# Neural Networks

doc. RNDr. lveta Mrázová, CSc.

DEPARTMENT OF THEORETICAL COMPUTER SCIENCE AND MATHEMATICAL LOGIC FACULTY OF MATHEMATICS AND PHYSICS, CHARLES UNIVERSITY IN PRAGUE

# Neural Networks:

Multi-layered Neural Networks: Analysis of Their Properties

doc. RNDr. Iveta Mrázová, CSc.

DEPARTMENT OF THEORETICAL COMPUTER SCIENCE AND MATHEMATICAL LOGIC FACULTY OF MATHEMATICS AND PHYSICS, CHARLES UNIVERSITY IN PRAGUE

## Neural Networks:

#### **Contents:**

- Multi-layered Neural Networks
- Multi-layered Neural Networks: Analysis of Their Properties
- Multi-layered Neural Networks: an Application Example

#### **Contents:**

- Multi-layered Neural Networks
  - Back-Propagation Training Algorithm
  - Strategies to Speed-up the Training Process
- Multi-layered Neural Networks: Analysis of Their Properties
  - Kolmogorov's Theorem
  - Function Approximation
  - The Complexity of Learning
  - The Number of Regions in the Feature Space
  - Vapnik-Chervonenkis Dimension
- Multi-layered Neural Networks: an Application Example
  - Internal Knowledge Representation and Pruning
  - Sensitivity Analysis and Feature Selection
  - Analysis of the World Bank Data

## Kolmogorov's Theorem – 1957

#### 13. Hilbert problem

~ Can every continuous function of *n* arguments be represented using a finite composition of functions of at most two arguments?

• Example:  $x \cdot y = \exp(\ln x + \ln y)$ 

<u>Theorem (Kolmogorov, 1957)</u>: Let  $f: [0,1]^n \rightarrow [0,1]$  be a continuous function. There exist functions of one argument g and  $\varphi_q$ , for q = 1, ..., 2n + 1 and constants  $\lambda_p$ , for p = 1, ..., n such that

$$f(x_1, \dots, x_n) = \sum_{q=1}^{2n+1} g\left(\sum_{p=1}^n \lambda_p \varphi_q(x_p)\right)$$

## Kolmogorov Networks

 to represent continuous functions of *n* variables

 $f(x_1, \dots, x_n) =$ 

$$= \sum_{q=1}^{2n+1} g(\sum_{p=1}^n \lambda_p \varphi_q(x_p))$$



I. MRÁZOVÁ: NEURAL NETWORKS

## Neural Networks:

#### **Contents:**

- Multi-layered Neural Networks
- Multi-layered Neural Networks: Analysis of Their Properties
- Multi-layered Neural Networks: an Application Example

#### **Contents:**

- Multi-layered Neural Networks
  - Back-Propagation Training Algorithm
  - Strategies to Speed-up the Training Process
- Multi-layered Neural Networks: Analysis of Their Properties
  - Kolmogorov's Theorem
  - Function Approximation
  - The Complexity of Learning
  - The Number of Regions in the Feature Space
  - Vapnik-Chervonenkis Dimension
- Multi-layered Neural Networks: an Application Example
  - Internal Knowledge Representation and Pruning
  - Sensitivity Analysis and Feature Selection
  - Analysis of the World Bank Data

## Function Approximation (1)

 Any continuous function can be reproduced exactly by a finite network of computing units, whereby the necessary primitive functions for each node exist

× the choice of the right transfer function!

The best possible approximation to a given function

× the choice of the right number of computing units with the considered transfer function

 $\Rightarrow$  Instead of an exact representation, try to approximate

## Function Approximation (2)

<u>Theorem</u>: Any continuous real function  $f: [0, 1] \rightarrow [0, 1]$  can be approximated using a network of threshold elements in such a way that the total approximation error E is lower than any given real number  $\varepsilon > 0$ :

$$E = \int_{0}^{1} \left| f(x) - \tilde{f}(x) \right| \, \mathrm{d}x < \varepsilon$$

where  $\tilde{f}$  denotes the network function.

