

# Lecture Notes

## The world of the neuron: from the perceptron to deep architectures

*Financial Machine Learning*

April 10, 2023

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

Oualid Missaoui  
AI4Quants Land

# Agenda

The world of the  
neuron: from the  
perceptron to deep  
architectures

- 1 Introduction
- 2 Feed-forward Neural Network
- 3 Learning neural networks: loss function
  - Regression
  - Classification
  - Backpropagation and dynamic programming
- 4 Issues of the back-propagation
- 5 References

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

- Starting point: linearly separable two-class  $(\omega_1, \omega_2)$  classification task
- $(\mathbf{x}_i, y_i), i = 1, \dots, N$ , with  $y_i \in \{-1, +1\}$ ,  $\mathbf{x}_i \in \mathbb{R}^D$
- We assume there is a hyperplane

$$\boldsymbol{\theta}_*^T \mathbf{x} = 0$$

such that

$$\begin{aligned}\boldsymbol{\theta}_*^T \mathbf{x} &> 0, & \text{if } \mathbf{x} \in \omega_1, \\ \boldsymbol{\theta}_*^T \mathbf{x} &< 0, & \text{if } \mathbf{x} \in \omega_2,\end{aligned}$$

The bias term of the hyperplane has been absorbed in  $\boldsymbol{\theta}_*$ .

- Goal: developing an algorithm that iteratively computes a hyperplane that classifies correctly all the patterns from both classes.

## Introduction

### Feed-forward Neural Network

### Learning neural networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

### Issues of the back-propagation

### References

- Let the available estimate at the current iteration step of the unknown parameters be  $\theta$ . Then there are two possibilities:
  - all points are classified correctly and a solution has been obtained, or
  - $\theta$  classifies correctly some of the points and the rest are misclassified
- Let  $\mathcal{Y}$  be the set of all misclassified samples
- The perceptron cost is defined as

$$\mathcal{J}(\theta) = - \sum_{i: \mathbf{x}_i \in \mathcal{Y}} y_i \theta^T \mathbf{x}_i \quad (1)$$

$$= \left( - \sum_{i: \mathbf{x}_i \in \mathcal{Y}} y_i \mathbf{x}_i^T \right) \theta \quad (2)$$

## Introduction

### Feed-forward Neural Network

### Learning neural networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

### Issues of the back-propagation

### References

# Perceptron Algorithm

- Starting from an arbitrary point  $\theta^{(0)}$ , the following iterative (batch) update

$$\theta^{(t)} = \theta^{(t-1)} + \mu_t \sum_{i: \mathbf{x}_i \in \mathcal{Y}} y_i \mathbf{x}_i \quad (3)$$

converges after a finite number of steps, where  $\mu_t$  is the user-defined step size.

- Stochastic version: considering one sample per iteration in a cyclic fashion

$$\theta^{(t)} = \begin{cases} \theta^{(t-1)} + \mu_t y_i \mathbf{x}_i, & \text{if } \mathbf{x}_i \text{ is misclassified by } \theta^{(t-1)}; \\ \theta^{(t-1)}, & \text{otherwise.} \end{cases}$$

Once all samples have been considered, we say that one epoch has been completed. If no convergence has been attained, all samples are reconsidered in a second epoch and so on.

## Introduction

### Feed-forward Neural Network

### Learning neural networks: loss function

Regression

Classification

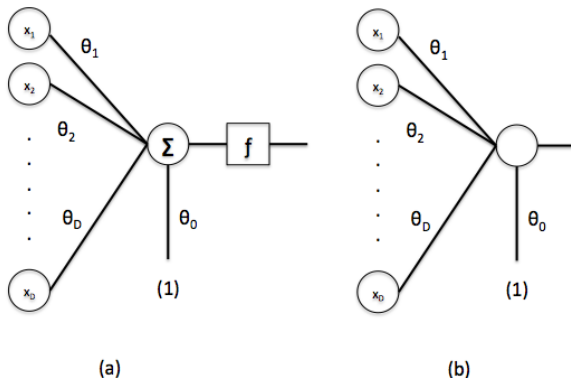
Backpropagation and dynamic programming

### Issues of the back-propagation

### References

# McCulloch-Pitts Neuron

The world of the neuron: from the perceptron to deep architectures



**Figure:** Schematic of a single hidden layer, feed-forward neural network (also called two layer feed-forward neural network).

The activation function is the Heaviside one,

$$f(z) = \begin{cases} 1 & \text{if } z > 0, \\ 0 & \text{if } z \leq 0. \end{cases}$$

Introduction

Feed-forward Neural Network

Learning neural networks: loss function

Regression

Classification

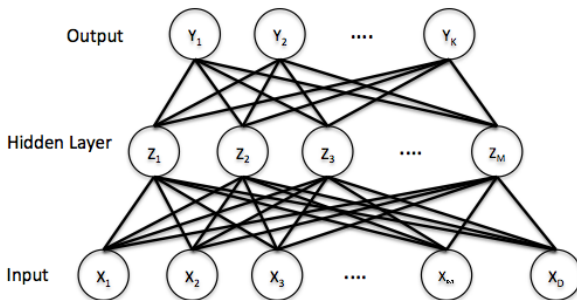
Backpropagation and dynamic programming

Issues of the back-propagation

References

# Feed-forward Neural Network

The world of the  
neuron: from the  
perceptron to deep  
architectures



**Figure:** Schematic of a single hidden layer, feed-forward neural network (also called two layer feed-forward neural network).

