

## ASYMPTOTICAL OPTIMALITY IN CLUSTER ANALYSIS

YURIJ S. KHARIN\* AND EUGENE E. ZHUK

*Department of Mathematical Modeling and Data Analysis, Belarusian State University, 4 Fr.Skariny av.,  
Minsk 220050, Belarus*

### SUMMARY

The problem of optimality and performance evaluation for cluster analysis procedures is investigated. For the situations where the classes are described by known or unknown prior probabilities and regular probability density functions with unknown parameters the asymptotic expansions of classification error probability are constructed. The results are illustrated for the case of well-known Fisher classification model. Copyright © 1999 John Wiley & Sons, Ltd.

KEY WORDS: cluster analysis; asymptotical optimality; 'plug-in' decision rule; classification error probability; asymptotic expansions

### I. INTRODUCTION: MATHEMATICAL MODEL

The problem of performance evaluation and optimality analysis for statistical classification decision rules remains one of the most important problems in statistical cluster analysis theory and its applications. This fact is caused not only by actual demands from practice but also by a special dominant role of this problem detected by Glick:<sup>1</sup> 'the task of estimating probabilities of correct classification confronts the statistician simultaneously with difficult distribution theory, questions intertwining sample size and dimension, problems of bias, variance, robustness and computation cost. But coping with such conflicting concerns enhances understanding of many aspects of statistical classification'.

Majority of results of performance evaluation and optimality analysis for statistical decision rules are concerned with the situation, where the training sample used for decision rule construction is classified and decision rules are being constructed by discriminant analysis methods. Review of these results can be found in the monographs of Lachenbruch,<sup>2</sup> Raudys<sup>3</sup> and McLachlan.<sup>4</sup>

In applied classification problems with poor statistical data the training samples are often unclassified, and methods of cluster analysis (self-training classification) are used for classification of training samples and also for classification of new-registered observations.

The main method of optimality analysis used up to now for these situations is Monte-Carlo simulation. Some analytical results are discussed in monographs of Bock<sup>5</sup> and Kharin.<sup>6</sup> This paper is devoted to optimality analysis for self-training statistical decision rules by asymptotic expansion method.

\* Correspondence to Y. S. Kharin, Department of Mathematical Modeling and Data Analysis, Belarusian State University, 4 Fr.Skariny av., Minsk 220050, Belarus.

Define now the mathematical model. Let

$$Q = \{q(x; \theta_*), x \in R^N : \theta_* \in \Theta \subseteq R^m\} \quad (1)$$

be a parametric family of probability density functions in  $R^N$  and  $L \geq 2$  different parameter values  $\{\theta_1^0, \dots, \theta_L^0\}$  be fixed. The set  $\{\theta_1^0, \dots, \theta_L^0\}$  determines  $L$  classes  $\{\Omega_1, \dots, \Omega_L\}$ . The class  $\Omega_i$  is described by the probability density function (p.d.f.)  $q(\cdot; \theta_i^0)$ ,  $i \in S = \{1, \dots, L\}$ .

Let a sample of  $n$  jointly independent random observations  $x_1, \dots, x_n$  from the classes  $\{\Omega_i\}_{i \in S}$  be registered in  $R^N$ . Introduce the notations:  $d_i^0 \in S$  is an unknown random index of the class to which the observation  $x_i$  belongs:

$$P\{d_i^0 = i\} = \pi_i^0, \quad i \in S \quad (2)$$

where  $\{\pi_i^0\}_{i \in S}$  ( $\sum_{i \in S} \pi_i^0 = 1$ ) are prior class probabilities;  $D^0 = (d_1^0, \dots, d_n^0)^T$  is the vector of true classification indices for the sample  $X = (x_1^T \dots x_n^T)^T$ , where  $T$  is the transposition symbol. The parameters  $\{\theta_i^0\}_{i \in S}$  and the prior probabilities  $\{\pi_i^0\}_{i \in S}$  are usually unknown. The cluster analysis problem consists in construction of a decision rule (DR)

$$d = d(x, X) : R^N \times R^{nN} \rightarrow S \quad (3)$$

for classifying a random observation  $x \in R^N$ . The DR (3) allows to calculate the estimate  $\hat{D} = (\hat{d}_1, \dots, \hat{d}_n)^T$  for the true classification vector of the sample  $X$ :  $\hat{d}_t = d(x_t, X)$ ,  $t = 1, \dots, n$ , and also to classify the new-registered observations  $x_{n+1}, x_{n+2}, \dots \in R^N$ .

