

PCA-BASED DATA MINING PROBABILISTIC AND FUZZY APPROACHES WITH APPLICATIONS IN PATTERN RECOGNITION

Luminita State

Dept. of Computer Science, University of Pitesti, Caderea Bastliei #45, Bucuresti – 1, Romania

Catalina Cocianu

Dept. of Computer Science, Academy of Economic Studies, Calea Dorobantilor #15-17, Bucuresti –1, Romania

Panayiotis Vlamos

Ionian University, Corfu, Greece

Viorica Stefanescu

Dept. of Mathematics, Academy of Economic Studies, Calea Dorobantilor #15-17, Bucuresti –1, Romania

Keywords: Principal component analysis, fuzzy clustering, supervised learning, cluster analysis, pattern recognition, data mining.

Abstract: The aim of the paper is to develop a new learning by examples PCA-based algorithm for extracting skeleton information from data to assure both good recognition performances, and generalization capabilities. Here the generalization capabilities are viewed twofold, on one hand to identify the right class for new samples coming from one of the classes taken into account and, on the other hand, to identify the samples coming from a new class. The classes are represented in the measurement/feature space by continuous repartitions, that is the model is given by the family of density functions $(f_h)_{h \in H}$, where H stands for the finite set of hypothesis (classes). The basis of the learning process is represented by samples of possible different sizes coming from the considered classes. The skeleton of each class is given by the principal components obtained for the corresponding sample.

1 PRINCIPAL COMPONENTS

The starting point for PCA is a n -dimensional random vector \mathbf{X} . There is available a sample $\mathbf{X}(1), \mathbf{X}(2), \dots, \mathbf{X}(T)$ from this random vector. No explicit assumptions on the probability density of the vectors are made in PCA, as long as the first – order and the second –order statistics are known or can be estimated from the sample. Also, no generative model is assumed for vector \mathbf{X} .

In the PCA transform, the vector \mathbf{X} is first centered by subtracting its mean, $\mathbf{X} = \mathbf{X} - E(\mathbf{X})$. In practice, the mean of the n -dimensional vector \mathbf{X} is estimated from the available sample. In the following, we assume that the vector \mathbf{X} is centered.

Next, \mathbf{X} is linearly transformed to another vector \mathbf{Y} with m elements, $m < n$, so that the redundancy induced by the correlations is removed. The transform consists in obtaining a rotated orthogonal coordinate system such that the elements of \mathbf{X} in the new coordinates become uncorrelated. At the same time, the variances of the projections of \mathbf{X} on the new coordinates become uncorrelated. At the same time, the variances of the projections of \mathbf{X} on the new coordinate axes are maximized so that the first axis corresponds to the maximal variance, the second axis corresponds to the maximal variance in the direction orthogonal to the first axis, and so on (Hyvarinen, 2001).

State L., Cocianu C., Vlamos P. and Stefanescu V. (2006).

PCA-BASED DATA MINING PROBABILISTIC AND FUZZY APPROACHES WITH APPLICATIONS IN PATTERN RECOGNITION.

In *Proceedings of the First International Conference on Software and Data Technologies*, pages 55-60

Copyright © SciTePress

In mathematical terms, consider a linear combination of the elements x_1, x_2, \dots, x_n of the vector \mathbf{X} , $y_1 = \sum_{i=1}^n w_{i1} x_i = \mathbf{W}_1^T \mathbf{X}$.

We look for a weight vector \mathbf{W}_1 maximizing the PCA criterion,

$$E(y_1^2) = E\left(\left(\mathbf{W}_1^T \mathbf{X}\right)^2\right) = \mathbf{W}_1^T \mathbf{S} \mathbf{W}_1 \quad (1)$$

$$\|\mathbf{W}_1\| = 1,$$

where the matrix \mathbf{S} is the covariance matrix of \mathbf{X} given for the zero-mean vector \mathbf{X} by the correlation matrix.