### Function Approximation (3)

#### <u> Proof - idea:</u>

• approximation of fby means of  $\varphi_N$ 



Inspinovano Riemannashin integnilem

## Function Approximation (4)

#### Proof (continued):

- Divide the interval [0, 1] into N equal segments selecting the points  $x_0, x_1, \dots, x_N \in [0, 1]; x_0 = 0, x_N = 1$
- Define a function  $\varphi_N$  as it follows:

 $\varphi_N(x) = \min\{f(x'); x' \in [x_i, x_{i+1}) \text{ for } x_i \le x < x_{i+1}\}$ 

• Further, consider  $\varphi_N$  an approximation of f so that the approximation error  $E_N$  is given by:

$$E_N = \int_0^1 |f(x) - \varphi_N(x)| \, \mathrm{d}x$$

## Function Approximation (5)

#### Proof (continued):

• Since  $f(x) \ge \varphi_N(x) \quad \forall x \in [0, 1], E_N$  corresponds to

$$E_N = \int_0^1 f(x) \, \mathrm{d}x - \int_0^1 \varphi_N(x) \, \mathrm{d}x \qquad \stackrel{\sim \text{ lower Riemann sum of the function } f}{}$$

- Since continuous functions are integrable  $\rightarrow$  the lower sum of f converges in the limit  $N \rightarrow \infty$  to the integral of f in the interval [0,1]
- Thus, it holds  $E_N \to 0$  when  $N \to \infty$ , hence for any real number  $\varepsilon > 0$  there exists an M such that  $E_N < \varepsilon \ \forall N \ge M$
- The function  $\varphi_N$  is therefore the desired approximation of f.

## Function Approximation (6)

#### Proof (continued):

- The function \(\varphi\_N(x)\) can be computed by a network of threshold units (~ a neural network)
  - $\varphi_N(x)$  is a step-wise function
  - in each of the N segments of the interval [0, 1]: [x<sub>0</sub>, x<sub>1</sub>), [x<sub>1</sub>, x<sub>2</sub>), ..., [x<sub>N-1</sub>, x<sub>N</sub>], φ<sub>N</sub>(x) has the respective value α<sub>1</sub>, ..., α<sub>N</sub>



## Udélám si costationi jemnou sit, abyen na judervalu [21] pohuy! usechang isznamni body (2mény y-smén), ve lutemich hudu mit treshold a budu na ném poravisiont Function Approximation (7)

#### Proof (continued):

This network can compute the step-wise function  $\varphi_N(x)$ :

- The single input to the network is x (from the interval [0,1])
- Each pair of units with the thresholds  $x_i$  and  $x_{i+1}$  guarantees that the unit with threshold  $x_i$  will be active when  $x_i \le x < x_{i+1}$ .
- The (linear) output unit adds all outputs of the previous layer of units and produces their (weighted) sum as a result
- The unit with the threshold  $x_N + \delta$ , where  $\delta$  is a small positive number, is used to recognize the case  $x_{N-1} \le x \le x_N$ .

This network computes the function  $\varphi_N(x)$ , that approximates the function f with the desired maximum error.

#### QED

## Function Approximation (8)

#### Corollary:

The theorem is valid also for any function  $f: [0, 1] \rightarrow (0, 1)$ and networks with the sigmoidal transfer function.



### Function Approximation (9)

#### Proof:

- The image of the function *f* has been limited to the interval
  (0, 1) in order to simplify the proof
- The function *f* can be approximated using the sketched network
- The transfer function of the units with the threshold  $x_i$  is given by  $s_c(x - x_i)$ , where c controls the slope of the function

$$s_c(x-x_i) = \frac{1}{1+e^{-c(x-x_i)}}$$



## Function Approximation (10)

#### Proof (continued):

The network can approximate the function  $\varphi_N(x)$  with an approximation error lower than any desired bound (> 0)

~ threshold functions can be approximated with any desired precision by a sigmoidal function parametrized with *c*)

