

## MAXIMUM COVERAGE METHOD FOR FEATURE SUBSET SELECTION FOR NEURAL NETWORK TRAINING

Štefan BOOR

*Institute of Information Engineering, Automation and Mathematics  
Faculty of Chemical and Food Technology  
Slovak University of Technology in Bratislava  
Radlinského 9  
812 37 Bratislava, Slovakia  
e-mail: stefan.boor@stuba.sk*

Communicated by Vladimír Kvasnička

**Abstract.** Every real object having certain properties can be described by a number of descriptors, visual or other, e.g. mechanical, chemical etc. A set of descriptors (features) characterizing a given object is described in the paper by a vector of descriptors, where each entry of the vector determines a value of some feature of the object. In general, it is important to describe the object as completely as possible, which means by a large number of descriptors. This paper deals with a problem of selection of a proper subset of descriptors, which have the most substantial influence on the properties of the object, so that irrelevant descriptors could be excluded. For this purpose, we introduce a new method, Maximum Coverage Method (MCM). This method has been combined with optimization by a classical genetic algorithm. The described method is used for a data pre-processing, with the resulting selected features serving as an input for a neural network.

**Keywords:** Neural network, cluster, coverage, significant, shift, prediction, correctness, eliminating, separation

### 1 BASIC DESCRIPTION OF THE USED NEURAL NETWORK

One of the basic problems in every area of science is a search for a relation or a function between the structure of its objects and their properties. Let there be

a set of objects  $\mathbf{A}$  and a set of their properties  $\mathbf{B}$ . Let there exist a function

$$\mathbf{F} : \mathbf{A} \rightarrow \mathbf{B} \quad (1)$$

which assigns a property  $\mathbf{y} \in \mathbf{B}$  to each object  $\mathbf{x} \in \mathbf{A}$ . With such a function, it is possible to predict a property of an unknown object characterized by its descriptors. The ultimate goal is to construct a function in an analytical form, which describes the relation of properties of objects to their structure. Unfortunately, in many cases this goal is either unreachable, or reachable only with great difficulties. This drawback can be removed relatively successfully by using nonlinear mapping determined by multilayer neural networks.

During classification and prediction of objects with multilayer neural networks we start from the assumption that there exists a set of classified objects  $\mathbf{A}$ , which will be consequently divided into a training set  $\mathbf{A}_{train}$  and a testing set  $\mathbf{A}_{test}$ , so that:

$$\mathbf{A} = \mathbf{A}_{train} \cup \mathbf{A}_{test}. \quad (2)$$

Let every object of the set  $\mathbf{A}$  be classified by  $n$  features (descriptors) which means that we can assign  $n$ -dimensional vectors  $\vec{\mathbf{d}}$  of descriptors to the objects of the set  $\mathbf{A}$ . These vectors are assigned both to objects from the set  $\mathbf{A}_{train}$  and from the set  $\mathbf{A}_{test}$ :

$$\begin{aligned} \vec{\mathbf{d}}_{train} &= \left( d_{train}^{(1)}, d_{train}^{(2)}, \dots, d_{train}^{(n)} \right) \\ \vec{\mathbf{d}}_{test} &= \left( d_{test}^{(1)}, d_{test}^{(2)}, \dots, d_{test}^{(n)} \right). \end{aligned} \quad (3)$$

Let us presume that for each object  $\mathbf{x} \in \mathbf{A}$  there exists a  $\mathbf{y} \in \mathbf{B}$  so that  $\mathbf{F}(\mathbf{x}) = \mathbf{y}$ , which means that to every vector  $\vec{\mathbf{d}}$  a value  $\mathbf{y} \in \mathbf{B}$  is assigned. Our ultimate goal is to find a function  $\mathbf{F}$ . The purpose of neural networks is to find such a function  $\mathbf{G}(\mathbf{x}, \mathbf{w})$ , which arguments are from the set  $\mathbf{A}_{train}$  and  $\mathbf{w}$  is such a parameter (or parameters) of mapping  $\mathbf{G}$  that functional values  $\mathbf{G}(\mathbf{x}, \mathbf{w})$  are as close as possible to the ideal results of the function  $\mathbf{F}(\mathbf{x})$ , i.e., to the required values. Let us define an objective function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{|\mathbf{A}_{train}|} (\mathbf{G}(\mathbf{x}_i, \mathbf{w}) - \mathbf{y}_i)^2. \quad (4)$$