[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

[Regression](#)

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)

- Let  $(\mathbf{x}_i, y_i), i = 1, \dots, N$ , be the set of training samples.
- We assume that the network comprises of  $L$  layers;  $L - 1$  hidden and one output layers.
- Each layer consists of  $k_r, r = 1, \dots, L$  neurons, thus

$$\mathbf{y}_i = [y_{i1}, \dots, y_{ik_L}]^T \in \mathbb{R}^{k_L}$$

- Let  $\theta_j^r$  denote the synaptic weights associated with the  $j$ th neuron and the  $r$ th layer, with  $j = 1, \dots, k_r$  and  $r = 1, 2, \dots, L$ , where the bias term is included in  $\theta_j^r$ , that is,

$$\theta_j^r = [\theta_{j0}^r, \theta_{j1}^r, \dots, \theta_{jk_{r-1}}^r]^T$$

- the synaptic weights link the respective neuron to all neurons in layer  $k_{r-1}$

[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

[Regression](#)

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)



A standard network would be defined as:

$$\mathbf{z}_{ij}^r = \boldsymbol{\theta}_j^{rT} \mathbf{y}_i^{r-1} \quad \text{Pre-activation}$$

$$\mathbf{y}_i^r = \sigma(\mathbf{z}_i^r) \quad \text{Post-activation}$$

Unrolling the function composition

$$\hat{f}(x) = \sigma(\theta \sigma(\theta \cdots \sigma(\theta x) \cdots))$$

$\sigma$  is the activation function

- for regression: the identity function
- for classification: the softmax function

$$\frac{e^{T_k}}{\sum_{l=1}^K e^{T_l}}$$

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

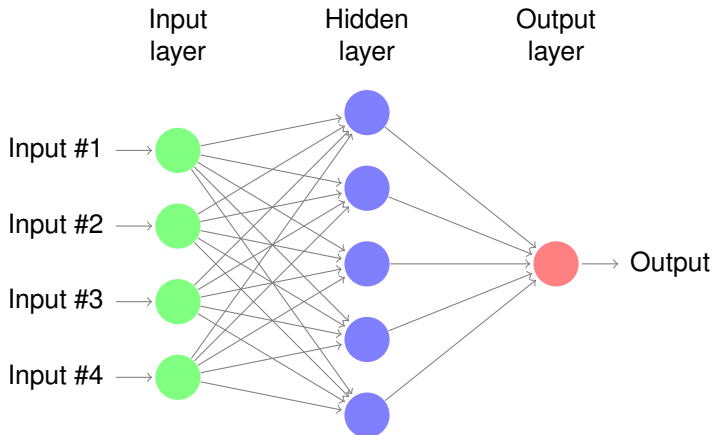
Issues of the  
back-propagation

References

- We will refer to the number of layers in a network as the *depth* of the network
- Networks with up to three (two hidden) layers, are known as *shallow*
- Those with more than three are called *deep* networks
- Network with more than 10 hidden layers are considered very deep neural networks.

# Feed-forward Neural Network

The world of the  
neuron: from the  
perceptron to deep  
architectures



[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

[Regression](#)

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)

# Activation functions

- step-function

$$\sigma(x) = \begin{cases} 1 & \text{if } z > 0, \\ 0 & \text{if } z \leq 0. \end{cases}$$

- Sigmoid (Fermi function)

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

- Hyperbolic tangent

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

- Rectified Linear unit (ReLU)

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

- Leaky rectified linear unit (Leaky ReLU)

$$\sigma(x) = \begin{cases} ax & \text{if } x \leq 0, \\ x & \text{otherwise.} \end{cases}$$

- Exponential linear units

$$\sigma(x) = \begin{cases} a \cdot (e^x - 1) & \text{if } x \leq 0, \\ x & \text{otherwise.} \end{cases}$$

# Illustration of Activation functions

The world of the  
neuron: from the  
perceptron to deep  
architectures

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

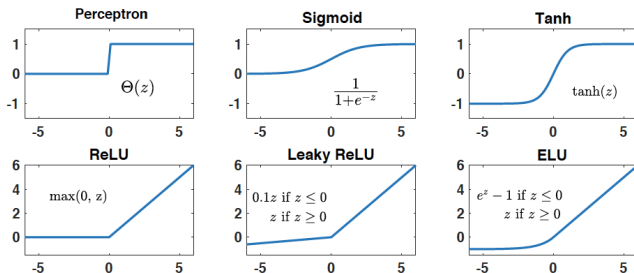
Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References



**Figure:** Possible non-linear activation functions for neurons.

Source: Pankaj Mehta, Ching-Hao Wang, Alexandre G. R. Day, and Clint Richardson *A high-bias, low-variance introduction to Machine Learning for physicists*

If the square error loss is adopted, we have

$$\mathcal{J}(\theta) = \sum_{i=1}^N \mathcal{J}_i(\theta)$$

, where

$$\mathcal{J}_i(\theta) = \frac{1}{2} \sum_{k=1}^{k_L} (\hat{y}_{ik} - y_{ik})^2$$

[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

**[Regression](#)**

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)

For classification, as is the case for logistic regression, the cost/loss function can be defined as the log-likelihood function

$$\begin{aligned}\mathcal{J}(\boldsymbol{\theta}) &= \log P(\mathcal{D}|\boldsymbol{\theta}) \\&= \log \prod_{i=1}^N P(\mathbf{x}_i|\boldsymbol{\theta}) \\&= \log \prod_{i=1}^N \sigma((\boldsymbol{\theta}^L)^T \mathbf{y}_i^{L-1})^{y_i} (1 - \sigma((\boldsymbol{\theta}^L)^T \mathbf{y}_i^{L-1}))^{1-y_i} \\&= \sum_{i=1}^N y_i \log (\sigma((\boldsymbol{\theta}^L)^T \mathbf{y}_i^{L-1})) + (1 - y_i) \log (1 - \sigma((\boldsymbol{\theta}^L)^T \mathbf{y}_i^{L-1}))\end{aligned}$$

The right hand side of the previous equation is known in statistics as the cross-entropy.