- The weights connecting the first layer of units to the output unit have been set in such a way that the sigmoid produces the desired values *α<sub>i</sub>* as a result
- Further it should be guaranteed that every input x produces a single 1 from the first layer to the output unit
  - → the first layer just finds out to which of the *N* segments of the interval [0, 1] the input *x* belongs

QED



## Function Approximation (11)

### The multidimensional case:

The network capable of approximating the function  $f: [0,1]^n \rightarrow (0,1)$  can be constructed using the same general idea as before in the one-dimensional case

- Extensions necessary for the two-dimensional case
  - Recognition of intervals in the *x* and *y* domains
    - 2 units left are used to test  $x_0 \le x < x_1$
    - 2 units right are used to test  $y_1 \le y < y_2$
  - The unit with the threshold 1.5 recognizes the conjunction of both conditions (for x and y)



## Function Approximation (12)



#### Extensions necessary for the 2D-case

- Recognition of intervals in the x and y domains
  - 2 units left are used to test  $x_0 \le x < x_1$
  - 2 units right are used to test  $y_1 \le y < y_2$
- The unit with the threshold 1.5 recognizes the conjunction of both conditions (for x and y)
- The "output" has the weight  $s_0^{-1}(\alpha_{12})$ , so the sigmoidal transfer function yields  $\alpha_{12}$ 
  - → this number corresponds to the desired approximation of the function f on  $[x_0, x_1) \times [y_1, y_2)$



#### Extensions necessary for the 2D-case

- Recognition of intervals in the x and y domains
  - 2 units left are used to test  $x_0 \le x < x_1$
  - 2 units right are used to test  $y_1 \le y < y_2$
- The unit with the threshold 1.5 recognizes the conjunction of both conditions (for x and y)
- The "output" has the weight  $s_0^{-1}(\alpha_{12})$ , so the sigmoidal transfer function yields  $\alpha_{12}$ 
  - → this number corresponds to the desired approximation of the function f on  $[x_0, x_1) \times [y_1, y_2)$

## Neural Networks:

#### **Contents:**

- Multi-layered Neural Networks
- Multi-layered Neural Networks: Analysis of Their Properties
- Multi-layered Neural Networks: an Application Example

#### **Contents:**

- Multi-layered Neural Networks
  - Back-Propagation Training Algorithm
  - Strategies to Speed-up the Training Process
- Multi-layered Neural Networks: Analysis of Their Properties
  - Kolmogorov's Theorem
  - Function Approximation
  - The Complexity of Learning
  - The Number of Regions in the Feature Space
  - Vapnik-Chervonenkis Dimension
- Multi-layered Neural Networks: an Application Example
  - Internal Knowledge Representation and Pruning
  - Sensitivity Analysis and Feature Selection
  - Analysis of the World Bank Data

## The Complexity of Learning

#### The satisfiability problem

D: Let V be a set of n logical variables, and let F be a logical expression in conjunctive normal form (conjunction of disjunctions of literals) which contains only variables from V.

The satisfiability problem consists in assigning truth values to the variables in V in such a way that the expression F becomes true.

**THEOREM:** The general learning problem for networks of threshold functions is NP-complete.

## The Complexity of Learning (2)

#### Proof - idea:

a network equivalent
 to 3SAT



## The Complexity of Learning (3)

#### Proof (continue):

1. 3SAT can be reduced to an instance of a learning problem for neural networks in polynomial time

A logical expression F in conjunctive normal form, which contains n variables can be transformed in polynomial time in the description of a network of the above type:

- For each variable  $x_i$  a weight  $w_i$  is defined
- The connections to the third layer are fixed according to the conjunctive normal form we are dealing with

## The Complexity of Learning (4)

### Proof (continue):

- This can be done (using a suitable coding) in polynomial time, because it holds for the number m of different possible disjunctions in a 3SAT formula that  $m \leq (2n)^3$
- If an instantiation A with logical values of the variables  $x_i$  exists, such that F becomes true, then there exist weights  $w_1, w_2, \ldots, w_n$ , that solve the learning problem

