

УДК 519.25

ОЦЕНИВАНИЕ КАЧЕСТВА ДИСКРИМИНАНТНЫХ ФУНКЦИЙ С РАЗБИЕНИЕМ НАБЛЮДЕНИЙ НА ОБУЧАЮЩИЕ И ПРОВЕРОЧНЫЕ ПОДВЫБОРКИ

А.П. Сарычев¹, Л.В. Сарычева²¹ *Институт технической механики НАНУ и ГКАУ, г. Днепропетровск*² *Национальный горный университет МОН Украины, г. Днепропетровск**Sarychev@prognoz.dp.ua, Sarycheval@ntmu.dp.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.

Введение

Решение задачи дискриминантного анализа в условиях структурной неопределенности по составу признаков предполагает принятие какого-либо способа сравнения дискриминантных функций (ДФ), построенных на различных множествах признаков. Два способа сравнения популярны в приложениях. Первый основан на разбиении наблюдений на обучающие и проверочные подвыборки. В этом способе обучающие подвыборки используются для оценивания коэффициентов ДФ, а проверочные – для оценивания ее качества классификации. Второй – способ скользящего экзамена, в котором в качестве проверочных выступают наблюдения, поочередно исключаемые из обучающей выборки. В литературе эти способы традиционно трактуются как эвристические приемы, хотя факт существования в них

оптимального множества признаков неоднократно подтверждался методом статистических испытаний. В рамках метода группового учета аргументов (МГУА) проведено аналитическое исследование этих способов [1–2].

1. Способ сравнения дискриминантных функций с разбиением наблюдений на обучающие и проверочные подвыборки

Пусть на этапе с номером $s (s=1,2,\dots, m)$ алгоритма полного перебора сочетаний признаков в ДФ может быть включено только s компонент из множества X , составляющих текущее анализируемое множество V . Пусть V соответствуют: 1) \mathbf{V}_I и \mathbf{V}_{II} – $(s \times n_I)$ - и $(s \times n_{II})$ -матрицы наблюдений из генеральных совокупностей P_I и P_{II} ; 2) \mathbf{v}_I и \mathbf{v}_{II} – $(s \times 1)$ -векторы математических ожиданий для P_I и P_{II} ; 3) Σ_V – ковариационная $(s \times s)$ -матрица. Рассмотрим оценку расстояния Махаланобиса, рассчитываемую с учетом разбиения наблюдений на обучающие и проверочные подвыборки. Вычислим оценки коэффициентов ДФ для множества компонент V на обучающей подвыборке A и используем их для оценивания расстояния Махаланобиса как отношение межгрупповой вариации к внутригрупповой вариации на проверочной подвыборке B :

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

В (1) $(s \times 1)$ -вектор $\hat{\mathbf{d}}_A$ представляет собой рассчитанную на подвыборке A фишеровскую оценку коэффициентов ДФ, которая построена в пространстве компонент множества V

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

где $(s \times 1)$ -векторы $\tilde{\mathbf{v}}_{IA}$ и $\tilde{\mathbf{v}}_{IIA}$ – оценки математических ожиданий \mathbf{v}_I и \mathbf{v}_{II} :

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

$(s \times s)$ -матрица \mathbf{S}_A – несмещенная оценка ковариационной матрицы Σ_V

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

В (4) $\mathbf{v}_{kA} (k = I, II)$ – $(s \times n_k)$ -матрицы, составленные из отклонений наблюдений \mathbf{V}_{kA} компонент множества V от оценок $\tilde{\mathbf{v}}_{kA}$

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

В (1) $(s \times 1)$ -векторы $\tilde{\mathbf{v}}_{IB}$ и $\tilde{\mathbf{v}}_{IIB}$ вычисляются аналогично (3), а $(s \times s)$ -матрица \mathbf{S}_B – аналогично (4)–(5); n_{IA} и n_{IIA} , n_{IB} и n_{IIB} – объемы обучающих и проверочных подвыборок соответственно такие, что выполняется $n_{IA} + n_{IB} = n_I$ и $n_{IIA} + n_{IIB} = n_{II}$. Используя (2), для (1) получаем

