

И. А. Суслина, О. В. Соколов

**СВЯЗЬ ЗАДАЧИ СЕЛЕКЦИИ РАЗРЕЖЕННОЙ
ПОДМАТРИЦЫ МАТРИЦЫ БОЛЬШОГО РАЗМЕРА
И БАЙЕСОВСКОЙ ЗАДАЧИ ПРОВЕРКИ ГИПОТЕЗ**

§1. ВВЕДЕНИЕ

Пусть наблюдения представлены в виде матрицы $\mathbf{Y} = \{Y_{ij}\}$ размера $N \times M$

$$Y_{ij} = s_{ij} + \xi_{ij}, \quad i = 1, \dots, N, \quad j = 1, \dots, M, \quad (1.1)$$

где $\{\xi_{ij}\}$ – независимые одинаково распределенные случайные величины и $s_{ij} \geq 0$, для всех $i \in \{1, \dots, N\}$, $j \in \{1, \dots, M\}$. Ошибки наблюдения ξ_{ij} распределены по стандартному гауссовскому закону.

Матрица размера $N \times M$ содержит $L = C_N^n \cdot C_M^m$ различных подматриц $\mathcal{C} = \mathcal{C}_k$, $k = 1, \dots, L$ размера $n \times m$. Пусть

$$\begin{aligned} \mathcal{C}_k = \{C = A \times B \subset \{1, \dots, N\} \times \{1, \dots, M\}, \\ \text{Card}(A) = n, \quad \text{Card}(B) = m\} \end{aligned}$$

представляет из себя набор индексов (i, j) элементов подматрицы \mathcal{C}_k , $k = 1, \dots, L$.

Рассмотрим $(L + 1)$ простую гипотезу:

$$\begin{aligned} H_0 : \quad Y_{ij} = \xi_{ij} &\quad \text{при любых } (i, j), \quad i \in \{1, \dots, N\}, \quad j \in \{1, \dots, M\}, \\ H_k : \quad Y_{ij} = \begin{cases} \xi_{ij} & \text{при } (i, j) \notin \mathcal{C}_k, \\ \xi_{ij} + s_{ij}, \quad s_{ij} \geq a > 0 & \text{при } (i, j) \in \mathcal{C}_k, \end{cases} & \quad (1.2) \\ k = 1, \dots, L. \end{aligned}$$

То есть, гипотеза H_0 состоит в том, что все элементы матрицы \mathbf{Y} центрированы, а гипотеза H_k – в том, что центрированы все элементы, кроме элементов подматрицы \mathcal{C}_k размера $n \times m$, а именно: $s_{ij} \geq a > 0$ для пар (i, j) из множества индексов \mathcal{C}_k и $s_{ij} = 0$ для пар $(i, j) \notin \mathcal{C}_k$.

Таким образом, перед нами стоит задача проверки гипотез, в частности, нахождения такого теста, с помощью которого можно будет по

Ключевые слова: байесовская задача проверки гипотез, байесовский риск, минимаксный риск, селекция, оптимальный байесовский тест, состоятельный тест.

исходным наблюдениям, представленным в виде матрицы \mathbf{Y} , определить, какая гипотеза будет принята. Качество теста характеризуется байесовским и минимаксным рисками.

Измеримые функции, определенные на множестве наблюдений $\{\mathbf{Y}\}$ со значениями в множестве $\{0, 1, \dots, L\}$, назовем тестами:

$$\psi = \psi(\mathbf{Y}) : \{\mathbf{Y}\} \rightarrow \{0, \dots, L\}.$$

Пусть \mathbf{P}_{H_k} – мера, соответствующая гипотезе H_k . Для каждого теста ψ , вводя обозначения $R_k(\psi) = \mathbf{P}_{H_k}(\psi \neq k)$, определим его риск. Положим

$$R(\psi) = \max_{k=0, \dots, L} \mathbf{P}_{H_k}(\psi(\mathbf{Y}) \neq k) = \max_{k=0, \dots, L} \{R_k(\psi)\}.$$

Минимаксный риск определяется соотношением

$$R_M = \inf_{\psi} R(\psi),$$

где инфимум берется по всем измеримым отображениям $\{\mathbf{Y}\}$ в $\{0, 1, \dots, L\}$.