The requirement that the function  $\mathbf{G}(\mathbf{x}, \mathbf{w})$  maps the objects  $\mathbf{x}_{train}$  as close as possible to the required values  $\mathbf{y}_{train}$  is realized by minimization of the function  $E(\mathbf{w})$  with respect to the parameters  $\mathbf{w}$ . Let  $\bar{\mathbf{w}}$  be an optimal value of a parameter (parameters)  $\mathbf{w}$  determined by the minimization process of the function  $E(\mathbf{w})$

$$\bar{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbf{W}} E(\mathbf{w}) \quad (5)$$

where  $\mathbf{W}$  is a set (space) of admissible values of the parameter  $\mathbf{w}$ . By the function  $\mathbf{G}(\mathbf{x}, \mathbf{w})$  (shortly  $\mathbf{G}(\mathbf{w})$ ) an adapted neural network [2] simulates the original ideal

function  $F(x)$  for the values of the argument from the training set  $A_{train}$  on the basis of the minimization criterion (5).

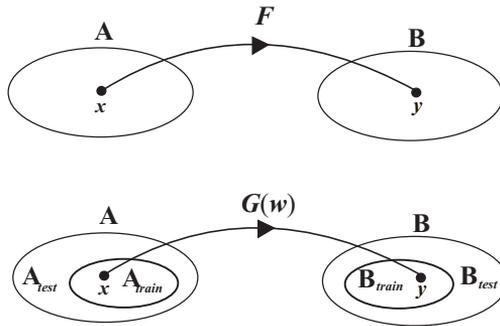


Fig. 1. Schematic description of the mapping  $F : A \rightarrow B$ . By narrowing of this mapping onto the subset  $A_{train}$  we get a new model mapping  $G(w)$ , whose function is determined by the parameter (parameters)  $w$

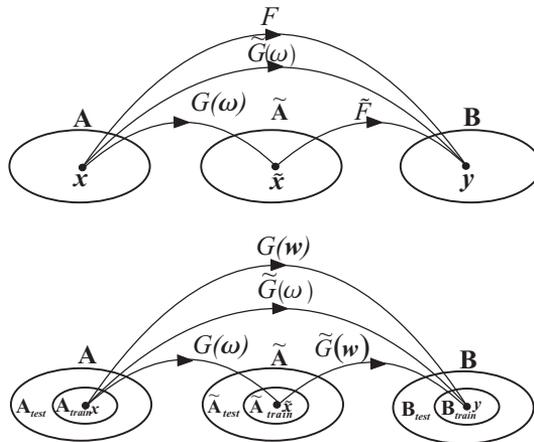


Fig. 2. If there exists such a function  $\tilde{G}(\omega) : A \rightarrow B$ , that  $G(\omega) : \tilde{A} \rightarrow A$  with a parameter (parameters)  $\omega$  and  $\tilde{G}(\omega) : \tilde{A} \rightarrow B$ , then the model mapping  $G(w)$  is replaceable by a (model) mapping  $\tilde{G}(\omega)$  with a parameter (parameters)  $\omega$ . Then the function  $F$  is replaced by three-layer neural network representing mapping  $\tilde{G}(\omega) : A \rightarrow B$ , where  $A$  is a set of target elements for the set  $A$  with a reduced number of descriptors.

The adapted neural net will be tested on the testing set, whether it approximates the ideal function learned from the training set. When the adapted function  $G(x, w)$  performs well also on the testing set, it is possible to use it also for prediction of properties of an unknown object  $x$  described by a corresponding vector of descriptors. The function  $G(x, w)$  which will be constructed suitably for the problem with

its parameter(s)  $\mathbf{w}$  determined from a minimization of the objective function  $E(\mathbf{w})$  so that it maps the training and testing set objects as close as possible to the required values, will be regarded as an analytical formula between the structure of studied objects and their properties, i.e., as the function  $\mathbf{F}(\mathbf{x})$  (Figure 1).

