

АСИМПТОТИЧЕСКИЙ АНАЛИЗ РИСКА НЕПАРАМЕТРИЧЕСКОЙ КЛАССИФИКАЦИИ В СЛУЧАЕ СУЩЕСТВЕННО ЗАВИСИМЫХ ПРИЗНАКОВ

Белорусский государственный университет

Введение. При решении разнообразных практических задач исследования сложных технических, производственных и экономических систем (установлении причинно-следственных связей, прогнозировании, управлении), приходится иметь дело с анализом многомерных статистических данных, которые представляют собой значения наблюдаемых и, как правило, взаимосвязанных характеристик (признаков) сложных систем. Для описания эмпирических данных, имеющих стохастическую природу, используются различные модели статистических зависимостей [1], [2].

В данной статье предполагается, что сложная система описывается случайным вектором признаков, компоненты которого связаны нелинейной статистической зависимостью. Особенностью рассматриваемой модели зависимости является то, что совместное распределение компонент вектора признаков близко к вырожденному, т. е. концентрируется на многообразии меньшей размерности, чем размерность пространства признаков. Рассматривается случай априорной непараметрической неопределенности, т. е. предполагается, что функциональный вид зависимости не известен. В общем случае допускается также неопределенность в разбиении вектора признаков на зависимые и независимые компоненты (эндогенные и экзогенные переменные). Если указанное разбиение считается известным, т. е. определена эндогенно-экзогенная структура модели, то рассматриваемые зависимости принимают вид нелинейных статистических зависимостей регрессионного типа.

Имеют место два режима функционирования сложной системы и, соответственно, два класса состояний, различающихся моделями статистических зависимостей для компонент случайного вектора признаков. Задача заключается в построении и исследовании решающего правила, предназначенного для оценки (прогнозирования) состояния сложной системы по имеющимся наблюдениям за системой (значениям вектора признаков). Для оценки состояния системы (классификации наблюдений) используются непараметрические подстановочные байесовские решающие правила (непараметрические классификаторы), получающиеся подстановкой в байесовское решающее правило [3-5] непараметрических ядерных оценок многомерных условных плотностей распределения наблюдений с фиксированным и адаптивным гауссовским ядром [6].

Аналитическое исследование рассматриваемых решающих правил проводится для случая моделью многомерной линейной регрессии с равномерно распределенными экзогенными переменными и гауссовскими ошибками наблюдения. С помощью метода асимптотических разложений в асимптотике растущего объема обучающей выборки и усиливающейся статистической зависимости компонент вектора признаков исследуется условный риск (средние потери) непараметрических подстановочных решающих правил при условии, что значения экзогенных переменных в регрессионных зависимостях являются заданными, и оценивается преимущество предлагаемого непараметрического классификатора с адаптивным ядром.

1. Модель наблюдений и непараметрическая оценка плотности с адаптивным гауссовским ядром. Пусть результаты наблюдения характеристик сложной системы в i -м эксперименте описываются случайным вектором признаков (наблюдений) $\mathbf{y}_i \in \mathfrak{R}^p$ ($p > 1$), который допускает в общем случае неизвестное разбиение на подвекторы:

$$\mathbf{y}_i = \begin{pmatrix} \mathbf{x}_i \\ \mathbf{z}_i \end{pmatrix} \in \mathfrak{R}^p, \quad \mathbf{x}_i = (x_{i1}, \dots, x_{iN})' \in \mathfrak{R}^N, \quad \mathbf{z}_i = (z_{i1}, \dots, z_{iM})' \in \mathfrak{Z} \subset \mathfrak{R}^M, \quad p=N+M, \quad i=1, \dots, n, \quad (1)$$

где $\mathfrak{Z} \subset \mathfrak{R}^M$ – ограниченная область в \mathfrak{R}^M .

Компоненты вектора $\mathbf{y}_i \in \mathfrak{R}^{N+M}$ связаны моделью статистической зависимости вида:

$$T(\mathbf{y}_i) = \mathbf{x}_i - f(\mathbf{z}_i) = \xi_i, i = 1, \dots, n, \quad (2)$$

где: $T(\cdot), f(\cdot)$ – неизвестные достаточно гладкие векторные функции; $\xi_i = (\xi_{i1}, \dots, \xi_{iN})' \in \mathfrak{R}^N$ – случайный вектор, с нулевым вектором математического ожидания и ковариационной матрицей $\Sigma \in \mathfrak{S}_N$, где \mathfrak{S}_N – семейство положительно определенных симметричных матриц размерности $N \times N$; случайные векторы $\mathbf{z}_i \in \mathbb{Z}$ и $\xi_i \in \mathfrak{R}^N$ являются статистически независимыми и имеют плотности распределения $p_z(\mathbf{z})$ и $p_\xi(\xi)$ соответственно; разбиение (1) в общем случае не известно. Для простоты записи случайный вектор и неслучайный аргумент соответствующей функции плотности распределения обозначаются одним и тем же символом.

С учетом сделанных предположений плотность распределения $p(\mathbf{y})$ случайного вектора $\mathbf{y}_i \in \mathfrak{R}^p$ ($i = 1, \dots, n$) имеет вид:

$$p(\mathbf{y}) = p_\xi(\mathbf{x} - f(\mathbf{z})) p_z(\mathbf{z}), \mathbf{x} \in \mathfrak{R}^N, \mathbf{z} \in \mathbb{Z} \subset \mathfrak{R}^M, \mathbf{y} \in \mathfrak{R}^{N+M}. \quad (3)$$

Относительно случайного вектора $\mathbf{y}_i \in \mathfrak{R}^p$ ($i = 1, \dots, n$), удовлетворяющего условиям (1)–(3), делается также предположение

$$\text{tr}(\Sigma) \rightarrow 0. \quad (4)$$

Условие (4) означает усиливающуюся статистическую зависимость компонент вектора \mathbf{y}_i при уменьшении дисперсий компонент вектора $\xi_i = (\xi_{i1}, \dots, \xi_{iN})' \in \mathfrak{R}^N$. При этом наблюдения $\{\mathbf{y}_i\} (i = 1, \dots, n)$ концентрируются в пространстве \mathfrak{R}^p «вблизи» некоторой N -мерной ($N < p$) гиперповерхности (многообразия) Γ , определяемой уравнением

$$T(\mathbf{y}) = \mathbf{0} \in \mathfrak{R}^N, \mathbf{y} \in \mathfrak{R}^p. \quad (5)$$

Тождество (5) можно интерпретировать как некоторое ожидаемое состояние системы в пространстве признаков \mathfrak{R}^p для определенного режима функционирования, а случайный вектор ξ_i – как отклонение системы от этого состояния в i -м эксперименте, обусловленное случайными и неконтролируемыми факторами.

Таким образом, распределение вероятностей случайного вектора $\mathbf{y}_i \in \mathfrak{R}^p$ ($i = 1, \dots, n$), описываемое соотношениями (1)–(3), при выполнении условия (4) приближается к вырожденному. На этом основании в работе [6] случайный вектор, удовлетворяющий условиям (1)–(4), называется случайным вектором с «существенно зависимыми» компонентами.

Для определенности и избежания громоздкости в аналитических исследованиях будем предполагать, что случайные векторы $\mathbf{z}_i \in \mathbb{Z} \subset \mathfrak{R}^M$ и $\xi_i \in \mathfrak{R}^N$ имеют соответственно равномерное в ограниченной области $\mathbb{Z} \subset \mathfrak{R}^M$ и N -мерное нормальное распределение с плотностями:

$$p_z(\mathbf{z}) = \frac{1}{\text{mes}\{\mathbb{Z}\}} \mathbf{I}_z(\mathbf{z}), p_\xi(\xi) = n_N(\xi | \mathbf{0}_N, \Sigma), \mathbf{z} \in \mathbb{Z}, \xi \in \mathfrak{R}^N, \quad (6)$$

где $n_N(\xi | \mathbf{0}_N, \Sigma)$ – функция плотности N -мерного нормального распределения с нулевым вектором математического ожидания и ковариационной матрицей $\Sigma \in \mathfrak{S}_N$; $\mathbf{I}_z(\mathbf{z})$ и $\text{mes}\{\mathbb{Z}\} < \infty$ – соответственно индикаторная функция и N -мерный объем (мера Лебега) области \mathbb{Z} .