It is well known from basic linear algebra that the solution of the PCA problem is given in terms of the unit-length eigen vectors of the matrix \mathbf{S} , $\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_n$. The ordering of the eigen vectors is such that the corresponding eigen values $\lambda_1, \lambda_2, \dots, \lambda_n$ satisfy $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. The solution maximizing $E(y_1^2) = \mathbf{W}_1^T \mathbf{S} \mathbf{W}_1$ is given by $\mathbf{W}_1 = \boldsymbol{\phi}_1$ and the first principal component of \mathbf{X} is $y_1 = \boldsymbol{\phi}_1^T \mathbf{X}$. The PCA criterion can be generalized to m principal components, $1 \leq m \leq n$. Let

$y_m = \sum_{i=1}^n w_{im} x_i = \mathbf{W}_m^T \mathbf{X}$ be the m -th principal component, with \mathbf{W}_m the corresponding unit norm weight vector. The solution maximizing $E(y_m^2) = \mathbf{W}_m^T \mathbf{S} \mathbf{W}_m$ under the constrain $E(y_m y_k) = 0$, for $1 \leq k < m$ is given by, (Hyvarinen, 2001) $\mathbf{W}_m = \boldsymbol{\phi}_m$ and the m -th principal component of \mathbf{X} is $y_m = \boldsymbol{\phi}_m^T \mathbf{X}$.

2 CLUSTER ANALYSIS

Cluster analysis is a data processing technique aiming to identify the natural grouping trends existing in a data collection, producing a set of overlapping clusters, the elements belonging to the same cluster sharing similar features. So far, there have been proposed a relatively small number of methods for testing the existence/inexistence of a natural grouping tendency in a data collection, most of them being based on arguments coming from mathematical statistics and heuristic graphical techniques (Panayirci and Dubes, 1983, Smith and Jain, 1984, Jain and Dubes, 1988, Tukey, 1977, Everitt, 1978).

The data are represented by p -dimensional vectors, $X = (x_1, \dots, x_p)^t$, whose components are the feature values of a specified attributes and the classification is performed against a certain given label set. The classification of a data collection $\mathfrak{S} = \{X_1, \dots, X_n\} \subset \mathfrak{R}^p$ corresponds to a labelling strategy of the objects of \mathfrak{S} .

In the fuzzy approaches, the clusters are represented as fuzzy sets $(u_i, 1 \leq i \leq c)$, $u_i : \mathfrak{S} \rightarrow [0, 1]$, where $u_{ik} = u_i(X_k)$ is the membership degree of X_k to the i -th cluster, $1 \leq i \leq c$, $1 \leq k \leq n$. A c -fuzzy partition is represented by the matrix $U = \|u_{ik}\| \in M_{c \times n}$. The number of labels c has to be selected in advance, the problem of finding the optimal c is usually referred as cluster validation.

The main types of label vectors are *crisp* N_c , *fuzzy* N_p , and *possibilistic* N_{poz} , defined as follows,

$$N_c = \left\{ y \mid y \in \mathfrak{R}^c, y = (y_1, y_2, \dots, y_c), y_i \in \{0, 1\}, \right. \\ \left. 1 \leq i \leq c, \sum_{i=1}^c y_i = 1 \right\} = \{e_1, e_2, \dots, e_c\}, \quad (2)$$

$$\text{where } (e_i)_j = \delta_{ij} = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}$$

$$N_p = \left\{ y \in \mathfrak{R}^c \mid y = (y_1, y_2, \dots, y_c), \forall i, y_i \in [0, 1], \right. \\ \left. \sum_{i=1}^c y_i = 1 \right\}, \quad (3)$$

$$N_{poz} = \left\{ y \in \mathfrak{R}^c \mid y = (y_1, y_2, \dots, y_c), \forall i, y_i \in [0, 1], \right. \\ \left. \exists j, y_j \neq 0 \right\}, \quad (4)$$

Obviously, $N_{poz} \supset N_p \supset N_c$. If we denote by $U = [U_1, \dots, U_n] = \|u_{ij}\|$ a partition of \mathfrak{S} , then, according to the types of label vectors, we get the c -partition types M_{pos} , M_p and M_c ,

$$M_{pos} = \left\{ U \mid U \in M_{c \times n}, U = [U_1, \dots, U_n], \right. \\ \left. \forall k, U_k \in N_{pos}, \forall i, \sum_{k=1}^n u_{ik} > 0 \right\} \quad (5)$$

$$M_p = \left\{ U \mid U \in M_{pos}, \forall k, U_k \in N_p \right\} \quad (6)$$

$$M_c = \left\{ U \mid U \in M_p, \forall k, U_k \in N_c \right\} \quad (7)$$

Note that $M_c \subset M_p \subset M_{pos}$.

