



Università degli  
Studi di Pavia

# *Deep Learning*

## *02-Artificial Neural Networks Basic Ideas, Notations and all that*

Marco Piastra

*This presentation can be downloaded at:*  
<http://vision.unipv.it/DL>

*Function approximation:  
Linear Combination*

# Function Approximation: linear combination

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

a.k.a. "single layer perceptron"

A first approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

i.e. this is a vector of dimension  $d$

*Note that, when the input is scalar, the approximator becomes*

$$\tilde{y} = wx + b$$

*i.e. a straight line*

# Function Approximation: linear combination

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

A first approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

## **dataset**

*A set of actual inputs and outputs is all we know about the target function*

$$D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \quad y^{(i)} = f^*(\mathbf{x}^{(i)}), \forall i$$

# Function Approximation: linear combination

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

A first approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

## **dataset**

*A set of actual inputs and outputs is all we know about the target function*

$$D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N, \quad y^{(i)} = f^*(\mathbf{x}^{(i)}), \forall i$$

*Three other fundamental aspects to be considered:*

- **representation:** *which parametric approximator for a given target function?*
- **evaluation:** *how do you tell that some parameter values are better than others?*
- **optimization:** *how can we learn optimal values for the parameters?*

# Function Approximation: linear combination

- Example: XOR

$$y = \text{XOR}(\mathbf{x}), \quad \mathbf{x} \in \{0, 1\}^2$$

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Dataset:

$$D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$$

$x_1$	$x_2$	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

this is our *dataset* ( $N = 4$ )

# Function Approximation: linear combination

## ■ Example: XOR

$$y = \text{XOR}(\mathbf{x}), \quad \mathbf{x} \in \{0, 1\}^2$$

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Dataset:

$$D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$$

Loss function (evaluation):

$$L(\mathbf{x}^{(i)}, y^{(i)}) := (\tilde{y}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

$$L(D) := \frac{1}{N} \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in D} L(\mathbf{x}^{(i)}, y^{(i)})$$

$x_1$	$x_2$	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

this is our *dataset* ( $N = 4$ )

Squared Error

Mean Squared Error (MSE)

# Function Approximation: linear combination

## ■ Example: XOR

$$y = \text{XOR}(\mathbf{x}), \quad \mathbf{x} \in \{0, 1\}^2$$

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Dataset:

$$D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$$

Optimization problem:

We need to find

$$(\mathbf{w}, b)^* = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} L(D)$$

i.e. the set of parameter values that minimizes loss w.r.t. to the dataset

$x_1$	$x_2$	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

this is our *dataset* ( $N = 4$ )



# Function Approximation: linear combination

- Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Loss function:

$$\begin{aligned} L(D) &:= \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}^{(i)}, y^{(i)}) \\ &= \frac{1}{N} \sum_{i=1}^N (\tilde{y}(\mathbf{x}^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{N} \sum_{i=1}^N ((\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)})^2 \end{aligned}$$

*Can we express this summation by using linear algebra?*

*As we will see later on, matrix representation may lead to a better **parallelization** of computations*

# Function Approximation: linear combination

- Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Loss function:

$$L(D) = \frac{1}{N} \sum_{i=1}^N ((\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)})^2$$

define:

$$\mathbf{X} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix} \quad \text{input data in matrix form (item index first)}$$

# Function Approximation: linear combination

## ■ Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Loss function:

$$L(D) = \frac{1}{N} \sum_{i=1}^N ((\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)})^2$$

define:

$$\hat{\mathbf{X}} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} & 1 \end{bmatrix} \quad \boldsymbol{\vartheta} := \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix} \quad \mathbf{y} := \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

The loss function becomes:

$$L(D) = \frac{1}{N} (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y})^2$$

**loss function in matrix form**  
— This is a positive-definite quadratic form

# Function Approximation: linear combination

- Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Loss function:

$$L(D) = \frac{1}{N} \sum_{i=1}^N ((\mathbf{w} \cdot \mathbf{x}^{(i)} + b) - y^{(i)})^2$$

define:

$$\hat{\mathbf{X}} := \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ x_1^{(N)} & \dots & x_d^{(N)} & 1 \end{bmatrix} \quad \boldsymbol{\vartheta} := \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix} \quad \mathbf{y} := \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

The loss function becomes:

$$L(D) = \frac{1}{N} (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y})^2$$

**loss function in matrix form**  
— This is a positive-definite quadratic form

# Function Approximation: linear combination

- Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Loss function:

$$L(D) = \frac{1}{N} (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y})^2$$

For XOR:

$$\hat{\mathbf{X}} := \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad \boldsymbol{\vartheta} := \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} \quad \mathbf{y} := \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

XOR

$x_1$	$x_2$	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

this is our *dataset* ( $N = 4$ )

# Function Approximation: linear combination

## ■ Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

*representation*

Loss function:

$$L(D) = \frac{1}{N} (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y})^2$$

*evaluation*

Optimization:

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) = 0$$

*optimization*

the loss function is convex:  
by solving this equation we can find  $\boldsymbol{\vartheta}^*$   
i.e. the optimal parameter values

# Function Approximation: linear combination

- Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Optimization:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) &= \frac{1}{N} \frac{\partial}{\partial \boldsymbol{\vartheta}} (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y})^2 \\ &= \frac{1}{N} \frac{\partial}{\partial \boldsymbol{\vartheta}} (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y})^T (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y}) = \frac{1}{N} \frac{\partial}{\partial \boldsymbol{\vartheta}} (\boldsymbol{\vartheta}^T \hat{\mathbf{X}}^T - \mathbf{y}^T) (\hat{\mathbf{X}} \boldsymbol{\vartheta} - \mathbf{y}) \\ &= \frac{1}{N} \frac{\partial}{\partial \boldsymbol{\vartheta}} (\boldsymbol{\vartheta}^T \hat{\mathbf{X}}^T \hat{\mathbf{X}} \boldsymbol{\vartheta} - \boldsymbol{\vartheta}^T \hat{\mathbf{X}}^T \mathbf{y} - \mathbf{y}^T \hat{\mathbf{X}} \boldsymbol{\vartheta} + \mathbf{y}^T \mathbf{y}) \\ &= \frac{1}{N} \frac{\partial}{\partial \boldsymbol{\vartheta}} (\boldsymbol{\vartheta}^T \hat{\mathbf{X}}^T \hat{\mathbf{X}} \boldsymbol{\vartheta} - 2\boldsymbol{\vartheta}^T \hat{\mathbf{X}}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \\ &= \frac{1}{N} (2\hat{\mathbf{X}}^T \hat{\mathbf{X}} \boldsymbol{\vartheta} - 2\hat{\mathbf{X}}^T \mathbf{y}) \end{aligned}$$