Согласно (3) и (6) плотность распределения случайного вектора $\mathbf{y}_i \in \mathfrak{R}^{N+M}$ определяется выражением:

$$\begin{aligned} p(\mathbf{y}) &= \frac{1}{\text{mes}\{\mathbb{Z}\}} \mathbf{I}_{\mathbb{Z}}(\mathbf{z}) n_N(\mathbf{x} | f(\mathbf{z}), \Sigma) = \\ &= \frac{1}{\text{mes}\{\mathbb{Z}\}} \mathbf{I}_{\mathbb{Z}}(\mathbf{z}) \left((2\pi)^N |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2} T'(\mathbf{y}) \Sigma^{-1} T(\mathbf{y})\right\} \right), \mathbf{x} \in \mathfrak{R}^N, \mathbf{z} \in \mathbb{Z} \subset \mathfrak{R}^M, \mathbf{y} \in \mathfrak{R}^{N+M}. \end{aligned} \quad (7)$$

Поскольку параметрический вид функций $T(\cdot), f(\cdot)$, а также само разбиение вектора $\mathbf{y} \in \mathfrak{R}^p$ на подвекторы не известны для оценивания плотности распределения $p(\mathbf{y})$ по случайной выборке $Y = (\mathbf{y}_i) \in \mathfrak{R}^{pn}$ будем использовать непараметрическую оценку плотности Розенблатта–Парзена с многомерным гауссовским ядром, определяемую по формуле [5]:

$$\hat{p}(\mathbf{y}) = \frac{1}{n} \sum_{j=1}^n n_N(\mathbf{y} | \mathbf{y}_j, h^2 H), \quad (8)$$

где H, h – управляемые параметры гауссовского ядра: $H \in \mathfrak{S}_{N+M}$ – матрица гауссовского ядра; $h \equiv h(n)$ – коэффициент сглаживания, удовлетворяющий условиям состоятельности оценки:

$$h(n) \rightarrow 0, nh(n) \rightarrow \infty \text{ при } n \rightarrow \infty. \quad (9)$$

При построении оценки (8) возникает проблема задания матрицы ядра H . Обычно [5] в качестве матрицы ядра используется либо фиксированная (например, единичная) матрица либо выборочная ковариационная матрица, вычисленная по всей выборке $Y = (\mathbf{y}_i) \in \mathfrak{R}^{pn}$.

В работе [6] при оценивании плотности распределения $p_{\mathbf{y}}(\mathbf{y})$ в условиях (1)–(4) показано преимущество выборочной оценки матрицы $H \in \mathfrak{S}_{N+M}$ и предлагается способ вычисления коэффициента сглаживания h для модели статистической зависимости (2) линейного вида (то есть, когда Γ является гиперплоскостью). В общем случае (в случае нелинейной гиперповерхности Γ) предлагается использовать оценку (8) с адаптивным гауссовским ядром. Адаптация ядра достигается за счет использования для каждого наблюдения $\mathbf{y}_i (i=1, \dots, n)$ своей матрицы ядра $H^{(i, m(i))}$ (локальной выборочной оценки ковариационной матрицы случайного вектора $\mathbf{y} \in \mathfrak{R}^p$), вычисленной в некоторой окрестности $S(i) \in \mathfrak{R}^p$ наблюдения \mathbf{y}_i . Оптимальный размер локальной окрестности для точки \mathbf{y}_i , определяемый количеством попавших в нее точек $m(i) - 1$ из выборки Y , находится из условия минимума статистики Андерсона [7], характеризующей степень множественной линейной зависимости компонент вектора $\mathbf{y} \in \mathfrak{R}^p$ в окрестности точки $\mathbf{y}_i (i=1, \dots, n)$. При этом само наблюдение \mathbf{y}_i в вычислении $H^{(i, m(i))}$ не используется.

2. Задачи исследования. Пусть сложная система характеризуется случайным вектором $\mathbf{y} \in \mathfrak{R}^p$, описываемым моделью (1)–(4), и имеют место два режима функционирования, которым соответствуют два класса состояний системы Ω_1 или Ω_2 . Номер класса состояния системы в i -м эксперименте описывается ненаблюдаемой случайной величиной $v_i = v(\mathbf{y}_i) \in S = \{1, 2\} (i = 1, \dots, n, \dots)$ с распределением вероятностей

$$P\{v_i = \alpha\} = \pi_\alpha > 0 (\alpha \in S), \pi_1 + \pi_2 = 1, \quad (10)$$

параметры $\{\pi_\alpha\}(\alpha \in S)$ являются априорными вероятностями классов состояний системы.

Классам $\{\Omega_\alpha\}$ в модели (2) соответствуют неизвестные функции $\{f_\alpha(\mathbf{z})\}$ (или $\{T_\alpha(\mathbf{y})\}$), удовлетворяющие условию

$$P(f_1(\mathbf{z}) = f_2(\mathbf{z})) = 0, \quad \forall \mathbf{z} \in \mathbb{Z}, \quad (11)$$

которое означает, что для различных классов состояний гиперповерхности $\{\Gamma_\alpha\}(\alpha \in S)$ различны.

Условная плотность распределения $p_\alpha(\mathbf{y})$ случайного вектора $\mathbf{y} \in \mathfrak{R}^p$ для класса состояния системы Ω_α имеет вид (7) при $f(\mathbf{z}) \equiv f_\alpha(\mathbf{z})(T(\mathbf{y}) \equiv T_\alpha(\mathbf{y})), \alpha \in S$.

Вероятностные характеристики классов $\{\pi_\alpha, p_\alpha(\mathbf{y})\}(\alpha \in S)$ не известны. Имеется классифицированная обучающая выборка наблюдений $Y = (\mathbf{y}_i) \in \mathfrak{R}^{pn}$, допускающая разбиение на подвыборки наблюдений из классов Ω_1 и Ω_2 : $Y = Y_1 \cup Y_2$, где $Y_\alpha = (\mathbf{y}_{\alpha i}) \in \mathfrak{R}^{pn_\alpha}$ – выборка наблюдений из класса Ω_α ($\alpha \in S, n = n_1 + n_2$).

Имеют место следующие задачи исследования:

1) построить непараметрическое решающее правило, предназначенное для принятия решения относительно класса состояния сложной системы по заданному вектору наблюдений $\mathbf{y} \in \mathfrak{R}^p$, удовлетворяющему условиям (1)–(4);

2) в случае известной эндогенно-экзогенной структуры модели с помощью метода асимптотических разложений в асимптотике растущего объема обучающей выборки и усиливающейся статистической зависимости компонент вектора признаков исследовать условный риск предлагаемого решающего правила при условии, что значения экзогенных переменных являются заданными.

3. Решающие правила и критерии оптимальности. Оптимальное в смысле минимума риска байесовское решающее правило (БРП) классификации наблюдений из классов $\{\Omega_\alpha\}$ с вероятностными характеристиками $\{\pi_\alpha, p_\alpha(\mathbf{y})\}(\alpha \in S)$ является нерандомизированным и имеет вид [5]:

$$d(\mathbf{y}) = \mathbf{1}(G(\mathbf{y})) + 1 = \begin{cases} 1, & \text{если } G(\mathbf{y}) < 0, \\ 2, & \text{если } G(\mathbf{y}) \geq 0, \end{cases} \quad (12)$$

где $\mathbf{1}(\cdot)$ – единичная функция Хевисайда, $G(\cdot)$ – байесовская дискриминантная функция, определяемая соотношениями:

$$G(\mathbf{y}) = c_2 p_2(\mathbf{y}) - c_1 p_1(\mathbf{y}), \quad c_1 = \pi_1(w_{12} - w_{11}), \quad c_2 = (1 - \pi_1)(w_{21} - w_{22}), \quad (13)$$

где $W = (w_{\alpha\beta})(\alpha, \beta \in S)$ – заданная матрица потерь.

Будем рассматривать альтернативные решающие правила, получаемые на основании БРП (11), (12), подстановкой в него вместо неизвестных истинных условных плотностей распределений $\{p_\alpha(\mathbf{y})\}$ их непараметрических оценок $\{\hat{p}_\alpha(\mathbf{y})\}$ вида (8), вычисленных по выборкам $\{Y_\alpha\}$ ($\alpha \in S$).

В рамках аналитического исследования риска подстановочных решающих правил будем предполагать, что модель зависимости (2) определяет в пространстве \mathfrak{R}^p некоторую гиперплоскость (многообразие) Γ размерности $N < p$:

$$T_\alpha(\mathbf{y}_i) \equiv T_\alpha \mathbf{y}_i = \mathbf{x}_i - B_\alpha \mathbf{z}_i = \xi_i, \quad i = 1, \dots, n, \quad (\alpha \in S) \quad (14)$$

где $T_\alpha = \left(\mathbf{I}_N \mid B_\alpha \right)$ – фиксированная $(N \times p)$ -матрица, B_α – неизвестная фиксированная $(N \times M)$ -матрица, удовлетворяющая условию разделимости классов типа (11):

$$P\left(B_1 \mathbf{z} = B_2 \mathbf{z}\right) = 0 \quad \forall \mathbf{z} \in \mathbb{Z}. \quad (15)$$

