

Neural Networks as Universal Approximators



Pratik Shah
Jan 20, 2023

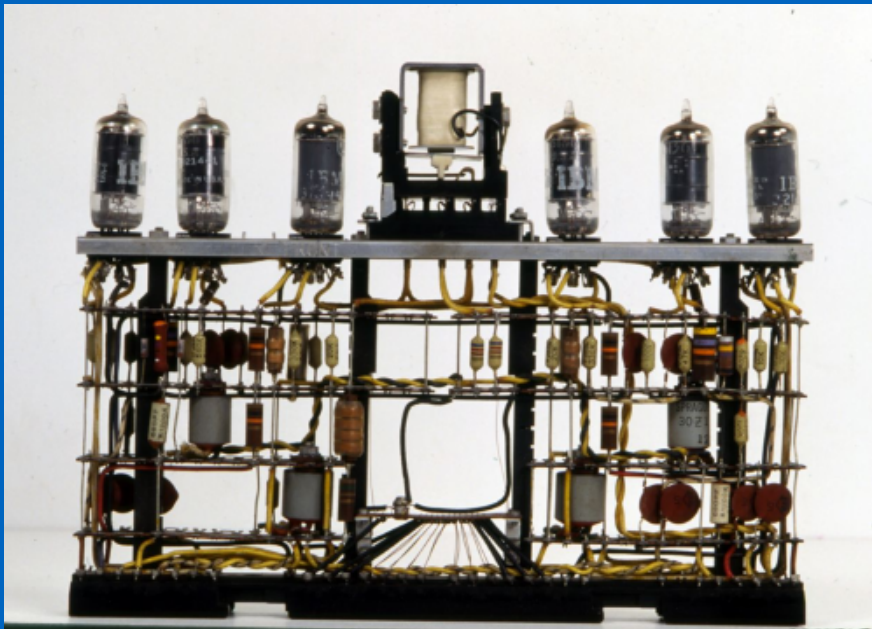
1943 - McCulloch and Pitts

In 1943, neurophysiologist Warren McCulloch and mathematician Walter Pitts wrote a paper on how neurons might work. In order to describe how neurons in the brain might work, they modeled a simple neural network using electrical circuits.

"A Logical Calculus of the Ideas Immanent in Nervous Activity."

1949 - "Organization of the Behaviour" Donald Hebb

1957 - Frank Rosenblatt - first paper - "perceptron"



IBM
704



1957-62 Perceptron - HYPE AI

1969 Perceptron, Minsky and Papert
(limitations and critique to AI)

1982-85 John Hopfield (Caltech)

Today...

- What is learning?
- What is the role of data in learning?

Data $\{(x_i, y_i)\}_{i \in I}$ $x_i \in \mathbb{R}^n, y_i \in \mathbb{R}^m$
 $\left\{ \begin{array}{l} \text{training} \\ \text{instances} \end{array} \right\} \begin{array}{l} \nearrow \text{I/P} \\ \nwarrow \text{O/P} \end{array}$

$$P(Y_i | X_i) \Rightarrow P(Y_i | X_i, \Theta)$$

↑ parameters



$$P(Y_i | X_i, S_T, \Theta)$$

↑ w, b

↑ training set

$$P(Y_i | X_i, r, S_T, \Theta)$$

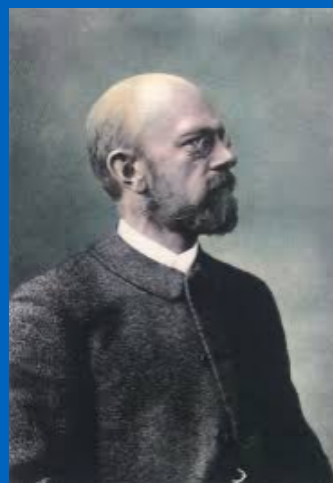
↑ seed

Corbonne, Paris

Further back in 1900

8th August

David
Hilbert



ICM

The 13th Hilbert's problem:

$$x^7 + ax^3 + bx^2 + cx + 1 = 0$$

- Is it possible to write its solution x , as a function of a, b and $c \rightarrow$ as a COMPOSITION of a finite number of two variable functions?

$$x = f(a, b, c)$$

$$\overbrace{\phi_2(\underbrace{\phi_1(a, b)}, c)}$$

$$f = \sum a_i \phi_i$$

linear combination

$$f = \phi_n \circ \phi_{n-1} \circ \dots \circ \phi_2 \circ \phi_1$$

composition of functions.