*all these terms are scalars*

# Function Approximation: linear combination

- Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

Optimization:

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) = \frac{1}{N} (2\hat{\mathbf{X}}^T \hat{\mathbf{X}} \boldsymbol{\vartheta} - 2\hat{\mathbf{X}}^T \mathbf{y})$$

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) = 0 \quad \implies \quad 2\hat{\mathbf{X}}^T \hat{\mathbf{X}} \boldsymbol{\vartheta} - 2\hat{\mathbf{X}}^T \mathbf{y} = 0$$

$$\hat{\mathbf{X}}^T \hat{\mathbf{X}} \boldsymbol{\vartheta} = \hat{\mathbf{X}}^T \mathbf{y}$$

$$\boldsymbol{\vartheta} = (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \mathbf{y} \quad \text{this is what we need}$$

this matrix is SQUARE  
and, typically, with actual datasets,  
is invertible (i.e. full rank)



# Function Approximation: linear combination

## Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

For XOR:

$$\boldsymbol{\vartheta} = (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \mathbf{y}$$

$$\hat{\mathbf{X}} := \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad \boldsymbol{\vartheta} := \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} \quad \mathbf{y} := \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$\hat{\mathbf{X}}^T \hat{\mathbf{X}} = \begin{bmatrix} 2 & 1 & 2 \\ 1 & 2 & 2 \\ 2 & 2 & 4 \end{bmatrix} \quad (\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} = \begin{bmatrix} 1 & 0 & 0.5 \\ 0 & 1 & 0.5 \\ 0.5 & 0.5 & 0.75 \end{bmatrix}$$

$$(\hat{\mathbf{X}}^T \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^T \mathbf{y} = \begin{bmatrix} 0 \\ 0 \\ 0.5 \end{bmatrix}$$

XOR	$x_1$	$x_2$	$x_1 \oplus x_2$
	0	0	0
	0	1	1
	1	0	1
	1	1	0

# Function Approximation: linear combination

- Loss minimization

Approximator: *linear combination*

$$\tilde{y} = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

For XOR:

$$\vartheta := \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0.5 \end{bmatrix}$$

hence the XOR linear approximator becomes:

$$\tilde{y} = 0.5$$

*What ???*

XOR	$x_1$	$x_2$	$x_1 \oplus x_2$
	0	0	0
	0	1	1
	1	0	1
	1	1	0

*Function approximation:  
Feed-Forward Neural Network*

# Feed-Forward Neural Network

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Second attempt: **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, b \in \mathbb{R}$$

— this is a non-linear scalar function, applied elementwise  
— i.e. this is a matrix of dimensions  $h \times d$

# Feed-Forward Neural Network

- Approximating a target function

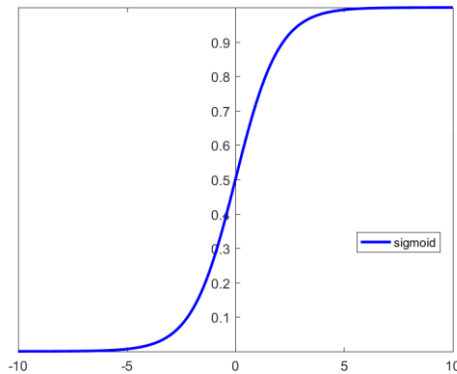
$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Second attempt: **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, b \in \mathbb{R}$$

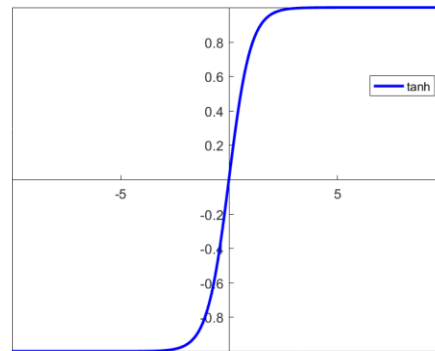
Popular choices for the non-linear function:

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



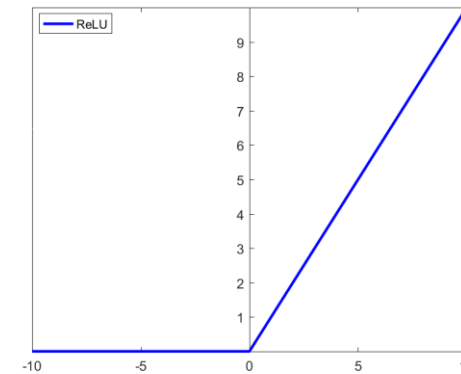
Sigmoid

$$g(x) = \tanh(x)$$



Hyperbolic Tangent

$$g(x) = \max(0, x)$$



ReLU

# Feed-Forward Neural Network

- Approximating a target function

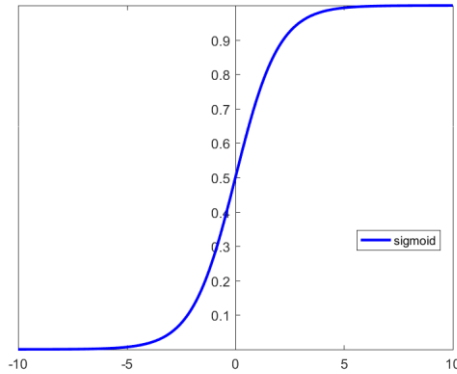
$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Second attempt: **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, \quad b \in \mathbb{R}$$

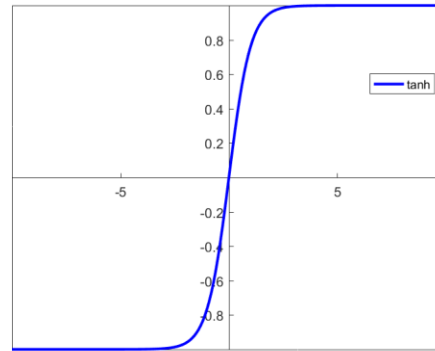
Popular choices for the non-linear function:

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



Sigmoid

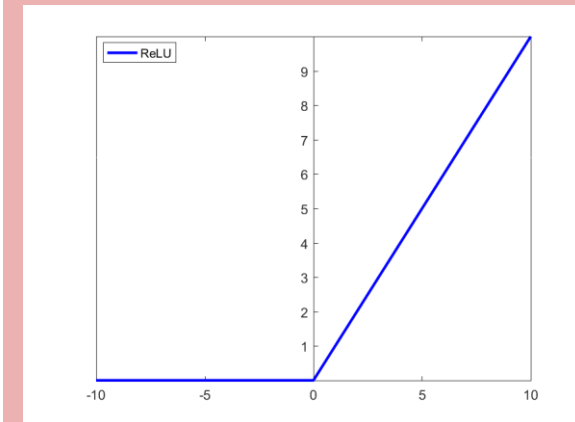
$$g(x) = \tanh(x)$$



Hyperbolic Tangent

*this is somewhat special...*

$$g(x) = \max(0, x)$$



ReLU

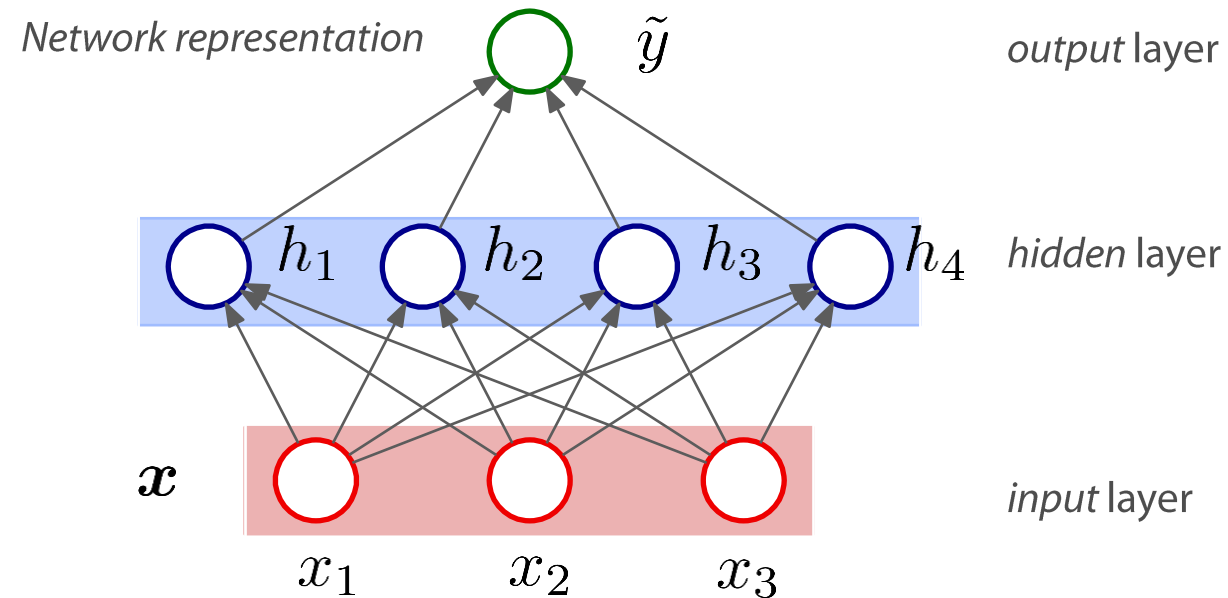
# Feed-Forward Neural Network

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Second attempt: **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, \quad b \in \mathbb{R}$$



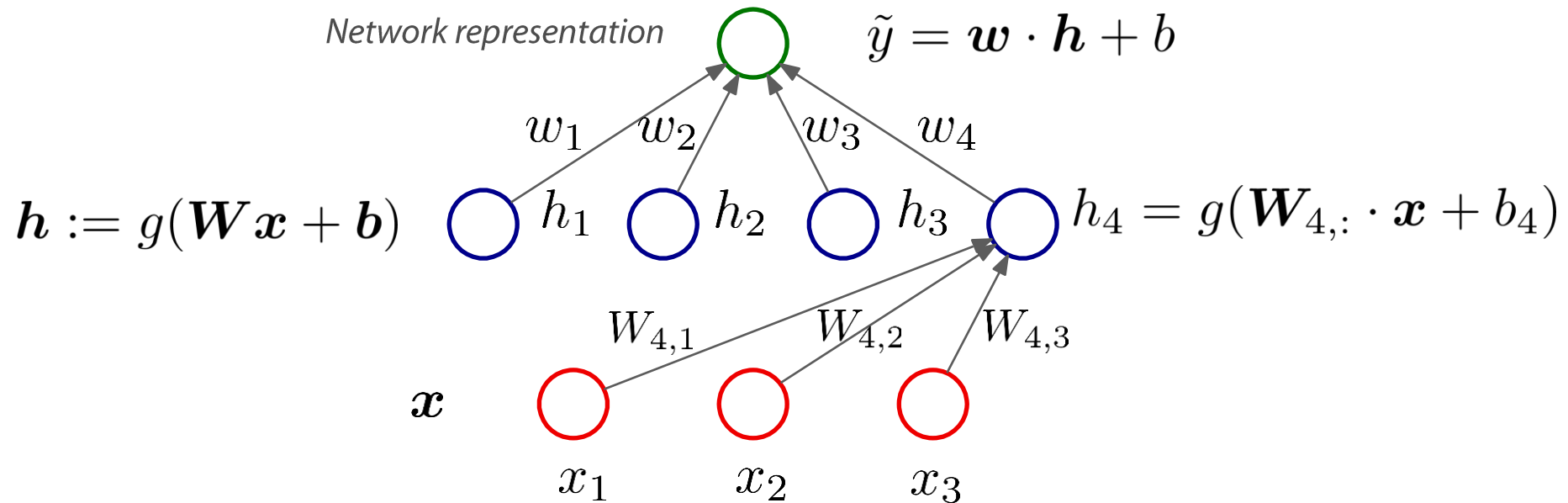
# Feed-Forward Neural Network

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Second attempt: **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, \quad b \in \mathbb{R}$$



NOTE: biases  $\mathbf{b}$  and  $b$  are NOT represented in the graph



# Universality of FF Neural Networks

- **Universal approximation theorem** (Cybenko, 1989; Hornik, 1991; Leshno et al. 1991)

For any target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d \quad (\text{which is continuous and Borel measurable})$$

and any  $\varepsilon > 0$  there exists parameters

$$h \in \mathbb{Z}^+, \mathbf{W} \in \mathbb{R}^{h \times d}, \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, b \in \mathbb{R}$$

*this is the dimension of the hidden layer: it is a parameter in the theorem*

such that the **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b$$

approximates the target function by less than  $\varepsilon$

$$\sup_{\mathbf{x}} | f^*(\mathbf{x}) - (\mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b) | < \varepsilon$$

