

Г. И. ПЕРЕТЯГИН

(Новосибирск)

ОТБОР ВЫДЕЛЯЮЩИХСЯ НАБЛЮДЕНИЙ И КРИТЕРИИ СДВИГА

При автоматизации экспериментальных исследований на стадии предварительной обработки поступающей информации часто возникает задача обнаружения элементов, выделяющихся из однородного (в некотором смысле) множества. В частности, к таковым относятся, например, задача отбора «информационных» фрагментов при обработке изображений, обнаружение «сигнальных» точек в статистически однородных случайных полях, отбраковка аномальных измерений и т. п.

1. В своей простейшей формулировке эти проблемы можно свести к следующему виду. В некоторой последовательности наблюдений, состоящей из n измеренных величин (x_1, x_2, \dots, x_n), имеется одно или несколько значений, не удовлетворяющих модели исследуемого процесса. Требуется обнаружить эти величины. Основная трудность здесь состоит в том, что измерения сопровождаются неустранимым шумом ($x_i = x_{0i} + \xi_i$) и проверка наблюдаемых величин должна основываться на анализе статистических характеристик процесса. Нужно отметить, что если проблема выделения одного выпадающего наблюдения допускает общее решение с точки зрения теории проверки гипотез [1], то обнаружение нескольких выделяющихся элементов (что более существенно для практики) до сих пор представляет собой «камень преткновения». Это связано с тем, что применяемые обычно последовательные процедуры отбора выпадающих наблюдений (начиная с экстремального) с теоретической точки зрения не вполне удовлетворительны — здесь принципиально нельзя учесть влияние всех выбросов на распределение критических статистик. Неудивительно, что эти статистики иногда искажаются выпадающими наблюдениями настолько, что выходят из своей критической зоны (подробнее об эффекте «маскирования» см. в работах [2, 3]). Очевидный путь повышения устойчивости критериев отбора выпадающих наблюдений состоит в том, чтобы основывать их на «центральной» части выборки. Цель данной работы и состоит в реализации такого подхода. Прежде всего находятся оптимальные решающие правила для фиксированного числа выделяющихся наблюдений. Принципы инвариантности и достаточности в применении к данным решающим правилам позволяют непосредственно найти оптимальные критерии сдвига в задачах с «мешающими» параметрами. Наконец, при неизвестном общем числе выпадающих наблюдений приводится оптимальная процедура отбора этих наблюдений.

2. Пусть X_1, X_2, \dots, X_n — независимые одинаково распределенные случайные величины, имеющие общую функцию распределения $F_\theta(X)$, где $\theta = \theta_0 \in \Theta$. Объем выборки n фиксирован. Альтернатива состоит в том, что некоторые из случайных величин $X_{v_1}, X_{v_2}, \dots, X_{v_k}$ имеют ту же самую общую функцию распределения $F_\theta(X)$, но $\theta_{v_j} \neq \theta_0$, $j = \overline{1, k}$. Будем считать, что заранее нельзя отдать предпочтение ни одному из C_n^k возможных наборов $(X_{v_1}, \dots, X_{v_k})$, удовлетворяющих альтернативному требованию задачи. В связи с этим естественно ограничиться лишь симметричными решающими правилами. Для этого сведем задачу к проверке многих гипотез H_0, H_1, \dots, H_M , $M = C_n^k$, и потребуем, чтобы гипотезы H_i , $i = \overline{1, M}$, были симметричными.

Обозначим через $f(X, \theta)$ плотность распределения случайной величины $X = (X_1, X_2, \dots, X_n)$, когда параметр равен θ . Кроме того, пусть $(\lambda_1^v, \lambda_2^v)$ обозначают всевозможные разбиения ($v=1, M$) множества

$$\bigcup_{i=1}^n \omega_i = \{x \in \mathbb{R}^n : x_i \in \omega_i\}$$

При альтернативных гипотезах $H_i^{(\lambda_1^i, \lambda_2^i)}$ ($i=1, M$) плотность совместного распределения компонент вектора X имеет вид

$$\begin{aligned} h_i(X) &= h(X|S_i) = \prod_{j=k+1}^n f(X_{ij}, \Theta_0) \prod_{l=1}^k f(X_{il}, \Theta_{il}) = \\ &= \prod_{j=1}^n f(X_j, \Theta_0) \prod_{l=1}^k \frac{f(X_{il}, \Theta_{il})}{f(X_{il}, \Theta_0)}. \end{aligned} \quad (1)$$

По данному выборочному значению $x = (x_1, x_2, \dots, x_n)$ требуется отдать предпочтение одной из гипотез H_0, H_1, \dots, H_M . Здесь нерандомизированное решающее правило может быть представлено измеримым разбиением выборочного пространства V_n на $M+1$ взаимно непересекающихся подмножеств $(\omega_0, \omega_1, \dots, \omega_M)$, где $\bigcup_{i=0}^M \omega_i = V_n$ и $\omega_i \cap \omega_j = \emptyset$, так что выбирается i -я альтернативная гипотеза H_i , когда $x \in \omega_i$. Поступая обычным способом, введем критическую функцию $\varphi(x) = (\varphi_0(x), \varphi_1(x), \dots, \varphi_M(x))$, где $\varphi_i(x)$ равна 1 в области ω_i и 0 вне этой области (при желании $\varphi(x)$ можно продолжить до рандомизированной критической функции [4]). Кроме того, пусть γ_v обозначает перестановку компонент n -мерного вектора X , соответствующую разбиению $(\lambda_1^v, \lambda_2^v)$ множества $\{1, 2, \dots, n\}$, и Γ — группа всех таких перестановок. Тогда требование симметрии будет означать, что решающая функция $\varphi(x)$ симметрична по Γ , т. е. замена $\gamma_v x$ на $\gamma_u x$ приводит к такому переупорядочению функций $\varphi_i(\gamma_v x)$, $i=1, M$, что в целом $\varphi(x)$ остается неизменной: $\varphi(\gamma_v x) = \varphi(\gamma_u x)$, $\gamma_v, \gamma_u \in \Gamma$. Нетрудно заметить, что в этом случае вероятность принятия гипотезы H_i при условии, что она верна, $P(H_i|H_i)$ не зависит от i для $i=1, M$. Принимая данные положения во внимание, будем рассматривать решающие правила $\varphi(x) = (\varphi_0(x), \varphi_1(x), \dots, \varphi_M(x))$, максимизирующие вероятности правильной классификации — $P(H_i|H_i)$, $i=1, M$ при условии, что $P(H_0|H_0) = 1 - \alpha$, $0 < \alpha < 1$.