## 2. 'PLUG-IN' DECISION RULES AND OPTIMALITY MEASURES

Note that if all class characteristics  $\{\pi_i^0, \theta_i^0\}_{i \in S}$  are *a priori* known, then the Bayesian DR (BDR)

$$d = d(x; \pi^0, \theta^0) = \arg \max_{i \in S} \{\pi_i^0 q(x; \theta_i^0)\} \quad (4)$$

where  $\pi^0 = (\pi_1^0, \dots, \pi_L^0)^T$ ,  $\theta^0 = (\theta_1^0, \dots, \theta_L^0)^T$ , classifies a random observation  $x$  from the class  $\Omega_{d^0}$  with the minimal value of classification error probability:

$$r_0 = P\{d \neq d^0\} = 1 - \int_{R^N} \max_{i \in S} \{\pi_i^0 q(x; \theta_i^0)\} dx \quad (5)$$

Now consider the situation, where  $\{\theta_i^0\}_{i \in S}$  (case *A*) or  $\{\pi_i^0, \theta_i^0\}_{i \in S}$  (case *B*) are unknown. Denote by  $\hat{\theta}$  and  $\{\hat{\pi}, \hat{\theta}\}$  any statistical estimators for  $\pi^0, \theta^0$  obtained by the sample  $X$  for cases *A* and *B*, respectively. Instead of the BDR (4) let us use the 'plug-in' DR (PDR):

$$d_A(x, X) = d(x; \pi^0, \hat{\theta}), \quad d_B(x, X) = d(x; \hat{\pi}, \hat{\theta}) \quad (6)$$

which are derived from the BDR  $d(x; \pi^0, \theta^0)$  by substitution of statistical estimators for unknown values of class characteristics.

As the optimality measures of PDR  $d_j(\cdot)$ ,  $j \in \{A, B\}$ , as in References 6 and 7, let us define:

(a) the classification error probability (CEP):

$$r_j(n) = P\{d_j(x, X) \neq d^0\}, \quad j \in \{A, B\}. \quad (7)$$

which means error probability at classifying a new-registered random observation  $x$  ( $x$  is independent on the random sample  $X$ );

(b) the relative bias of the CEP:

$$\kappa_j(n) = \frac{r_j(n) - r_0}{r_0} > 0, \quad j \in \{A, B\} \quad (8)$$

where  $r_0 > 0$  is defined in (5); the smaller  $\kappa_j(n)$  is, the more effective PDR  $d_j(\cdot)$  is ( $j \in \{A, B\}$ ). If for increasing sample size  $n$

$$\lim_{n \rightarrow \infty} \kappa_j(n) = 0 \quad (9)$$

then the PDR  $d_j(\cdot, X)$  is asymptotically optimal;

(c) the  $\delta$ -admissible sample size:

$$n_j^* = n_j^*(\delta) = \min\{n : \kappa_j(n) \leq \delta\}, \quad j \in \{A, B\} \quad (10)$$

where  $\delta > 0$  is any predetermined value; if  $n \geq n_j^*(\delta)$ , then the inequality  $\kappa_j(n) \leq \delta$  holds ( $j \in \{A, B\}$ ).

Note that the optimality measures (7)–(10) make possible to investigate symptotical optimality of PDR (6) with respect to the sample size  $n$  by means of risk asymptotic expansion method, which was proposed by Kharin.<sup>6</sup>

### 3. OPTIMALITY OF 'PLUG-IN' DECISION RULES BASED ON *ML*-ESTIMATES

Let in  $R^N$  with prior probabilities  $\pi_1^0, \pi_2^0 = 1 - \pi_1^0$  random observations from  $L = 2$  classes  $\Omega_1, \Omega_2$  are being registered (note that the results of this section also can be generalized for  $L > 2$  classes).