Нас будут интересовать асимптотические свойства минимаксного риска в случае разреженных подматриц $\mathcal{C}_k = \mathcal{C}_k(N, M, n, m, a)$. Предположим, что

$$\begin{aligned} N &\rightarrow \infty, & M &\rightarrow \infty, & n &\rightarrow \infty, & m &\rightarrow \infty, \\ p = n/N &\rightarrow 0, & q = m/M &\rightarrow 0. \end{aligned} \tag{1.3}$$

Назовем тест состоятельный в минимаксном смысле, если $R(\psi) \rightarrow 0$ при выполнении условий (1.3).

Будем искать граничные значения для a : нижние границы – условия, при которых минимаксный риск стремится к единице, а также верхние границы – условия, при которых минимаксный риск стремится к нулю. При изучении верхних границ будет построен состоятельный в минимаксном смысле тест.

Аналогичные задачи в несколько иной постановке изучались, в частности в [1–5]. Мы существенно используем результаты работы [1]. В ней достаточно подробно представлена история изучения подобных задач. Результаты, полученные для гауссовского случая в работе [2], верны только по порядку. Поэтому результаты, приведенные в ней для случая, когда ошибки наблюдения ξ_{ij} распределены не по стандартному гауссовскому закону, не представляются возможным использовать

в наших рассмотрениях. В работе [3] показано, что проверка гипотезы

$$H_0 : \quad Y_{ij} = \xi_{ij} \quad \text{при любых } (i, j), \quad i \in \{1, \dots, N\}, \quad j \in \{1, \dots, M\},$$

против альтернативы: существует такая подматрица $\mathcal{C} = \mathcal{C}_k$, $k = 1, \dots, L$, размера $n \times m$ исходной матрицы $\mathbf{Y} = \{Y_{ij}\}$ размера $N \times M$, что

$$Y_{ij} = \begin{cases} \xi_{ij} & \text{при } (i, j) \notin \mathcal{C}_k, \\ \xi_{ij} + s_{ij}, \quad s_{ij} \geq a > 0 & \text{при } (i, j) \in \mathcal{C}_k, \end{cases}$$

возможна в более широких условиях, чем изучаемая нами задача (см. [1], глава 2.2).

§2. ОПТИМАЛЬНЫЙ БАЙЕСОВСКИЙ ТЕСТ

Рассмотрим байесовскую задачу проверки k простых гипотез. Пусть $\{\mathcal{X}, \mathcal{A}, \mathcal{P}\}$ – статистический эксперимент и пусть $\mathbf{P}_1, \dots, \mathbf{P}_k$ – вероятностные меры, $\mathbf{P}_l \in \mathcal{P}$ соответствует l -ой гипотезе о наблюдении $x \in \mathcal{X}$. Пусть $\pi = \{\pi_1, \dots, \pi_k\}$ – априорные вероятности соответствующих гипотез, где $\pi_l \geq 0$, $l = 1, \dots, k$, $\sum_{l=1}^k \pi_l = 1$. Для теста $\psi = \psi(x)$ – измеримого отображения $\mathcal{X} \rightarrow \{1, \dots, k\}$ рассмотрим вероятности ошибок $R_l(\psi) = \mathbf{P}_l(\psi \neq l) = 1 - \mathbf{P}_l(\psi = l)$. Их среднее значение (смесь по мере π) назовем байесовским риском теста:

$$R_B(\psi) = \sum_{l=1}^k \pi_l R_l(\psi) = 1 - \sum_{l=1}^k \pi_l \mathbf{P}_l(\psi = l).$$

Будем называть тест $\psi^*(x)$ *оптимальным*, если $R_B(\psi^*) = \inf_{\psi} R_B(\psi) = R_B$, где инфимум берется по всем измеримым отображениям \mathcal{X} в $\{1, \dots, k\}$. Хорошо известно, что для любой априорной меры $\pi = \{\pi_0, \dots, \pi_L\}$, $\sum_{k=0}^L \pi_k = 1$, байесовский риск не превосходит минимаксного риска. То есть,

$$R_B = \inf_{\psi} \sum_{k=0}^L \pi_k R_k(\psi) \leq R_M = \inf_{\psi} R(\psi) = \inf_{\psi} \max_{k=0, \dots, L} R_k(\psi). \quad (2.1)$$

Будем предполагать, что при всех $l = 1, \dots, k$ существуют $f_l(x)$ – плотности мер \mathbf{P}_l по отношению к общей мере \mathbf{P} .