3. Сформулируем теперь байесовскую задачу решения о выборе одной из гипотез H_0, H_1, \dots, H_M . При этом будем предполагать, что θ_0 известно. Прежде всего введем функцию $L_j(H_i)$, определяющую потери, связанные с принятием гипотезы H_j , когда верна гипотеза H_i . Примем ее равной

$$L_j(H_i) = 1 - \delta_{ij} = \begin{cases} 1, & i \neq j; \\ 0, & i = j. \end{cases} \quad (2)$$

Тогда для данного правила множественного решения $\varphi(x)$ математическое ожидание потерь (функция риска) есть вектор $(R(H_0, \varphi), R(H_1, \varphi), \dots, R(H_M, \varphi))$,

где

$$R(H_i, \varphi) = \int_{V_n} \sum_{j=0}^M L_j(H_i) \varphi_j(x) h(x|s_i) dx = 1 - E_{s_i} \varphi_i(x). \quad (3)$$

Используя байесовский принцип, будем считать, что p_i — априорные вероятности гипотез H_i , а истинное распределение имеет плотность $h(x|s_i)$. В этом случае байесовский риск относительно $p = (p_0, p_1, \dots, p_M)$ есть величина

$$r(p, \varphi) = \sum_{i=0}^M p_i R(H_i, \varphi) = 1 - \sum_{i=0}^M p_i E_{s_i} \varphi_i(x). \quad (4)$$

Проблема выбора $\varphi(x)$, минимизирующей данный риск, эквивалентна проблеме выбора $\varphi(x)$ для максимизации вероятности правильной классификации [4]. Известное байесовское решение относительно данного p задается следующим правилом: принять гипотезу H_j , если

$$p_j h_j(x) = \max_i p_i h_i(x), \quad i = 0, 1, \dots, M. \quad (5)$$

Другими словами, каждое правило $\varphi(x)$, для которого $\{\varphi_j(x) = 1, \varphi_i(x) = 0, j \neq i\}$, если $p_j h_j(x)$ больше любого другого $p_i h_i(x)$, имеет минимальный байесовский риск. (Здесь нужно отметить, что байесовское правило разрешает определенную рандомизацию среди тех j , для которых $p_j h_j(x) = \max_i p_i h_i(x)$, $i = 1, M$; следует лишь удовлетворить требованию $\sum_{i=1}^M \varphi_i(x) = 1$.) Принимая во внимание симметрию задачи, будем брать априорные вероятности симметричными, исходя из того, что вероятность принятия гипотезы H_i не должна зависеть от перестановок наблюдений:

$$p_i = p(H_i) = p(H_j) = p_j (j = 1, M).$$

В этом случае симметричные априорные распределения должны давать равные веса всем гипотезам H_i , $i = 1, M$. Поэтому приходим к представлению априорных распределений в виде $\hat{p}_0 = p(H_0) = 1 - Mp^0$, $\hat{p}_i = p(H_i) = p^0$ для всех $i = 1, M$ и $0 \leq p^0 \leq 1/M$. Используя правило (5), найдем байесовскую стратегию для распределения $\hat{p}(H_i)$.

Пусть $(\lambda_1^v, \lambda_2^v)$ — произвольное разбиение множества $\{1, 2, \dots, n\}$, а $(\lambda_1^k, \lambda_2^k)$ — такое разбиение, что как только $x_{v_i} = \gamma_k x_i$, то

$$\max_v \prod_{j=1}^k \frac{f_j(x_{v_j})}{f_0(x_{v_j})} = \prod_{j=1}^k \frac{f_j(x_{k_j})}{f_0(x_{k_j})} = W.$$

Тогда каждое решающее правило $\hat{\varphi}(x)$, для которого

$$\hat{\varphi}_0(x) = 1, \quad \text{если } W < \frac{1 - Mp^0}{p^0};$$

$$\hat{\varphi}_v(x) = 1, \quad \text{если } \prod_{i=1}^k \frac{f_i(x_{v_i})}{f_0(x_{v_i})} = W > \frac{1 - Mp^0}{p^0} \quad (v = 1, M), \quad (6)$$

является байесовским относительно $\hat{p}(H_i)$ и максимизирует вероятность выбора правильного решения.

Если число выпадающих наблюдений (k) неизвестно, то можно использовать минимаксное решающее правило, соответствующее наи-

менее благоприятному априорному распределению $\hat{p}^*(H_i)$, для которого

$$\min_{\varphi} r(\hat{p}^*, \varphi) = \max_{\hat{p}} \min_{\varphi} r(\hat{p}, \varphi).$$

В этом случае

$$r(\hat{p}, \varphi) = 1 - \sum_{i=0}^M \hat{p}_i E_{s_i} \varphi_i(x) = 1 - (1 - Mp^0) E_{s_0} \varphi_0(x) - p^0 \sum_{i=1}^M E_{s_i} \varphi_i(x),$$

где $0 < p^0 < 1/M$, $M = C_n^k$. Выбрав для определенности $p^0 = \varepsilon/M$, где $0 < \varepsilon < 1$, заметим, что $r(\hat{p}, \varphi)$ возрастает по переменной p^0 с ростом k , достигая своего максимального значения при $k = n/2$. Отсюда следует, что наименее благоприятное априорное распределение \hat{p}^* соответствует либо $k = n/2$, либо некоторому $k_{\max} < n/2$. Отметим, что решающее правило (6) будет симметричным байесовским правилом и относительно наименее благоприятного априорного распределения.