The training sample  $X = (x_1^T; \dots; x_n^T)^T$  is a random sample of size  $n$  from the mixture of two distributions:

$$q_{\pi^0}(x; \theta^0) = \pi_1^0 q(x; \theta_1^0) + (1 - \pi_1^0) q(x; \theta_2^0) \quad (11)$$

Because  $X$  is a unclassified sample and we want to avoid non-uniqueness, derived by denotations of classes, we assume, that  $\theta_2^0 \succ \theta_1^0$  (here  $\succ$  is the symbol of lexicographic comparison). In this model  $q_{\pi^0}(\cdot; \theta^0)$  is an element of the family of mixtures:

$$Q_{\pi^0} = \{q_{\pi^0}(x; \theta), x \in R^N; \theta = (\theta_1^T; \theta_2^T)^T \in \Theta^2 \subseteq R^{2m}, \theta_2 \succ \theta_1\} \quad (12)$$

At first, let us investigate the case  $A$  (the prior probability  $\pi_1^0$  of the class  $\Omega_1$  is assumed to be known). As an estimator  $\tilde{\theta}$  we use the maximum likelihood estimator (MLE):

$$\tilde{\theta} = \arg \max_{\theta} \sum_{i=1}^n \ln q_{\pi^0}(x_i; \theta) \quad (13)$$

Note that to solve the multiextremum problem (13) the different numerical methods may be used (see, for example, References 8 and 9).

Assume that the family of p.d.f.s (1) satisfies the following regularity conditions:

(C<sub>1</sub>)  $\theta_*$  is an identifiable parameter of p.d.f.  $q(\cdot; \theta_*)$ :

$$E_{\theta_*} \{\ln q(x; \theta_*)\} > E_{\theta_*} \{\ln q(x; \theta_{**})\}, \quad \theta_* \neq \theta_{**}$$

where for any function  $f(x; \theta_{**})$

$$E_{\theta_*} \{f(x; \theta_{**})\} = \int_{R^N} f(x; \theta_{**}) q(x; \theta_*) dx$$

(C<sub>2</sub>) For any compact  $K \subset \Theta$  and any points  $\theta_1^0, \theta_2^0 \in K$  some neighbourhoods  $U_{\theta_1^0}, U_{\theta_2^0} \subset K$  exist, such that for some  $a, c > 1$ ,  $b > 2$ , any neighbourhood  $U \subset U_{\theta_1^0}$  and any  $\theta_1 \in U_{\theta_1^0}$ ,  $\theta_2 \in U_{\theta_2^0}$  the functions

$$|\ln q(x; \theta_k)|^a, \quad \left( \sup_{\theta' \in U} |\ln q(x; \theta')| \right)^a$$

$$\left| \frac{\partial^2 \ln q(x; \theta_k)}{\partial \theta_{ki} \partial \theta_{kj}} \right|^b, \quad \left| \frac{\partial \ln q(x; \theta_k)}{\partial \theta_{ki}} \cdot \frac{\partial \ln q(x; \theta_*)}{\partial \theta_{*j}} \right|^b$$

$$\left| \frac{\partial^3 \ln q(x; \theta_k)}{\partial \theta_{ki} \partial \theta_{kj} \partial \theta_{kt}} \right|^c, \quad \left| \frac{\partial \ln q(x; \theta_s)}{\partial \theta_{st}} \cdot \frac{\partial^2 \ln q(x; \theta_k)}{\partial \theta_{ki} \partial \theta_{kj}} \right|^c$$

are uniformly integrable with respect to p.d.f.  $q(x; \theta_*)$ ,  $x \in R^N$  ( $\theta_* \in K$ ;  $k, s \in S = \{1, 2\}$ ;  $i, j, t = \overline{1, m}$ ).

(C<sub>3</sub>)  $E_{\theta_k^0} \{\nabla_{\theta_k^0} \ln q(x; \theta_k^0)\} = 0_m$ ,  $\theta_k^0 \in \Theta$ ,  
where  $0_m$  is  $m$ -vector, all elements of which are equal to 0;

(C<sub>4</sub>) Fisher information matrices

$$H_k = E_{\theta_k^0} \{-\nabla_{\theta_k^0}^2 \ln q(x; \theta_k^0)\}, \quad k \in S$$

$$J = J(\theta^0) = E_{\theta^0} \{-\nabla_{\theta^0}^2 \ln q_{\pi^0}(x; \theta^0)\}$$

where

$$E_{\theta^0} \{g(x)\} = \int_{R^N} g(x) q_{\pi^0}(x; \theta^0) dx$$

are positively defined, so that the minimal eigenvalues of these matrices are separated from zero.

