

УДК 519.25

DISCRIMINANT FUNCTIONS QUALITY ESTIMATION ON THE BASIS OF TRAINING AND TESTING SAMPLES

Alexander Sarychev¹ and Lyudmyla Sarycheva²

¹ *Institute of Technical Mechanics of the National Academy of Sciences of Ukraine,
15 Leshko-Popel St., Dnipropetrovs'k, 49005, Ukraine*

² *National Mining University of Ukraine, 19 K. Marks Ave., Dnipropetrovs'k, 49027, Ukraine
Sarychev@prognoz.dp.ua, Sarycheval@nmu.org.ua*

Обґрунтовано спосіб порівняння дискримінантних функцій з розбиттям вибірок спостережень на навчальні й перевірні підвибірки. Отримано умови існування оптимальної множини ознак, які залежать від параметрів генеральних сукупностей і обсягів вибірок. Виявлено закономірності спрощення оптимальної дискримінантної функції при зменшенні обсягів вибірок і при збільшенні дисперсій ознак.

Ключові слова: метод групового урахування аргументів, невизначеність за складом ознак, критерій якості лінійної дискримінантної функції.

The way of comparison of discriminant functions with dividing samples observations on training and testing subsamples is proved. Conditions of existence of optimum set of features which depend on parameters of general sets and volumes samples are received. Laws of simplification of optimum discriminant function at decrease of volumes samples and at increase of dispersions of features are revealed.

Keywords: Group Method of Data Handling, uncertainty on structure of features, criterion of linear discriminant function quality.

Обоснован способ сравнения дискриминантных функций с разбиением выборок наблюдений на обучающие и проверочные подвыборки. Получены условия существования оптимального множества признаков, которые зависят от параметров генеральных совокупностей и объемов выборок. Выявлены закономерности упрощения оптимальной дискриминантной функции при уменьшении объемов выборок и при увеличении дисперсий признаков.

Ключевые слова: метод группового учета аргументов, неопределенность по составу признаков, критерий качества линейной дискриминантной функции.

Introduction

The decision of task of the discriminant analysis in conditions of structural uncertainty on structure of features assumes acceptance of any way of comparison of discriminant functions which are constructed on various sets of features. Two ways of comparison are popular in practice. The first way is based on dividing of observations on training and testing subsamples. In this way training subsamples are used for estimation coefficients of discriminant functions, and testing subsamples are used for estimation its qualities of classification. The second way is sliding examination. In this way observations which are serially excluded from training subsamples are used as testing observations. In the literature these ways are traditionally treated as heuristic methods though the fact of existence in them of optimum set of features repeatedly proved by a method of statistical tests. In the Group Method of Data Handling (GMDH) analytical research of these two ways is carried out [1-4]. For the decision of a task of the discriminant analysis in conditions of structural uncertainty

except for a way of comparison discriminant functions it is required to specify algorithm of generation of various combinations of the features included in discriminant functions. It is supposed, that as such method is chosen the complete sorting-out of all possible combinations of features.

1. Way of comparison of discriminant functions on the basis of training and testing subsamples

Suppose that at the step with number s ($s = 1, 2, \dots, m$) of algorithm complete sorting-out of all possible sets of features only s components from the set X can be included in the discriminant function and these features form the current set V . In the following we suppose that \mathbf{V}_I and \mathbf{V}_{II} are $(s \times n_I)$ and $(s \times n_{II})$ matrices of observations from general sets P_I and P_{II} , \mathbf{v}_I and \mathbf{v}_{II} are s -dimensional column vectors of the mathematical expectations in the sets P_I and P_{II} , Σ_V is covariance $(s \times s)$ matrix of the sets P_I and P_{II} .

Let's consider the estimation of Mahalanobis distance that is constructed with account of dividing of observations on training and testing subsamples. We shall calculate estimations of coefficients discriminant function for set the component V on the training subsample A and it is used them for estimation Mahalanobis distances as the relation of an intergroup variation to an intragroup variation on testing subsample B :

$$D_{AB}^2(V) = \frac{\hat{\mathbf{d}}_A^T (\tilde{\mathbf{v}}_{IB} - \tilde{\mathbf{v}}_{IIB}) (\tilde{\mathbf{v}}_{IB} - \tilde{\mathbf{v}}_{IIB})^T \hat{\mathbf{d}}_A}{\hat{\mathbf{d}}_A^T \mathbf{S}_B \hat{\mathbf{d}}_A} \tag{1}$$