Adaptation of a neural network (function  $\mathbf{G}(x, w)$ ) is a rather difficult and computationally demanding process. Its complexity is influenced by several factors. One of them is the length of a vector of descriptors assigned to the objects from the set  $\mathbf{A}$ . This length corresponds to the number of features which describe these objects and on the basis of which they are classified. The features are in many cases introduced “ad-hoc” without a more detailed study of their interdependences and significance for the classification of the objects. In that case the function  $\mathbf{G}(x, w)$  might have superfluous number of parameters  $\mathbf{w}$ , which is proportional to the size of a neural network and to the complexity of its adaptation process. The problem of excluding of descriptors, which do not have a substantial influence on the classification of objects can be transformed onto a problem of finding a function  $\mathbf{G}(x, \omega)$  (shortly  $\mathbf{G}(\omega)$ ) with a parameter (parameters)  $\omega$ , which will map a set  $\mathbf{A}$  onto a set  $\tilde{\mathbf{A}}$  so that every vector of descriptors (3a), resp. (3b) (Figure 3) will be mapped onto the vector

$$\begin{aligned}\tilde{\mathbf{d}}_{train} &= \left( \tilde{d}_{train}^{(1)}, \tilde{d}_{train}^{(2)}, \dots, \tilde{d}_{train}^{(m)} \right) \\ \tilde{\mathbf{d}}_{test} &= \left( \tilde{d}_{test}^{(1)}, \tilde{d}_{test}^{(2)}, \dots, \tilde{d}_{test}^{(m)} \right)\end{aligned}\quad (6)$$

while  $m < n$ . Then the functions  $\mathbf{F}$  and  $\mathbf{G}(\mathbf{w})$  are replaced by functions  $\tilde{\mathbf{F}}$  and  $\tilde{\mathbf{G}}(\mathbf{w})$ , respectively (Figure 2), while the function  $\mathbf{G}(\omega)$  reduces the number of descriptors of the objects from the set  $\mathbf{A}$ . This function will map an  $n$ -dimensional space  $\mathbf{R}^n$  onto an  $r$ -dimensional space  $\mathbf{R}^r$ :

$$\mathbf{G}(\omega) : \mathbf{R}^n \rightarrow \mathbf{R}^r. \quad (7)$$

The subject of our following consideration will be construction of a function  $\mathbf{G}(\omega)$  and finding its parameters  $\omega$  so that the conditions

$$\forall \mathbf{x} \in \mathbf{A}_{train} : \tilde{\mathbf{G}}(\omega) : \mathbf{x}_{train} \rightarrow \mathbf{y}_{train} \quad (8)$$

$$\forall \mathbf{x} \in \mathbf{A}_{test} : \tilde{\mathbf{G}}(\omega) : \mathbf{x}_{test} \rightarrow \mathbf{y}_{test} \quad (9)$$

will be satisfied as fully as possible.

It is obvious that the requirement (9) is very demanding and it can be satisfied practically only to a certain precision. For the description of the function  $\mathbf{G}(\omega)$  and finding its parameters we shall use a new MCM method (Maximum coverage method). The problem of excluding superfluous descriptors was solved by the KNN method (K Nearest Neighbour) [1].

## 2 THE BASIC PRINCIPLES OF MAXIMUM COVERAGE METHOD

Entities from a set  $\mathbf{A}$  are entities from an  $n$ -dimensional left-open and right-closed interval  $\mathbf{J} = ]a^{(i)}, b^{(i)}]^n, i = 1, \dots, n$ , classified by an  $n$ -dimensional vector of descriptors  $\vec{d}_j = (d_j^{(1)}, d_j, \dots, d_j^{(n)})$  for each entity  $x_j \in \mathbf{A}, j = 1 \dots m$ .  $\mathbf{J} = ]a^{(i)}, b^{(i)}]^n = ]a^{(1)}, b^{(1)}] \times ]a^{(2)}, b^{(2)}] \times \dots \times ]a^{(n)}, b^{(n)}] = J^{(1)} \times J^{(2)} \times \dots \times J^{(n)}$ . The number of entities  $\mathbf{x} \in \mathbf{A}$  is  $m$  and the set  $\mathbf{B}$  is a set of attributes of entities  $\mathbf{x} \in \mathbf{A}$ , and  $|\mathbf{B}| < m$ . The construction of intervals and end points  $A(a^1, \dots, a^n), B(b^1, \dots, b^n)$  of the interval  $\mathbf{J}$  is as follows:

1. Make the sequence  $d_1^{(i)} < d_2^{(i)} < \dots < d_{m^{(i)}}^{(i)} \forall i = 1, \dots, n$ , when  $m^{(i)} \leq m$ .

$$m^{(i)} = m \Leftrightarrow \forall \mathbf{x}_p, \mathbf{x}_q \in \mathbf{A}; x_p^{(i)} \neq x_q^{(i)}, \forall p \neq q, p, q \in \{1, 2, \dots, m\}, m = |\mathbf{A}|$$

2.  $a_1^{(i)} < a_2^{(i)} < \dots < a_{m^{(i)}-1}^{(i)} \forall i = 1, \dots, n; a_j^{(i)} = \frac{d_{j+1}^{(i)} + d_j^{(i)}}{2} \forall j = 1, 2, \dots, m^{(i)} - 1, i = 1, \dots, n$ .

3. Let's put  $a^i = a_0^{(i)} = d_1^{(i)} - \frac{d_2^{(i)} - d_1^{(i)}}{2}; b^i = a_{m^{(i)}}^{(i)} = d_{m^{(i)}}^{(i)} + \frac{d_{m^{(i)}}^{(i)} - d_{m^{(i)}-1}^{(i)}}{2}, i = 1, \dots, n$ .

This construction assigns just one  $n$ -dimensional interval  $C$  to  $\forall \mathbf{x} \in \mathbf{A}$ .

**Definition 1.** Let  $\bar{N} = \{1, \dots, n\}$ . The set

$$\mathbf{D}^{(i)} = \left\{ a_j^{(i)} \in ]a^{(i)}, b^{(i)}]; i \in \bar{N}; j = 1, 2, \dots, m_i - 1 \right\} \tag{10}$$

will be called a partition  $\mathbf{D}^{(i)}$  of interval  $\mathbf{J}^{(i)}$  for index  $j$  of the element  $x_j \in \mathbf{A}, j = 1, \dots, m, i = 1, \dots, n$ .

**Definition 2.** Let  $\bar{N} = \{1, \dots, n\}$ ,

$$\mathbf{J}_j^{(i)} = ]a_{j-1}^i, a_j^i]; i \in \bar{N}; j = 1, 2, \dots, m_i \tag{11}$$

Then the set

$$X^{(i)} = \{ \mathbf{J}_j^{(i)}; j = 1, 2, \dots, m_i; i \in \bar{N} \} \tag{12}$$

will be called a disjunctive coverage of the interval  $\mathbf{J}^{(i)} = ]a^{(i)}, b^{(i)}]$ .

**Definition 3.** Let us define the set  $covC$  as a disjunctive coverage of an  $n$ -dimensional interval  $\mathbf{J} = ]a^{(i)}, b^{(i)}]^n$  with the following attribute:

$$covC = \{ C \subset \mathbf{J}; C = X_1 \times X_2 \times \dots \times X_n, \forall X_1 \in X^{(1)}, X_2 \in X^{(2)}, \dots, X_n \in X^{(n)} \} \tag{13}$$

**Definition 4.** The set  $clC \in covC$  will be called the *cluster* of  $\mathbf{A}$  if to each entity  $\mathbf{x} \in C$ , which is in  $\mathbf{A}$ , the same value  $\mathbf{y} \in \mathbf{B}$  is assigned.

A non empty set  $clC \subset covC$  with the previous attribute will be called the set of *clusters*.

When preparing data for neural network entry, the following attributes are important for its optimal training and prediction:

- the correctness of data
- eliminating the descriptors not significant for classification of  $\mathbf{y} \in \mathbf{B}$
- optimal separation of the set  $\mathbf{A}$  to  $\mathbf{A}_{train} \cup A_{test}$ .