Let us construct asymptotic expansion of the CEP  $r_A(n)$  at  $n \rightarrow +\infty$ . Denote by

$$\Gamma = \{x: G(x; \theta^0) = 0\} \subset R^N \quad (14)$$

the Bayesian discriminant hypersurface, where

$$G(x; \theta^0) = (1 - \pi_1^0) q(x; \theta_2^0) - \pi_1^0 q(x; \theta_1^0) \quad (15)$$

$$\alpha = \frac{1}{2} \int_{\Gamma} (\nabla_{\theta^0} G(x; \theta^0))^T J^{-1} \nabla_{\theta^0} G(x; \theta^0) |\nabla_x G(x; \theta^0)|^{-1} d\mathcal{S}_{N-1} \geq 0 \quad (16)$$

$\mathbf{1}(z) = \{1 \text{ if } z \geq 0; 0, \text{ if } z < 0\}$  is the unit function.

#### Theorem 1

If the conditions (C<sub>1</sub>)–(C<sub>4</sub>) are satisfied, p.d.f.s  $\{q(x; \theta_k^0)\}_{k \in S}$  have derivatives w.r.t.  $x$  ( $x \in R^N$ ,  $k \in S$ ), and surface integral (16) is finite:  $\alpha < +\infty$ , then the CEP of PDR

$$d_A(x, X) = \mathbf{1}(G(x; \tilde{\theta})) + 1$$

allows the asymptotic expansion:

$$r_A(n) = r_0 + \alpha n^{-1} + O(n^{-3/2}) \quad (17)$$

where

$$r_0 = 1 - \pi_1^0 - \int_{R^s} \mathbf{1}(G(x; \theta^0)) G(x; \theta^0) dx \quad (18)$$

*Proof.* Under the regularity conditions (C<sub>1</sub>)–(C<sub>4</sub>) using Chibisov<sup>10</sup> stochastic expansion for the random deviation  $\Delta\theta = \hat{\theta} - \theta^0$  of MLE (13) we find, that  $\Delta\theta$  has the third-order moments, and the following asymptotic expansions are true at  $\tau = 1/\sqrt{n} \rightarrow 0$ :

for the bias:

$$E\{\Delta\theta\} = \mathbf{1}_{2m} O(\tau^3) \quad (19)$$

for the covariance matrix of  $\hat{\theta}$ :

$$E\{\Delta\theta(\Delta\theta)^\top\} = J^{-1} \tau^2 + \mathbf{1}_{2m \times 2m} O(\tau^3) \quad (20)$$

for third-order moments ( $k, l, s \in S$ ;  $i, j, t = \overline{1, m}$ ):

$$E\{(\hat{\theta}_{ki} - \theta_{ki}^0)(\hat{\theta}_{lj} - \theta_{lj}^0)(\hat{\theta}_{st} - \theta_{st}^0)\} = O(\tau^3) \quad (21)$$

where  $\mathbf{1}_{2m}$ ,  $\mathbf{1}_{2m \times 2m}$  are  $(2m)$ -vector and  $(2m \times 2m)$ -matrix, all elements of which are equal to 1.

Now expression (17) is obtained by applying of Taylor formula to the CEP (7):

$$r_A(n) = 1 - \pi_1^0 - E \left\{ \int_{R^s} \mathbf{1}(G(x; \hat{\theta})) G(x; \theta^0) dx \right\}$$

in the neighbourhood of  $\theta^0$  with respect to  $\Delta\theta$ , and of the relations (19)–(21).  $\square$

### Corollary

Under the conditions of Theorem 1 the decision rule  $d_A(x, X)$  is asymptotically optimal:  $\kappa_A(n) \rightarrow 0$ , if  $n \rightarrow +\infty$ .

Now let us consider the situation, where both the parameter values  $\theta_1^0, \theta_2^0$  and the prior probability  $\pi_1^0$  are unknown. In this case the composed vector  $\eta^0 = (\theta_1^{0\top}; \theta_2^{0\top}; \pi_1^0)^\top$  of unknown parameters is  $(2m+1)$ -dimensional vector and its MLE  $\hat{\eta} = (\hat{\theta}_1^\top; \hat{\theta}_2^\top; \hat{\pi}_1)^\top$  has the form

$$\hat{\eta} = \arg \max_{\substack{\theta_1, \theta_2 \\ 0 < \pi_1 < 1}} \sum_{t=1}^n \ln q_\pi(x_t; \theta) \quad (22)$$

### Theorem 2

If the conditions of Theorem 1 hold, then

$$r_B(n) = r_0 + (\alpha + \beta + \gamma)n^{-1} + O(n^{-3/2}) \quad (23)$$

where

$$\beta = \frac{\pi_1^0}{2(1 - \pi_1^0)} \int_{\Gamma} (q(x; \theta_1^0))^2 |\nabla_x G(x; \theta^0)|^{-1} d\mathcal{S}_{N-1} > 0$$

and  $\gamma = O(r_0) > 0$ .

