

Л. В. РЫКОВА

НЕПАРАМЕТРИЧЕСКИЙ КРИТЕРИЙ  
КОЛМОГОРОВА — СМИРНОВА ДЛЯ ВЫБОРОК ИЗ КОНЕЧНЫХ  
СОВОКУПНОСТЕЙ

(Представлено академиком А. Н. Колмогоровым 22 X 1973)

Для проверки однородности двух выборок Н. В. Смирновым <sup>(1)</sup> были введены статистики

$$\begin{aligned} D_{m,n}^+ &= \sup_x [F_n(x) - G_m(x)], \quad D_{m,n}^- = -\inf_x [G_m(x) - F_n(x)], \\ D_{m,n} &= \sup_x |F_n(x) - G_m(x)|, \end{aligned} \tag{1}$$

где  $F_n(x)$  и  $G_m(x)$  — эмпирические функции распределения (э.ф.р.), построенные по результатам двух выборок  $\xi_1, \dots, \xi_n$  и  $\eta_1, \dots, \eta_m$  объема  $n$  и  $m$  соответственно. При этом предполагалось, что  $\xi_1, \dots, \xi_n$  ( $\eta_1, \dots, \eta_m$ ) — взаимно независимые случайные величины с функцией распределения  $F(x)$  ( $G(x)$ ). С помощью (1) можно проверять гипотезу о совпадении функций распределения  $F(x)$  и  $G(x)$ .

В данной работе исследуются свойства аналогичных статистик, которыми можно пользоваться при изучении конечной совокупности без предположения взаимной независимости наблюдений.

Пусть  $\mathfrak{P}$  — конечная совокупность объектов  $O_i$ , каждому из которых сопоставляется скалярная величина  $X_i$ ,  $i=1, \dots, N$ . Предположим, что  $X_i \neq X_j$ ,  $i \neq j$ , и объекты занумерованы в порядке возрастания величин  $X_i$ , т. е.  $X_1 < X_2 < \dots < X_N$ . Совокупность  $\mathfrak{P}$  будем характеризовать э.ф.р.  $F_N(x) = l/N$ , где  $l$  — число таких объектов, у которых  $X_i < x$ .

Из совокупности  $\mathfrak{P}$  берутся две случайные выборки без возвращения объема  $n_1$  и  $n_2$ . Результатам выборок сопоставим э.ф.р.  $F_{n_1}(x) = k_1/n_1$  и  $F_{n_2}(x) = k_2/n_2$  соответственно, где  $k_1$  ( $k_2$ ) равно числу объектов, попавших в первую (вторую) выборку, у которых  $X_i \leq x$ . Ниже будет показано, что при  $n_i \rightarrow \infty$ ,  $(n_i/N) \rightarrow v_i$ ,  $0 < v_i < 1$ ,  $i=1, 2$ ,

$$P \left\{ \sup_x \left( \frac{N-1}{N} \frac{n_1 n_2}{n_1 + n_2} \right)^{1/2} |F_{n_1}(x) - F_{n_2}(x)| < z \right\} \rightarrow K(z) = \sum_k (-1)^k e^{-2k^2 z^2}. \tag{2}$$

Случай  $v_2=0$  ( $v_1=0$ ) идентичен рассмотренному ранее <sup>(2)</sup>.

Определим величину

$$\rho_{n_1, n_2, N} = \sup_x \left( \frac{N-1}{N} \frac{n_1 n_2}{n_1 + n_2} \right)^{1/2} |F_{n_1}(x) - F_{n_2}(x)|. \tag{3}$$

Поскольку  $F_{n_1}(x)$  и  $F_{n_2}(x)$  постоянны на  $[X_i, X_{i+1}]$ ,  $i=1, \dots, N-1$ , и имеют скачки  $1/n_1$  и  $1/n_2$  соответственно, в точках  $X_i$ ,  $i=1, \dots, N$ , то отображение  $X_i$  в точку  $i/N$  отрезка  $[0, 1]$  не изменит величины (3). Поэтому распределение величины (3) совпадает с распределением

$$\rho_{n_1, n_2, N} = \max_{1 \leq i \leq N} \left( \frac{N-1}{N} \frac{n_1 n_2}{n_1 + n_2} \right)^{1/2} \left| \frac{x_1(i)}{n_1} - \frac{x_2(i)}{n_2} \right|,$$

где  $\kappa_i(l)$  равно числу объектов в  $i$ -й выборке,  $i=1, 2$ , номера которых меньше  $l$ .