*(on any compact subset of  $\mathbb{R}^d$ )*

*This theorem holds with any of the non-linear functions seen before*

# Universality of FF Neural Networks

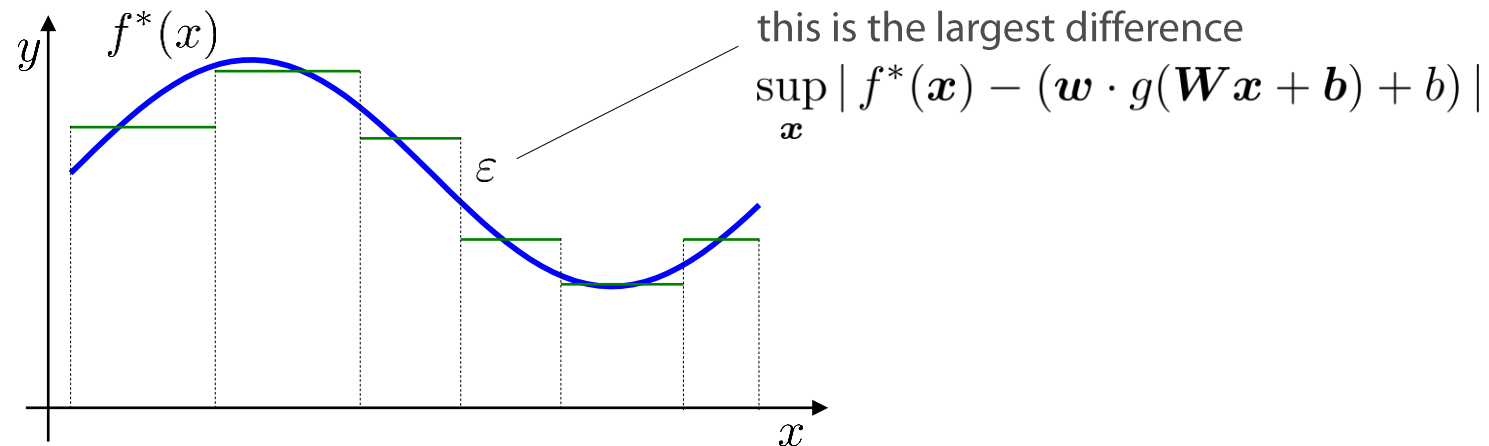
- **Universal approximation theorem** (Cybenko, 1989; Hornik, 1991; Leshno et al. 1991)

## Intuitive rationale

Any continuous target function

$$y = f^*(x), \quad x \in \mathbb{R}$$

can be approximated arbitrarily well by a stepwise function



for simplicity, assume  $x$  is *scalar* (hence  $\mathbf{W}$  is *vector*)

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}x + \mathbf{b}) + b$$

# Universality of FF Neural Networks

- **Universal approximation theorem** (Cybenko, 1989; Hornik, 1991; Leshno et al. 1991)

## Intuitive rationale

Consider the *step function* as the non-linearity

$$\tilde{y} = \mathbf{w} \cdot \text{step}(\mathbf{W}x + \mathbf{b}) + b$$

then, by expanding the approximator

$$\tilde{y} = w_1 \text{step}(W_1x + b_1) + \dots + w_h \text{step}(W_hx + b_h) + b$$

where each step occurs at

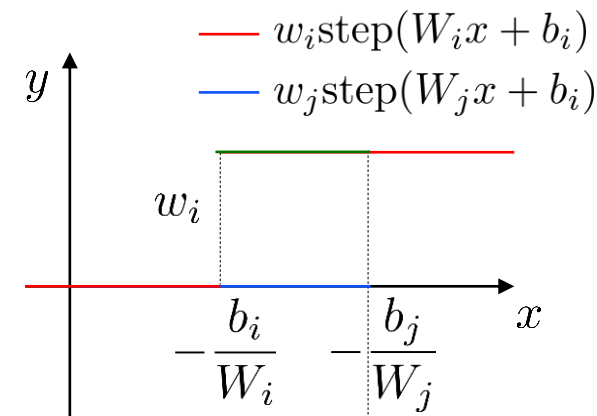
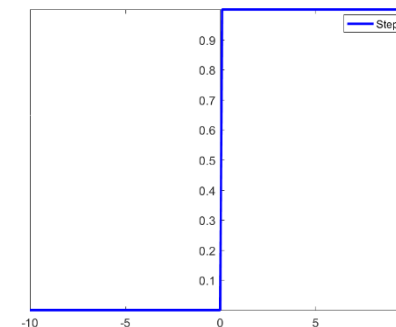
$$W_i \cdot x + b_i = 0 \implies W_i \cdot x = -b_i \implies x = -\frac{b_i}{W_i}$$

Consider *pairs* of steps  $i$  and  $j$  and impose

$$-\frac{b_i}{W_i} < -\frac{b_j}{W_j}, \quad W_i, W_j > 0, \quad w_i = -w_j$$

in this way we can construct  $\frac{h}{2}$  such function steps

$$g(x) = \text{step}(x)$$



# *Learning Feed-Forward Neural Networks*

# Learning with FF Neural Networks

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Second attempt: **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, \quad b \in \mathbb{R}$$

**Optimization problem (learning)**

Given a dataset  $D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ ,  $y^{(i)} = f^*(\mathbf{x}^{(i)})$ ,  $\forall i$

*/ the dimension of the hidden layer is pre-defined*

we want to find parameter values  $\mathbf{W} \in \mathbb{R}^{h \times d}$ ,  $\mathbf{w}, \mathbf{b} \in \mathbb{R}^h$ ,  $b \in \mathbb{R}$

that *minimize* the loss function  $L(D) := \frac{1}{N} \sum_D (\tilde{y}^{(i)} - y^{(i)})^2$

where:  $\tilde{y}^{(i)} := \mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b$

# Learning with FF Neural Networks

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Second attempt: **(shallow) feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, \quad b \in \mathbb{R}$$

## Difficulty

In general, *minimizing* the loss function

$$L(D) = \frac{1}{N} \sum_D ((\mathbf{w} \cdot g(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2$$

cannot be done directly, since

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) = 0$$

cannot be solved analytically

— this loss function is not convex, in general

We need to find another way...

