

NUMERICAL METHOD OF CONSTRUCTION THE QUADRATIC DISCRIMINATION CRITERION MINIMIZING THE MAXIMUM OF VARIATION

V.Y. SHELEKHOVA, Y.P. YURACHKOVSKI*

The problem of adequate regression function choice is considered. The concept of normalized quadratic discrimination criterion (NQDC) is introduced. Numerical method based on the theory of linear operators eigenvalues disturbances for the purpose of construction the NQDC with minimum of its maximum variation is offered. Examples illustrating offered method applications are given.

1. Introduction

Discrimination criteria are intended to choose from a given finite set of known functions (with unknown parameters) the one which coincides with unknown "true" dependence kept in the sample of observations and distorted with additive noise. The value of discrimination criterion is used as a measure of distance from any function to the "true" dependence. The quality of discrimination criterion can be compared by means of different approaches based on the various definitions of discrimination criterion optimality. Here we suppose, that the criterion minimizing the maximum of variation is the optimal one.

2. Assumptions and Problem Statement

Let's suppose that the output of a random experiment, carried out in some points $v \in V \subset \mathbb{R}^m$ is the observed random value $y(v, \omega)$, $\omega \in \Omega$, where Ω is the space of elementary events. We also suppose that random values $y(v, \omega)$ are independent and identically distributed and they have an unknown mathematical expectation $\bar{y}(v)$, unknown independent from v finite variation σ^2 and equal to zero coefficients of asymmetry and excess:

*Cybernetics Institute of Ukraine Academy of Sciences, 252028 Kiev, Ukraine-USSR

$$Ey(v, \omega) = \bar{y}(v) \quad (1)$$

$$E(y(v, \omega) - \bar{y}(v))^2 = Dy(v, \omega) = \sigma^2 < \infty \quad (2)$$

$$E(y(v, \omega) - \bar{y}(v))^3 / \sigma^3 = 0 \quad (3)$$

$$E(y(v, \omega) - \bar{y}(v))^4 / \sigma^4 - 3 = 0 \quad (4)$$

where E means operator of mathematical expectation.

The number of experiments we designate as n and the points in which the experiment was made are w_1, \dots, w_n . Let $f_\alpha(\cdot, \cdot)$, $f_\beta(\cdot, \cdot)$ be two functions determined on the sets $V \times \Theta_\alpha$, $V \times \Theta_\beta$ correspondingly. Furthermore, let these functions be linear by the parameters $\Theta_\alpha \in \Theta_\alpha$, $\Theta_\beta \in \Theta_\beta$ correspondingly, i.e.

$$f_\alpha(\cdot, \Theta_\alpha) = \Theta_\alpha^T \phi(\cdot) \quad (5)$$

$$f_\beta(\cdot, \Theta_\beta) = \Theta_\beta^T \phi(\cdot) \quad (6)$$

where

$$\phi(\cdot) = [\phi_1(\cdot), \dots, \phi_{m_\alpha}(\cdot)]^T \quad \phi(\cdot) = [\phi_1(\cdot), \dots, \phi_{m_\beta}(\cdot)]^T$$

Function $f_\alpha(\cdot, \cdot)$, $f_\beta(\cdot, \cdot)$ are called structures of regression models or simply structures and instead of $f_\alpha(\cdot, \cdot)$, $f_\beta(\cdot, \cdot)$ we shall write f_α , f_β accordingly.

Definition 1. Structure f_γ is adequate to random value $y(\cdot, \omega)$ on the set $V_0 \subset V$ if there exists such a value $\Theta_\gamma^0 \in \Theta_\gamma$ that $f_\gamma(v, \Theta_\gamma^0) = \bar{y}(v)$, $v \in V_0$.

In terms defined above our problem is to find an adequate (on some set) structure among f_α and f_β if one of them is adequate.

Function $\bar{y}(v)$ is unknown, therefore a definition can't be applied for solving this problem. Below we shall use statistical estimations of structure which may be adequate.

