

# Deep Learning theory - Lecture 2

[0,1]

for a given  $P, \mathcal{C}, \mathcal{F}_0, \mathcal{A}$  we study

$$\Delta(\mathcal{D}_n) = \mathcal{R}(\mathcal{A}(\mathcal{D}_n)) - \inf_{f: X \rightarrow Y} \mathcal{R}(f)$$

Population Risk decomposition:

$$\underbrace{\mathcal{R}(f_{\hat{s}}) - \inf_{f: X \rightarrow Y} \mathcal{R}(f)}_{\text{population risk}} = \underbrace{\mathcal{R}(f_{\hat{s}}) - \inf_{s \in \Theta} \mathcal{R}(f_s)}_{\text{estimation err.}} + \underbrace{\inf_{s \in \Theta} \mathcal{R}(f_s) - \inf_{f: X \rightarrow Y} \mathcal{R}(f)}_{\text{approximation err.}}$$

ANALYSIS ( $\mathcal{F}$ )

Question: What is an (artificial) neural network?

What is the corresponding  $\mathcal{F}_0$ ? How big is it?

"Definition" of neural network:

Neural networks are a class of functions combining (usually pointwise) nonlinear operations - the neurons - through linear maps/transformations / connections - the network

Different connection structures / nonlinearities give different NN.

# 1) Single (hidden) layer neural networks (shallow)

Def: Let  $X = \mathbb{R}^d$ . A single layer neural network of width  $n \in \mathbb{N}$  is a function of the form:

$$\hat{f}_\Theta(x) = \alpha \cdot \sum_{j=1}^n a_j \sigma(w_j \cdot x + b_j)$$

where

- $\alpha \in \mathbb{R}$  is the scaling of the network

- $a_j \in \mathbb{R}^{d'}$  are the output weights

- $w_j \in \mathbb{R}^d$  are the input weights

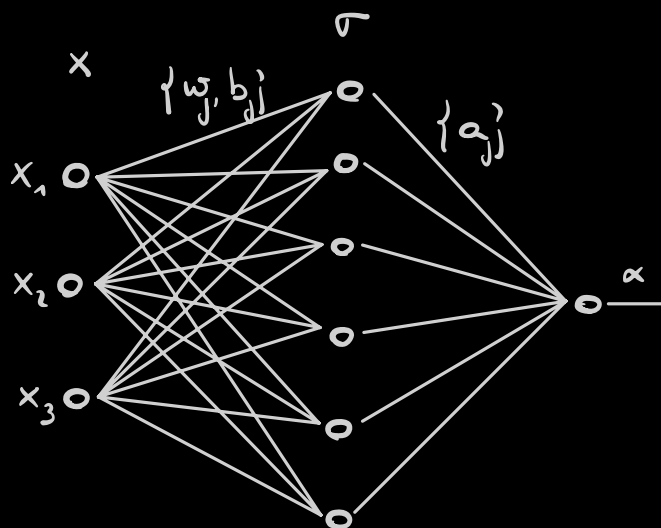
- $b_j \in \mathbb{R}$  are the biases.

- $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  is the activation function.

we will define  $\Theta = (a_j, w_j, b_j)_j \in \mathbb{R}^{(1+d+1)n} =: \Theta$

Note .. Most of the time we will set  $d' = 1$ .

- The above network can be represented as a graph:



here  $d=3, d'=1, n=6$

$$f(x) = \alpha \sum_j a_j \sigma(w_j \cdot x + b_j)$$

- Often (e.g. today)  $\alpha = 1$ , but we will use it when  $n \rightarrow \infty$ .
- The network can also be written using matrix notation, defining the action of  $\sigma$  on  $\mathbb{R}^n$  to be componentwise:

$$\hat{f}_\sigma(x) = \alpha \vec{a} \cdot \sigma(W\vec{x} + \vec{b})$$

where  $\vec{a} = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} \in \mathbb{R}^n$ ,  $\vec{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \in \mathbb{R}^n$ ,  $W = \begin{pmatrix} w_{1,1} & \dots & w_{1,d} \\ \vdots & & \vdots \\ w_{n,1} & \dots & w_{n,d} \end{pmatrix} \in \mathbb{R}^{n \times d}$