[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

[Regression](#)

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)

## Formal setup

- Gradient descent scheme:

$$\theta_j^r(t+1) = \theta_j^r(t) + \Delta\theta_j^r$$

where

$$\Delta\theta_j^r = -\mu \frac{\partial \mathcal{J}}{\partial \theta_j^r} \Big|_{\theta_j^r(t)}$$

- If the square error loss is adopted, we have

$$\mathcal{J}(\theta) = \sum_{i=1}^N \mathcal{J}_i(\theta)$$

, where

$$\mathcal{J}_i(\theta) = \frac{1}{2} \sum_{k=1}^{k_L} (\hat{y}_{ik} - y_{ik})^2$$

- For the neurons at the layer  $r$ , we denote its output by  $\mathbf{y}_i^r = [1, y_{i1}^r, \dots, y_{ik_r}^r]^T$ 
  - at the output layer,  $r = L$ ,  $y_{im}^L = \hat{y}_{im}$ ,  $m = 1, 2, \dots, k_L$
  - for  $r = 1$ , we have  $y_{im}^0 = x_{im}$ ,  $m = 1, 2, \dots, k_0$ ; that is,  $y_{im}^0$  are set equal to the input feature values.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References



## Gradient of loss function

- Let  $z_{ij}^r$  denote the output of the linear combiner of the  $j$ th neuron in the  $r$ th layer

$$z_{ij}^r = \sum_{m=0}^{k_{r-1}} \theta_{jm}^r y_{im}^{r-1} = \boldsymbol{\theta}_j^{rT} \mathbf{y}_i^{r-1}$$

and

$$\mathbf{z}_i^r = [z_{i1}^r, z_{i2}^r, \dots, z_{iK_r}^r]^T$$

- It follows that

$$\mathbf{y}_i^r = \sigma(\mathbf{z}_i^r)$$

- We can now write that:

$$\frac{\partial \mathcal{J}_i}{\partial \boldsymbol{\theta}_j^r} = \frac{\partial \mathcal{J}_i}{\partial z_{ij}^r} \frac{\partial z_{ij}^r}{\partial \boldsymbol{\theta}_j^r} = \frac{\partial \mathcal{J}_i}{\partial z_{ij}^r} \mathbf{y}_i^{r-1}$$

- Let us now define

$$\delta_{ij}^r = \frac{\partial \mathcal{J}_i}{\partial z_{ij}^r}$$

- Then we have

$$\Delta \boldsymbol{\theta}_j^r = -\mu \sum_i^N \delta_{ij}^r \mathbf{y}_i^{r-1}, r = 1, 2, \dots, L$$

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

## Computation of $\delta_{ij}^r$

Here is where the heart of the back-propagation algorithm beats.  
For the computation of the gradients,  $\delta_{ij}^r$ , one starts at the last layer  $r = L$ , and proceeds backwards toward  $r = 1$  which justifies the name given to the algorithm.

- $r = L$ , we have that  $\delta_{ij}^L = \frac{\partial \mathcal{J}_i}{\partial z_{ij}^L}$ .

- For the squared error loss function,

$$J_n = \frac{1}{2} \sum_{k=1}^{k_L} \left( f(z_{ik}^L) - y_{ik} \right)^2$$

- Hence,  $\delta_{ij}^L = (\hat{y}_{ij} - y_{ij}) f'(z_{ij}^L)$

- $r < L$ : due to the successive dependence between the layers, the value of  $z_{ij}^{r-1}$  influences all the values  $z_{ik}^r, k = 1, 2, \dots, k_r$  of the next layer.

Employing the chain rule for differentiation, we get

- $\delta_{ij}^{r-1} = \frac{\partial \mathcal{J}_i}{\partial z_{ij}^{r-1}} = \sum_{k=1}^{k_r} \frac{\partial \mathcal{J}_i}{\partial z_{ik}^r} \frac{\partial z_{ik}^r}{\partial z_{ij}^{r-1}} = \sum_{k=1}^{k_r} \delta_{ik}^r \frac{\partial z_{ik}^r}{\partial z_{ij}^{r-1}}$

- $\frac{\partial z_{ik}^r}{\partial z_{ij}^{r-1}} = \frac{\partial \left( \sum_{m=1}^{k_{r-1}} \theta_{km}^r y_{im}^{r-1} \right)}{\partial z_{ij}^{r-1}} = \theta_{kj}^r f'(z_{ij}^{r-1})$  since  $y_{im}^{r-1} = f(z_{ij}^{r-1})$

- Thus

$$\delta_{ij}^{r-1} = \left( \sum_{k=1}^{k_r} \delta_{ik}^r \theta_{kj}^r \right) f'(z_{ij}^{r-1})$$

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

## Alternative representation of $\delta_i^r$

We can clearly see that  $\delta_{ij}^{r-1}$  can be rewritten in a matrix format:

$$\delta_{ij}^{r-1} = \left( \sum_{k=1}^{k_r} \delta_{ik}^r \theta_{kj}^r \right) f'(z_{ij}^{r-1})$$

where

$$\delta_i^{r-1} = \Sigma'(\mathbf{z}_i^{r-1})(\boldsymbol{\theta}^r)^T \delta_i^r$$

where  $\Sigma'(\mathbf{z}_i^{r-1})$  is a square matrix whose diagonal entries are the values  $\sigma'(z_{ij}^{r-1})$ , and whose off-diagonal entries are zero.