Предложение 2.1. *Оптимальный тест в байесовской проверке k простых гипотез имеет вид:*

$$\psi^*(x) = \operatorname{argmax}_{1 \leq l \leq k} \pi_l f_l(x).$$

Таким образом, тест $\psi^*(x) = l_0$, если $\pi_{l_0} f_{l_0}(x) \geq \pi_l f_l(x)$ при $l \in \{1, \dots, L\}$, $l \neq l_0$. В случае равенства $\pi_{l_0} f_{l_0}(x) = \pi_l f_l(x)$ при $l \in \{1, \dots, L\}$, $l \neq l_0$, применяется дополнительная процедура розыгрыша.

Доказательство. Нужно найти минимум следующего выражения:

$$R_B(\psi) = \sum_{i=1}^{i=k} \mathbf{P}_i(\psi \neq i) \pi_i,$$

где риск минимизируется по множеству всех тестов $\psi : \mathcal{X} \rightarrow \{1, \dots, k\}$.

Имеем

$$R_B(\psi) = \sum_{i=1}^{i=k} \mathbf{P}_i(\psi \neq i) \pi_i = 1 - \sum_{i=1}^{i=k} \mathbf{P}_i(\psi = i) \pi_i.$$

Таким образом, задача сводится к нахождению максимума выражения:

$$W = \sum_{i=1}^{i=k} \mathbf{P}_i(\psi = i) \pi_i. \quad (2.2)$$

В силу того, что $f_i(x)$ является плотностью вероятностной меры \mathbf{P}_i , выражение (2.2) примет вид:

$$W = \sum_{i=1}^{i=k} \int_{\mathcal{X}_i} \pi_i f_i(x) dx, \quad (2.3)$$

где $\mathcal{X}_i = \{x \in \mathcal{X} : \psi(x) = i\}$, $i = 1, \dots, k$. Легко видеть, что сумма (2.3) достигает максимума, если

$$\mathcal{X}_i = \{x \in \mathcal{X} : \pi_i f_i(x) \geq \pi_j f_j(x), \forall j = 1, \dots, k\}.$$

□

Замечание 2.1. Если мы хотим связать байесовскую задачу проверки гипотез (1.2) с селекцией, то естественно предполагать, что априорные вероятности π_k , $k = 1, \dots, L$, совпадают.

Предложение 2.2. *Обозначим $M = \frac{mna}{2} + \frac{\log L}{a} + \frac{\log \frac{p}{1-p}}{a}$. Рассмотрим априорную меру $\pi = (\pi_0, \dots, \pi_L)$, $\pi_0 = p$, $\pi_k = \frac{1-p}{L}$, $k = 1, \dots, L$,*

где $0 < p < 1$ – вероятность того, что $s_{ij} = 0$ для любых i, j . Тогда тест

$$\tilde{\psi}^{\text{оп}}(\mathbf{Y}) = \begin{cases} 0, & \max_k \left\{ \sum_{C_k} Y_{ij} \right\} < M, \\ l > 0, & \max_k \left\{ \sum_{C_k} Y_{ij} \right\} \geq M, \quad \operatorname{argmax}_k \left\{ \sum_{C_k} Y_{ij} \right\} = C_l, \end{cases} \quad (2.4)$$

является оптимальным байесовским тестом для проверки гипотез (1.2), если $s_{ij} = 0$ или $s_{ij} = a$.

Доказательство. Из предложения 2.1 следует, что оптимальный байесовский тест в данной задаче должен иметь вид:

$$\psi(\mathbf{Y}) = \operatorname{argmax}_{0 \leq k \leq L} \pi_k f_k(Y),$$

где $f_0(Y)$ соответствует матрице \mathbf{Y} с центрированными элементами $Y_{ij} \sim N(0, 1)$, а каждая функция $f_k(Y)$, $k = 1, \dots, L$, соответствует матрице \mathbf{Y} , в которой $s_{ij} = 0$ при $(i, j) \notin C_k$ и $s_{ij} = a > 0$ при $(i, j) \in C_k$. Аргумент максимума не изменится, если разделить все выражения на $f_0(Y)$ – положительную константу. Следовательно, оптимальный тест можно искать в виде:

$$\psi(\mathbf{Y}) = \operatorname{argmax}_{0 \leq k \leq L} \pi_k \frac{f_k(Y)}{f_0(Y)},$$

где $\frac{f_k(Y)}{f_0(Y)} = \frac{d\mathbf{P}_k}{d\mathbf{P}_0}$ – отношение правдоподобия для гипотезы H_k . Здесь $f_0(Y)$ соответствует гипотезе $H_0 : s_{ij} = 0$ при всех допустимых (i, j) , а $f_k(Y)$ соответствует гипотезе $H_k : s_{ij} = a > 0$, $(i, j) \in C_k$, $s_{ij} = 0$, $(i, j) \notin C_k$.