**Теорема 1.** Двумерный процесс  $(\kappa_1(l), \kappa_2(l))$  является двумерным марковским процессом с вероятностями переходов

$$\begin{aligned} p_1 &= P\{\kappa_1(l+1) = k_1 + 1, \kappa_2(l+1) = k_2 | \kappa_1(l) = k_1, \kappa_2(l) = k_2\} = (n_1 - k_1)/(N - l), \\ p_2 &= P\{\kappa_1(l+1) = k_1, \kappa_2(l+1) = k_2 + 1 | \kappa_1(l) = k_1, \kappa_2(l) = k_2\} = (n_2 - k_2)/(N - l), \\ p &= P\{\kappa_1(l+1) = k_1, \kappa_2(l+1) = k_2 | \kappa_1(l) = k_1, \kappa_2(l) = k_2\} = 1 - p_1 - p_2. \end{aligned} \quad (4)$$

Доказательство теоремы состоит в проверке совпадения вероятностей отбора объектов с номерами  $l_1, l_2, \dots, l_{n_1}$  и  $m_1, m_2, \dots, m_{n_2}$  в первую и вторую выборки соответственно и вероятностей скачков процесса с вероятностями переходов (4) в тех же точках.

Поскольку для любых  $l_1, l_2$

$$P\{\kappa_1(l_1, l_2) = d_1, \kappa_2(l_1, l_2) = d_2\} = h_{N, l_2 - l_1}^{n_1, d_1} h_{N - n_1, l_2 - l_1 - d_1}^{n_2, d_2}, \quad (5)$$

где  $h_{N, M}^{n, m} = \binom{M}{m} \binom{N-M}{n-m} / \binom{N}{n}$  — гипергеометрические вероятности,

$\kappa_i(l, m) = \kappa_i(m) - \kappa_i(l)$ , то процесс  $(\kappa_1(l), \kappa_2(l))$  естественно называть двумерным гипергеометрическим процессом.

**Лемма 1.** Пусть  $m, n, M, N$  — целые числа,  $m \leq M \leq N$ ,  $m \leq n \leq N$ ,  $b(n, N) = n(N-n)/N$ ,  $0 < \underline{c} < \Delta < \bar{c} < 1$ ,

$$\nu_+ = \overline{\lim}_{n \rightarrow \infty} (n/N) < 1, \quad |M/N - \Delta| < 1/N, \quad \frac{|m-n|M/N|^3}{b^2(n, N)} = \alpha_{N, M}^{n, m} \rightarrow 0$$

при  $n \rightarrow \infty$ .

Тогда для любого  $\varepsilon > 0$

$$h_{N, M}^{n, m} = \frac{\exp\{- (m-n\Delta)^2 / [2b(n, N)\Delta(1-\Delta)]\}}{[2\pi b(n, N)\Delta(1-\Delta)]^{1/2}} (1 + \theta_{N, M}^{n, m}),$$

где

$$|\theta_{N, M}^{n, m}| < \gamma_{N, M}^{n, m} / [\Delta(1-\Delta)(1-\nu_+)(1-\varepsilon)]^2,$$

$$\gamma_{N, M}^{n, m} = \min(\alpha_{N, M}^{n, m}, n^{-1/2}).$$

Рассмотрим процесс  $\beta_{n_1, n_2}(l_1/N, l_2/N) = (\xi_{n_1, n_2}^{(1)}(l_1/N), \xi_{n_1, n_2}^{(2)}(l_2/N))$ , где

$$\xi_{n_1, n_2}^{(i)}\left(\frac{l}{N}\right) = \left( \frac{N-1}{N} \frac{n_1 n_2}{n_1 + n_2} \right)^{1/2} \left[ \frac{\kappa_i(l)}{n_i} - \frac{l}{N} \right].$$