Similarly, since  $\delta_{ij}^L = (\hat{y}_{ij} - y_{ij}) f'(z_{ij}^L)$ , we can rewrite the error at the final layer as:

$$\delta_i^L = \Sigma'(\mathbf{z}_i^L) \nabla_y \mathcal{J}$$

Combining the above, we get

$$\delta_i^r = \Sigma'(\mathbf{z}_i^r)(\boldsymbol{\theta}^{r+1})^T \cdots \Sigma'(\mathbf{z}_i^{L-1})(\boldsymbol{\theta}^L)^T \Sigma'(\mathbf{z}_i^L) \nabla_y \mathcal{J}$$

[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

[Regression](#)

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)

# Naive/Brute force computation of the gradient

The world of the  
neuron: from the  
perceptron to deep  
architectures

$$\delta_i^0 = \Sigma'(z_i^0)(\theta^1)^T \Sigma'(z_i^1)(\theta^2)^T \dots \Sigma'(z_i^{L-1})(\theta^L)^T \Sigma'(z_i^L) \nabla_y \mathcal{J}$$

$$\delta_i^1 = \Sigma'(z_i^1)(\theta^2)^T \Sigma'(z_i^2)(\theta^3)^T \dots \Sigma'(z_i^{L-1})(\theta^L)^T \Sigma'(z_i^L) \nabla_y \mathcal{J}$$

$\vdots$

$$\delta_i^r = \Sigma'(z_i^r)(\theta^{r+1})^T \Sigma'(z_i^{r+1})(\theta^{r+2})^T \dots \Sigma'(z_i^{L-1})(\theta^L)^T \Sigma'(z_i^L) \nabla_y \mathcal{J}$$

$$\delta_i^{r+1} = \Sigma'(z_i^{r+1})(\theta^{r+2})^T \Sigma'(z_i^{r+2})(\theta^{r+3})^T \dots \Sigma'(z_i^{L-1})(\theta^L)^T \Sigma'(z_i^L) \nabla_y \mathcal{J}$$

$\vdots$

$$\delta_i^{L-1} = \Sigma'(z_i^{L-1})(\theta^L)^T \Sigma'(z_i^L) \nabla_y \mathcal{J}$$

$$\delta_i^L = \Sigma'(z_i^L) \nabla_y \mathcal{J}$$

Clearly, computing the gradients independently would have been very costly and inefficient due to the redundant operations.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Back-propagation algorithm

The world of the  
neuron: from the  
perceptron to deep  
architectures

The following back-propagation algorithm computes the gradient of the loss function with respect of all the network parameters efficiently

**Input** : Set the corresponding activation for the input layer

**Feedforward** : For each hidden layer  $2 \leq r \leq L$ , compute  $\mathbf{z}_i^r$

**Output error** : Compute the vector  $\delta_i^L = \Sigma'(\mathbf{z}_i^L) \nabla_{\mathbf{y}} \mathcal{J}$

**Backpropagate the error** : For each  $r = L - 1, L - 2, \dots, 2$ , compute  $\delta_i^{r-1} = \Sigma'(\mathbf{z}_i^{r-1})(\theta^r)^T \delta_i^r$  (hence the name bac-propagation)

**Output** : The gradient of the cost function is given by  $\delta_{ij}^r \mathbf{y}_i^{r-1}$

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Gradient descent and its generalization

The world of the neuron: from the perceptron to deep architectures

In the simplest gradient descent (GD) algorithm, we update the parameters as follows. Initialize the parameters to some value  $\theta_0$  and iteratively update the parameters to the equation

$$\begin{aligned}\mathbf{v}_t &= \eta_t \nabla_{\theta} \mathcal{J}(\theta_t) \\ \theta_{t+1} &= \theta_t - \mathbf{v}_t\end{aligned}$$

where  $\nabla_{\theta} \mathcal{J}(\theta)$  is the gradient of  $\mathcal{J}(\theta)$  w.r.t  $\theta$  and  $\eta_t$  is the learning rate that controls the magnitude of the step we would take in the direction of the gradient at time step  $t$ .