3 C-MEANS MODEL

The most popular classification methods are the *c-means* algorithms. The variational problem corresponding to *c-means* model is given by

$$\min_{(U,V)} \left\{ J_m(U,V;w) = \sum_{i=1}^c \sum_{k=1}^n u_{ik}^m D_{ik}^2 + \sum_{i=1}^c w_i \sum_{k=1}^n (1-u_{ik})^m \right\}$$

(8) where

$U \in M_c / M_p / M_{pos}$, $V = (v_1, \dots, v_c) \in M_{c \times p}$, v_i is the centroid of the i -th cluster, $w = (w_1, \dots, w_c)^T$ is the penalties vector corresponding to the cluster system, $m \geq 1$ is the fuzzification degree, and $D_{ik}^2 = \|x_k - v_i\|^2$.

Let us denote by (\hat{U}, \hat{V}) a solution of (8). Then,

1. The crisp model:

$$(U, V) \in M_c \times M_{c \times p}; w_i = 0, 1 \leq i \leq c,$$

$$\hat{u}_{ik} = \begin{cases} 1, & D_{ik} \leq D_{ij}, i \neq j \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

$$\hat{v}_i = \frac{\sum_{k=1}^n \hat{u}_{ik} x_k}{\sum_{k=1}^n \hat{u}_{ik}}; \quad 1 \leq i \leq c, 1 \leq k \leq n \quad (10)$$

2. The fuzzy model:

$$(U, V) \in M_p \times M_{c \times p}; m > 1, w_i = 0, 1 \leq i \leq c$$

$$\hat{u}_{ik} = \left[\sum_{j=1}^c \left(\frac{D_{ik}}{D_{jk}} \right)^{\frac{2}{m-1}} \right]^{-1} \quad (11)$$

$$\hat{v}_i = \frac{\sum_{k=1}^n u_{ik}^m x_k}{\sum_{k=1}^n u_{ik}^m}; \quad 1 \leq i \leq c, 1 \leq k \leq n \quad (12)$$

3. The possibilistic model:

$$(U, V) \in M_{pos} \times M_{c \times p}; \forall i, w_i > 0$$

$$\hat{u}_{ik} = \left[1 + \left(\frac{D_{ik}^2}{w_i} \right)^{\frac{1}{m-1}} \right]^{-1}$$

$$\hat{v}_i = \frac{\sum_{k=1}^n u_{ik}^m x_k}{\sum_{k=1}^n u_{ik}^m}; \quad 1 \leq i \leq c, 1 \leq k \leq n \quad (13)$$

The general scheme of a cluster procedure φ is,

$t \leftarrow 0$

repeat

$t \leftarrow t+1$

$$U_t \leftarrow F_\varphi(V_{t-1})$$

$$V_t \leftarrow G_\varphi(U_{t-1})$$

until $(t = \text{Tor} \|V_t - V_{t-1}\| \leq \varepsilon)$

$$(U, V) \leftarrow (U_t, V_t)$$

where c is the given number of clusters, T is upper limit on the number of iterations, m is the weight parameter, $1 \leq m < \infty$, C is the terminal condition, w is the system of weights $\forall i, w_i > 0$, $V_0 = (v_{1,0}, \dots, v_{c,0}) \in M_{c \times p}$ is the initial system of centroids and F_φ, G_φ are the updating functions.

4 PCA-BASED ALGORITHM FOR EXTRACTING SKELETON INFORMATION

In the following a new learning by examples PCA-based algorithm for extracting skeleton information from data to assure both good recognition performances, and generalization capabilities, is developed. Here the generalization capabilities are viewed twofold, on one hand to identify the right class for new samples coming from one of the classes taken into account and, on the other hand, to identify the samples coming from a new class. The classes are represented in the measurement/feature space by continuous repartitions, that is the model is given by the family of density functions $(f_h)_{h \in H}$, where H stands for the finite set of hypothesis (classes).

The basis of the learning process is represented by samples of possible different sizes coming from the considered classes. The skeleton of each class is given by the principal components obtained for the corresponding sample. The recognition algorithm identifies the class whose skeleton is the "nearest" to the tested example, where the closeness degree is expressed in terms of the amount of disturbance determined by the decision of allotting it to the corresponding class. The model is presented as follows. Let X_1, X_2, \dots, X_N be a series of n -dimensional vectors coming from a certain class C . The sample covariance matrix is

$$\Sigma_N = \frac{1}{N-1} \sum_{i=1}^N (X_i - \mu_N)(X_i - \mu_N)^T, \quad (14)$$

where $\mu_N = \frac{1}{N} \sum_{i=1}^N X_i$.