Our goal is to introduce a method which will answer the following three problems.

**Definition 5.** The descriptor  $d^{(i)}$  will be called non-significant if there is a set  $covC$ , for which:  $|clC|$  is minimal and  $\mathbf{D}^{(i)} = \emptyset$  for each entity  $C \in clC$ .

Figures 3 a) and 3 b) show points with two different values before and after constructing clusters, when descriptor  $d^{(2)}$  is non-significant, i.e., the attributes of entities  $\mathbf{x} \in \mathbf{A}$  in clusters can be described exactly by  $d^{(1)}$ .

**Definition 6.** Let the vectors of descriptors  $\vec{\mathbf{d}}_i = (d_i^{(1)}, d_i^{(2)}, \dots, d_i^{(n)})$  and  $\vec{\mathbf{d}}_j = (d_j^{(1)}, d_j^{(2)}, \dots, d_j^{(n)})$  characterize the attributes  $\mathbf{y}_i$ , respectively  $\mathbf{y}_j$ . Then the data consisting of entities in sets  $\mathbf{A}$ ,  $\mathbf{B}$  are correct only provided that for every two vectors of descriptors the equality  $\vec{\mathbf{d}}_i = \vec{\mathbf{d}}_j$  implies  $\mathbf{y}_i = \mathbf{y}_j$ . Otherwise, the data are incorrect.

**Remark.** This data regularity definition means that data are regular only if there exists a function  $F : \mathbf{A} \rightarrow \mathbf{B}$ .

### 3 CONSTRUCTION OF CLUSTERS AND PURPOSE FUNCTION

To the entity  $C_1 \in covC$  a value  $\mathbf{y}_1 \in \mathbf{B}$  of a randomly chosen entity  $\mathbf{x} \in \mathbf{A}$  is assigned only if  $\mathbf{x} \in C_1$ . Let us denote by  $c_1$  the number of entities  $\mathbf{x} \in C_1$  for which the value of the entity  $\mathbf{x}$  does not equal  $\mathbf{y}_1$ .

To the entity  $C_i \in covC$  a value  $\mathbf{y}_i \in \mathbf{B}$  of the entity  $\mathbf{x} \in \mathbf{A} - \bigcup_{j=1}^{i-1} C_j$  is assigned. Let  $c_i$  be the number of entities  $\mathbf{x} \in C_i$  whose value does not equal  $\mathbf{y}_i$ .

$C_k$  is the last cluster only if  $\mathbf{A} - \bigcup_j = 1^k C_j = \emptyset$ . The optimal construction of the set  $covC$ , as well as the clusters will be reached by optimization of the following function

$$f(\omega) = \sum_{i=1}^n c_i + \frac{\sum_{i=1}^n |\mathbf{D}^{(i)}|}{n \times \max_{i \in \{1, \dots, n\}} (m^{(i)})} \tag{14}$$

provided that  $a_j^{(i)} = \omega_j^{(i)} a_j^{(i)}$ ,  $i = 1, \dots, n, j = 1, \dots, m^{(i)}$ , where  $\omega_j^{(i)}$  is a binary value, whose zero value means that the boundary point  $a_j^{(i)}$  is deleted and the neighbouring intervals are joined. This in turn may affect the value of  $c_i$ .

