

ПРОВЕРКА ГИПОТЕЗ ОБ ОДНОРОДНОСТИ И СИММЕТРИЧНОСТИ РАСПРЕДЕЛЕНИЙ ДЛЯ МНОГОМЕРНЫХ ДАННЫХ

Н.К. БАКИРОВ

Аннотация. Рассматриваются задачи проверки непараметрических гипотез для многомерных данных.

Ключевые слова: проверка непараметрических гипотез для многомерных данных.

1. ВВЕДЕНИЕ

В настоящей работе изучаются новые критерии проверки непараметрических гипотез для многомерных данных: об однородности распределений двух случайных векторов и симметричности многомерного распределения. Построенные тестовые статистики имеют простую структуру и инвариантны к линейным преобразованиям данных, после надлежащей нормировки их распределения слабо сходятся при нулевой гипотезе к распределениям типа омега-квадрат. Построенные критерии состоятельны против широкого класса альтернатив при минимальных моментных ограничениях.

2. ПРОВЕРКА СИММЕТРИЧНОСТИ

Пусть X_1, X_2, \dots, X_n — повторная выборка с общей функцией распределения (ф.р.) $F(x)$. Для проверки гипотезы о симметричности:

$$H_0 : 1 - F(x + 0) - F(-x) = 0 \quad \text{для всех } x \in R^1 \quad (1)$$

в одномерном случае используют известные критерии: w^2 , Колмогоровского типа, Уотсона-Дарлингга, Хилла-Рао, знаковых статистик и др. [1]. Асимптотические уровни значимости могут быть найдены в предположении, что ф.р. $F(x)$ непрерывна. К сожалению, асимптотическое распределение многомерных аналогов этих статистик зависит от ф.р. F , так что они не являются подобными. Задачи проверки симметричности многомерных распределений рассматривались, в частности, в работах [2]–[10]. В настоящем параграфе мы строим тесты для проверки многомерной симметричности ф.р. с заданным асимптотическим уровнем значимости, инвариантные к линейным преобразованиям данных.

Итак, пусть X_1, X_2, \dots, X_n , — повторная выборка $X_k \in R^d$, $X_1 \neq 0, a.s.$, $E|X_1| < \infty$. Мы рассматриваем гипотезы о симметричности распределения выборки с центром в нуле:

H_{01} : *диагональная симметричность* или $X_1 \stackrel{D}{=} -X_1$;

H_{02} : *сферическая симметричность* или $X_1 \stackrel{D}{=} CX_1$ для всех ортогональных матриц C ;

H_{03} : *эллиптическая симметричность*, то есть случайный вектор MX_1 сферически симметричен для некоторой положительно определенной матрицы M , а также варианты этих гипотез с неизвестными центрами симметрии.

N.K. BAKIROV, TESTING HOMOGENEITY AND SYMMETRY FOR MULTIVARIATE DATA.

© БАКИРОВ Н.К. 2009.

Поступила 01 июня 2009 г.

Все наши построения базируются на использовании эмпирических характеристических функций:

$$f_n(t) = \frac{1}{n} \sum_{i=1}^n \exp\{i(t, X_i)\}, \quad (2)$$

где (\cdot, \cdot) — скалярное произведение в R^d . В работах [11]-[12] они были применены для проверки гипотезы о независимости двух многомерных распределений.

§ 1. Диагональная симметричность

Пусть $f(t)$ — характеристическая функция случайного вектора X_1 . Легко видеть, что

$$H_{01} \iff f(t) = f(-t), \forall t \in R^d \iff J(f, \varphi) \stackrel{\text{def}}{=} \int_{R^d} |f(t) - f(-t)|^2 \varphi(t) dt = 0$$

для любой весовой функции $\varphi(t)$ такой, что интеграл существует и $\varphi(t) > 0$ почти всюду (п.в.). Мы будем предполагать, что $E|X_1| < \infty$.

Рассмотрим статистику

$$J(f_n, \varphi_0) = \int_{R^d} |f_n(t) - f_n(-t)|^2 \varphi_0(t) dt, \quad \varphi_0(t) = |t|^{-1-d}.$$

Заметим, что $|f(t) - f(-t)| \leq E|X_1||t|$, поэтому $J(f, \varphi_0)$ имеет интегрируемую особенность в нуле. По аналогичным причинам интеграл $J(f_n, \varphi_0)$ существует. Учитывая тождество

$$\int_{R^d} (1 - \exp\{i(t, X)\})|t|^{-d-\gamma} dt = C(d, \gamma)|X|^\gamma, \quad \forall \gamma \in (0, 2), \quad (3)$$

где $C(d, \gamma)$ — положительная константа, мы можем легко вычислить

$$J(f_n, \varphi_0) = 2C(d, 1)n^{-2} \left(\sum_{i,j=1}^n |X_i + X_j| - \sum_{i,j=1}^n |X_i - X_j| \right)$$

(интеграл (3) сходится в нулевой точке в смысле главного значения, при необходимости, ввиду симметричности функции $\sin(t, X)$, мы можем в (3) заменить $\exp\{i(t, X)\}$ на $\cos(t, X)$ и считать несобственный интеграл (3) сходящимся в обычном смысле). Используя формулу 3.241.4 [13] и рассматривая интеграл (3) как повторный, можно показать, что $C(d, 1) = \pi^{\frac{d+1}{2}} / \Gamma\left(\frac{d+1}{2}\right)$. Далее, по закону больших чисел для статистик Мизеса [1]

$$\begin{aligned} J(f_n, \varphi_0) &\xrightarrow[n \rightarrow \infty]{P} 2C(d, 1) (E|X_1 + X_2| - E|X_1 - X_2|) = \\ &= 2C(d, 1) \int_{R^d} |f(t) - f(-t)|^2 \varphi_0(t) dt. \end{aligned} \quad (4)$$

Так что, для всех альтернатив к гипотезе H_{01} мы имеем $J(f_n, \varphi_0) \xrightarrow[n \rightarrow \infty]{P} \infty$.

Аналогично

$$G(f_n) \stackrel{\text{def}}{=} 2 \int_{R^d} (1 - f_n(2t)) \varphi_0(t) dt = 4C(d, 1)n^{-1} \sum_{i=1}^n |X_i| \xrightarrow[n \rightarrow \infty]{P}$$

$$\xrightarrow[n \rightarrow \infty]{P} 4C(d, 1)E|X_1| = 2 \int_{R^d} (1 - f(2t)) \varphi_0(t) dt = G(f).$$

Заметим, что $G(f) \neq 0$, если $X_1 \not\equiv 0$ п.в. Рассмотрим теперь асимптотическое поведение статистики $J(f_n, \varphi_0)$ при нулевой гипотезе H_{01} . Обозначим

$$\xi_n(t) = \sqrt{n}(f_n(t) - f_n(-t)).$$

Нетрудно вычислить, что при $H_{01}, \forall t, s$

$$E\xi_n(t) = 0, \quad E\xi_n(t)\overline{\xi_n(s)} = 2(f(t-s) - f(t+s)). \quad (5)$$

Конечномерные распределения случайного процесса (сл.пр.) $\xi_n(t)$ сходятся при $n \rightarrow \infty$ к конечномерным распределениям гауссовского (сл.пр.) $\xi(t)$ с моментами (5) в силу многомерной центральной предельной теоремы (ЦПТ). Из (5) следует, что

$$E|\xi_n(t) - \xi_n(s)|^2 = E|\xi(t) - \xi(s)|^2 = 4(1 - f(t-s)) + 2(2f(t+s) - f(2t) - f(2s)) \leq 8|t-s|E|X_1|. \quad (6)$$

Поэтому [14], гауссовский (сл.пр.) $\xi(t)$ имеет модификацию с непрерывными траекториями. С другой стороны,

$$E|\xi_n(t)|^2 = E|\xi(t)|^2 = 2(1 - f(2t)) \leq 4 \min(1, |t|E|X_1|). \quad (7)$$

Следовательно, $\int_{R^d} E|\xi(t)|^2 \varphi_0(t) dt < \infty$, поэтому по теореме Фубини мы можем определить $Q \stackrel{\text{def}}{=} \int_{R^d} |\xi(t)|^2 \varphi_0(t) dt$ (для непрерывной модификации $\xi(t)$) как римановский интеграл.

Обозначим $A_1 = \{t : |t| < \varepsilon\}$, $A_2 = \{t : |t| > 1/\varepsilon\}$, $A_3 = \{t : \varepsilon \leq |t| \leq 1/\varepsilon\}$. Применяя (7), мы можем записать

$$E \int_{A_2} (|\xi_n(t)|^2 + |\xi(t)|^2) \varphi_0(t) dt \leq 8 \int_{|t|>1/\varepsilon} \varphi_0(t) dt = 8w_d \varepsilon, \quad (8)$$

где w_d — площадь единичной сферы в R^d . С другой стороны, используя (7), запишем

$$\begin{aligned} E \int_{A_1} (|\xi_n(t)|^2 + |\xi(t)|^2) \varphi_0(t) dt &= 4 \int_{|t|<\varepsilon} (1 - f(2t)) \varphi_0(t) dt = \\ &= 4E|X_1| \delta, \quad \delta = \int_{|t|<\varepsilon|X_1|} (1 - \cos 2t_1) \varphi_0(t) dt, \end{aligned}$$

здесь $t = (t_1, t_2, \dots, t_d)$ и мы использовали замену переменных. Очевидно, $0 \leq \delta \leq 2C(d, 1)$ и $\delta \xrightarrow{\text{П.В.}} 0$, если $\varepsilon \rightarrow 0$. Итак, по теореме Лебега об ограниченной сходимости

$$4E|X_1| \delta \xrightarrow{\varepsilon \rightarrow 0} 0. \quad (9)$$

Далее, для любого разбиения $\Delta_k, k = 1, 2, \dots, N$ множества A_3 ($\max_k \text{diam } \Delta_k = \tau$) и $t_k \in \Delta_k$ имеем

$$\int_{A_3} |\xi_n(t)|^2 \varphi_0(t) dt = \sum_{k=1}^N |\xi_n(t_k)|^2 \int_{\Delta_k} \varphi_0(t) dt + \alpha_1, \quad (10)$$