- For  $\eta_t$  sufficiently small, the method will converge to a local minimum
- The smaller the  $\eta_t$ , the more steps we have to take to reach the local minimum
- If  $\eta_t$  is too large, we can overshoot the minimum and the algorithm becomes unstable (it either oscillates or even move away from the minimum)

Introduction

Feed-forward Neural Network

Learning neural networks: loss function

Regression

Classification

Backpropagation and dynamic programming

Issues of the back-propagation

References

# Understanding Gradient descent

The world of the neuron: from the perceptron to deep architectures

In Newton's method, we choose the step  $\mathbf{v}_t$  for the parameters in such a way as to minimize a second-order Taylor expansion to the energy function

$$\mathcal{J}(\boldsymbol{\theta} + \mathbf{v}) \approx \mathcal{J}(\boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) \mathbf{v} + \frac{1}{2} \mathbf{v}^T \mathbf{H}(\boldsymbol{\theta}) \mathbf{v}$$

where  $\mathbf{H}(\boldsymbol{\theta})$  is the Hessian matrix of the second derivatives. Ideally, using the second order approximation, we can find  $\mathbf{v}$  that makes  $\mathcal{J}(\boldsymbol{\theta} + \mathbf{v})$  minimal, which translates to minimizing

$$\mathcal{J}(\boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) \mathbf{v} + \frac{1}{2} \mathbf{v}^T \mathbf{H}(\boldsymbol{\theta}) \mathbf{v}$$

with respect to  $\mathbf{v}$ , which leads to

$$\mathbf{v}_{\text{opt}} = \mathbf{H}(\boldsymbol{\theta})^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta})$$

and hence Newton's method

$$\begin{aligned} \mathbf{v}_t &= \mathbf{H}(\boldsymbol{\theta})^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}_t) \\ \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \mathbf{v}_t \end{aligned}$$

Introduction

Feed-forward Neural Network

Learning neural networks: loss function

Regression

Classification

Backpropagation and dynamic programming

Issues of the back-propagation

References

Considering the one dimensional case and a quadratic loss function, the second order expansion becomes

$$\mathcal{J}(\theta + v) \approx \mathcal{J}(\theta) + \partial_{\theta} \mathcal{J}(\theta)v + \frac{1}{2} \partial_{\theta}^2 \mathcal{J}(\theta)v^2$$

From a gradient descent perspective, the optimal learning rate  $\eta_{\text{opt}}$  can be reliably derived from the Newton's method updating step

$$v = \eta_t \partial_{\theta} \mathcal{J}(\theta) = [\partial_{\theta}^2 \mathcal{J}(\theta)]^{-1} \partial_{\theta} \mathcal{J}(\theta)$$

which gives

$$\eta_{\text{opt}} = [\partial_{\theta}^2 \mathcal{J}(\theta)]^{-1}$$

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References



It can be shown that there are four qualitatively different regimes possible

- $\eta < \eta_{\text{opt}}$ : multiple small steps to reach the solution
- $\eta = \eta_{\text{opt}}$ : one single step is required
- $\eta_{\text{opt}} < \eta < 2\eta_{\text{opt}}$ : ascillation across both sides of the potential before eventually converging to the minimum
- $\eta > 2\eta_{\text{opt}}$ : the gradient descent diverges

source:

LeCun Yann, Leon Bottou, Genevieve B Orr, and Klaus Robert Muller (1998b), *Efficient backprop*, in Neural networks: Tricks of the trade (Springer) pp. 9-50.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

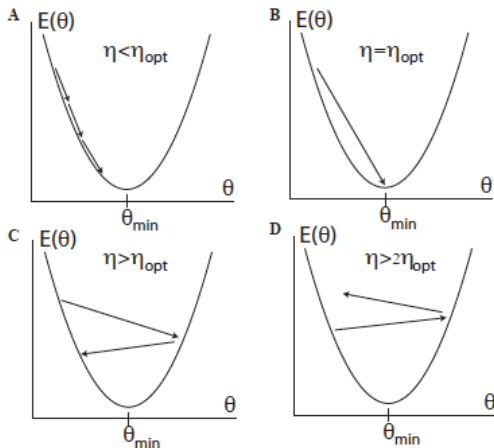
Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Understanding Gradient descent

The world of the neuron: from the perceptron to deep architectures



**Figure:** Effect of learning rate on convergence

Source: Pankaj Mehta, Ching-Hao Wang, Alexandre G. R. Day, and Clint Richardson *A high-bias, low-variance introduction to Machine Learning for physicists*

Introduction

Feed-forward Neural Network

Learning neural networks: loss function

Regression

Classification

Backpropagation and dynamic programming

Issues of the back-propagation

References

# Limitations of the simplest gradient descent algorithm

The world of the  
neuron: from the  
perceptron to deep  
architectures

- The gradient descent finds local minima of the cost function
- Gradients are computationally expensive to calculate for large datasets
- The gradient descent is very sensitive to choices of the learning rates
- The gradient descent treats all directions in parameter space uniformly
- The gradient descent is sensitive to initial conditions
- The gradient descent can take exponential time to escape saddle points, even with random initialization

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Stochastic Gradient Descent (SGD)

The world of the  
neuron: from the  
perceptron to deep  
architectures

- Traditionally, SGD is a variation of the GD where the gradient is approximated on one data point at a time. The stochasticity is incorporated through the choice of the order of the data points presented to the gradient descent.
- This can be generalized to a subset of the data instead of one data point at each iteration: we call this subset *minibatch*
- The size of the minibatches is almost always much smaller than the total number of data points
- At each gradient descent step we approximate the gradient using a single minibatch, and update the parameters  $\theta$
- if  $N$  is the total number of the data points, and  $M$  is the size of the minibatch, then the number of SGD steps/iterations is  $\frac{N}{M}$
- A full iteration over all  $N$  data points is called an *epoch*

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Gradient Descent with momentum (GDM)

The world of the neuron: from the perceptron to deep architectures

$$\begin{aligned}\mathbf{v}_t &= \gamma \mathbf{v}_{t-1} + \eta_t \nabla_{\theta} \mathcal{J}(\theta_t) \\ \theta_{t+1} &= \theta_t - \mathbf{v}_t\end{aligned}$$

where  $\gamma$  is the momentum parameter with  $0 \leq \gamma \leq 1$ .