1957 Arnold (Kolmogorov)

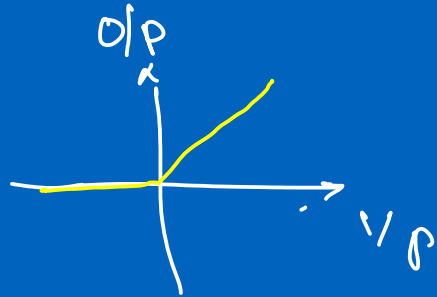
"Can every continuous function of 3 variables be expressed as a composition of finitely many continuous functions of two variables?"

$$f(\underline{x, y, z}) = x^3 + xy^2z + zx = 0$$

$$f(x, y, z) = \Phi_n(\overline{\Phi_{n-1}(x)}, \overline{y})$$

Neuron:

• ReLU



• tanh, Sigmoid

$[-1, +1]$ $[0, 1]$

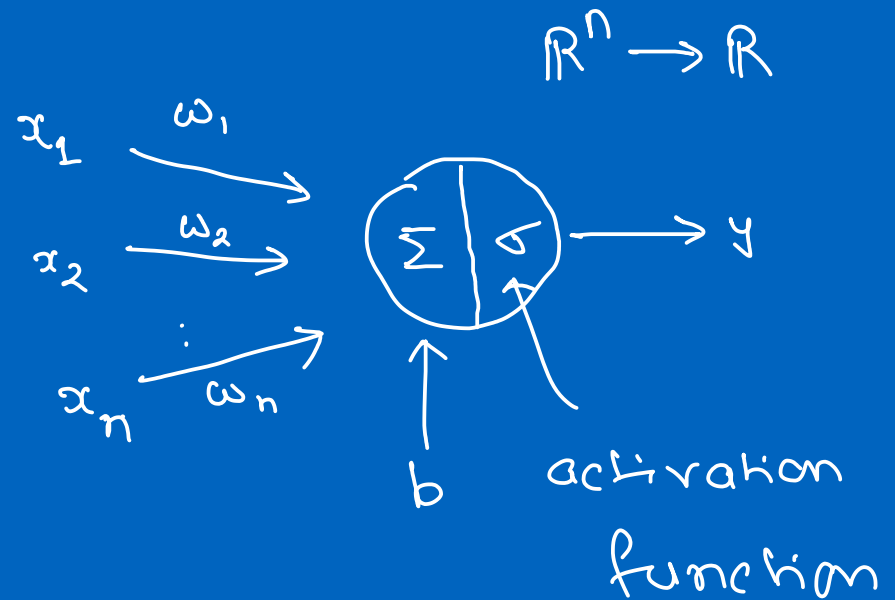
• Nonlinearity!