Для H_0 отношение правдоподобия равно 1, найдем $\frac{f_k(Y)}{f_0(Y)}$ при $k = 1, \dots, L$. Все $\{\xi_{ij}\}$ независимы и распределены по стандартному нормальному закону. Следовательно,

$$\frac{f_k(Y)}{f_0(Y)} = \frac{\prod_{(i,j) \notin C_k} \exp(-Y_{ij}^2/2) \times \prod_{(i,j) \in C_k} \exp(-(Y_{ij} - a)^2/2)}{\prod_{\substack{1 \leq i \leq N, \\ 1 \leq j \leq M}} \exp(-Y_{ij}^2/2)}.$$

Приходим к выражению:

$$\frac{f_k(Y)}{f_0(Y)} = \exp\left(-\frac{a^2 nm}{2} + a \sum_{(i,j) \in C_k} Y_{ij}\right).$$

Откуда

$$\begin{aligned} \pi_0 \frac{f_0(Y)}{f_0(Y)} &= p, \\ \pi_k \frac{f_k(Y)}{f_0(Y)} &= \frac{1-p}{L} \exp\left(-\frac{a^2 nm}{2} + a \sum_{(i,j) \in C_k} Y_{ij}\right), \quad k = 1, \dots, L. \end{aligned}$$

В результате использования оптимального байесовского теста гипотеза H_0 принимается при условии:

$$\max_k \sum_{C_k} Y_{ij} < \frac{mna}{2} + \frac{\log(L)}{a} + \frac{\log \frac{p}{1-p}}{a}.$$

Если это условие не выполнено, в выражении $\pi_k \frac{f_k(Y)}{f_0(Y)}$ нужно максимизировать сумму $\sum_{(i,j) \in C_k} Y_{ij}$. Следовательно, тест должен принимать гипотезу H_l , если $\operatorname{argmax}_k \sum_{C_k} Y_{ij} = C_l$. Таким образом, тест (2.4) является оптимальным байесовским тестом. \square

Заметим, что при выполнении условий (1.3)

$$\frac{mna}{2} + \frac{\log(L)}{a} + \frac{\log \frac{p}{1-p}}{a} \sim \frac{mna}{2} + \frac{\log(L)}{a} = M_1.$$

Далее мы будем рассматривать и анализировать тест

$$\psi^{\text{оп}}(\mathbf{Y}) = \begin{cases} 0, & \max_k \left\{ \sum_{C_k} Y_{ij} \right\} < M_1, \\ l > 0, & \max_k \left\{ \sum_{C_k} Y_{ij} \right\} \geq M_1, \quad \operatorname{argmax}_k \left\{ \sum_{C_k} Y_{ij} \right\} = C_l. \end{cases} \quad (2.5)$$

Замечание 2.2. Заметим, что из соображений симметрии $R_i(\psi^{\text{оп}}) = R_j(\psi^{\text{оп}})$ при $1 \leq i, j \leq L$.

§3. ВЕРХНИЕ ГРАНИЦЫ

Пусть выполнено (1.3). Укажем условия, при которых минимаксный риск $R_M \rightarrow 0$, а также покажем, что при этих условиях тест (2.5) является состоятельным, то есть, его минимаксный риск стремится к нулю. Заметим, что, если $R_M \rightarrow 0$, то байесовский риск $R_B = \inf_{\psi} R_B(\psi) \rightarrow 0$ для любого набора априорных вероятностей $\pi = \{\pi_1, \dots, \pi_k\}$, $p_i \geq 0$, $\sum_{i=1}^k p_i = 1$.