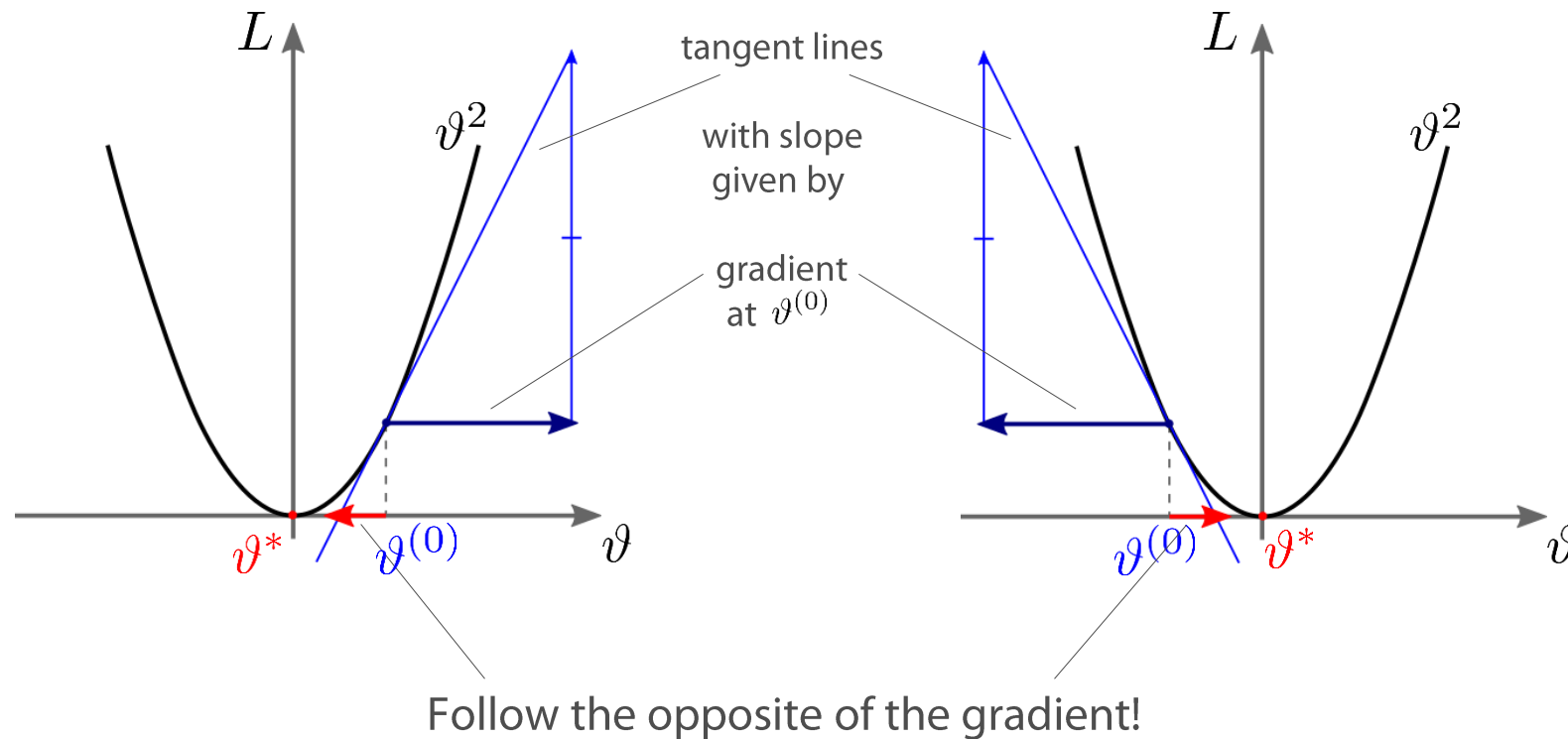
# Gradient Descent (GD): intuition

- Optimization problem

$$\vartheta^* := \operatorname{argmin}_{\vartheta} L(D, \vartheta)$$

Just making the dependence explicit

- Minimizing a generic function



# Gradient Descent (GD): intuition

## ■ Optimization problem

$$\boldsymbol{\vartheta}^* := \operatorname{argmin}_{\boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta})$$

Just making the dependence explicit

## ■ Iterative method

Step in the method

1. Initialize  $\boldsymbol{\vartheta}^{(0)}$  at random
2. Update  $\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta}^{(t-1)})$
3. Unless some termination criterion has been met, go back to step 2.

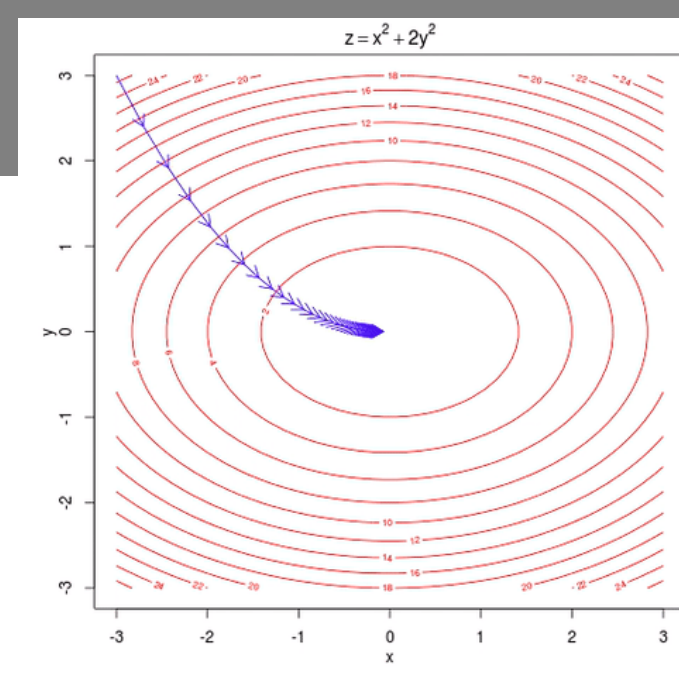
where

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta}) := \frac{1}{N} \sum_D \frac{\partial}{\partial \boldsymbol{\vartheta}} L(\hat{y}^{(i)}, y^{(i)}, \boldsymbol{\vartheta})$$

$$\eta \ll 1$$

The gradient of the loss over the dataset  $D$  is the average of gradients over each data item

A learning rate, it is arbitrary (i.e., an hyperparameter)





# Gradient Descent (GD): convergence

- **Convergence**

When  $L(D, \boldsymbol{\vartheta})$  is *convex, derivable*, and its gradient is *Lipschitz continuous*, that is

$$\left\| \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta}_1) - \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta}_2) \right\| \leq C \|\boldsymbol{\vartheta}_1 - \boldsymbol{\vartheta}_2\|, \quad C > 0$$