## The Complexity of Learning (5)

### Proof (continue):

- It is sufficient to set the weights w<sub>i</sub> = 1, if x<sub>i</sub> = 1; and w<sub>i</sub> = 0, if x<sub>i</sub> = 0.
  (in both cases, we thus choose w<sub>i</sub> = x<sub>i</sub>)
- Similarly in the opposite way: if there exist weights  $w_1, w_2, ..., w_n$ , that solve the learning problem, then the instantiation  $x_i = 1$  for  $w_i \ge 0.5$  and  $x_i = 0$  otherwise, is a valid instantiation that makes F true

## The Complexity of Learning (6)

### Proof (continue):

- 2. Further, we must show that the learning problem belongs to the class NP (its solution can be checked in polynomial time)
  - If the weights  $w_1, w_2, ..., w_n$  are given, then a single run of the network can be used to check if the output F is equal to 1
  - The number of computation steps is directly proportional to the number *n* of variables and to the number *m* of disjunctive clauses (which is bounded by the polynomial (2n)<sup>3</sup>)

## The Complexity of Learning (7)

### Proof (continue):

- The time required to check an instantiation is therefore bounded by a polynomial in n
- The given learning problem thus belongs to the class NP

#### QED

#### Remark:

For some special types of simple neural networks, the learning problem can be solved in polynomial time (by means of linear programming algorithms)

## Neural Networks:

#### **Contents:**

- Multi-layered Neural Networks
- Multi-layered Neural Networks: Analysis of Their Properties
- Multi-layered Neural Networks: an Application Example

#### **Contents:**

- Multi-layered Neural Networks
  - Back-Propagation Training Algorithm
  - Strategies to Speed-up the Training Process
- Multi-layered Neural Networks: Analysis of Their Properties
  - Kolmogorov's Theorem
  - Function Approximation
  - The Complexity of Learning
  - The Number of Regions in the Feature Space
  - Vapnik-Chervonenkis Dimension
- Multi-layered Neural Networks: an Application Example
  - Internal Knowledge Representation and Pruning
  - Sensitivity Analysis and Feature Selection
  - Analysis of the World Bank Data

## The XOR Problem (1)



- Input vector  $\vec{x} = (x_1, x_2)$  is a point in the two-dimensional space;  $f_j$  defines the labeling for all points  $(x_1, x_2)$ , where  $x_i \in \{0, 1\}$
- The perceptron can compute functions for which a hyperplane can separate the points labeled by 0 from the points labeled by 1

## The XOR Problem (2)

#### 16 Boolean functions of two variables



## The Number of Regions in the Feature Space (1)

 The capacity of a neuron depends on the dimension of the weight space and the number of cuts with separating hyperplanes

#### $\rightarrow$ <u>Question</u>:

How many regions are defined by m cutting hyper-planes of the dimension n - 1 in an n-dimensional space?

- we consider only hyperplanes going through the origin
- $\rightarrow$  The intersection of *l* hyperplanes;  $l \leq n$  is of dimension n l

# The Number of Regions in the Feature Space (2)

#### <u>2-dimensional case:</u>

m lines going through the origin define at most  $2 \cdot m$  different regions

<u>3 – dimensional case:</u>

each new cut increases the number of regions two times

■ <u>in general:</u>

n cuts with (n - 1)-dimensional hyperplanes in n-dimensional space define at most  $2^n$ different regions



+ 2 new regions



The Number of Regions in the Feature Space (3)

### Theorem:

Let R(m, n) denote the number of different regions defined by m separating hyperplanes of the dimension n - 1 in an n-dimensional space.

We set R(1, n) = 2 for  $n \ge 1$  and  $R(m, 0) = 0 \forall m \ge 1$ .