Объектом исследования в данной статье являются решающие правила, использующие непараметрические оценки плотности $\left\{ \hat{p}_\alpha^{(l)}(\mathbf{y}) \right\} (l=0,1,2, \alpha \in S)$, различающиеся выбором матрицы ядра H : оценки $\left\{ \hat{p}_\alpha^{(1)}(\mathbf{y}) \right\}$ имеют место, если $H \equiv H_1 \in \mathfrak{S}_{N+M}$ – произвольная фиксированная матрица; в случае оценок $\left\{ \hat{p}_\alpha^{(2)}(\mathbf{y}) \right\}$ для выборочных наблюдений $\left\{ \mathbf{y}_{\alpha i} \right\} (i=1, \dots, n_\alpha)$ используются матричные статистики $H^{(i,m(\alpha,i))} \in \mathfrak{S}_{N+M}$, вычисленные по $m(\alpha,i)$ наблюдениям из локальной окрестности $S(\alpha,i)$ точки $\mathbf{y}_{\alpha i}$ в соответствии с описанной в п. 1 процедурой. В теоретических исследованиях будем также использовать оценки $\left\{ \hat{p}_\alpha^{(0)}(\mathbf{y}) \right\}$, в которых матрица ядра совпадает с ковариационной матрицей $H_0 \in \mathfrak{S}_{N+M}$ случайного вектора $\mathbf{y} \in \mathfrak{R}^p$. Для вычисления перечисленных оценок плотностей справедливы следующие формулы:

$$\begin{aligned} \hat{p}_\alpha^{(l)}(\mathbf{y}) &= \frac{1}{n_\alpha} \sum_{i=1}^{n_\alpha} n_N(\mathbf{y} \mid \mathbf{y}_{\alpha i}, h^2 H_l) \quad (l=0,1), \quad \hat{p}_\alpha^{(2)}(\mathbf{y}) = \frac{1}{n_\alpha} \sum_{i=1}^{n_\alpha} n_N(\mathbf{y} \mid \mathbf{y}_{\alpha i}, h^2 H^{(i,m(\alpha,i))}), \quad (16) \\ H^{(i,m(\alpha,i))} &= \frac{1}{m(\alpha,i) - 1} \sum_{\mathbf{y}_{\alpha j} \in S(\alpha,i) \setminus \mathbf{y}_{\alpha i}} (\mathbf{y}_{\alpha j} - \mathbf{y}^{(\alpha,i)}) (\mathbf{y}_{\alpha j} - \mathbf{y}^{(\alpha,i)})', \quad \mathbf{y}^{(\alpha,i)} = \frac{1}{m(\alpha,i) - 1} \sum_{\mathbf{y}_{\alpha j} \in S(\alpha,i) \setminus \mathbf{y}_{\alpha i}} \mathbf{y}_{\alpha j}. \end{aligned}$$

Согласно (12), подстановочные решающие правила $\hat{d}^{(l)}(\mathbf{y}) (l=0,1,2)$ имеют вид:

$$\hat{d}^{(l)}(\mathbf{y}) = \mathbf{1}(\hat{G}^{(l)}(\mathbf{y})) + 1, \quad \hat{G}^{(l)}(\mathbf{y}) = c_2 \hat{p}_2^{(l)}(\mathbf{y}) - c_1 \hat{p}_1^{(l)}(\mathbf{y}). \quad (17)$$

При решении задачи 2 в качестве критериев оптимальности подстановочных решающих правил $\hat{d}^{(l)}(\mathbf{y}) (l=0,1,2)$ с учетом вида зависимости (2) компонент вектора $\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ \mathbf{z} \end{pmatrix} \in \mathfrak{R}^{N+M}$, ($\mathbf{x} \in \mathbb{X} = \mathfrak{R}^N$, $\mathbf{z} \in \mathbb{Z} \subset \mathfrak{R}^M$) будем использовать следующие функционалы риска (средних потерь):

– условный риск (условное математическое ожидание потерь) в фиксированной точке $\mathbf{z} \in \mathbb{Z}$:

$$r_n^{(l)}(\mathbf{z}) = E_Y \left\{ R_n^{(l)}(\mathbf{z}, Y) \right\}, \quad R_n^{(l)}(\mathbf{z}, Y) = \sum_{\alpha \in S} \pi_\alpha \int_{\mathbb{X}} w(\alpha, \hat{d}^{(l)}(\mathbf{y})) p_\xi(\mathbf{x} - f_\alpha(\mathbf{z})) dx; \quad (18)$$

– условный ε -риск в фиксированной точке $\mathbf{z} \in \mathbb{Z}$:

$$r_n^{(l)}(\varepsilon, \mathbf{z}) = E_Y \left\{ R_n^{(l)}(\varepsilon, \mathbf{z}, Y) \right\} (l=1,2), \quad R_n^{(l)}(\varepsilon, \mathbf{z}, Y) = \sum_{\alpha \in S} \pi_\alpha \int_{T(\varepsilon, \mathbf{z})} w(\alpha, \hat{d}^{(l)}(\mathbf{y})) p_\xi(\mathbf{x} - f_\alpha(\mathbf{z})) dx,$$

где $T(\varepsilon, \mathbf{z}) \subset \mathbb{X}$ – ограниченная область, удовлетворяющая условию

$$\left| r_n^{(l)}(\varepsilon, \mathbf{z}) - r_n^{(l)}(\mathbf{z}) \right| \leq \varepsilon, \quad (19)$$

здесь величина $\varepsilon > 0$ задает точность приближения $r_n^{(l)}(\varepsilon, \mathbf{z})$ к $r_n^{(l)}(\mathbf{z})$. Согласно (19) использование условного ε -риска обеспечивает возможность оценки условного риска классификации с заданной

точностью, определяемой величиной ε ($0 < \varepsilon < 1$). При $\varepsilon \rightarrow 0$ точность вычисления условного риска возрастает.

Заметим, что в случае «антиединичной» матрицы потерь W (когда $w_{\alpha\beta} = 1 - \delta_{\alpha\beta}$, где $\delta_{\alpha\beta}$ – символ Кронекера) функционалы (17), (18) имеют смысл условных вероятностей ошибок классификации.

Можно показать [4], что условие (19) справедливо, если $P(\mathbf{x} \in T(\varepsilon, z)) \geq 1 - \varepsilon / (c_1 + c_2) \forall \mathbf{z} \in \mathbb{Z}$. В этом случае имеет место выражение для условного риска:

$$r_n^{(l)}(\mathbf{z}) = \pi_1(w(1, 1)) + \pi_2(w(2, 1)) - \int_{\mathbb{X}} E_Y \left\{ \mathbf{1}(\tilde{G}^{(l)}(\mathbf{y})) \right\} \tilde{G}(\mathbf{y}) d\mathbf{x} \quad (l=1,2), \quad (20)$$

где $\tilde{G}(\mathbf{y}) = c_2 p_{\xi}(\mathbf{x} - f_2(\mathbf{z})) - c_1 p_{\xi}(\mathbf{x} - f_1(\mathbf{z}))$ – дискриминантная функция БРП (12), (13) с учетом рассматриваемой модели зависимости компонент вектора наблюдений (2).

4. Асимптотические разложения риска подстановочных решающих правил. В условиях рассматриваемой модели зависимости компонент вектора $\mathbf{y} \in \mathfrak{Y}^p$ методом асимптотического разложения исследуем условные риски описанных выше подстановочных решающих правил $\{\tilde{d}^{(l)}(\mathbf{y})\}$. Сформулируем вначале вспомогательные утверждения.

Будем использовать обозначения: $H_0 \in \mathfrak{S}_{N+M}$ и $Q \in \mathfrak{S}_M$ – ковариационные матрицы случайных векторов $\mathbf{y} \in \mathfrak{Y} = \mathfrak{R}^p$ и $\mathbf{z} \in \mathbb{Z}$ соответственно; $b_{n,\alpha}^{(l)}(\mathbf{y})$ и $v_{n,\alpha}^{(l)}(\mathbf{y})$ – смещение и вариация оценки плотности $\tilde{p}_{\alpha}^{(l)}(\mathbf{y})$ вида (16) соответственно; $\tilde{G}^{(l)}(\mathbf{y})$ – дискриминантная функция ПБРП $\tilde{d}^{(l)}(\mathbf{y})$, определяемого соотношениями (17); $g^{(l)}(\mathbf{y}) = \tilde{G}^{(l)}(\mathbf{y}) - G(\mathbf{y})$ – случайное отклонение оценки дискриминантной функции при использовании ПБРП $\tilde{d}^{(l)}(\mathbf{y})$ от дискриминантной функции вида (13) БРП; $T_{\alpha} = (\mathbf{I}_N \mid -B_{\alpha})$ – блочная $(N \times (N+M))$ -матрица, \mathbf{I}_N – единичная $(N \times N)$ -матрица; $\{c_{\alpha}\} (\alpha \in S)$ определяются по формуле (13); $n_1 = n_2 (n = 2n_1)$;

$$k_l = \frac{1}{(4\pi)^{(N+M)/2} \sqrt{\det(H_e)}} \quad (l=0,1), \quad k_2 = \frac{1}{(4\pi)^{(N+M)/2} \sqrt{\det(\Sigma) \det(Q)}}. \quad (21)$$

Л е м м а 1. Если выполняются условия (14), (15), (6), то для фиксированного значения $\mathbf{y} \in \mathfrak{Y} = \mathfrak{R}^p$ случайные отклонения $\{g^{(l)}(\mathbf{y})\} (l=0,1,2)$ имеют моменты до третьего порядка включительно, причем

$$E \left\{ (g_l(\mathbf{y}))^2 \right\} = \frac{k_l}{n_0 h^{N+M}} \sum_{\alpha \in S} c_{\alpha}^2 p_{\alpha}(\mathbf{y}) + (c_1 b_{n,1}^{(l)}(\mathbf{y}) - c_0 b_{n,0}^{(l)}(\mathbf{y}))^2, \quad (22)$$

$$E \left\{ (g_l(\mathbf{y}))^3 \right\} = o_1, \quad o_1 = o(h^4 + n^{-1} h^{-(N+M)}), \quad (23)$$