*Proof.* Is conducted analogously to the proof of Theorem 1.  $\square$

It is seen from Theorems 1 and 2 that

$$r_B(n) = r_A(n) + (\beta + \gamma)n^{-1} + O(n^{-3/2})$$

where the term  $(\beta + \gamma)n^{-1}$  is derived by random errors of estimation of the prior probability  $\pi_1^0$  in the case *B*.

Note that the DR  $d_B(x, X)$  also is the asymptotically optimal decision rule.

#### 4. OPTIMALITY OF DECISION RULES BASED ON LIKELIHOOD FUNCTION

Let us investigate the case of  $L \geq 2$  classes and define composite probability distribution of the random sample  $X = (x_1^T; \dots; x_n^T)^T \in R^{nN}$  and the random true classification vector  $D^0 = (d_1^0, \dots, d_n^0)^T \in S^n$ :

$$p(X, D^0 | \pi^0, \theta^0) = \prod_{i=1}^n \pi_{d_i^0}^0 q(x_i; \theta_{d_i^0}^0) \quad (24)$$

The logarithmic likelihood function, which corresponds to (24), has the form

$$l(\pi, \theta, D) = n^{-1} \ln p(X, D | \pi, \theta) = n^{-1} \sum_{i=1}^n \ln (\pi_{d_i} q(x_i; \theta_{d_i})) \quad (25)$$

As in Section 3, let us consider the following two cases.

(A) Prior probability vector  $\pi^0 = (\pi_1^0, \dots, \pi_L^0)^T$  is known. In this case we need to estimate the vectors  $D^0$  and  $\theta^0 = (\theta_1^0; \dots; \theta_L^0)^T$ :

$$l(\pi^0, \theta, D) \rightarrow \max_{\theta \in \Theta^L, D \in S^n} \quad (26)$$

(B) All class characteristics  $\pi^0$  and  $\theta^0$  are unknown, then we need to solve the following problem:

$$l(\pi, \theta, D) \rightarrow \max_{\substack{\theta \in \Theta^L, D \in S^n, \\ \pi: 0 < \pi_i < 1, \pi_1 + \dots + \pi_L = 1}} \quad (27)$$

Denote by

$$W_n(\pi, \theta) = \max_{D \in S^n} l(\pi, \theta, D) = n^{-1} \sum_{i=1}^n f(x_i; \pi, \theta) \quad (28)$$

the statistical estimator for the functional

$$W(\pi, \theta) = \int_{R^N} f(x; \pi, \theta) q_{\pi^0}(x; \theta^0) dx \quad (29)$$

where

$$f(x; \pi, \theta) = \max_{i \in S} \ln(\pi_i q(x; \theta_i))$$

$$q_{\pi^0}(x; \theta^0) = \sum_{i \in S} \pi_i^0 q(x; \theta_i^0)$$

The optimization problems (26) and (27) can be rewritten in the form of PDR:

$$d_A(x, X) = \arg \max_{i \in S} \{\pi_i^0 q(x; \hat{\theta}_i)\}, \quad \hat{\theta} = \arg \max_{\theta} W_n(\pi^0, \theta) \quad (30)$$

$$d_B(x, X) = \arg \max_{i \in S} \{\hat{\pi}_i q(x; \hat{\theta}_i)\}, \quad \{\hat{\pi}, \hat{\theta}\} = \arg \max_{\pi, \theta} W_n(\pi, \theta) \quad (31)$$

Note, that the PDR (30) and (31) are formal notations of the problems (26) and (27). Unlike Section 3 the processes of unknown class characteristic estimation and classification of the sample  $X$  take place simultaneously. But expressions (30) and (31) are more convenient for investigation than expressions (26) and (27).

As in Section 3, to avoid non-uniqueness assume that the components of the estimates  $\hat{\pi}, \hat{\theta}$  and  $\tilde{\theta}$  are lexicographically ordered and correspond to the true values  $\pi^0$  and  $\theta^0$ .