$$x = [x_1 \ x_2 \ \dots \ x_n]^T$$

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

$$\omega = [\omega_1 \ \omega_2 \ \dots \ \omega_n]^T$$

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

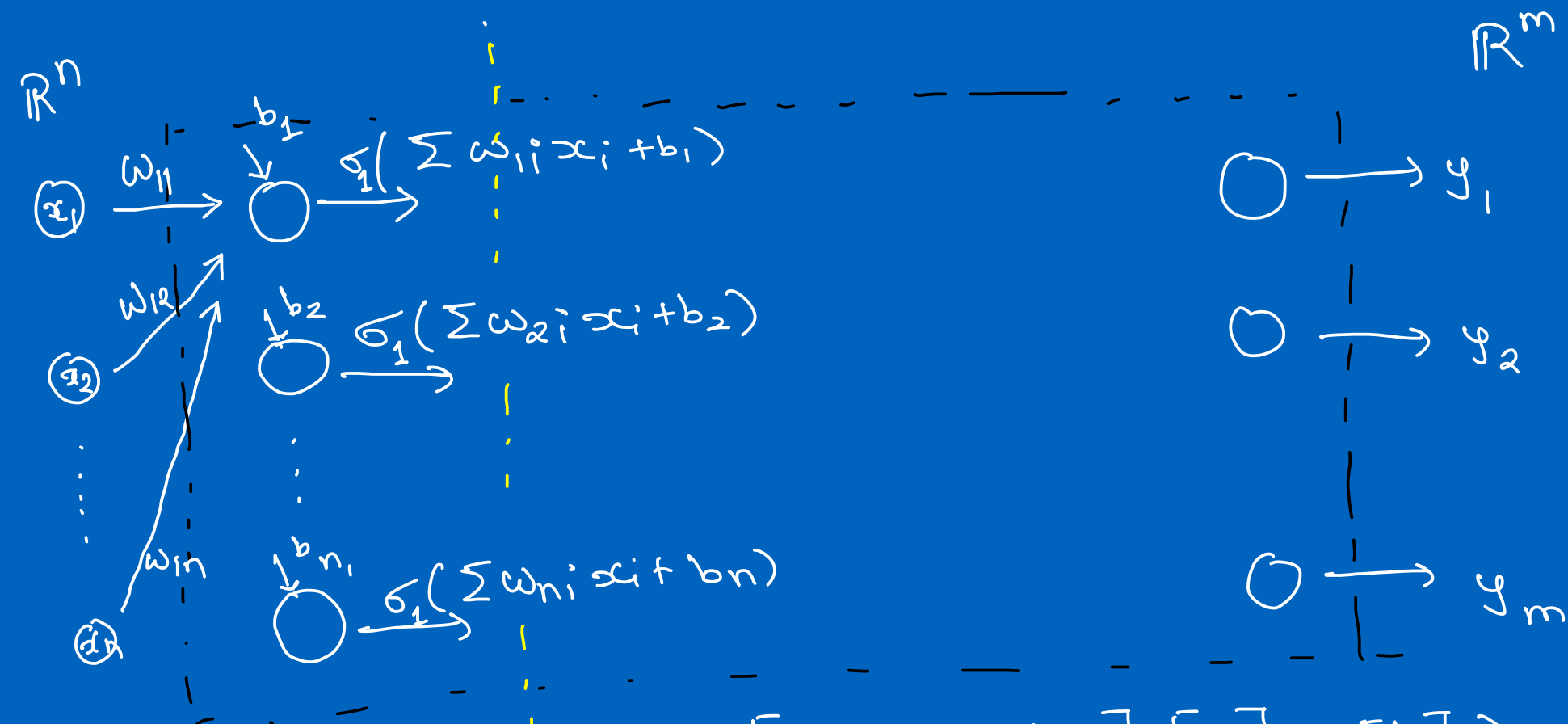


$$y =$$

$$\sigma \left(\sum_{i=1}^n \omega_i x_i + b \right)$$

↖ affine function

$$\sigma(\omega^T x + b)$$