- momentum helps the gradient descent algorithm gain speed in directions with persistent but small gradients even in the presence of stochasticity,
- while suppressing oscillations in highcurvature directions.
- This becomes especially important in situations where the landscape is shallow and flat in some directions and narrow and steep in others

Introduction

Feed-forward Neural Network

Learning neural networks: loss function

Regression

Classification

Backpropagation and dynamic programming

Issues of the back-propagation

References

# Nesterov Accelerated Gradient (NAG)

The world of the  
neuron: from the  
perceptron to deep  
architectures

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

$$\begin{aligned}\mathbf{v}_t &= \gamma \mathbf{v}_{t-1} + \eta_t \nabla_{\theta} \mathcal{J}(\boldsymbol{\theta}_t + \gamma \mathbf{v}_{t-1}) \\ \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \mathbf{v}_t\end{aligned}$$

where  $\gamma$  is the momentum parameter with  $0 \leq \gamma \leq 1$ .

One of the major advantages of NAG is that it allows for the use of a larger learning rate than GDM for the same choice of  $\gamma$ .

## Second order: RMSprop

$$\begin{aligned} \mathbf{g}_t &= \nabla_{\theta} \mathcal{J}(\theta_t) \\ \mathbf{s}_t &= \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2 \\ \theta_{t+1} &= \theta_t - \eta_t \frac{\mathbf{g}_t}{\mathbf{s}_t + \epsilon} \end{aligned}$$

where  $\beta$  controls the averaging time of the second moment and is typically taken to be about  $\beta = 0.9$ ,  $\eta_t$  is a learning rate typically chosen to be  $10^{-3}$ , and  $\epsilon \approx 10^{-8}$  is a small regularization constant to prevent divergences.

- learning rate is reduced in directions where the gradient is consistently large.
- This greatly speeds up the convergence by allowing us to use a larger learning rate for flat directions.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

## Second order: ADAM (adaptive moment estimation)

The world of the  
neuron: from the  
perceptron to deep  
architectures

$$\begin{aligned} \mathbf{g}_t &= \nabla_{\theta} \mathcal{J}(\theta_t) \\ \mathbf{m}_t &= \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \\ \mathbf{s}_t &= \beta_2 \mathbf{s}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2 \\ \hat{\mathbf{m}}_t &= \frac{\mathbf{m}_t}{1 - \beta_1^t} \\ \hat{\mathbf{s}}_t &= \frac{\mathbf{s}_t}{1 - \beta_2^t} \end{aligned}$$

$$\sigma_t^2 = \hat{\mathbf{s}}_t - (\hat{\mathbf{m}}_t)^2 \theta_{t+1} = \theta_t - \eta_t \frac{\hat{\mathbf{m}}_t}{\sqrt{\sigma_t^2 + \hat{\mathbf{m}}_t^2 + \epsilon}}$$

where  $\beta_1$  and  $\beta_2$  set the memory lifetime of the first and second moment and are typically taken to be 0.9 and 0.99

- adapting our step size so that we cut off large gradient directions (and hence prevent oscillations and divergences)
- measuring gradients in terms of a natural length scale, the standard deviation

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References



# Remarks on the back-propagation

- Stopping criterion:
  - value of the cost function gets smaller than a preselected threshold
  - the gradient values become small, i.e. the values of the weights do not change much from iteration to iteration
- the choice of the step-size (or the learning rate) is very critical
  - it has to be small to guarantee convergence
  - but not too small; otherwise convergence speed slows down
- multiple local minima where the algorithm can be trapped due to the high non-linearity of the neural network problem: reinitialization
- Stochastic versus batch learning scheme

[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

[Regression](#)

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)

# Universal approximation property of feed-forward neural networks

The world of the neuron: from the perceptron to deep architectures

## Theorem

Let  $g(x)$  be a continuous function defined in a compact subset  $S \subset \mathbb{R}^I$  and any  $\epsilon > 0$ . Then there is a two layer neural network  $\hat{g}(x)$  with  $K(\epsilon)$  hidden nodes so that

$$|g(x) - \hat{g}(x)| < \epsilon \quad \forall x \in S$$

## Proof.

The proof can be found in:

George Cybenko. *Approximation by superpositions of a sigmoidal function*. Mathematics of control, signals and systems, 2(4):303-314, 1989.



Introduction

Feed-forward Neural Network

Learning neural networks: loss function

Regression

Classification

Backpropagation and dynamic programming

Issues of the back-propagation

References

- A deep architecture can represent certain functions (exponentially) more compactly
- Example: Boolean functions
  - a Boolean circuit is a sort of feed-forward network where hidden units are logic gates (i.e. AND, OR or NOT functions of their arguments)
  - any Boolean function can be represented by a *single hidden layer* Boolean circuit
    - however, it might require an exponential number of hidden units
  - it can be shown that there are Boolean functions which
    - require an exponential number of hidden units in the single layer case
    - require a polynomial number of hidden units if we can adapt the number of layers

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Issues of the back-propagation

The world of the  
neuron: from the  
perceptron to deep  
architectures

- Initialization
- Overfitting
- Scaling of the inputs
- Number of hidden units and layers
- Multiple Minima