Для перехода к классической постановке теории проверки гипотез нужно положить $\frac{1 - Mp^0}{p^0} = C_\alpha$; в этом случае риск $R(H_0, \varphi)$ называется уровнем значимости критерия $\alpha = 1 - E_{s_0} \varphi_0(x) = R(H_0, \varphi)$. Если определить ρ как число разбиений $(\lambda_1^v, \lambda_2^v)$, для которых $\prod_{j=1}^k \frac{f_j(x_{v_j})}{f_0(x_{v_j})} = W$, и найти C_α такое, что $E_{s_0} \varphi_0(x) = 1 - \alpha$, то из соотношения (6) непосредственно следует, что верна гипотеза H_0 при $W < C_\alpha$ и гипотезы H_v при $\prod_{i=1}^k \frac{f_i(x_{v_i})}{f_0(x_{v_i})} = W > C_\alpha$:

$$\begin{aligned} \hat{\varphi}_0(x) &= 1, \text{ если } W < C_\alpha; \\ \hat{\varphi}_i(x) &= \frac{1}{\rho}, \text{ если } \prod_{i=1}^k \frac{f_i(x_{j_i})}{f_0(x_{j_i})} = \max_v \prod_{i=1}^k \frac{f_i(x_{v_i})}{f_0(x_{v_i})} = W > C_\alpha. \end{aligned} \quad (6')$$

(При $W = C_\alpha$ возможна рандомизация: $\varphi_0(x) = \eta$ и $\varphi_j(x) = (1 - \eta)/\rho$.) Данное правило является симметричным допустимым, в силу того что оно симметрично по Γ и байесовское относительно $\hat{p}(H_i)$, где $p^0 = 1/(C_\alpha + M)$. Более того, пусть $\tilde{\varphi}(x) = (\tilde{\varphi}_0(x), \tilde{\varphi}_1(x), \dots, \tilde{\varphi}_M(x))$ есть некоторое другое симметричное решающее правило, имеющее $E_{s_0} \tilde{\varphi}_0(x) \geq E_{s_0} \hat{\varphi}_0(x)$. Тогда, поскольку правило (6') байесовское относительно $\hat{p}(H_i)$, $r(\hat{p}, \tilde{\varphi}) \geq r(\hat{p}, \hat{\varphi})$ и $r(\hat{p}, \tilde{\varphi}) - r(\hat{p}, \hat{\varphi}) = \sum_{i=0}^M \hat{p}(H_i) (E_{s_i} \hat{\varphi}_i(x) - E_{s_i} \tilde{\varphi}_i(x)) \geq 0$. По предположению $E_{s_0} \tilde{\varphi}_0(x) \geq E_{s_0} \hat{\varphi}_0(x)$ и $\hat{p}(H_i) > 0$ получаем $\sum_{i=1}^M E_{s_i} \hat{\varphi}_i(x) \geq \sum_{i=1}^M E_{s_i} \tilde{\varphi}_i(x)$. Учитывая, что $\sum_{i=1}^M E_{s_i} \varphi_i(x)$ характеризует мощность критерия $\varphi(x)$, резюмируем вышеизложенное следующим образом.

В классе всех симметричных правил, удовлетворяющих условиям $P(H_0|H_0) = 1 - \alpha$ и $P(H_i|H_i) = P(H_j|H_j)$, $j = 1, M$, каждый критерий $\varphi(x)$ вида (6') для некоторого C_α ($0 \leq \eta(x) \leq 1$) является наиболее мощным объема α для проверки гипотезы H_0 против H_i , $i = 1, M$. (Отметим, что данное предложение — обобщение результата Фергюсона [1], относящегося к проверке одного выпадающего наблюдения.)

Если альтернативные гипотезы сложные, то для семейств распределений $h_i(x) = h(x|s_i)$ с монотонным отношением правдоподобия

данное предложение приводит к равномерно наиболее мощным (РНМ) симметричным критериям. Именно если существуют функции $T_v(x)$ такие, что отношение $h(x|s_v)/h(x|s_0)$ является неубывающей функцией от $T_v(x)$ и $\frac{h(x|s_v)}{h(x|s_\mu)} \geq 1$, когда соответственно $T_v(x)/T_\mu(x) \geq 1$, то правило (6') сводится к следующему:

$$\begin{aligned}\hat{\varphi}_0(x) &= 1 \text{ при } \max_v T_v(x) < C_\alpha; \\ \hat{\varphi}_j(x) &= 1 \text{ при } T_j(x) = \max_v T_v(x) > C_\alpha.\end{aligned}\tag{7}$$

Данное правило является РНМ симметричным.