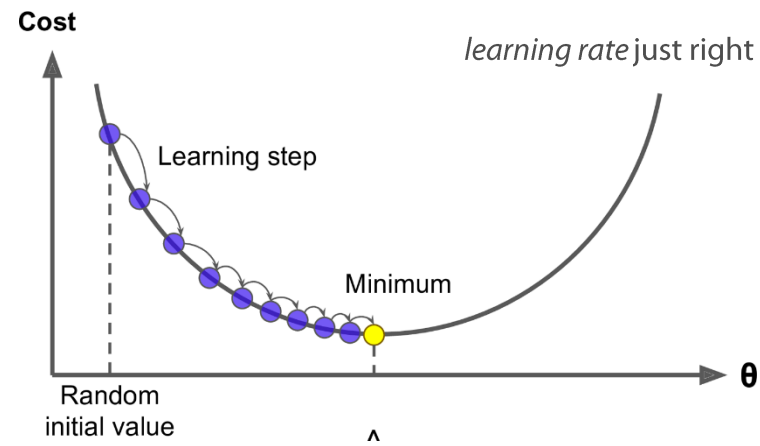
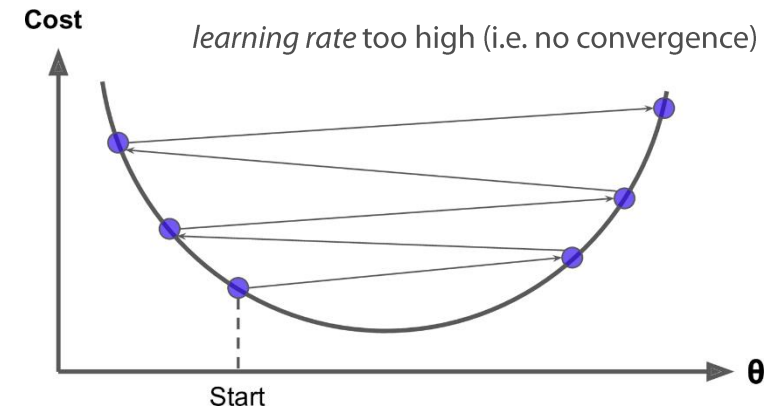
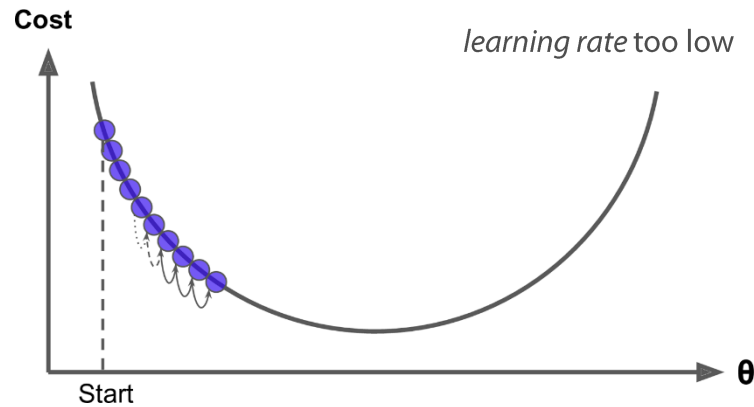
the gradient descent method converges to the optimal  $\boldsymbol{\vartheta}^*$  for  $t \rightarrow \infty$   
provided that  $\eta \leq 1/C$

When  $L(D, \boldsymbol{\vartheta})$  is *derivable* but not *convex*, and its gradient is *Lipschitz continuous*,  
the gradient descent method converges to a local minimum of  $L(D, \boldsymbol{\vartheta})$   
under the same conditions

# Gradient Descent (GD): practicalities

- *Convergence in practice*

The choice of the *learning rate*  $\eta$  is crucial

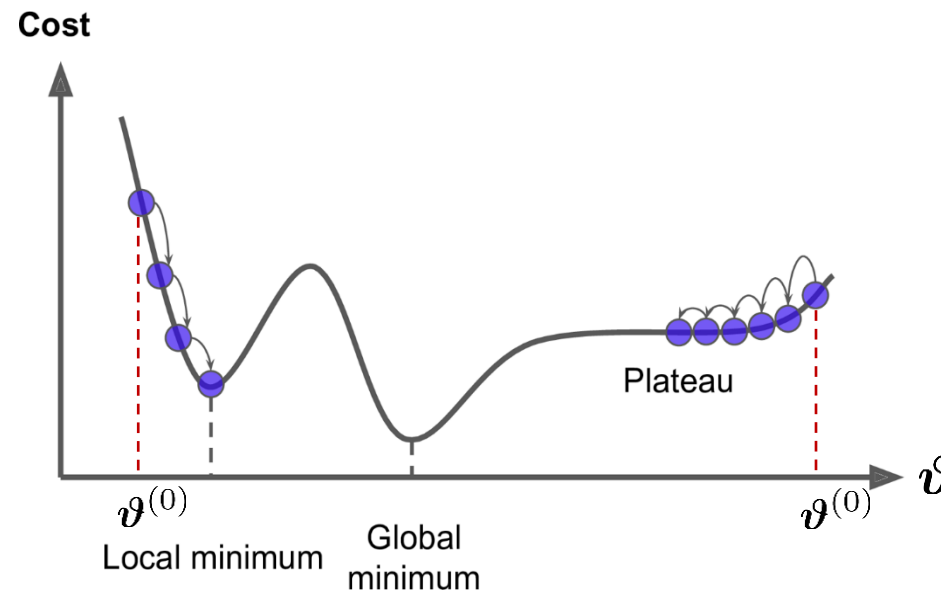


Images from <https://www.safaribooksonline.com/library/view/hands-on-machine-learning/9781491962282/ch04.html>

# Gradient Descent (GD): practicalities

- *Convergence in practice*

When  $L(D, \vartheta)$  is not convex, the **initial estimate**  $\vartheta^{(0)}$  is crucial



The outcome of the method will depend on which  $\vartheta^{(0)}$  is picked

Image from <https://www.safaribooksonline.com/library/view/hands-on-machine-learning/9781491962282/ch04.html>

*Learning  
Feed-Forward Neural Networks  
(contd.)*

# Gradient Descent for FF Neural Networks

Recall that the *item-wise* loss for a specific data item in the dataset is

$$L(\tilde{y}^{(i)}, y^{(i)}) := (\tilde{y}^{(i)} - y^{(i)})^2$$

then

$$L(D) = \frac{1}{N} \sum_D L(\tilde{y}^{(i)}, y^{(i)})$$

and the gradient of the loss function is

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) &= \frac{\partial}{\partial \boldsymbol{\vartheta}} \frac{1}{N} \sum_D L(\tilde{y}^{(i)}, y^{(i)}) \\ &= \frac{1}{N} \sum_D \frac{\partial}{\partial \boldsymbol{\vartheta}} L(\tilde{y}^{(i)}, y^{(i)}) \end{aligned}$$

