Deep Learning

A course about theory & practice

Artificial Neural Networks Basic ideas, notation and all that

Marco Piastra



Approximating a target function

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

a.k.a. "single layer perceptron"

A first approximator: *linear combination*

$$ilde{y} = oldsymbol{w} \cdot oldsymbol{x} + b, \;\; oldsymbol{w} \in \mathbb{R}^d, b \in \mathbb{R}$$
 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

Approximating a target function

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

A first approximator: *linear combination*

$$ilde{y} = \boldsymbol{w} \cdot \boldsymbol{x} + b, \ \ \boldsymbol{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

dataset

A set of input and output pairs (**data items**) is what know about the target function

$$D := \{(\boldsymbol{x}^{(i)}, y^{(i)})\}_{i=1}^{N}, \quad y^{(i)} = f^{*}(\boldsymbol{x}^{(i)}), \forall i$$
A set of data items

$$Item \underline{index}$$

Approximating a target function

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

A first approximator: *linear combination*

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

dataset

A set of input and output pairs (**data items**) is what know about the target function

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

Three fundamental aspects:

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

Deep Learning 2024–2025

• Example: XOR $y = XOR(\boldsymbol{x}), \ \boldsymbol{x} \in \{0,1\}^2$

Approximator: *linear combination*

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

Dataset:

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

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 ($N=4$)			



Approximator: *linear combination*

$$ilde{y} = oldsymbol{w} \cdot oldsymbol{x} + b, \hspace{1em} oldsymbol{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

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

Dataset:

Example: XOR

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

Loss function (evaluation):

Deep Learning 2024-2025

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 ($N=4$)			



Dataset:

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

Optimization problem:

We need to find

parameter values that minimize the loss w.r.t. to the dataset

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

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 ($N=4$)			

Deep Learning 2024-2025

Artificial Neural Networks [8]

Loss minimization

Approximator: *linear combination*

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

Loss function:

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

Can we express this summation by using linear algebra?

Matrix representation leads to a better **parallelization** of computations (more on this, later on)

Deep Learning 2024-2025

Loss minimization

Approximator: *linear combination*

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

Loss function:

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

define:

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

Loss minimization

Approximator: *linear combination*

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

Loss function:

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

define:

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

The loss function becomes:

$$L(D) = rac{1}{N} (\hat{m{X}} m{artheta} - m{y})^2$$

 $herefore This is a positive-definite quadratic form$

Loss minimization

Approximator: *linear combination*

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

Loss function:

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

define:

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

The loss function becomes:

$$L(D) = rac{1}{N} (\hat{m{X}} m{artheta} - m{y})^2$$

 $herefore This is a positive-definite quadratic form$

Loss minimization

Approximator: *linear combination*

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

Loss function:

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

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 ($N=4$)			(-4)

For XOR:

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

Loss minimization

Approximator: *linear combination*

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

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

Loss function:

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

Optimization:

this loss function is <u>convex</u>: by solving this equation, we can find $artheta^*$ i.e. the optimal parameter values representation

evaluation

optimization



Deep Learning 2024-2025

Artificial Neural Networks [14]

Loss minimization

Approximator: *linear combination*

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

Optimization:

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

Loss minimization

Approximator: *linear combination*

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

Optimization:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) &= \frac{1}{N} (2\hat{\boldsymbol{X}}^T \hat{\boldsymbol{X}} \boldsymbol{\vartheta} - 2\hat{\boldsymbol{X}}^T \boldsymbol{y}) \\ \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) &= 0 \implies 2\hat{\boldsymbol{X}}^T \hat{\boldsymbol{X}} \boldsymbol{\vartheta} - 2\hat{\boldsymbol{X}}^T \boldsymbol{y} = 0 \\ \hat{\boldsymbol{X}}^T \hat{\boldsymbol{X}} \boldsymbol{\vartheta} &= \hat{\boldsymbol{X}}^T \boldsymbol{y} \\ \boldsymbol{\vartheta} &= (\hat{\boldsymbol{X}}^T \hat{\boldsymbol{X}})^{-1} \hat{\boldsymbol{X}}^T \boldsymbol{y} \quad \text{this is what we need} \\ & \text{this matrix is SQUARE and SYMMETRIC} \\ & \text{and, typically, with actual datasets} \\ & \text{is invertible (i.e. full rank)} \end{aligned}$$

Loss minimization

Approximator: *linear combination*

$$ilde{y} = oldsymbol{w} \cdot oldsymbol{x} + b, \hspace{0.2cm} oldsymbol{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

