

## PERCEPTRONS AND MULTILAYER PERCEPTRONS

Vincent Barra

LIMOS, UMR 6158 CNRS, Université Clermont Auvergne

PERCEPTRON

MULTILAYER PERCEPTRONS

# THRESHOLD LOGIC UNIT

## Mc Culloch and Pitts, 1943

First mathematical model for a neuron

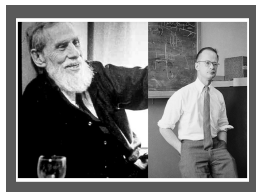
For  $x$  boolean vector,  $w, b \in \mathbb{R}$ :

$$f(x) = \mathbb{1}_{\{w \sum_i x_i + b \geq 0\}}$$

and in particular

- ▶  $OR(x, y) = \mathbb{1}_{\{x+y-0.5 \geq 0\}}$
- ▶  $AND(x, y) = \mathbb{1}_{\{x+y-1.5 \geq 0\}}$
- ▶  $NOT(x) = \mathbb{1}_{\{-x+0.5 \geq 0\}}$

Any Boolean function can be build with such units.



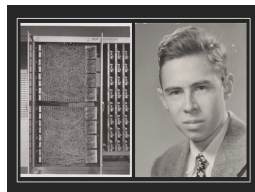
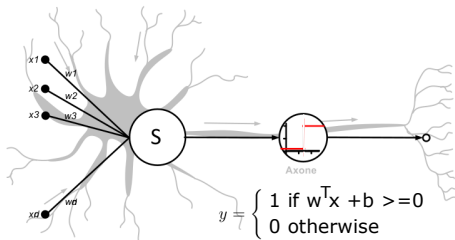
## PERCEPTRON

## Rosenblatt 1957

Generalization:  $w, x \in \mathbb{R}^d, b \in \mathbb{R}$

$$f(x) = \mathbb{1}_{\{w^T x + b \geq 0\}}$$

Relation to biology



# PERCEPTRON

## A more general view

$$f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

where

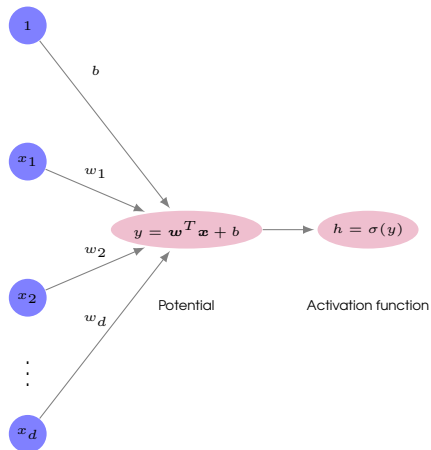
- ▶  $\mathbf{w}$ : synaptic weights
- ▶  $b$ : bias
- ▶  $\mathbf{w}^T \mathbf{x}$  : post synaptic potential
- ▶  $\sigma$ : activation function

# REPRESENTING THE PERCEPTRON

## Graphical representations

1 "Neural" representation

2



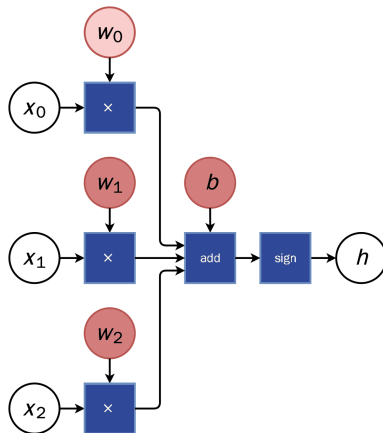
# REPRESENTING THE PERCEPTRON

## Graphical representations

1 "Neural" representation

2 **Computational graph**

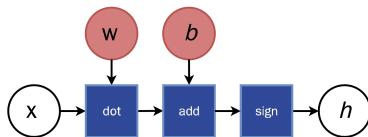
- ▶ white nodes: inputs and outputs
- ▶ red nodes: model parameters
- ▶ blue nodes: operations



# REPRESENTING THE PERCEPTRON

## Graphical representations

- 1 "Neural" representation
- 2 **Computational graph**
  - ▶ white nodes: inputs and outputs
  - ▶ red nodes: model parameters
  - ▶ blue nodes: operations

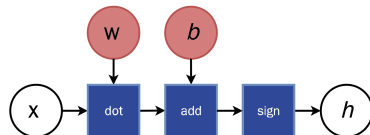




# REPRESENTING THE PERCEPTRON

## Basic brick

This unit is the basic brick of all neural networks



# LEARNING THE PERCEPTRON

## Problem statement

How to build the model ?