In formula (1), vector $\hat{\mathbf{d}}_A$ is an estimate of the coefficients of the Fisher function that is calculate on training subsamples A

$$\hat{\mathbf{d}}_A = \mathbf{S}_A^{-1} (\tilde{\mathbf{v}}_{IA} - \tilde{\mathbf{v}}_{IIA}), \tag{2}$$

where $\tilde{\mathbf{v}}_{IA}$ and $\tilde{\mathbf{v}}_{IIA}$ are estimate of the mathematical expectation \mathbf{v}_I and \mathbf{v}_{II} :

$$\tilde{\mathbf{v}}_{kA} = (n_{kA})^{-1} \sum_{i=1}^{n_{kA}} \mathbf{V}_{kiA}, \quad k = I, II; \tag{3}$$

the matrix \mathbf{S}_A is an unbiased estimate of covariance matrix Σ_V

$$\mathbf{S}_A = (n_{IA} - n_{IIA} - 2)^{-1} [\mathbf{v}_{IA} \mathbf{v}_{IA}^T + \mathbf{v}_{IIA} \mathbf{v}_{IIA}^T]. \tag{4}$$

In formula (4) \mathbf{v}_{kA} ($k = I, II$) are matrices of deviations of observations \mathbf{V}_{kA} from estimate $\tilde{\mathbf{v}}_{kA}$

$$\mathbf{v}_{kA} = [\mathbf{V}_{k1A} - \tilde{\mathbf{v}}_{kA}, \mathbf{V}_{k2A} - \tilde{\mathbf{v}}_{kA}, \dots, \mathbf{V}_{kn_kA} - \tilde{\mathbf{v}}_{kA}]. \tag{5}$$

In formula (5) vectors $\tilde{\mathbf{v}}_{IB}$ and $\tilde{\mathbf{v}}_{IIB}$ calculated analogues (3), and matrix \mathbf{S}_B calculated analogues (4)–(5); n_{IA} and n_{IIA} , n_{IB} and n_{IIB} are volume of training and testing subsamples respectively, and it is true $n_{IA} + n_{IB} = n_I$ and $n_{IIA} + n_{IIB} = n_{II}$. Using (2), we obtain for $D_{AB}^2(V)$:

$$D_{AB}^2(V) = \frac{(\tilde{\mathbf{v}}_{IA} - \tilde{\mathbf{v}}_{IIA})^T \mathbf{S}_A^{-1} (\tilde{\mathbf{v}}_{IB} - \tilde{\mathbf{v}}_{IIB}) (\tilde{\mathbf{v}}_{IB} - \tilde{\mathbf{v}}_{IIB})^T \mathbf{S}_A^{-1} (\tilde{\mathbf{v}}_{IA} - \tilde{\mathbf{v}}_{IIA})}{(\tilde{\mathbf{v}}_{IA} - \tilde{\mathbf{v}}_{IIA})^T \mathbf{S}_A^{-1} \mathbf{S}_B \mathbf{S}_A^{-1} (\tilde{\mathbf{v}}_{IA} - \tilde{\mathbf{v}}_{IIA})}. \quad (6)$$

Let $\tau_V^2 = (\mathbf{v}_I - \mathbf{v}_{II})^T \Sigma_V^{-1} (\mathbf{v}_I - \mathbf{v}_{II})$ be the Mahalanobis distance for the set V , $r = n_{IA} + n_{IIA} - 2$, $c_A^{-1} = (n_{IA}^{-1} + n_{IIA}^{-1})$, $c_B^{-1} = (n_{IB}^{-1} + n_{IIB}^{-1})$.

Theorem. For mathematical expectation of random variable $D_{AB}^2(V)$, we have

$$E\{D_{AB}^2(V)\} = \left(\tau_V^2 - \frac{\tau_V^2 [s - (r-1)/(r-s)] c_A^{-1}}{\tau_V^2 + s c_A^{-1}} + c_B^{-1} \frac{r-1}{r-s} \right) \frac{r-s}{r-1}. \quad (7)$$

The validity of theorem follows from the validity of the following: 1) the estimates obtained on subsamples A and B are independent; 2) the estimate (3) and estimate (4) are independent; 3) matrix \mathbf{S}_A is random ($s \times s$) matrix which has the Wishart distribution with r degrees of freedom.

Definition 1. The optimal set components (set features) is defined as the set V_{OPT} for which

$$V_{OPT} = \arg \max_{V \subseteq X} E\{D_{AB}^2(V)\}. \quad (8)$$

Definition 2. Optimal discriminant function with respect to the number and composition of the components is defined as the Fisher discriminant function constructed on the set of components V_{OPT} .

We proved that optimal set of components exist and formulated the conditions under which the optimal discriminant function is simplified in number of the features included in it. For this purpose, it was investigated $E\{D_{AB}^2(V)\}$ depending on composition of set V .

