

Математическое моделирование

УДК 519.24

ИСПОЛЬЗОВАНИЕ СПЕЦИАЛЬНЫХ ФУНКЦИЙ ЭРМИТА ДЛЯ ИССЛЕДОВАНИЯ МОЩНОСТНЫХ СВОЙСТВ КРИТЕРИЯ ГРАББСА

Л. К. Ширяева

Самарский государственный экономический университет,
443090, Россия, Самара, ул. Советской Армии, 141.

E-mail: shiryeva_lk@mail.ru

Рассмотрен случай, когда в нормально распределенной выборке имеется выброс. Получено новое представление для плотности распределения вероятностей студентизированного отклонения выброса от среднего по выборке, основанное на использовании специальных функций Эрмита с отрицательным целым значком. На основе этого представления найдены интегральные соотношения для мер мощности Дэвида критерия Граббса в случае, когда статистика критерия является статистикой отношения правдоподобия. Найдена величина максимально возможного отклонения вероятности обнаружения присутствия одиночного выброса в выборке от вероятности его точного обнаружения. Определена область критических значений статистики Граббса, в которой меры мощности Дэвида, предназначенные для точного обнаружения выброса, совпадают. Выполнены модельные расчёты функции мощности критерия Граббса для случая нормально распределенных выборок с выбросом, отличающимся от остальных наблюдений сдвигом вправо. Результаты вычислений оказались близки к теоретически ожидаемым.

Ключевые слова: выброс, критерий Граббса проверки на один выброс, функция Эрмита, меры мощности Дэвида для критерия Граббса, нормальный закон распределения.

Введение. Пусть $X_1, X_2, \dots, X_{n-1}, X_n$ — случайная выборка из n значений нормально распределенной случайной величины X ; $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n-1)} \leq X_{(n)}$ — построенный по ней упорядоченный вариационный ряд; $X_{(j)}$ — j -тая порядковая статистика ($j = 1, 2, \dots, n$). Проверяемая нулевая гипотеза H_0 состоит в том, что наблюдения $X_1, X_2, \dots, X_{n-1}, X_n$ являются независимыми случайными величинами с нормальным распределением $N(a, \sigma^2)$. В качестве конкурирующей гипотезы H_1 рассмотрим случай, когда какие-либо $(n-1)$ из n наблюдений имеют одинаковое $N(a, \sigma^2)$ распределение, а одно из них — выброс X_{out} — имеет распределение $N(a + \lambda\sigma, \nu\sigma^2)$. Параметр сдвига $\lambda \geq 0$ характеризует среднее (ожидаемое) смещение выброса вправо, а параметр масштаба $\nu > 0$ определяет изменение дисперсии выброса относительно остальных наблюдений. Гипотеза H_1 моделирует ситуацию, когда вероятность «засорения» выборки аномальными наблюдениями весьма

Людмила Константиновна Ширяева (к.ф.-м.н., доц.), доцент, каф. математической статистики и эконометрики.

мала [1], а сам выброс отличается от «обычных» наблюдений сдвигом вправо своего математического ожидания и измененной дисперсией. В частности, для случая $\lambda > 0$ и $\nu = 1$ получаем модель, учитывающую только сдвиг вправо математического ожидания выброса, а для $\lambda = 0$ и $\nu \neq 1$ модель учитывает только изменение дисперсии выброса.

Статистика критерия Граббса для проверки на один верхний выброс имеет вид [2]

$$G_n = (X_{(n)} - \bar{X}) / S, \quad (1)$$

где $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ и $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$.

Гипотеза H_0 отвергается, если наблюдаемое значение статистики Граббса превысит критическое $G_{n;\alpha}^{cr}$, отвечающее выбранному уровню значимости α .

Для исследования мощностных свойств критерия Граббса будем использовать следующие меры мощности Дэйвида для статистики Граббса [1]:

$$\begin{aligned} P_1 &= P(G_n > G_{n;\alpha}^{cr} | H_1), \\ P_2 &= P((X_{out} - \bar{X})/S > G_{n;\alpha}^{cr} | H_1), \\ P_3 &= P(\{G_n > G_{n;\alpha}^{cr}\} \cap \{X_{(n)} = X_{out}\} | H_1), \\ P_4 &= P(\{(X_{out} - \bar{X})/S > G_{n;\alpha}^{cr}\} \cap \{(X_{(n-1)} - \bar{X})/S < G_{n;\alpha}^{cr}\} | H_1). \end{aligned}$$

Мера мощности P_1 является, по сути, «классической» функцией мощности критерия Граббса, ибо она равна вероятности не допустить ошибку второго рода. Поэтому мера P_1 особенно подходит для задачи выявления выборок с аномальными наблюдениями. Меры P_2 – P_4 могут быть использованы для обнаружения выброса в выборке. Заметим также, что вероятность выявления выборки с аномальным наблюдением не совпадает с вероятностью обнаружения выброса, так как [1]

$$P_1 \geq P_2 \geq P_3 \geq P_4. \quad (2)$$

В работе [3] найдены формулы для вычисления мер мощности P_1 – P_4 в случае, когда для выборки объёма n ($n \geq 3$) справедлива гипотеза H_1 :

$$P_1(t) = \begin{cases} 1, & t \leq \frac{1}{\sqrt{n}}; \\ 1 - \int_{-\frac{n-1}{\sqrt{n}}}^t F_{n-1}^{(1)}(\rho_n(t, x)) f_{\bar{T}_n}(x) dx, & \frac{1}{\sqrt{n}} < t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}; \end{cases} \quad (3)$$

$$P_2(t) = \begin{cases} 1, & t < -\frac{n-1}{\sqrt{n}}; \\ \int_t^{\frac{n-1}{\sqrt{n}}} f_{\bar{T}_n}(\xi) d\xi, & -\frac{n-1}{\sqrt{n}} \leq t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}; \end{cases} \quad (4)$$

$$P_3(t) = \begin{cases} \int_{\frac{1}{\sqrt{n}}}^{\frac{n-1}{\sqrt{n}}} F_{n-1}^{(1)}(g_n(x)) f_{\bar{T}_n}(x) dx, & t \leq \frac{1}{\sqrt{n}}; \\ \int_t^{\frac{n-1}{\sqrt{n}}} F_{n-1}^{(1)}(g_n(x)) f_{\bar{T}_n}(x) dx, & \frac{1}{\sqrt{n}} < t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}; \end{cases} \quad (5)$$

$$P_4(t) = \begin{cases} 0, & t < -\frac{1}{\sqrt{n}}; \\ \int_t^{\frac{n-1}{\sqrt{n}}} F_{n-1}^{(1)}(\rho_n(t, x)) f_{\tilde{T}_n}(x) dx, & -\frac{1}{\sqrt{n}} \leq t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}. \end{cases} \quad (6)$$