где в силу (6),(7)

$$\begin{aligned} E|\alpha_1| &\leq \sum_{k=1}^N \int_{\Delta_k} E||\xi_n(t)|^2 - |\xi_n(t_k)|^2| \varphi_0(t) dt \leq \\ &\leq 8 \left(\tau E|X_1| + \sqrt{2\tau E|X_1|} \right) \int_{A_3} \varphi_0(t) dt \xrightarrow{\tau \rightarrow 0} 0. \end{aligned}$$

Представление, аналогичное (10), справедливо и для случайной величины (сл.в.) $\int_{A_3} |\xi(t)|^2 \varphi_0(t) dt$. Далее, в силу сходимости конечномерных распределений сл.пр. $\xi_n(t)$ к конечномерным распределениям сл.пр. $\xi(t)$ имеем, что

$$\sum_{k=1}^N |\xi_n(t_k)|^2 \int_{\Delta_k} \varphi_0(t) dt \xrightarrow[n \rightarrow \infty]{D} \sum_{k=1}^N |\xi(t_k)|^2 \int_{\Delta_k} \varphi_0(t) dt. \quad (11)$$

Объединяя теперь (8)–(11), получаем

$$Q = \int_{R^d} |\xi(t)|^2 \varphi_0(t) dt = D - \lim_{n \rightarrow \infty} \int_{R^d} |\xi_n(t)|^2 \varphi_0(t) dt,$$

$$Q = D - \lim_{\substack{\varepsilon \rightarrow 0 \\ \tau \rightarrow 0}} Q_N, \quad Q_N = \sum_{k=1}^N |\xi(t_k)|^2 \int_{\Delta_k} \varphi_0(t) dt \quad (12)$$

для некоторого разбиения Δ_k множества A_3 с $\tau \xrightarrow{\varepsilon \rightarrow 0} 0$. Поэтому, в частности, сл.в. Q есть слабый предел квадратичных форм Q_N от централизованных гауссовских сл.в. и, следовательно, мы можем применить неравенство из [15]:

$$P \left\{ \frac{Q}{EQ} > \left(\Phi^{-1} \left(1 - \frac{\alpha}{2} \right) \right)^2 \right\} \leq \alpha, \quad \forall \alpha \leq 0.21515\dots$$

Определим тестовую статистику равенством

$$T_n = \frac{\sum_{i,j=1} |X_i + X_j| - \sum_{i,j=1} |X_i - X_j|}{2 \sum_{i=1} |X_i|} = \frac{nJ(f_n, \varphi_0)}{G(f_n)}.$$

Выше было показано, что при любой альтернативе

$$\frac{T_n}{n} \xrightarrow[n \rightarrow \infty]{P} \text{const} > 0,$$

а при нулевой гипотезе H_{01}

$$T_n \xrightarrow[n \rightarrow \infty]{D} Q,$$

где Q есть слабый предел квадратичных форм от централизованных гауссовских сл.в., $EQ = 1$.

Таким образом, справедлива следующая теорема.

Теорема 1. *Критерий, отвергающий гипотезу H_{01} , при*

$$T_n \geq \Lambda, \quad \Lambda = \left(\Phi^{-1} \left(1 - \frac{\alpha}{2} \right) \right)^2 \quad (13)$$

состоятелен против всех альтернатив, он имеет асимптотический уровень значимости не более чем $\alpha, \forall \alpha \leq 0.21515\dots$

Рассмотрим теперь асимптотическое поведение мощности β_n критерия (13). Запишем

$$1 - \beta_n = P\{T_n \leq \Lambda\} = P\{U_n - EU_n \geq \nu + (1 - \Lambda) \frac{2}{n-1} \tilde{X}\}, \quad (14)$$

где

$$U_n = \frac{2}{n(n-1)} \sum_{i < j} (|X_i - X_j| - |X_i + X_j|), \quad \tilde{X} = \frac{1}{n} \sum_{i=1}^n |X_i|,$$

$$\nu = -EU_n = E(|X_1 + X_2| - |X_1 - X_2|),$$

заметим, что $\nu > 0$ при альтернативах и в случае $P\{|X_1| = 0\} > 0$ мы определяем мощность β_n как правую часть (14).

Пусть выполнено условие

C1) $\forall H > 0 : E \exp\{H|X_1|^2\} < \infty$,

тогда по экспоненциальному неравенству Чебышева

$$P \left\{ \frac{2(1 - \Lambda)\tilde{X}}{n-1} > \frac{\nu}{\sqrt{n}} \right\} \leq P \left\{ \sum_{k=1}^n |X_k| \geq \frac{(n-1)\sqrt{n}\nu}{2(1 + \Lambda)} \right\} \leq$$

$$\leq \exp \left\{ -\frac{(n-1)\sqrt{n}\nu}{2(1 + \Lambda)} \right\} (E \exp\{|X_1|\})^n \leq e^{-Cn\sqrt{n}} \quad (15)$$

для некоторой положительной константы C . С другой стороны, при условии С1) по теореме о больших отклонениях Р. Дасгупты для невырожденных U -статистик (см. теорему 1.6.4. в [1]) для любой последовательности $\gamma_n \xrightarrow[n \rightarrow \infty]{} 0$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \ln P\{U_n - EU_n > \nu + \gamma_n\} = u\left(\frac{\nu}{2\sigma}\right), \quad (16)$$

где $\sigma^2 = D\psi$, $\psi = |X_1 + X_2| - |X_1 - X_2|$,

$$u(x) = \ln \inf_{t \geq 0} (e^{-tn} E e^{t(\nu - \psi)}).$$

Объединяя (15) и (16), получаем

Предложение 1. Для всех альтернатив, удовлетворяющих С1)

$$\lim_{n \rightarrow \infty} \frac{1}{n} \ln(1 - \beta_n) = u\left(\frac{\nu}{2\sigma}\right).$$

Заметим, что $-2u(\frac{\nu}{2\sigma})$ есть показатель Ходжеса-Лемана критерия (13) и

$$u\left(\frac{\nu}{2\sigma}\right) = -\frac{\nu^2}{8\sigma^2}(1 + o(1)),$$

когда $\nu \rightarrow 0$, $\sigma > const > 0$ (см. [1], § 1.2).

Замечание 1. Определим коэффициент диагональной симметричности как:

$$SYM = \frac{T_n}{n} = \frac{\sum_{i,j=1} |X_i + X_j| - \sum_{i,j=1} |X_i - X_j|}{2n \sum_{i=1} |X_i|}.$$

Ясно, что $0 \leq SYM \leq 1$ и $SYM = 1$ тогда и только тогда, когда сл.в. $X_i \equiv const \neq 0, \forall i$. Некоторые из его свойств даны в (4), теореме 1 и предложении 1. Заметим, что коэффициент SYM инвариантен к ортогональным преобразованиям данных и изменению масштаба.

Замечание 2. Если требование $E|X_1| < \infty$ для нулевой гипотезы и альтернатив не выполнено, то тогда можно заметить, что гипотезы H_{01} и H_{02} эквивалентны аналогичным гипотезам с повторной выборкой Y_1, Y_2, \dots, Y_n , где $Y_k = X_k |X_k|^{-1} \arctan |X_k|$, $k = 1, 2, \dots, n$, (если сл.в. X_1 не имеет моментов) или $Y_k = X_k |X_k|^{\gamma-1}$ (если мы знаем, что $E|X_1|^\gamma < \infty$). Заметим, что в первом случае соответствующее значение T_n не инвариантно к изменению масштаба.

Замечание 3. Наши рассуждения остаются справедливыми и для коэффициента

$$T'_n = \frac{\sum_{i,j=1} |X_i + X_j|^\gamma - \sum_{i,j=1} |X_i - X_j|^\gamma}{2 \sum_{i=1} |X_i|^\gamma}, \quad \gamma \in (0, 2).$$

Выясним теперь вопрос: "которое γ лучше?" в некотором специальном случае и в некотором специальном смысле. А именно, рассмотрим локальную ($\nu \rightarrow 0$) эффективность по Ходжесу-Леману, сдвиговые альтернативы с $X_1 \stackrel{D}{=} N(\theta, E)$, где $\theta = EX_1$ — параметр сдвига, E — единичная матрица. Заметим, что сл.в. $|X_1 - \theta|^2$ имеет распределение с d -степенями свободы. Можно подсчитать, что в данном случае при $|\theta| \rightarrow 0$ показатель Ходжеса-Лемана есть

$$\frac{\nu^2}{4\sigma^2} = |\theta|^4 G(\gamma)(1 + o(1)),$$

где

$$G(\gamma) = \frac{\nu^2}{8d^2} \left(\frac{\Gamma(\frac{d}{2})\Gamma(\frac{d}{2} + \nu)}{\Gamma^2(\frac{d+\nu}{2})} - 1 \right)^{-1}.$$

Ниже мы доказываем, что $G(\gamma)$ возрастает по γ при $\gamma \in (0, 1]$ так, что ответ такой: чем больше γ , тем лучше.

Лемма 1. $G'(\gamma) > 0, \forall \gamma \in (0, 1], \forall d$.

Доказательство. Это неравенство эквивалентно следующему:

$$1 - \frac{\Gamma^2(n+x)}{\Gamma(n)\Gamma(n+2x)} \geq x \left(\frac{\Gamma'(n+2x)}{\Gamma(n+2x)} - \frac{\Gamma'(n+x)}{\Gamma(n+x)} \right), \quad (17)$$

где $n = d/2, x = \gamma/2$. Известна следующая формула [16]:

$$\Psi(1+z) \stackrel{\text{def}}{=} \frac{\Gamma'(1+z)}{\Gamma(1+z)} = -C_0 + \sum_{m=1}^{\infty} \frac{z}{m(m+z)}, \quad (18)$$

где C_0 — абсолютная константа. Теперь (17) может быть переписано как

$$1 - \exp\left\{ \int_n^{n+x} \Psi(z) dz - \int_{n+x}^{n+2x} \Psi(z) dz \right\} \geq x (\Psi(n+2x) - \Psi(n+x)),$$

или в силу (18) для $v = m + n + x - 1$

$$1 - \prod_{m=1}^{\infty} \left(1 - \frac{x^2}{v^2} \right) \geq x^2 \sum_{m=1}^{\infty} \frac{1}{v(v+x)}. \quad (19)$$

Применяя здесь неравенство $\prod(1 - a_i) \leq 1 - \sum a_i + \sum_{i < j} a_i a_j$ для $a_i \geq 0$, мы видим, что (19) верно, если

$$\sum_{m=1}^{\infty} \frac{1}{v^2(v+x)} + \frac{x}{2} \sum_{m=1}^{\infty} \frac{1}{v^4} \geq \frac{x}{2} \left(\sum_{m=1}^{\infty} \frac{1}{v^2} \right)^2. \quad (20)$$

Далее, оценивая суммы соответствующими интегралами, мы можем записать, что

$$\begin{aligned} \sum_{m=1}^{\infty} \frac{1}{v^2(v+x)} &\geq \frac{1}{(n+x)^2(n+2x)} + \int_{n+x+1}^{\infty} \frac{dz}{z^2(z+x)} \geq \\ &\geq \frac{1}{(n+x)^2(n+2x)} + \frac{1}{2(n+x+1)^2} - \frac{x}{3(n+x+1)^3}, \\ \sum_{m=1}^{\infty} \frac{1}{v^4} &\geq \frac{1}{(n+x)^4} + \frac{1}{3(n+x+1)^3}, \quad \sum_{m=1}^{\infty} \frac{1}{v^2} \leq \frac{1}{(n+x)^2} + \frac{1}{n+x}. \end{aligned}$$

Применяя эти неравенства, мы можем доказать (20).