Let $y_i = (w_i, \omega)$, $i = 1, \dots, n$,

$$Y = [y_1, \dots, y_n]^T$$

$$\bar{Y} = [\bar{y}(w_1), \dots, \bar{y}(w_n)]^T$$

Definition 2. Value

$$cr(f_\gamma) = Y^T F_\gamma Y \tag{7}$$

is called normalized quadratic discrimination criterion (NQDC) if for matrix F_γ the following properties are satisfied:

- 1) $F = F^T$
- 2) $F \geq 0$
- 3) if structure f_γ is adequate on the set $\{w_1, \dots, w_n\}$ then $\bar{Y}^T F_\gamma \bar{Y} = 0$.
- 4) $tr F_\gamma = 1$.

Statement 1. If structure f_β is adequate on the set $\{w_1, \dots, w_n\}$, then for any F_α and F_β

$$E cr(f_\alpha) - E cr(f_\beta) \geq 0 \tag{8}$$

Proof. It's easy to obtain the following equality

$$E cr(f_\gamma) = E[Y^T F_\beta Y] = \bar{Y}^T F_\gamma \bar{Y} + \sigma^2 \tag{9}$$

for any γ . In addition $\bar{Y}^T F_\beta \bar{Y} = 0$ and $\bar{Y}^T F_\alpha \bar{Y} \geq 0$. Statement (8) follows from (9) and these inequalities.

So we see that

$$\hat{f} = \arg \min_{f_\gamma \in \{f_\alpha, f_\beta\}} cr(f_\gamma) \tag{10}$$

may be used as estimation of the adequate structure on the set $\{w_1, \dots, w_n\}$. Matrix F_γ in (7) and in (10) is not unique. Arbitrariness in its choice we shall use to optimize discrimination criterion quality.

Let

$$X = [x_{ij}] = [\phi_j(w_i)] \quad Z = [z_{ij}] = [\phi_j(w_i)]$$

The criterion $d(\alpha/\beta)$ of discrimination criterion quality assumes a maximum of variation of random value $cr(f_\alpha)$ calculated in supposition that structure f_β is adequate on the set $\{w_1, \dots, w_n\}$ and $\|\Theta_\beta^0\|^2 = \Theta_\beta^{0T} \Theta_\beta^0 = \eta^2 > 0$, i.e.

$$d(\alpha/\beta) = \max_{\bar{Y} \in \mathcal{Y}_\beta} D \text{ cr}(f_\alpha)$$

where

$$\mathcal{Y}_\beta = \{\bar{Y} : \bar{Y} = Z\Theta_\beta, \|\Theta_\beta\|^2 = \eta^2 > 0\}$$

We define also

$$F_D^*(\alpha/\beta) = \arg \min_{\{F_\alpha\}} d(\alpha/\beta)$$

Our main purpose is to create the numerical algorithm of matrix $F_D^*(\alpha/\beta)$ construction.

3. Method and Mathematical Substantiation

Let N be such a matrix that

$$N^T N = I$$

$$N N^T = I - X(X^T X)^{-1} X^T$$

and

$$N^T Z Z^T N = U F U^T \quad (11)$$

is the spectral factorization in which

$$F = \text{diag}(\gamma_1, \dots, \gamma_n) = \begin{bmatrix} \gamma_1 & & 0 \\ & \ddots & \\ 0 & & \gamma_n \end{bmatrix}, \quad \gamma_1 \leq \dots \leq \gamma_n$$