Здесь t – критическое значение статистики (1), отвечающее некоторому уровню значимости α ($0 \leq \alpha \leq 1$); $\rho_n(t, x) = \left(t + \frac{x}{n-1}\right) / \sqrt{\frac{n-1}{n-2} \left(1 - \frac{n}{(n-1)^2} x^2\right)}$, $|x| < \frac{n-1}{\sqrt{n}}$; $g_n(x) = \rho_n(x, x)$; $F_m^{(1)}(t) = P(G_m < t | H_0)$.

Закон распределения случайной величины G_m в условиях справедливой гипотезы H_0 известен [4, 5]:

$$F_m^{(1)}(t) = \begin{cases} 0, & t \leq \frac{1}{\sqrt{m}}, \\ 1, & t > \frac{m-1}{\sqrt{m}}, \quad m \geq 2; \\ m \int_{\frac{1}{\sqrt{m}}}^t F_{m-1}^{(1)}(g_m(x)) f_{T_m}(x) dx, & \frac{1}{\sqrt{m}} < t \leq \frac{m-1}{\sqrt{m}}, \quad m \geq 3, \end{cases} \quad (7)$$

где $f_{T_m}(x) = \frac{1}{m-1} \sqrt{\frac{m}{\pi}} \Gamma\left(\frac{m-1}{2}\right) / \Gamma\left(\frac{m-2}{2}\right) \left(1 - \frac{m}{(m-1)^2} x^2\right)^{\frac{m-4}{2}}$, $|x| < \frac{m-1}{\sqrt{m}}$; $\Gamma(x) = \int_0^{+\infty} \xi^{x-1} e^{-\xi} d\xi$ – гамма-функция.

Для вычисления мер мощности по формулам (3)–(6) следует знать закон распределения случайной величины

$$\tilde{T}_n = (X_{out} - \bar{X})/S. \quad (8)$$

Случайная величина \tilde{T}_n является студентизированным отклонением выброса X_{out} от среднего, найденным по выборке объёма n . В работе [3] доказана следующая теорема о законе её распределения.

ТЕОРЕМА 1. Пусть наблюдения $X_1, X_2, \dots, X_{n-1}, X_{out}$ являются независимыми случайными величинами, причём X_1, X_2, \dots, X_{n-1} имеют нормальное распределение $N(a, \sigma^2)$, а выброс X_{out} имеет нормальный закон распределения $N(a + \lambda\sigma, \nu\sigma^2)$, где $\lambda \geq 0$ и $\nu > 0$. Тогда для $n \geq 3$ плотность распределения вероятностей случайной величины \tilde{T}_n имеет вид

$$f_{\tilde{T}_n}(t) = \begin{cases} K_n \left[\frac{(n-1)^2}{n} - t^2\right]^{\frac{n-4}{2}} I_n(t), & |t| < \frac{n-1}{\sqrt{n}}, \\ 0, & |t| \geq \frac{n-1}{\sqrt{n}}, \end{cases} \quad (9)$$

где

$$K_n = \frac{(n-1)^2}{n\sqrt{2\pi}\Gamma\left(\frac{n-2}{2}\right)} \left(\frac{\eta}{2}\right)^{\frac{n-2}{2}} e^{-\frac{\mu^2}{2}}, \quad \eta = \frac{1 + \nu(n-1)}{n}, \quad \mu = \lambda\sqrt{\frac{n-1}{n\eta}},$$

$$I_n(t) = \int_0^\infty y^{\frac{n-3}{2}} e^{t\mu\sqrt{y}-0,5q(t)y} dy, \quad q(t) = \eta\frac{(n-1)^2}{n} + (1-\eta)t^2.$$

В данной работе найдено новое представление для плотности $f_{\bar{T}_n}$ через специальную функцию Эрмита и показано, что это представление совпадает с формулой (9). Найденное представление было использовано для получения интегральных представлений для мер мощности критерия Граббса, используя функции Эрмита. Отдельно рассмотрен случай, когда статистика критерия является статистикой отношения правдоподобия.

1. Вывод основных соотношений для плотности вероятностей случайной величины \bar{T}_n . Чтобы вывести новое представление для плотности распределения вероятностей случайной величины (8), докажем два вспомогательных утверждения.

ЛЕММА 1. Пусть W и Z — независимые случайные величины с плотностями распределения вероятностей f_W и f_Z и областями значений \mathbb{R}_+ и \mathbb{R} соответственно. Тогда $\forall r \in \mathbb{R}$ случайная величина $U = r\sqrt{W} - Z$ имеет плотность распределения вероятностей

$$f_U(u) = \int_0^\infty f_W(x)f_Z(r\sqrt{x} - u)dx. \quad (10)$$

Доказательство. Интегральная функция распределения случайной величины $U = r\sqrt{W} - Z \forall r \in \mathbb{R}$ следующая:

$$F_U(u) = P(U < u) = P(r\sqrt{W} - Z < u) = P(Z > r\sqrt{W} - u).$$

По условию W и Z — случайные величины с областями значений \mathbb{R}_+ и \mathbb{R} соответственно, поэтому

$$F_U(u) = P(\{r\sqrt{W} - u < Z < \infty\} \cap \{W > 0\}).$$

Случайные величины W и Z независимы, следовательно,

$$F_U(u) = \int_0^\infty f_W(x)dx \int_{r\sqrt{x}-u}^\infty f_Z(z)dz. \quad (11)$$