- ▶ Input: Learning set  $Z = \{(\mathbf{x}_i, y_i), i \in \llbracket 1 \cdots n \rrbracket, \mathbf{x}_i \in \mathbb{R}^{d+1}, y_i \in \mathbb{R}\}$
- ▶ Unknown:  $\mathbf{w} \in \mathbb{R}^{d+1}$

# LEARNING THE PERCEPTRON

## Problem statement

How to build the model ?

- ▶ Input: Learning set  $Z = \{(\mathbf{x}_i, y_i), i \in [1 \cdots n], \mathbf{x}_i \in \mathbb{R}^{d+1}, y_i \in \mathbb{R}\}$
- ▶ Unknown:  $\mathbf{w} \in \mathbb{R}^{d+1}$

## Key Idea

For each  $\mathbf{x}_i \in Z$ :

- ▶ expected output:  $y_i$
- ▶ computed output:  $h_i = \sigma(\mathbf{w}^T \mathbf{x}_i) = f_{\mathbf{w}}(\mathbf{x})$

If  $\mathcal{L} : \mathbb{R}^{d+1} \times \mathbb{R}^{d+1} \rightarrow \mathbb{R}$  is a loss function

$$\hat{\mathbf{w}} = \mathit{Arg} \min_{\mathbf{w}} \sum_{(\mathbf{x}, y) \in Z} \mathcal{L}(f_{\mathbf{w}}(\mathbf{x}), y)$$

## EXAMPLES OF LOSS FUNCTIONS

### Binary classification (-1/1)

- 1 Characteristic function:  $\mathcal{L}(f_{\mathbf{w}}(x), y) = \mathbb{1}_{yf_{\mathbf{w}}(x) \leq 0}$
- 2 Logistic loss :  $\mathcal{L}(f_{\mathbf{w}}(x), y) = \ln \left( 1 + e^{-yf_{\mathbf{w}}(x)} \right)$
- 3 binary cross-entropy:  $\mathcal{L}(f_{\mathbf{w}}(x), y) = -(y \log(f_{\mathbf{w}}(x)) + (1 - y) \log(1 - f_{\mathbf{w}}(x)))$

### Regression

- 1 Hinge loss :  $\mathcal{L}(f_{\mathbf{w}}(x), y) = (1 - yf_{\mathbf{w}}(x))_+ = \max(0, 1 - yf_{\mathbf{w}}(x))$
- 2 MSE ( $L_2$  loss) :  $\mathcal{L}(f_{\mathbf{w}}(x), y) = \|f_{\mathbf{w}}(x) - y\|^2$
- 3 Huber loss :  $\mathcal{L}(f_{\mathbf{w}}(x), y) = \begin{cases} \frac{1}{2\epsilon}(f_{\mathbf{w}}(x) - y)^2 & \text{if } |f_{\mathbf{w}}(x) - y| \geq \epsilon \\ 0 & \text{otherwise} \end{cases}$
- 4 Vapnik loss:  $\mathcal{L}(f_{\mathbf{w}}(x), y) = \begin{cases} 0 & \text{if } |f_{\mathbf{w}}(x) - y| \leq \epsilon \\ |f_{\mathbf{w}}(x) - y| - \epsilon & \text{otherwise} \end{cases}$

## FIRST TRAINING ALGORITHM

Here,  $\sigma(x) \in \{-1, 1\}$

Given a training set

$$Z = \{(\mathbf{x}_i, y_i), i \in [1 \dots n], \mathbf{x}_i \in \mathbb{R}^{d+1}, y_i \in \{-1, 1\}\}$$

this linear operator can be trained for a binary classification problem.

