quad u = \sum_{i=1}^2 w_i x_i$$

activation function : typical example is the sigmoid, tanh

Linear regression :

$$f(x; \beta) = x \cdot \beta$$



$w$  = our parameter

$b$  = bias

or ... .. . . . . . . . . . .

$$g(x_1 w_1 + x_2 w_2 + b) = y$$

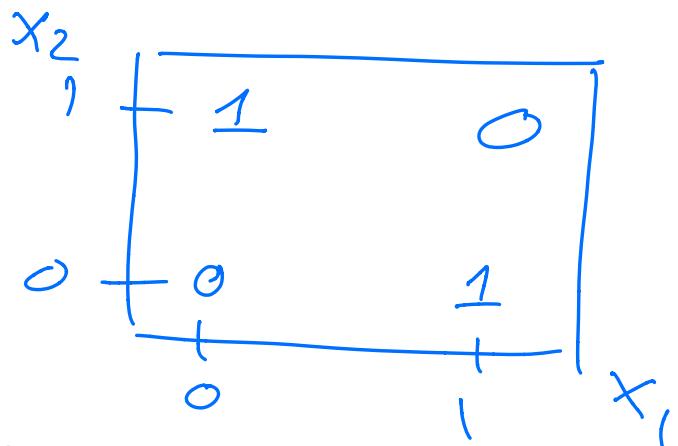
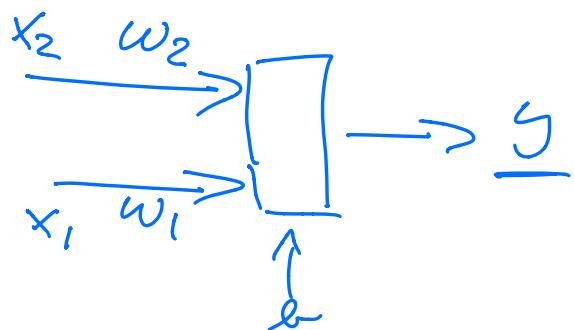
$$y \in [0, 1]$$

XOR-model (Exclusive OR)