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

Теорема. Для математического ожидания случайной величины $D_{AB}^2(V)$ выполняется

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

где $\tau_V^2 = (\mathbf{v}_I - \mathbf{v}_{II})^T \Sigma_V^{-1} (\mathbf{v}_I - \mathbf{v}_{II})$ – расстояние Махаланобиса для множества V ; $r = n_{IA} + n_{IIA} - 2$, $c_A^{-1} = (n_{IA}^{-1} + n_{IIA}^{-1})$, $c_B^{-1} = (n_{IB}^{-1} + n_{IIB}^{-1})$.

Справедливость теоремы следует из следующих утверждений: 1) оценки, полученные на подвыборках A и B , независимы; 2) оценки математических ожиданий (3) и ковариационной матрицы (4) независимы; 3) \mathbf{S}_A – случайная $(s \times s)$ -матрица, имеет распределение Уишарта с r степенями свободы.

Определение 1. Оптимальным множеством компонент (признаков) называется множество V_{OPT} :

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

Определение 2. Оптимальной по количеству и составу компонент называется фишеровская дискриминантная функция, построенная на множестве V_{OPT} .

Доказано существование оптимального множества признаков в способе с разбиением наблюдений на обучающую и проверочную подвыборки, и сформулированы условия, при выполнении которых оптимальная ДФ упрощается по числу входящих в нее компонент.

С этой целью исследована зависимость $E\{D_{AB}^2(V)\}$ от состава множества V . Множество компонент X может быть разбито на непересекающиеся подмножества $X = \overset{\circ}{X} \cup \overset{\circ}{R} \cup \tilde{R} = \overset{\circ}{V} \cup \tilde{R}$: 1) $\overset{\circ}{X} \neq \emptyset$ (\emptyset – пустое множество) – множество компонент (m – их число), для математических ожиданий которых

выполнено $\overset{\circ}{\chi}_{Ih} \neq \overset{\circ}{\chi}_{IIh}, h=1,2,\dots,m$; 2) $\overset{\circ}{R}$ – множество компонент, для математических ожиданий которых выполнено $\overset{\circ}{\rho}_{Ih} = \overset{\circ}{\rho}_{IIh}, h=1,2,\dots,l$, где l – их число, и каждая компонента из множества $\overset{\circ}{R}$ статистически зависит хотя бы от одной компоненты из множества $\overset{\circ}{X}$ (множество $\overset{\circ}{R}$ может быть пустым); 3) \tilde{R} – множество компонент, для математических ожиданий которых выполнено $\tilde{\rho}_{Ih} = \tilde{\rho}_{IIh}, h=1,2,\dots,\tilde{l}$, где \tilde{l} – их число, и каждая компонента из множества \tilde{R} статистически не зависит от любой из компонент множества $\overset{\circ}{X}$ (множество \tilde{R} может быть пустым).

Сформулированы в виде лемм соотношения между расстоянием Махаланобиса для множества компонент $\overset{\circ}{V} = \overset{\circ}{X} \cup \overset{\circ}{R}$ и расстоянием Махаланобиса для произвольного текущего анализируемого множества компонент $V \subseteq X$ [1–4]. Для случая известных параметров генеральных совокупностей из сформулированных лемм следует: 1) любая компонента из множества $\overset{\circ}{X}$ необходима в том смысле, что ее включение в текущее множество компонент V увеличивает расстояние Махаланобиса τ_V^2 ; 2) любая компонента из множества $\overset{\circ}{R}$ необходима в том смысле, что ее включение в V увеличивает расстояние Махаланобиса τ_V^2 ; 3) любая компонента из множества \tilde{R} избыточна в том смысле, что ее включение в множество V не увеличивает расстояния Махаланобиса τ_V^2 .