*Moral: we must be capable to compute the gradient on each data item*

# Gradient Descent for FF Neural Networks

Suppose we can compute the four *item-wise gradients*, w.r.t. to the parameters:

$$\frac{\partial}{\partial \mathbf{W}} L(\tilde{y}^{(i)}, y^{(i)}) \quad \frac{\partial}{\partial \mathbf{b}} L(\tilde{y}^{(i)}, y^{(i)}) \quad \frac{\partial}{\partial \mathbf{w}} L(\tilde{y}^{(i)}, y^{(i)}) \quad \frac{\partial}{\partial b} L(\tilde{y}^{(i)}, y^{(i)})$$

we can then apply a *gradient descent* method

## ■ Gradient Descent

1. Assign initial values to the four parameters
2. Update the four parameters by adding

$$\mathbf{W}^{(0)}, \mathbf{b}^{(0)}, \mathbf{w}^{(0)}, b^{(0)}$$

$$\Delta \mathbf{W} = -\eta \frac{1}{N} \sum_D \frac{\partial}{\partial \mathbf{W}} L(\tilde{y}^{(i)}, y^{(i)}) \quad \Delta \mathbf{b} = -\eta \frac{1}{N} \sum_D \frac{\partial}{\partial \mathbf{b}} L(\tilde{y}^{(i)}, y^{(i)})$$

$$\Delta \mathbf{w} = -\eta \frac{1}{N} \sum_D \frac{\partial}{\partial \mathbf{w}} L(\tilde{y}^{(i)}, y^{(i)}) \quad \Delta b = -\eta \frac{1}{N} \sum_D \frac{\partial}{\partial b} L(\tilde{y}^{(i)}, y^{(i)})$$

3. Unless complete, return to step 2.

# Computing Gradients

All we need to apply the descent method is computing the item-wise gradients

For instance:

$$\begin{aligned}\frac{\partial}{\partial \mathbf{W}} L(\tilde{y}^{(i)}, y^{(i)}) &= \frac{\partial}{\partial \mathbf{W}} (\tilde{y}^{(i)} - y^{(i)})^2 \\ &= \frac{\partial}{\partial \mathbf{W}} ((\mathbf{w} \cdot g(\mathbf{W} \mathbf{x}^{(i)} + \mathbf{b}) + b) - y^{(i)})^2\end{aligned}$$

(similar expressions hold for the other three gradients)

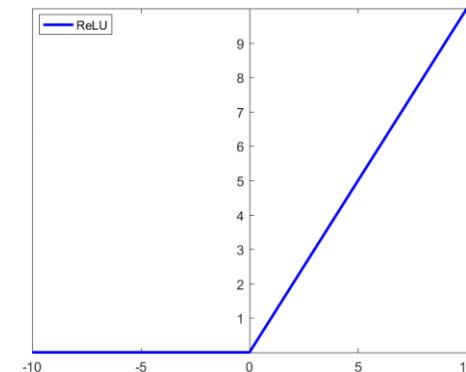
Assume

$$g(x) = \text{ReLU}(x) := \max(0, x)$$

i.e., the non-linearity is ReLU

Easy, huh?

$$g(x) = \max(0, x)$$



# Function Approximation: FF Neural Networks

- Loss minimization

Approximator:  
*(shallow) feed-forward neural network*

$$\tilde{y} = \mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b$$

Optimal values for XOR and  $h = 2$  :

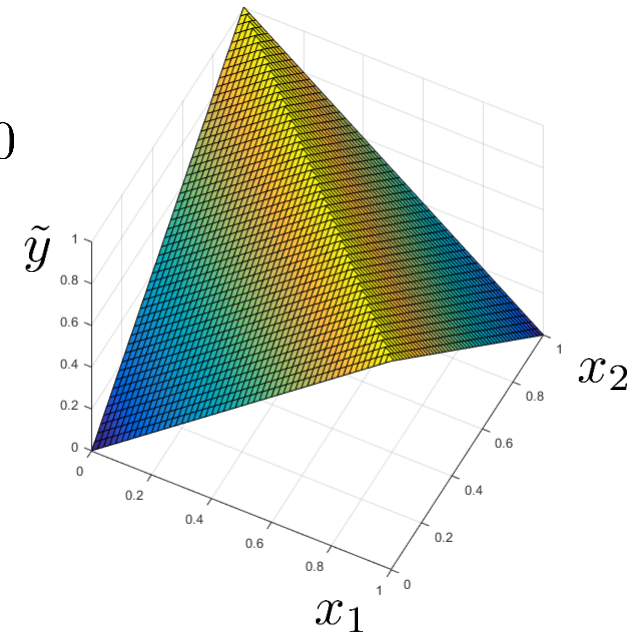
*dimension of the hidden layer*

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

XOR

$x_1$	$x_2$	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

$b = 0$





# *Stochastic and Mini-Batch Gradient Descent*

# Function Approximation: FF Neural Networks

- Loss minimization

Approximator:

*(shallow) feed-forward neural network*

$$\tilde{y} = \mathbf{w} \cdot \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b}) + b$$

In this case our dataset was tiny... ( )  $N = 4$

What if the dataset was very large?

XOR

$x_1$	$x_2$	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

/  
this is our dataset

# Stochastic Gradient Descent (SGD): intuition

- *Objective*

$$\boldsymbol{\vartheta}^* := \operatorname{argmin}_{\boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta})$$

- *Iterative method*

1. Initialize  $\boldsymbol{\vartheta}^{(0)}$  at random
2. Pick a data item  $(\boldsymbol{x}^{(i)}, y^{(i)}) \in D$  with uniform probability
3. Update  $\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta^{(t)} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(\tilde{y}^{(i)}, y^{(i)}, \boldsymbol{\vartheta}^{(t-1)})$
4. Unless some termination criterion has been met, go back to step 2.

$$\eta^{(t)} \ll 1$$

Note that the *learning rate* may vary across iterations...

# Stochastic Gradient Descent for FF Neural Networks

With very large datasets, the sum in:

$$\Delta \boldsymbol{\vartheta} = -\eta \frac{1}{N} \sum_D \frac{\partial}{\partial \boldsymbol{\vartheta}} L(\tilde{y}^{(i)}, y^{(i)})$$

may take very long to compute (and this must be repeated at each iteration)

## ■ **Stochastic Gradient Descent (SGD)** (i.e. "you don't actually need to sum up them all")

1. Assign initial values to the four parameters  $\mathbf{W}^{(0)}$ ,  $\mathbf{b}^{(0)}$ ,  $\mathbf{w}^{(0)}$ ,  $b^{(0)}$
2. Pick up a data item  $(\mathbf{x}^{(i)}, y^{(i)})$  from  $D$  with uniform probability and update the four parameters (with  $\eta \ll 1.0$ ,  $\eta \rightarrow 0$  as iterations progress)

$$\Delta \mathbf{W} = -\eta \frac{\partial}{\partial \mathbf{W}} L(\tilde{y}^{(i)}, y^{(i)}) \qquad \Delta \mathbf{b} = -\eta \frac{\partial}{\partial \mathbf{b}} L(\tilde{y}^{(i)}, y^{(i)})$$

$$\Delta \mathbf{w} = -\eta \frac{\partial}{\partial \mathbf{w}} L(\tilde{y}^{(i)}, y^{(i)}) \qquad \Delta b = -\eta \frac{\partial}{\partial b} L(\tilde{y}^{(i)}, y^{(i)})$$