Linear fit

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

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



Define a cost function

$$C(w, b) = \frac{1}{4} \sum_{i=1}^4 (y_i - \hat{y}_i)^2$$

$$= \frac{1}{4} \sum_{i=1}^4 (t_i - \hat{y}_i)^2$$

↑ target      ↓ output from model

$$\beta = [b, w_1, w_2] \quad (\oplus)$$

Design Matrix

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

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

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

$$b = 1/2 \quad w_1 = w_2 = 0$$

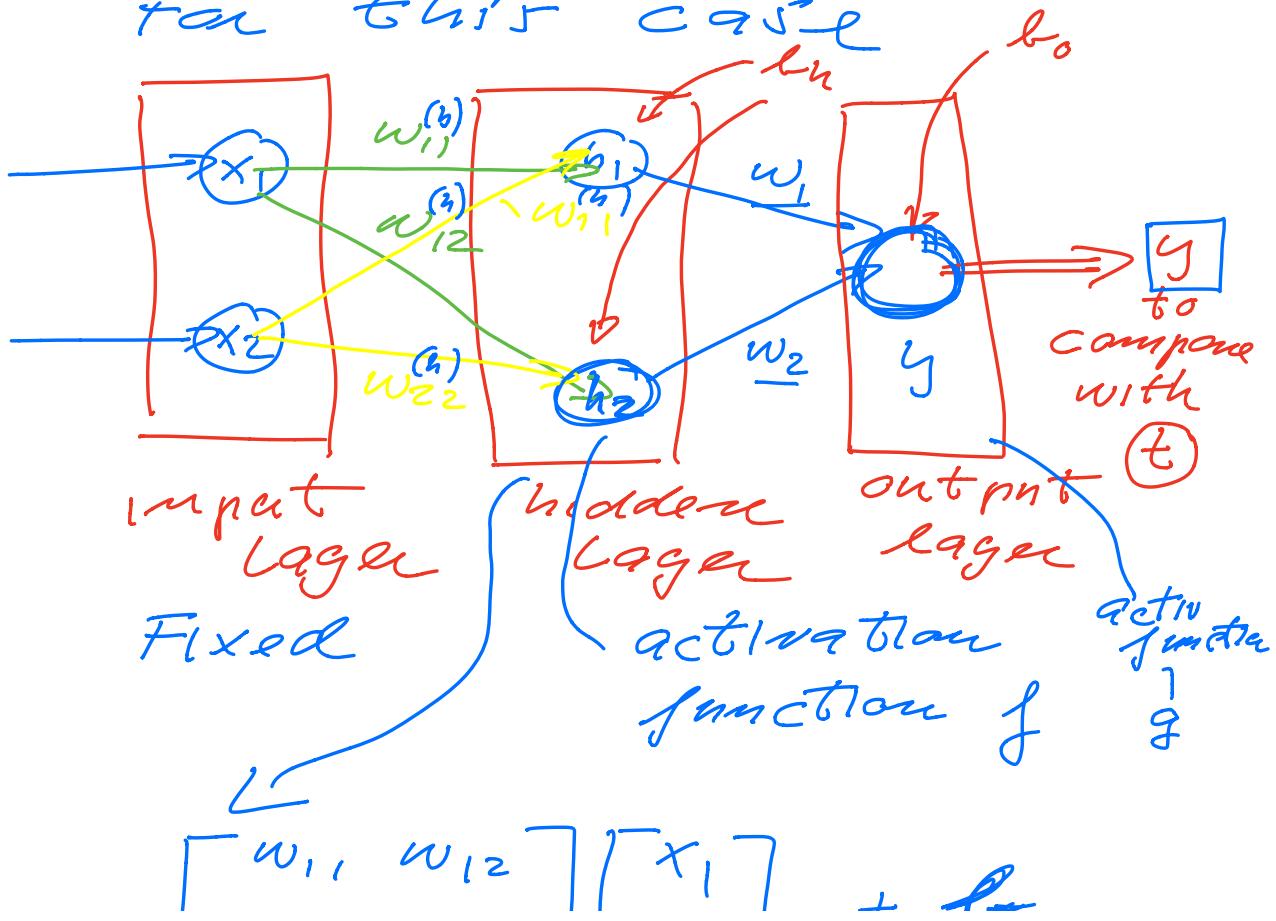
$$y = X\beta = 1/2$$

all predictions give 1/2

We want  $\underline{[0 \ 1 \ 1]}$

Linear regression fails.  
We can make a network  
with a hidden layer  
and two neurons/nodes.  
 $\Rightarrow$  This gives correct  
results

For this case



$$[w_{z1} \ w_{z2}] [x_2] \rightarrow$$

$$W \quad X$$

(i) Feed Forward stage  
with randomly  
defined values for  
 $W$  and  $b$ .

$$h_1 = w_{11}x_1 + w_{12}x_2$$

$$h_2 = w_{21}x_1 + w_{22}x_2$$

$$Wx + b = h \rightarrow$$

$$f(h) = f(Wx + b)$$

$\rightarrow$  output layer

$$g(f(Wx + b))$$

$$= g(w_1 h_1 + w_2 h_2)$$

$$= y$$

to be compared with  
target  $t$

How do we fit the targets  $t$  with our model outputs  $y$ ?

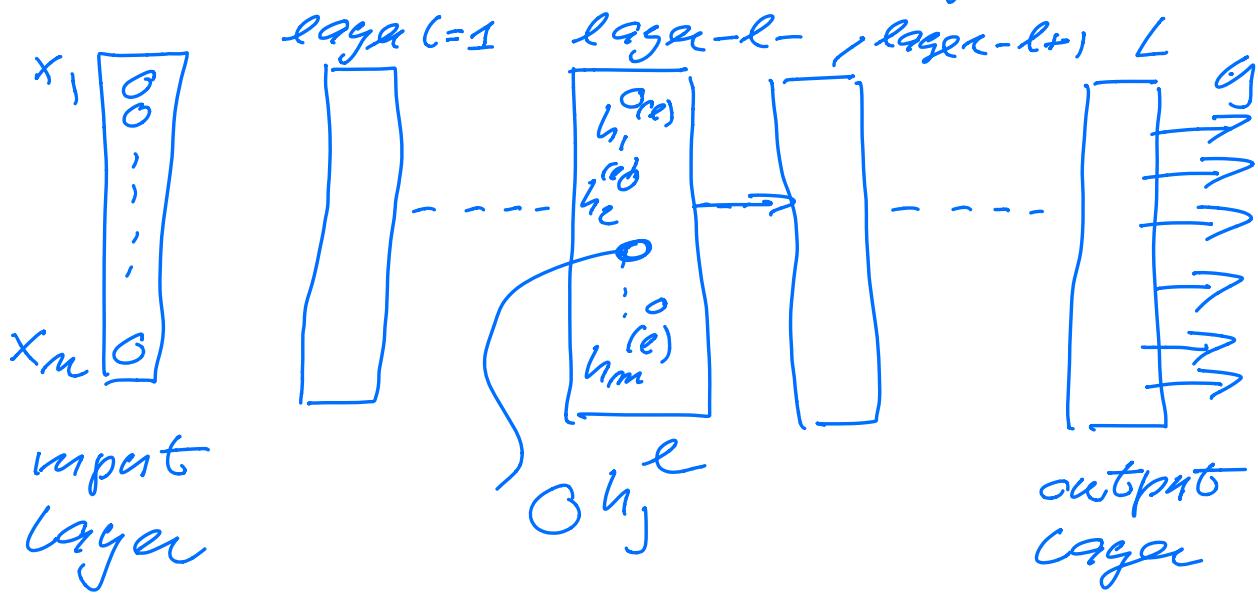
We need an algorithm to find weights and biases.