We denote by $\lambda_1^N \geq \lambda_2^N \geq \dots \geq \lambda_n^N$ the eigen values and by $\psi_1^N, \dots, \psi_n^N$ the corresponding orthonormal eigen vectors of Σ_N .

If X_{N+1} is a new sample, then, for the series $X_1, X_2, \dots, X_N, X_{N+1}$, we get

$$\Sigma_{N+1} = \Sigma_N + \frac{1}{N+1} (X_{N+1} - \mu_N)(X_{N+1} - \mu_N)^T - \frac{1}{N} \Sigma_N \quad (15)$$

Lemma. In case the eigen values of Σ_N are pairwise distinct, the following first order approximations hold,

$$\lambda_i^{N+1} = \lambda_i^N + (\psi_i^N)^T \Delta \Sigma_N \psi_i^N = (\psi_i^N)^T \Sigma_{N+1} \psi_i^N \quad (16)$$

$$\psi_i^{N+1} = \psi_i^N + \sum_{j=1, j \neq i}^n \frac{(\psi_j^N)^T \Delta \Sigma_N \psi_i^N}{\lambda_i^N - \lambda_j^N} \psi_j^N \quad (17)$$

Proof Using the perturbation theory, we get, $\Sigma_{N+1} = \Sigma_N + \Delta \Sigma_N$ and, $\psi_i^{N+1} = \psi_i^N + \Delta \psi_i^N$, $\lambda_i^{N+1} = \lambda_i^N + \Delta \lambda_i^N$, $1 \leq i \leq n$. Then,

$$\Delta \Sigma_N = \frac{1}{N+1} (\mathbf{X}_{N+1} - \mu_N)(\mathbf{X}_{N+1} - \mu_N)^T - \frac{1}{N} \Sigma_N \quad (18)$$

$$\begin{aligned} & (\Sigma_N + \Delta \Sigma_N)(\psi_i^N + \Delta \psi_i^N) = \\ & = (\lambda_i^N + \Delta \lambda_i^N)(\psi_i^N + \Delta \psi_i^N) \end{aligned} \quad (19)$$

Using first order approximations, from (19) we get,

$$\begin{aligned} & \lambda_i^N \psi_i^N + \Sigma_N \Delta \psi_i^N + \Delta \Sigma_N \psi_i^N \cong \\ & \cong \lambda_i^N \psi_i^N + \lambda_i^N \Delta \psi_i^N + \Delta \lambda_i^N \psi_i^N \end{aligned} \quad (20)$$

hence,

$$\begin{aligned} & (\psi_i^N)^T \Sigma_N \Delta \psi_i^N + (\psi_i^N)^T \Delta \Sigma_N \psi_i^N \cong \\ & \cong \lambda_i^N (\psi_i^N)^T \Delta \psi_i^N + \Delta \lambda_i^N \|\psi_i^N\|^2 \end{aligned} \quad (21)$$

Using $\lambda_i^N (\psi_i^N)^T = (\psi_i^N)^T \Sigma_N$ we obtain ,

$$\begin{aligned} & \lambda_i^N (\psi_i^N)^T \Delta \psi_i^N + (\psi_i^N)^T \Delta \Sigma_N \psi_i^N \cong \\ & \cong \lambda_i^N (\psi_i^N)^T \Delta \psi_i^N + \Delta \lambda_i^N \end{aligned} \quad (22)$$

hence $\Delta \lambda_i^N = (\psi_i^N)^T \Delta \Sigma_N \psi_i^N$ that is,

$$\lambda_i^{N+1} = \lambda_i^N + (\psi_i^N)^T \Delta \Sigma_N \psi_i^N = (\psi_i^N)^T \Sigma_{N+1} \psi_i^N \quad (23)$$