$\mathbf{w}^0 = \mathbf{0}$

$k = 0$

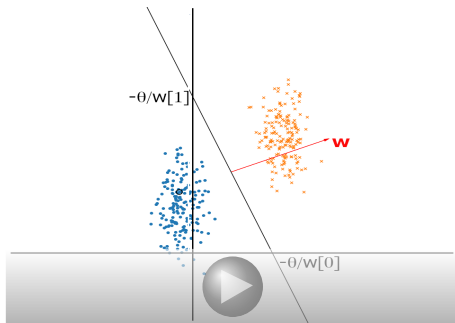
**while**  $\exists i$  such that

$y_i((\mathbf{w}^k)^T \mathbf{x}_i) \leq 0$  **do**

$\mathbf{w}^{k+1} = \mathbf{w}^k + y_i \mathbf{x}_i$

$k = k + 1$

**end**



## FIRST TRAINING ALGORITHM

Convergence iff:

- ▶ Points lie in a sphere of radius  $R$ :  $(\forall i \in \llbracket 1 \cdots n \rrbracket) \|\mathbf{x}_i\| \leq R$
- ▶ The two classes can be separated by a margin:

$$\exists \tilde{\mathbf{w}}, \|\tilde{\mathbf{w}}\| = 1 \exists \gamma > 0, (\forall i \in \llbracket 1 \cdots n \rrbracket) y_i(\tilde{\mathbf{w}}^T \mathbf{x}_i) \geq \gamma/2$$

If so, the perceptron stops as soon as it finds a separating hyperplane.

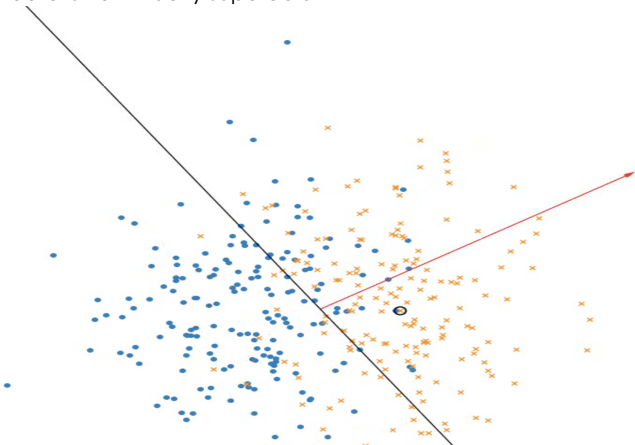
## FIRST TRAINING ALGORITHM

Convergence iff:

- ▶ Points lie in a sphere of radius  $R$ :  $(\forall i \in \llbracket 1 \cdots n \rrbracket) \|\mathbf{x}_i\| \leq R$
- ▶ The two classes can be separated by a margin:

$$\exists \tilde{\mathbf{w}}, \|\tilde{\mathbf{w}}\| = 1 \exists \gamma > 0, (\forall i \in \llbracket 1 \cdots n \rrbracket) y_i (\tilde{\mathbf{w}}^T \mathbf{x}_i) \geq \gamma/2$$

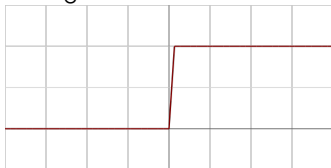
If so, the perceptron stops as soon as it finds a separating hyperplane. But what if the data is non linearly separable ?



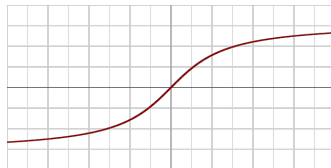
## SECOND TRAINING ALGORITHM

One possible solution: minimize the amount of errors.

- 1 Change  $\sigma$  function to make it differentiable



→



- 2 Error

$$\ell(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in Z} \mathcal{L}(f_{\mathbf{w}}(\mathbf{x}), y)$$

- 3 Minimize the error w.r.t  $\mathbf{w}$ .