Лемма доказана.

Замечание. Аналогично можно показать, что $\forall \gamma \in (1, 2) \exists d_0 = d_0(\gamma) : \forall d > d_0$
 $G(\gamma) = \sup\{G(s) | s \leq \gamma\}$.

§ 2. Сферическая симметричность

Гипотеза о сферической симметричности распределения выборки состоит в следующем:
 $H_{02} : X_1 \stackrel{D}{=} CX_1$ для всех ортогональных матриц C .

В одномерном случае ($d = 1$) $H_{01} \equiv H_{02}$, поэтому мы предполагаем, что $d \geq 2$. Как легко видеть,

$$H_{02} \iff f(t) = f(Ct), \text{ для всех } t \in R^d \text{ и ортогональных матриц } C \iff$$

$$\iff R(f, \varphi) \stackrel{\text{def}}{=} \frac{1}{2} \int \int_{R^d} |f(t) - f(C^T t)|^2 \varphi(t) dt d\chi(C) = 0$$

для любой весовой функции $\varphi(t)$ такой, что интеграл существует и $\varphi(t) > 0$ п.в., здесь $\chi(C)$ — нормализованная мера Хаара на множестве ортогональных матриц, $\int d\chi(C) = 1$.

Пусть

$$f_n(t) = n^{-1} \sum_{k=1}^n \exp\{i(t, X_k)\}, \quad f_n(t, C) = n^{-1} \sum_{k=1}^n \exp\{i(t, CX_k)\}$$

— выборочные характеристические функции. Рассмотрим статистику

$$\begin{aligned} R(f_n, \varphi_0) &= \frac{1}{2} \int \int_{R^d} |f_n(t) - f_n(t, C)|^2 \varphi_0(t) dt d\chi(C) = \\ &= C(d, 1) n^{-2} \sum_{i,j=1}^n (g(|X_i|, |X_j|) - |X_i - X_j|), \end{aligned}$$

где

$$g(u, v) = \frac{w_{d-1}}{w_d} \int_{-1}^1 \sqrt{u^2 - 2uvx + v^2} (1 - x^2)^{\frac{d-3}{2}} dx, \quad w_1 \stackrel{\text{def}}{=} 2$$

и мы использовали формулу (3) § 1 вместе с равенством: $\forall a, b \in R^d$

$$\begin{aligned} \int |a - Cb| d\chi(C) &= \int \sqrt{|a|^2 - 2(a, Cb) + |b|^2} d\chi(C) = \\ &= E \sqrt{|a|^2 - 2|a||b|\xi + |b|^2} = g(|a|, |b|), \end{aligned}$$

где ξ — первая координата случайного вектора равномерно распределенного на единичной сфере в R^d . Заметим, что $g(u, v) \leq |u| + |v|$ и

$$g(u, u) = \mu_d |u|, \quad \mu_d = \frac{2^{d-1} \Gamma(\frac{d}{2}) \Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(d - \frac{1}{2})},$$

здесь $\Gamma(\cdot)$ есть гамма-функция и мы использовали формулу: $w_d = 2\pi^{d/2} / \Gamma(\frac{d}{2})$. Для $u \neq v$ функция $g(u, v)$ может быть вычислена при помощи формулы 3.665.2 из [13].

По закону больших чисел для статистик Мизеса мы имеем, также как и в § 1,

$$R(f_n, \varphi_0) \xrightarrow[n \rightarrow \infty]{P} R(f, \varphi_0). \quad (21)$$

Далее, при нулевой гипотезе

$$\begin{aligned} nER(f_n, \varphi_0) &= C(d, 1) Eg(|X_1|, |X_1|) = \\ &= P - \lim_{n \rightarrow \infty} V_n, \quad V_n = C(d, 1) \frac{\mu_d}{n} \sum_{k=1}^n |X_k|. \end{aligned}$$

Окончательно определим тестовую статистику формулой

$$Q_n = \frac{nR(f_n, \varphi_0)}{V_n} = \frac{\sum_{i,j=1}^n (g(|X_i|, |X_j|) - |X_i - X_j|)}{\mu_d \sum_{k=1}^n |X_k|}.$$

Можно показать, как и ранее, что при $n \rightarrow \infty$ статистика Q_n слабо сходится к квадратичной форме от центрированных гауссовских сл.в., $EQ = 1$. С другой стороны, в силу (21) мы имеем для любой альтернативы к гипотезе H_{02} , что: $\frac{1}{n} Q_n \xrightarrow[n \rightarrow \infty]{P} const > 0$.

Таким образом, справедлива следующая теорема.

Теорема 2. *Критерий, отвергающий гипотезу H_{02} , при $Q_n \geq \Lambda$ состоятелен против всех альтернатив, он имеет асимптотический уровень значимости не более чем α , $\forall \alpha \leq 0.21515\dots$*

§ 3. Проверка симметричности с неизвестными центром симметрии

Для неизвестного параметра $a \stackrel{\text{def}}{=} EX_1$ мы рассматриваем гипотезы:

$$H_{01}^* : X_1 - a \stackrel{D}{=} -(X_1 - a),$$

$$H_{02}^* : X_1 - a \stackrel{D}{=} C(X_1 - a) \text{ для всех ортогональных матриц } C.$$

Идея состоит в применении методов §§ 1-2 к центрированным сл.в. $X_k - \bar{X}$, $k = 1, 2, \dots, X_n$, где $\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k$ — выборочное среднее. Предположим дополнительно, что

C2) $E|X_1|^2 < \infty$, матрица $R = Cov(X_1, X_1)$ невырождена.

C3) $\int_{R^d} |t|^2 f^2(t) \varphi_0(t) dt < \infty$.

Рассмотрим подробнее условие C3). Известно равенство (см. [17], гл. 1, § 5, стр. 63)

$$\int_{S_\rho} e^{i(s, X)} ds = \left(\frac{2\pi\rho}{|X|} \right)^{d/2} |X| J_{\frac{d-2}{2}}(\rho|X|),$$

где ρ — радиус сферы S_ρ с центром в нуле в R^d , $d \geq 2$, $J_m(x)$ — функция Бесселя первого рода m -го порядка. Так что для $Z = X_1 - X_2$

$$\begin{aligned} \int_{R^d} |t|^2 f^2(t) \varphi_0(t) dt &= \int_0^\infty \left(\int_{S_y} f^2(t) dt \right) \frac{dy}{y^{d-1}} = \\ &= E \int_0^\infty \left(\frac{2\pi y}{|Z|} \right)^{d/2} |Z| J_{\frac{d-2}{2}}(y|Z|) \frac{dy}{y^{d-1}} = C_0 E|Z|^{-1}, \end{aligned}$$

где (см. [13], форм. 6.511, 6.561.14)

$$C_0 = (2\pi)^{d/2} \int_0^\infty y^{1-d/2} J_{\frac{d-2}{2}}(y) dy = \frac{2\pi^{\frac{d+1}{2}}}{\Gamma(\frac{d-1}{2})} = C(d, 1)(d-1).$$

Другими словами, условие C3) эквивалентно свойству $E|X_1 - X_2|^{-1} < \infty$, которое не является неестественным в многомерном случае: $d > 1$.

Рассмотрим сначала гипотезу о диагональной симметричности. Пусть

$$\hat{R} = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X})(X_k - \bar{X})^T$$

— выборочная ковариационная матрица и $\tilde{X}_i = \hat{R}^{-1/2}(X_i - \bar{X})$ — нормированные выборочные значения.

Обозначим через D_n тестовую статистику:

$$D_n = \frac{\sum_{i,j=1}^n (|\tilde{X}_i + \tilde{X}_j| - |\tilde{X}_i - \tilde{X}_j|)}{2 \sum_{k=1}^n |\tilde{X}_k| + 4 \frac{d-1}{n-1} \sum_{i<j} |\tilde{X}_i - \tilde{X}_j|^{-1} - \frac{2}{n-1} \sum_{i,j=1}^n |\tilde{X}_i - \tilde{X}_j|}.$$

Можно показать, что $D_n \xrightarrow[n \rightarrow \infty]{P} \infty$ при всех альтернативах удовлетворяющих C2), C3).

Справедлива следующая теорема.

Теорема 3. При условиях C2), C3) критерий, отвергающий гипотезу H_{01}^* , при $D_n \geq \Lambda$ состоятелен против всех альтернатив, он имеет асимптотический уровень значимости не более чем α , $\forall \alpha \leq 0.21515\dots$

Рассмотрим гипотезу о сферической симметричности. Обозначим

$$S_n = \frac{\sum_{i,j=1}^n [g(|X_i - \bar{X}|, |X_j - \bar{X}|) - |X_i - X_j|]}{\mu_d \sum_{i=1}^n |X_i - \bar{X}| + Tr \hat{R} \frac{2(d-1)}{d(n-1)} \sum_{i<j} |X_i - X_j|^{-1} - \frac{1}{n-1} \sum_{i,j=1}^n |X_i - X_j|},$$

где функция g была определена ранее.

Справедлива теорема.