It is possible to divide set of components X into the following nonintersecting subsets $X = \overset{\circ}{X} \cup \overset{\circ}{R} \cup \tilde{R} = \overset{\circ}{V} \cup \tilde{R}$: so that 1) $\overset{\circ}{X} \neq \emptyset$ (where \emptyset is the empty set) is the set of components whose mathematical expectation satisfy $\overset{\circ}{\chi}_{Ih} \neq \overset{\circ}{\chi}_{IIh}$, $h = 1, 2, \dots, m$, where m is their number; 2) \tilde{R} is the set of components whose mathematical

expectation satisfy $\overset{\circ}{\rho}_{Ih} = \overset{\circ}{\rho}_{IIIh}, h = 1, 2, \dots, \overset{\circ}{l}$, where $\overset{\circ}{l}$ is their number and each component in $\overset{\circ}{R}$ depends statistically on the least one components in the set $\overset{\circ}{X}$ (the set $\overset{\circ}{R}$ may be empty); 3) \tilde{R} is the set of components whose mathematical expectation satisfy $\tilde{\rho}_{Ih} = \tilde{\rho}_{IIIh}, h = 1, 2, \dots, \tilde{l}$, where \tilde{l} is number and each component each component in \tilde{R} is statistically independent from each set $\overset{\circ}{X}$ (the set \tilde{R} may be empty). Relationship between the Mahalanobis distance for the set components $\overset{\circ}{V} = \overset{\circ}{X} \cup \overset{\circ}{R}$ and the Mahalanobis distance for a current analyzed set of components $V \subseteq X$ is formulated in the form of lemmas [1-4].

In case of known parameters of general sets P_I and P_{II} it follows from the stated lemmas that: 1) every component from set $\overset{\circ}{X}$ is necessary in the sense that its inclusion into the current set of components V increase the Mahalanobis distance τ_V^2 ; 2) every component from the set $\overset{\circ}{R}$ is necessary in the sense that its inclusion into the current set of components V increase the Mahalanobis distance τ_V^2 ; 3) every components from the set \tilde{R} is redundant in the sense, that its inclusion into the current set V does not increase the Mahalanobis distance τ_V^2 .

2. Conditions of reduction (simplification) of optimum discriminant function

As a rule, in practical applications parameters of general populations are unknown; however they can be estimated as statistical estimates on training samples of observations of limited volume. It is known, that if we use constructed rule of classification to the training sample, then estimate of recognition quality will be overstated by mathematical expectation in comparison with the same evaluation of quality on data, independent of training data.

The way for comparison of the discriminant functions based on dividing of the initial data sample on training and testing subsamples give not overstated estimates of recognition quality. Experience of practical applications and test investigations of this way on basis of method of statistical test show that in this way: 1) on increase of size of observations samples increases the number of components in the set, on which the best quality of recognition is attained, and on decrease of size of observations samples the number of components in such set decreases; 2) on increase of the Mahalanobis distance τ_X^2 between general populations (from which observation samples were obtain) the number of components increases in the set, on which the best quality of recognitions is attained, and on decrease of this distance the number of components in such set decreases.

Our analytical investigations confirm these empirically determined regularities about the existence of the discriminant function optimal by the number

and composition of components. Let's formulate the conditions of reduction (simplification) optimal discriminant function for a special case of an independent feature. Let the set of V is those, that is carried out $\overset{\circ}{X} = V \cup \overset{\circ}{x}$, where $\overset{\circ}{x} \in \overset{\circ}{X}$ (one feature is missed). Taking into account (7), we receive

$$\Delta(V) = E\{D_{AB}^2(\overset{\circ}{X})\} - E\{D_{AB}^2(V)\} =$$

$$\left(\tau_{\overset{\circ}{X}}^2 - \frac{\tau_{\overset{\circ}{X}}^2 \left[\overset{\circ}{m} - (r-1) / (r - \overset{\circ}{m}) \right] c_A^{-1}}{\tau_{\overset{\circ}{X}}^2 + \overset{\circ}{m} c_A^{-1}} + c_B^{-1} \frac{r-1}{r - \overset{\circ}{m}} \right) \frac{r - \overset{\circ}{m}}{r-1} -$$

$$- \left(\tau_V^2 - \frac{\tau_V^2 \left[(\overset{\circ}{m}-1) - (r-1) / (r - \overset{\circ}{m} + 1) \right] c_A^{-1}}{\tau_V^2 + (\overset{\circ}{m}-1) c_A^{-1}} + c_B^{-1} \frac{r-1}{r - \overset{\circ}{m} + 1} \right) \frac{r - \overset{\circ}{m} + 1}{r-1}. \quad (9)$$