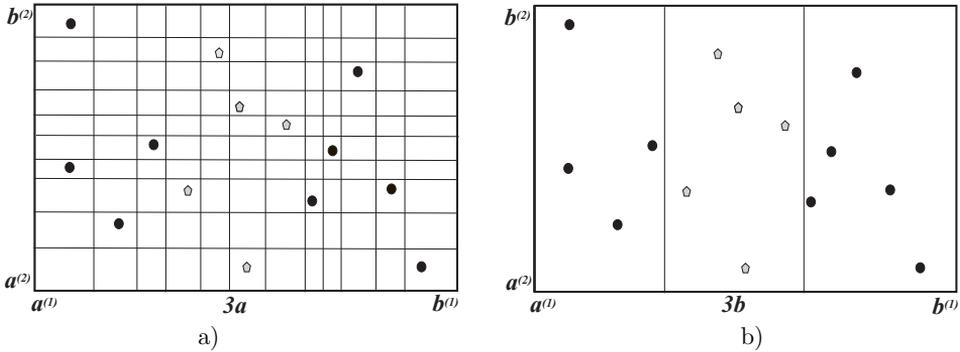


Fig. 3. Randomly generated points on 2-dimensional interval  $[a^{(1)}, b^{(1)}] \times [a^{(2)}, b^{(2)}]$  with attributes  $y_1 = , y_2 =$  that in some cases can be sorted to clusters, where each cluster does not depend on the  $i^{\text{th}}$  descriptor, i.e.  $d^{(i)} \in [a^{(i)}, b^{(i)}]$  without any dividing points of this interval (3a). In the first stage of the optimization process each element  $x_i \in A$  is assigned to the corresponding cluster with attribute  $y_i \in B$ . In the first phase of the optimisation process, each element  $x_i \in A$  is assigned to a cluster with a corresponding attribute  $y_i \in B$ . After the optimization (3b), the clusters formed by objects  $x \in A$ , which have the same attributes  $y \in B$ , are of the maximum size. This figure shows that the entities from set  $x \in A$  may be sorted to clusters with identical attributes, where  $d^{(2)}$  is non-significant from the view of adaptation process of the neural network. The attributes  $y \in B$  are exactly determined by the descriptor  $d^{(1)}$ .

#### 4 USING THE GENETIC ALGORITHMS (GA)

Let the population [3, 4]  $P$  be composed of  $p$  randomly generated mapping set vectors (chromosomes)  $\omega^{(i)} \in \{0, 1\}^{m^{(i)}-1} \forall i \in \overline{N}$ . Let the objective function be given by equation. During iterations the population  $P$  changes through an application of crossover and mutation to its vectors. The selection of vectors to the process of crossover and further selection is realized by the so-called roulette, which means by a random selection of a chromosome with the probability proportional to its fitness  $F$  defined using the objective function by the formula

$$F(\omega) = \frac{\frac{1}{f(\omega)}}{\sum_{\omega \in P} \frac{1}{f(\omega)}}. \tag{15}$$

The crossover performs an exchange of randomly selected parts of two vectors chosen for a crossover process. The *crossover* operator is applied to 20 % of selected vectors; the remaining 80 % are not affected by the operation.

The *mutation* operator is applied to all mapping set vectors created by the *crossover*. The *mutation* randomly changes single entries of the vector, each with probability 0.002 %. After the operations of *crossover* and *mutation* are realized, a newly created population replaces the “parental” population of chromosomes (vectors of descriptors). Each iteration consists in calculation of objective function values



tor  $d^{(2)}$  is the same and maximal, thus the clusters are practicably assigned with descriptor  $d^{(1)}$ .

## 6 ILLUSTRATIVE EXAMPLE

In [3] the authors dealt with construction of the neural network for the prediction of  $^{13}\text{C}$  NMR (nuclear magnetic resonance) of chemical shifts of acyclic alkanes – saturated hydrocarbons. The chemical shift is a local property of carbon atoms; therefore, the referred carbon atom will be described by its closest neighbourhood (Figure 5). The used description of the neighbourhood consists in a creation of a vector describing a frequency of presence of some in advance chosen fragments in the alkane, i.e. how often the pertaining fragment appears in the alkane with the specified carbon atom.

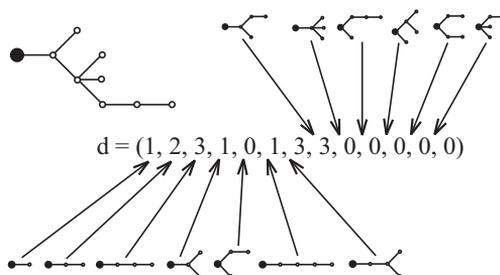


Fig. 5. Diagram of the description of a carbon atom denoted by a heavy dot in the molecule of the acyclic hydrocarbon. The neighborhood of this atom is described by a vector of descriptors  $\vec{d}$ , containing 13 integers corresponding to the “frequencies of appearance” of the fragments shown at the ends of the arrows.