Теорема 4. При условиях C2), C3) критерий, отвергающий гипотезу H_{02}^* , при $S_n \geq \Lambda$ состоятелен против всех альтернатив, он имеет асимптотический уровень значимости не более чем α , $\forall \alpha \leq 0.21515\dots$

§ 4. Эллиптическая симметричность

Пусть выполнено условие

$$C4) E|X_1|^4 < \infty.$$

Рассмотрим гипотезу об эллиптической симметричности распределения выборки или эквивалентно:

H_{03} : распределение $R^{-1/2}X_1$ сферически симметрично, $EX_1 = 0$.

Здесь $R = Cov(X_1, X_1)$. Обозначим $\widehat{R} = \frac{1}{n} \sum_{k=1}^n X_k X_k^T$, $Y_i = \widehat{R}^{-1/2} X_i$ и

$$g_n(t) = n^{-1} \sum_{k=1}^n \exp\{i(t, Y_k)\}, \quad g_n(t, C) = n^{-1} \sum_{k=1}^n \exp\{i(t, CY_k)\}.$$

Рассмотрим статистику

$$\begin{aligned} P(g_n, \varphi_0) &= \frac{1}{2} \int \int_{R^d} |g_n(t) - g_n(t, C)|^2 \varphi_0(t) dt d\chi(C) = \\ &= C(d, 1)n^{-2} \sum_{i,j=1}^n (g(|Y_i|, |Y_j|) - |Y_i - Y_j|), \end{aligned}$$

где функция $g(u, v)$ определена в § 2. Мы анализируем статистику $P(g_n, \varphi_0)$. Также, как и ранее, сделаем только некоторые замечания.

Во-первых, для $V_k = R^{-1/2}X_k$ и $\Delta = \widehat{R}^{-1/2}R^{1/2} - E$ по формуле Тейлора и закону больших чисел

$$\begin{aligned} \sqrt{n}(g_n(t) - g_n(t, C)) &\approx \frac{1}{\sqrt{n}} \sum_{k=1}^n [e^{i(t, V_k)} - e^{i(t, CV_k)} + i(t, \Delta V_k)e^{i(t, V_k)} - \\ &- i(t, C\Delta V_k)e^{i(t, CV_k)}] \approx \frac{1}{\sqrt{n}} \sum_{k=1}^n [e^{i(t, V_k)} - e^{i(t, CV_k)}] + \\ &+ (t, \sqrt{n}\Delta f'(t)) - (C^T t, \sqrt{n}\Delta f'(C^T t)). \end{aligned}$$

С другой стороны, для $\delta \stackrel{def}{=} \widehat{R} - R = R^{1/2}(\widehat{V} - E)R^{1/2}$, где $\widehat{V} = n^{-1} \sum_{k=1}^n V_k V_k^T$ мы имеем в силу C2) асимптотическое равенство $\widehat{R}^{1/2} = R^{1/2} + L(\delta) + o(\delta)$, $\delta \rightarrow 0$ для некоторой симметричной матрицы $L(\delta)$ линейной по δ . Возводя его в квадрат, получаем $\delta = R^{1/2}L(\delta) + L(\delta)R^{1/2}$ или

$$\widehat{V} - E = R^{-1/2}L(\delta) + L(\delta)R^{-1/2}. \quad (22)$$

При справедливости H_{03} имеем $f(t) = g(|t|^2/2)$ для некоторой функции $g : R^1 \rightarrow R^1$. Далее, так как $\Delta = -\widehat{R}^{-1/2}(\widehat{R}^{1/2} - R^{1/2}) = -R^{-1/2}L(\delta) + o(\delta)$, то мы имеем в силу (22) для любых $t \in R^d$

$$\begin{aligned} (t, \sqrt{n}\Delta f'(t)) &= (t, \sqrt{n}\Delta t)g'(|t|^2/2) = -\frac{\sqrt{n}}{2}(t, (\widehat{V} - E)t)g'(|t|^2/2) + \\ &+ o(\delta) = o(\delta) - \frac{1}{2\sqrt{n}} \sum_{k=1}^n (t, (V_k V_k^T - E)f'(t)). \end{aligned}$$

Это поможет нам проанализировать главную часть разности $g_n(t) - g_n(t, C)$. Также, как и ранее, мы можем вывести при справедливости H_{03} и условий C2), C3), C4) что

$$nP(g_n, \varphi_0) \xrightarrow[n \rightarrow \infty]{D} Q,$$

где Q есть квадратичная форма от центрированных гауссовских сл.в. и для $Z = R^{-1/2}X_1$

$$EQ = \frac{1}{2} \int \int_{R^d} |\Delta_1(t) - \frac{1}{2}\Delta_2(t)|^2 \varphi_0(t) dt d\chi(C),$$

где

$$\begin{aligned} \Delta_1(t) &= \exp\{i(t, Z)\} - \exp\{i(t, CZ)\}, \\ \Delta_2(t) &= (t, (ZZ^T - E)f'(t)) - (C^T t, (ZZ^T - E)f'(C^T t)). \end{aligned}$$

Рассуждая, как и ранее, получаем

$$E \int \int_{R^d} |\Delta_1(t)|^2 \varphi_0(t) dt d\chi(C) = 2C(d, 1)\mu_d E|Z| \quad (23)$$

и

$$\begin{aligned} (t, f'(t)) &= |t|^2 g'(|t|^2/2), \quad \Delta_2(t) = ((t, Z)^2 - (t, CZ)^2) g'(|t|^2/2), \\ E(\Delta_2(t))^2 &= (g')^2 (2E(t, Z)^4 - 2E(t, Z)^2 (t, CZ)^2). \end{aligned}$$

Рассмотрим слагаемые в последнем выражении. В силу сферической симметричности ф.р. сл.в. $Z = (Z_1, Z_2, \dots, Z_d)$

$$E(t, Z)^4 = |t|^4 E(|t|/t, Z)^4 = |t|^4 E Z_1^4 = |t|^4 E|Z|^4 \frac{3}{d(d+2)},$$

здесь последнее равенство выполнено для всех сферически симметричных распределений.

С другой стороны,

$$\int (t, CZ)^2 d\chi(C) = |t|^2 |Z|^2 E\xi^2 = \frac{1}{d} |t|^2 |Z|^2, \quad (24)$$

где ξ есть первая координата случайного вектора, равномерно распределенного на единичной сфере в R^d . Итак,

$$\begin{aligned} E \int (\Delta_2(t))^2 d\chi(C) &= 2(g')^2 (|t|^4 E Z_1^4 - \frac{1}{d} |t|^2 E|Z|^2 (t, Z)^2) = \\ &= 2(t, f'(t))^2 \varrho_d E|Z|^4, \end{aligned} \quad (25)$$

где $\varrho_d = 2(d-1)d^{-2}(d+2)^{-1}$. Пусть Y есть независимая копия сл.в. Z , обозначим

$$J = \int_{R^d} (t, f'(t))^2 \varphi_0(t) dt, \quad v(x, y) = |x - y|^{-3} (|x|^2 |y|^2 - (x, y)^2)$$

и предположим выполненным условие

$$C5) \quad E|X_1|^2 |X_1 - X_2|^{-1} < \infty.$$

Теперь нам потребуется следующая лемма.

Лемма 2. $J = C(d, 1) E v(Z, Y)$.

Доказательство. Нетрудно доказать следующую цепочку равенств:

$$\begin{aligned} J &= E \int_{R^d} (t, Z)(t, Y) \exp\{i(t, Z - Y)\} \varphi_0(t) dt = \\ &= -E \int_{R^d} \frac{\partial^2}{\partial \alpha \partial \beta} \Big|_{\alpha=\beta=1} (1 - \exp\{i(t, \alpha Z - \beta Y)\}) \varphi_0(t) dt = \\ &= -E \frac{\partial^2}{\partial \alpha \partial \beta} \Big|_{\alpha=\beta=1} C(d, 1) |\alpha Z - \beta Y| = C(d, 1) E v(Z, Y), \end{aligned}$$

здесь мы использовали формулу (3). Все интегралы понимаются в смысле их главного значения (на бесконечности), $E v(Z, Y) < \infty$ в силу C5).

Лемма доказана.

Следствие. В силу леммы 4 и (25)

$$E \int \int_{R^d} |\Delta_2(t)|^2 \varphi_0(t) dt d\chi(C) = 2C(d, 1) \varrho_d E v(Z, Y) E |Z|^4. \quad (26)$$

Можно также показать, что

$$\begin{aligned} & Re E \int \int_{R^d} \Delta_1(t) \overline{\Delta_2(t)} \varphi_0(t) dt d\chi(C) = \\ &= E \int \int_{R^d} [\exp\{i(t, Z)\} - \exp\{i(t, CZ)\}] ((t, Z)^2 - (t, CZ)^2) \\ & \quad g'(|t|^2/2) \varphi_0(t) dt d\chi(C) = \\ &= 2E \int_{R^d} \exp\{i(t, Z)\} ((t, Z)^2 - \frac{1}{d}|t|^2|Z|^2) g'(|t|^2/2) \varphi_0(t) dt, \end{aligned} \quad (27)$$

здесь мы применили формулу (24). Действуя, как ранее, мы можем преобразовать (27):

$$\begin{aligned} & E \int_{R^d} \exp\{i(t, Z)\} (t, Z)^2 g'(|t|^2/2) \varphi_0(t) dt = \\ &= E \int_{R^d} \exp\{i(t, Z)\} (t, Z) (Z, f'(t)) \varphi_0(t) dt = \\ &= E(Z, Y) \int_{R^d} \frac{\partial}{\partial \alpha} \Big|_{\alpha=1} (1 - \exp\{i(t, \alpha Z - Y)\}) \varphi_0(t) dt = C(d, 1) v_1(Z, Y), \end{aligned} \quad (28)$$

где $v_1(x, y) = (x, y)(|x|^2 - (x, y))|x - y|^{-1}$. Рассуждая аналогично, получаем

$$\begin{aligned} & E \int_{R^d} \exp\{i(t, Z)\} |t|^2 |Z|^2 g'(|t|^2/2) \varphi_0(t) dt = \\ &= E |Z|^2 \int_{R^d} \exp\{i(t, Z)\} (t, f'(t)) \varphi_0(t) dt = C(d, 1) v_2(Z, Y), \end{aligned} \quad (29)$$

где $v_2(x, y) = |x|^2((x, y) - |y|^2)|x - y|^{-1}$.