According to the above mentioned lemmas for Mahalanobis distances of sets V and $\overset{\circ}{X}$ the ratio $\tau_V^2 = \tau_{\overset{\circ}{X}}^2 - \gamma^2$ is carried out, where $\gamma^2 = \sigma_{\overset{\circ}{x}}^{-2} (\chi_{\text{I}} - \chi_{\text{II}})^2$ is the component of Mahalanobis distance, that caused by the missed independent feature $\overset{\circ}{x} \in \overset{\circ}{X}$. In view of it, having limited to accuracy $(1/n)$, neglecting members of the order $(1/n^2)$, we receive

$$\Delta(V) = \frac{1}{\left(\tau_{\overset{\circ}{X}}^2 + \overset{\circ}{m} c_A^{-1} \right) \cdot \left[\left(\tau_{\overset{\circ}{X}}^2 - \gamma^2 \right) + (\overset{\circ}{m}-1) c_A^{-1} \right]} \cdot \left\{ - \left(\tau_{\overset{\circ}{X}}^2 \cdot \frac{r - \overset{\circ}{m} + 1}{r-1} + \frac{r - \overset{\circ}{m}}{r-1} \cdot \overset{\circ}{m} \cdot c_A^{-1} \right) \cdot (\gamma^2)^2 + \right.$$

$$\left. + \tau_{\overset{\circ}{X}}^2 \cdot \left(\tau_{\overset{\circ}{X}}^2 \cdot \frac{r - \overset{\circ}{m} + 2}{r-1} + 2 \cdot \frac{r - \overset{\circ}{m}}{r-1} \cdot \overset{\circ}{m} \cdot c_A^{-1} \right) \cdot \gamma^2 - \left(\tau_{\overset{\circ}{X}}^2 \right)^2 \cdot \left(\tau_{\overset{\circ}{X}}^2 \cdot \frac{1}{r-1} + \frac{r - \overset{\circ}{m}}{r-1} \cdot c_A^{-1} \right) \right\}. \quad (10)$$

The value $\Delta(V)$ can be both positive, and negative. If $\Delta(V) > 0$ the feature $\overset{\circ}{x}$ is necessary for including in discriminant function. If the $\Delta(V) < 0$ the $\overset{\circ}{x}$ should not be included in discriminant function as it will lead to decreasing of value D_{AB}^2 , i. e. addition of an feature $\overset{\circ}{x}$ does not improve quality discriminant function by considered criterion. The condition $\Delta(V) < 0$ is a condition of a reduction (simplification) of discriminant function that is optimal by quantity and structure of features. This condition represents a condition of negative definiteness of a quadratic

trinomial relatively γ^2 in braces (10). Reduction of discriminant function is possible when value γ^2 below then threshold value

$$(\gamma^2)_{por} = \tau_{\overset{\circ}{X}}^2 \cdot \frac{\left(\frac{\tau_{\overset{\circ}{X}}^2}{r-1} \right) + c_A^{-1}}{\tau_{\overset{\circ}{X}}^2 \left(\frac{r-m+1}{r-1} \right) + m c_A^{-1}} \tag{11}$$

In figure dependences of threshold value (11) from volume samples n for a set of Mahalanobis distance $\tau_{\overset{\circ}{X}}^2$ ($\tau_{\overset{\circ}{X}}^2 = 6, 8, \dots, 18$) are submitted at fixed $m = 6$.