## SECOND TRAINING ALGORITHM

### Gradient

At a local minimum the gradient is null: 
$$\sum_{(\mathbf{x}, y) \in Z} \nabla_{\mathbf{w}} \mathcal{L}(f_{\mathbf{w}}(\mathbf{x}), y) = \mathbf{0}$$

## SECOND TRAINING ALGORITHM

### Gradient

At a local minimum the gradient is null: 
$$\sum_{(\mathbf{x}, y) \in Z} \nabla_{\mathbf{w}} \mathcal{L}(f_{\mathbf{w}}(\mathbf{x}), y) = \mathbf{0}$$

### Gradient Descent Algorithm

1 Initialization:  $\mathbf{w} = \mathbf{w}_0, k = 0$

2 While (non stop)

$$2.1 \mathbf{g}_k = \frac{1}{|Z|} \sum_{(\mathbf{x}, y) \in Z} \nabla_{\mathbf{w}} \mathcal{L}(f_{\mathbf{w}_k}(\mathbf{x}), y)$$

$$2.2 \mathbf{w}_{k+1} = \mathbf{w}_k - \eta \mathbf{g}_k$$

$$2.3 k = k + 1$$

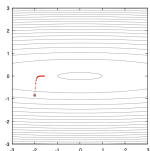
### Additional resource

See Slides “toy example” and “Optimization for deep Learning”.

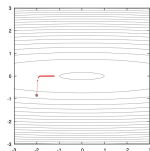
## SECOND TRAINING ALGORITHM

Algorithm parameters:

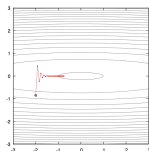
- ▶ stopping criterion
- ▶  $\eta$ : learning rate
- ▶ Weight initialization



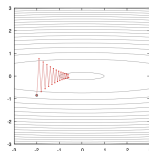
$$\eta = 10^{-2}$$



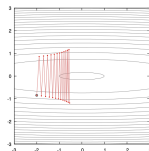
$$\eta = 2.10^{-2}$$



$$\eta = 4.10^{-2}$$



$$\eta = 5.10^{-2}$$

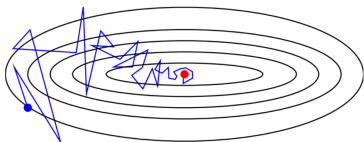
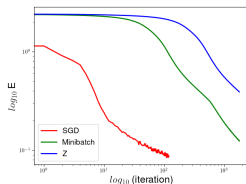


$$\eta = 5.310^{-2}$$

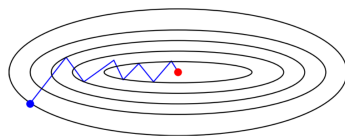
## SECOND TRAINING ALGORITHM

### Different learning strategies

- ▶ Compute the error over all  $Z$ : real gradient descent
- ▶ Compute the error on one example only: stochastic gradient descent (SGD)
- ▶ Compute the error on a batch of example: batch learning (minibatch)



SGD



minibatch

But...

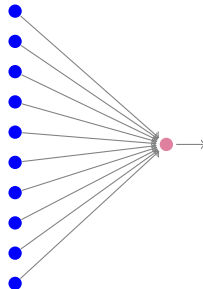
If we want to accurately classify the data (and allow a good generalization property), we need to find something else...