Определим основные величины, необходимые в дальнейших рассмотрениях:

$$\begin{aligned} B &= B_{n, m, N, M} = \min\{A_1, A_2, A\} \\ \text{где } A &= \frac{a\sqrt{nm}}{\sqrt{2(n \log(p^{-1}) + m \log(q^{-1}))}}, \\ A_1 &= \frac{a\sqrt{m}}{\sqrt{2}(\sqrt{\log(n)} + \sqrt{\log(N-n)})}, \\ A_2 &= \frac{a\sqrt{n}}{\sqrt{2}(\sqrt{\log(m)} + \sqrt{\log(M-m)})}. \end{aligned} \tag{3.1}$$

Нам потребуется следующее замечание.

Замечание 3.1. Функция распределения суммы $Y = \sum_{i=1}^n X_i$ – независимых случайных величин, распределенных по стандартному нормальному закону, имеет вид $F_Y(t) = \mathbf{P}(Y < t) = \Phi(t/\sqrt{n})$, где $\Phi(\cdot)$ – функция распределения стандартного нормального закона. Кроме того, при выполнении условий (1.3) справедливо соотношение $\log(C_N^n) \sim n \log(N/n)$.

Предложение 3.1. Пусть выполнено (1.3) и $B = B_{n, m, N, M}$, определенное (3.1), удовлетворяет условию

$$\liminf B_{n, m, N, M} > 1. \tag{3.2}$$

Тогда тест, определенный (2.5), является состоятельным, то есть $R(\psi^{op}) = \max_{k=0, \dots, L} \mathbf{P}_{H_k}(\psi(\mathbf{Y}) \neq k) \rightarrow 0$.

Доказательство. Так как

$$R(\psi^{op}) = \max_{k=0, \dots, L} R_k(\psi^{op}) = \max(R_0(\psi^{op}), R_1(\psi^{op})),$$

нужно доказать, что $R_0(\psi^{\text{оп}})$ и $R_1(\psi^{\text{оп}})$ одновременно стремятся к нулю.

I. Оценка риска $R_1(\psi^{\text{оп}})$. Покажем, что при выполнении условий предложени $R_1(\psi^{\text{оп}}) \rightarrow 0$. Очевидно, что

$$\begin{aligned} R_1(\psi^{\text{оп}}) &= \mathbf{P}_{H_1}(\psi^{\text{оп}} \neq 1) = \mathbf{P}_{H_1}\left(\max_k \sum_{C_k} Y_{ij} > \sum_{C_1} Y_{ij}\right) \\ &\quad + \mathbf{P}_{H_1}\left(\max_k \sum_{C_k} Y_{ij} = \sum_{C_1} Y_{ij}, \sum_{C_1} Y_{ij} < \frac{mna}{2} + \frac{\log L}{a}\right) \quad (3.3) \\ &\leq \mathbf{P}_1 + \mathbf{P}_2, \end{aligned}$$

где $\mathbf{P}_1 = \mathbf{P}_{H_1}\left(\max_k \sum_{C_k} Y_{ij} > \sum_{C_1} Y_{ij}\right)$, $\mathbf{P}_2 = \mathbf{P}_{H_1}\left(\sum_{C_1} Y_{ij} < \frac{mna}{2} + \frac{\log L}{a}\right)$.

В доказательстве теоремы 2.1 в [1] показывается, что $\mathbf{P}_1 \rightarrow 0$ при выполнении условий (1.3) и (3.2). Значит это верно и в условиях предложения.

Рассмотрим теперь \mathbf{P}_2 – второе слагаемое в последнем выражении (3.3). Для него справедливо равенство

$$\begin{aligned} \mathbf{P}_2 &= \mathbf{P}_{H_1}\left(\sum_{C_1} Y_{ij} < \frac{mna}{2} + \frac{\log L}{a}\right) = \mathbf{P}_{H_1}\left(\sum_{C_1} (\xi_{ij} + a) < \frac{mna}{2} + \frac{\log L}{a}\right) \\ &= \mathbf{P}_{H_1}\left(\sum_{C_1} \xi_{ij} + mna < \frac{mna}{2} + \frac{\log L}{a}\right) \\ &= \mathbf{P}_{H_1}\left(\sum_{C_1} \xi_{ij} < -\frac{mna}{2} + \frac{\log L}{a}\right). \end{aligned}$$

Используя замечание 3.1, имеем

$$\begin{aligned} \mathbf{P}_{H_1}\left(\sum_{C_1} \xi_{ij} < -\frac{mna}{2} + \frac{\log L}{a}\right) &= \Phi\left(-\frac{a\sqrt{mn}}{2} + \frac{\log L}{a\sqrt{mn}}\right) \\ &= \Phi\left(-\frac{a\sqrt{mn}}{2}\left(1 - \frac{2\log L}{a^2 mn}\right)\right). \end{aligned}$$