Объединяя (23), (26)–(29) получаем, что при гипотезе H_{03} для $Y_i = \widehat{R}^{-1/2} X_i$

$$EQ = C(d, 1)P - \lim_{n \rightarrow \infty} R_n,$$

$$\begin{aligned} R_n = R_n(Y_1, Y_2, \dots, Y_n) &= \frac{\mu d}{n} \sum_{i=1}^n |Y_i| + \frac{1}{2n(n-1)} \sum_{i < j} v(Y_i, Y_j) \frac{\varrho_d}{n} \sum_{i=1}^n |Y_i|^4 - \\ & - \frac{2}{n(n-1)} \sum_{i < j} (v_1(Y_i, Y_j) - \frac{1}{d} v_2(Y_i, Y_j)). \end{aligned}$$

Окончательно определим тестовую статистику формулой $G_n = nP(g_n, \varphi_0)/R_n$.

Итак, верна следующая теорема.

Теорема 5. Пусть выполнены условия C2), C4) и C5). Критерий, отвергающий гипотезу H_{03} , при $G_n \geq \Lambda$ состоятелен против всех альтернатив, он имеет асимптотический уровень значимости не более чем $\alpha, \forall \alpha \leq 0.21515\dots$

Рассмотрим гипотезу:

H_{03}^* : распределение $R^{-1/2}(X_1 - \theta)$ сферически симметрично при некотором неизвестном $\theta \in R^d$.

Обозначим $\widehat{R} = \frac{1}{n} \sum_{k=1}^n (X_k - \overline{X})(X_k - \overline{X})^T$, $Z_i = \widehat{R}^{-1/2}(X_i - \overline{X})$,

$$G_n^* = \frac{\frac{1}{n} \sum_{i,j=1}^n (g(|Z_i|, |Z_j|) - |Z_i - Z_j|)}{R_n^*},$$

$$R_n^* = R_n(Z_1, Z_2, \dots, Z_n) + \frac{2(d-1)}{n(n-1)} \sum_{i < j} |Z_i - Z_j|^{-1} + \frac{4}{n(n-1)} \sum_{i < j} v_3(Z_i, Z_j),$$

где $v_3(x, y) = ((x, y) - |y|^2)|x - y|^{-1}$.

Теорема 6. Пусть выполнены условия $C2), C4)$ и $C5)$. Критерий, отвергающий гипотезу H_{03}^* , при $G_n^* \geq \Lambda$ состоятелен против всех альтернатив, он имеет асимптотический уровень значимости не более чем $\alpha, \forall \alpha \leq 0.21515\dots$

3. ПРОВЕРКА ОДНОРОДНОСТИ

Пусть X_1, X_2, \dots, X_n и $Y_1, Y_2, \dots, Y_m, X_i, Y_i \in R^d$ — две независимые повторные выборки $X_1, Y_1 \neq 0, a.s.$ Проверка гипотезы H_0 об однородности (совпадении) распределений X -ов и Y -ов в случае $d = 1$ традиционно производится на основе критериев типа Колмогорова-Смирнова, w^2 , хи-квадрат, Манна-Уитни, Вилкоксона и других. В многомерном случае возникают трудности с определением асимптотических уровней значимости их аналогов, эти аналоги не инвариантны к линейным преобразованиям исходных данных. В настоящем параграфе работы предлагаются алгоритмы проверки однородности распределений двух повторных выборок в различных постановках и изучаются их свойства.

Рассмотрим следующий коэффициент:

$$НОМ = H_{m,n} = 1 - \frac{S_1 + S_2}{2S_3}, \quad (30)$$

где

$$S_1 = \frac{1}{n^2} \sum_{i,j=1}^n |X_i - X_j|, S_2 = \frac{1}{m^2} \sum_{i,j=1}^m |Y_i - Y_j|, S_3 = \frac{1}{mn} \sum_{i=1}^n \sum_{j=1}^m |X_i - Y_j|,$$

здесь $|\cdot|$ — евклидова норма в R^d и в случае $S_3 = 0$ (что влечет $S_1 = S_2 = 0$) мы полагаем по определению, что $H_{m,n} = 0$. Коэффициент (30) может быть представлен через выборочные характеристические функции:

$$f_n(t) = \frac{1}{n} \sum_{k=1}^n \exp\{i(t, X_k)\}, \quad g_m(t) = \frac{1}{m} \sum_{k=1}^m \exp\{i(t, Y_k)\},$$

(где (\cdot) — скалярное произведение в R^d), а именно:

$$H_{m,n} = \frac{\int_{R^d} |f_n(t) - g_m(t)|^2 \frac{dt}{|t|^{1+d}}}{\int_{R^d} (1 - |f_n(t)|^2 + 1 - |g_m(t)|^2 + |f_n(t) - g_m(t)|^2) \frac{dt}{|t|^{1+d}}}.$$

Справедлива следующая теорема.

Теорема 7. 1) Пусть $E(|X_1| + |Y_1|) \ln^\delta(1 + |X_1| + |Y_1|) < \infty$, для некоторого $\delta > 1$, тогда почти, наверное

$$H \stackrel{def}{=} \lim_{m,n \rightarrow \infty} H_{m,n} = 1 - \frac{E|X_1 - X_2| + E(|Y_1 - Y_2|)}{2E|X_1 - Y_1|};$$

2) $0 \leq H \leq 1, H = 1 \iff X_1 \equiv C_0 = const, Y_1 \equiv C_1 = const, C_0 \neq C_1,$

$H = 0 \iff$ выборки X -ов и Y -ов однородны;

3) $0 \leq H_{m,n} \leq 1, H_{m,n} = 1 \iff X_i \equiv C_0 = const, Y_j \equiv C_1 = const, C_0 \neq C_1, \forall i, j$, если $m = n$, то тогда $H_{n,n} = 0 \iff$ выборки X -ов и Y -ов совпадают (без учета их порядка), если $m \neq n$, то тогда $H_{m,n} = 0 \iff X_i = Y_j, \forall i, j$;

4) В условиях пункта 1) при справедливости нулевой гипотезы H_0 об однородности X -ов и Y -ов имеет место слабая сходимость распределений:

$$D - \lim_{m,n \rightarrow \infty} \frac{2mn}{m+n} H_{m,n} = Q, \quad EQ = 1,$$

где Q — неотрицательная квадратичная форма от центрированных гауссовских случайных величин.

Доказательство. Пункт 1) следует из закона больших чисел для одновыборочных и двух-выборочных статистик Мизеса [16, стр. 69, 73]. Второй и третий пункты можно легко доказать, используя интегральное представление для $H_{m,n}$ через выборочные характеристические функции приведенное выше, и, пункт 4) доказывается так же как это делалось выше в §1.

Замечание 1. Для квадратичных форм Q , упомянутых в пункте 4), равномерно по X_1, Y_1 имеет место неравенство: $P\{Q \geq \Lambda\} \leq \alpha \forall \alpha \leq 0,21515\dots$, где $\Lambda = (\Phi^{-1}(1 - \frac{\alpha}{2}))^2$, $\Phi(\cdot)$ — ф.р. нормальной (0,1) сл.в., $\Phi^{-1}(\cdot)$ — обратная функция. Таким образом, критерий, отвергающий гипотезу H_0 об однородности X -ов и Y -ов при

$$\frac{2mn}{m+n} H_{m,n} \geq \Lambda, \quad (31)$$

состоятелен против всех альтернатив, удовлетворяющих условиям пункта 1), он имеет асимптотический уровень значимости, равный $\alpha, \forall \alpha \leq 0,21515\dots$

Замечание 2. Задачу проверки гипотезы об однородности распределений одномерных случайных величин X_k, Y_k , принимающих конечное число значений, можно свести к проверке однородности распределений выборок:

$$X'_k = \begin{pmatrix} \chi(X_k = a_1) \\ \chi(X_k = a_2) \\ \dots \\ \chi(X_k = a_d) \end{pmatrix}, Y'_k = \begin{pmatrix} \chi(Y_k = a_1) \\ \chi(Y_k = a_2) \\ \dots \\ \chi(Y_k = a_d) \end{pmatrix} \quad k = 1, 2, \dots, n,$$

здесь $\chi(\cdot)$ — индикатор соответствующего множества и $\{a_i\}_{i=1}^d$ — вероятные значения случайных величин X_k, Y_k . Рассуждая, как и ранее, можно получить формулу

$$H_{m,n} = \frac{\sum_{k=1}^d (p_k - q_k)^2}{2 - 2 \sum_{k=1}^d p_k q_k},$$

где $p_k, (q_k), k \geq 1$ — выборочные частоты значений a_k для выборки X -ов (Y -es). Действительно, пусть $t = (t_1, t_2, \dots, t_d)$, тогда

$$f_n(t) = \frac{1}{n} \sum_{j=1}^n e^{i(t, X'_j)} = \sum_{k=1}^d p_k e^{it_k}, \quad g_m(t) = \frac{1}{m} \sum_{j=1}^m e^{i(t, Y'_j)} = \sum_{k=1}^d q_k e^{it_k},$$

$$|f_n(t) - g_m(t)|^2 = \sum_{k=1}^d (p_k - q_k)^2 + \sum_{k \neq l} (p_k - q_k)(p_l - q_l) e^{i(t_k - t_l)} =$$

$$= \sum_{k \neq l} (p_k - q_k)(p_l - q_l) (e^{i(t_k - t_l)} - 1),$$

следовательно, числитель $H_{m,n}$ равен

$$\begin{aligned} 2S_3 - S_1 - S_2 &= \frac{1}{C(d, 1)} \int_{R^d} \frac{|f_n(t) - g_m(t)|^2}{|t|^{d+1}} dt = -r \sum_{k \neq l} (p_k - q_k)(p_l - q_l) = \\ &= r \sum_{k=1}^d (p_k - q_k)^2 \end{aligned}$$

для некоторой константы r и, с другой стороны,

$$\begin{aligned} 1 - f_n(t)\overline{g_m(t)} &= 1 - \sum_{k,l=1}^d p_k q_l e^{i(t_k - t_l)} = 1 - \sum_{k=1}^d p_k q_k - \sum_{k \neq l}^d p_k q_l e^{i(t_k - t_l)} = \\ &= \sum_{k \neq l}^d p_k q_l (1 - e^{i(t_k - t_l)}), \end{aligned}$$

следовательно знаменатель $H_{m,n}$ равен

$$2S_3 = 2Re \frac{1}{C(d,1)} \int_{R^d} \frac{1 - f_n(t)\overline{g_m(t)}}{|t|^{d+1}} dt = 2r \sum_{k \neq l}^d p_k q_l = 2r \left(1 - \sum_{k=1}^d p_k q_k \right).$$

Мы видим, что вклад относительно малых p_k, q_k в $H_{m,n}$ будет относительно мал. В этом — одно из их отличий $\frac{2mn}{m+n} H_{m,n}$ от хорошо известного варианта статистики хи-квадрат критерия

$$\frac{2mn}{m+n} \sum_{k=1}^d \frac{(p_k - q_k)^2}{p_k + q_k}$$

для проверки однородности ([20], стр. 88).

