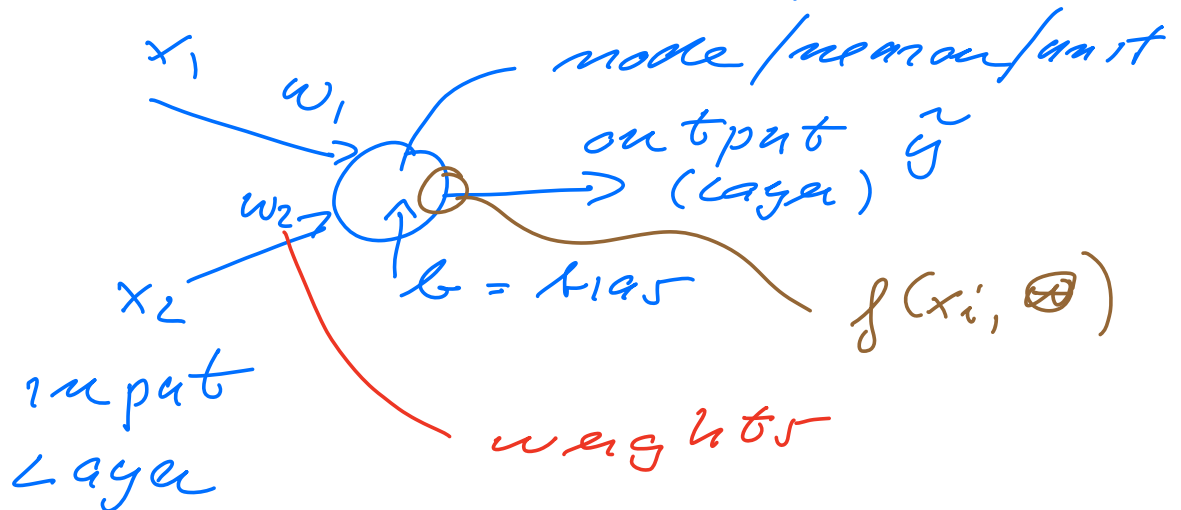# Lecture November 29

Mathematics of feed forward
Neural Networks (FFNN) /
Multi-Layer Perceptron (MLP)

- Simple perceptron
  model.

  - activation function,

$x_1$
$w_1$

node/neuron/unit
output $\tilde{y}$
(layer)

$w_2$

$b = $ bias

$f(x_i, \Theta)$

$x_2$

input
Layer

weights

$$\tilde{y}_i = f(x_i; \Theta) = w^T x_i + b$$
activation

$$\Theta = \{ w, b \} \qquad w^T = [w_1 \; w_2]$$

Targets $\{$ $y_i \quad (= t_i)$
outputs $\{$

$$D = \{ (x_0, y_0), (x_1 \, y_1) \ldots (x_{m-1} \, y_{m-1}) \}$$

Examples: XOR, OR, AND
gates

## OR- gate



| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

$y = 0$ •

$y = 1$ O

linear reg
logistic
reg.

## AND - gate

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

## XOR - Gate

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |



$$\tilde{y}_i = x_i \cdot w_1 + x_i \cdot w_2 + \underbrace{b}_{\text{intercept}}$$

## Linear regression

$$X = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 3}$$

$$\beta \rightarrow \Theta = \begin{bmatrix} b \\ w_1 \\ w_2 \end{bmatrix}$$

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

$$\hat{e} = (x^T x)^{-1} x^T y$$

OR-Gate
$$= [1/4, 1/2, 1/2]^T$$

$$\tilde{y} = x\hat{e} = [1/4, 3/4, 3/4, 5/4]$$

$$\tilde{y}^2 \leq 1/2 \quad \text{then} \quad \tilde{y} = 0$$

$$\tilde{y} > 1/2 \quad \text{then} \quad \tilde{y}^2 = 1$$

Linear/Logistic are ck

## AND gate

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

$$\tilde{y} = [-1/4, 1/4, 1/4, 3/4]^T$$

$$\tilde{y} \leq 0.5 \quad \text{then} \quad \tilde{y}^2 = 0$$

$$\text{else} \quad \tilde{y}^2 = 1$$

## XOR - gate

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

$$\tilde{y} = [\tfrac{1}{2}, \tfrac{1}{2}, \tfrac{1}{2}, \tfrac{1}{2}]^T$$

With one hidden layer and
two hidden neurons/nodes
we can reproduce the xor
output