- There are many possible choices for activation  $\sigma$ :

a) Rectified Linear Unit (ReLU):

$$\sigma(z) = z_+ := \max(0, z) = z \cdot \mathbb{1}(z \geq 0)$$

b) Sigmoid

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

c) Hyperbolic tangent

$$\sigma(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Less common activations include:

d) step function:  $\sigma(z) = \mathbb{1}(z \geq 0) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$

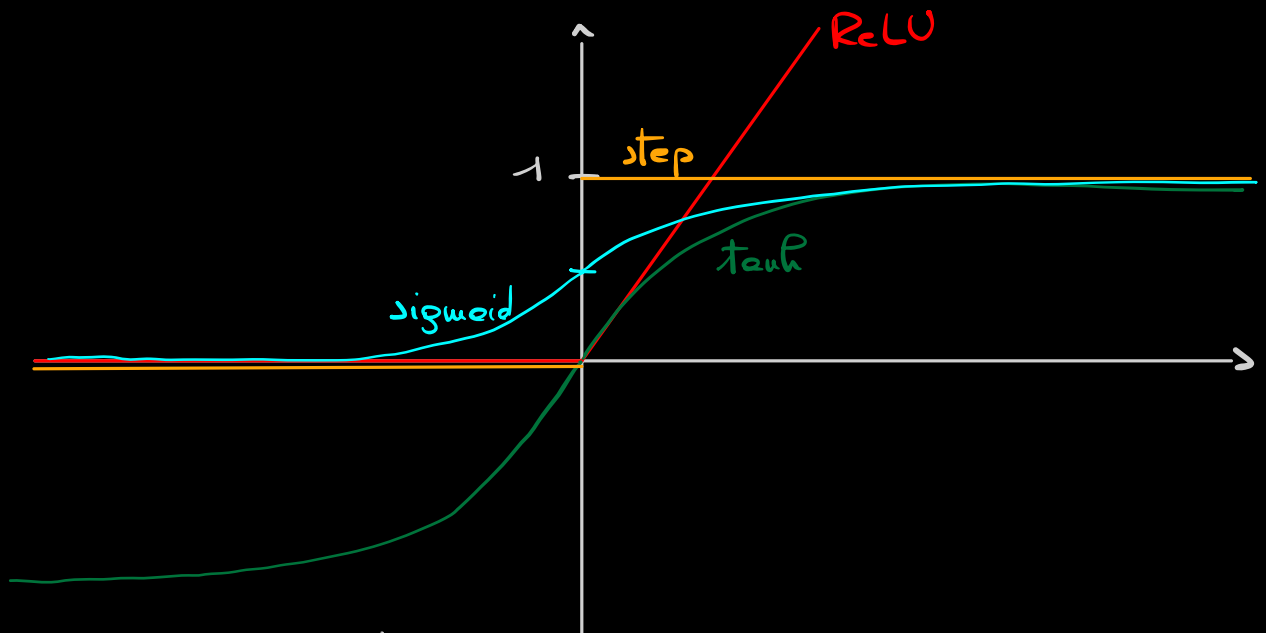
e) identity:  $\sigma(z) = z$

f) smooth ReLUs:  $\sigma(z) = x \cdot \Phi(x)$ ,  $\sigma(z) = x \cdot (1 + e^{-z})^{-1}$

g) gaussian density:  $\sigma(z) = e^{-x^2}$

h) squared:  $\sigma(z) = z^2$

d) (... anything ...)