Отметим также, что статистика

$$H_{m,n} = \frac{\sum_{k=1}^{\infty} (p_k - q_k)^2}{2 - 2 \sum_{k=1}^{\infty} p_k q_k}$$

применима для проверки однородности, когда случайные величины X_k, Y_k изменяются на счетном множестве значений.

§ 5. Асимптотика функции мощности

Оценим теперь мощность β критерия (31) при увеличении объемов выборок. Пусть $m = n, d = 1$ и выполнено условие Г. Крамера: $\exists L > 0 : E \exp\{L(|X_1| + |Y_1|)\} < \infty$. Имеем

$$\begin{aligned} 1 - \beta &= P\{n(2S_3 - S_1 - S_2) \leq 2S_3 \Lambda\} \leq P\{(n - \Lambda)(2S_3 - S_1 - S_2) \leq \\ &\leq 2\Lambda \log n\} + P\{S_1 \geq \log n\} + P\{S_2 \geq \log n\}. \end{aligned} \quad (32)$$

Оценим два последних слагаемых. Ясно, что $S_1 \leq \frac{2}{n} \sum_{i=1}^n |X_i|$ (аналогичное соотношение выполнено для S_2), следовательно, в силу экспоненциального неравенства для сумм независимых случайных величин, удовлетворяющих условию Г. Крамера ([18], стр. 81): $\exists n_0 : \forall n \geq n_0$

$$P\{S_1 \geq \log n\} + P\{S_2 \geq \log n\} \leq \exp\{-n\sqrt{\log n}\}.$$

Обозначим единым символом C все положительные константы, зависящие только от размерности данных d и моментов $E|X_1|, E|Y_1|$. Далее, имеем

$$C(2S_3 - S_1 - S_2) = \int_{R^d} |f_n(t) - g_n(t)|^2 \frac{dt}{|t|^{1+d}} \geq \Delta - S_n, \quad (33)$$

где

$$\begin{aligned} \Delta &= \int_{R^d} |f(t) - g(t)|^2 \frac{dt}{|t|^{1+d}}, \quad S_n = \frac{1}{n} \sum_{k=1}^n \xi_k, \quad E\xi_k = 0, \\ \xi_k &= 2Re \int_{R^d} (\overline{g}(t) - \overline{f}(t))(e^{i(t, X_k)} - f(t) - e^{i(t, Y_k)} + g(t)) \frac{dt}{|t|^{1+d}} \end{aligned}$$

(черта сверху означает комплексное сопряжение). Итак, для некоторой неслучайной последовательности $\alpha_n \xrightarrow{n \rightarrow \infty} 0$

$$1 - \beta \leq \exp\{-n\sqrt{\log n}\} + P\{S_n \geq \Delta(1 + \alpha_n)\}.$$

Легко видеть, что

$$\int_{R^d} (1 - \exp\{i(t, X)\}) \frac{dt}{|t|^{1+d}} \equiv C|X|.$$

Интеграл понимается в смысле главного значения

$$\xi_k^2/\Delta + |\xi_k| \leq C(|X_k| + E|X_k| + |Y_k| + E|Y_k|),$$

поэтому, в частности, для случайных величин ξ_k условие Г. Крамера выполнено, следовательно ([19], стр. 208)

$$\nu = \limsup_{n \rightarrow \infty} \frac{1}{n} \ln(1 - \beta) \leq -h(\Delta) = -\sup_x (\Delta x - \ln Ee^{x\xi_1}).$$

В соответствии с общей теорией больших уклонений $h(\Delta) < \infty$ для достаточно малых Δ функция $h(\cdot)$ неотрицательна, выпукла и $h(0) = 0$ ([19], стр. 204).

Пусть теперь $\Delta \rightarrow 0$ и $\tau^* = \sup\{\tau | E \exp\{\tau|X_1| + \tau|Y_1|\}\} > 0$ равномерно для достаточно малых Δ , тогда $h(\Delta) = (1 + o(1))\Delta^2/(2\sigma^2)$ (см. [19], стр. 204, 208), здесь $\sigma^2 \stackrel{def}{=} E\xi_1^2$.

Рассмотрим далее простые сдвиговые альтернативы $H_1 : g(t) = e^{i\theta t} f(t), \forall t$. Предположим, что существует плотность распределения $p(x)$, (соответствующая характеристической функции $f(t)$), принадлежащая $L_2(R^1)$, положительная, абсолютно непрерывная, обладающая ненулевой, конечной фишеровской информацией I и пусть при $\theta \rightarrow 0$

$$-\ln \int_{-\infty}^{\infty} \sqrt{p(x)p(x+\theta)} dx \sim \frac{\theta^2 I}{8}.$$

Из результатов работы [20](стр. 88) следует, что

$$eff = \limsup_{\theta \rightarrow 0} \left(-\frac{8\nu}{\theta^2 I} \right) \leq 1,$$

здесь eff — локальная, относительная асимптотическая эффективность критерия (31) по Ходжесу-Леману. Далее, при $\theta \rightarrow 0$

$$(1 + o(1))\Delta^2 = (\theta^2 \int_{-\infty}^{\infty} |f(t)|^2 dt)^2 = (2\pi\theta^2 \int_{-\infty}^{\infty} p^2(x) dx)^2,$$

$$(1 + o(1))\sigma^2 = 8\theta^2 E(\text{Im} \int_{R^d} t^{-1} \bar{f}(t)(e^{itX_1} - f(t)) dt)^2 = \frac{8}{3}\pi^2\theta^2,$$

следовательно, в рассматриваемом случае для критерия (31)

$$eff \geq 6 \left(\int_{-\infty}^{\infty} p^2(x) dx \right)^2 \left(\int_{-\infty}^{\infty} \frac{(p'(x))^2}{p(x)} dx \right)^{-1},$$

поэтому, например, для гауссовских распределений $eff \geq \frac{3}{2\pi} = 0,47\dots$, для распределения Лапласа: $p(x) = \lambda \exp\{-2\lambda|x|\}$, $eff \geq 0,375$ и для логистического: $p(x) = e^x(1 + e^x)^{-2}$, $eff \geq 0,5$. Отметим здесь, что $eff = 1$ для критериев Колмогорова-Смирнова и w^2 , более того, они асимптотически оптимальны по Ходжесу-Леману (см. [20]).

§ 6. Проверка внутренней однородности выборки

Рассмотрим задачу проверки гипотезы H_{01} о внутренней однородности как таковой выборки из независимых случайных векторов $X_0, X_1, X_2, \dots, X_n$, $E|X_j| < \infty$, при H_{01} распределения X_k -х совпадают. В параметрической постановке она рассматривалась, например, А.Д. Бернштейном (против альтернатив, сближающихся с нулевой гипотезой при $n \rightarrow \infty$). Обозначим $f_k(t)$ характеристические функции случайных векторов X_k . Ясно, что

$$H_{01} \iff f_k(t) = f_l(t), \forall k, l, t \iff$$

$$U(f) \stackrel{\text{def}}{=} \frac{\int_0^1 x(1-x) \int_{R^d} |f(t, x) - g(t, x)|^2, \frac{dt}{|t|^{1+d}} dx}{\int_0^1 \int_{R^d} (1 - \operatorname{Re} f(t, x) \overline{g(t, x)}, \frac{dt}{|t|^{1+d}} dx} = 0,$$

где

$$f_n(t, x) = \frac{1}{[nx] + 1} \sum_{k=0}^{[nx]} f_k(t) \quad g(t, x) = \frac{1}{n - [nx]} \sum_{k=[nx]+1}^n f_k(t),$$

здесь $[\cdot]$ — целая часть числа, таким образом, $U(f) \neq 0$ для всех альтернатив для фиксированного n .

Проверка гипотезы H_{01} может быть произведена на основе статистики $H_n = U(f_n)$, полученной подстановкой в $U(f)$ вместо $f(t, x)$ и $g(t, x)$ соответственно их выборочных значений:

$$f_n(t, x) = \frac{1}{[nx] + 1} \sum_{k=0}^{[nx]} \exp\{i(t, X_k)\}, \quad g_n(t, x) = \frac{1}{n - [nx]} \sum_{k=[nx]+1}^n \exp\{i(t, X_k)\}$$

и заменой в $U(f_n)$ интегралов на соответствующие интегральные суммы:

$$H_n = \frac{\frac{1}{n} \sum_{k=1}^n x_k(1-x_k)(2S_3(k) - S_1(k) - S_2(k))}{\frac{1}{n} \sum_{k=1}^n S_3(k)},$$

где

$$S_1(k) = \frac{1}{k^2} \sum_{i,j=0}^{k-1} |X_i - X_j|, \quad S_2(k) = \frac{1}{(n-k+1)^2} \sum_{i,j=k}^n |X_i - X_j|,$$

$$S_3(k) = \frac{1}{k(n-k+1)} \sum_{i=0}^{k-1} \sum_{j=k}^n |X_i - X_j|, \quad x_k = \frac{k}{n}.$$

Ясно, что $0 \leq H_n \leq 0,5$ с вероятностью 1 и что коэффициент H_n инвариантен к изометрическим преобразованиям и изменению масштаба данных, так что он может рассматриваться как мера внутренней однородности выборки. Можно также заметить, что $H_n = 0 \iff X_1 = X_2 = \dots = X_n$, п.в. При справедливости нулевой гипотезы H_{01} можно показать (также, как и в [18], § 1, § 2), что $D - \lim_{n \rightarrow \infty} nH_n = Q^*$, $EQ^* = 1$, где Q^* — неотрицательная квадратичная форма от центрированных гауссовских случайных величин, так что критерий, отвергающий H_{01} при

$$nH_n \geq \Lambda, \tag{34}$$

имеет асимптотический уровень значимости не более чем α , $\forall \alpha \leq 0,21515\dots$