Кроме того, снова используя замечание 3.1, условие (3.2), а также равенство $C_N^n \cdot C_M^m = L$, получаем

$$\frac{a^2 mn}{2 \log L} \sim 1 + \delta, \quad \delta > 0, \quad (3.4)$$

где величина $\delta = \delta(N, M, n, m, a)$ равномерно отделена от нуля при выполнении (1.3). Следовательно,

$$\Phi\left(-\frac{a\sqrt{mn}}{2}\left(1 - \frac{2\log L}{a^2mn}\right)\right) \sim \Phi\left(-\frac{a\sqrt{mn}}{2}\frac{\delta}{(1+\delta)}\right).$$

Из условия (3.2) следует, что $a\sqrt{mn} \rightarrow \infty$. Отсюда получаем, что $\mathbf{P}_2 \rightarrow 0$, условия (1.3), (3.2) влечут $R_1(\psi^{\text{оп}}) \rightarrow 0$.

II. Оценка риска R_0 . Используя хорошо известную асимптотику

$$\Phi(x) \sim \frac{1}{\sqrt{2\pi}} \frac{e^{-x^2/2}}{(-x)} \quad \text{при } x \rightarrow -\infty, \quad (3.5)$$

покажем, что $R_0(\psi^{\text{оп}}) \rightarrow 0$. Имеем

$$\begin{aligned} R_0(\psi^{\text{оп}}) &= \mathbf{P}_{H_0}(\psi^{\text{оп}} \neq 0) = \mathbf{P}_{H_0}\left(\max_{k=1, \dots, L} \sum_{(i,j) \in C_k} Y_{ij} \geqslant \frac{mna}{2} + \frac{\log L}{a}\right) \\ &\leqslant \sum_{k=1}^L \mathbf{P}_{H_0}\left(\sum_{(i,j) \in C_k} Y_{ij} \geqslant \frac{mna}{2} + \frac{\log L}{a}\right) \\ &= L \cdot \Phi\left(-\left(\frac{a\sqrt{mn}}{2} + \frac{\log L}{a\sqrt{mn}}\right)\right). \end{aligned}$$

Из (3.5) следует, что при выполнении (1.3) для достаточно больших n, m справедливо неравенство:

$$L \cdot \Phi\left(-\left(\frac{a\sqrt{mn}}{2} + \frac{\log L}{a\sqrt{mn}}\right)\right) \leqslant \exp\left(-\left(\frac{a\sqrt{mn}}{2} + \frac{\log L}{a\sqrt{mn}}\right)^2 / 2 + \log L\right). \quad (3.6)$$

Рассмотрим правую часть неравенства (3.6). Из соотношения (3.4) получаем

$$\begin{aligned} \frac{mna^2}{8} - \frac{\log L}{2} + \frac{\log^2 L}{2mna^2} &\sim \frac{1+\delta}{4} \log L - \frac{\log L}{2} + \frac{\log L}{4(1+\delta)} \\ &= \frac{\log L}{4} \frac{\delta^2}{1+\delta} \rightarrow \infty. \end{aligned}$$

Следовательно, в условиях предложения 3.1, имеем $R_0(\psi^{\text{оп}}) \rightarrow 0$. \square

Следствие 3.1. При выполнении условий (1.3), (3.2) минимаксный риск $R_M = \inf_{\psi} R(\psi) \leqslant R(\psi^{\text{оп}}) \rightarrow 0$. Обозначим через $R_B(\pi)$ байесовский риск, возникающий при априорном распределении π . Имеем

$R_B(\pi) \leq R_M$, следовательно, $R_B(\pi) \rightarrow 0$ при любом априорном распределении π .

Замечание 3.2. В параграфе 3 все рассмотрения и доказательства проводились в предположении, что $s_{ij} = 0$ или $s_{ij} = a$. Легко видеть (сравни с [1], замечание 2.1), что все рассмотрения и доказательства переносятся на случай $s_{ij} = 0$ или $s_{ij} \geq a$.

§4. НИЖНИЕ ГРАНИЦЫ

Укажем условия, при которых минимаксный риск R_M стремится к 1.

