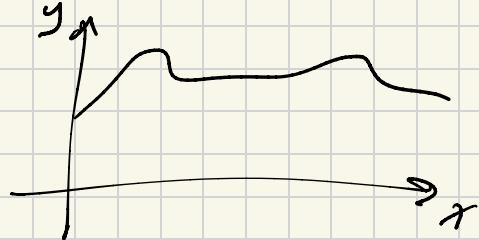
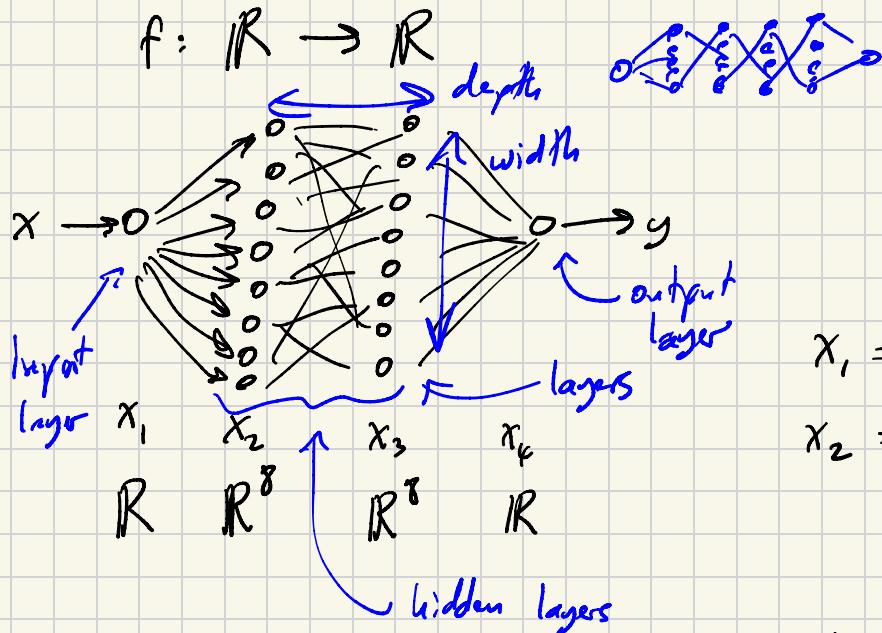