For XOR:

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

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

$$\hat{\boldsymbol{X}} := \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad \boldsymbol{\vartheta} := \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} \quad \boldsymbol{y} := \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$
$$\hat{\boldsymbol{X}}^T \hat{\boldsymbol{X}} = \begin{bmatrix} 2 & 1 & 2 \\ 1 & 2 & 2 \\ 2 & 2 & 4 \end{bmatrix} \quad (\hat{\boldsymbol{X}}^T \hat{\boldsymbol{X}})^{-1} = \begin{bmatrix} 1 & 0 & 0.5 \\ 0 & 1 & 0.5 \\ 0.5 & 0.5 & 0.75 \end{bmatrix} \quad (\hat{\boldsymbol{X}}^T \hat{\boldsymbol{X}})^{-1} \hat{\boldsymbol{X}}^T \boldsymbol{y} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.5 \end{bmatrix}$$

()

0

Loss minimization

Approximator: *linear combination*

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

For XOR:

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

hence the XOR linear approximator becomes:

$$\tilde{y} = 0.5$$

What ??? A constant function?

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

Approximating a target function

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

Second attempt: (shallow) feed-forward neural network

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

this is a matrix of dimensions $h \times d$
this is a non-linear scalar function, applied *elementwise*

Deep Learning 2024–2025

Approximating a target function

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

Second attempt: (shallow) feed-forward neural network

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

Popular choices for the non-linear function:



Deep Learning 2024-2025

Approximating a target function

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

Second attempt: (shallow) feed-forward neural network

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

Popular choices for the non-linear function:

this is somewhat special...



Approximating a target function

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

Second attempt: (shallow) feed-forward neural network

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



NOTE: <u>biases</u> $oldsymbol{b}$ and $oldsymbol{b}$ are NOT represented in the graph

Approximating a target function

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

Second attempt: (shallow) feed-forward neural network

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



NOTE: <u>biases</u> $oldsymbol{b}$ and $oldsymbol{b}$ 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

$$y=f^*(oldsymbol{x}), \hspace{0.2cm} oldsymbol{x}\in \mathbb{R}^d$$
 (which is

(which is continuous and Borel measurable)

and any $\ \varepsilon > 0$ there exists parameters

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

this is the dimension of the hidden layer: it is a <u>parameter</u> in the theorem

such that the (shallow) feed-forward neural network

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

approximates the target function by less than $\, arepsilon \,$

$$\sup_{oldsymbol{x}} |f^*(oldsymbol{x}) - (oldsymbol{w} \cdot g(oldsymbol{W} oldsymbol{x} + oldsymbol{b}) + b)| < \varepsilon$$
 (on any compact subset of \mathbb{R}^d)

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

Deep Learning 2024-2025

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 now that x is *scalar* (hence W is *vector*)

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

Deep Learning 2024-2025

Artificial Neural Networks [26]

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;

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

then, by expanding the scalar product:

$$\tilde{y} = w_1 \operatorname{step}(W_1 x + b_1) + \dots + w_h \operatorname{step}(W_h x + b_h) + b$$

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

$$-rac{b_i}{W_i}<-rac{b_j}{W_j}, \ \ W_i, W_j>0, \ \ w_i=-w_j$$
 in this way we can construct $rac{h}{2}$ such function steps



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

----- Step

Artificial Neural Networks [27]

Deep Learning 2024-2025

Learning with Feed-Forward Neural Networks

Deep Learning 2024–2025

Artificial Neural Networks [28]

Learning with FF Neural Networks

Approximating a target function

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

Second attempt: (shallow) feed-forward neural network

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

Optimization problem (learning)

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

/ the dimension of the hidden layer is pre-defined

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

that minimize the loss function
$$\,L(D):=rac{1}{N}\sum_{D}{(ilde{y}^{(i)}-y^{(i)})^2}$$

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

Deep Learning 2024–2025

Learning with FF Neural Networks

Approximating a target function

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

Second attempt: (shallow) feed-forward neural network

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

Difficulty

In general, *minimizing* the loss function

$$L(D) = \frac{1}{N} \sum_{D} ((\boldsymbol{w} \cdot g(\boldsymbol{W}\boldsymbol{x}^{(i)} + \boldsymbol{b}) + b) - y^{(i)})^2$$
this loss function is not convex, in general

cannot be done directly, since

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

cannot be solved analytically

We need to find another way...