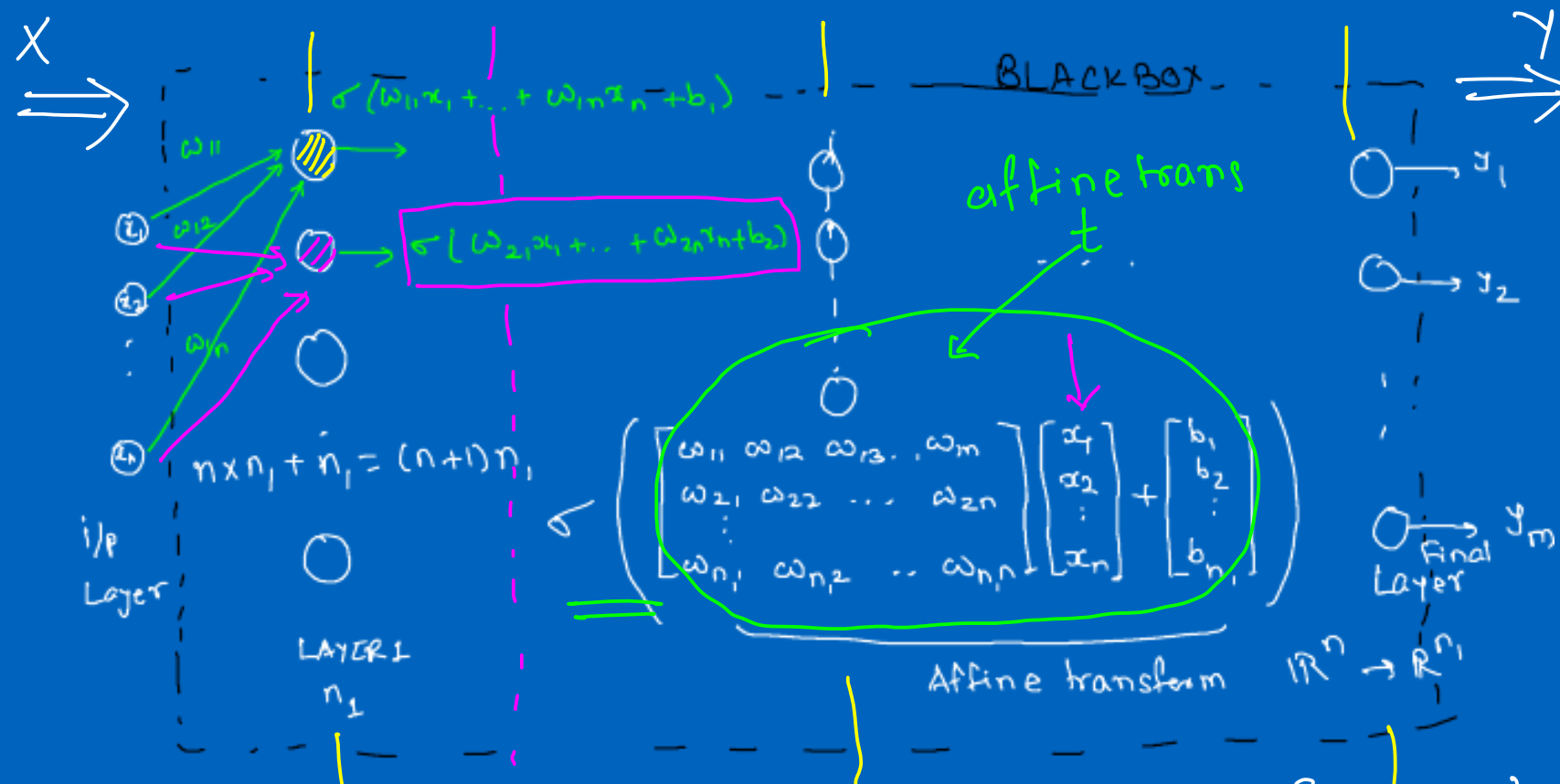
n_1
Layer.1

$$\sigma_1 \left(\begin{bmatrix} \omega_{11} & \omega_{12} & \dots & \omega_{1n} \\ \omega_{21} & \omega_{22} & \dots & \omega_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_{n_1 1} & \dots & \omega_{n_1 n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{n_1} \end{bmatrix} \right)$$