## Stacking linear classifiers

A linear classifier of the form

$$f : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$$

$$\mathbf{x} \mapsto \sigma(\mathbf{w}^T \mathbf{x} + b)$$



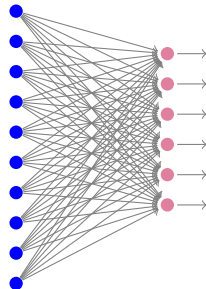
## Stacking linear classifiers

A linear classifier of the form

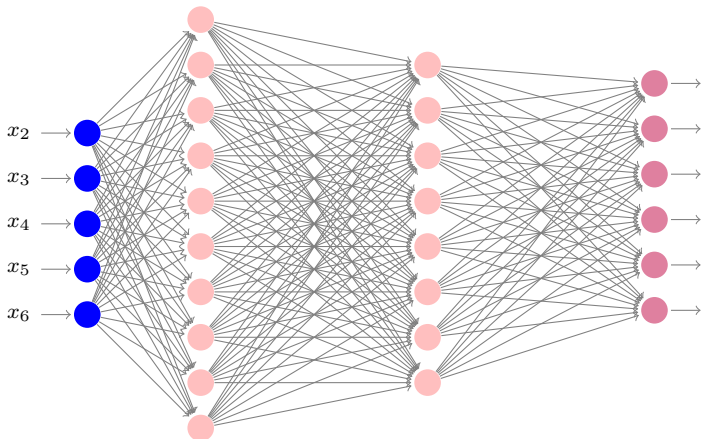
$$f : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$$

$$\mathbf{x} \mapsto \sigma(\mathbf{w}^T \mathbf{x} + b)$$

can naturally be component-wise extended to any function  $f : \mathbb{R}^{d+1} \rightarrow \mathbb{R}^c$



And even...





The general structure can be defined using  $\mathbf{x}^{(0)} = \mathbf{x}$  and

$$(\forall l \in \llbracket 1 \cdots L \rrbracket) \quad \mathbf{x}^{(l)} = \sigma(\mathbf{w}^{(l)T} \mathbf{x}^{(l-1)} + b^{(l)})$$

This is a *Multilayer Perceptron (MLP)*.

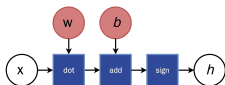
# BUILDING COMPLEX NEURAL NETWORKS

$$h = \sigma(\mathbf{w}^T x + b)$$

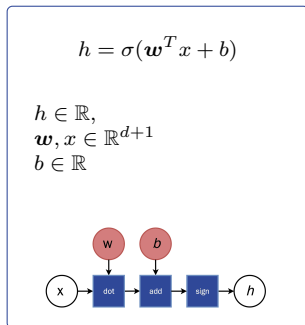
$$h \in \mathbb{R},$$

$$\mathbf{w}, x \in \mathbb{R}^{d+1}$$

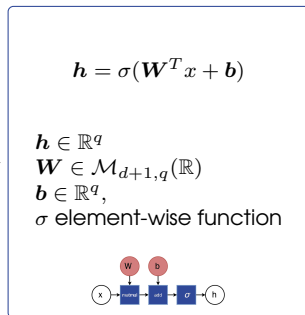
$$b \in \mathbb{R}$$



## BUILDING COMPLEX NEURAL NETWORKS



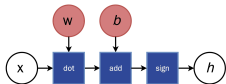
Parallel composition →



## BUILDING COMPLEX NEURAL NETWORKS

$$h = \sigma(\mathbf{w}^T x + b)$$

$$\begin{aligned} h &\in \mathbb{R}, \\ \mathbf{w}, x &\in \mathbb{R}^{d+1} \\ b &\in \mathbb{R} \end{aligned}$$



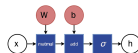
Parallel composition →



100x speed up

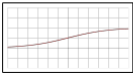
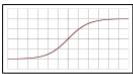
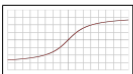
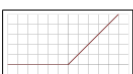
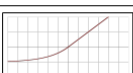
$$h = \sigma(\mathbf{W}^T x + b)$$

$$\begin{aligned} h &\in \mathbb{R}^q \\ \mathbf{W} &\in \mathcal{M}_{d+1,q}(\mathbb{R}) \\ \mathbf{b} &\in \mathbb{R}^q, \\ \sigma &\text{ element-wise function} \end{aligned}$$



$h$  is the output of a layer.

$\sigma$  has to be non linear (otherwise equivalent to a perceptron).