4. Вообще говоря, распределение результатов наблюдений, получаемых в эксперименте, зависит дополнительно от некоторых «мешающих» параметров. В этом случае в качестве семейств с монотонным отношением правдоподобия обычно берут распределения, инвариантные относительно подходящим образом выбранной группы преобразований [5]. Для иллюстрации рассмотрим следующую задачу. Пусть (x_1, \dots, x_n) — независимые случайные переменные, относительно распределения которых можно выдвинуть $(C_n^k + 1)$ гипотез:

$$\begin{aligned}H_0: x_j &\sim \mathcal{N}(\mu, \sigma^2), j = \overline{1, n}; \\ H_v: x_{v_i} &\sim \mathcal{N}(\mu + \eta_i, \sigma^2), i = \overline{1, k}, \text{ и } x_{v_j} \sim \mathcal{N}(\mu, \sigma^2), j = \overline{k+1, n},\end{aligned}$$

где μ, σ^2 — неизвестные параметры, а $(\eta_1, \eta_2, \dots, \eta_k)$ — неизвестные положительные числа. Нетрудно заметить, что проблема инвариантна к изменениям параметров сдвига (μ) и масштаба (σ) в том смысле, что умножение всех x_i на постоянную λ и добавление к ним одинаковой константы оставляет каждую из гипотез неизменной. Естественно потребовать, чтобы решающее правило $\varphi(x)$ для выбора одной из $(M+1)$ гипотез было инвариантным относительно изменений параметров сдвига и масштаба:

$$\varphi_i(x_1, x_2, \dots, x_n) = \varphi_i(\lambda x_1 + u, \lambda x_2 + u, \dots, \lambda x_n + u)$$

для всех $i = 0, 1, \dots, M$, u и λ . Для таких правил вероятности $P(H_i | H_i)$ не зависят от μ и σ^2 .

Выделение класса решающих правил, инвариантных относительно изменений положения и масштаба наблюдений, означает выделение в исходном выборочном пространстве некоторого класса \mathcal{A}_1 инвариантных событий A таких, что для любых действительных u и λ имеет место эквивалентность [5]

$$\{(x_1, x_2, \dots, x_n) \in A\} \Leftrightarrow \{(\lambda x_1 - u, \lambda x_2 - u, \dots, \lambda x_n - u) \in A\}.$$

Если определить

$$\widehat{f}(x_1, x_2, \dots, x_n) = \int_0^\infty \int_{-\infty}^\infty f(\lambda x_1 - u, \lambda x_2 - u, \dots, \lambda x_n - u) \lambda^{n-2} du d\lambda$$

и

$$\widehat{f}_0(x_1, x_2, \dots, x_n) = \int_0^\infty \int_{-\infty}^\infty f_0(\lambda x_1 - u, \lambda x_2 - u, \dots, \lambda x_n - u) \lambda^{n-2} du d\lambda,$$

то

$$\frac{\widehat{f}(x_1, x_2, \dots, x_n)}{\widehat{f}_0(x_1, x_2, \dots, x_n)} = \frac{\widehat{f}(\lambda x_1 - u, \lambda x_2 - u, \dots, \lambda x_n - u)}{\widehat{f}_0(\lambda x_1 - u, \lambda x_2 - u, \dots, \lambda x_n - u)}$$

и для любого $A \subset \mathcal{A}_1$

$$P(A) = \int_A \frac{\hat{f}(x)}{\hat{f}_0(x)} f_0(x) dx.$$

В нашем примере имеем сложную нулевую гипотезу, состоящую из плотностей вида

$$h_0(x | \mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 + \sum_{i=k+1}^n (x_{v_i} - \mu^2) \right\},$$

Так как распределение инвариантной статистики критерия не зависит от параметров сдвига и масштаба, то для упрощения дальнейших вычислений положим $\mu = 0$ и $\sigma^2 = 1$. Тогда

$$\begin{aligned} \hat{h}_0(x | 0, 1) &= \frac{1}{2} n^{-\frac{1}{2}} \pi^{-\frac{n}{2} + \frac{1}{2}} \Gamma\left(\frac{n}{2} - \frac{1}{2}\right) \left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^{-\frac{n}{2} + \frac{1}{2}}, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \\ \hat{h}_v(x | \eta, 1) &= n^{-\frac{1}{2}} (2\pi)^{-\frac{n}{2} + \frac{1}{2}} \left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^{-\frac{n}{2} + \frac{1}{2}} \exp \left\{ -\frac{1}{2} \left(\sum_{i=1}^k \eta_i^2 - \frac{1}{2} \left(\sum_{i=1}^k \eta_i \right)^2 \right) \right\} \int_0^\infty e^{-\frac{1}{2}(\lambda_1^2 - 2\lambda_1 t)} \lambda_1^{n-2} d\lambda_1, \\ t &= \frac{\sum_{i=1}^k \eta_i (x_{v_i} - \bar{x})}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2}}. \end{aligned}$$

Следовательно,

$$\frac{\hat{h}_v(x | \eta, 1)}{\hat{h}_0(x | 0, 1)} = C(\eta) \int_0^\infty \exp \left\{ -\frac{1}{2} \lambda^2 + \lambda \frac{\sum_{i=1}^k \eta_i (x_{v_i} - \bar{x})}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2}} \right\} \lambda^{n-2} d\lambda.$$

Если $\eta_1 = \eta_2 = \dots = \eta_k = \eta$, то полученное отношение будет монотонно возрастающей функцией от $T_v(x) = \frac{\sum_{i=1}^k (x_{v_i} - \bar{x})}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2}}$. В соответствии с (7)

равномерно наиболее мощным симметричным объемом α решающим правилом, инвариантным относительно изменений параметров сдвига и масштаба, будет следующее:

$$\hat{\Phi}_0(x) = \begin{cases} 1, & \text{если } \max_v T_v(x) < C_\alpha; \\ 0, & \text{если } \max_v T_v(x) \geq C_\alpha; \end{cases}$$

$$\hat{\varphi}_j(x) = \begin{cases} 1, & \text{если } T_j(x) = \max_v T_v(x) \geq C_\alpha; \\ 0 & \text{в остальных случаях.} \end{cases}$$

Константа C_α выбирается так, чтобы $P(H_0|H_0) = 1 - \alpha$. Данное правило имеет максимальную величину $P(H_i|H_i)$ для $i \neq 0$ в классе всех симметричных правил независимо от $\eta > 0$. Можно заметить, что

$$\max_v T_v(x) = \max_v \frac{\sum_{i=1}^k (x_{v_i} - \bar{x})}{S} = \frac{\sum_{i=n-k+1}^n (x_{(i)} - \bar{x})}{S},$$

где $x_{(1)} < x_{(2)} < \dots < x_{(n)}$ — порядковые статистики выборки. Кроме того, не существует РНМ симметричного инвариантного критерия, если $\eta_1 \neq \eta_2 \neq \dots \neq \eta_k$.