Then for  $n \ge 1$  and m > 1:

R(m,n) = R(m-1,n) + R(m-1,n-1)

The Number of Regions in the Feature Space (4)

### **Proof** (by induction on <u>m</u>):

1. m = 2 and n = 1: The formula is valid, because R(2,1) = R(1,1) + R(1,0) = 2 + 0 = 2

2. m = 2 and  $n \ge 2$ :  $R(2,n) = 4 \implies$  valid, because R(2,n) = R(1,n) + R(1,n-1) = 2 + 2 = 4

3. m + 1 hyperplanes of dimension n - 1 are given in the *n*-dimensional space and in general position  $(n \ge 2)$ :

• The first m hyperplanes define R(m, n) regions in the n-dimensional space

## The Number of Regions in the Feature Space (5)

#### Proof (continue):

- (m + 1)-st hyperplane intersects the first m hyperplanes in m hyperplanes of dimension n 2
- These m hyperplanes (of dimension n-2) divide the (n-1)-dimensional space into R(m, n-1) regions
- After the cut with the hyperplane (m + 1), exactly R(m, n 1) new regions have been created
- $\rightarrow$  The new number of regions is therefore:

• 
$$R(m+1,n) = R(m,n) + R(m,n-1)$$

QED

## Number of Regions in the Feature Space (6)

• A useful alternative for R(m, n):

$$R(m,n) = 2 \sum_{i=0}^{n-1} \binom{m-1}{i}$$

- X With growing *n*, the number of Boolean functions grows significantly quicker than the number of regions formed by hyperplanes in a general position
  - this number can be in general larger than the number of threshold functions over binary inputs

## Number of Regions in the Feature Space (7)

### Example:

n	# Boolean functions	# threshold functions	# regions
1	4	2	2
2	16	14	14
3	256	104	128
4	65536	1882	3882
5	$4.3 \times 10^{9}$	94572	412736

Number of Regions in the Feature Space (8)

#### **Consequences:**

- Learnability problems ~ if the number of input vectors is too high, the network might be not able to form enough regions with the given number of hidden neurons
- Generalization
  - ~ expected number of correctly classified examples
- Over-fitting
  - ~ erroneous interpolation of patterns outside of the training set
- Vapnik-Chervonenkis dimension (VC-dimension)
  - ~ finite VC-dimension  $\rightarrow$  "the class of concepts" is learnable

## Neural Networks:

#### **Contents:**

- Multi-layered Neural Networks
- Multi-layered Neural Networks: Analysis of Their Properties
- Multi-layered Neural Networks: an Application Example

#### **Contents:**

- Multi-layered Neural Networks
  - Back-Propagation Training Algorithm
  - Strategies to Speed-up the Training Process
- Multi-layered Neural Networks: Analysis of Their Properties
  - Kolmogorov's Theorem
  - Function Approximation
  - The Complexity of Learning
  - The Number of Regions in the Feature Space
  - Vapnik-Chervonenkis Dimension
- Multi-layered Neural Networks: an Application Example
  - Internal Knowledge Representation and Pruning
  - Sensitivity Analysis and Feature Selection
  - Analysis of the World Bank Data

## Vapnik-Chervonenkis Dimension (VC-dimension) (1)

#### **Definition:**

Let  $C = \{f_i\}$  be a set of functions (concept class).

The set of *m* training patterns  $\{t_k\}_{k=1,...,m}$  can be shattered by means of *C*, if for each of the  $2^m$  possible labelings of these patterns with 1/0, there exists at least one function, that satisfies this labeling.

#### **Definition:**

The VC-dimension V of a set of functions C is defined as the biggest m, for which a set of m training patterns exists that can be shattered.

## Vapnik-Chervonenkis Dimension (VC-dimension) (2)

If there exists for any m a set of m training patterns, that can be shattered by means of C, the VC-dimension of C is infinite

 $\rightarrow$  Such a problem is not " <u>LEARNABLE</u>"

- In general, the VC-dimension of a set of functions does not depend on the number of parameters
- VC-dimension impacts adequate generalization
  - The network can have many parameters, but it should have a small VCdimension → better generalization
  - High VC-dimension correlates with worse generalization

## Vapnik-Chervonenkis Dimension (VC-dimension) (3)

#### Example:

1. The VC-dimension of a set of linear indicator functions

 $Q(\vec{z}, \alpha) = \Theta\{\sum_{p=1}^{n} \alpha_p z_p + \alpha_0\}$  in the *n*-dimensional space is n + 1,