Hidden Layer
— # hidden layers
    — # neurons in a layer

—    W :   $w_{12}$, $w_{11}$, $w_{21}$ $w_{22}$
     b :   $b_1$, $b_2$
—    activation functions
—    output layer
     — # output nodes
     — activation functions
Universal approximation

theorem:

Given a function $F \in \mathbb{C}[0,1]^d$
and a parameter $\varepsilon > 0$,
then there is a one layer
( = 1 hidden layer) neural
network $f(x; W, b)$

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

with $W \in \mathbb{R}^{n \times m}$ and
$b \in \mathbb{R}^m$ for which

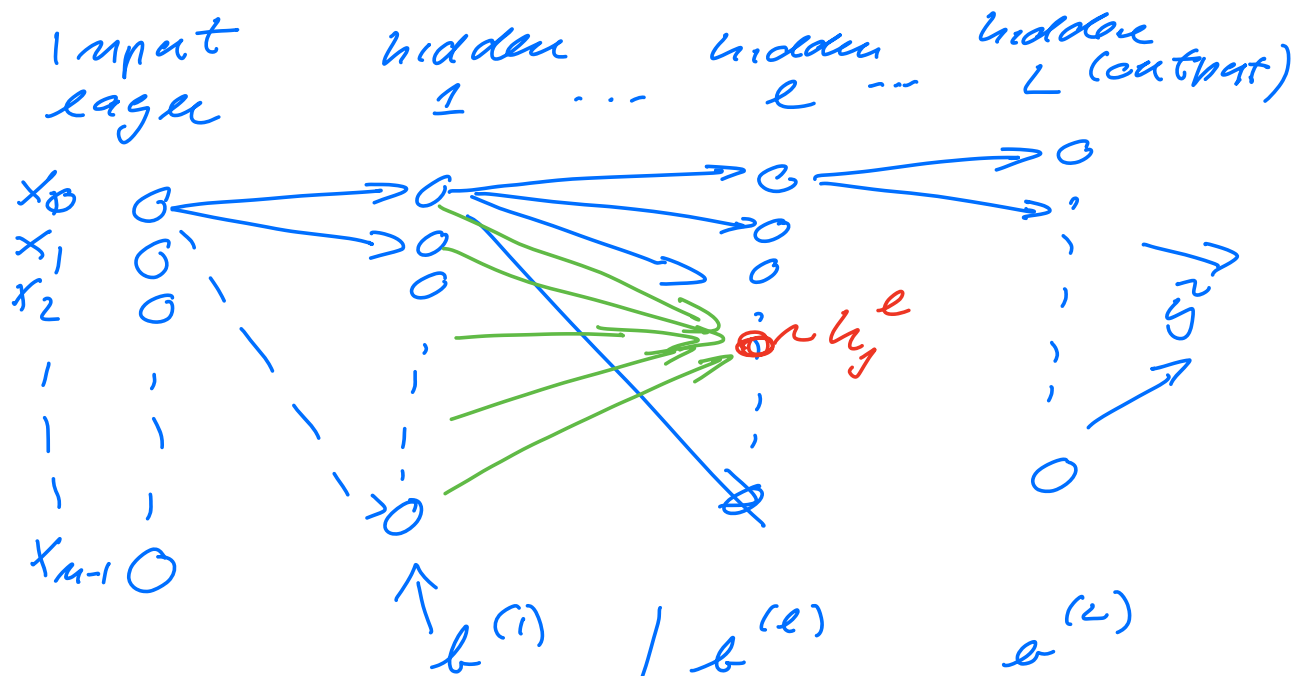$$|f(x; \Theta) - F(x)| < \varepsilon$$

for all $x \in [0,1]^d$ (unit cube)

practicalities: How do we
find $\Theta$ ?