# Function approximation with neural networks



$$x_1 = x$$

$$x_2 = \sigma(w_1 x_1 + b_1)$$

$\mathbb{R} \rightarrow \mathbb{R}^8$

$$x_i \in \mathbb{R}$$

$$w_1 \quad 8 \times 1$$

$$b_1 \quad 8$$

$$w_2 \quad 8 \times 8$$

$$b_2 \quad 8$$

$$w_3 \quad 1 \times 8$$

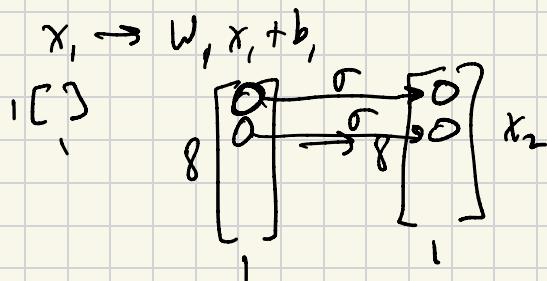
$$b_3 \quad 1$$

parameters

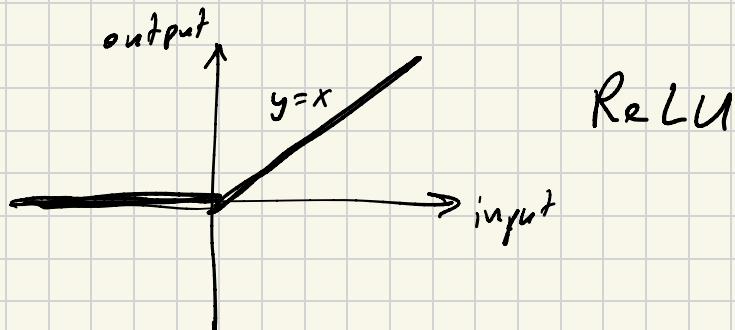
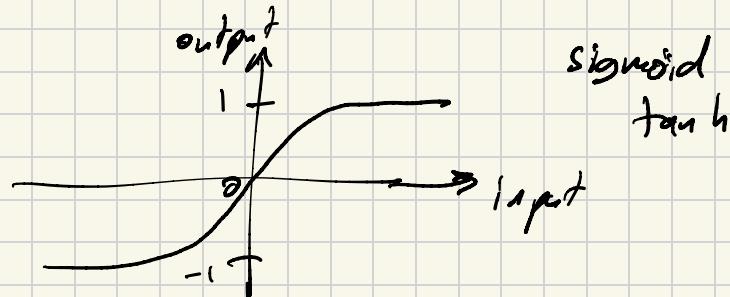
$$x_3 = \sigma(w_2 x_2 + b_2)$$

$$x_4 = w_3 x_3 + b_3$$

$$y = x_4$$



$\sigma$  = activation function

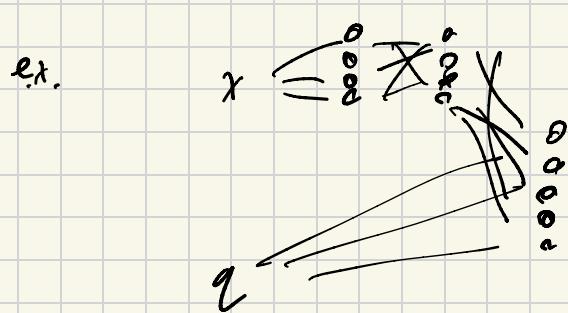


$$x \in \mathbb{R} \rightarrow \text{NN} \rightarrow y = f(x)$$

$$\theta \in \mathbb{R}^{q_f}$$

$$\text{parameters} = (\omega_1, b_1, \omega_2, b_2, \omega_3, b_3)$$

$$v_\theta^f$$



# Universal approximation theorems

[https://en.wikipedia.org/wiki/Universal\\_approximation\\_theorem](https://en.wikipedia.org/wiki/Universal_approximation_theorem)

Single hidden layer, arbitrary width (1990s)

$$\mathbb{R}^n \rightarrow \mathbb{R}^m$$

**Universal approximation theorem** — Let  $C(X, \mathbb{R}^m)$  denote the set of **continuous functions** from a subset  $X$  of a Euclidean  $\mathbb{R}^n$  space to a Euclidean space  $\mathbb{R}^m$ . Let  $\sigma \in C(\mathbb{R}, \mathbb{R})$ . Note that  $(\sigma \circ x)_i = \sigma(x_i)$ , so  $\sigma \circ x$  denotes  $\sigma$  applied to each component of  $x$ .

Then  $\sigma$  is not **polynomial** if and only if for every  $n \in \mathbb{N}$ ,  $m \in \mathbb{N}$ , **compact**  $K \subseteq \mathbb{R}^n$ ,  $f \in C(K, \mathbb{R}^m)$ ,  $\varepsilon > 0$  there exist  $k \in \mathbb{N}$ ,  $A \in \mathbb{R}^{k \times n}$ ,  $b \in \mathbb{R}^k$ ,  $C \in \mathbb{R}^{m \times k}$  such that

$$\sup_{x \in K} \|f(x) - g(x)\| < \varepsilon$$

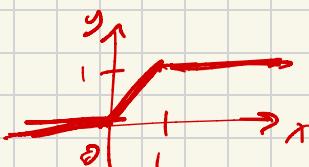
where  $g(x) = C \cdot (\sigma \circ (A \cdot x + b))$   $\leftarrow$  NN



Cybenko, G. (1989). "Approximation by superpositions of a sigmoidal function". Mathematics of Control, Signals, and Systems. 2 (4): 303–314. doi:[10.1007/BF02551274](https://doi.org/10.1007/BF02551274)

Hornik, Kurt (1991). "Approximation capabilities of multilayer feedforward networks". Neural Networks. 4 (2): 251–257. doi:[10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T)