In order to create training and a testing set, this vector was enlarged by an experimentally measured value of the chemical shift. The objects of the set **A** from Figure 2 are now represented by vectors of descriptors  $\vec{d}$  and the set **B** is determined by the corresponding experimentally measured chemical shifts. A well adapted neural network approximates the function **F** (Figure 1). A question arises: “Do all the descriptors have an influence on the chemical shift of protons in carbon atom?” We shall try to get the answer by a maximum coverage method (MCM). The population was composed of 16 randomly generated mapping chromosomes, where No. 1 was generated with probability 0.24. The crossover process was applied to 20% of the population and the mutation process was carried out with probability 0.021. After 464 iterations the value of the optimized function (14) is 5.32867133 and the resulting chromosome provides the following results.

1. Descriptors  $d^{(5)}$ ,  $d^{(9)}$ ,  $d^{(10)}$ ,  $d^{(11)}$ ,  $d^{(13)}$  are non-significant in the meaning of Definition 5.

1	1	1	1	1	1								
1	1	0	1	0	0	0	0	1					
1	1	1	1	1	0	1							
1	1	1	0	0	0	1							
1	0	0	0	0	0	0	0	0	0	0	0	1	
1	1	1	1	0	1								
1	1	0	1	0	1								
1	1	0	0	0	0	0	0	1					
1	0	0	1										
1	0	0	0	0	0	0	1						
1	0	0	0	0	0	0	0	1					
1	1	0	0	1	0	1	1	0	1				
1	0	0	0	0	0	0	0	0	1				

Fig. 6. The resulting chromosome, consisting of binary mapping 0, 1 determining the points of division  $\mathbf{D}^{(i)}$  of the interval  $\mathbf{J}^{(i)}$ . The best result shows that  $\mathbf{J}^{(5)}, \mathbf{J}^{(9)}, \mathbf{J}^{(10)}, \mathbf{J}^{(11)}, \mathbf{J}^{(13)}$  do not have (save the boundary points) any points of division, which means that the mapping of the elements  $\mathbf{x} \in \mathbf{A}$  does not depend on the corresponding descriptors.

2. Descriptors  $d^{(2)}, d^{(8)}$  also have a small significance for the clustering, since the whole interval is divided by one dividing point only.

### 7 COMPARISON OF RESULTS OF KNN AND MCM METHODS

Let us define weights  $w^{(i)}$  of significance of descriptors as

$$w^{(i)} = \frac{m_i - 2}{n_i - 2}, n_i > 2 \tag{16}$$

where  $m_i$  is the number of number of intervals  $\mathbf{J}^{(i)}$  reduced according to the chromosome, where 0 entries do not define division of the interval and  $n_i$  is the number of all possible dividing points of the given interval.

$w^{(i)}$	1	2/6	4/5	2/5	0	3/4	1/2	1/7	0	0	0	1/2	0
$\omega$	1	1	0	1	1	0	1	1	0	1	0	0	1

Fig. 7

The first row of the table shows weights of significance of descriptors calculated from results of the MCM method. The second row provides a “mask” vector of method KNN [4]. Descriptor with a value 1 of “mask” vector means its significance, descriptors with values equal to 0 are labeled as insignificant.

By comparison of the two methods we arrived at two descriptors selected as significant by both methods; MCM selected altogether 3 insignificant descriptors and KNN selected 5 insignificant descriptors. It is necessary to mention that a KNN method indicates sharp significance or insignificance only. The MCM method indicates weights of significance, while by comparing results we can say that  $d$  is selected

as insignificant by both methods. The results of both methods do not fit totally, but they are very similar.