5. Для поиска соответствующих тестов (при произвольных η_i) рассмотрим задачу с точки зрения общей линейной гипотезы. Именно будем считать, что n наблюдаемых величин (x_1, \dots, x_n) являются линейными комбинациями m неизвестных постоянных $\theta_1, \theta_2, \dots, \theta_m$ и «ошибок» $\xi_1, \xi_2, \dots, \xi_n$:

$$x_i = p_{i1}\theta_1 + p_{i2}\theta_2 + \dots + p_{im}\theta_m + \xi_i, \quad i = \overline{1, n},$$

где $\{p_{ij}\} = P$ — матрица известных постоянных коэффициентов. Целью анализа является получение выводов относительно $\{\xi_i\}$: соответствуют ли они основным предположениям ($\xi_i \sim N(0, \sigma^2)$, σ^2 неизвестно) либо необходимо «исключение» соответствующих x_{jl} (при $\xi_{jl} \sim N(\eta_l, \sigma^2)$, $l = \overline{1, k}$).

Для представления проблемы в канонической форме введем выборочное пространство V_n наблюдений $x = (x_1, \dots, x_n)$ и будем считать, что вектор $E(x) = P\theta$ при гипотезе H_0 лежит в подпространстве V_m , порождаемом вектором-столбцом матрицы $P(E(x) = \eta_{(m)})$; при альтернативных гипотезах $H_v, v = \overline{1, C_n^k}$, вектор $E(x)$ может принадлежать некоторому из C_n^k возможных подпространств V_{m+k}^v размерности $(m+k)$. Обозначим через \hat{x} проекцию x на V_m и через \hat{x}_1 — проекцию x на $V_{n-m} = V_m^\perp$ (ортогональное дополнение к V_m). Тогда $x = \hat{x} + \hat{x}_1$, поскольку каждый вектор $v \in V_n$ допускает единственное представление вида $v = v_1 + v_2$, где $v_1 \in V_m$, $v_2 \in V_m^\perp$. Вектор \hat{x} является наилучшим приближением вектора x в подпространстве V_m в том смысле, что на нем достигается минимум функции $\|x - \hat{x}\|^2$ [4]. (Здесь V_n разлагается в прямую сумму подпространств V_m и V_{n-m} : $V_n = V_m \oplus V_{n-m}$).

Рассмотрим структуру подпространства V_{n-m} . Согласно предположениям о сдвиге некоторых наблюдений $x_{v_1}, x_{v_2}, \dots, x_{v_k}$ относительно выбранной модели, будем идентифицировать с каждой гипотезой H_v векторное подпространство $V_{kv} \subset V_{n-m}$ и полагать, что вектор $E(\hat{x}^\perp) = \eta_{(m)}^\perp$ лежит в этом подпространстве. Обозначим \hat{x}_{kv}^\perp — проекцию вектора x^\perp на V_{kv} — через V_{kv}^\perp — ортогональное дополнение V_{kv} до V_{n-m} . По критерию факторизации статистика $U(\hat{x}, \|x^\perp\|^2)$ достаточна для $N(\eta_{(m)}, \sigma^2 I_n)$, так как при H_0

$$\frac{1}{2\sigma^2} \|x - \eta_{(m)}\|^2 = \frac{1}{2\sigma^2} (\|\hat{x} - \eta_{(m)}\|^2 + \|\hat{x}^\perp\|^2).$$

В то же время при H_v достаточной статистикой для $\eta_{(m)}^\perp$ будет $T_v(x) =$

$= \hat{x}^\perp_{kv}$. Нетрудно заметить, что функции

$$\psi(U, T_v) = \psi_v(x) = \frac{\|\hat{x}^\perp_{kv}\|^2}{\|\hat{x}^\perp_{kv}\|^2}$$

монотонно возрастают по T_v при фиксированной U и независимы от U при H_0 . Симметричным критерием проверки линейной гипотезы H_0 при $H_v: E(x^\perp) \in V_{kv}$ будет критерий с критической областью S_k , определяемой неравенством

$$\max_v \psi_v(x) = \max_v \frac{\|\hat{x}^\perp_{kv}\|^2}{\|\hat{x}^\perp_{kv}\|^2} \geq F_\alpha(k, n - m - k),$$

где $F_\alpha(k, n - m - k) = 100\alpha$ -процентная точка распределения переменной $\max_v \psi_v(x)$. Здесь критическая область S_k является подобной размера α , поскольку функции $\psi_v(x)$ независимы от U и $P_{\theta, \sigma^2}(S_k | U = U_0) = \alpha$ для всех $(\theta_1, \theta_2, \dots, \theta_m)$ и σ^2 (вероятность отклонения гипотезы H_0 равна α на каждой из поверхностей $U = U_0$) [5]. Таким образом, РНМ симметричным критерием в классе подобных является критерий $\varphi(x)$ вида

$$\begin{aligned} \hat{\varphi}_0(x) &= 1, \text{ если } \max_v \psi_v(x) < F_\alpha; \\ \hat{\varphi}_j(x) &= 1, \text{ если } \psi_j(x) = \max_v \psi_v(x) \geq F_\alpha. \end{aligned} \tag{8}$$