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

- Usually, starting values for weights are chosen to be random values near zero.
- Hence the model starts out nearly linear, and becomes nonlinear as the weights increase.
- Use of exact zero weights leads to zero derivatives and the algorithm never moves.
- Starting instead with large weights often leads to poor solutions

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

- Often neural networks have too many weights and will overfit the data.
- An early stopping rule can be used to avoid overfitting
  - A validation dataset is useful for determining when to stop, since we expect the validation error to start increasing.
  - Since the weights start at a highly regularized (linear) solution, this has the effect of shrinking the final model toward a linear model.
- a more explicit method for regularization is weight decay, which consists of adding a penalty to the loss function

$$\mathcal{J}(\theta) + \lambda\Omega(\theta)$$

where

$$\Omega(\theta) = \sum_{rjk} (\theta_{jk}^r)^2$$

- This is analogous to ridge regression

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

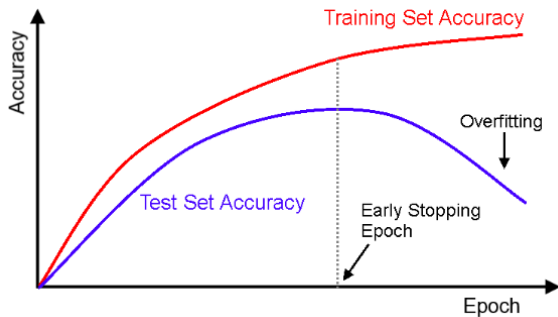
Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Overfitting

The world of the neuron: from the perceptron to deep architectures



**Figure:** Early stopping

Introduction

Feed-forward Neural Network

Learning neural networks: loss function

Regression

Classification

Backpropagation and dynamic programming

Issues of the back-propagation

References

- The loss function is nonconvex, possessing many local minima.
- As a result, the final solution obtained is quite dependent on the choice of starting weights.
- One must at least try a number of random starting configurations and choose the solution giving lowest (penalized) error.
- Probably a better approach is to use the average predictions over the collection of networks as the final prediction.
- This is preferable to averaging the weights, since the nonlinearity of the model implies that this averaged solution could be quite poor.
- Another approach is via bagging, which averages the predictions of networks training from bootstraps of the training data.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References



## Number of hidden units and layers

- Generally speaking it is better to have too many hidden units than too few.
- With too few hidden units, the model might not have enough flexibility to capture the non linearities in the data
- with too many hidden units, the extra weights can be shrunk toward zero if appropriate regularization is used
- Typically the number of hidden units is somewhere in the range of 5 to 100, with the number increasing with the number of inputs and number of training cases.
- It is most common to put down a reasonably large number of units and train them with regularization.
- Some researchers use cross-validation to estimate the optimal number
- Choice of the number of hidden layers is guided by background knowledge and experimentation
- Each layer extracts features of the input
- Use of multiple hidden layers allows construction of hierarchical features at different level of resolution.

Why training is hard:

- First hypothesis: optimization is harder (underfitting)
  - vanishing gradient problem
- Second hypothesis: overfitting
  - we are exploring a space of complex functions
  - deep nets usually have lots of parameters
  - Might be in a high variance/low bias situation

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

- Depending on the problem, one or the other situation will tend to dominate
- If first hypothesis (underfitting): better optimize
  - use better optimization methods
  - use non saturating activation functions
- If second hypothesis (overfitting): use better regularization
  - Implicit regularization using SGD
  - Dropout
  - Batch Normalization

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# The problem of vanishing or exploding gradients

The world of the  
neuron: from the  
perceptron to deep  
architectures

- This problem is especially pronounced in neural networks that try to capture long-range dependencies, such as Recurrent Neural Networks for sequential data.
- We can illustrate this problem by considering a simple network with one neuron in each layer
- We further assume that all weights are equal, and denote them by  $\theta$

$$\delta_j^1 = \delta_j^L \prod_{j=0}^{L-1} \theta \sigma'(z_j) = \delta_j^L (\theta)^L \prod_{j=0}^{L-1} \sigma'(z_j)$$

where  $\delta_j^L$  is the error in the  $L$ -th topmost layer, and  $(\theta)^L$  is the weight to the power  $L$ .

- Let us now also assume that the magnitude  $\sigma'(z_j)$  is fairly constant and we can approximate  $\sigma'(z_j) \approx \sigma'_0$

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# The problem of vanishing or exploding gradients (cont.)

The world of the  
neuron: from the  
perceptron to deep  
architectures

Under these simplifying assumptions:

$$\delta_j^1 = \delta_j^L (\theta \sigma'_0)^L$$

for large  $L$ , the error  $\delta_j^1$  has very different behavior depending on the value of  $\theta \sigma'_0$

- if  $\theta \sigma'_0 > 1$ , the errors and the gradient blow up
- if  $\theta \sigma'_0 < 1$ , the errors and the gradient vanish
- only when the weights satisfy  $\theta \sigma'_0 \approx 1$  and the neurons are not saturated will the gradient stay well behaved for deep networks
- if any activation unit saturates in either direction, the errors and gradients vanish

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

- SGD acts as an implicit regularizer by introducing stochasticity (from the use of mini-batches) that prevents overfitting.
- it is important that the weight initialization is chosen randomly,
- another important thing is to appropriately choose the learning rate or step-size by searching over five logarithmic grid points. If the best performance occurs at the edge of the grid, repeat this procedure until the optimal learning rate is in the middle of the grid parameters.
- it is common to center or whiten the input data
- Another important form of regularization that is often employed in practice is Early Stopping.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