где

$$b_{n,\alpha}^{(0)}(\mathbf{y}) = -\frac{h^2}{2} (N - \mathbf{y}' T_{\alpha}' \Sigma^{-1} T_{\alpha} \mathbf{y}) p_{\alpha}(\mathbf{y}) + o(h^2), \quad b_{n,\alpha}^{(2)}(\mathbf{y}) = b_{n,\alpha}^{(0)}(\mathbf{y}) + o(h^2), \quad (24)$$

$$b_{n,\alpha}^{(1)}(\mathbf{y}) = -\frac{h^2}{2} \text{tr} \left(T_{\alpha} H_l T_{\alpha}' \Sigma^{-1} (\mathbf{I}_N - T_{\alpha} \mathbf{y} \mathbf{y}' T_{\alpha}' \Sigma^{-1}) \right) p_{\alpha}(\mathbf{y}) + o(h^2), \quad (25)$$

$$v_{n,\alpha}^{(0)}(\mathbf{y}) = \frac{p_{\alpha}(\mathbf{y}) n^{-1} h^{-(N+M)}}{(4\pi)^{(N+M)/2} \sqrt{|H_1| |\Sigma|}} + b_{n,\alpha}^{(0)}(\mathbf{y}) + o_1, \quad v_{n,\alpha}^{(2)}(\mathbf{y}) = v_{n,\alpha}^{(0)}(\mathbf{y}) + o_1, \quad (26)$$

$$v_{n,\alpha}^{(1)}(\mathbf{y}) = \frac{p_\alpha(\mathbf{y})n^{-1}h^{-(N+M)}}{(4\pi)^{(N+M)/2}\sqrt{|H_1|}} + b_{n,\alpha}^{(l)}(\mathbf{y}) + o_1. \quad (27)$$

Доказательство. Асимптотические разложения (24)–(27) были получены при сделанных в лемме 1 предположениях в работе [8]. Пусть $\zeta_{\alpha l}$ – случайная величина с плотностью распределения $p_{\alpha l}(\zeta_{\alpha l})$ такая, что для фиксированного значения $\mathbf{y} \in \mathbb{Y}$:

$$\zeta_{\alpha l} \equiv \zeta_{\alpha l}(\mathbf{y}) = \widehat{p}_\alpha^{(l)}(\mathbf{y}) - p_\alpha(\mathbf{y}) \quad (l=0,1,2, \alpha \in S). \quad (28)$$

Тогда для отклонений $\{g_l(\mathbf{y})\} (l=0,1,2)$ справедливо представление

$$g_l(\mathbf{y}) = c_2 \zeta_{2l} - c_1 \zeta_{1l}. \quad (29)$$

С учетом (29) искомые моменты случайных величин $\{g_l(\mathbf{y})\}$ выражаются через моменты первого, второго и третьего порядка для взаимно независимых по построению случайных величин $\{\zeta_{\alpha l}\}$. С учетом (24)–(27) имеем: $E\{\zeta_{\alpha l}\} = b_{n,\alpha}^{(l)}(\mathbf{y}) = O(h^2)$, $E\{\zeta_\alpha^2\} = v_{n,\alpha}^{(l)}(\mathbf{y}) = O(h^4 + n^{-1}h^{-(N+M)})$, $E\{\zeta_\alpha^3\} = o_1$. На основании (25) и свойств математического ожидания случайных величин получаем формулы (22), (23). Лемма 1 доказана.

Введем обозначения: $G_1(\cdot; \mathbf{z}) : \mathbf{x} \rightarrow \mathfrak{R}^1$ – байесовская дискриминантная функция для фиксированной точки $\mathbf{z} \in \mathbb{Z}$;

$$\Gamma_z = \{\mathbf{x} : G_1(\mathbf{x}; \mathbf{z}) = 0\} \quad (30)$$

– соответствующая дискриминантная гиперповерхность; $T(\varepsilon, \mathbf{z}) \subset \mathbb{X}$ – ограниченная область, представляющая собой объединение двух гиперэллипсоидов:

$$T(\varepsilon, \mathbf{z}) = T^{(1)}(\varepsilon, \mathbf{z}) \cup T^{(2)}(\varepsilon, \mathbf{z}), \quad T^{(\alpha)}(\varepsilon, \mathbf{z}) \equiv \left\{ \mathbf{x} : (\mathbf{x} - B_\alpha \mathbf{z})' \Sigma^{-1} (\mathbf{x} - B_\alpha \mathbf{z}) \leq u_{\varepsilon, \alpha}^2 \right\} \quad (\varepsilon > 0), \quad (31)$$

где $u_{\varepsilon, \alpha}^2$ – квантиль уровня $1-\varepsilon$ для χ^2 -распределения с N степенями свободы;

$l_\alpha(\mathbf{x}, \mathbf{z}) = \sqrt{(\mathbf{x} - B_\alpha \mathbf{z})' \Sigma^{-1} (\mathbf{x} - B_\alpha \mathbf{z})}$ ($\alpha \in S$) – расстояние Махаланобиса от точки $\mathbf{x} \in \mathbb{X}$ до центра гиперэллипсоида $T^{(\alpha)}(\varepsilon, \mathbf{z})$;

$$l_\alpha(\mathbf{z}) = (-1)^\alpha \frac{\Delta^2(\mathbf{z})}{2} - \ln \left(\frac{c_1}{c_2} \right); \quad (32)$$

$\Delta(\mathbf{z}) = \sqrt{\mathbf{z}' B' \Sigma^{-1} B \mathbf{z}}$ ($B = B_2 - B_1$) – расстояние Махаланобиса между центрами гиперэллипсоидов $T^{(0)}(\varepsilon, \mathbf{z})$ и $T^{(1)}(\varepsilon, \mathbf{z})$; $\Delta_1(\mathbf{z}) = \mathbf{z}' B' \Sigma^{-2} B \mathbf{z}$.

Представим матрицу $\Sigma^{-1} = (\bar{\sigma}_{ij})$ в блочном виде, а векторы $\boldsymbol{\beta}(\mathbf{z}) = (\beta_k(\mathbf{z})) = \Sigma^{-1} B \mathbf{z} \in \mathfrak{R}^N$, $\mathbf{b}_\alpha(\mathbf{z}) = (b_{\alpha j}(\mathbf{z})) = B_\alpha \mathbf{z} \in \mathfrak{R}^N$ ($N > 1$) разобьем на подвекторы:

$$\Sigma^{-1} = \begin{pmatrix} \tilde{\Sigma}^{-1} & \tilde{\boldsymbol{\sigma}}_N \\ \tilde{\boldsymbol{\sigma}}_N' & \tilde{\sigma}_{NN} \end{pmatrix}, \quad \boldsymbol{\beta}(\mathbf{z}) = \begin{pmatrix} \tilde{\boldsymbol{\beta}}(\mathbf{z}) \\ \tilde{\beta}_N(\mathbf{z}) \end{pmatrix}, \quad \mathbf{b}_\alpha(\mathbf{z}) = \begin{pmatrix} \mathbf{b}_i(\mathbf{z}) \\ \mathbf{b}_{iN}(\mathbf{z}) \end{pmatrix}.$$

Л е м м а 2. Если модель зависимости компонент случайного вектора $\mathbf{y} \in \mathbb{Y} = \mathfrak{R}^p$ определяется соотношениями (14), (15), (6) и область $T(\varepsilon, \mathbf{z}) \subset \mathfrak{X}$ имеет вид (31), то область пересечения $\tilde{T}^{(\alpha)}(\varepsilon, \mathbf{z}) \equiv \Gamma_z \cap T^{(\alpha)}(\varepsilon, \mathbf{z})$ ($\varepsilon > 0$) является гиперэллипсоидом в пространстве $\tilde{\mathfrak{X}} = \mathfrak{R}^{N-1}$ $\left(\tilde{\mathbf{x}} = (x_1, \dots, x_{N-1})' \in \tilde{\mathfrak{X}} \right)$ и описывается выражением:

$$\tilde{T}^{(\alpha)}(\varepsilon, \mathbf{z}) = \left\{ \tilde{\mathbf{x}} : \left(\tilde{\mathbf{x}} - \mu^{(\alpha)}(\mathbf{z}) \right)' F(\mathbf{z}) \left(\tilde{\mathbf{x}} - \mu^{(\alpha)}(\mathbf{z}) \right) \leq \tilde{u}_{\varepsilon, \alpha}^2 \right\}, \quad (33)$$

где

$$F(\mathbf{z}) = \Sigma_1 + \mathbf{a}(\mathbf{z})\mathbf{a}(\mathbf{z})' \in \mathfrak{S}_{N-1}, \quad (34)$$

$$\Sigma_1 = \tilde{\Sigma}^{-1} - \frac{1}{\bar{\sigma}_{NN}} \tilde{\sigma}_N \tilde{\sigma}_N', \quad \mathbf{a}(\mathbf{z}) = \left(\beta_N(\mathbf{z}) \sqrt{\bar{\sigma}_{NN}} \right)^{-1} \left(\bar{\sigma}_{NN} \tilde{\beta}_N(\mathbf{z}) + \beta_N(\mathbf{z}) \tilde{\sigma}_N \right),$$

и имеет место:

$$\mu^{(\alpha)}(\mathbf{z}) = -\sqrt{\bar{\sigma}_{NN}} \left(\frac{c(\mathbf{z})}{\beta_N(\mathbf{z})} + \beta_{\alpha N}(\mathbf{z}) \right) F(\mathbf{z}) \mathbf{a}(\mathbf{z}) - \left(\Sigma_1 - \frac{1}{\beta_N(\mathbf{z})} \tilde{\beta}_N(\mathbf{z}) \tilde{\sigma}_N' \right) \tilde{\mathbf{b}}_\alpha(\mathbf{z}), \quad (35)$$

$$\tilde{u}_{\varepsilon, \alpha}^2(\mathbf{z}) = u_{\varepsilon, \alpha}^2 - l_\alpha^2(\mathbf{z}) / \Delta^2(\mathbf{z}) \quad (\alpha \in S). \quad (36)$$

