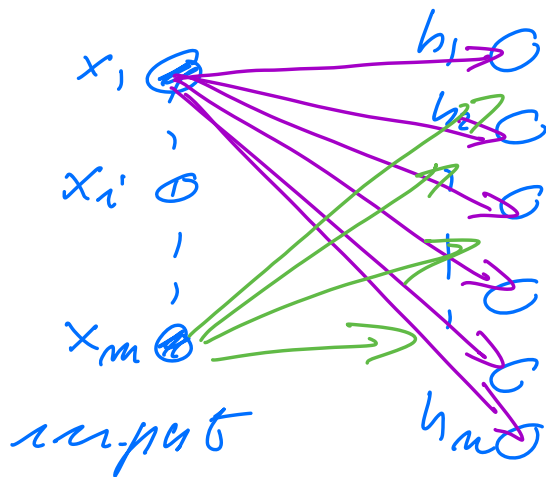


# Lecture October 15

## Typical NN

- $m$  input values
- First hidden layer  
 $n$  - nodes/neurons



weight  
matrix

$$W \in \mathbb{R}^{m \times n}$$

$$b \in \mathbb{R}^n$$

$$\uparrow b$$

$$z(x) = W^T x + b$$

$$x \mapsto \sigma(z(x))$$

activation function

$$= [a_1(x), a_2(x), \dots, a_n(x)]$$

With  $L$ -layers

$$z_1 : \mathbb{R}^{m_1} \rightarrow \mathbb{R}^{m_2}$$

$$z_l : \mathbb{R}^{m_{l-1}} \rightarrow \mathbb{R}^{m_l} \text{ for } 2 \leq l \leq L$$

↑  
output  
layer

Our model

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

$$f(x; \Theta) = \sigma_L(z_L(\dots \sigma_1(z_1(x))))$$

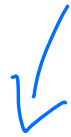
Math of activation  
functions

consider a simple NN  
where  $x, w, b$  are scalar  
quantities,  $L = 2$

$$f(x; \Theta) = \sigma_2(w_2 \sigma_1(w_1 x + b_1) + b_2)$$

in Backpropagation +  
gradient optimization

$$\partial_{w_1} f(x; \theta) \quad \text{and} \quad \partial_{w_2} f(x; \theta)$$



$$\frac{\partial f}{\partial w_1} = \nabla_2' (w_2 \nabla_1' (\overbrace{w_1 x + b_1}^{z_1} + b_2)) \\ \otimes \underline{w_2} \nabla_1' (w_1 x + b_1) x$$

$L$ -layers

$$\partial_{w_1} f(x; \theta) = \left[ \prod_{l=2}^L w_l \right] \\ \otimes \left[ \prod_{l=1}^L \nabla_l' (z_l) \right] x$$

standard activation  
function  $\sigma$

$$\nabla(z) = \frac{1}{1 + e^{-z}} \quad \text{sigmoid}$$

tanh

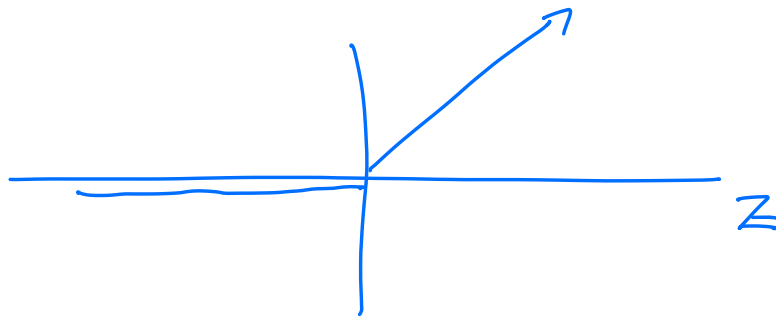
when  $|z| \gg 0$ , then

$\nabla'(z)$  can become small,  
 $\Rightarrow$  vanishing gradients

$$\text{ReLU} = \nabla(z) = \max(0, z)$$

$$\nabla'(z) = 1 \quad \text{for } z > 0$$

$$\text{Leaky ReLU} = \nabla(z) = \begin{cases} \alpha \cdot z & z < 0 \\ z & z \geq 0 \end{cases}$$
$$\alpha \sim 10^{-2} - 10^{-3}$$



$$\nabla'(z) = \begin{cases} \alpha & \text{for } z < 0 \\ 1 & \text{for } z \geq 0 \end{cases}$$

$$\nabla(z) = \text{ELU} = \int \alpha(e^z - 1) \quad z < 0$$

$$|z| \quad z \geq 0$$

$$\sigma'(z) = \begin{cases} \alpha e^z & z < 0 \\ 1 & \text{for } z \geq 0 \end{cases}$$

$$\sigma(z) = \tanh(z)$$

$$\sigma'(z) = 1 - (\tanh(z))^2$$

$$\sigma(z) = \text{sigmoid} = \frac{1}{1 + e^{-z}}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

Classical approximations

$$y = F(x) \quad (\text{continuous})$$

$$x \in [0, 1]^d$$

Given  $F(x)$  and  $\varepsilon > 0$   
there exists a function

$$f(x; \epsilon)$$

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

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

Examples

$$f(x; \epsilon) = \sum_j \epsilon_j x^j$$

Theorem: Stone-Weierstrass

$$F \in C([0, 1]^d), \text{ for each}$$

$\epsilon$  there exists a polynomial

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

$$\forall x \in [0, 1]^d$$

Example:

$$f(x; \epsilon) = \sum_j \epsilon(\beta_j) e^{i 2\pi x \beta_j}$$

$$\left[ \int_0^1 (F(x) - f(x; \epsilon))^2 dx \right]^{1/2}$$

$$\left[ \frac{1}{[c,1]^d} \right]$$

$$< \varepsilon$$

For NNs [Cybenko, 1989]

- Let  $\sigma$  be any continuous sigmoidal function

$$\sigma(z) \rightarrow \begin{cases} 1 & \text{as } z \rightarrow \infty \\ 0 & \text{as } z \rightarrow -\infty \end{cases}$$

Given an  $F \in C([c,1]^d)$  and  $\varepsilon > 0$ , there is a one layer neural network

$$f(x; W, b) = \sigma(x; \theta)$$

with  $W \in \mathbb{R}^{m \times n}$

and  $b \in \mathbb{R}^m$  for

which

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

for all  $x \in [a, b]^d$

or:

Any continuous  $F(x)$

for  $[a, b]^d$  can be

approximated by a one  
layer sigmoidal network  
to arbitrary accuracy.

Called Universal  
approximation theorem

Hornik (1991) refined  
the theorem by letting  
any non-constant  
bounded function  
to be included,

- The theorems do not



say anything about  
the number of nodes  
or the values of the  
weights and biases

- does not mean that  
any NN can be used  
to compute exactly  
any function,

We get an approximation  
that is as good as we  
want.

- The  $F(x)$  functions  
are continuous  
functions.