6. Конструкция прямой суммы подпространств, примененная выше, является теоретической основой дисперсионного анализа, позволяющей эффективно осуществить ортогональное разложение гипотез [6]. Для более полной редукции проблемы отбора выпадающих наблюдений в рассматриваемой постановке введем последовательность «группированных» гипотез: $H_1^* = H^{i_1}$ [в выборке (x_1, \dots, x_n) имеется один выпадающий элемент — x_{i_1}], $H_2^* = H^{i_1} \cap H^{i_2}$ (два элемента — x_{i_1} и x_{i_2}), \dots ; $H_k^* = H^{i_1} \cap H^{i_2} \cap \dots \cap H^{i_k}$. В соответствии с гипотезами $H^{i_l}, l = 1, k$, введем k подпространств $W_{i_1}, W_{i_2}, \dots, W_{i_k}$, попарно ортогональных и порождающих все $V_{kv}: V_{kv} = \bigoplus_{j=1}^k W_{v_j}$. Очевидно, что $V_n = V_m \oplus V_{n-m}$ и операторами проектирования на V_m и V_{n-m} являются соответственно $P(P^t P)^{-1} P^t$ и $(I - P(P^t P)^{-1} P^t)$. Пусть $B_0 = V_{n-m}^v$ и $B_r = \left(\bigoplus_{i=1}^r W_{v_i} \right)^\perp$, так что $B_k = V_{kv}^\perp$ и $B_{k+1} = \emptyset$. Подпространства B_i образуют убывающую последовательность подпространств $V_{n-m}^v = B_0 \supset B_1 \supset \dots \supset B_k \supset B_{k+1} = \emptyset$. Обозначая через E_i проектор на подпространство B_i , получаем «убывающую» последовательность проекционных операторов

$$I - P(P^t P)^{-1} P^t = E_0 \supset E_1 \supset \dots \supset E_k \supset E_{k+1} = 0.$$

Поскольку

$$V_{n-m}^v = \bigoplus_{i=0}^k (B_i \ominus B_{i+1}) = \left(\bigoplus_{i=1}^k W_{v_i} \right) \oplus V_{kv}^\perp, \tag{9}$$

а оператор ортогонального проектирования на подпространство $(B_i \ominus \ominus B_{i+1})$ равен разности $(E_i - E_{i+1})$, то разложение (9) можно записать

в терминах проекционных операторов:

$$I - P(P^T P)^{-1} P^T = \sum_{i=0}^k (E_i - E_{i+1}).$$

Это дает разложение вектора X на сумму ($k+2$) взаимно ортогональных векторов:

$$\begin{aligned} x &= P(P^T P)^{-1} P^T x + (I - P(P^T P)^{-1} P^T)x = P(P^T P)^{-1} P^T x + \\ &+ \sum_{i=0}^k (E_i - E_{i+1})x = \hat{x}_{V_m} + \hat{x}_{V_{kv}^\perp} + \sum_{j=1}^k \hat{x}_{W_{v_j}}, \end{aligned}$$

откуда

$$\|x\|^2 = \|\hat{x}_{V_m}\|^2 + \|\hat{x}_{V_{kv}^\perp}\|^2 + \sum_{j=1}^k \|\hat{x}_{W_{v_j}}\|^2.$$

Нетрудно показать, что проекционные операторы равны

$$E_r = \Lambda_r^v P (P^T \Lambda_r^v P)^{-1} P^T \Lambda_r^v,$$

где

$$\Lambda_1^v = \begin{pmatrix} v_1 & & \\ 1_1 & \ddots & 0 \\ & \ddots & \vdots \\ & 0 & 1 \\ & & \ddots & 1 \end{pmatrix}_{v_1}; \quad \Lambda_2^v = \begin{pmatrix} v_1 & & v_2 & & \\ 1_1 & \ddots & 0 & & \\ & \ddots & \ddots & \ddots & \\ & 0 & \ddots & 1_0 & \\ & & & & \ddots & 1 \end{pmatrix}_{v_2} \text{ и т. д.}$$

Используя это обстоятельство, получаем разложение квадратичной формы

$$\begin{aligned} \|x\|^2 &= x^T x = x^T (P(P^T P)^{-1} P^T) x + x^T (I - P(P^T P)^{-1} P^T - \Lambda_1^v + \\ &+ \Lambda_1^v P (P^T \Lambda_1^v P)^{-1} P^T \Lambda_1^v) x + x^T (\Lambda_1^v - \Lambda_1^v P (P^T \Lambda_1^v P)^{-1} P^T \Lambda_1^v - \Lambda_2^v + \\ &+ \Lambda_2^v P (P^T \Lambda_2^v P)^{-1} P^T \Lambda_2^v) x + \dots + x^T (\Lambda_{k+1}^v - \Lambda_{k+1}^v P (P^T \Lambda_{k+1}^v P)^{-1} \times \\ &\times P^T \Lambda_{k+1}^v) x = \|\hat{x}_{V_m}\|^2 + \|\hat{x}_{V_{kv}^\perp}\|^2 + Z_{v_1}^2 + Z_{v_2}^2 + \dots + Z_{v_k}^2 = \\ &= Z_m^2 + Z_{n-m-k}^2 + \sum_{i=1}^k Z_{v_i}^2, \end{aligned} \tag{10}$$

на котором основывается применение следующей процедуры отбора выпадающих наблюдений.