The first order approximations of the orthonormal eigen vectors of Σ_{N+1} can be derived using the

expansion of each vector $\Delta \psi_i^N$ in the basis represented by the orthonormal eigen vectors of Σ_N ,

$$\Delta \psi_i^N = \sum_{j=1}^n b_{i,j} \psi_j^N, \quad (24)$$

where

$$b_{i,j} = (\psi_j^N)^T \Delta \psi_i^N. \quad (25)$$

Using the orthonormality, we get,

$$\begin{aligned} 1 & = \|\psi_i^N + \Delta \psi_i^N\|^2 \cong \|\psi_i^N\|^2 + 2(\psi_i^N)^T (\Delta \psi_i^N) = \\ & = 1 + 2(\psi_i^N)^T (\Delta \psi_i^N), \end{aligned} \quad (26)$$

that is

$$b_{i,i} = (\psi_i^N)^T \Delta \psi_i^N = 0 \quad (27)$$

Using (19), the approximation,

$$\Sigma_N \Delta \psi_i^N + \Delta \Sigma_N \psi_i^N \cong \lambda_i^N \Delta \psi_i^N + \Delta \lambda_i^N \psi_i^N. \quad (28)$$

holds for each $1 \leq i \leq n$.

For $1 \leq j \neq i \leq n$, from (28) we obtain the following equations,

$$\begin{aligned} & (\psi_j^N)^T \Sigma_N \Delta \psi_i^N + (\psi_j^N)^T \Delta \Sigma_N \psi_i^N \cong \\ & \cong \lambda_i^N (\psi_j^N)^T \Delta \psi_i^N + \Delta \lambda_i^N (\psi_j^N)^T \psi_i^N \end{aligned} \quad (29)$$

$$\begin{aligned} & (\psi_j^N)^T \Sigma_N \Delta \psi_i^N + (\psi_j^N)^T \Delta \Sigma_N \psi_i^N \cong \\ & \cong \lambda_i^N (\psi_j^N)^T \Delta \psi_i^N \end{aligned} \quad (30)$$

$$\begin{aligned} & \lambda_j^N (\psi_j^N)^T \Delta \psi_i^N + (\psi_j^N)^T \Delta \Sigma_N \psi_i^N \cong \\ & \cong \lambda_i^N (\psi_j^N)^T \Delta \psi_i^N \end{aligned} \quad (31)$$

From (31) we get,

$$(\lambda_i^N - \lambda_j^N) (\psi_j^N)^T \Delta \psi_i^N = (\psi_j^N)^T \Delta \Sigma_N \psi_i^N \quad (32)$$

$$b_{i,j} = (\psi_j^N)^T \Delta \psi_i^N = \frac{(\psi_j^N)^T \Delta \Sigma_N \psi_i^N}{\lambda_i^N - \lambda_j^N} \quad (33)$$

Consequently, the first order approximation of the eigen vectors of Σ_{N+1} are,

$$\psi_i^N + \Delta \psi_i^N = \psi_i^N + \sum_{j=1, j \neq i}^n \frac{(\psi_j^N)^T \Delta \Sigma_N \psi_i^N}{\lambda_i^N - \lambda_j^N} \psi_j^N \quad (34)$$

The skeleton of C is represented by the set of estimated principal components $\psi_1^N, \dots, \psi_n^N$. When the example X_{N+1} is included in C , then the new skeleton is $\psi_1^{N+1}, \dots, \psi_n^{N+1}$. The skeleton disturbance induced by the decision that X_{N+1} has to be allotted to C is measured by

$$D = \frac{1}{n} \sum_{k=1}^n d(\psi_k^N, \psi_k^{N+1}) \quad (35)$$

The crisp classification procedure identifies for each example the closest cluster in terms of the measure (35). Let $H = \{C_1, C_2, \dots, C_M\}$. In order to protect against misclassifications of samples coming from new classes not belonging to H , a threshold $T > 0$ is imposed, that is the example X_{N+1} is allotted to one of C_j for which

$$D = \frac{1}{n} \sum_{k=1}^n d(\psi_{k,j}^N, \psi_{k,j}^{N+1}) = \min_{1 \leq p \leq M} \frac{1}{n} \sum_{k=1}^n d(\psi_{k,p}^N, \psi_{k,p}^{N+1}) \quad (36)$$

and $D < T$, where the skeleton of C_j is $\psi_{1,j}^N, \dots, \psi_{n,j}^N$.

The classification of samples for which the resulted value of D is larger than T is postponed and the samples are kept in a new possible class CR. The reclassification of elements of CR is then performed followed by the decision concerning to either reconfigure the class system or to add CR as a new class in H .