Предложение 4.1. Пусть выполнено (1.3) и $B = B_{n,m,N,M}$, определенное (3.1), удовлетворяет условию

$$\limsup B_{n,m,N,M} < 1. \quad (4.1)$$

Тогда

$$\mathbf{P}_{H_1} \left(\max_k \sum_{C_k} Y_{ij} > \sum_{C_1} Y_{ij} \right) \rightarrow 1.$$

Доказательство. Рассмотрим задачу чистой селекции, которая аналогочна данной, но исключает возможность реализации гипотезы H_0 . То есть, в матрице \mathbf{Y} обязательно есть подматрица C_k , такая что $s_{ij} = a > 0$ для $(i,j) \in C_k$ и $s_{ij} = 0$ для $(i,j) \notin C_k$. Теорема 2.2 в [1] утверждает, что при выполнении условий (1.3), (4.1)

$$\inf_{\psi} \left\{ \max_{k=1,\dots,L} \mathbf{P}_{H_k} (\psi(\mathbf{Y}) \neq k) \right\} \rightarrow 1,$$

где инфимум берется по всем измеримым отображениям $\{\mathbf{Y}\}$ в множество $\{1, \dots, L\}$.

Рассмотрим тест

$$\psi^+(\mathbf{Y}) = \operatorname{argmax}_{k=1,\dots,L} \left\{ \sum_{(i,j) \in C_k} Y_{ij} \right\}.$$

В силу симметрии $\mathbf{P}_{H_i} (\psi^+(\mathbf{Y}) \neq i) = \mathbf{P}_{H_j} (\psi^+(\mathbf{Y}) \neq j)$, $1 \leq i, j \leq L$, следовательно,

$$\begin{aligned} \inf_{\psi} \left\{ \max_{k=1,\dots,L} \mathbf{P}_{H_k} (\psi(\mathbf{Y}) \neq k) \right\} &\leq \max_{k=1,\dots,L} \mathbf{P}_{H_k} (\psi^+(\mathbf{Y}) \neq k) \\ &= \mathbf{P}_{H_1} (\psi^+(\mathbf{Y}) \neq 1) = \mathbf{P}_{H_1} \left(\max_k \sum_{C_k} Y_{ij} > \sum_{C_1} Y_{ij} \right). \quad \square \end{aligned}$$

Предложение 4.2. *Обозначим $R = mna/2$. Тест*

$$\psi^u(\mathbf{Y}) = \begin{cases} 0, & \max_{k \in \{1, \dots, L\}} \left\{ \sum_{C_k} Y_{ij} \right\} \leq R, \\ l, & \operatorname{argmax}_k \left\{ \sum_{C_k} Y_{ij} \right\} = C_l, \quad \max_{k \in \{1, \dots, L\}} \left\{ \sum_{C_k} Y_{ij} \right\} > R, \end{cases}$$

является оптимальным байесовским тестом при априорной мере $\pi = (\pi_0, \dots, \pi_L)$, $\pi_k = 1/(L+1)$, $k = 0, \dots, L$.

Доказательство. Из предложения 2.1 следует, что оптимальный байесовский тест в данной задаче должен иметь вид:

$$\psi(\mathbf{Y}) = \operatorname{argmax}_{0 \leq k \leq L} \pi_k f_k(\mathbf{Y}),$$

где $f_0(Y)$ соответствует матрице \mathbf{Y} с центрированными элементами $Y_{ij} \sim N(0, 1)$, а $f_k(Y)$, $k = 1, \dots, L$, соответствует матрице \mathbf{Y} , в которой $\mathbf{E}(Y_{ij}) = 0$ при $(i, j) \notin C_k$ и $\mathbf{E}(Y_{ij}) = a > 0$ при $(i, j) \in C_k$.

Следовательно, оптимальный тест можно искать в виде:

$$\psi(\mathbf{Y}) = \operatorname{argmax}_{0 \leq k \leq L} \frac{f_k(Y)}{f_0(Y)},$$

где $\frac{f_k(Y)}{f_0(Y)} = \frac{d\mathbf{P}_k}{d\mathbf{P}_0}$ – отношение правдоподобия для гипотезы H_k . Для

H_0 отношение правдоподобия равно единице, отношения $\frac{f_k(Y)}{f_0(Y)}$ для $k = 1, \dots, L$ были найдены при доказательстве предложения 2.2.