For each nonlinearity  $\sigma$ , we denote by

$$\overline{F}_{\sigma, n} := \left\{ x \sum_{j=1}^n a_j \cdot \sigma(w_j \cdot x + b_j) : (a_j, w_j, b_j) \in \mathbb{R}^{1+d+1} \forall j \right\}$$

the space of single layer NN of width  $n$  with activation  $\sigma$ ,

$$\overline{F}_{\sigma} = \bigcup_{n=1}^{\infty} \overline{F}_{\sigma, n}$$

Example: let  $X \subseteq \mathbb{R}$

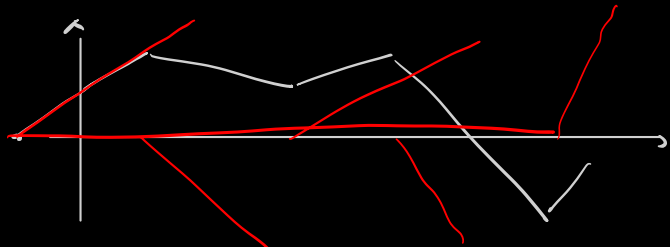
open/closed

$$\text{let } PL_n(X) = \left\{ f \in C(X) : f(x) = (\tilde{a}_n \cdot x + \tilde{c}_n) \mathbb{1}(x \in [\tilde{b}_n, \tilde{b}_{n+1}]), \begin{matrix} \tilde{b}_n \in X \\ \tilde{a}_n, \tilde{c}_n \in \mathbb{R} \\ n \in \{1, \dots, n\} \end{matrix} \right\}$$

Lemma 2.2: For every  $h \in PL_n(X)$  there exists a single layer NN  $f$  of width  $n$  with ReLU nonlinearity such that  $f(x) = h(x)$   
In other words,  $PL_n(X) \subseteq \overline{F}_{\text{ReLU}, n}$

Proof: Set  $w_j = 1 \quad \forall j \in \{1, \dots, n\}$ .

Then we can write



$$f_\theta(x) = \sum_{j=1}^n a_j (x + b_j) \mathbb{1}(x + b_j \geq 0) = \sum_{j=1}^n (a_j x + a_j b_j) \mathbb{1}(x \geq -b_j)$$

So, for  $k=1$ :  $-b_1 = \tilde{b}_1$ ,  $a_1 = \tilde{a}_1$

$$k \rightarrow k+1: \text{ solve } \sum_{j=1}^{k+1} (a_j x + a_j b_j) = \tilde{a}_{k+1} x + \tilde{c}_{k+1} \implies \begin{cases} a_{k+1} = \tilde{a}_{k+1} - \sum_{j=1}^k a_j \\ b_{k+1} = \frac{1}{a_{k+1}} \left( \tilde{c}_{k+1} - \sum_{j=1}^k a_j b_j \right) \end{cases}$$

□

- Sometimes a function is applied to the output of a network:

$$\hat{f}_\theta(x) = \left( 1 + \exp \left( -\alpha \sum_{j=1}^n a_j \sigma(w_j x + b_j) \right) \right)^{-1}$$

# Approximation error for single layer NNs

Which functions can we approximate arbitrarily well?

Def: a class of functions  $\mathcal{F}$  is a universal approximator over a compact set  $X$  if for every  $g \in C(X)$  there exists  $f \in \mathcal{F}$  with  $\|f - g\|_\infty \leq \varepsilon$ .

Consider first the univariate case. Throughout, we set wlog.  $X = [0, 1]$   
We start with intuitive, constructive pf.

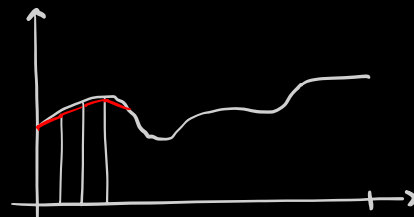
Prop 2.1: Let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be  $p$ -Lipschitz. Then for every  $\varepsilon > 0$  there is a single-layer NN  $f_\varepsilon$  with  $\sigma(z) = z_+$  of width  $n = \lceil p/\varepsilon \rceil$  such that  $\|g - f_\varepsilon\|_\infty \leq \varepsilon$

Proof: By the above lemma we only need to show that  $\exists h \in \mathcal{PL}_n(X)$  such that  $\|g - h\|_\infty \leq \varepsilon$