Name	Graph	$f$	$f'$
<i>Logistic / sigmoid</i>		$f(x) = \frac{1}{1+e^{-x}}$	$f'(x) = f(x)(1-f(x))$
<i>tanh</i>		$f(x) = \frac{2}{1+e^{-2x}} - 1$	$f'(x) = 1 - f^2(x)$
<i>atan</i>		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2+1}$
<i>ReLU</i>		$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{if } x \neq 0 \\ 1 & \text{if } x = 0 \end{cases}$
<i>Linear exponential</i>		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$

# LEARNING THE MLP

## Expanding the gradient descent

- ▶ At step  $k$  of the gradient descent, need to evaluate

$$\nabla_{\theta} \mathcal{L}(f_{\theta}(\mathbf{x}), y)$$

- ▶ Evaluation of the total derivatives  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_j}$  and  $\frac{\partial \mathcal{L}}{\partial \mathbf{b}_j}$ ,  $j \in \llbracket 1 \dots L \rrbracket$

⇒ Automatic differentiation on the computational graph

## LEARNING THE MLP

## Expanding the gradient descent

- ▶ At step  $k$  of the gradient descent, need to evaluate

$$\nabla_{\theta} \mathcal{L}(f_{\theta}(\mathbf{x}), y)$$

- ▶ Evaluation of the total derivatives  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_j}$  and  $\frac{\partial \mathcal{L}}{\partial \mathbf{b}_j}$ ,  $j \in \llbracket 1 \dots L \rrbracket$

⇒ Automatic differentiation on the computational graph

## Chain Rule

Let  $g : \mathbb{R} \rightarrow \mathbb{R}^m$  and  $f : \mathbb{R}^m \rightarrow \mathbb{R}$

$$f \circ g(x) = f(\mathbf{u}) = y \text{ where } \mathbf{u} = g(x) = (g_1(x) \dots g_m(x))^T = (u_1 \dots u_m)$$

Chain rule:

$$\frac{dy}{dx} = \sum_{j=1}^m \frac{\partial y}{\partial u_j} \underbrace{\frac{du_j}{dx}}_{\text{recursive}}$$

# LEARNING THE MLP

## Automatic differentiation

- ▶ MLP = composition of differentiable functions
- ▶ The total derivatives of the loss can be evaluated backward, by applying the chain rule recursively over its computational graph.



# LEARNING THE MLP

## Automatic differentiation

- ▶ MLP = composition of differentiable functions
- ▶ The total derivatives of the loss can be evaluated backward, by applying the chain rule recursively over its computational graph.

## Automatic differentiation

- 1 Forward pass: values are all computed from inputs to outputs
- 2 Backward pass: the total derivatives are computed by walking through all paths from outputs to parameters in the computational graph and accumulating the terms.

# LEARNING THE MLP

## Automatic differentiation

- ▶ MLP = composition of differentiable functions
- ▶ The total derivatives of the loss can be evaluated backward, by applying the chain rule recursively over its computational graph.

## Automatic differentiation

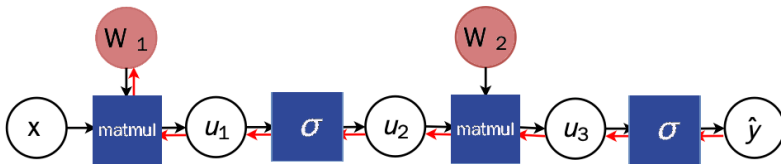
- 1 Forward pass: values are all computed from inputs to outputs
- 2 Backward pass: the total derivatives are computed by walking through all paths from outputs to parameters in the computational graph and accumulating the terms.

## Additional resource

See Slides "backpropagation" and "Vanishing gradient".