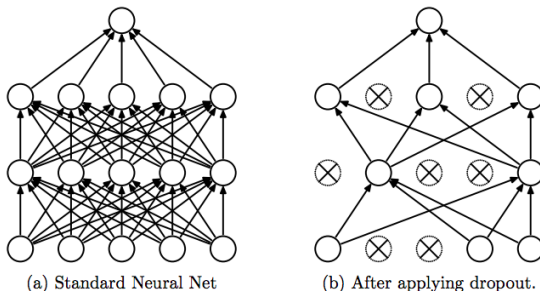
Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Drop-out

- "Cripple" neural network by removing hidden units stochastically
  - each hidden unit is set to 0 with probability  $p$  ( e.g. 0.5)



**Figure:** Drop out illustration

Source: Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, Ruslan Salakhutdinov. 2014. *Dropout: A simple way to prevent neural networks from overfitting*

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

- "Cripple" neural network by removing hidden units stochastically
- This averaging of weights is similar in spirit to the Bagging procedure

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

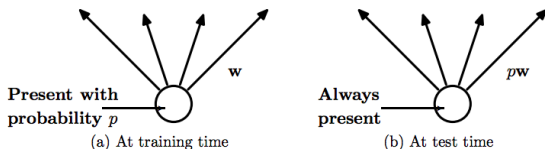
Issues of the  
back-propagation

References



# Drop-out during testing

The world of the  
neuron: from the  
perceptron to deep  
architectures



**Figure:** Drop out during testing

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

## Drop-out formal setup

This section describes the dropout neural network model. Consider a neural network with  $L$  hidden layers. Let  $r \in \{1, \dots, L\}$  index the hidden layers of the network. A standard network would be defined as:

$$\begin{aligned}z_{ij}^r &= \theta_j^{rT} \mathbf{y}_i^{r-1} \\ \mathbf{y}_i^r &= \sigma(\mathbf{z}_i^r)\end{aligned}$$

With drop-out, the feed-forward operation becomes

$$\begin{aligned}q_j^r &\sim \text{Bernoulli}(p) \\ \tilde{\mathbf{y}}_i^{r-1} &= \mathbf{q}^r * \mathbf{y}_i^{r-1} \\ z_{ij}^r &= \theta_j^{rT} \tilde{\mathbf{y}}_i^{r-1} \\ \mathbf{y}_i^r &= \sigma(\mathbf{z}_i^r)\end{aligned}$$

where  $*$  denotes an element-wise product.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

Dropout Decreases the Reliance of Neural Networks on Spurious Interaction Effects.

Source: B. Lengerich, E. P. Xing, R. Caruana, *On Dropout, Overfitting, and Interaction Effects in Deep Neural Networks*, arXiv 2007.00823, July 3, 2020

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

- The basic inspiration behind Batch Normalization is the long-known observation that training in neural networks works best when the inputs are centered around zero with respect to the bias.
- The reason for this is that it prevents neurons from saturating and gradients from vanishing in deep nets.
- In the absence of such centering, changes in parameters in lower layers can give rise to saturation effects in higher layers, and vanishing gradients.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References

# Batch Normalization

The world of the neuron: from the perceptron to deep architectures

The idea of Batch Normalization is to introduce additional new *BatchNorm* layers that standardize the inputs by the mean and variance of the mini-batch. Consider a layer  $r$  with  $K_r$  neurons whose inputs are  $(\mathbf{z}_1^r, \dots, \mathbf{z}_{K_r}^r)$ . We standardize each dimension so that

$$\mathbf{z}_k^r \longrightarrow \hat{\mathbf{z}}_k^r = \frac{\mathbf{z}_k^r - \mathbb{E}[\mathbf{z}_k^r]}{\sqrt{\text{Var}[\mathbf{z}_k^r]}}$$

where the mean and variance are taken over all samples in the mini-batch.

To avoid unwanted changes in the representational power of the neural network, we introduce new parameters that can additionally shift and scale the normalized input

$$\hat{\mathbf{z}}_k^r \longrightarrow \tilde{\mathbf{z}}_k^r = \gamma_k^r \hat{\mathbf{z}}_k^r + \beta_k^r$$

[Introduction](#)

[Feed-forward Neural Network](#)

[Learning neural networks: loss function](#)

[Regression](#)

[Classification](#)

[Backpropagation and dynamic programming](#)

[Issues of the back-propagation](#)

[References](#)

- ❶ Michael Nielsen, *Neural Networks and Deep Learning*, Online, 2019
- ❷ Ian Goodfellow, Yoshua Bengio, and Aaron Courville, *Deep Learning*, The MIT Press, 2016
- ❸ Charu C. Aggarwal *Neural Networks and Deep Learning*, Springer, 2018
- ❹ Martin Anthony and Peter L. Bartlett *Neural Network Learning: Theoretical Foundations*, Cambridge University Press, 2009
- ❺ Sergios Theodoridis, *Machine Learning: A Bayesian and optimization perspective*, Elsevier, 2015.
- ❻ Christopher M. Bishop, *Pattern Recognition and Machine Learning*, 4Myeloma Press, 2010.

Introduction

Feed-forward Neural  
Network

Learning neural  
networks: loss function

Regression

Classification

Backpropagation and  
dynamic programming

Issues of the  
back-propagation

References