We choose  $h$  as the linear interpolator of  $\{g(u\varepsilon/\lceil p/\varepsilon \rceil)\}_{u=0}^{\lceil p/\varepsilon \rceil}$

$$\text{Let } w_j = \frac{\varepsilon}{p}, \quad \tilde{b}_j = j \frac{\varepsilon}{p}.$$

$$a_0 = g(0) \quad \tilde{a}_j = \frac{g(b_j) - g(b_{j-1})}{\varepsilon}$$



$$\text{so that } f_\varepsilon(x) = \sum_{j=0}^{n-1} a_j \mathbb{1}(x \geq b_j) = \sum_{j=0}^{n-1} g(b_j) \mathbb{1}(x \in (b_j, b_{j+1}])$$

$$\Rightarrow |g(x) - f_\varepsilon(x)| \leq |g(x) - g(b_j)| + |g(b_j) - f_\varepsilon(b_j)| + |f_\varepsilon(b_j) - f_\varepsilon(x)|$$
$$x \in (b_j, b_{j+1}) \quad \leq p(x - b_j) + 0 + 0 \leq p \frac{\varepsilon}{p} = \varepsilon.$$

Alternatively we could have used the result

Lemma 2.2 Let  $X$  be a compact subset of  $\mathbb{R}$ . Then  $PL$  is dense in  $C(X)$  in the  $\|\cdot\|_\infty$  topology.

Proof: see course of real analysis (e.g. Rudin)

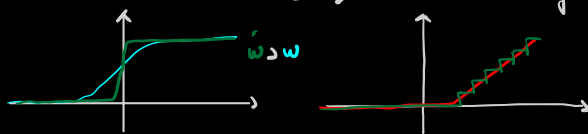
Cor. 2.3: For every compact set  $X \subseteq \mathbb{R}$ , single layer neural networks with ReLU nonlinearity are universal approximators

Remark: While the above result is more general (approximation of  $C(X)$  instead of Lipschitz) Prop 2.1 gives bounds on the size of the network low complexity  $\leftrightarrow$  low width

Exercise: Let  $X$  be compact, then  $\overline{F}_{\text{sigmoid}}$  is univ. approx in  $X$ .

Idea: By Cor 2.3  $\exists h_u \in \overline{F}_{\text{ReLU}, n}$  with  $\|h_u - g\|_\infty \leq \frac{\varepsilon}{2}$

Then, approximate each neuron of  $h_u$  (with error  $\frac{\varepsilon}{2n}$ ) with  $\sum a_j \sigma(w_j x + b_j)$  for  $w_j \gg 1, a_j \ll 1$



Lemma 2.4 Let  $X \subseteq \mathbb{R}$  be compact.  $\overline{F}_{\mathbb{1}_{\geq 0}}$  is a universal approximator over  $X$

PF: This follows from the fact that

$$\mathbb{1}(x \in [b_1, b_2]) = \mathbb{1}(x \geq b_1) \cdot \mathbb{1}(x \leq b_2) = \mathbb{1}(x - b_1 \geq 0) \mathbb{1}(-x + b_2 \geq 0)$$

together with the density of simple functions in  $C(X)$

Remark: Polynomials of bounded degree are closed under linear transf.

$$p \in \mathcal{P}(n) \implies a p(a'x + b') + b \in \mathcal{P}(n)$$

and are not dense in  $C(X)$  for fixed  $n \implies$  no univ. approx

Remark: The compactness of  $X$  is necessary: for any finite relu network  $F$ , if  $X = \mathbb{R}$   $\|\cos(x) - F(x)\|_\infty \geq 1$ .

Remark: Why the  $\|\cdot\|_\infty$  topology? Because we want to bound  $R(\sigma)$  without knowing  $P$ .

Q: How about higher dimensions?  $d \geq 1$

Idea: We would like to approximate with sums.

$$\prod_e \mathbb{1}(x_e \in A_e) \quad A_e = \bigcup_{j=1}^d [b_{e,j}, b'_{e,j})$$

However, this cannot be done (at least directly) with shallow nets.