### 8 TESTING OF THE MCM

The results presented below method were obtained from 30 test runs with the following set of parameter values:

- 16 randomly generated mapping chromosomes
- No. 1 was generated with probability 0.24
- the crossover process was applied to 20% of the population
- the mutation process was carried out with probability 0.021

where  $\bar{w}^{(i)}$  is the average of weights  $w^{(i)}$ . The average number of cycles necessary to obtain results was 549,56 and the average value of the optimized function is 5.333566434. It is apparent that in the meaning of Definition 5 the descriptors  $d^{(5)}, d^{(9)}, d^{(13)}$  are non-significant for all the runs,  $d^{(10)}, d^{(11)}$  were non-significant in 8 tests. Every test resulted in at least 4 nonsignificant descriptors. The best obtained result is shown in Figure 6.

$i$	1	2	3	4	<u>5</u>	6	7	8	<u>9</u>	<u>10</u>	<u>11</u>	12	<u>13</u>
$\bar{w}^{(i)}$	1	0.45	0.64	0.30	<u>0</u>	0.75	0.5	0.23	<u>0</u>	<u>0.19</u>	<u>0.04</u>	0.36	<u>0</u>

Fig. 8

### 9 CONCLUSIONS

The main result of the presented paper is a new method of selection of classifiers of elements, or in other words, exclusion of non-significant features from a set of descriptors, which enter as input values into a neural network classifier. The method solves significance of a property, and provides a possibility of reduction of the input vector of descriptors with a suitable setting of parameters. Moreover, the method solves correctness of data and selection of elements into the testing set.

The advantages of the proposed MCM method as compared with the well established KNN method are as follows:

1. The design of the object function is proposed in such a way that it shows the number of non-correct elements in the data set, which the KNN method cannot provide. If  $f(\omega) < 1$ , data are regular. When for example  $f(\omega) = 5.333566434$  then there exist 5 copies of the vector of values of descriptors  $\vec{d}$ , which differ in their functional values  $\mathbf{y}$ .
2. It follows from the principles of the MCM method (Figures 3 a), 3 b)) that there exists a function  $\mathbf{F} : \mathbf{A} \rightarrow \mathbf{B}$ , in an equivalent form  $\mathbf{G}(\omega) : \mathbf{R}^{n-p} \rightarrow \mathbf{R}^r$ ,

where  $p$  is the number of non-significant features from a set of descriptors. The function  $\mathbf{G}(\omega)$ , which is able to hold the mapping  $\mathbf{F} : \mathbf{A} \rightarrow \mathbf{B}$  while dropping out non-significant descriptors. This means that the neural network shall work with a smaller number of descriptors leaving out 38 % of non-significant descriptors. This reduction will allow faster learning of the neural network with the same resulting quality of prediction.

3. The MCM method not only allows us to reduce the number of descriptors (unlike the KNN [4] method), but also, unlike KNN, enables the study of neural network optimisation for weighted input data.

### Acknowledgement

This work was supported by Scientific Grant Agency Vega of Slovak Republic under grants 1/0071/09 and by Slovak Research and Development Agency APVV-0446-07.

### REFERENCES

- [1] KVASNIČKA, V.—POSPÍCHAL, J.—HESEK, D.: Augmented Simulated Annealing Algorithm for the TSP. Central European Journal for Operations Research and Economics, Vol. 2, 1993, pp. 307–317.
- [2] MITCHEL, M.: An Introduction to Genetic Algorithms. MIT Press, Boston, MA 1996.
- [3] SVOZIL, D.—POSPÍCHAL, J.—KVASNIČKA, V.: Neural-Network Prediction of Carbon-13 NMR Chemical Shifts of Alkanes. J. Chem. Inf. Comp. Sci., Vol. 35, 1995, pp. 924–928.
- [4] BOOR, Š.—KVASNIČKA, V.—POSPÍCHAL, J.: Feature Subset Selection by GA Combined with KNN Method. Mendel 2000, PC-DIR, Brno, 2000, ISBN 80-214-1609-2, pp. 407–412.



**Štefan Boor** graduated in mathematic and descriptive geometry from Faculty of Natural Sciences, Comenius University in Bratislava in 1980, and received the Ph.D. in theoretical and computational chemistry from Faculty of Chemical and Food Technologies, Slovak University of Technology, Bratislava in 2010. Now he is with the Department of Mathematics at the same Faculty as Assistant Professor. He has co-authored journal and conference papers related to fuzzy logic, evolutionary algorithms and neural networks.