Д о к а з а т е л ь с т в о. С учетом (14), (15), (6) дискриминантная функция $G_1(\mathbf{x}; \mathbf{z})$ имеет вид:

$$G_1(\mathbf{x}; \mathbf{z}) = \mathbf{x}' \Sigma^{-1} B \mathbf{z} + c(\mathbf{z}), \quad c(\mathbf{z}) = \mathbf{z}' B' \Sigma^{-1} (B_2 + B_1) \mathbf{z} + \ln(c_1 / c_2). \quad (37)$$

Согласно (30) и (37) компоненты составного вектора $\mathbf{x} = (\tilde{\mathbf{x}}, x_N)' \in \Gamma_z$ связаны соотношением:

$$x_N = -\frac{1}{\beta_N(\mathbf{z})} \left(\tilde{\mathbf{x}}' \tilde{\beta}_N(\mathbf{z}) + c(\mathbf{z}) \right).$$

Поэтому квадратичная форма $l_\alpha^2(\mathbf{x}, \mathbf{z})$ допускает представление:

$$l_\alpha^2(\mathbf{x}, \mathbf{z}) = \left(\tilde{\mathbf{x}} - \mu^{(\alpha)}(\mathbf{z}) \right)' F(\mathbf{z}) \left(\tilde{\mathbf{x}} - \mu^{(\alpha)}(\mathbf{z}) \right) + d_\alpha^2(\mathbf{z}), \quad (38)$$

где $F(\mathbf{z})$, $\mu^{(\alpha)}(\mathbf{z})$ имеют вид (34), (35), а

$$\begin{aligned} d_\alpha^2(\mathbf{z}) = & \left(\tilde{\mathbf{b}}_\alpha(\mathbf{z}) \right)' \tilde{\mathbf{b}}_\alpha(\mathbf{z}) + 2 \left(\frac{c(\mathbf{z})}{\beta_N(\mathbf{z})} + b_{\alpha N}(\mathbf{z}) \right) \left(\tilde{\sigma}_N \right)' \tilde{\mathbf{b}}_\alpha(\mathbf{z}) + \\ & + \bar{\sigma}_{NN} \left(\frac{c(\mathbf{z})}{\beta_N(\mathbf{z})} + b_{\alpha N}(\mathbf{z}) \right)^2 - \left(\mu^{(\alpha)}(\mathbf{z}) \right)' F(\mathbf{z}) \mu^{(\alpha)}(\mathbf{z}). \end{aligned}$$

Из (31), (33), (38) следует, что $\tilde{u}_{\varepsilon, \alpha}^2(\mathbf{z}) = u_\varepsilon^2 - d_\alpha^2(\mathbf{z})$. Заметим, что $d_\alpha^2(\mathbf{z})$ – расстояние Махаланобиса между центрами гиперэллипсоидов $\tilde{T}^{(\alpha)}(\varepsilon, \mathbf{z})$ и $T^{(\alpha)}(\varepsilon, \mathbf{z})$, а $\tilde{u}_{\varepsilon, \alpha}^2(\mathbf{z})$ и u_ε^2 соответствен-

но «радиусы» этих гиперэллипсоидов. Поэтому $d_\alpha^2(\mathbf{z}) = l_\alpha^2(\mathbf{x}^*, \mathbf{z})$, где $\mathbf{x}^* \in \Gamma_\varepsilon$ – общий центр областей $\{\tilde{T}^{(\alpha)}(\varepsilon, \mathbf{z})\} (\alpha \in S)$. Решая систему из двух уравнений

$$\begin{cases} l_2^2(\mathbf{x}^*, \mathbf{z}) - l_1^2(\mathbf{x}^*, \mathbf{z}) = -2 \ln(c_1 / c_2), \\ l_2^2(\mathbf{x}^*, \mathbf{z}) + l_1^2(\mathbf{x}^*, \mathbf{z}) = \Delta^2(\mathbf{z}), \end{cases}$$

получаем $d_\alpha^2(\mathbf{z}) = l_\alpha^2(\mathbf{z}) / \Delta^2(\mathbf{z})$ ($\alpha \in S$), что влечет (36). Лемма 2 доказана.

Докажем основные утверждения. Будем использовать обозначения: $r(\varepsilon, \mathbf{z})$, $r_n^{(l)}(\varepsilon, \mathbf{z})$ ($r(\mathbf{z})$, $r_n^{(l)}(\mathbf{z})$) – условный ε -риск (условный риск) в точке $\mathbf{z} \in \mathbb{Z} \subset \mathfrak{R}^M$ соответственно БРП (12) и ПБРП (17) ($l = 0, 1, 2$); $m = 3 - \arg \min_{\alpha \in S} \{c_\alpha\}$, где $\{c_\alpha\}$ определяются по формулам (13); $p_{\alpha l}^{(k)}(\zeta_{\alpha l})$ – k -я производная функции плотности $p_{\alpha l}(\zeta_{\alpha l})$ ($k = 0, 1, 2$), где $\zeta_{\alpha l} \equiv \zeta_{\alpha l}(\mathbf{y})$, ($l = 0, 1, 2$, $\alpha \in S$) определяются (28); $p_{g_l}(v)$ ($v \in (-\infty, \infty)$) – плотность распределения случайной величины $g_l(\mathbf{y})$ вида (29) при фиксированном $\mathbf{y} \in \mathfrak{R}^p$.

Для $t \in \mathfrak{R}^1$ и фиксированной точки $\mathbf{z} \in \mathbb{Z} \subset \mathfrak{R}^M$ определим функции:

$$Q_l(t, \mathbf{z}) = \int_{\mathfrak{X}} E \left\{ \mathbf{1} \left(G_1(\mathbf{x}; \mathbf{z}) + t g^{(l)}(\mathbf{y}) \right) \right\} G_1(\mathbf{x}; \mathbf{z}) d\mathbf{x} \quad (l = 1, 2). \quad (39)$$

Т е о р е м а 1. Пусть модель статистической зависимости компонент вектора $\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ \mathbf{z} \end{pmatrix} \in \mathbb{Y} = \mathfrak{R}^p$ ($\mathbf{x} \in \mathfrak{X} = \mathfrak{R}^N$, $\mathbf{z} \in \mathbb{Z} \subset \mathfrak{R}^M$, $\text{mes}\{\mathbb{Z}\} < \infty$) определяется плотностью распределения вида (7), в которой $T_\alpha(\mathbf{y}) = T_\alpha \mathbf{y} = \mathbf{x} - B_\alpha \mathbf{z}$, где B_α ($\alpha \in S$) – фиксированная $(N \times M)$ -матрица, удовлетворяющая условию (15). Тогда, если

$$\left| p_{\alpha l}^{(k)}(\zeta_{\alpha l}) \right| < \infty \quad \forall \mathbf{y} \in \mathbb{Y}, \quad k = 0, 1, 2, \quad l = 1, 2, \quad \alpha \in S, \quad (40)$$

то для риска подстановочных решающих правил $\{\tilde{d}^{(l)}(\mathbf{y})\} (l = 1, 2)$ вида (17) при $n \rightarrow \infty$ справедливы асимптотические разложения:

$$r_n^{(l)}(\mathbf{z}) = r(\mathbf{z}) + \frac{1}{2} \text{mes}\{\mathbb{Z}\} \left(\frac{k_l}{2n_0 h^{N+M}} \sum_{\alpha \in S} c_\alpha \bar{\alpha}_l(\mathbf{z}) + \bar{\beta}_n^{(l)}(\mathbf{z}) \right) + o_1, \quad (41)$$

$$\bar{\alpha}_l(\mathbf{z}) = \int_{\Gamma_z} c_\alpha p_\alpha(\mathbf{y}) |G_1(\mathbf{x}; \mathbf{z})|^{-1} d\gamma_{N-1}, \quad (42)$$

$$\bar{\beta}_n^{(l)}(\mathbf{z}) = \int_{\Gamma_z} (c_1 b_{n,1}^{(l)}(\mathbf{y}) - c_0 b_{n,0}^{(l)}(\mathbf{y}))^2 |G_1(\mathbf{x}; \mathbf{z})|^{-1} d\gamma_{N-1}, \quad (43)$$

интегралы в (42), (43) есть поверхностные интегралы по $(N-1)$ -мерной поверхности $\Gamma_z = \{\mathbf{x} : G_1(\mathbf{x}; \mathbf{z}) = 0\}$, $d\gamma_{N-1}$ – элемент поверхности Γ_z , $|G_1(\mathbf{x}; \mathbf{z})|$ – норма вектора первых производных от байесовской дискриминантной функции для фиксированной точки $\mathbf{z} \in \mathbb{Z}$.