Легко проверить, что  $E\beta_{n_1, n_2}(l_1/N, l_2/N) = 0$ ,

$$S = \|a_{ij}\| = \left( \begin{array}{cc} \frac{(N-n_1)n_2}{N(n_1+n_2)} \frac{l_1}{N} \left(1 - \frac{l_2}{N}\right), & -\frac{n_1 n_2}{(n_1+n_2)N} \frac{l_1}{N} \left(1 - \frac{l_2}{N}\right) \\ -\frac{n_1 n_2}{(n_1+n_2)N} \frac{l_1}{N} \left(1 - \frac{l_2}{N}\right), & \frac{(N-n_2)n_1}{(n_1+n_2)N} \frac{l_1}{N} \left(1 - \frac{l_2}{N}\right) \end{array} \right),$$

где

$$a_{ij} = E\xi_{n_1, n_2}^{(i)}(l_1/N) \xi_{n_1, n_2}^{(j)}(l_2/N), \quad i, j = 1, 2.$$

**Теорема 2.** При  $n_i \rightarrow \infty$ ,  $i=1, 2$ ,  $\overline{\lim}_{n_i \rightarrow \infty} (n_i/N) = \nu_i < 1$  конечномерные распределения процесса  $\beta_{n_1, n_2}(l_1/N, l_2/N)$  сходятся к конечномерным распределениям

лением двумерного гауссовского процесса  $\beta(t, s) = (\beta_1(t), \beta_2(s))$  с нулевым математическим ожиданием и ковариационной матрицей

$$Q = \begin{pmatrix} \frac{v_2(1-v_1)t(1-s)}{v_1+v_2} & -\frac{v_1v_2t(1-s)}{v_1+v_2} \\ -\frac{v_1v_2t(1-s)}{v_1+v_2} & \frac{v_1(1-v_2)t(1-s)}{v_1+v_2} \end{pmatrix}, \quad t < s.$$

При доказательстве теоремы надо воспользоваться марковскими свойствами процесса  $\beta_{n_1, n_2}(l_1/N, l_2/N)$ , равенством (5), леммой 1 и методом интегральных сумм ((3), стр. 188).

Лемма 2. Для любых  $t_1$  и  $t_2$  из  $[0, 1]$ , всех  $N$  и  $n_i$  таких, что  $0 < v_i \leq n_i/N \leq \bar{v} < 1$  существует такая  $H = H(v, \bar{v}) > 0$ , что

$$\mu_i^{(t)} = E|\xi_{n_1, n_2}^{(t)}(t_1) - \xi_{n_1, n_2}^{(t)}(t_2)|^4 \leq H|t_1 - t_2|^2, \quad i=1, 2. \quad (6)$$

Неравенство (6) получается как следствие гипергеометричности распределения приращений процесса  $\kappa_i(l)$ .

Выполнение условий теоремы 2 и леммы 2 дает возможность воспользоваться теоремой (4) о том, что для всех непрерывных на  $C_{[0, 1]}$  функционалов  $f$  распределение  $f(\beta_{n_1, n_2}(l_1/N, l_2/N))$  будет сходиться к распределению  $f(\beta(t, t))$ ,  $C_{[0, 1]}$  — пространство всех непрерывных на  $[0, 1]$  функций с метрикой

$$\rho(x, y) = \sup_{0 \leq t \leq 1} |x(t) - y(t)|.$$

В частности, распределение

$$f(\beta_{n_1, n_2}(l_1/N, l_2/N)) = \max_{1 \leq i \leq N} |\xi_{n_1, n_2}^{(1)}(l_1/N) - \xi_{n_1, n_2}^{(2)}(l_2/N)| = \rho_{n_1, n_2, N}$$

будет сходиться к распределению

$$f(\beta(t, t)) = \max_{0 \leq t \leq 1} |\beta_1(t) - \beta_2(t)| = \max_{0 \leq t \leq 1} |\beta_0(t)|.$$