\* Pinkus, Allan (1999). "Approximation theory of the MLP model in neural networks". Acta Numerica. 8: 143–195. doi:[10.1017/S0962492900002919](https://doi.org/10.1017/S0962492900002919)



*→ deep learning revolution*

## Fixed width, arbitrary depth (2010s)

**Universal approximation theorem** (*L<sub>1</sub> distance, ReLU activation, arbitrary depth, minimal width*). For any Bochner–Lebesgue p-integrable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and any  $\epsilon > 0$ , there exists a **fully-connected ReLU** network  $F$  of width exactly  $d_m = \max\{n + 1, m\}$ , satisfying

$$\int_{\mathbb{R}^n} \|f(x) - F(x)\|^p dx < \epsilon.$$

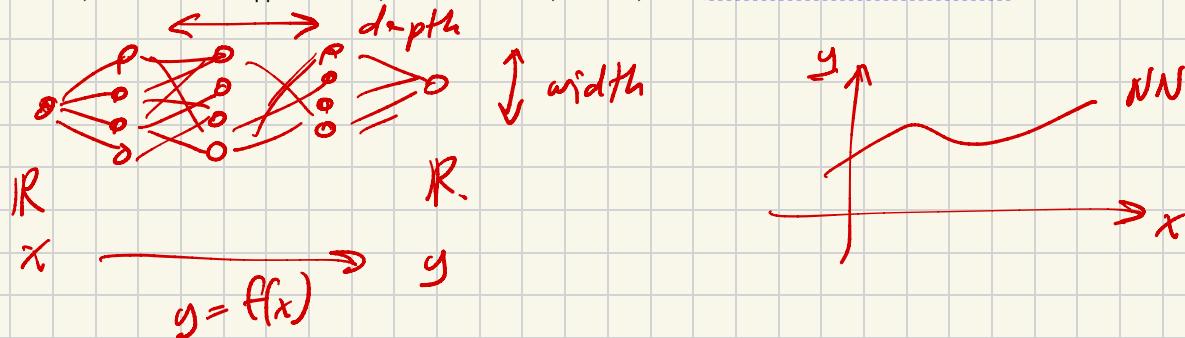
*NN*

Moreover, there exists a function  $f \in L^p(\mathbb{R}^n, \mathbb{R}^m)$  and some  $\epsilon > 0$ , for which there is no **fully-connected ReLU** network of width less than  $d_m = \max\{n + 1, m\}$  satisfying the above approximation bound.

Lu, Zhou; Pu, Hongming; Wang, Feicheng; Hu, Zhiqiang; Wang, Liwei. "The Expressive Power of Neural Networks: A View from the Width". NeurIPS 2017. [arXiv:1709.02540](https://arxiv.org/abs/1709.02540)

Park, Yun, Lee, Shin, Sejun, Chulhee, Jaeho, Jinwoo. "Minimum Width for Universal Approximation". ICLR 2021. [arXiv:2006.08859](https://arxiv.org/abs/2006.08859)

R. DeVore, B. Hanin, and G. Petrova, "Neural network approximation", Acta Numerica 30, 327-444, 2021. [doi:10.1017/S0962492921000052](https://doi.org/10.1017/S0962492921000052)