Продифференцировав (11), получим плотность случайной величины U :

$$f_U(u) = \frac{dF_U(u)}{du} = \int_0^\infty f_W(x)f_Z(r\sqrt{x} - u)dx,$$

что и требовалось доказать. \square

ЛЕММА 2. Пусть случайные величины W и Z являются независимыми, причём величина W имеет распределение $\chi^2(n-1)$, а Z имеет нормальное распределение $N(a_0, 1)$. Тогда $\forall t \in \left(-\frac{n-1}{\sqrt{n}}; \frac{n-1}{\sqrt{n}}\right)$ плотность распределения вероятностей случайной величины $V_n(t) = \beta_n(t)\sqrt{W} - Z$ в точке $v = 0$ следующая:

$$f_{V_n(t)}(0) = \frac{A_n}{h^{\frac{n-1}{2}}(t)} [(n-1)^2 - nt^2]^{\frac{n-1}{2}} H_{-n+1}\left(-\frac{a_0 t}{\sqrt{2h(t)}}\right), \quad (12)$$

где

$$\beta_n(t) = \frac{ct}{\sqrt{(n-1)^2 - nt^2}}, \quad c = \text{const} > 0;$$

$$A_n = c^{-n+1} \sqrt{\frac{2}{\pi}} e^{-\frac{a_0^2}{2}} \cdot \frac{\Gamma(n-1)}{\Gamma(\frac{n-1}{2})}; \quad h(t) = \frac{(n-1)^2 + (c^2 - n)t^2}{c^2};$$

$$H_k(z) = \frac{1}{\Gamma(-k)} \int_0^\infty e^{-\xi^2 - 2z\xi} \xi^{-k-1} d\xi, \quad k < 0.$$

Доказательство. Обозначим плотности случайных величин W и Z через $f_{\chi_{n-1}^2}$ и f_Z соответственно. По условию W и Z — независимые случайные величины с областями значений \mathbb{R}_+ и \mathbb{R} соответственно. Следовательно, $\forall t \in \left(-\frac{n-1}{\sqrt{n}}; \frac{n-1}{\sqrt{n}}\right)$ плотность случайной величины $V_n(t) = \beta_n(t)\sqrt{W} - Z$ можно найти по формуле (10):

$$f_{V_n(t)}(v) = \int_0^\infty f_{\chi_{n-1}^2}(x) f_Z(\beta_n(t)\sqrt{x} - v) dx.$$

Здесь плотность величины W следует вычислять по формуле [6]

$$f_{\chi_k^2}(x) = \begin{cases} \frac{x^{\frac{k-2}{2}} e^{-\frac{x}{2}}}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})}, & x \geq 0, \\ 0, & x < 0, \end{cases} \quad (13)$$

при $k = n - 1$, а плотность величины Z — по формуле

$$f_Z(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-a_0)^2}{2}}. \quad (14)$$

Далее, в точке $v = 0$ плотность величины $V_n(t)$ есть

$$f_{V_n(t)}(0) = \int_0^\infty f_{\chi_{n-1}^2}(x) f_Z(\beta_n(t)\sqrt{x}) dx. \quad (15)$$

Сделав под знаком интеграла в (15) замену переменных $\xi = \sqrt{\frac{1+\beta_n^2(t)}{2}}\sqrt{x}$, приведём его к виду

$$f_{V_n(t)}(0) = \frac{2f_Z(0)\Gamma(n-1)}{(1 + \beta_n^2(t))^{\frac{n-1}{2}} \Gamma(\frac{n-1}{2})} H_{-n+1}\left(-\frac{a_0\beta_n(t)}{\sqrt{2(1 + \beta_n^2(t))}}\right),$$

где $H_k(z) = \frac{1}{\Gamma(-k)} \int_0^\infty e^{-\xi^2 - 2z\xi} \xi^{-k-1} d\xi$ — функция Эрмита с отрицательным целым значком ($k < 0$) [7].

Положим

$$h(t) = \frac{(n-1)^2 + (c^2 - n)t^2}{c^2}, \quad A_n = c^{-n+1} \sqrt{\frac{2}{\pi}} e^{-\frac{a_0^2}{2}} \cdot \frac{\Gamma(n-1)}{\Gamma(\frac{n-1}{2})},$$

где $c = \text{const} > 0$.

Легко убедиться, что $\forall c > 0$ справедливы равенства

$$\frac{a_0 \beta_n(t)}{\sqrt{2(1 + \beta_n^2(t))}} = \frac{a_0 t}{\sqrt{2h(t)}}, \quad \frac{1}{(\beta_n(t) + 1)^{\frac{n-1}{2}}} = \frac{c^{-n+1}}{h^{\frac{n-1}{2}}(t)} [(n-1)^2 - nt^2]^{\frac{n-1}{2}}.$$

Отсюда

$$f_{V_n(t)}(0) = \frac{A_n}{h^{\frac{n-1}{2}}(t)} [(n-1)^2 - nt^2]^{\frac{n-1}{2}} H_{-n+1}\left(-\frac{a_0 t}{\sqrt{2h(t)}}\right),$$

что и требовалось доказать. \square

Используя соотношения (10) и (12), можно найти новое представление для плотности вероятностей случайной величины \tilde{T}_n .

ТЕОРЕМА 2. Пусть выполняются условия теоремы 1. Тогда для $n \geq 3$ плотность распределения вероятностей случайной величины \tilde{T}_n имеет вид

$$f_{\tilde{T}_n}(t) = \begin{cases} B_n \cdot \frac{(n-1)^2}{\sqrt{[(n-1)^2 - nt^2]^3}} \cdot f_{V_n(t)}(0), & |t| < \frac{n-1}{\sqrt{n}}, \\ 0, & |t| \geq \frac{n-1}{\sqrt{n}}, \end{cases} \quad (16)$$

где $B_n = \sqrt{\frac{2n}{\eta}} \cdot \frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n-2}{2})}$, $\eta = \frac{1+\nu(n-1)}{n}$; $V_n(t) = \beta_n(t)\sqrt{W} - Z$; $f_{V_n(t)}(0)$ вычисляется по формуле (12) при $c = \sqrt{\frac{n}{\eta}}$; случайные величины W и Z являются независимыми; W имеет распределение $\chi^2(n-1)$, Z имеет нормальное распределение $N(a_0, 1)$, $a_0 = \mu$, $\mu = \lambda\sqrt{\frac{n-1}{n\eta}}$.