It is easy to prove (Yurachkovski) that any matrix F_α may be represented in the form

$$F_\alpha = N A A^T N^T \quad (12)$$

where A is such arbitrary matrix that

$$\text{tr} A A^T = I$$

As it was proved by Seber (1980)

$$D \text{ cr}(f_\alpha) = D[Y^T F_\alpha Y] = 2\sigma^4 \text{tr} F_\alpha^2 + 4\sigma^2 \bar{Y}^T F_\alpha \bar{Y}$$

therefore

$$d(\alpha/\beta) = \max_{\bar{Y} \in \mathcal{Y}_\beta} D \text{ cr}(f_\alpha) = 2\sigma^4 \text{tr} F_\alpha^2 + 4\sigma^2 \bar{Y}^T F_\alpha \bar{Y} \quad (13)$$

$$= 2\sigma^4 \text{tr} F_\alpha^2 + 4\sigma^2 \max_{\bar{Y} \in \mathcal{Y}_\beta} \bar{Y}^T F_\alpha \bar{Y} \quad (14)$$

$$= 2\sigma^4 \text{tr} F_\alpha^2 + 4\sigma^2 \max_{\|\Theta_\beta\|^2 = \eta^2} \Theta^T Z^T F_\alpha Z \Theta \quad (15)$$

$$= 2\sigma^4 \text{tr} F_\alpha^2 + 4\sigma^2 \eta^2 \lambda_{-1}[Z^T F_\alpha^2 Z] \quad (16)$$

where $\lambda_{-1}[C]$ designates the greatest eigenvalue of matrix C .

From equalities (11), (12) and $\lambda_{-1}[CC^T] = \lambda_{-1}[C^T C]$ follows that

$$d(\alpha/\beta) = 2\sigma^4 \text{tr}(AA^T)^2 + 4\sigma^2 \eta^2 \lambda_{-1}[AA^T U \Gamma U^T AA^T].$$

Let $h = \sigma^2/(2\eta^2)$ then

$$F_D^*(\alpha/\beta) = N(AA^T)_D^* N^T$$

where

$$(AA^T)_D^* = \arg \min_{\{AA^T: \text{tr} AA^T = 1\}} \left(\lambda_{-1}[AA^T U \Gamma U^T AA^T] + h \cdot \text{tr}(AA^T)^2 \right)$$

or

$$F_D^*(\alpha/\beta) = NU(BB^T)_D^* U^T N^T$$

where

$$(BB^T)_D^* = \arg \min_{\{BB^T: \text{tr} BB^T = I\}} \left(\lambda_{-1}[BB^T \Gamma BB^T] + h \text{tr}(BB^T)^2 \right)$$

and

$$BB^T = UAA^T U^T$$

We shall accomplish the search of $(BB^T)_D^*$ by means of the gradient projection method. This method consists of two procedures:

1) finding the direction of the greatest decreasing of the minimized function

$$\Phi(B) = \lambda_{-1}[BB^T \Gamma BB^T] + h \operatorname{tr}(BB^T)^2$$

2) projection of obtained direction vector on the permissible domain

$$B = \{BB^T : \operatorname{tr} BB^T = I\}$$

The direction of the greatest decreasing of function $\Phi(B)$ is given by matrix $-\partial\Phi(B)/\partial B$. It is easy to see that

$$\frac{\partial\Phi(B)}{\partial B} = \left[\frac{\partial\Phi(B)}{\partial b_{ij}} \right] = \left[\frac{\partial\lambda_{-1}[BB^T \Gamma BB^T]}{\partial b_{ij}} \right] + h \left[\frac{\partial\operatorname{tr}(BB^T)^2}{\partial b_{ij}} \right]$$

If multiplicity of $\lambda_{-1}[BB^T \Gamma BB^T]$ is one then

$$\frac{\partial\lambda_{-1}[BB^T \Gamma BB^T]}{\partial b_{ij}} =$$

$$\lim_{\varepsilon \rightarrow 0} \frac{\lambda_{-1}[(B + \varepsilon G_{ij})(B + \varepsilon G_{ij})\Gamma(B + \varepsilon G_{ij})(B + \varepsilon G_{ij})^T] + \lambda_{-1}[BB^T \Gamma BB^T]}{\varepsilon}$$

where G_{ij} is a matrix whose element with number (i, j) is equal to one and the rest of the elements are equal to zero.