Layer 1

$$t_1: \mathbb{R}^n \rightarrow \mathbb{R}^{n_1}$$

$$\text{Output} \rightarrow \sigma_1 \circ t_1(x)$$



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

$$\downarrow n_1$$

$$\sigma_1(Wx + b) \in \mathbb{R}^{n_1}$$

$$\vdots$$

$$\downarrow$$

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

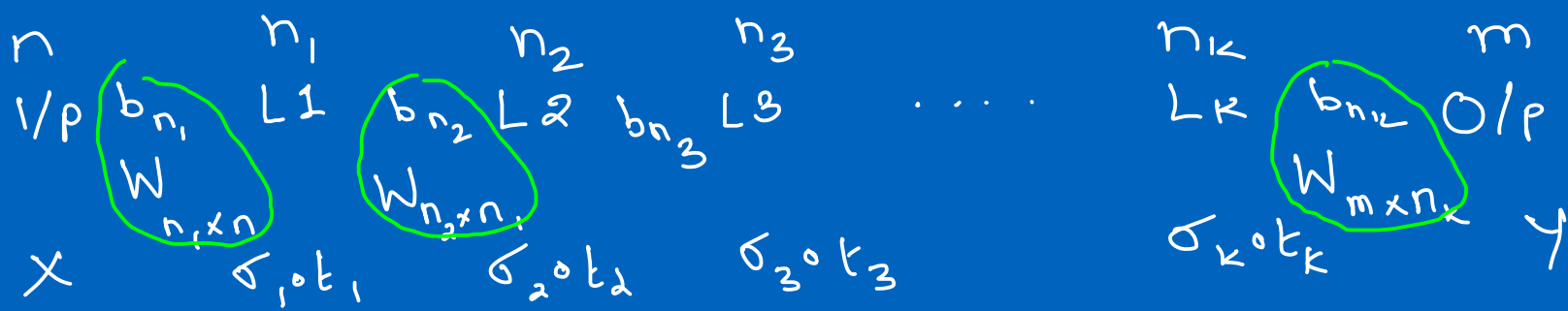
$f(x)$ = $t_n \circ \dots \circ \sigma_3 \circ t_3 \circ \sigma_2 \circ t_2 \circ \sigma_1 \circ t_1(x)$

σ -net $\rightarrow f$ $\sigma = \{\sigma_1, \dots, \sigma_k\}$

$$f(x^{(i)}, \theta) \sim y^{(i)}$$

$$\min_{\theta} \sum_{i=1}^n (f(x^{(i)}, \theta) - y^{(i)})^2$$

$$\theta = (w^i, b^i)$$



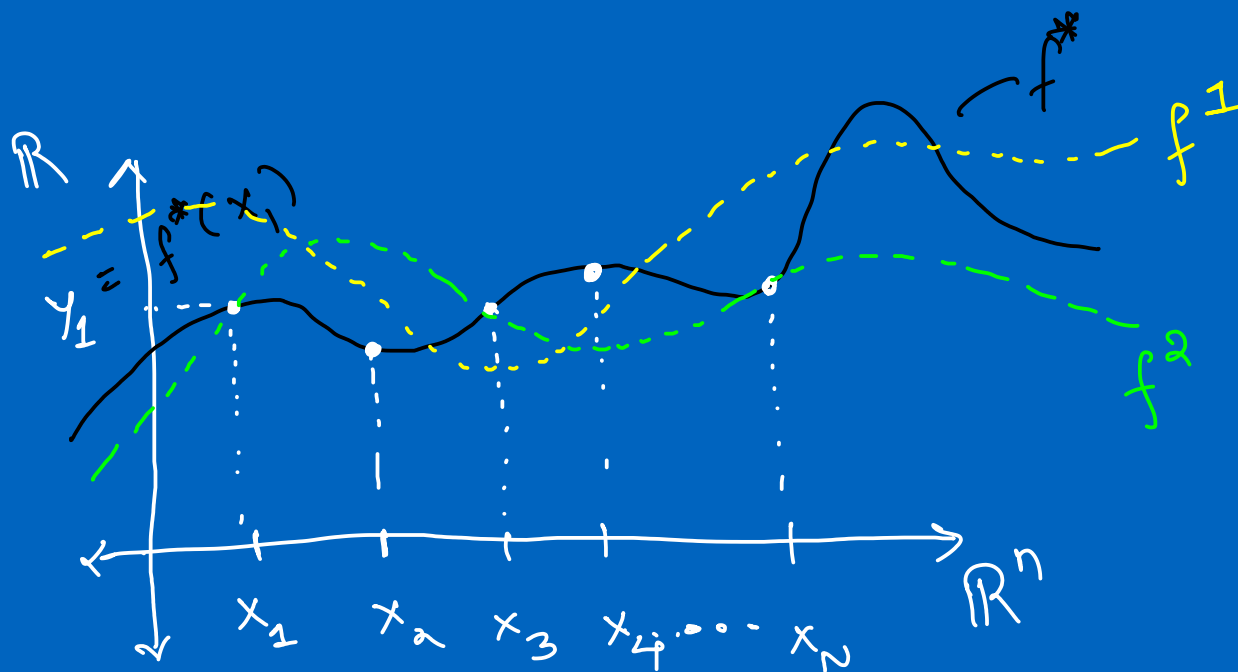
$$Y = \underbrace{\sigma_k \circ t_k \circ \sigma_{k-1} \circ t_{k-1} \circ \dots \circ \sigma_2 \circ t_2 \circ \sigma_1 \circ t_1}_{\text{set of activation functions}}(X).$$