Доказательство. Нетрудно проверить, что при $n \geq 3$ справедливы следующие равенства:

$$X_{\text{out}} - \bar{X} = \frac{n-1}{n} (X_{\text{out}} - \bar{X}^*), \quad S^2 = \frac{n-2}{n-1} S^{*2} + \frac{(X_{\text{out}} - \bar{X}^*)^2}{n}, \quad (17)$$

где $\bar{X}^* = \frac{1}{n-1} \sum_{i=1}^{n-1} X_i$ и $S^{*2} = \frac{1}{n-2} \sum_{i=1}^{n-1} (X_i - \bar{X}^*)^2$ вычисляются по выборке объёма $(n-1)$, не содержащей выброса.

С учётом (17) для случайной величины \tilde{T}_n получим

$$\tilde{T}_n = \frac{n-1}{\sqrt{n}} \frac{\text{sign}(X_{\text{out}} - \bar{X}^*)}{\sqrt{\frac{n(n-2)S^{*2}}{(n-1)(X_{\text{out}} - \bar{X}^*)^2} + 1}}. \quad (18)$$

Рассмотрим случайные величины

$$Y = (n-2)S^{*2}/\sigma^2, \quad Z = \sqrt{\frac{n-1}{n\eta}}(X_{\text{out}} - \bar{X}^*)/\sigma. \quad (19)$$

Согласно теореме Фишера [6] случайная величина Y имеет распределение $\chi^2(n-2)$. Легко проверить, что случайная величина Z имеет нормальное распределение $N(\mu, 1)$.

С учётом введенных обозначений формула (18) примет вид

$$\tilde{T}_n = \frac{n-1}{\sqrt{n}} \frac{\text{sign}(Z)}{\sqrt{1 + \eta^{-1}Y/Z^2}}. \quad (20)$$

Пусть $t < 0$. Найдём интегральную функцию $F_{\tilde{T}_n}(t) = P(\tilde{T}_n < t)$. С учётом (20) получаем

$$F_{\tilde{T}_n}(t) = P\left(\{Z < 0\} \cap \left\{-\frac{1}{\sqrt{1 + \eta^{-1}Y/Z^2}} < \frac{\sqrt{nt}}{n-1}\right\}\right),$$

$$F_{\tilde{T}_n}(t) = P\left(\{Z < 0\} \cap \left\{\frac{Y}{\eta Z^2} < \frac{(n-1)^2}{nt^2} - 1\right\}\right).$$

Заметим, что для $t \leq -\frac{n-1}{\sqrt{n}}$ событие $\left\{\frac{Y}{\eta Z^2} < \frac{(n-1)^2}{nt^2} - 1\right\}$ становится невозможным, следовательно,

$$F_{\tilde{T}_n}(t) = 0, \quad t \leq -\frac{n-1}{\sqrt{n}}. \quad (21)$$

Далее $\forall t \in (-\frac{n-1}{\sqrt{n}}; 0)$ имеем

$$F_{\tilde{T}_n}(t) = P\left(\left\{Z < \sqrt{\frac{n}{\eta}} \cdot \frac{t}{\sqrt{(n-1)^2 - nt^2}} \sqrt{Y}\right\}\right).$$

Как и ранее (см. лемму 2), обозначим $\beta_n(t) = \frac{ct}{\sqrt{(n-1)^2 - nt^2}}$, $c = \text{const} > 0$.

Положим $c = \sqrt{\frac{n}{\eta}}$, тогда $F_{\tilde{T}_n}(t) = P\left(\left\{Z < \beta_n(t)\sqrt{Y}\right\} \cap \{0 < Y < \infty\}\right)$.

Согласно теореме Фишера случайные величины Z и Y являются независимыми [6], следовательно,

$$F_{\tilde{T}_n}(t) = \int_0^\infty f_Y(y) dy \int_{-\infty}^{\beta_n(t)\sqrt{y}} f_Z(z) dz. \quad (22)$$

Здесь плотность $f_Z(z)$ случайной величины Z вычисляется по формуле (14) при условии $a_0 = \mu$, а плотность $f_Y(y)$ случайной величины Y вычисляется согласно (13) для $k = n - 2$.

Аналогичным образом можно получить для случая $t \geq 0$:

$$F_{\tilde{T}_n}(t) = \begin{cases} \frac{1}{2} - \Phi(\mu) + \int_0^\infty f_Y(y) dy \int_0^{\sqrt{y}\beta_n(t)} f_Z(z) dz, & 0 \geq t < \frac{n-1}{\sqrt{n}}, \\ 1, & t \geq \frac{n-1}{\sqrt{n}}, \end{cases} \quad (23)$$

где $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_0^z e^{-\frac{\xi^2}{2}} d\xi$ — функция Лапласа.