# LEARNING THE MLP

Example: derivatives with respect to  $W_1$



- 1 Forward pass:  $u_1, u_2, u_3$  and  $\hat{y}$  computed by traversing the graph, given  $x, W_1$  and  $W_2$
- 2 Backward pass :

$$\begin{aligned} \frac{d\hat{y}}{dW_1} &= \frac{\partial \hat{y}}{\partial u_3} \frac{\partial u_3}{\partial u_2} \frac{\partial u_2}{\partial u_1} \frac{\partial u_1}{\partial W_1} \\ &= \frac{\partial \sigma(u_3)}{\partial u_3} \frac{\partial W_2^T u_2}{\partial u_2} \frac{\partial \sigma(u_1)}{\partial u_1} \frac{\partial W_1^T u_1}{\partial W_1} \end{aligned}$$

Evaluating the partial derivatives requires the intermediate values computed forward

# UNIVERSAL APPROXIMATION

## Theorem (Cybenko 1989; Hornik et al, 1991)

Let  $\sigma$  be a bounded, non-constant continuous function.  
 Let  $I_d$  denote the  $d$ -dimensional hypercube, and  $C(I_d)$  denote the space of continuous functions on  $I_d$ .

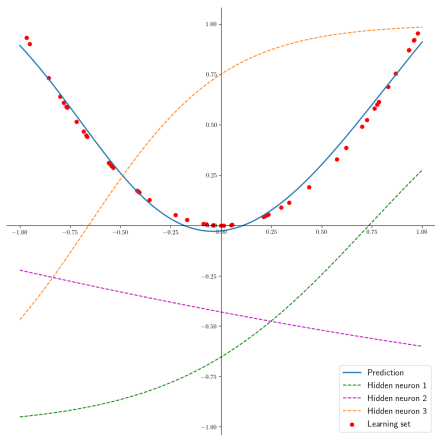
$(\forall f \in C(I_d))(\forall \epsilon > 0)(\exists q > 0, v_i, \mathbf{w}_i, b_i, i \in \llbracket 1 \dots q \rrbracket)$  such that

$$F(\mathbf{x}) = \sum_{i=1}^q v_i \sigma(\mathbf{w}_i^T \mathbf{x} + b_i)$$

satisfies

$$\sup_{\mathbf{x} \in I_d} |f(\mathbf{x}) - F(\mathbf{x})| < \epsilon$$

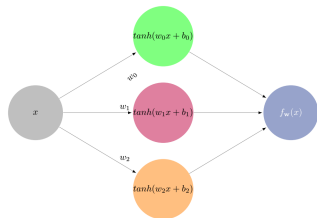
# UNIVERSAL APPROXIMATION



$$f(x) = x^2, |Z| = 50$$

## A simple example

- ▶  $|Z|$  points uniformly sampled (red) over the definition set
- ▶ 1 hidden layer MLP, 3 neurons.
- ▶  $\tanh$  activation function, and linear output neurons
- ▶ network output : blue curve
- ▶ hidden neurons outputs: dashed curves



# UNIVERSAL APPROXIMATION

## Properties

- ▶ Guarantees that a single hidden layer network can represent any classification problem in which the boundary is locally linear (smooth)
- ▶ Does not inform about good/bad architectures, nor how they relate to the optimization procedure
- ▶ Generalizes to any non-polynomial (possibly unbounded) activation function, including the ReLU

# UNIVERSAL APPROXIMATION

## Theorem (Barron, 1992)

Let a one-hidden layer MLP with  $q$  hidden neurons,  $p$  inputs and  $|Z| = n$ . The mean integrated square error between the estimated network  $\hat{F}$  and the target function  $f$  is bounded by

$$O\left(\frac{C_f^2}{q} + \frac{qp}{n} \log(n)\right)$$

where  $C_f$  measures the global smoothness of  $f$ .

## Properties

- ▶ Combines approximation and estimation errors.
- ▶ Provided enough data, guarantees that adding more neurons will result in a better approximation

## EFFECT OF DEPTH

### Theorem (Montúfar et al, 2014)

A MLP with ReLU as activation functions,  $p$  inputs,  $L$  hidden layers with  $q \geq p$  neurons can compute functions having  $\Omega\left(\left(\frac{q}{p}\right)^{(L-1)p} q^p\right)$  linear regions (asymptotic lower bound).

### Properties

- ▶ The number of linear regions of deep models grows exponentially in  $L$  and polynomially in  $q$ .
- ▶ Even for small values of  $L$  and  $q$ , deep rectifier models are able to produce substantially more linear regions than shallow rectifier models.