σ - set of activation functions $\sigma_i \in \sigma$

σ -net $Y := f(x, \overbrace{W, b}^{\Theta})$ $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$

$C(\mathbb{R}^n, \mathbb{R}^m) := \left\{ \begin{array}{l} \text{Set of all continuous functions} \\ \text{from } \mathbb{R}^n \rightarrow \mathbb{R}^m \end{array} \right\}$

Universal Approximation:



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

$$D = \{ (x_i, y_i) \}_{i \in \{1, \dots, N\}}$$

$$\sigma\text{-net } f^1 : \mathbb{R}^n \rightarrow \mathbb{R}$$

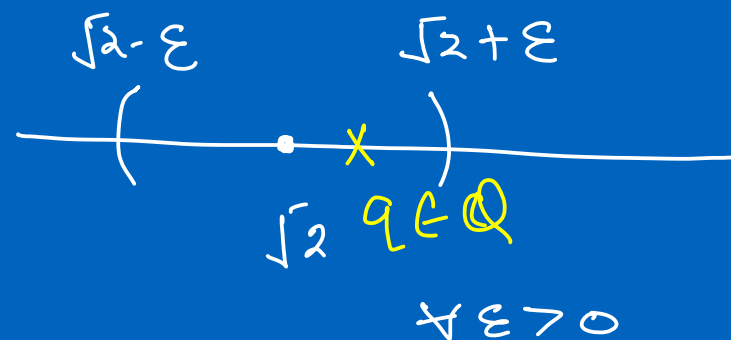
$$d(f^1, f^*) = \sum_{i=1}^N (f^1(x_i) - f^*(x_i))^2$$

$$\mathbb{Q} \subset \mathbb{R}$$

rational $\sqrt{2}$
irrationals

• Floating point arithmetic

\mathbb{Q} is dense in \mathbb{R} .