Deep Learning 2024-2025

Artificial Neural Networks [30]

Gradient Descent (GD): intuition

Optimization problem

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

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

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

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

 $\eta \ll 1$

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

Deep Learning 2024-2025

Artificial Neural Networks [32]

Gradient Descent (GD): convergence

Convergence

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

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

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

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

Gradient Descent (GD): practicalities

Convergence in practice



The choice of the *learning rate* η is crucial

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

Deep Learning 2024-2025

Artificial Neural Networks [34]

Gradient Descent (GD): practicalities

Convergence in practice

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



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

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

Deep Learning 2024-2025

Artificial Neural Networks [35]

Learning Feed-Forward Neural Networks (contd.)

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

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

Moral: we need to compute the gradient on each data item

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

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

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

then we can apply a *gradient descent* method

Gradient Descent

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

$$\Delta \boldsymbol{W} = -\eta \frac{1}{N} \sum_{D} \frac{\partial}{\partial \boldsymbol{W}} L(\tilde{y}^{(i)}, y^{(i)}) \qquad \Delta \boldsymbol{b} = -\eta \frac{1}{N} \sum_{D} \frac{\partial}{\partial \boldsymbol{b}} L(\tilde{y}^{(i)}, y^{(i)})$$
$$\Delta \boldsymbol{w} = -\eta \frac{1}{N} \sum_{D} \frac{\partial}{\partial \boldsymbol{w}} L(\tilde{y}^{(i)}, y^{(i)}) \qquad \Delta \boldsymbol{b} = -\eta \frac{1}{N} \sum_{D} \frac{\partial}{\partial \boldsymbol{b}} L(\tilde{y}^{(i)}, y^{(i)})$$

3. Unless complete, return to step 2.

Deep Learning 2024–2025

Computing Gradients

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

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

(similar expressions hold for the other three gradients)

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



Assume

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

i.e., the non-linearity is ReLU Easy, huh?

Deep Learning 2024–2025

Function Approximation: FF Neural Networks

Loss minimization

Approximator: (shallow) feed-forward neural network

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

XOR x_1 x_2 $x_1 \oplus x_2$ 00001101101

Optimal values for XOR and h = 2

dimension of the hidden layer

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



Artificial Neural Networks [40]

Deep Learning 2024-2025

Stochastic and Mini-Batch Gradient Descent

Function Approximation: FF Neural Networks

Loss minimization

Approximator: (shallow) feed-forward neural network

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

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

What if the dataset was <u>very</u> 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 $(\pmb{x}^{(i)}, y^{(i)}) \in D$ 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 $\,oldsymbol{W}^{(0)},\,oldsymbol{b}^{(0)},\,oldsymbol{w}^{(0)},\,b^{(0)}$
 - 2. Pick up a data item $(x^{(i)}, y^{(i)})$ from D with uniform probability and update the four parameters (with $\eta \ll 1.0, \eta \rightarrow 0$ as iterations progress)

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

3. Unless complete, return to step 2.

Deep Learning 2024–2025

Artificial Neural Networks [44]

<u>Stochastic</u> Gradient Descent (SGD): convergence

Convergence

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

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

the stochastic gradient descent method converges to the optimal ϑ^* for $t \to \infty$ provided that

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

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

Speed of Convergence

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

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

Define *accuracy* as an upper bound:

optimal value current parameter estimate $|L(artheta^*) - L(ilde{artheta})| <
ho$

descent (SGD)

 $N\,$ size of the dataset

[Bottou & Bousquet, 2008] $\rac{q}{}$ r	number of (scalar) para	meters in $artheta$
Algorithm	Cost per iteration	lterations to reach accuracy $ ho$	Time to reach accuracy ρ
<i>Gradient descent</i> (GD)	$\mathcal{O}(N q)$	$\mathcal{O}\left(\log \frac{1}{\rho}\right)$	$\mathcal{O}\left(Nq\log\frac{1}{\rho}\right)$
Stochastic aradient			

 $\mathcal{O}(q)$

Deep Learning 2024–2025

Artificial Neural Networks [46]

 $\left(q\frac{1}{q}\right)$

Qualitative comparison of GD methods



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

Deep Learning 2024–2025

Artificial Neural Networks [47]

Mini-batch Gradient Descent (MBGD): intuition

Objective

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

- Iterative method
 - 1. Initialize $\theta^{(0)}$ 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.

where:

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

This method has the same convergence properties of SGD

Deep Learning 2024-2025

Mini-batch Gradient Descent for FF Neural Networks

Mini-batch Gradient Descent (MBGD)

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

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

3. Unless complete, return to step 2.

This method has the same convergence properties of SGD