According to the eigenvalue perturbation theory the last limit may be easily calculated. It is equal to

$$v_{-1}^T (G_{ij} B^T \Gamma BB^T + B G_{ij}^T \Gamma BB^T + B^T \Gamma G_{ij} B^T + BB^T \Gamma B G_{ij}^T) v_{-1}$$

where v_{-1} is matrix $BB^T \Gamma BB^T$ eigenvector corresponding to the greatest eigenvalue. Thus

$$\frac{\partial\lambda_{-1}[BB^T \Gamma BB^T]}{\partial b_{ij}} = 2v_{-1} v_{-1}^T BB^T \Gamma B + 2\Gamma BB^T v_{-1} v_{-1}^T B \quad (17)$$

It is easy to show that

$$h \frac{\partial \text{tr}(\mathbf{B}\mathbf{B}^T)^2}{\partial b_{ij}} = 4h\mathbf{B}\mathbf{B}^T\mathbf{B} \quad (18)$$

From (13), (14) we obtain that

$$\frac{\partial \Phi(\mathbf{B})}{\partial \mathbf{B}} = 2\mathbf{v}_{-1}\mathbf{v}_{-1}^T\mathbf{B}\mathbf{B}^T\mathbf{\Gamma}\mathbf{B} + 2\mathbf{\Gamma}\mathbf{B}\mathbf{B}^T\mathbf{v}_{-1}\mathbf{v}_{-1}^T\mathbf{B} + 4h\mathbf{B}\mathbf{B}^T\mathbf{B} \quad (19)$$

Denote $\mathbf{B}_k\mathbf{B}_k^T$ the approximation of k -th iteration. Then in procedure 1 of $(k + 1)$ -th iteration consists in calculation of

$$\widetilde{\mathbf{B}}_{k+1} = \mathbf{B}_k - s \frac{\partial \Phi(\mathbf{B})}{\partial \mathbf{B}} \Big|_{\mathbf{B}=\mathbf{B}_k} = \mathbf{B}_k - s\Phi'(\mathbf{B}_k)$$

where $s > 0$ is the step of $(k + 1)$ -th iteration.

Procedure 2 consists of projecting the matrix $\widetilde{\mathbf{B}}_{k+1}\widetilde{\mathbf{B}}_{k+1}^T$ on domain \mathcal{B} , i.e.

$$\begin{aligned} \mathbf{B}_{k+1}\mathbf{B}_{k+1}^T &= \widetilde{\mathbf{B}}_{k+1}\widetilde{\mathbf{B}}_{k+1}^T / \text{tr}\widetilde{\mathbf{B}}_{k+1}\widetilde{\mathbf{B}}_{k+1}^T = \\ &= \frac{(\mathbf{B}_k - s\Phi'(\mathbf{B}_k))(\mathbf{B}_k - s\Phi'(\mathbf{B}_k))^T}{\text{tr}(\mathbf{B}_k - s\Phi'(\mathbf{B}_k))(\mathbf{B}_k - s\Phi'(\mathbf{B}_k))^T} \end{aligned} \quad (20)$$

If multiplicity of eigenvalue $\lambda_{-1}[\mathbf{B}\mathbf{B}^T\mathbf{\Gamma}\mathbf{B}\mathbf{B}^T]$ is not one formula (13) is not correct and hence formula (16) does not lead to the step of gradient projection method. If the only multiplicity of $\lambda_{-1}[\mathbf{B}\mathbf{B}^T\mathbf{\Gamma}\mathbf{B}\mathbf{B}^T]$ isn't one after one or more steps performed by formula (16) algorithm comes to the point where multiplicity of eigenvalue is one and the using of formula (14) becomes well-grounded.