Объединяя соотношения (21), (22) и (23), получим

$$F_{\tilde{T}_n}(t) = \begin{cases} 0, & t \leq -\frac{n-1}{\sqrt{n}}; \\ \int_0^\infty f_Y(y)dy \int_{-\infty}^{\beta_n(t)\sqrt{y}} f_Z(z)dz, & -\frac{n-1}{\sqrt{n}} < t < 0; \\ \frac{1}{2} - \Phi(\mu) + \int_0^\infty f_Y(y)dy \int_0^{\beta_n(t)\sqrt{y}} f_Z(z)dz, & 0 \leq t < \frac{n-1}{\sqrt{n}}; \\ 1, & t \geq \frac{n-1}{\sqrt{n}}. \end{cases}$$

Продифференцировав $F_{\tilde{T}_n}(t)$, найдём плотность распределения вероятностей случайной величины \tilde{T}_n :

$$f_{\tilde{T}_n}(t) = \begin{cases} 0, & |t| \geq \frac{n-1}{\sqrt{n}}; \\ \beta'_n(t) \int_0^\infty \sqrt{y} f_Y(y) f_Z(\beta_n(t)\sqrt{y}) dy, & |t| < \frac{n-1}{\sqrt{n}}, \end{cases} \quad (24)$$

где

$$\beta'_n(t) = \sqrt{\frac{n}{\eta}} \cdot \frac{(n-1)^2}{[(n-1)^2 - nt^2]^{\frac{3}{2}}}. \quad (25)$$

Используя соотношение (13), легко проверить, что

$$\sqrt{y} f_Y(y) = \frac{\sqrt{2}\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n-2}{2})} f_{\chi_{n-1}^2}(y). \quad (26)$$

Тогда соотношение (24) с учётом (25) и (26) примет вид

$$f_{\tilde{T}_n}(t) = \begin{cases} B_n \cdot \frac{(n-1)^2}{\sqrt{[(n-1)^2 - nt^2]^3}} \cdot \int_0^\infty f_{\chi_{n-1}^2}(y) f_Z(\beta_n(t)\sqrt{y}) dy, & |t| < \frac{n-1}{\sqrt{n}}, \\ 0, & |t| \geq \frac{n-1}{\sqrt{n}}, \end{cases}$$

где $B_n = \sqrt{\frac{2n}{\eta}} \cdot \frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n-2}{2})}$.

Положим, что имеется случайная величина W , распределённая по закону $\chi^2(n-1)$; кроме того, будем считать, что W и Z — независимые случайные величины и величина Z определена согласно (19). Тогда $\forall t \in (-\frac{n-1}{\sqrt{n}}; \frac{n-1}{\sqrt{n}})$ и $n \geq 3$ можно определить случайную величину $V_n(t) = \beta_n(t)\sqrt{W} - Z$ с плотностью $f_{V_n(t)}$.

Используя формулу (10), получаем, что интеграл в (27) равен

$$\int_0^\infty f_W(y) f_Z(\beta_n(t)\sqrt{y}) dy = f_{V_n(t)}(0).$$

Поэтому соотношение (27) примет вид

$$f_{\tilde{T}_n}(t) = \begin{cases} B_n \cdot \frac{(n-1)^2}{\sqrt{[(n-1)^2 - nt^2]^3}} \cdot f_{V_n(t)}(0), & |t| < \frac{n-1}{\sqrt{n}}, \\ 0, & |t| \geq \frac{n-1}{\sqrt{n}}, \end{cases}$$

что и требовалось доказать. \square

Теперь, используя леммы 2 и 3, можно получить представление плотности $f_{\tilde{T}_n}$ через функцию Эрмита.

Для этого найдём по формуле (12) плотность $f_{V_n(t)}(0)$ для случая, когда $c = \sqrt{\frac{n}{\eta}}$ и $a_0 = \mu$, и подставим полученное выражение в (16). Легко убедиться, что в результате формула (16) примет вид

$$f_{\tilde{T}_n}(t) = \begin{cases} D_n q^{-\frac{n-1}{2}}(t) \left[1 - \frac{nt^2}{(n-1)^2}\right]^{\frac{n-4}{2}} H_{-n+1}\left(-\frac{t\mu}{\sqrt{2q(t)}}\right), & |t| < \frac{n-1}{\sqrt{n}}, \\ 0, & |t| \geq \frac{n-1}{\sqrt{n}}, \end{cases} \quad (28)$$

где $D_n = \frac{2\Gamma(n-1)}{\Gamma(\frac{n-2}{2})} \cdot \frac{e^{-\frac{\mu^2}{2}}}{\sqrt{\pi}} \cdot \left(\sqrt{\frac{n}{\eta}} \cdot (n-1)\right)^{n-2}$.

Покажем, что соотношения (9) и (28) для плотности $f_{\tilde{T}_n}(t)$ совпадают.

Для $k = -n + 1$ можно записать [7]

$$H_{-n+1}\left(-\frac{t\mu}{\sqrt{2q(t)}}\right) = \frac{1}{\Gamma(n-1)} \int_0^\infty e^{-\xi^2 + \frac{2t\mu}{\sqrt{2q(t)}}\xi} \xi^{n-2} d\xi.$$

Вводя новую переменную интегрирования $y = 2\xi^2/q(t)$, можно представить последнее соотношение в виде

$$H_{-n+1}\left(-\frac{t\mu}{\sqrt{2q(t)}}\right) = \frac{q^{\frac{n-1}{2}}(t)}{2^{\frac{n+1}{2}}\Gamma(n-1)} \int_0^\infty e^{t\mu\sqrt{y}-0,5q(t)y} y^{\frac{n-3}{2}} dy.$$

Поскольку

$$I_n(t) = \int_0^\infty e^{t\mu\sqrt{y}-0,5q(t)y} y^{\frac{n-3}{2}} dy$$

(см. теорему 1), то равенство (29) принимает вид

$$H_{-n+1}\left(-\frac{t\mu}{\sqrt{2q(t)}}\right) = \frac{q^{\frac{n-1}{2}}(t)}{2^{\frac{n+1}{2}}\Gamma(n-1)} I_n(t). \quad (30)$$

С учётом (30) соотношение (28) можно привести к виду

$$f_{\tilde{T}_n}(t) = \begin{cases} \frac{D_n}{2^{\frac{n+1}{2}}\Gamma(n-1)} \left(\frac{\sqrt{n}}{n-1}\right)^{n-4} \left[\frac{(n-1)^2}{n} - t^2\right]^{\frac{n-4}{2}} I_n(t), & |t| < \frac{n-1}{\sqrt{n}}, \\ 0, & |t| \geq \frac{n-1}{\sqrt{n}}. \end{cases} \quad (31)$$

Если вернуться к обозначению (см. теорему 1)

$$K_n = \frac{(n-1)^2}{n\sqrt{2\pi}\Gamma(\frac{n-2}{2})} \left(\frac{\eta}{2}\right)^{\frac{n-2}{2}} e^{-\frac{\mu^2}{2}},$$

то легко убедиться в справедливости равенства

$$K_n = \frac{D_n}{2^{\frac{n+1}{2}}\Gamma(n-1)} \left(\frac{\sqrt{n}}{n-1}\right)^{n-4}.$$

Поэтому соотношение (31) принимает вид соотношения (9). Откуда следует, что соотношения (9) и (28) совпадают.