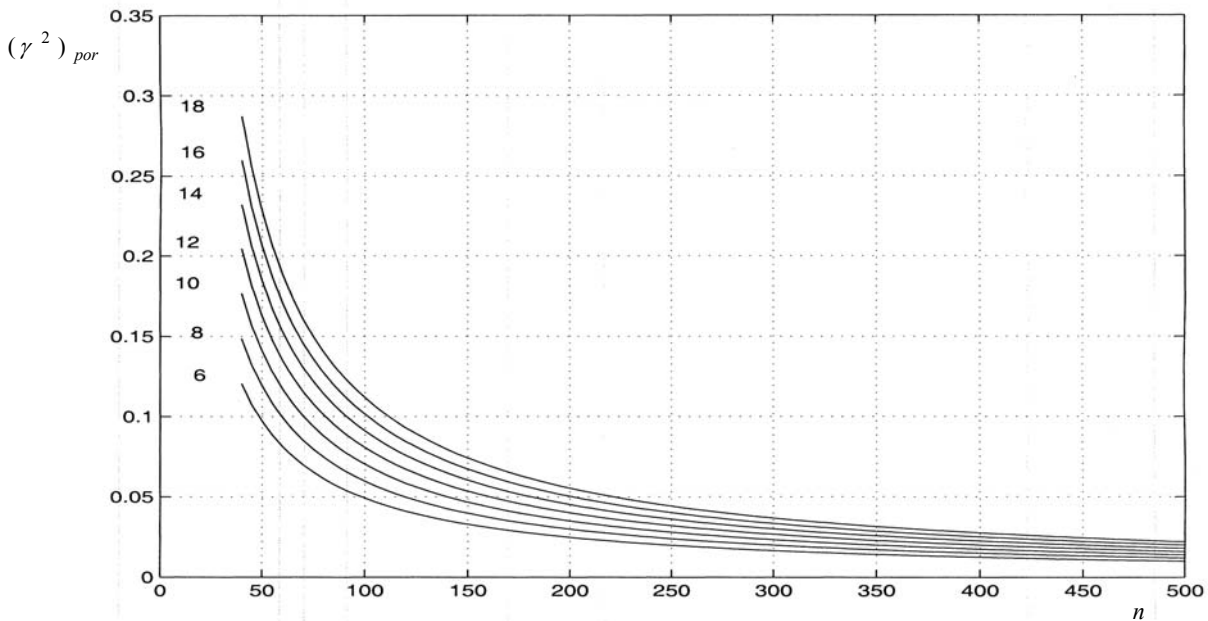


Fig. – Dependences of threshold value $(\gamma^2)_{por}$ on volume of subsamples n

Let's note, that in asymptotic at $n \rightarrow \infty$ ($r \rightarrow \infty, c_A^{-1} \rightarrow 0$) the condition of the reduction is not carried out, i. e. $V_{OPT} = \overset{\circ}{X}$.

3. Conclusions

The method for comparison of the discriminant functions based on dividing of the initial data sample on training and testing subsamples is proved. In spite of successful use of this way in practice and repeated confirmation of its efficiency by the method of statistical test, it was supposed traditionally as heuristic method. Conditions of

reduction (simplification) of discriminant function which is optimal by structure of features are revealed. It is obtained that this conditions depend on volumes samples and parameters of general sets, i.e. on mathematical expectations and covariance matrices of features. It was shown that under condition of structural uncertainty and the absence of a priori estimates of parameters of general sets this method make it possible to solve the problem of search of the discriminant function of optimal complexity.

References

1. Sarychev A. P. Circuit of Discriminant Analysis with Learning and Checking Subsamplings of Observations. *Automatica (Ukraine)*, 1: 32–41, 1990 (*Journal of Automation and Information Sciences*. – Scripta Technika Inc. – Vol. 23, 1: 29–39, 1990).
2. Miroshnichenko L. V., Sarychev A. P. Scheme of the Sliding Exam for Search of the Optimal Set of Characters in the Problem of Discriminant Analysis. *Automatica (Ukraine)*, 1: 35–44, 1992 (*Journal of Automation and Information Sciences*. – Scripta Technika Inc., Vol. 25, 1: 33–42, 1992).
3. Sarychev A. P. The Solution of the Discriminant Analysis Task in Conditions of Structural Uncertainty on Basis of the Group Method of Data Handling. *Problemy Upravlenia i Informatiki (Ukraine)*, 3: 100–112, 2008 (*Journal of Automation and Information Sciences*. – Begell House Inc., Vol. 40, 6: 27–40, 2008).
4. Sarychev A. P. Identification of Structural-Uncertain Systems States. The book. Dnipropetrovs'k, Ukraine, NAS Ukraine and NSA Ukraine, Institute of Technical Mechanics, 268, 2008.
5. Sarychev A. P., L. V. Sarycheva. The Optimal Set Features Determination in Discriminant Analysis by the Group Method of Data Handling // *Systems Analysis and Modeling Simulation (SAMS)*, Overseas Publishers Association. – 1998. – Vol. 31. – P. 153–167.
6. Sarychev A. P. S-Scheme of Sliding Examination for Optimal Set Features Determination in Discriminant Analysis by the Group Method of Data Handling / A. P. Sarychev // *System Analysis and Modelling Simulation (SAMS)*, Taylor & Francis. – 2003. – Vol. 43. – № 10. – P. 1351–1362.
7. Sarychev A. P., Sarycheva L. V. Quality Estimation of Discriminant Functions by Sliding Examination // *The Forth Workshop on Inductive Modelling (IWIM-2011)* : July 4–8 2011, Kyiv, Ukraine : Proc. – Kyiv : IRTC ITS, 2011. – ISBN 978-966-02-6078-8. – P. 104–108.