2. Условие редукции (упрощения) оптимальной дискриминантной функции

В практических приложениях параметры генеральных совокупностей, как правило, неизвестны, но могут быть получены как статистические оценки по обучающим выборкам наблюдений ограниченного объема. Известно, что если применить построенное правило классификации к обучающей выборке, то оценка качества распознавания будет завышена по математическому ожиданию по сравнению с той же оценкой качества на независимых от обучения данных. Способ с разбиением наблюдений на обучающие и проверочные подвыборки дает незавышенные оценки качества распознавания. Опыт практических применений и тестовые исследования на основе метода статистических испытаний показывают, что в этой схеме: 1) с увеличением объема выборок наблюдений увеличивается количество компонент во множестве V_0 , на котором достигается наилучшее качество распознавания, а с уменьшением объема

выборки наблюдений количество компонент в V_0 уменьшается; 2) с увеличением расстояния Махаланобиса τ_X^2 между генеральными совокупностями (из которых получены выборки наблюдений) увеличивается количество компонент во множестве V_0 , а с уменьшением этого расстояния количество компонент в таком множестве уменьшается. Проведенные аналитические исследования объясняют эти эмпирически установленные закономерности.

Сформулируем условие редукции (упрощения) оптимальной ДФ для частного случая независимого признака. Пусть множество V таково, что выполняется $\overset{\circ}{X} = V \cup \overset{\circ}{x}$, где $\overset{\circ}{x} \in \overset{\circ}{X}$ (в ДФ пропущен один признак). Учитывая (7), получаем

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

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

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

В соответствии с вышеупомянутыми леммами для расстояний Махаланобиса множеств V и $\overset{\circ}{X}$ выполняется соотношение: $\tau_V^2 = \tau_{\overset{\circ}{X}}^2 - \gamma^2$, где $\gamma^2 = \sigma_x^{-2} (\chi_{I} - \chi_{II})^2$ – составляющая расстояния Махаланобиса, обусловленная пропущенным признаком $\overset{\circ}{x} \in \overset{\circ}{X}$ (при условии, что $\overset{\circ}{x}$ статистически не зависим с другими компонентами из $\overset{\circ}{X}$). С учетом этого, ограничившись точностью $(1/n_k)$ и пренебрегая членами порядка $(1/n_k^2)$, получаем

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

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

Величина $\Delta(V)$ может быть как положительной, так и отрицательной. Если величина $\Delta(V) > 0$, то признак $\overset{\circ}{x}$ необходимо включать в ДФ. Если величина $\Delta(V) < 0$, то признак $\overset{\circ}{x}$ не следует включать в ДФ, поскольку это приведет к уменьшению величины D_{AB}^2 , т.е. добавление признака $\overset{\circ}{x} \in \overset{\circ}{X}$ не улучшает качество ДФ по рассматриваемому критерию.

Условие $\Delta(V) < 0$ является условием редукции (упрощения) ДФ, оптимальной по количеству и составу признаков. Оно представляет собой условие отрицательной определенности квадратичного трехчлена относительно γ^2 в фигурных скобках (10). Пороговым значением для γ^2 , ниже которого возможна редукция ДФ, является значение:

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

На рисунке представлены зависимости порогового значения (11) от объема выборок n для набора расстояний Махаланобиса $\tau_{\overset{\circ}{X}}^2$ ($\tau_{\overset{\circ}{X}}^2 = 6, 8, \dots, 18$) при $m = 6$.