In case of fuzzy classification, the value of the membership degree of X_{N+1} to each cluster of $HT = H \cup \{CR\}$ is computed as follows. Let

$$d(X_{N+1}, C_i) = \frac{1}{n} \sum_{k=1}^n d(\psi_{k,i}^N, \psi_{k,i}^{N+1}), \quad 1 \leq i \leq M \quad \text{and}$$

$$d(X_{N+1}, CR) = \begin{cases} \frac{1}{n} \sum_{k=1}^n d(\psi_{k,R}^N, \psi_{k,R}^{N+1}), & \text{if } CR \neq \emptyset \\ \infty, & \text{otherwise} \end{cases},$$

where the skeleton of CR is represented by the set of estimated principal components $\psi_{1,R}^N, \dots, \psi_{n,R}^N$.

$$\mu_{C_i}(X_{N+1}) = 1 - \frac{d(X_{N+1}, C_i)}{S}, \quad 1 \leq i \leq M \quad (37)$$

$$\mu_{CR}(X_{N+1}) = \begin{cases} 1 - \frac{d(X_{N+1}, CR)}{S}, & \text{if } d(X_{N+1}, CR) \neq \infty \\ 0, & \text{otherwise} \end{cases} \quad (38),$$

where

$$S = \begin{cases} \sum_{C \in HT} d(X_{N+1}, C), & \text{if } d(X_{N+1}, CR) \neq \infty \\ \sum_{C \in H} d(X_{N+1}, C), & \text{if } d(X_{N+1}, CR) = \infty \end{cases} \quad (39)$$

5 EXPERIMENTAL RESULTS AND CONCLUDING REMARKS

Several tests were performed on simulated data and they pointed out very successful performance of the proposed classification strategy.

A series of tests were performed on 4-dimensional simulated data coming from 5 classes each of them having 50 examples. The classes are represented by normal repartitions $C_i \sim (\mu_i, \Sigma_i)$, $1 \leq i \leq 5$, where

$$\mu_1 = [10 \ 11 \ 2 \ -12]$$

$$\Sigma_1 = \begin{bmatrix} 3.5944 & 2.0100 & 1.4720 & 0.6460 \\ 2.0100 & 3.7500 & 3.1000 & 1.5350 \\ 1.4720 & 3.1000 & 3.0600 & 0.8010 \\ 0.6460 & 1.5350 & 0.8010 & 1.8369 \end{bmatrix}$$

$$\mu_2 = [12 \ -5 \ 8 \ 13]$$

$$\Sigma_2 = \begin{bmatrix} 1.5566 & 0.7755 & 0.5230 & 0.3745 \\ 0.7755 & 1.7766 & 0.9305 & 0.6720 \\ 0.5230 & 0.9305 & 2.4586 & 1.0050 \\ 0.3745 & 0.6720 & 1.0050 & 2.3961 \end{bmatrix}$$

$$\mu_3 = [-10 \ 0 \ 9 \ 11]$$

$$\Sigma_3 = \begin{bmatrix} 1.7300 & 1.3740 & 0.0200 & 0.6000 \\ 1.3740 & 2.5144 & 0.0120 & 0.4600 \\ 0.0200 & 0.0120 & 3.0725 & 0.3250 \\ 0.6000 & 0.4600 & 0.3250 & 2.3144 \end{bmatrix}$$

$$\mu_4 = [-3 \ 14 \ 3 \ -11.5]$$

$$\Sigma_4 = \begin{bmatrix} 2.3618 & 0.5825 & 0.3814 & 0.9805 \\ 0.5825 & 2.4038 & 0.9920 & 0.6029 \\ 0.3814 & 0.9920 & 2.3724 & 0.4250 \\ 0.9805 & 0.6029 & 0.4250 & 4.1054 \end{bmatrix}$$

$$\mu_5 = [-7 \ -10.5 \ -14 \ 11.5]$$

$$\Sigma_5 = \begin{bmatrix} 1.8017 & 0.2860 & 0.6330 & 0.6540 \\ 0.2860 & 2.4436 & 0.0240 & 0.0210 \\ 0.6330 & 0.0240 & 1.6301 & 0.3454 \\ 0.6540 & 0.0210 & 0.3454 & 2.9441 \end{bmatrix}$$

The classification criterion is: allot X_{N+1} to C_{j_i} if

$$D = \min_{1 \leq l \leq t} \frac{1}{m_{j_l}} \sum_{k=1}^{m_{j_l}} d(\psi_{j_l}^k, \psi_{j_l, N+1}^k) \quad (37)$$

In order to evaluate the generalization capacities, 100 new examples were generated for each distribution. The results are presented in Table 1.