4. Algorithm Description

The gradient search is used in the algorithm. Therefore, the speed of the algorithm depends considerably on the initial approximation of matrix \mathbf{B} (or $\mathbf{B}\mathbf{B}^T$), direction of motion and step at every iteration, and halt criteria. *Initial approximation of matrix $\mathbf{B}\mathbf{B}^T$* we obtain as a weighted sum of two matrix, i.e.

$$\mathbf{B}\mathbf{B}^T = p_\lambda(\mathbf{B}\mathbf{B}^T)_\lambda + p_h(\mathbf{B}\mathbf{B}^T)_h$$

where $(\mathbf{B}\mathbf{B}^T)_\lambda$ and $(\mathbf{B}\mathbf{B}^T)_h$ are matrices on which minimum of functionals $\lambda_{-1}[\mathbf{B}\mathbf{B}^T\Gamma\mathbf{B}\mathbf{B}^T]$ and $\text{tr}[(\mathbf{B}\mathbf{B}^T)^2]$ are achieved correspondingly. Note that this matrices are diagonal. Elements of matrix $(\mathbf{B}\mathbf{B}^T)$ depend on n_2 – number of zeroes in $\gamma = \text{diag}(\gamma_1, \dots, \gamma_n) : (\mathbf{B}\mathbf{B}^T)_\lambda = \text{diag}(1/n_z, \dots, 1/n_z, 0, \dots, 0)$, if n_z isn't zero ; otherwise, elements $(\mathbf{B}\mathbf{B}^T)_\lambda$ are calculated from equations

$$\begin{cases} (\mathbf{B}\mathbf{B}^T)_\lambda^i \gamma_i^{1/2} = (\mathbf{B}\mathbf{B}^T)_\lambda^j \gamma_j^{1/2} & \text{for every } i, j = 1, \dots, n \\ \sum_{i=1}^n (\mathbf{B}\mathbf{B}^T)_\lambda^i = 1 \end{cases}$$

Hence, i -th element of $\text{diag}(\mathbf{B}\mathbf{B}^T)_\lambda$ is

$$(\mathbf{B}\mathbf{B}^T)_\lambda^i = \frac{1}{\gamma_i^{1/2} \left(\sum_{j=1}^n 1/\gamma_j^{1/2} \right)}$$

Diagonal elements of matrix $(\mathbf{B}\mathbf{B}^T)_h^i$ are equal to $1/n$. Weights p_λ and p_h depend on $\{\gamma_i\}$ and h .

Choice of direction and step. At any point of iteration algorithm, motion is conducted towards the antigradient $-\partial\Phi(\mathbf{B}\mathbf{B})/\partial\mathbf{B}\mathbf{B}$. The norm of $\partial\Phi(\mathbf{B}\mathbf{B})/\partial\mathbf{B}\mathbf{B}$ depends linearly on the values $\{\gamma_i\}$ and h , so that if the value $\|\partial\Phi(\mathbf{B}\mathbf{B})/\partial\mathbf{B}\mathbf{B} \cdot \text{step}\|$ isn't very large (or small) it's necessary to normalize values $\{\gamma_i\}$ and h .

The motion is conducted until Φ_k decreases, where Φ_k is the value of functional Φ at the k -th iteration. The increasing of Φ_k is possible either when the step exceeds distance to extremum or when antigradient of functional Φ_k is considerably changed. The last case we shall determine below as a leap over junction of functional Φ . In both cases the motion should be continued, the step decreased beforehand. Note, that the iteration process can be conducted along the functional junction. In this case the rapid step decrease will result in unnecessary growth of calculation time. To obtain available value of eigenvalue $\lambda_{-1}[\mathbf{B}\mathbf{B}^T\Gamma\mathbf{B}\mathbf{B}^T]$ Jacobi algorithm is used, and to get corresponding eigenvector $v_{-1}[\mathbf{B}\mathbf{B}^T\Gamma\mathbf{B}\mathbf{B}^T]$ power method is used.

Stop criterion. During the iteration process the minimal value of Φ (we denote it Φ_{opt}) and the matrix in which it's achieved are saved. The criterion of the end of calculations consists in unchangeability of Φ_{opt} during N times leaps over junctions.