7. Будем считать, что независимые переменные $Z_{(i_1)} < Z_{(i_2)} < \dots < Z_{(i_k)}$ выбраны таким образом, что они соответствуют k «подозреваемым» наблюдениям x_{i_1}, \dots, x_{i_k} , имеющим максимальные отклонения от их оцененных значений (это можно сделать последовательно). Нетрудно заметить, что $E(Z_{(ij)}) = \eta_j^\perp$ и $\max_v \psi_v(x) = \sum_{j=1}^k Z_{(ij)}^2 / Z_{n-m-k}^2$.

Если число выпадающих наблюдений неизвестно, но не больше k , то нужно решить, к какому из k подпространств $V_{rv} = \bigoplus_{j=1}^r W_{v_j}$ принадлежит вектор $E(Z_k^v) = (E(Z_{v_1}), E(Z_{v_2}), \dots, E(Z_{v_k}))$. Для поиска соответствую-

щей процедуры формализуем эту задачу известным образом [7], считая

$$\begin{aligned} H_k^* &= \hat{H}^{i_1} : \eta_1^\perp \neq 0 \quad (\eta^\perp \in V_{k_i}); \\ H_{k-1}^* &= \hat{H}^{i_2} : \eta_1^\perp = 0, \quad \eta_2^\perp \neq 0 \quad (\eta^\perp \in V_{(k-1)i}); \\ &\vdots \\ H_1^* &= \hat{H}^{i_k} : \eta_1^\perp = \eta_2^\perp = \dots = \eta_{k-1}^\perp = 0, \quad \eta_k^\perp \neq 0 \quad (\eta^\perp \in W_{i_1}); \\ &\hat{H}^{i_{k+1}} : \eta_1^\perp = \eta_2^\perp = \dots = \eta_k^\perp = 0. \end{aligned} \quad (11)$$

Данной последовательности гипотез противопоставим последовательность нулевых гипотез:

$$\begin{aligned} \hat{H}_0^{i_1} : \eta_1^\perp &= 0; \\ \hat{H}_0^{i_1 i_2} : \eta_1^\perp = \eta_2^\perp &= 0; \\ &\vdots \\ \hat{H}_0^{i_1 i_2 \dots i_k} : \eta_1^\perp = \eta_2^\perp = \dots = \eta_k^\perp &= 0. \end{aligned} \quad (12)$$

Ясно, что если некоторая гипотеза $H_0^{i_1 \dots i_l}$ верна, то предшествующие ей гипотезы также должны быть верны. Если же она неверна, то неверны и все последующие гипотезы: $(\hat{H}_0^{i_1} \supset \hat{H}_0^{i_1 i_2} \supset \dots \supset \hat{H}_0^{i_1 i_2 \dots i_k})$. Кроме того, семейства множеств, определяемые гипотезами (11) и (12), связаны соотношениями:

$$\begin{aligned} \hat{H}_0^{i_1} &= \hat{H}^{i_2} : \eta_1^\perp \in W_{i_2} \oplus W_{i_3} \oplus \dots \oplus W_{i_k}; \\ \hat{H}_0^{i_1 i_2} &= \hat{H}^{i_3} : \eta_1^\perp \in W_{i_3} \oplus W_{i_4} \oplus \dots \oplus W_{i_k}; \\ &\vdots \\ \hat{H}_0^{i_1 i_2 \dots i_k} &= \hat{H}^{i_{k+1}} : \eta_1^\perp \in \emptyset. \end{aligned} \quad (13)$$

Здесь будем ограничивать вероятности ошибок, связанных с принятием решений о том, что коэффициенты $\eta_i^\perp, i = 1, k$, отличны от нуля, в то время как на самом деле они равны нулю. С этой целью зададим каждой из нулевых гипотез определенный уровень значимости:

$$\begin{aligned} \pi_1 &= P \{ \text{отвергнуть } \hat{H}_0^{i_1} | \hat{H}_0^{i_1} \}; \\ \pi_1 + \pi_2 &= P \{ \text{отвергнуть } \hat{H}_0^{i_1 i_2} | H_0^{i_1 i_2} \}; \\ &\vdots \\ \pi_1 + \pi_2 + \dots + \pi_k &= P \{ \text{отвергнуть } \hat{H}_0^{i_1 i_2 \dots i_k} | \hat{H}_0^{i_1 i_2 \dots i_k} \} = \\ &= P \{ \text{отвергнуть } \hat{H}^{i_{k+1}} | \hat{H}^{i_{k+1}} \}, \end{aligned} \quad (14)$$

где $\pi_i \geq 0$ и $\sum_{i=1}^k \pi_i \leq 1$. В такой постановке вероятность отклонения каждой последующей нулевой гипотезы, когда она верна, не меньше, чем аналогичная вероятность для предшествующей гипотезы. Использование этих соотношений дает определенные значения вероятностям принятия гипотезы \hat{H}^{i_1} , когда на самом деле верна одна из последующих гипотез $\hat{H}^{i_{k+1}}, \dots, \hat{H}^{i_k}$:

$$\begin{aligned} \pi_1 &= P \{ \text{принять } \hat{H}^{i_1} | \hat{H}^{i_2} \cup \dots \cup \hat{H}^{i_{k+1}} \}; \\ \pi_2 &= P \{ \text{принять } \hat{H}^{i_2} | \hat{H}^{i_3} \cup \dots \cup \hat{H}^{i_{k+1}} \}; \\ &\vdots \\ \pi_k &= P \{ \text{принять } \hat{H}^{i_k} | \hat{H}^{i_{k+1}} \}. \end{aligned} \quad (15)$$

Как было показано выше, критерий (8) с областью отклонения S_k является оптимальным для проверки гипотезы $\hat{H}_0^{i_1} : \eta_1^\perp = 0$ (дополнительной