(i.e., it can shatter at most n + 1 patterns)



Vapnik-Chervonenkis Dimension (VC-dimension) (4)  $\theta(x) = \begin{cases} 1 & \text{for } x \ge 0 \\ 0 & \text{for } x < 0 \end{cases}$ 

- 2. VC-dimension of the set of functions  $f_{\alpha}(z) = \theta(\cos \alpha z), \ \alpha \in R$  is infinite
  - The points  $z_1 = 10^{-1}, ..., z_m = 10^{-m}$  can be shattered by means of the functions from this set
  - To shatter these patterns into two classes (+1/-1) given by the sequence  $\delta_1, \ldots, \delta_m$ ;  $\delta_i \in \{0, 1\}$  it is sufficient to choose the value of the parameter as

$$\alpha = \left(1 + \sum_{i=1}^{m} (1 - \delta_i) \ 10^i\right) \cdot \pi$$

## Vapnik-Chervonenkis Dimension (VC-dimension) (5)

For the set of functions  $f_{\alpha}(z) = \theta(\cos \alpha z), \ \alpha \in R$ 

• The points  $z_1 = 10^{-1}, ..., z_m = 10^{-m}$  can be shattered by this set of functions

- To shatter these patterns into two classes (+1/-1) given by the sequence  $\delta_1, \dots, \delta_m$ ;  $\delta_i \in \{0, 1\}$  it is sufficient to choose the value of the parameter  $\alpha = \left(1 + \sum_{i=1}^m (1 \delta_i) \ 10^i\right) \cdot \pi$
- e.g., for  $\delta_1 = 1$ ,  $\delta_2 = 0$ ,  $\delta_3 = 1$ ,  $\alpha = \pi (1 + 0 \cdot 10^1 + 1 \cdot 10^2 + 0 \cdot 10^3) = 101 \cdot \pi$
- $\cos \alpha z_1 = \cos 10.1 \, \pi \approx 0.9511 > 0$
- $\cos \alpha z_2 = \cos 1.01 \, \pi \approx -0.9995 < 0$
- $\cos \alpha z_3 = \cos 0.101 \, \pi \approx 0.9501 > 0$

## Vapnik-Chervonenkis Dimension (VC-dimension) (6)



when choosing a suitable coefficient  $\alpha$ , it is possible to approximate any function bounded in the interval (+1; -1) for any number *m* of selected points by *cos*  $\alpha z$ 

## Vapnik-Chervonenkis Dimension (VC-dimension) (7)

<u>The problem of "overfitting"</u> ~ the network learns also the noise



## Vapnik-Chervonenkis Dimension (VC-dimension) (8)

 For the network with W weights and N neurons and with the required limit for the generalization error ε, the number P of training patterns necessary for good generalization is:

$$P \ge \left(\frac{W}{\varepsilon}\right) \log_2\left(\frac{N}{\varepsilon}\right)$$

- A multi-layered network with 1 hidden layer cannot generalize well, if there were less than  $W/\varepsilon$  randomly chosen training patterns, i.e.,  $P \ge W/\varepsilon$ 
  - To achieve the accuracy of at least 90 %, it is necessary to provide at least 10 · W patterns

## Neural Networks:

#### **Contents:**

- Multi-layered Neural Networks
- Multi-layered Neural Networks: Analysis of Their Properties
- Multi-layered Neural Networks: an Application Example

#### **Contents:**

- Multi-layered Neural Networks
  - Back-Propagation Training Algorithm
  - Strategies to Speed-up the Training Process
- Multi-layered Neural Networks: Analysis of Their Properties
  - Kolmogorov's Theorem
  - Function Approximation
  - The Complexity of Learning
  - The Number of Regions in the Feature Space
  - Vapnik-Chervonenkis Dimension
- Multi-layered Neural Networks: an Application Example
  - Internal Knowledge Representation and Pruning
  - Sensitivity Analysis and Feature Selection
  - Analysis of the World Bank Data