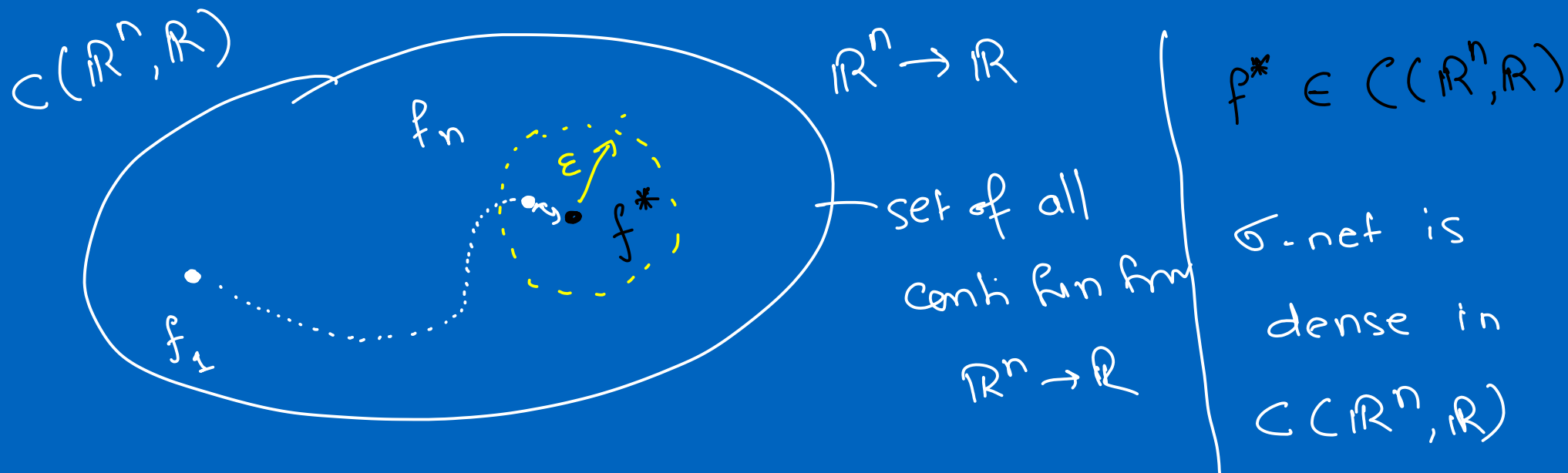


$$d(a, b) = |a - b|$$

minimize
w, b

$$e_f(w, b) = \sum_{i=1}^N (f(x_i, w, b) - y_i)^2$$

$$w_{i+1} \leftarrow w_i - \underbrace{\eta}_{\text{learning rate}} \nabla_{w, b} e_f(w, b)$$



Reference	Function class	Activation ρ	Upper / lower bounds
Lu et al. (2017)	$L^1(\mathbb{R}^{d_x}, \mathbb{R})$	ReLU	$d_x + 1 \leq w_{\min} \leq d_x + 4$
	$L^1(\mathcal{K}, \mathbb{R})$	ReLU	$w_{\min} \geq d_x$
Hanin and Sellke (2017)	$C(\mathcal{K}, \mathbb{R}^{d_y})$	ReLU	$d_x + 1 \leq w_{\min} \leq d_x + d_y$
Johnson (2019)	$C(\mathcal{K}, \mathbb{R})$	uniformly conti. [†]	$w_{\min} \geq d_x + 1$
Kidger and Lyons (2020)	$C(\mathcal{K}, \mathbb{R}^{d_y})$	conti. nonpoly [‡]	$w_{\min} \leq d_x + d_y + 1$
	$C(\mathcal{K}, \mathbb{R}^{d_y})$	nonaffine poly	$w_{\min} \leq d_x + d_y + 2$
	$L^p(\mathbb{R}^{d_x}, \mathbb{R}^{d_y})$	ReLU	$w_{\min} \leq d_x + d_y + 1$
Ours (Theorem 1)	$L^p(\mathbb{R}^{d_x}, \mathbb{R}^{d_y})$	ReLU	$w_{\min} = \max\{d_x + 1, d_y\}$
Ours (Theorem 2)	$C([0, 1], \mathbb{R}^2)$	ReLU	$w_{\min} = 3 > \max\{d_x + 1, d_y\}$
Ours (Theorem 3)	$C(\mathcal{K}, \mathbb{R}^{d_y})$	ReLU+STEP	$w_{\min} = \max\{d_x + 1, d_y\}$
Ours (Theorem 4)	$L^p(\mathcal{K}, \mathbb{R}^{d_y})$	conti. nonpoly [‡]	$w_{\min} \leq \max\{d_x + 2, d_y + 1\}$

ICLR 2020, Min width for Universal approximation

Park, Lee