Доказательство. Выражение (18) для условных рисков $\{r_n^{(l)}(\mathbf{z})\}$ (при условии, что $\mathbf{z} \in \mathbb{Z}$ – фиксированное значение) с учетом (20) и (39) допускает представление:

$$r_n^{(l)}(\mathbf{z}) = \pi_1(w(1, 1)) + \pi_2(w(2, 1)) - \text{mes}\{\mathbb{Z}\} Q_l(1, \mathbf{z}), \mathbf{z} \in \mathbb{Z} (l=1,2). \quad (44)$$

Аналогично для условного риска БРП (12) имеем:

$$r(\mathbf{z}) = \pi_1(w(1, 1)) + \pi_2(w(2, 1)) - \text{mes}\{\mathbb{Z}\} Q_l(0, \mathbf{z}), \forall l=0,1,2. \quad (45)$$

Аналогично [4] можно показать, что если выполняется условие (40), то функция $Q_l(t, \mathbf{z})$ ($\mathbf{z} \in \mathbb{Z}$) трижды дифференцируема по t в $(0,1)$, при этом:

$$Q_l^{(3)}(t, \mathbf{z}) = \int_{\mathbb{X}} Q_l^*(t, \mathbf{z}) G_1(\mathbf{x}; \mathbf{z}) d\mathbf{x} \quad (l=1,2), \quad Q_l^*(t, \mathbf{z}) = \int_{-\infty}^{\infty} v^3 p_{g_l}(v) \delta^{(2)}(G_1(\mathbf{x}; \mathbf{z}) + tv) d\mathbf{x} \quad (l=1,2), \quad (46)$$

($\delta^{(2)}(\cdot)$ – производная 2-го порядка дельта-функции) и при $n \rightarrow \infty$:

$$|Q_l^{(3)}(t, \mathbf{z})| < \infty, (l=0,1,2). \quad (47)$$

На этом основании применим к функции $Q_l(t, \mathbf{z})$ ($\mathbf{z} \in \mathbb{Z}$) формулу Тейлора в окрестности точки $t=0$ ($Q_l^{(k)}(\cdot, \mathbf{z})$ – производная k -го порядка функции $Q_l(\cdot, \mathbf{z})$):

$$Q_l(t, \mathbf{z}) = Q_l(0, \mathbf{z}) + Q_l^{(1)}(0, \mathbf{z})t + Q_l^{(2)}(0, \mathbf{z})\frac{t^2}{2} + Q_l^{(3)}(t_1, \mathbf{z})\frac{t^3}{6} \quad (z \in \mathbb{Z}), \text{ где } 0 < t_1 < t. \quad (48)$$

На основании (39), а также свойств дельта-функции и ее производной $u\delta(u) = 0$, $u\delta^{(1)}(u) = -\delta(u)$ [9] получаем:

$$Q_0(0, \mathbf{z}) = \int_{\mathbb{X}} E\{\mathbf{1}(G_1(\mathbf{x}; \mathbf{z}))\} G_1(\mathbf{x}; \mathbf{z}) d\mathbf{x},$$

$$Q_l^{(1)}(0, \mathbf{z}) = 0, \quad Q_l^{(2)}(0, \mathbf{z}) = - \int_{\mathbb{X}} E\left\{\left(g^{(l)}(\mathbf{y})\right)^2\right\} \delta(G(\mathbf{y})) d\mathbf{x}. \quad (49)$$

Используя лемму 1 и формулы (46), (47), получаем $Q_l^{(3)}(t_1, \mathbf{z}) = o_1$ для $t_1 \in (0,1)$ ($l=0,1,2$). На основании леммы 3.3 из [4], преобразуем выражение для $Q_l^{(2)}(0, \mathbf{z})$ в (49) к виду:

$$Q_l^{(2)}(0, \mathbf{z}) = - \int_{\Gamma_z} E\left\{\left(g^{(l)}(\mathbf{y})\right)^2\right\} |\nabla_x G_1(\mathbf{x}; \mathbf{z})|^{-1} d\gamma_{N-1}. \quad (50)$$

Подставляя (49), (50) в (48) и затем, используя полученное выражение в (44), с учетом (45), получаем:

$$r_n^{(l)}(\mathbf{z}) = r(\mathbf{z}) + \frac{1}{2} \text{mes}\{\mathbb{Z}\} \int_{\Gamma_z} E\left\{\left(g^{(l)}(\mathbf{y})\right)^2\right\} |\nabla_x G_1(\mathbf{x}; \mathbf{z})|^{-1} d\gamma_{N-1} + o_1 \quad (l=1,2), \quad (51)$$

откуда на основании (22), (23) следует (41)–(43). Теорема 1 доказана.

Поскольку поверхность $\Gamma_z = \{\mathbf{x} : G_1(\mathbf{x}; \mathbf{z}) = 0\}$ имеет неограниченную площадь, то интегралы в (42), (43) являются неограниченными. Вследствие чего, вместо разложения (41) будем использовать разложение условного ε -риска:

$$r_n^{(l)}(\varepsilon, \mathbf{z}) = r(\varepsilon, \mathbf{z}) + \frac{1}{2} \text{mes}\{\mathbb{Z}\} \left(\frac{k_l}{2n_0 h^{N+M}} \sum_{\alpha \in S} c_\alpha \alpha_\alpha(\mathbf{z}) + \beta_n^{(l)}(\mathbf{z}) \right) + o_1 \quad (l=1,2), \quad (52)$$

в котором интегралы $\{\alpha_i(\mathbf{z})\}$, $\beta_n^{(l)}(\mathbf{z})$ отличаются от интегралов $\{\bar{\alpha}_i(\mathbf{z})\}$, $\bar{\beta}_n^{(l)}(\mathbf{z})$ областью интегрирования, имеющей с учетом (30), (31) вид: $\Gamma_{\mathbf{z}} \cap T(\varepsilon, \mathbf{z}) = (\Gamma_{\mathbf{z}} \cap T^{(1)}(\varepsilon, \mathbf{z})) \cup (\Gamma_{\mathbf{z}} \cap T^{(2)}(\varepsilon, \mathbf{z}))$.

Введем обозначения: $\eta_\alpha \in \mathfrak{R}^N$ ($\alpha \in S$) – случайный гауссовский вектор с математическим ожиданием и ковариационной матрицей $\mu_\eta(\alpha, \mathbf{z})$, $\Sigma_\eta(\alpha, \mathbf{z})$, определяемыми по формулам:

$$\mu_\eta(\alpha, \mathbf{z}) = -\frac{l_\alpha(\mathbf{z})}{\Delta^2(\mathbf{z})} \boldsymbol{\beta}^*(\mathbf{z}), \quad \Sigma_\eta(\alpha, \mathbf{z}) = \mathbf{I}_N - \frac{1}{\Delta^2(\mathbf{z})} \boldsymbol{\beta}^*(\mathbf{z}) (\boldsymbol{\beta}^*(\mathbf{z}))', \quad \boldsymbol{\beta}^*(\mathbf{z}) = \Sigma^{-\frac{1}{2}} B \mathbf{z} \in \mathfrak{R}^N; \quad (53)$$

$$\tau_{m,\varepsilon}^{(l)}(\mathbf{z}, A) = \int_{\Theta(\varepsilon)} (\boldsymbol{\tau}' A \boldsymbol{\tau})^k n_N(\boldsymbol{\tau} | \mu_\eta(\alpha, \mathbf{z}), \Sigma_\eta(\alpha, \mathbf{z})) d\boldsymbol{\tau}, \quad \Theta(\varepsilon) = \{\boldsymbol{\tau} : \boldsymbol{\tau}' \boldsymbol{\tau} \leq u_\varepsilon^2\}; \quad (54)$$

$$\tilde{\Psi}^{(l)} = \Psi_2^{(l)} - \Psi_1^{(l)}, \quad \Psi_\alpha^{(l)} = T_\alpha H_l T_\alpha' \quad (l=0,1,2), \quad C = \text{tr} \left(\Psi^{(1)} + 2(-1)^m \frac{\ln(c_1/c_2)}{N} \Psi_m^{(1)} \Sigma^{-1} \right), \quad (55)$$

где $\{T_\alpha\}$, $\Sigma \in \mathfrak{Z}_N$ – определенные выше фиксированные матрицы, $m = 3 - \arg \min_{\alpha \in S} \{c_\alpha\}$.