Consider  $\sigma(z) = \cos(z)$ . then Or can it?

$$\begin{aligned} \sigma(z) \cdot \sigma(y) &= \cos(z) \cos(y) \\ &= \cos(z+y) + \cos(z-y) \\ &= \sigma(z+y) + \sigma(z-y) \end{aligned}$$

So, at least in principle, using the above property of cosines and the Fourier approximation of step functions:

$$\begin{aligned} g(x) &\approx \sum_j \alpha_j \prod_e \mathbb{1}(x_e \in [b_{j,e}, b'_{j,e})) \approx \sum_j \alpha_j \prod_e \sum_k \beta_k \cos(\gamma_k x_e) \\ \text{algebra} &\rightarrow \approx \sum_j \alpha_j \sum_k \beta_k \prod_e \cos(\gamma_k x_e) \approx \sum_j \alpha_j \sum_k \beta \sum_m \delta_m \cos(\sum_e \gamma'_e x_e) \end{aligned}$$



$$\text{expressivity} \approx \sum_{\delta} \alpha_{\delta} \sum_{\alpha} \beta_{\alpha} \sum_{\omega} \delta_{\omega} \sum_{\eta} \eta_{\eta} \mathbb{1}[\eta'_{\omega} \cdot x \geq b_{\omega}]$$

In practice, this is painful, but the heavy lifting was done for us:

Thm 2.5 (Stone-Weierstrass): Let  $\mathcal{F} \subseteq C(X)$  for compact  $X \subseteq \mathbb{R}^d$  satisfy:

a) for every  $x \in X$ , there exists  $f \in \mathcal{F}$  such that  $f(x) \neq 0$

b) for every pair  $x, x' \in X$  with  $x \neq x'$  there exists  $f \in \mathcal{F}$  with

$$f(x) \neq f(x') \quad (\mathcal{F} \text{ separates points})$$

c)  $\mathcal{F}$  is closed under pointwise multiplication ( $\mathcal{F}$  is an algebra)

then  $\mathcal{F}$  is a universal approximator.

Lemma 2.6  $\mathcal{F}_{\cos}$  is universal

Pf: a) each  $f \in \mathcal{F}_{\cos}$  is continuous (finite sum of cont. functions)

b)  $\cos(0 \cdot x) = 1 \quad \forall x \in X$

c)  $x \neq x' \Rightarrow f(z) = \cos\left(\frac{(z-x') \cdot (x-x')}{\|x-x'\|^2}\right)$  satisfies  $\begin{cases} f(x') = 1 \\ f(x) = 0 \end{cases}$

d) already checked.  $\square$

Thm 2.7 Suppose  $\sigma \in C(\mathbb{R})$  is sigmoidal:  $\begin{cases} \lim_{z \rightarrow -\infty} \sigma(z) = 0 \\ \lim_{z \rightarrow \infty} \sigma(z) = 1 \end{cases}$   
then  $\mathcal{F}_{\sigma}$  is universal.

Also,  $\mathcal{F}_{\text{ReLU}}$  is universal

Pf (sketch): By Lemma 2.7 we have there exists  $n \in \mathbb{N}$ ,

$$h_n(x) = \sum_{j=1}^n \tilde{a}_j \cos(\tilde{\omega}_j \cdot x + \tilde{b}_j) \in \overline{\mathcal{F}}_{\cos}$$

$$\text{with } \|h - g\|_{\infty} \leq \frac{\varepsilon}{2}$$

Then, since  $h_{n,j}(x) = \tilde{a}_j \cos(\tilde{\omega}_j \cdot x + \tilde{b}_j) \in C(X)$ , by exercise we have  $\exists f_{n,j} \in \mathcal{F}_{\text{sigmoid}} : \|f_{n,j} - h_{n,j}\| \leq \frac{\varepsilon}{2n}$

$$\implies \text{for } f(x) = \sum_j f_{n,j}(x) \in \mathcal{F}_{\text{sigmoid}}$$

$$\|f_n - g\|_{\infty} \leq \|f_n - h_n\| + \|h_n - g\| \leq \sum_{j=1}^n \frac{\varepsilon}{2n} + \frac{\varepsilon}{2} \leq \varepsilon. \quad \square$$

Note: the algebra condition does not hold for polynomials of bounded degree. In fact.

Thm (Leshno, 1993):  $\overline{\mathcal{F}}_{\sigma}$  is universal iff  $\sigma \in C(\mathbb{R})$  is not a polynomial