Let's mark the additional peculiarities of the given algorithm. On some steps of the gradient method it may occur that matrix $\mathbf{B}\mathbf{B}^T$ becomes

diagonal, in this case the further search will remain within the space of the diagonal matrix. To avoid this it is necessary sometimes to check the diagonality of BB^T and, when diagonality takes place, it is necessary to noise up the matrix elements:

$$(BB^T)_{k+1}^{ij} = (BB^T)_k^{ij} + \varepsilon_k,$$

where ε_k is the monotonously decreasing sequence. Apart from the checking up of diagonality it's necessary to watch if BB^T is a singular matrix. In algorithm a Cholesky triangle factorization procedure for positive definite symmetric matrix is used. Noising up of diagonal elements makes it possible to leave the class of the singular matrix, because for any singular matrix $A(A + \xi I)$ is a nonsingular one for any small $\xi > 0$.

To increase the speed of calculations eigenvector v_{-1} of the k -th iteration is handed as an initial approximation of v_{-1} at the $(k+1)$ -th iteration. When doing this we must bear in mind that the junction of functional means that λ_{-1} is considerably changed. Thus it can occur that v_{-1} from the previous iteration is saved as eigenvector, but it does not corresponds to λ_{-1} . It takes place, for example, for $h = 0$ eigenvectors of matrix $BB^T \Gamma BB^T$ are orthonormal basis e_1, \dots, e_n . If v_{-1} is e_i at the k -th iteration and leaps over the junction it becomes e_j at the $(k+1)$ -th iteration ($i \neq j$), power method will not detect the error. Thus we must change the initial approximation of v_{-1} , when Φ_k increases:

$$v_{-1}^k = v_{-1}^{k+1} + \delta, \quad \delta = \text{const}$$

Note, that every time the elements are noised up we must normalize v_{-1} and BB^T by formulas $\|v_{-1}\| = 1$, $\text{tr} BB^T = I$.

In the algorithm we use following values and constants:

$$p_\lambda/p_h = \text{tr} \gamma/h, \quad gh = \max(\gamma_1, \dots, \gamma_n h), \quad \text{step}_k/\text{step}_{k+1} = 1.02,$$

$$\{\varepsilon_k\} : \varepsilon_0 = 0.01, \quad \text{and } \varepsilon_{k+1} = \varepsilon_k/2, \quad \xi_k = 10\varepsilon_k, \quad \delta = 0.01.$$

5. Examples

We have set up in all the examples the same matrix $\Gamma = \text{diag}(1, 4, 9, 16)$.

1. $h = 0$.

In this case we obviously obtained the diagonal matrix

$$(\mathbf{B}\mathbf{B}^T)_D^* = \text{diag}(0.48, 0.24, 0.16, 0.12).$$

2. $h = 0.1$.

$$(\mathbf{B}\mathbf{B}^T)_D^* = \begin{bmatrix} 0.480478 & 0.000020 & 0.000031 & -0.000068 \\ & 0.239764 & -0.000176 & -0.000087 \\ & & 0.159956 & -0.000018 \\ & & & 0.119803 \end{bmatrix}.$$

3. $h = 100$.

It's easy to show that $(\mathbf{B}\mathbf{B}^T)_D^* \rightarrow \text{diag}(1/n, \dots, 1/n)$ when $h \rightarrow \infty$. In this example matrix $(\mathbf{B}\mathbf{B}^T)_D^*$ is close to the diagonal one. Really, the modules of all nondiagonal elements are less than 10^{-6} . In our calculations

$$(\mathbf{B}\mathbf{B}^T)_D^* = \text{diag}(0.258929, 0.258929, 0.258929, 0.223214).$$

References

- Yurachkovski Yu.P.**, *Analytical construction of optimal quadratic discrimination criterion.*— Soviet Journal of Automation and Information Sciences, v.21, No.1.
- Seber J.** (1980): *Linear regression analyses.*— Moscow, Mir, (in Russian).