Т е о р е м а 2. Пусть в условиях теоремы 1 для фиксированного значения $\mathbf{z} \in \mathbb{Z}$ область $T(\varepsilon, \mathbf{z}) \subset \mathbb{X}$ ($\varepsilon > 0$) имеет вид (31), тогда для функционала условного ε -риска подстановочных решающих правил $\{\hat{d}^{(l)}(\mathbf{y})\}$ ($l=1,2$), определяемых (17), при $n \rightarrow \infty$ справедливы асимптотические разложения:

$$r_n^{(l)}(\varepsilon, \mathbf{z}) = r(\varepsilon, \mathbf{z}) + \frac{k_l (c_1 + c_2) \text{mes}\{\tilde{T}^{(m)}(\varepsilon, \mathbf{z})\} \text{mes}\{\mathbb{Z}\}}{2\Delta_1(\mathbf{z})} n^{-1} h^{-(N+M)} + \frac{c_m}{4\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \frac{l_m(\mathbf{z})}{\Delta^2(\mathbf{z})} \right\} \beta^{(l)}(\mathbf{z}, A) h^4 + o_1, \quad (56)$$

где $k_1, k_2 > 0$ определяются (21), $o_1 = o(h^4 + n^{-1} h^{-(N+M)})$, $A = \Sigma^{-\frac{1}{2}} \tilde{\Psi}^{(l)} \Sigma^{-\frac{1}{2}}$,

$$\beta^{(1)}(\mathbf{z}, A) = C^2 \tau_{m,\varepsilon}^{(0)}(\mathbf{z}, A) + 2C \tau_{m,\varepsilon}^{(1)}(\mathbf{z}, A) + \tau_{m,\varepsilon}^{(2)}(\mathbf{z}, A) \geq 0, \quad (57)$$

$$\beta^{(2)}(\mathbf{z}, A) = 4 \left(\ln \left(\frac{c_1}{c_2} \right) \right)^2 \tau_{m,\varepsilon}^{(0)}(\mathbf{z}, A) \geq 0, \quad (58)$$

Д о к а з а т е л ь с т в о. Согласно (31) $P\{\mathbf{x} \in T^{(\alpha)}(\varepsilon, \mathbf{z})\} = 1 - \varepsilon$ (здесь \mathbf{x} – случайный вектор с плотностью распределения $n_N(\mathbf{x} | B_\alpha \mathbf{z}, \Sigma)$) и следовательно справедливо (20). В условиях теоремы $\forall \mathbf{x} \in \Gamma_{\mathbf{z}}$ справедливо: $c_1 n_N(\mathbf{x} | B_1 \mathbf{z}, \Sigma) = c_0 n_N(\mathbf{x} | B_0 \mathbf{z}, \Sigma)$. Поэтому из (37) для $\mathbf{x} \in \Gamma_{\mathbf{z}}$ и $\alpha \in S$ следует:

$$|\nabla_{\mathbf{x}} G_1(\mathbf{x}; \mathbf{z})|^{-1} = 1 / \Delta_1(\mathbf{z}) c_\alpha P_\alpha(\mathbf{y}). \quad (59)$$

Так как области $\tilde{T}^{(\alpha)}(\varepsilon, \mathbf{z}) \equiv \Gamma_{\mathbf{z}} \bigcap T^{(\alpha)}(\varepsilon, \mathbf{z})$ – концентрические гиперэллипсоиды и, согласно (32), (36), $\tilde{T}^{(3-m)}(\varepsilon, \mathbf{z}) \subseteq \tilde{T}^{(m)}(\varepsilon, \mathbf{z})$ (где $m = 1 - \arg \min_{\alpha \in S} \{c_{\alpha}\}$), то $\Gamma_{\mathbf{z}} \bigcap T(\varepsilon, \mathbf{z}) = \tilde{T}^{(m)}(\varepsilon, \mathbf{z})$. Поэтому на основании (42) и (57) имеем:

$$\alpha_{\alpha}(\mathbf{z}) = \text{mes} \left\{ \tilde{T}^{(m)}(\varepsilon, \mathbf{z}) \right\} / \Delta_1(\mathbf{z}) \quad \forall \alpha \in S. \quad (60)$$

Заметим, что из (42) и (60) следует неограниченность $\{\bar{\alpha}_{\alpha}(\mathbf{z})\}$. На основании (24), (25), (43), (52) получаем выражения для $\beta^{(1)}(\mathbf{z}, A)$ и $\beta^{(2)}(\mathbf{z}, A)$:

$$\beta^{(1)}(\mathbf{z}, A) = \frac{h^4 c_m}{4\Delta_1(\mathbf{z})} \left(C^2 J_{m,\varepsilon}^{(0)}(\mathbf{z}, A) + 2C J_{m,\varepsilon}^{(1)}(\mathbf{z}, A) + J_{m,\varepsilon}^{(2)}(\mathbf{z}, A) \right) + o(h^4) \geq 0, \quad (61)$$

$$\beta^{(2)}(\mathbf{z}, A) = \frac{h^4 c_m}{4\Delta_1(\mathbf{z})} \left(\ln \left(\frac{c_1}{c_2} \right) \right)^2 J_{m,\varepsilon}^{(0)}(\mathbf{z}, A) + o(h^4) \geq 0, \quad (62)$$

где для области $\mathbb{T}(m, \varepsilon, \mathbf{z}) \equiv \tilde{T}^{(m)}(\varepsilon, \mathbf{z})$ и $k = 0, 1, 2$:

$$\begin{aligned} J_{m,\varepsilon}^{(k)}(\mathbf{z}, A) &= \int_{\mathbb{T}(m,\varepsilon,\mathbf{z})} \left((\mathbf{x} - B_m \mathbf{z})' \Sigma^{-1/2} A \Sigma^{-1/2} (\mathbf{x} - B_m \mathbf{z}) \right)^k n_N(\mathbf{x} | B_m \mathbf{z}, \Sigma) d\gamma_N = \\ &= \Delta_1(\mathbf{z}) n_1 \left(0 | l_m(\mathbf{z}), \Delta^2(\mathbf{z}) \right) \tau_{m,\varepsilon}^{(k)}(\mathbf{z}, A), \quad n_1 \left(0 | l_m(\mathbf{z}), \Delta^2(\mathbf{z}) \right) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \frac{l_m(\mathbf{z})}{\Delta^2(\mathbf{z})} \right\}. \end{aligned} \quad (63)$$

Подставляя выражения (60)–(63) в (52), получаем (56)–(58). Теорема 2 доказана.

В разложение (56) входят три главных члена: первый характеризует условный ε -риск БРП классификации наблюдений из области $T(\varepsilon, \mathbf{z}) \subset \mathbb{X}$ ($\varepsilon > 0$), второй и третий характеризуют приращение условного ε -риска, обусловленное соответственно вариацией и смещением оценок плотностей $\{\hat{p}_{\alpha}^{(l)}(\mathbf{y})\}$ ($l = 1, 2$).

С л е д с т в и е 1. В условиях теоремы 2 имеет место соотношение:

$$r_n^{(2)}(\varepsilon, \mathbf{z}) = r_n^{(0)}(\varepsilon, \mathbf{z}) + o_1. \quad (64)$$

Д о к а з а т е л ь с т в о. Полагая в разложении для $r_n^{(1)}(\varepsilon, \mathbf{z})$ $H_1 = H_0$, получаем (64).

Сравним по точности классификации решающие правила, использующие оценки плотности с фиксированным и переменным ядром в предположении, что точность аппроксимации риска возрастает, а зависимость компонент вектора признаков $\mathbf{y} \in \mathfrak{R}^p$ усиливается (имеет место условие (4)). Обозначим $\bar{\sigma} = \max_{j=1,\dots,N} \{\bar{\sigma}_{jj}\}$ и рассмотрим асимптотику:

$$\varepsilon \rightarrow \infty, \quad \bar{\sigma} \rightarrow \infty \quad (\text{tr}(\Sigma) \rightarrow 0), \quad |Q| \rightarrow \infty, \quad |H_0| \rightarrow 0, \quad n \rightarrow \infty. \quad (65)$$

При этом будем предполагать, что $\Delta(\mathbf{z}) \rightarrow \infty$ и $u_{\varepsilon,\alpha}^2 \rightarrow \infty$ с такой скоростью, что для $\tilde{u}_{\varepsilon,m}^2(\mathbf{z})$, определяемого (32), (36), выполняются ограничения: $0 < \tilde{u}_{\varepsilon,m}^2(\mathbf{z}) < \infty$ ($\mathbf{z} \in \mathbb{Z}$). В противном случае: при $u_{\varepsilon,\alpha}^2 \rightarrow 0$ уменьшается «степень пересечения» как классов, так и областей $\{T^{(\alpha)}(\varepsilon, \mathbf{z})\}$, и точ-

ность оценивания риска с помощью функционалов $\{r_n^{(l)}(\varepsilon, \mathbf{z})\}$ убывает. С другой стороны, при $u_{\varepsilon, \alpha}^2 \rightarrow \infty$ интегралы $\{\alpha_\alpha(\mathbf{z})\}$ в разложении (52) становятся неограниченными.

С л е д с т в и е 2. Пусть выполняются условия теоремы 2 и область $\mathbb{Z} = Z_1 \times \dots \times Z_M \subset \mathfrak{R}^M$ – ограниченный параллелепипед, причем компонента $z_j \in Z_j$ вектора $\mathbf{z} = (z_j) \in \mathbb{Z}$ равномерно распределена на интервале Z_j ($j=1, \dots, M$). Тогда функционалы $r_n^{(1)}(\varepsilon, \mathbf{z})$ и $r_n^{(2)}(\varepsilon, \mathbf{z})$ в асимптотике (65) при условии, что $0 < \tilde{u}_{\varepsilon, m}^2(\mathbf{z}) < \infty (\mathbf{z} \in \mathbb{Z})$ связаны соотношением:

$$r_n^{(1)}(\varepsilon, \mathbf{z}) = r_n^{(2)}(\varepsilon, \mathbf{z}) + \rho_{\varepsilon, \sigma}(\mathbf{z}) + o_2, \quad \rho_{\varepsilon, \sigma}(\mathbf{z}) = O(\bar{\sigma}^4) \geq 0, \quad o_2 = o(\bar{\sigma}^4 h^4). \quad (66)$$