 — Back propagation algo-
  rithm (chain rule)
   $\Rightarrow$ Gradients to be used
    in optimization of
    $W$ and $b$

# Back propagation algorithm

- Regression

$$C(\Theta) = \frac{1}{2} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2$$

Input layer | hidden 1 ... | hidden $\ell$ ... | hidden $L$ (output)



$x_0$
$x_1$
$x_2$

$x_{M-1}$

$\sim h_j^\ell$

$b^{(1)}$      $b^{(\ell)}$      $b^{(L)}$

hidden layer $-\ell-$

input to node $j$:

$z_j^\ell$

$a_{i-1}^{\ell-1} w_{i-1j}$

$h_c^{\ell-1}$   $a_i^{\ell-1} w_{ij}$   $h_j^\ell$

$x_0 w_0$
$x_1 w_1$

$x_i w_i$

$$z_j^{\ell} = \sum_{i=0}^{M_{\ell-1}-1} w_{ij}^{\ell} \, a_i^{\ell-1} + b_j^{\ell} \qquad z = \sum_i w_i' x_i' + \ell$$

outputs from layer $\ell-1$ and node $i$

inputs to layer $-\ell-$

$$z^{\ell} = (W^{\ell})^T a^{\ell-1} + b$$

This is modulated by an activation $f(z_j^{\ell})$ and produces an output

$$a_j^{\ell} = f(z_j^{\ell})$$

$$a^{\ell} = f(z^{\ell})$$

with two hidden layers

First hidden layer

$$z_j^1 = \sum_i w_{ij}^1 x_i + t_j^1$$

output $\qquad a_j^1 = f(z_j^1)$ $\overset{\text{input}}{\overset{\text{layer}}{}}$

$$a^1 = f(z^1)$$

2nd hidden = output layer
(last = output)

$$z_j^2 = \sum_i w_{ij}^2 a_i^1 + b_j^2$$

output $= \tilde{y} = a^2 = f(z^2)$

$$= f(f(z'))$$

3-layer : $f(f(f(z')))$

we have defined

$$a_j^\ell = f(z_j^\ell)$$

typical function

$$f(z_j^\ell) = \frac{1}{1+\exp(-z_j^\ell)}$$

we need to define :

$$\partial z_1^\ell \qquad\qquad \ell - 1$$

$$\frac{\partial v}{\partial w_{ij}^{\ell}} = a_i$$

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

$$\frac{\partial a_j^{\ell}}{\partial z_j^{\ell}} = \frac{\partial f(z_j^{\ell})}{\partial z_j^{\ell}}$$

$$= f(z_j^{\ell})(1 - f(z_j^{\ell}))$$

Specialize to output
layer $\ell = L$

my output i's $\tilde{y}_i = a_i^L$

$$\frac{\partial C(\theta)}{\partial w_{jk}^L} = \frac{\partial}{\partial w_{jk}^L}\left[\frac{1}{2}\sum_i (\tilde{y}_i - y_i)^2\right]$$

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

$\smile w_{jk}$

chain - rule !

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

$$= a_j^L (1 - a_j^L) a_k^{L-1}$$

$$\frac{\partial C(e)}{\partial w_{jk}^L} = (a_j^L - y_j)(1 - a_j^L)$$
$$\times a_k^{L-1} a_j^L$$

Define :

$$\delta_j^L = \underbrace{a_j^L (1 - a_j^L)}_{f'(z_j^L)} \underbrace{(a_j^L - y_j)}_{\frac{\partial C}{\partial a_j^L}}$$

$$\boxed{\frac{\partial C}{\partial \ _j^L} = \delta_j^L a_k^{L-1}}$$

$$\frac{\partial}{\partial w_{jk}} \quad \bigg]$$

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

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

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

$$\delta_j^L = f'(z_j^L) \frac{\partial C}{\partial a_j^L}$$

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

all for final layer.

Need to Back propagate

$$L \to \ell$$

$$\delta_j^\ell = \frac{\partial C}{\partial z_j^\ell} = $$

$$\ell + 1$$

$$\sum_k \frac{\partial c}{\partial z_k^{\ell+1}} \frac{\partial z_k^{\ell+1}}{\partial z_j^{\ell}}$$

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

$$z_j^{\ell+1} = \sum_i w_{ij}^{\ell+1} \underline{a_i^{\ell}} + b_j^{\ell+1}$$

$$\underline{f(z_i^{\ell})}$$

$$\frac{\partial z_k^{\ell+1}}{\partial z_j^{\ell}} = w_{kj}^{\ell+1} f'(z_j^{\ell})$$

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

Now can set up gradients
for each layer:

$$\ell = L-1, L-2, \ldots 1$$

$$w_{jk}^{\ell} \leftarrow w_{jk}^{\ell} - \eta \, \delta_j^{\ell} a_k^{\ell+1}$$

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

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