Имеем

$$\frac{f_0(Y)}{f_0(Y)} = 1, \tag{4.2}$$

$$\frac{f_k(Y)}{f_0(Y)} = \exp\left(-\frac{a^2 nm}{2} + a \sum_{(i,j) \in C_k} Y_{ij}\right), \quad k = 1, \dots, L. \tag{4.3}$$

Легко видеть, что $\frac{f_0(Y)}{f_0(Y)} = 1$ является максимумом при условии

$$\sum_{C_k} Y_{ij} < R, \quad k \in \{1, \dots, L\},$$

то есть, в этом случае оптимальный тест должен принимать значение 0. При

$$\max_{k \in \{1, \dots, L\}} \left\{ \sum_{(i,j) \in C_k} Y_{ij} \right\} > R$$

максимальное отношение правдоподобия соответствует максимальной сумме $\sum_{C_k} Y_{ij}$. При $\max_{k \in \{1, \dots, L\}} \left\{ \sum_{C_k} Y_{ij} \right\} = R$ максимальное отношение правдоподобия достигается по крайней мере при двух гипотезах. Будем считать, что в этом случае принимается H_0 . \square

Предложение 4.3. *Пусть выполнены условия (1.3) и (4.1). Тогда минимаксный риск $R_M \rightarrow 1$.*

Доказательство. Рассмотрим априорную меру $\pi = (\pi_0, \dots, \pi_L)$, $\pi_k = 1/(L+1)$, $k = 0, \dots, L$. Из предложения 4.2 следует, что оптимальным байесовским тестом при данной априорной мере является тест ψ^u и $R_B(\psi^u) = R_B$. Согласно (2.1) $R_B(\psi^u) = R_B \leq R_M$. Следовательно, достаточно доказать, что в условиях предложения $R_B(\psi^u) \rightarrow 1$. Определим байесовский риск теста ψ^u :

$$R_B(\psi^u) = \sum_{k=0}^L \frac{1}{L+1} R_k(\psi^u) \geq \frac{L}{L+1} R_1(\psi^u). \quad (4.4)$$

Из предложения 4.1 следует

$$R_1(\psi^u) = \mathbf{P}_{H_1}(\psi^u \neq 1) \geq \mathbf{P}_{H_1} \left(\max_k \sum_{C_k} Y_{ij} > \sum_{C_1} Y_{ij} \right) \rightarrow 1. \quad (4.5)$$

\square

ЛИТЕРАТУРА

1. C. Butucea, Y. I. Ingster, I. A. Suslina, *Sharp variable selection of a sparse submatrix in a high-dimensional noisy matrix*. — ESAIM: Probab. Statist. **19** (2015), 115–134.
2. M. Kolar, S. Balakrishnan, A. Rinaldo, A. Singh, *Minimax localization of structural information in large noisy matrices*. — In: Advances in Neural Information Processing Systems 24, J. Shawe-Taylor et al. eds., pp. 909–917, NIPS, 2011.
3. C. Butucea, Y. I. Ingster, *Detection of a sparse submatrix of a high-dimensional noisy matrix*. — Bernoulli **19**, No. 5B (2013), 2652–2688.
4. C. Butucea, G. Gayraud, *Sharp detection of smooth signals in a high-dimensional sparse matrix with indirect observations*. — Ann. Inst. H. Poincaré: Probab. Statist. **52**, No. 4 (2016), 1564–1591.

5. X. Sun, A. B. Nobel, *On the maximal size of large-average and ANOVA-fit submatrices in a Gaussian random matrix*. — Bernoulli **19**, No. 1 (2013), 275–294.

Suslina I. A., Sokolov O. V. The connection between the selection problem for a sparse submatrix of a large-size matrix and the Bayes problem of hypothesis testing.

We associate the selection problem for a sparse submatrix of a matrix of large dimension and the problem of testing the hypothesis of the existence of a sparse submatrix possessing the required properties with the Bayesian hypothesis testing problem.

С.-Петербургский Национальный
исследовательский Университет
информационных технологий,
механики и оптики,
Кронверкский проспект, дом 49,
197101, Санкт-Петербург, Россия
E-mail: isuslina@mail.ru

Поступило 8 ноября 2017 г.

МакКинзи и Компания
улица Лесная, дом 5, строение "С",
125047, Москва, Россия
E-mail: olesokolov@gmail.com