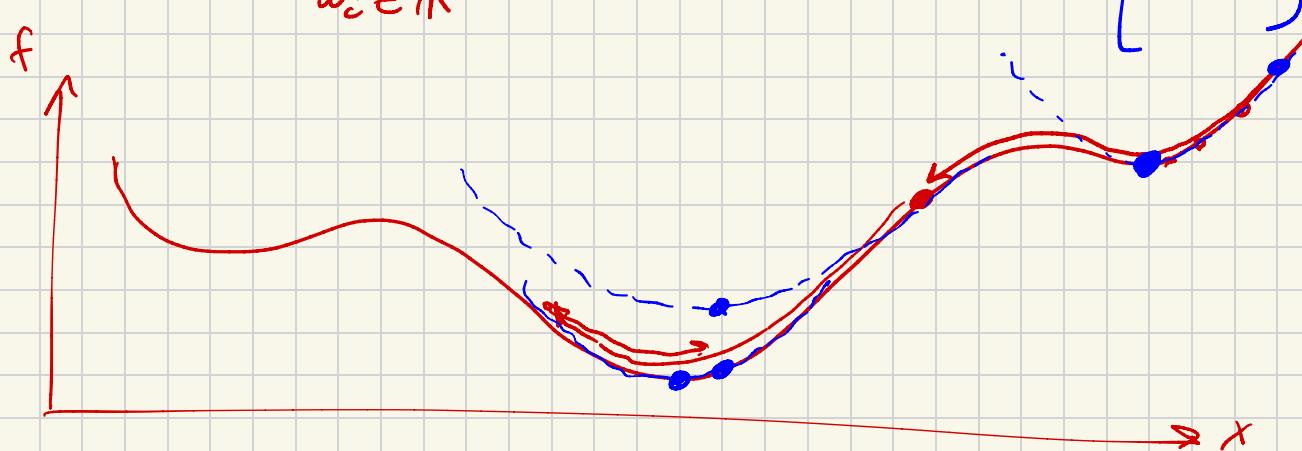
$$\frac{\partial f}{\partial x} \rightarrow x.\text{grad}$$

$$\begin{matrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{matrix} \quad \left. \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} \right\} \rightarrow$$

$$\nabla_x f \quad f \in \mathbb{R}$$

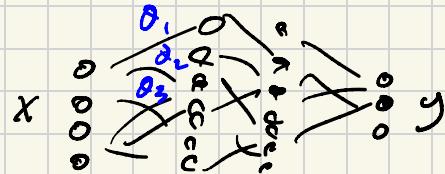
$$\nabla(\nabla f)$$

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



$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$f_\theta(x)$  is a NN  
 ↗ params



$$x \in \mathbb{R}^k$$

$$y \in \mathbb{R}^3$$

$$x_2 = \sigma(w_1 x_1 + b_1)$$

$\equiv$

$\theta$

sample set  $\{x_i\}$  compute the  $y_i = f(x_i)$

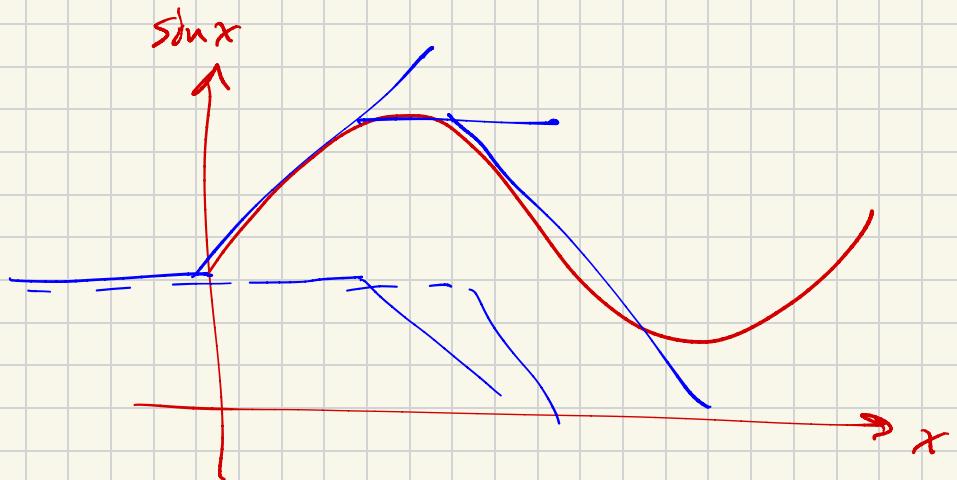
training data  $\{(x_i, y_i)\}$

Loss function  $L_\theta = \sum_i \|f_\theta(x_i) - y_i\|^2$

$$\theta \mapsto L_\theta$$

$$\mathbb{R}^{q_1} \rightarrow \mathbb{R}$$

$$\theta^{k+1} = \theta^k - \alpha \nabla_{\theta} L(\theta^k)$$



$$y_3 = \text{ReLU}(w_2 x_2 + b_2)$$

$\uparrow$   
 weight  
 $\uparrow$   
 bias