To investigate optimality of the procedures (26) and (27), let us construct the asymptotic expansions of the CEP for the PDR (30) and (31). Introduce the following notations:

$$\rho_{ij} = \max_{\substack{i \neq j \\ i, j \in S}} \int_{R^N} \sqrt{q(x; \theta_i^0) q(x; \theta_j^0)} dx \quad (32)$$

$$F_{ij}(x) = \pi_i^0 q(x; \theta_i^0) - \pi_j^0 q(x; \theta_j^0), \quad \Gamma_{ij} = \{x: F_{ij}(x) = 0\} \subset R^N (i \neq j \in S)$$

$$B_{kl ij} = \pi_k^0 \pi_l^0 \int_{\Gamma_{ij}} (\nabla_{\theta_k} q(x; \theta_k^0))^T I_{kl}(\theta^0) \nabla_{\theta_l} q(x; \theta_l^0) |\nabla_x F_{ij}(x)|^{-1} d\mathcal{S}_{N-1}$$

where  $I_{kl}(\theta^0)$  is the  $(k, l)$ th  $(m \times m)$ -block of the matrix  $I(\theta^0) = (I_{kl}(\theta^0))_{k, l \in S}$ :

$$I(\theta) = V^{-1}(\theta) \int_{R^N} \nabla_{\theta} f(x; \pi^0, \theta) (\nabla_{\theta} f(x; \pi^0, \theta))^T q_{\pi^0}(x; \theta^0) dx V^{-1}(\theta) \quad (33)$$

$$V(\theta) = \nabla_{\theta}^2 W(\pi^0, \theta)$$

For the case of two classes ( $L = 2$ ) and unknown  $\pi_1^0$  let us introduce

$$\tilde{z} = \frac{1}{2} \int_{\Gamma_{12}} (\nabla_{\eta^0} F(x; \eta^0))^T \tilde{I}(\eta^0) \nabla_{\eta^0} F(x; \eta^0) |\nabla_x F(x; \eta^0)|^{-1} d\mathcal{S}_{N-1} \quad (34)$$

where

$$F(x; \eta^0) = (1 - \pi_1^0)q(x; \theta_2^0) - \pi_1^0 q(x; \theta_1^0)$$

$$\tilde{I}(\eta^0) = \tilde{V}^{-1}(\eta^0) \int_{R^N} \nabla_{\eta^0} f(x; \pi^0, \theta^0) (\nabla_{\eta^0} f(x; \pi^0, \theta^0))^T q_{\pi^0}(x; \theta^0) dx \tilde{V}^{-1}(\eta^0)$$

$$\tilde{V}(\eta^0) = \nabla_{\eta}^2 W(\pi, \theta)|_{\eta=\eta^0}, \quad \eta^0 = (\theta_1^0; \theta_2^0; \pi_1^0)^T$$

### Theorem 3

If under the regularity conditions (C<sub>1</sub>)–(C<sub>4</sub>) the following asymptotics takes place:

$$\rho_n \rightarrow 0, \quad n \rightarrow +\infty, \quad (35)$$

then the estimators  $\tilde{\theta}, \tilde{\pi}, \hat{\theta}$  from (30) and (31) are strongly consistent:

$$\tilde{\theta} \xrightarrow{P-1} \theta^0, \quad \tilde{\pi} \xrightarrow{P-1} \pi^0, \quad \hat{\theta} \xrightarrow{P-1} \theta^0 \quad (36)$$

If in addition p.d.f.  $q(x; \theta_*)$ ,  $\theta_* \in \Theta$ , is differentiable w.r.t.  $x \in R^N$ , the matrices  $V(\theta^0)$ ,  $\tilde{V}(\eta^0)$  are non-singular and surface integrals in (32) and (34) are finite, then the CEP functionals  $r_A(n), r_B(n)$  of the PDR (30) and (31) have the form

$$r_A(n) = r_0 + \frac{1}{2} \sum_{j=2}^L \sum_{l=1}^{j-1} (B_{jllj} + B_{lljj} - 2B_{ljjl}) n^{-1} + o(n^{-1}) \quad (37)$$

$$r_B(n) = r_0 + \tilde{\alpha} n^{-1} + o(n^{-1}) \quad (L=2) \quad (38)$$

where the CEP  $r_0$  of the BDR (4) is determined in (5).