Поскольку  $E\beta_0(t) = 0$ ,  $E\beta_0(t)\beta_0(s) = t(1-s)$ ,  $t < s$ , то  $\beta_0(t)$  является условным винеровским процессом на  $[0, 1]$  с закрепленными концами  $\beta_0(0) = \beta_0(1) = 1$ , а распределение  $\max_{0 \leq t \leq 1} |\beta_0(t)|$  является известным распределением Колмогорова  $K(z)$ . Таким образом, верна

Теорема 3. При  $n_i \rightarrow \infty$ ,  $\lim_{n_i \rightarrow \infty} (n_i/N) = v_i$ ,  $0 \leq v_i < 1$ ,  $i=1, 2$ , распределение  $\rho_{n_1, n_2, N}$  сходится к распределению  $\beta_0 = \max_{0 \leq t \leq 1} |\beta_0(t)|$ , т. е. справедливо соотношение (2).

Статистика  $\rho_{n_1, n_2, N}$  может быть использована в следующих математических моделях:

а) Объекты совокупности  $\Phi$  характеризуются скалярными величинами  $X_i = X_i(s)$ , вообще говоря зависящими от времени (хранения)  $s$ ,  $s \geq 0$ . Для проверки гипотезы о неизменности характеристик, т. е. гипотезы  $H_0 = \{X_i(s_1) = X_i(s_2), i \in [1, 2, \dots, N]\}$ , в два момента времени  $s_1 < s_2$  берутся две случайные выборки без возвращения объемов  $n_1, n_2$ . Если  $n_i \gg 1$ ,  $i=1, 2$ ,  $N \gg 1$ ,  $n_i/N \leq v_i < 1$ , то при тестовой статистике  $\rho_{n_1, n_2, N}$  критическое множество с уровнем значимости, близким к  $\alpha$ , имеет вид

$$\rho_{n_1, n_2, N} > z_{1-\alpha}, \quad K(z_{1-\alpha}) = 1 - \alpha. \quad (7)$$

б) Выясняется вопрос о различии влияния двух типов воздействия на характеристики объектов  $O_i$  конечной совокупности  $\Phi$ ,  $i=1, \dots, N$ . Для

этого объекты  $\mathfrak{P}$  разбиваются на три группы объемом  $n_1$ ,  $n_2$  и  $N - (n_1 + n_2)$  на основе случайных выборок без возвращения объема  $n_1$  и  $n_2$ . На объекты, составляющие первую выборку, оказывается воздействие первого типа, на объекты, составляющие вторую выборку, оказывается воздействие второго типа. За нулевую гипотезу принимается идентичность воздействия обоих типов или отсутствие его. Тестовая статистика  $\rho_{n_1, n_2, N}$  имеет то же критическое множество (7). Отбор объектов для проверки идентичности воздействий можно проводить с постоянной вероятностью. При этом объемы выборок  $n_1$  и  $n_2$  будут случайными величинами.

Заметим, что относительно предельных распределений статистик

$$\rho_{n_1, n_2, N}^+ = \sup_x [F_{n_1}(x) - F_{n_2}(x)], \quad \rho_{n_1, n_2, N}^- = -\inf_x [F_{n_2}(x) - F_{n_1}(x)],$$

$$T = \int_{-\infty}^{\infty} [F_{n_1}(x) - F_{n_2}(x)]^2 dH_{n_1+n_2}(x),$$

где

$$H_{n_1+n_2}(x) = \varphi_1(n_1, n_2, N) F_{n_1}(x) + \varphi_2(n_1, n_2, N) F_{n_2}(x),$$

будут справедливы утверждения, аналогичные теореме 3, вследствие непрерывности этих статистик, рассматриваемых как функционалы от траекторий гипергеометрических процессов.

Московский государственный университет  
им. М. В. Ломоносова

Поступило  
4 X 1973

#### ЦИТИРОВАННАЯ ЛИТЕРАТУРА

<sup>1</sup> Н. В. Смирнов, Бюлл. Московск. унив., сер. А, 2, 3 (1939). <sup>2</sup> Ю. К. Беляев, Л. В. Рыкова, ДАН, 210, № 6 (1973). <sup>3</sup> В. Феллер, Введение в теорию вероятностей, 1, М., 1967. <sup>4</sup> И. И. Гихман, А. В. Скороход, Введение в теорию случайных процессов, «Наука», 1965.