3. Unless complete, return to step 2.

# Stochastic Gradient Descent (SGD): convergence

- **Convergence**

When  $L(D, \boldsymbol{\vartheta})$  is *convex, derivable*, and its gradient is *Lipschitz continuous*, that is

$$\left\| \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta}_1) - \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta}_2) \right\| \leq C \|\boldsymbol{\vartheta}_1 - \boldsymbol{\vartheta}_2\|, \quad C > 0$$

the stochastic gradient descent method converges to the optimal  $\boldsymbol{\vartheta}^*$  for  $t \rightarrow \infty$  provided that

$$\eta^{(t)} \leq \frac{1}{Ct} \quad \text{Note that } \eta^{(t)} \rightarrow 0 \text{ for } t \rightarrow \infty$$

When  $L(D, \boldsymbol{\vartheta})$  is *derivable*, and its gradient is *Lipschitz continuous* but not convex the stochastic gradient descent method converges to a local minimum of  $L(D, \boldsymbol{\vartheta})$  under the same conditions

# Speed of Convergence

Perhaps surprisingly, **stochastic gradient descent** shares the same properties and could be faster than GD ...

Consider a generic loss function  $L(\vartheta)$  which is *convex* in the parameter  $\vartheta$

Define *accuracy* as an upper bound:

$$|L(\vartheta^*) - L(\tilde{\vartheta})| < \rho$$

optimal value                      current parameter estimate

[from Bottou & Bousquet, 2008]

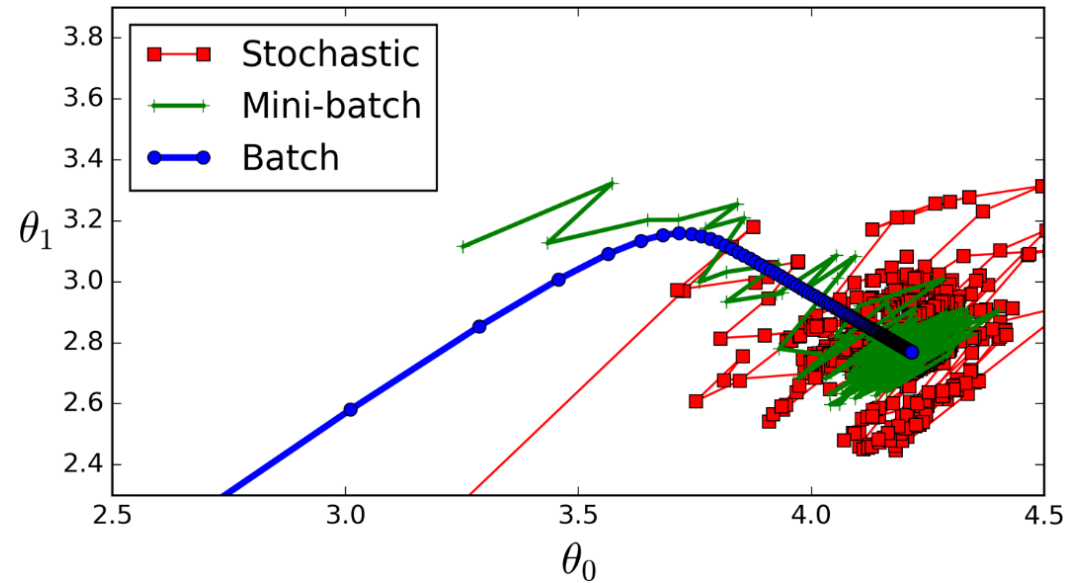
$N$  size of the dataset

$q$  number of (scalar) parameters in  $\vartheta$

Algorithm	Cost per iteration	Iterations to reach accuracy $\rho$	Time to reach accuracy $\rho$
Gradient descent (GD)	$\mathcal{O}(N q)$	$\mathcal{O}\left(\log \frac{1}{\rho}\right)$	$\mathcal{O}\left(N q \log \frac{1}{\rho}\right)$
Stochastic gradient descent (SGD)	$\mathcal{O}(q)$	$\mathcal{O}\left(\frac{1}{\rho}\right)$	$\mathcal{O}\left(q \frac{1}{\rho}\right)$

# Qualitative comparison of GD methods

Typical traces  
of the three methods  
(batch = GD)



In general:

- GD is more regular but slower (with large datasets)
- SGD is faster (with large datasets) but noisy
- MBGD is often the right compromise in practice...

Image from <https://www.safaribooksonline.com/library/view/hands-on-machine-learning/9781491962282/ch04.html>

# Mini-batch Gradient Descent (MBGD): intuition

- *Objective*

$$\boldsymbol{\vartheta}^* := \operatorname{argmin}_{\boldsymbol{\vartheta}} L(D, \boldsymbol{\vartheta})$$

- *Iterative method*

1. Initialize  $\boldsymbol{\theta}^{(0)}$  at random
2. Pick a mini batch  $B \subseteq D$  with uniform probability
3. Update  $\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta^{(t)} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$
4. Unless some termination criterion has been met, go back to step 2.

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}) := \frac{1}{|B|} \sum_B \frac{\partial}{\partial \boldsymbol{\vartheta}} L(\hat{y}^{(i)}, y^{(i)}, \boldsymbol{\vartheta})$$

This method has the same convergence properties of SGD



# Mini-batch Gradient Descent for FF Neural Networks

## ▪ **Mini-batch Gradient Descent (MBGD)**

1. Assign initial values to the four parameters  $\mathbf{W}^{(0)}$ ,  $\mathbf{b}^{(0)}$ ,  $\mathbf{w}^{(0)}$ ,  $b^{(0)}$
2. Pick a *mini-batch*  $B \subseteq D$  with uniform probability and update the four parameters (with  $\eta \ll 1.0$ ,  $\eta \rightarrow 0$  as iterations progress)

$$\Delta \mathbf{W} = -\eta \frac{1}{|B|} \sum_B \frac{\partial}{\partial \mathbf{W}} L(\tilde{y}^{(i)}, y^{(i)}) \quad \Delta \mathbf{b} = -\eta \frac{1}{|B|} \sum_B \frac{\partial}{\partial \mathbf{b}} L(\tilde{y}^{(i)}, y^{(i)})$$

$$\Delta \mathbf{w} = -\eta \frac{1}{|B|} \sum_B \frac{\partial}{\partial \mathbf{w}} L(\tilde{y}^{(i)}, y^{(i)}) \quad \Delta b = -\eta \frac{1}{|B|} \sum_B \frac{\partial}{\partial b} L(\tilde{y}^{(i)}, y^{(i)})$$

3. Unless complete, return to step 2.

This method has the same convergence properties of SGD