*Proof.* The strong consistency of  $\tilde{\theta}, \tilde{\pi}$  and  $\hat{\theta}$  (expressions (36)) is proved by applying of the strong law of large numbers to (28):

$$W_n(\pi, \theta) \xrightarrow{P-1} W(\pi, \theta)$$

and the continuity theorems from Reference<sup>11</sup> under asymptotics (35). The asymptotic expansions (37) and (38) are constructed as in the proof of Theorem 1.  $\square$

### Corollary

The 'plug-in' decision rules (30) and (31) are asymptotically optimal. If according to the conditions of Theorem 3,

$$b_{lj} = \inf_{x \in R^N} \left| \nabla_x \ln \frac{q(x; \theta_l^0)}{q(x; \theta_j^0)} \right| \geq r, \quad r = O(1), \quad l \neq j \in S \quad (39)$$

then

$$r_k(n) = r_0 + o(n^{-1}), \quad k \in \{A, B\} \quad (40)$$



*Proof.* The asymptotic optimality of the PDR (30) and (31) follows from (37), (38) and (9). Under the conditions (C<sub>1</sub>)–C<sub>4</sub> and (39) by means of Cauchy-Schwarz inequality we obtain

$$\max\{|B_{ljlj}|, B_{llj}, B_{jjl}\} < p_1 \rho_+^{1/2}, \quad l \neq j \in S$$

$$\bar{x} < p_2 \rho_-^{1/2}$$

where  $p_1, p_2 < +\infty$  are some positive constants. From last inequalities and from the expressions (35), (36) and (38) we have (40).  $\square$

Note (see Reference 6) that the asymptotics (35) is of practical value, because under 'large overlapping of classes' (when (35) is disturbed) the value  $r_0$  is large and all DRs are not advisable for applications.

### 5. FISHER MODEL AND ILLUSTRATION OF OBTAINED RESULTS

Let us illustrate the obtained results for the case with two equiprobable classes ( $L=2$ ):  $\pi_1^0 = \pi_2^0 = 0.5$ , when the family (1) is Gaussian:

$$q(x; \theta_i^0) = n_N(x | \theta_i^0, \Sigma), \quad i \in S = \{1, 2\} \quad (41)$$

where

$$n_N(x | \theta_*, \Sigma) = (2\pi)^{-N/2} (\det(\Sigma))^{-0.5} \exp(-0.5(x - \theta_*)^T \Sigma^{-1} (x - \theta_*))$$

is  $N$ -variate Gaussian p.d.f. with mathematical mean vector  $\theta_*$  and non-singular covariance ( $N \times N$ )-matrix  $\Sigma$  ( $\det(\Sigma) > 0$ ). Relation (41) defines the so-called Fisher<sup>12</sup> model.

Denote by

$$\Delta = \sqrt{(\theta_1^0 - \theta_2^0)^T \Sigma^{-1} (\theta_1^0 - \theta_2^0)}$$

the Mahalanobis interclass distance. The CEP  $r_0$  of the BDR (4) has the form

$$r_0 = \Phi\left(-\frac{\Delta}{2}\right)$$

where  $\Phi(\cdot)$  is the standard Gaussian distribution function with the p.d.f.:

$$\varphi(z) = \Phi'(z) = (2\pi)^{-1/2} \exp(-z^2/2), \quad z \in R$$

For the PDR based on the MLE (13) and (22) by the results of Section 3 we evaluated main terms in the CEP asymptotic expansions (17) and (23) and determined the  $\delta$ -admissible sample size (10):

$$n_A^*(\delta) = \left[ \frac{1}{2\delta} \left( N + \frac{\Delta^2}{4} + \frac{\Delta^3}{8} \sqrt{\frac{\pi}{2}} e^{-\Delta^2/8} \right) \right] + 1 \quad (42)$$

$$n_B^*(\delta) = n_A^*(\delta) + \left[ \frac{1}{\delta} \left( \frac{1}{2} + \frac{2}{\Delta^2} + \sqrt{\frac{\pi}{2}} \frac{1}{\Delta} e^{-\Delta^2/8} \right) \right] + 1$$

Table 1. Cluster analysis of Iris data

Sample size, $n$	Case A			Case B		
	Evaluated $\Delta$	$\hat{r}_A$	$\hat{\gamma}_n$	Evaluated $\Delta$	$\hat{r}_B$	$\hat{\gamma}_n$
10	1.713	0.248	0.200	1.692	0.285	0.300
20	1.913	0.193	0.150	1.889	0.209	0.200
30	1.909	0.186	0.167	1.872	0.199	0.200
40	1.882	0.185	0.175	1.882	0.192	0.175

where  $[z]$  means the entire part of the value  $z$ . If  $n = n_j^*(\delta)$ , then for the CEP  $r_j(n)$  the following approximation is true:

$$r_j(n) \approx \hat{r}_j, \quad \hat{r}_j = \hat{r}_j(n) = (1 + \delta)r_0, \quad j \in \{A, B\}. \quad (43)$$