Обозначим теперь $U_0(f) = \int_0^1 x(1-x) \int_{R^d} |f(t, x) - g(t, x)|^2, \frac{dt}{|t|^{1+d}} dx$. Для последовательности альтернатив:

$$H_{1n} : nU_0(f) \xrightarrow{n \rightarrow \infty} \infty, \sup_{i,j,n} E|X_i - X_j| < \infty$$

имеем $E \frac{1}{n} \sum_{k=1}^n S_3(k) \leq C$, поэтому аналогично (32), (33) мы можем получить: $\forall K > 0$

$$1 - \beta \leq P\left\{n \int_0^1 x(1-x) \int_{R^d} |f_n(t, x) - g_n(t, x)|^2 \frac{dt}{|t|^{1+d}} dx \leq C\Lambda K\right\} + \\ + \frac{C}{K} \leq P\{nS_n^* \geq nU_0(f) - C\Lambda K\} + \frac{C}{K},$$

где β — мощность критерия (34),

$$S_n^* = 2Re \int_0^1 x(1-x) \int_{R^d} \xi_n(t)(\bar{g}(t) - \bar{f}(t)) \frac{dt}{|t|^{1+d}} dx \\ \xi_n(t) = f_n(t, x) - f(t, x) - g_n(t, x) + g(t, x).$$

По неравенству Коши-Буняковского

$$E(S_n^*)^2 \leq 4U_0(f) \int_0^1 x(1-x) \int_{R^d} E|\xi_n(t)|^2 \frac{dt}{|t|^{1+d}} dx \leq CU_0(f)/n,$$

поэтому $\beta \xrightarrow{n \rightarrow \infty} 1$ и, значит, критерий (34) состоятелен против последовательности альтернатив H_{1n} .

Пример. Рассмотрим проверку нулевой гипотезы $H_0 : X_k = \xi_k$, $k \geq 1$ где ξ_k , $k \geq 1$ повторная выборка против альтернативы $H_1 : X_k = \xi_k + F(\frac{k}{n})$, $k \geq 1$, где $F(x)$ — функция, локально интегрируемая по Риману, и $F(x)$ не есть константа почти всюду по отношению к мере Лебега. В этом случае

$$U_0(f) \xrightarrow{n \rightarrow \infty} \\ \xrightarrow{n \rightarrow \infty} \int_0^1 x(1-x) \int_{R^d} \left| \frac{1}{x} \int_0^x e^{i(t, F(s))} ds - \frac{1}{1-x} \int_x^1 e^{i(t, F(s))} ds \right|^2 \frac{|f(t)|^2 dt}{|t|^{1+d}} dx,$$

что не равно нулю, так как иначе мы имели бы для всех достаточно малых $|t|$, что интегранд равен нулю или после элементарных преобразований

$$\int_0^x e^{i(t, F(s))} ds = x \int_0^1 e^{i(t, F(s))} ds$$

почти всюду по отношению к мере Лебега. Производная по x интеграла в левой части здесь существует и равна интегранду или $e^{i(t, F(x))} = C(t)$, что невозможно, так как $F(x)$ не есть константа.

Таким образом, рассматриваемый критерий состоятелен против всех сдвиговых альтернатив с ф.р. F . К примеру, мы можем тестировать линейный тренд данных.

§ 7. Проверка некоторых линейных гипотез

1. Однородность к сдвигам. Рассмотрим гипотезу $H_0 : F_Y(x) = F_X(x - \theta)$, $\forall x \in R^d$, для некоторого $\theta \in R^d$ другими словами, распределения X -ов и Y -ов принадлежат одному сдвиговому семейству распределений или эквивалентны $X - EX \stackrel{D}{=} Y - EY$. Пусть имеются две независимые повторные выборки X_1, X_2, \dots, X_n и Y_1, Y_2, \dots, Y_n равных объемов. В соответствии с используемым нами методом рассмотрим эмпирические характеристические функции

$$f_n(t) = \frac{1}{n} \sum_{k=1}^n \exp\{i(t, X_k - \bar{X})\}, \quad g_n(t) = \frac{1}{n} \sum_{k=1}^n \exp\{i(t, Y_k - \bar{Y})\}.$$

Пусть выполнены условия:

С) $E|X_1|^2 + E|Y_1|^2 < \infty$, матрицы $Cov(X_1, X_1)$ и $Cov(Y_1, Y_1)$ невырождены, обозначим $Z_1 = X_1 - EX_1 - Y_1 + EY_1$, в случае $d \neq 1$ $E|R^{-1/2}(Z_1)|^{-1} < \infty$ и если $d = 1$, то тогда сл.в. Z_1 имеет непрерывную плотность распределения $p(x)$ такую, что $p(0) \neq 0$.

Обозначим $f(t)$ х.ф. сл.в. $X_1 - EX_1$ и рассмотрим случайный процесс

$$\xi_n(t) = \sqrt{n}(f_n(t) - g_n(t)).$$

Нетрудно подсчитать, что при справедливости H_0

$$\begin{aligned} E|\xi_n(t)|^2 &= 2(1 - |f(t)|^2) + 2n \left(|f(t)|^2 - \left| f\left(t - \frac{t}{n}\right) f^{n-1}\left(-\frac{t}{n}\right) \right|^2 \right) \xrightarrow{n \rightarrow \infty} \\ &\xrightarrow{n \rightarrow \infty} 2(1 - |f(t)|^2) + 2Re[(t, f'(t))f(-t) - f(t)(t, f'(-t))] + |f(t)|^2(t, Rt) = V(t), \end{aligned}$$

где $R = Cov(X_1, X_1)$. Обозначим $X_1^* = R^{-1/2}(X_1 - EX_1)$, $Y_1^* = R^{-1/2}(Y_1 - EY_1)$, рассуждая далее, как и в § 1 (проверка симметричности), мы можем подсчитать $\forall d \neq 1$

$$\begin{aligned} &\frac{1}{C(d, 1)} \int_{R^d} \frac{V(t)}{(t, Rt)^{\frac{d+1}{2}}} dt = 2detR^{-1/2} E|X_1^* - Y_1^*| + \\ &+ 2E \frac{1}{C(d, 1)} \int_{R^d} \frac{i(t, X_1 - EX_1 - Y_1 + EY_1) e^{i(t, X_1 - EX_1 - Y_1 + EY_1)}}{(t, Rt)^{\frac{d+1}{2}}} dt + \\ &+ \frac{2detR^{-1/2}}{C(d, 1)} \int_{R^d} \frac{|t|^2 |f(R^{-1/2}t)|^2}{|t|^{d+1}} dt = 2detR^{-1/2} E|X_1^* - Y_1^*| \left(1 - \frac{C^*(d, 1)}{C(d, 1)} \right) + \\ &+ 2 \frac{C_0(d, 1) detR^{-1/2}}{C(d, 1)} E|X_1^* - Y_1^*|^{-1} = 2(d-1) detR^{-1/2} E|X_1^* - Y_1^*|^{-1} = \\ &= P - \lim_{n \rightarrow \infty} U_n, \end{aligned}$$

где $C_0(d, 1) = (d-1)C(d, 1)$ (см. §1 (проверка симметричности)) и

$$C^*(d, 1) = \int_{R^d} \frac{s_1 \sin s_1}{|s|^{d+1}} ds = C(d, 1),$$

что может быть показано также, как и при выводе формулы для $C(d, 1)$:

$$\begin{aligned} U_n &= 2(d-1) \frac{detR^{-1/2}}{n^2} \sum_{i,j=1}^n |X_i^* - Y_j^*|^{-1}, \quad X_i^* = \widehat{R}^{-1/2}(X_i - \bar{X}), \\ Y_j^* &= \widehat{R}^{-1/2}(Y_j - \bar{Y}), \quad \widehat{R} = \frac{1}{2n} \sum_{i=1}^n \{ (X_i - \bar{X})(X_i - \bar{X})^T + (Y_i - \bar{Y})(Y_i - \bar{Y})^T \}. \end{aligned}$$

В случае $d = 1$ мы имеем

$$\frac{1}{C(d, 1)} \int_{R^d} \frac{V(t)}{(t, Rt)^{\frac{d+1}{2}}} dt = \frac{2detR^{-1/2}}{C(d, 1)} \int_{R^d} |f(R^{-1/2}t)|^2 dt = \frac{4\pi detR^{-1/2}}{C(d, 1)} p(0),$$

где $p(x)$ — плотность распределения $X_1^* - Y_1^*$. Рассмотрим следующую тестовую статистику:

$$\begin{aligned} Q_n &= \frac{n}{U_n} \times \frac{1}{C(d, 1)} \int_{R^d} \frac{|f_n(t) - g_n(t)|^2}{(t, \widehat{R}t)^{\frac{d+1}{2}}} dt = \\ &= n \times \frac{\frac{2}{n^2} \sum_{i=1}^n \sum_{k=1}^n |X_i^* - Y_k^*| - \frac{1}{n^2} \sum_{i,j=1}^n |X_i^* - X_j^*| - \frac{1}{n^2} \sum_{k,l=1}^n |Y_k^* - Y_l^*|}{\frac{2(d-1)}{n^2} \sum_{i=1}^n \sum_{k=1}^n |X_i^* - Y_k^*|^{-1}}, \end{aligned}$$

где в случае $d = 1$ мы заменяем знаменатель на $4\pi \widehat{p}(0)$ для некоторой состоятельной оценки плотности $p(x)$ в нулевой точке. Рассуждая, как и ранее, можно показать, что $Q_n \xrightarrow{n \rightarrow \infty} Q$ слабо, где Q — квадратичная форма от центрированных гауссовских случайных величин,

$EQ = 1$. Соответственно, асимптотический уровень значимости будет иметь заданное значение $\alpha, \forall \alpha \leq 0.21515 \dots$ и тест будет состоятельным против всех альтернатив, подчиненных условиям С).