$\hat{H}^{i_1 \dots i_l} = H^*$. Для проверки гипотезы $\hat{H}_0^{i_1 \dots i_l} : \eta_{i_1}^\perp = 0$ с уровнем значимости α в предположении, что $\eta_{i_1}^\perp = \dots = \eta_{i_{l-1}}^\perp = 0$, состоит в ее отклонении, если

$$t_l^2 = \frac{Z_{(i_l)}^2}{Z_{n-m-k}^2 + Z_{(i_1)}^2 + \dots + Z_{(i_{l-1})}^2} \geq \frac{t_l^2(\alpha)}{n-m-k+l+1}. \quad (16)$$

Здесь вероятность отклонения гипотезы $\hat{H}_0^{i_1 \dots i_l}$ не зависит от «мешающих» параметров $\theta_1, \theta_2, \dots, \theta_m, \sigma^2$ и $\eta_{i+1}^\perp, \dots, \eta_k^\perp$. Тем самым, если S_l — подобная область размера α_l для проверки гипотезы $\hat{H}_0^{i_1 \dots i_l} : \eta_{i_l}^\perp = 0$, или формально

$$P\{S_l \mid \eta_{i_1}^\perp = \eta_{i_2}^\perp = \dots = \eta_{i_{l-1}}^\perp = \eta_{i_l}^\perp = 0\} = \alpha_l$$

для всех $\Theta = (\Theta_1, \Theta_2, \dots, \Theta_m)^T, \sigma^2$, то S_l выделяет условную вероятность α_l на почти всех комбинациях значений достаточных для этих параметров статистик \hat{x}_{V_m} и $Z_{n-m-k+l}^2$ при $\eta_{i_1}^\perp = \eta_{i_2}^\perp = \dots = \eta_{i_l}^\perp = 0$.

Нерандомизированная статистическая процедура для данной задачи состоит в следующем [7]. Фиксируется набор $(k+1)$ попарно непересекающихся областей $R_1, R_2, \dots, R_k, R_{k+1}$ в выборочном пространстве V_n . Если выборочная точка попадает в R_j , то принимается гипотеза $\hat{H}^{i_j} = H_{k-j+1}^*$. Приписывание уровней значимости (15) приводит к тому, что эти области становятся «подобными» в том смысле, что при $\eta_{i_1}^\perp = \eta_{i_2}^\perp = \dots = \eta_{i_l}^\perp = 0$ вероятности попадания выборочной точки в R_1, R_2, \dots, R_l равны соответственно $\pi_1, \pi_2, \dots, \pi_l$ независимо от $(\Theta_1, \Theta_2, \dots, \Theta_m), \sigma^2$ и $\eta_{i+1}^\perp, \eta_{i+2}^\perp, \dots, \eta_k^\perp$. Пусть T_l^* — область, определяемая соотношением (16) при $\alpha_l = \pi_l / (1 - \pi_1 - \pi_2 - \dots - \pi_{l-1})$. Тогда согласно [7], какими бы ни были R_1, R_2, \dots, R_{l-1} , наилучший выбор области R_l состоит в том, что она должна являться частью множества T_l^* , не содержащей точек из R_1, R_2, \dots, R_{l-1} . При таком выборе ($R_l = W_{i_l}$) вероятность

$$\begin{aligned} P\{R_l \mid \eta_{i_1}^\perp = \eta_{i_2}^\perp = \dots = \eta_{i_{l-1}}^\perp = 0\} &= (1 - \pi_1 - \pi_2 - \dots \\ &\dots - \pi_{l-1}) P\{T_l^* \mid \eta_{i_1}^\perp = \dots = \eta_{i_l}^\perp = 0\} \end{aligned}$$

не зависит от выбора R_1, \dots, R_{l-1} и для каждого значения $\eta_{i_l}^\perp$ вероятность $P\{R_l \mid \eta_{i_1}^\perp, \eta_{i_2}^\perp = \dots = \eta_{i_{l-1}}^\perp = 0\}$ принимает максимальное значение. Оптимальная процедура состоит, по сути дела, в поочередной проверке гипотез $\eta_{i_1}^\perp = 0, \eta_{i_2}^\perp = 0, \dots$ до тех пор, пока либо какая-то из гипотез не будет отвергнута, либо будут приняты все гипотезы вплоть до $\hat{H}_0^{i_1 i_2 \dots i_k} = \hat{H}^{i_{k+1}}$. Если все α_l фиксированы и равны α , то $\pi_1 = \alpha, \pi_2 = \alpha(1-\alpha), \dots, \pi_k = \alpha(1-\alpha)^k, \pi_{k+1} = (1-\alpha)^{k+1}$. Наконец, удобные для практических расчетов аппроксимации распределения статистик t_l^2 можно получить, используя результаты статьи [8].

ЛИТЕРАТУРА

1. Ferguson T. S. On the rejection of outliers.— Proc. of the Fourth Berkeley Symposium on Mathematical Statistics and Probability. California, University of California, Press, 1961, vol. 1, p. 253—287.
2. Tietjen G. L., Moore R. H. Some Grubbs-type statistics for the detection of several outliers.— "Technometrics", 1972, vol. 14, N 3, p. 583—597.
3. Кирчук В. С., Луценко Б. Н. Исключение недостоверных данных.— «Автометрия», 1970, № 6, с. 47—54.
4. Рао С. Р. Линейные статистические методы и их применение. М., «Наука», 1968.
5. Леман Е. Л. Проверка статистических гипотез. М., «Наука», 1964.
6. Барра Ж.-Р. Основные понятия математической статистики. М., «Мир», 1974.
7. Андерсон Т. Статистический анализ временных рядов. М., «Мир», 1976.
8. Большев Л. Н. О критериях исключения резко выделяющихся наблюдений.— Труды Ин-та прикладной математики Тбилисского гос. ун-та. Вып. II. Тбилиси, 1969, с. 159—177.

Поступила в редакцию 21 января 1976 г.