Table 1: Results on new simulated examples.

Class	C ₁	C ₂	C ₃	C ₄	C ₅
Number of correct classified examples	100	100	96	99	100
Number of misclassified examples	0	0	4 allotted to C ₂	1 allotted to C ₁	0
The mean value of D in case of correct classifications	0.08	0.05	0.75	0.21	0.14
The maximum value of D in case of correct classified examples	0.41	0.19	1.85	0.55	0.53

The evaluation of the generalization capacities in case of examples coming from new classes was performed on 1000 samples generated from $N(\mu, \Sigma)$, where $\mu = [0 \ 11 \ -9 \ -9.5]$ and

$$\Sigma = \begin{bmatrix} 8.2725 & 3.1080 & 1.7925 & 1.3680 \\ 3.1080 & 6.8986 & 1.8390 & 2.5561 \\ 1.7925 & 1.8390 & 6.0422 & 1.6410 \\ 1.3680 & 2.5561 & 1.6410 & 5.2261 \end{bmatrix}$$

The admissibility criterion for allotting a sample to a certain class is given by the maximum value of D corresponding to correct classifications. The results showed that about 975 examples were classified in CR, that is the algorithm managed to detect the intruded examples.

REFERENCES

Al Sultan K.S., Selim,S.Z., 1993. Global Algorithm for Fuzzy Clustering Problem, *Patt.Recogn.* 26,1375-1361
 Bensaid,A., Hall,L.O.,Bezdek,J.C., Clarke,L.P., 1996 Partially supervised clustering for image segmentation, *Patt.Recog.*,29(5),859-871.
 Bezdek,J.C., Keller,J., Krisnapuram,R., Pal,N.K., 2005. *Fuzzy Models and Algorithms for Pattern Recognition and Image Processing*, Springer Verlag
 Bezdek, J.C., 1981, *Pattern Recognition with Fuzzy Objective Function Algorithms*, Plenum Press
 Bensaid, H., Bezdek,J.C., Clarke, L.P., Silbiger, M.L., Arrington, J.A., Murtagh, R.F., 1996. Validity-Guided (Re)clustering with applications to image segmentation, *IEEE Trans. Fuzzy Systems*, 4, 112-123

Clark,M., Hall,L.O., Goldgof,D.B., Clarke,L.P., Velthuisen,R.P., Silbiger,M.S., 1994. MRI Segmentation using Fuzzy Clustering Techniques, *IEEE Engineering in Medicine and Biology*, 1994
 Everitt, B. S., 1978. *Graphical Techniques for Multivariate Data*, North Holland, NY
 Gath,J., Geva,A.B., 1989. Unsupervised optimal fuzzy clustering, *IEEE Trans. Pattern.Anal.Machine Intell.*, 11, 773-781
 A. Hyvarinen, J. Karhunen, E. Oja, 2001. *Independent Component Analysis*, John Wiley & Sons
 Huang,C., Shi,Y. 2002. *Towards Efficient Fuzzy Information Processing*, Physica-Verlag, Heidelberg
 Jain,A.K., Dubes,R., 1988. *Algorithms for Clustering Data*, Prentice Hall,Englewood Cliffs, NJ.
 Jin,Y., 2003. *Advanced Fuzzy Systems Design and Applications*, Physica-Verlag, Heidelberg
 Krishnapuram, R., Keller,J.M., 1993. A possibilistic approach to clustering, *IEEE Trans. Fuzzy Syst.*, 1(2)
 Li,C., Goldgof,D.B., Hall, L.O., 1993. Automatic segmentation and tissue labelling of MR brain images, *IEEE Transactions on Medical Imaging*, 12(4),1033-1048
 Pal,N.,R., Bezdek,J.C., 1995. On Cluster validity for the Fuzzy c-Means Model, *IEEE Trans. On Fuzzy Syst.*, Vol. 3,no.3
 Smith,S.P., Jain,A.K., 1984. Testing for uniformity in multidimensional data, *IEEE Trans. Patt. Anal. and Machine Intell.*, 6(1),73-81
 Wu,K-L., Yang,M-S, 2005. A Cluster validity index for fuzzy clustering, *Patt.Recog.Lett.* 26, 1275-1291
 Zahid,N., Abouelala,O., Limouri,M., Essaid,A., 1999. Unsupervised fuzzy clustering, *Patt.Recog.Lett.*, 20,1