2. Интегральные представления для мер мощности в случае, когда статистика критерия является статистикой отношения правдоподобия. Предположим, что выброс X_{out} имеет сдвиг вправо в математическом ожидании и ту же дисперсию, что и остальные наблюдения в выборке. Это предположение часто считают правдоподобным приближением к действительности, поэтому его исследование наиболее интересно с практической точки зрения. Известно, что в такой ситуации статистика G_n становится статистикой отношения правдоподобия [5]. В этом случае параметр масштаба $\nu = 1$, в то время как параметр сдвига $\lambda > 0$. Поэтому соотношение (28) примет вид

$$f_{\tilde{T}_n}(t) = \begin{cases} \frac{\sqrt{n}}{n-1} d_n \left[1 - \frac{nt^2}{(n-1)^2} \right]^{\frac{n-4}{2}} H_{-n+1} \left(-\frac{t\mu}{\sqrt{2q_0}} \right), & |t| < \frac{n-1}{\sqrt{n}}, \\ 0, & |t| \geq \frac{n-1}{\sqrt{n}}, \end{cases} \quad (32)$$

где $d_n = \frac{\Gamma(n-1)}{\Gamma(\frac{n-2}{2})} \cdot \frac{2}{\sqrt{\pi}} e^{-\frac{\mu^2}{2}}$, $\mu = \lambda \sqrt{\frac{n-1}{n}}$, $q_0 = \frac{(n-1)^2}{n}$.

Подставив (32) в (3)–(6) и выполнив замену переменных $x = \frac{n-1}{\sqrt{n}} \sin \varphi$, получим следующие интегральные представления для мер мощности:

$$P_1(t) = \begin{cases} 1, & t \leq \frac{1}{\sqrt{n}}; \\ 1 - d_n \int_{-\pi/2}^{\arcsin \frac{\sqrt{nt}}{n-1}} F_{n-1}^{(1)}(r_n(t, \varphi)) \times \\ \quad \times \cos^{n-3} \varphi H_{-n+1}(\theta_n(\varphi)) d\varphi, & \frac{1}{\sqrt{n}} < t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}; \end{cases} \quad (33)$$

$$P_2(t) = \begin{cases} 1, & t < -\frac{n-1}{\sqrt{n}}; \\ d_n \int_{\arcsin \frac{\sqrt{nt}}{n-1}}^{\pi/2} \cos^{n-3} \varphi H_{-n+1}(\theta_n(\varphi)) d\varphi, & -\frac{n-1}{\sqrt{n}} \leq t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}; \end{cases} \quad (34)$$

$$P_3(t) = \begin{cases} d_n \int_{\arcsin \frac{1}{n-1}}^{\pi/2} F_{n-1}^{(1)}(s_n(\varphi)) \times \\ \quad \times \cos^{n-3} \varphi H_{-n+1}(\theta_n(\varphi)) d\varphi, & t \leq \frac{1}{\sqrt{n}}; \\ d_n \int_{\arcsin \frac{\sqrt{nt}}{n-1}}^{\pi/2} F_{n-1}^{(1)}(s_n(\varphi)) \times \\ \quad \times \cos^{n-3} \varphi H_{-n+1}(\theta_n(\varphi)) d\varphi, & \frac{1}{\sqrt{n}} < t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}; \end{cases} \quad (35)$$

$$P_4(t) = \begin{cases} 0, & t < -\frac{1}{\sqrt{n}}; \\ d_n \int_{\arcsin \frac{\sqrt{nt}}{n-1}}^{\pi/2} F_{n-1}^{(1)}(r_n(t, \varphi)) \times \\ \quad \times \cos^{n-3} \varphi H_{-n+1}(\theta_n(\varphi)) d\varphi, & -\frac{1}{\sqrt{n}} \leq t \leq \frac{n-1}{\sqrt{n}}; \\ 0, & t > \frac{n-1}{\sqrt{n}}. \end{cases} \quad (36)$$

Здесь $r_n(t, \varphi) = \sqrt{\frac{n-2}{n(n-1)}} \left(\frac{\sqrt{nt}}{\cos \varphi} + \operatorname{tg} \varphi \right)$, $s_n(\varphi) = \sqrt{\frac{(n-2)n}{n-1}} \operatorname{tg} \varphi$, $|\varphi| < \pi/2$;
 $\theta_n(\varphi) = -\frac{\mu}{\sqrt{2}} \sin \varphi$.

Из соотношений (33) и (36) следует, что $\forall t \in \left(\frac{1}{\sqrt{n}}, \frac{n-1}{\sqrt{n}} \right)$

$$P_1(t) = \delta_n(t, \lambda) + P_4(t).$$

Здесь

$$\delta_n(t, \lambda) = 1 - d_n \int_{-\pi/2}^{\pi/2} F_{n-1}^{(1)}(r_n(t, \varphi)) \cos^{n-3} \varphi H_{-n+1}(\theta_n(\varphi)) d\varphi \quad (37)$$

— величина отклонения.

Из условия (2) вытекает, что

$$P_1(t) - P_2(t) \leq P_1(t) - P_3(t) \leq P_1(t) - P_4(t) \equiv \delta_n(t, \lambda).$$

Таким образом, величина $\delta_n(t, \lambda)$ на заданном уровне значимости α определяет максимально возможное отклонение вероятности обнаружить присутствие одиночного выброса в выборке (мера P_1) от вероятности его точного обнаружения (меры P_2 – P_4).