Д о к а з а т е л ь с т в о. Пусть V_{N-1} – объем $(N-1)$ -мерного параллелепипеда. С учетом (30)–(33) имеем:

$$\text{mes}(\tilde{T}^{(\alpha)}(\varepsilon, \mathbf{z})) = V_{N-1} \sqrt{|\Sigma|} \tilde{u}_{\varepsilon, \alpha}^{N-1}(\mathbf{z}) L(\mathbf{z}, \Sigma), \quad L(\mathbf{z}, \Sigma) = \bar{\sigma}_{NN} / \sqrt{(1 + \mathbf{a}'(\mathbf{z}) \Sigma^{-1} \mathbf{a}(\mathbf{z}))}. \quad (67)$$

На основании выражений (56)–(58) и (67) справедливо соотношение:

$$r_n^{(1)}(\varepsilon, \mathbf{z}) = r_n^{(2)}(\varepsilon, \mathbf{z}) + \tilde{\rho}_{\varepsilon, \sigma}(\mathbf{z}) n^{-1} h^{-(N+M)} + \rho_{\varepsilon, \sigma}(\mathbf{z}) h^4 + o_1,$$

где $(a = C^2 - 4(\ln(c_1/c_2))^2 \geq 0)$:

$$\tilde{\rho}_{\varepsilon, \sigma}(\mathbf{z}) = \frac{\text{mes}\{\mathbb{Z}\} (c_1 + c_2) V_{N-1} \tilde{u}_{\varepsilon, \alpha}^{N-1}(\mathbf{z}) L(\mathbf{z}, \Sigma) \sqrt{|\mathcal{Q}||\Sigma|} - \sqrt{|H_1|}}{2\Delta_1(\mathbf{z}) \sqrt{|\mathcal{Q}||H_1|}},$$

$$\rho_{\varepsilon, \sigma}(\mathbf{z}) = \frac{c_m n_1 (0 | l_m(\mathbf{z}), \Delta^2(\mathbf{z}))}{4} (a \tau_{m, \varepsilon}^{(0)}(\mathbf{z}, A) + 2C \tau_{m, \varepsilon}^{(1)}(\mathbf{z}, A) + \tau_{m, \varepsilon}^{(2)}(\mathbf{z}, A)).$$

В асимптотике (65) при $0 < \tilde{u}_{\varepsilon, m}^2(\mathbf{z}) < \infty (\mathbf{z} \in \mathbb{Z})$ имеем: согласно (34), (67) $L(\mathbf{z}, \Sigma) = O(1)$, $\text{mes}\{\mathbb{Z}\} = O(|\mathcal{Q}|)$; $|H_0| = |\mathcal{Q}| \times |\Sigma| \rightarrow 0$ и $1/\Delta_1(\mathbf{z}) = o(1/\Delta^2(\mathbf{z}))$. Вследствие чего, $\tilde{\rho}_{\varepsilon, \sigma}(\mathbf{z}) \rightarrow 0$.

При тех же условиях с учетом (53)–(55) получаем:

$$\tau_{m, \varepsilon}^{(0)}(\mathbf{z}, A) \rightarrow P\{\eta_m \eta_m' \leq u_\varepsilon^2\} \leq 1, \quad \tau_{m, \varepsilon}^{(1)}(\mathbf{z}, A) \rightarrow E\{\eta_m A \eta_m'\} = \text{tr}\{K_m(\mathbf{z}) A\},$$

$$\tau_{m, \varepsilon}^{(2)}(\mathbf{z}, A) \rightarrow E\{(\eta_m A \eta_m')^2\} = \text{tr}^2\{K_m(\mathbf{z}) A\} + 2 \text{tr}\{(K_m(\mathbf{z}) A)^2\}.$$

где: $\eta_\alpha \in \mathfrak{R}^N (\alpha \in S)$ – случайный гауссовский вектор с математическим ожиданием $\mu_\eta(\alpha, \mathbf{z})$ и ковариационной матрицей $\Sigma_\eta(\alpha, \mathbf{z})$; $A = \Sigma^{-\frac{1}{2}} \tilde{\Psi}^{(l)} \Sigma^{-\frac{1}{2}}$, $K_m(\mathbf{z}) = \Sigma_\eta(\alpha, \mathbf{z}) + \mu_\eta(\alpha, \mathbf{z}) \mu_\eta'(\alpha, \mathbf{z}) \in \mathfrak{S}_N$.

На основании свойств следа матрицы: $\text{tr}(K_m(\mathbf{z}) A) + \text{tr}^2(K_m(\mathbf{z}) A) + \text{tr}(K_m(\mathbf{z}) A)^2 = O(\text{tr}^4(\Sigma^{-1}))$,

что влечет $\rho_{\varepsilon, \sigma}(\mathbf{z}) = O(\bar{\sigma}^4)$. Следствие 2 доказано.

Как известно, качество непараметрических ядерных оценок плотности, а, следовательно, и точность непараметрических классификаторов критическим образом зависит от выбора коэффи-

циентов сглаживания $h = h(n)$, удовлетворяющих (9). Применительно к используемым в статье ядерным оценкам плотности данная проблема рассматривается в [6]. В асимптотике усиливающейся зависимости компонент вектора $\mathbf{y} \in \mathcal{R}^p$ (4) исследуем влияние выбора матрицы ядра и коэффициентов сглаживания на точность рассматриваемых подстановочных решающих правил.

Обозначим $\Delta r_n(\varepsilon, \mathbf{z}) = r_n^{(1)}(\varepsilon, \mathbf{z}) - r_n^{(2)}(\varepsilon, \mathbf{z})$ ($\mathbf{z} \in \mathbb{Z}$), $\beta = \bar{\sigma}h$ и рассмотрим последовательности однотипных ситуаций, возникающих, если при $\bar{\sigma} \rightarrow \infty$ и $h \rightarrow 0$: а) $\beta \rightarrow 0$, б) $\beta \rightarrow \lambda$ ($0 < \lambda < \infty$), в) $\beta \rightarrow \infty$. С учетом следствия 2 справедливо следующее утверждение.

С л е д с т в и е 3. В условиях теоремы 2 для описанных выше ситуаций имеет место:

$$\text{а) } \Delta r_n(\varepsilon, z) \rightarrow 0, \text{ б) } \Delta r_n(\varepsilon, z) \rightarrow \lambda > 0, \text{ в) } \Delta r_n(\varepsilon, z) \rightarrow \infty.$$

Таким образом, подстановочное решающее правило вида (17), использующее оценки плотности $\{\hat{p}_\alpha^{(1)}(\mathbf{y})\}$ с фиксированным ядром в ситуациях б), в) (возникающих в случае «малого» объема выборки) проигрывает по точности решающему правилу, использующему оценки $\{\hat{p}_\alpha^{(2)}(\mathbf{y})\}$ с адаптивным ядром, и этот проигрыш обусловлен увеличением в указанных ситуациях смещения оценок плотности с фиксированным ядром, определяемым выражением (25).

Литература

1. Айвазян, С.А. Прикладная статистика. Исследование зависимостей / С.А. Айвазян, И.С. Енюков, Л.Д. Мешалкин. М., 1985.
2. Харин, Ю.С. Эконометрическое моделирование / Ю.С. Харин, В.И. Малюгин, А.Ю. Харин. Минск, 2003.
3. Прикладная статистика. Классификация и снижение размерности / С.А. Айвазян [и др.]. М., 1989.
4. Харин, Ю.С. Робастность в статистическом распознавании образов. Мн., 1992.
5. Фукунага, К. Введение в статистическую теорию распознавания образов. М., 1979.
6. Малюгин, В.И. Об оценивании плотности случайных векторов с существенно зависимыми компонентами // Вестник БГУ. Сер. 1. – 1985. – № 2. – С. 41–44.
7. Андерсон, Т. Введение в многомерный статистический анализ данных. М., 1963.
8. Малюгин, В.И. Учет зависимости признаков в задачах непараметрической классификации // Динамика систем. Управление, оптимизация, адаптация. Горький, 1983. – С. 143-160.
9. Гельфанд, И.М. Обобщенные функции и действия над ними / И.М. Гельфанд, Г.Е. Шилков. М., 2007. – 408 с.

Аннотация

Рассматривается задача классификации многомерных наблюдений в пространстве существенно зависимых признаков с помощью подстановочных решающих правил, основанных на непараметрических оценках плотности с фиксированным и адаптивным гауссовским ядром. Аналитическое исследование решающих правил проводится для случая, когда существенная зависимость признаков описывается моделью многомерной линейной регрессии с равномерно распределенными экзогенными переменными и гауссовскими ошибками наблюдения. С помощью асимптотических разложений условного риска проводится сравнительный анализ альтернативных решающих правил и показывается преимущество предлагаемого непараметрического классификатора с адаптивным ядром в условиях усиливающейся зависимости признаков и растущего объема обучающей выборки.

V. I. MALUGIN

ASYMPTOTIC ANALYSIS OF THE RISK OF NONPARAMETRIC CLASSIFICATION IN THE CASE OF ESSENTIAL DEPENDENT FEATURES

The problem of the classification of multivariate observations in the case of essential dependent features is considered. The decision “plug-in” rules based on the nonparametric kernel density estimators with fixed and adaptive Gaussian kernel are suggested. Analytical investigations of decision rules are conducted for the case when the essential dependent of the features described by the multivariate linear regression models with univariate distributed exogenous variables and Gaussian disturbances. Alternative decision rules are compared by means of asymptotic expansions of the conditional risk. The preference of the adaptive decision rule in the asymptotic of strengthening dependents of features and growing learning sample size is established.