Если условие С) не выполнено, то тогда можно рассмотреть альтернативный вариант для проверки H_0 , используя статистику:

$$Q_{m,n} = \inf_{A \in R^d} Q_{m,n}(A), \quad Q_{m,n}(A) = \frac{2mn}{m+n} \times$$

$$\frac{\frac{2}{mn} \sum_{i=1}^n \sum_{k=1}^m |X_i - A - Y_k| - \frac{1}{n^2} \sum_{i,j=1}^n |X_i - X_j| - \frac{1}{m^2} \sum_{k,l=1}^m |Y_k - Y_l|}{\frac{1}{n^2} \sum_{i,j=1}^n |X_i - X_j| + \frac{1}{m^2} \sum_{k,l=1}^m |Y_k - Y_l|}.$$

Можно заметить, что минимум $Q_{m,n}(A)$ по A легко найти, перебирая значения $A = X_i - Y_j$ и $Q_{m,n}(A) \leq Q_{m,n}(\theta) \xrightarrow[m,n \rightarrow \infty]{D} Q, EQ = 1$, где Q есть неотрицательная квадратичная форма от центрированных гауссовских случайных величин. Таким образом, мы получаем верхнюю оценку для асимптотического уровня значимости:

$$\lim_{m,n \rightarrow \infty} P\{Q_{m,n} \geq (\Phi^{-1}(1 - \frac{\alpha}{2}))^2\} \leq \lim_{m,n \rightarrow \infty} P\{Q_{m,n}(\theta) \geq (\Phi^{-1}(1 - \frac{\alpha}{2}))^2\} \leq$$

$$\leq P\{Q \geq (\Phi^{-1}(1 - \frac{\alpha}{2}))^2\} = \alpha \forall \alpha \leq 0,21515 \dots$$

Соответствующий критерий будет состоятельным против всех альтернатив, подчиненных условию $E(|X_1| + |Y_1|) < \infty$. Действительно, рассмотрим

$$R_{m,n}(A) \stackrel{def}{=} C(d, 1) \times$$

$$\times \left(\frac{2}{mn} \sum_{i=1}^n \sum_{k=1}^m |X_i - A - Y_k| - \frac{1}{n^2} \sum_{i,j=1}^n |X_i - X_j| - \frac{1}{m^2} \sum_{k,l=1}^m |Y_k - Y_l| \right),$$

которая выпукла и ограничена: $R_{m,n}(A) \leq 2|A| + \frac{4}{n} \left(\sum_{i=1}^n |X_i| + \sum_{j=1}^m |Y_j| \right)$. Применяя закон больших чисел, мы получаем

$$\lim_{m,n \rightarrow \infty} \inf_{A \in R^d} R_{m,n}(A) = \inf_{A \in R^d} C(d, 1) (2|X_1 - A - Y_1| - E|X_1 - X_2| -$$

$$- E|Y_1 - Y_2|) = \inf_{A \in R^d} \int_{R^d} \frac{|e^{-i(t,A)} f(t) - g(t)|^2}{|t|^{d+1}} dt,$$

правая часть здесь неотрицательна и равна нулю только при справедливости H_0 , таким образом, для всех альтернатив $Q_{m,n} \xrightarrow[m,n \rightarrow \infty]{} \infty$, *a.s.* что означает состоятельность критерия (здесь $f(t), g(t)$ есть соответствующие характеристические функции).

2. Принадлежность сдвига-масштабному семейству распределений. Пусть X, Y есть случайные вектора со значениями в R^d . Рассмотрим гипотезу $H_0 : Y \stackrel{D}{=} A + BX$ для некоторых неслучайных вектора A и матрицы B . Рассмотрим сначала одномерный случай: пусть $E|X_1|^2 + E|Y_1|^2 < \infty$. Рассуждая, как и ранее, рассмотрим тестовую статистику

$$L_{m,n} = \inf_{A \in R^1, B > 0} L_{m,n}(A, B), \quad L_{m,n}(A, B) = \frac{2mn}{m+n} \times$$

$$\times \frac{\frac{2}{mn} \sum_{i=1}^n \sum_{k=1}^m |Y_k - A - BX_i| - \frac{1}{n^2} \sum_{i,j=1}^n B|X_i - X_j| - \frac{1}{m^2} \sum_{k,l=1}^m |Y_k - Y_l|}{\frac{1}{n^2} \sqrt{\frac{\widehat{S}_Y}{\widehat{S}_X}} \sum_{i,j=1}^n |X_i - X_j| + \frac{1}{m^2} \sum_{k,l=1}^m |Y_k - Y_l|},$$

где $\widehat{S}_Y, \widehat{S}_X$ есть соответствующие выборочные дисперсии. Заметим, что $L_{m,n}(A, B)$ есть выпуклая функция от A, B , что облегчает вычисление инфимума. Можно показать, что критерий, отвергающий H_0 в случае

$$L_{m,n} \geq (\Phi^{-1}(1 - \frac{\alpha}{2}))^2,$$

имеет асимптотический уровень значимости менее чем $\alpha \forall \alpha \leq 0,21515\dots$ и является состоятельным против всех альтернатив, подчиненных условию $E(|X_1| + |Y_1|) < \infty$.

Рассмотрим теперь случай $d \neq 1$ и специальный случай матрицы B : $B = bB'$ для некоторой положительной константы b и ортогональной матрицы B' , этот случай соответствует ситуации, когда одна из выборок подвергается сдвигам, вращениям и изменению масштаба (что соответствует различным способам регистрации данных). Обозначим $\widehat{S}_Y, \widehat{S}_X$ выборочные дисперсии:

$$\widehat{S}_Y = \frac{1}{m} \sum_{k=1}^m (Y_k - \bar{Y})(Y_k - \bar{Y})^T, \quad \widehat{S}_X = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{Y})(X_i - \bar{Y})^T$$

и

$$M_{m,n} = \inf_{A \in R^d, b > 0, B'} M_{m,n}(A, B), \quad M_{m,n}(A, B) = \frac{2mn}{m+n} \times$$

$$\times \frac{\frac{2}{mn} \sum_{i=1}^n \sum_{k=1}^m |Y_k - A - bB'X_i| - \frac{1}{n^2} \sum_{i,j=1}^n b|X_i - X_j| - \frac{1}{m^2} \sum_{k,l=1}^m |Y_k - Y_l|}{\frac{1}{n^2} \sqrt{\frac{\text{Tr} \widehat{S}_Y}{\text{Tr} \widehat{S}_X}} \sum_{i,j=1}^n |X_i - X_j| + \frac{1}{m^2} \sum_{k,l=1}^m |Y_k - Y_l|}.$$

Можно показать, что критерий, отвергающий H_0 в случае

$$M_{m,n} \geq (\Phi^{-1}(1 - \frac{\alpha}{2}))^2,$$

имеет асимптотический уровень значимости менее чем $\alpha \forall \alpha \leq 0,21515\dots$ и является состоятельным против всех альтернатив, подчиненных условию $E(|X_1| + |Y_1|) < \infty$.

СПИСОК ЛИТЕРАТУРЫ

1. Никитин Я.Ю. *Асимптотическая эффективность непараметрических критериев*. М.: Наука. 1995. 238 с.
2. M.L. Puri, P.K. Sen *On the theory of rank order tests for location in the multivariate one-sample problem* // Ann. Math. Statist. V. 38. 1967. P. 1216–1228.
3. M. Huskova *Asymptotic distribution of rank statistics used for multivariate symmetry* // J. Multiv. Analysis. V. 1. No. 1. 1971. P. 461–484.
4. J. Mottonen, T.P. Hettmansperger, H. Oja, J. Tienari *On the efficiency of affine invariant multivariate rank test* // J. Multiv. Analysis. V. 66. 1998. P. 118–132.
5. K-T. Fang, L-X. Zhu, P.M. Bentler *A necessary test of goodness of fit for symmetry* // J. Multiv. Analysis. V. 45. 1993. P. 34–55.
6. V.I. Koltchinskii, L. Li *Testing for spherical symmetry of a multivariate distributions* // J. Multiv. Analysis. V. 65. 1998. P. 218–244.
7. J.C. Lee, T.C. Chang, P.R. Krishnaiah *Approximation of the distribution s of the likelihood ratio statistics for testing certain structures of the covariance matrix of real multivariate normal populations*. // APL TR 75–167, Aerospace Research Laboratory, Wright-Patterson, Ohio.
8. S. Csorgo, C.R. Heathcote *Testing for symmetry* // Biometrika. V. 74, No. 1. 1987. P. 177–184.
9. L. Baringhaus // Ann. Statist. V. 19, No. 2. 1991. P. 899–917.

10. G. Neuhaus, L-X. Zhu *Permutation test for reflected symmetry* // J. Multiv. Analysis. V. 67. 1998. P. 129–153.
11. Bakirov N.K., Rizzo M.L. , Szekely G.J. *A multivariate nonparametric test of independence*// Journal of Multivariate Analysis. 2006. V.97. Issue 8. P.1742-1756.
12. Bakirov N.K., Rizzo M.L. , Szekely G.J. *Measuring and Testing Dependence by Correlation of Distances* // The Annals of Mathematical Statistics. 2007. V.35. No. 6. P. 2769-2794.
13. Градштейн И.С., Рыжик И.М. *Таблицы интегралов, сумм, рядов и произведений, 4-е изд.* М.: Физматгиз. 1962. 1100 с.
14. R.M. Dudley *Gaussian processes on several parameters* // Ann. Math. Statist. 1965. V. 36, No. 3. P. 771–788.
15. N.K. Bakirov, G.J. Szekely *Extremal properties for Gaussian quadratic forms*"Probability theory and related fields // Probability theory and related fields. 2003. V. 126. No. 2. P. 184–202.
16. M. Abramovitz, I. Stegun *Handbook of mathematical functions.* // National Bureau of Standards. 1964.
17. Гихман И.И., Скороход А.В. *Введение в теорию случайных процессов.* М.: Наука. 1965. 654 с.
18. Петров В.В. *Предельные теоремы для сумм независимых случайных величин.* М.: Наука. 1987. 320 с.
19. Боровков А.А. *Теория вероятностей, 2-е изд.* М.: Наука. 1986. 432 с.
20. Никитин Я.Ю. *Об асимптотической эффективности по Ходжесу-Леману непараметрических критериев согласия и однородности* // Теория вероятностей и ее применения. 1987. Т. 32, № 1. С. 82–91.

Наиль Кутлужанович Бакиров,
 Институт математики с ВЦ УНЦ РАН,
 ул. Чернышевского, 112,
 450008, г. Уфа, Россия
 E-mail: bakirovnk@rambler.ru