Из формулы (37) следует также, что $\delta_n(t, \lambda)$ — убывающая функция аргумента t . Поскольку с ростом критических значений уровень значимости критерия уменьшается, то переход, например, от пятипроцентного уровня значимости к однопроцентному приведёт лишь к уменьшению величины $\delta_n(t, \lambda)$.

Легко проверить, что в области

$$Q = \left[t_f \leq t \leq \frac{n-1}{\sqrt{n}}, \arcsin \sqrt{\frac{n-2}{2(n-1)}} \leq \varphi \leq \pi/2 \right], \quad \left(t_f = \sqrt{\frac{(n-1)(n-2)}{2n}} \right)$$

имеем

$$r_n(t, \varphi) \geq \frac{n-2}{\sqrt{n-1}}, \quad s_n(\varphi) \geq \frac{n-2}{\sqrt{n-1}}.$$

Следовательно, с учётом (7) получим

$$F_{n-1}^{(1)}(r_n(t, \varphi)) = 1, \quad (t, \varphi) \in Q;$$

$$F_{n-1}^{(1)}(s_n(\varphi)) = 1 \quad \left(\arcsin \sqrt{\frac{n-2}{2(n-1)}} \leq \varphi \leq \pi/2 \right).$$

Поэтому для критических значений t , удовлетворяющих условию $t_f \leq t \leq \frac{n-1}{\sqrt{n}}$, третья и четвёртая меры мощности совпадают со второй:

$$P_3(t) = P_4(t) = P_2(t).$$

n	t_f	$\alpha(n; t_f)$	n	t_f	$\alpha(n; t_f)$
4	0,87	0,8453	19	2,84	0,0109
5	1,10	0,6806	20	2,92	0,0079
6	1,29	0,5334	21	3,01	0,0058
7	1,46	0,4109	22	3,09	0,0042
8	1,62	0,3126	23	3,17	0,0030
9	1,76	0,2356	24	3,25	0,0022
10	1,90	0,1763	25	3,32	0,0016
11	2,02	0,1312	26	3,40	0,0011
12	2,14	0,0972	27	3,47	0,0008
13	2,25	0,0717	28	3,54	0,0006
14	2,36	0,0527	29	3,61	0,0004
15	2,46	0,0387	30	3,68	0,0003
16	2,56	0,0283	31	3,75	0,0002
17	2,66	0,0206	32	3,81	0,0002
18	2,75	0,0150	33	3,88	0,0001

Уровни значимости $\alpha(n; t_f)$ критерия, соответствующие критическим значениям t_f , можно найти из условия

$$\alpha(n; t_f) = P(G_n > t_f | H_0) = 1 - F_n^{(1)}(t_f),$$

где значение $F_n^{(1)}(t_f)$ может быть найдено численно по формуле (7).

В таблице приведены результаты численных расчётов уровней значимости $\alpha(n; t_f)$ для критических значений t_f нулевого распределения статистики G_n в случае выборок объёмов n от 4 до 33.

Из таблицы видно, что на пятипроцентном уровне значимости для числа наблюдений $4 \leq n \leq 14$ получаем

$$P_2(G_{n;0,05}^{cr}) = P_3(G_{n;0,05}^{cr}) = P_4(G_{n;0,05}^{cr}).$$

Для выборок объёмом $15 \leq n \leq 19$ имеем

$$P_2(G_{n;0,01}^{cr}) = P_3(G_{n;0,01}^{cr}) = P_4(G_{n;0,01}^{cr}).$$

Для $n \geq 20$ можно записать

$$P_2(G_{n;\alpha < 0,005}^{cr}) = P_3(G_{n;\alpha < 0,005}^{cr}) = P_4(G_{n;\alpha < 0,005}^{cr}).$$

Обычно исследователь использует «стандартные» уровни значимости ($0,01 \leq \alpha \leq 0,05$), поэтому можно быть уверенным, что для выборок объёмом $n \geq 20$ имеет место строгое неравенство

$$P_2(G_{n;\alpha}^{cr}) > P_3(G_{n;\alpha}^{cr}) > P_4(G_{n;\alpha}^{cr}), \quad 0,01 \leq \alpha \leq 0,05.$$

3. Численное моделирование. Полученные интегральные соотношения для мер (33)–(36) могут быть использованы для исследования чувствительности критерия Граббса к наличию в выборке одного выброса. Цель расчётов — продемонстрировать возможности практического применения мер мощности. Для этого был разработан алгоритм вычисления мер мощности по формулам (33)–(36). Определенные интегралы в мерах P_1 – P_4 вычислялись приближенно по формуле Симпсона, при этом значения функции распределения $F_{n-1}^{(1)}(x)$ вычислялись рекурсивно по формуле (7), плотности $f_{\tilde{T}_n}(t)$ — по формуле (28). Алгоритм вычисления мер мощности P_1 – P_4 был реализован на языке программирования Object Pascal.

Численные расчёты мер P_1 – P_4 , а также отклонения δ_n были выполнены для случая нормально распределенной выборки, в которой присутствовал выброс с параметрами $\lambda > 0$ и $\nu = 1$.

На рис. 1 представлены результаты численных расчётов меры $P_1(t_{cr})$ по формуле (33) для выборок объёмом n от 5 до 100 и значений параметра λ от 1 до 5. Величина t_{cr} была выбрана равной критическому значению статистики G_n на стандартном уровне значимости $\alpha = 0,05$, т. е. $t_{cr} = G_{n;0,05}^{cr}$. Мера $P_1(t_{cr})$, таким образом, равна вероятности не совершить ошибку второго рода на пятипроцентном уровне значимости. Как и следовало ожидать в соответствии с общей теорией вопроса, для выборки фиксированного объёма с ростом параметра λ наблюдалось увеличение мощности критерия. Из рис. 1 видно, как мера $P_1(t_{cr})$ возрастает с ростом параметра сдвига λ для разных значений n . Вероятность не совершить ошибку второго рода, т. е. обнаружить присутствие аномального наблюдения в выборке объёма n , мала для $\lambda \leq 2$ и близка к 1 для $\lambda \geq 4$.