- (ii) Back Propagation algo,  
application of chain  
rule, New values of  
weights and biases,
- (iii) New FF pass, New  
Back prop.
- (iv) Repeat (i)-(iii) till  
achieved accuracy.

In Essence :

- Back propagation
- in each layer we  
need gradient descent

MLP / FFNN : Basic def's:



Regression :

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

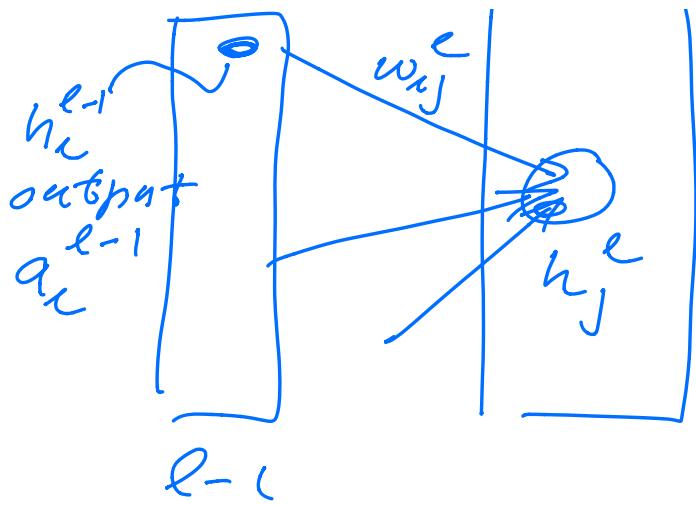
weights in a layer  $w_{ij}^l$

inputs from layer  $l-1$  to node  $m$  in  $l$  :  $a_i^{l-1}$

bias in layer  $l$  for mode  $j$  :  $b_j^l$

input to mode  $j$  in layer

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



- new parameters :
  - # hidden layers
  - # nodes in hidden layers
  - # different activation functions
  - ( normal approach is to have same activation for hidden nodes and different activation function for output layer )

$\rightarrow$  output from node  $h_j^l$

$$a_j^l = f(z_j^l) \stackrel{\text{sigmoid}}{=}$$

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

Derivatives & Back prop.

$$\frac{\partial C(w, b)}{\partial w} \rightarrow \frac{\partial C}{\partial b}$$

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

$$\frac{\partial z_j^e}{\partial a_i^{l-1}} = w_{ji}^e$$

$$\begin{aligned} \frac{\partial a_j^e}{\partial z_j^e} &= a_j^e(1-a_j^e) \\ &= f(z_j^e)(1-f(z_j^e)) \end{aligned}$$

$l = L = \text{output layer}$ .

$$C(w^L, b^L) = \frac{1}{n} \sum (y_i - t_i)^2$$

$$= \frac{1}{2} \sum_i (\hat{q}_i^L - t_i)^2$$

$$\frac{\partial C}{\partial w_{jk}^L} = (\hat{q}_j^L - t_j) \frac{\partial q_j^L}{\partial w_{jk}^L}$$

$$\frac{\partial q_j^L}{\partial w_{jk}^L} = \frac{\partial q_j^L}{\partial z_j^L} \frac{\partial z_j^L}{\partial w_{jk}^L} =$$

$$q_j^L (1 - q_j^L) q_k^{L-1}$$

$$\frac{\partial C(w, b)}{\partial w_{jk}^L} = (\hat{q}_j^L - t_j) \times q_j^L (1 - q_j^L) q_k^{L-1}$$

$$\begin{aligned} \delta_j^L &= q_j^L (1 - q_j^L) (\hat{q}_j^L - t_j) \\ &= f'(z_j^L) \frac{\partial C}{\partial q_j^L} \end{aligned}$$

$\overbrace{\dots}^L - L - L - 1$

$$\frac{\partial \tilde{y}}{\partial w_{jk}^l} = \delta_j \, a_k$$

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

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

1

## Parameters

$$C(w, b) = \frac{1}{2} \sum (y_i - t_i)^2$$

- $w, b$
- # hidden layer
- # nodes hidden layer
- activation function  
for hidden layer
- Learning rates &

for grad descent method

- activation func for output layer
- type of gradient method
- hyperparameter  
 $\lambda \quad \lambda \|\mathbf{W}\|_2^2$