In practice, this fact allows to estimate the risk (7) when the classification process has finished and when the Mahalanobis distance  $\Delta$  has been evaluated. As a numerical example let us use the well known Fisher<sup>12</sup> Iris data. Four samples ( $n = 10; 20; 30; 40; N = 4$ ) from two ( $L = 2$ ) equiprobable classes (Iris versicolor, Iris virginica) were used. The experimental results for the case A (prior class probabilities are assumed to be known:  $\pi_1^0 = \pi_2^0 = 0.5$ ) and for the case B ( $\{\pi_i^0\}_{i \in S}$  are assumed to be unknown) are presented in Table 1.

Here  $\hat{\gamma}_n$  is the experimental frequency of error decisions:

$$\hat{\gamma}_n = \frac{1}{n} \sum_{i=1}^n (1 - \delta_{d_i, d_i^0})$$

Note, in real situation, when the true classification vector  $D^0 = (d_1^0, \dots, d_n^0)^T$  is unknown, as the CEP estimates we recommend to use  $\hat{r}_j, j \in \{A, B\}$ .

Now let us illustrate the results obtained in Section 4 for the PDR based on the likelihood function. Suppose that the Mahalanobis distance  $\Delta$  is increasing:  $\Delta \rightarrow +\infty, n \rightarrow +\infty$ , and the asymptotics (35) takes place (because of  $\rho_- = \exp(-\Delta^2/8)$  under Fisher model (41)). All conditions of Corollary of Theorem 3 are satisfied:

$$b_{12} = \sqrt{(\theta_2^0 - \theta_1^0)^T \Sigma^{-1} \Sigma^{-1} (\theta_2^0 - \theta_1^0)} \geq v, \quad v = O(1)$$

and according to (40) the CEP  $r_A(n)$  and  $r_B(n)$  of the PDR (30) and (31) have no terms of order  $O(n^{-1})$  in their asymptotic expansions:

$$r_j(n) = r_0 + o(n^{-1}), \quad r_0 = \Phi\left(-\frac{\Delta}{2}\right), \quad j \in \{A, B\}$$

#### ACKNOWLEDGEMENTS

These investigations were partially supported by Belarussian National Grants F40-263, MP94-03.

#### REFERENCES

1. N. Glick, 'Additive estimators for probabilities of correct classification', *Pattern Recognition*, **1**, 211-222 (1978).
2. P. Lachenbruch, *Discriminant Analysis*, Hafner Press, New York, 1975.

3. Sh. Raudys, 'Small sample effects in classification problems', in *Statistical Problems of Control*, Vol. 18, IMC, Vilnius, 1976 (in Russian).
4. G. J. McLachlan, *Discriminant Analysis and Statistical Pattern Recognition*, Wiley, New York, 1992.
5. H. H. Bock, 'Probabilistic aspects in cluster analysis', *Proc. 13th Conf. of the Gesellschaft für Klassifikation*, Springer, Berlin, 1989.
6. Yu. S. Kharin, *Robustness in Statistical Pattern Recognition*, Kluwer Academic Publishers, Dordrecht, 1996.
7. Yu. S. Kharin and E. E. Zhuk, 'Robustness in statistical pattern recognition under "contaminations" of training samples', in *Proc. 12th IAPR Int. Conf. on Pattern Recognition*, Vol. 2, IEEE Comp. Press, Washington, DC, 1994.
8. J. E. Dennis, Jr., and R. B. Schnabel, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Prentice-Hall, Englewood cliffs, NJ, 1983.
9. L. Redner and J. Walker, 'Mixture densities, maximum likelihood and the EM algorithm', *SIAM Rev.* **26** (2), (1984).
10. D. M. Chibisov, 'Asymptotic expansion for some estimate family including *ML*-estimates', *Teor. Ver. Prim.*, **18** (4), 689–701 (1973).
11. A. A. Borovkov, *Mathematical Statistics: Parameters Estimation and Hypotheses Testing*, Nauka, Moscow, 1984 (in Russian).
12. R. A. Fisher, 'The use of multiple measurements in taxonomic problems', *Ann. Eugen.*, **7**, 179–188 (1936).