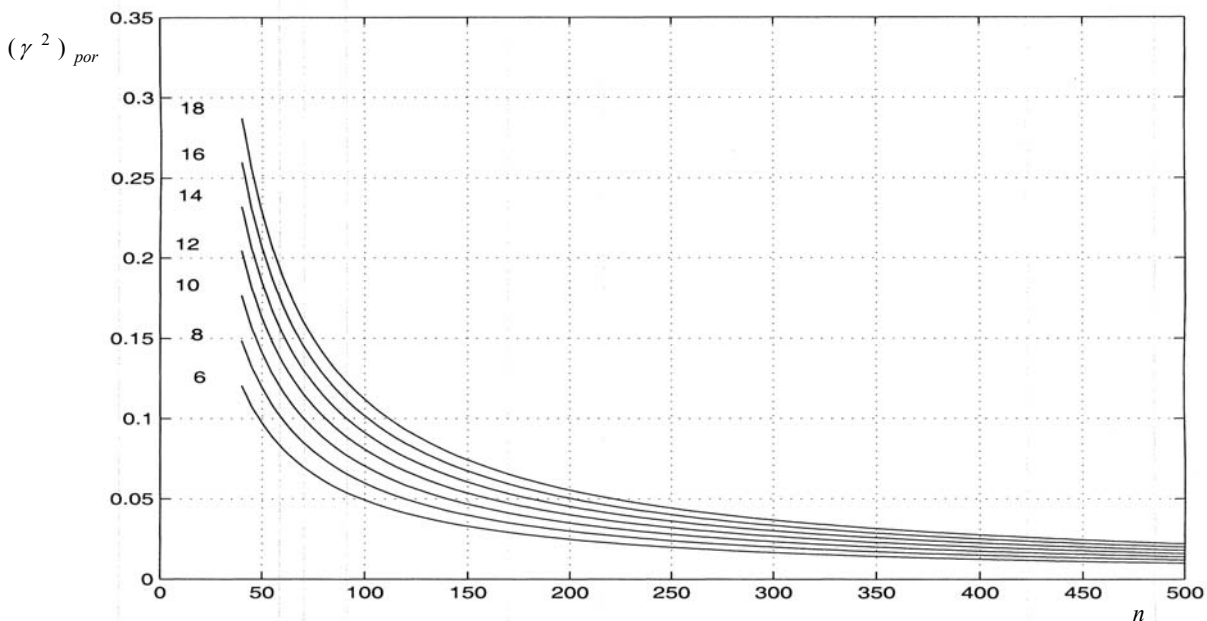


Рисунок 1 – Зависимости порогового значения $(\gamma^2)_{por}$ от объема выборок n

Отметим, что в асимптотике при $n_A \rightarrow \infty$ ($r \rightarrow \infty$) условие редукции не выполняется, поскольку $(\gamma^2)_{por} = 0$, а $\gamma^2 > 0$, т.е. $V_{OPT} = \overset{\circ}{X}$.

3. Выводы

Обоснован способ сравнения дискриминантных функций с разбиением выборок наблюдений на обучающие и проверочные подвыборки. Несмотря на успешное применение этого способа на практике и неоднократное подтверждение его работоспособности методом статистических испытаний, он традиционно считался эвристическим приёмом.

Получены условия существования оптимального множества признаков, зависящие от параметров генеральных совокупностей и объемов выборок. Выявлены закономерности упрощения оптимальной дискриминантной функции при уменьшении объемов выборок и при увеличении дисперсий признаков. Показано, что в условиях структурной неопределенности и отсутствия априорных оценок ковариационной матрицы признаков применение этого способа позволяет решать задачу поиска дискриминантной функции оптимальной сложности.

Литература

1. Сарычев А. П. Схема дискриминантного анализа с обучающими и проверочными подвыборками наблюдений / А. П. Сарычев // Автоматика. – 1990. – № 1. – С. 32–41.
2. Мирошниченко Л. В. Схема скользящего экзамена для поиска оптимального множества признаков в задаче дискриминантного анализа / Л. В. Мирошниченко, А. П. Сарычев // Автоматика. – 1992. – № 1. – С. 35–44.
3. Сарычев А. П. Решение задачи дискриминантного анализа в условиях структурной неопределенности на основе метода группового учета аргументов / А. П. Сарычев // Проблемы управления и информатики. – 2008. – № 3. – С. 100–112.
4. Сарычев А. П. Идентификация состояний структурно-неопределенных систем / А. П. Сарычев – Днепропетровск: Институт технической механики НАН Украины и НКА Украины, 2008. – 268 с.