Видно также, что вероятность обнаружить присутствие выброса в выборке из 20 наблюдений практически не отличается от вероятности обнаружить его присутствие в выборке из 100 наблюдений.

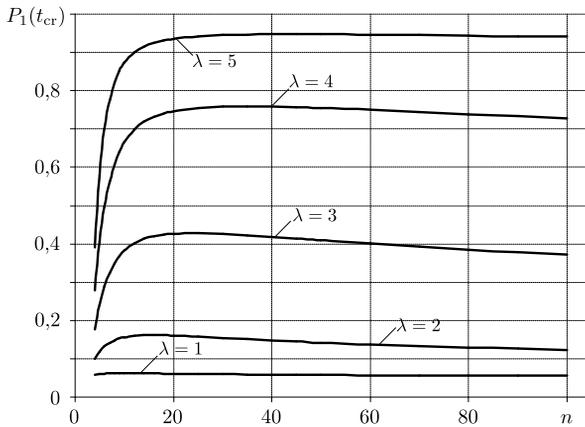


Рис. 1. Графики зависимости первой меры мощности $P_1(t_{cr})$ от объёма выборки n для значений λ от 1 до 5 ($t_{cr} = G_{n;0,05}^{cr}$)

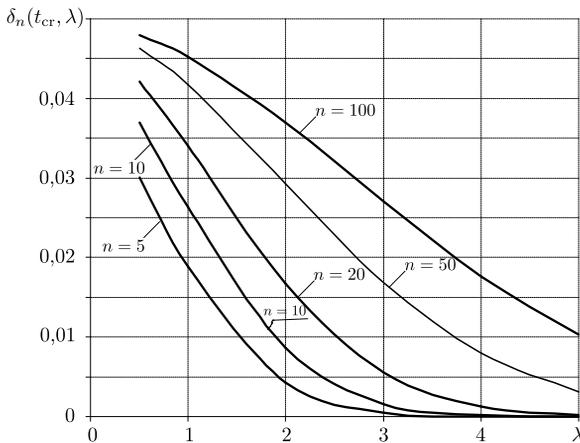


Рис. 2. Графики зависимости отклонения $\delta_n(t_{cr}, \lambda)$ от λ для объёмов выборок n от 5 до 100 ($t_{cr} = G_{n;0,05}^{cr}$)

На рис. 2 представлены результаты численных расчётов по формуле (37) величины отклонения $\delta_n(t_{cr}, \lambda) = P_1(t_{cr}) - P_4(t_{cr})$ для случая, когда в нормально распределенной выборке объёма n имеется выброс с параметром сдвига λ . Видно, что при фиксированных уровне значимости α и объёме выборки n величина отклонения δ_n является убывающей функцией параметра сдвига λ . Если не менять параметр сдвига λ и уровень значимости α , то увеличение числа наблюдений n будет приводить к росту отклонения δ_n .

Из рис. 2 видно также, что выборка объёма $n = 20$ может считаться более предпочтительной для исследователя, чем, например, объёма $n = 50$ или $n = 100$; при переходе от $n = 20$ к $n \geq 50$ мера P_1 меняется незначительно (см. рис. 1), однако отклонения от нее мер P_2-P_4 могут существенно вырасти.

Заключение. Получено новое представление для закона распределения студентизированного отклонения выброса от среднего в нормально распределенной выборке, основанное на использовании функции Эрмита с отрицательным целым значком. Это представление было использовано для получения интегральных представлений мер мощности Дэйвида P_1-P_4 критерия Граббса в случае, когда статистика критерия является статистикой отношения правдоподобия. При этом нулевой гипотезой H_0 служило предположение о том, что выборка из n наблюдений случайно извлечена из нормальной $N(a; \sigma^2)$ генеральной совокупности. Конкурирующая гипотеза H_1 состояла в том, что в выборке имеется одно anomальное наблюдение X_{out} с распределением $N(a + \lambda\sigma; \sigma^2)$. Определена область критических значений статистики Граббса, в которой меры мощности Дэйвида P_2-P_4 , предназначенные для обнаружения выброса, совпадают. Получена формула для вычисления максимально возможного отклонения мер мощности Дэйвида, предназначенных для обнаружения выброса, от классической функции мощности критерия. Выполнены модельные расчёты классической функции мощности критерия для случая нормально распределенной выборки с выбросом, отличающимся от остальных наблюдений сдвигом вправо. Результаты вычислений оказались близки к теоретически ожидаемым.

Автор выражает искреннюю благодарность профессору Олегу Александровичу Репину за помощь и поддержку, оказанные им при работе со специальными функциями.

БИБЛИОГРАФИЧЕСКИЙ СПИСОК

1. David H. A. Order Statistics: 2nd ed. / Wiley Series in Probability and Mathematical Statistics. New York: John Wiley & Sons, 1981. xiii+360 pp.; русск. пер.: Дэвид Г. Порядковые статистики. М.: Наука, 1979. 336 с.
2. Grubbs F. E. Sample criteria for testing outlying observations // *Ann. Math. Statistics*, 1950. Vol. 21, no. 1. Pp. 27–58.
3. Ширяева Л. К. Вычисление мер мощности критерия Граббса проверки на один выброс // *Сиб. журн. индустр. матем.*, 2010. Т. 13, № 4. С. 141–154. [Shiryayeva L. K. Calculation of power measures of the Grubbs test for one outlier // *Sib. Zh. Ind. Mat.*, 2010. Vol. 13, no. 4. Pp. 141–154].
4. Zhang J., Keming Y. The null distribution of the likelihood-ratio test for one or two outliers in a normal sample // *Test*, 2006. Vol. 15, no. 1. Pp. 141–150.
5. Barnett V., Lewis T. Outliers in statistical data: 3rd ed. / Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics. Chichester: John Wiley & Sons, 1994. xviii+584 pp.
6. Айвазян С. А., Мхитарян В. С. Теория вероятностей и прикладная статистика. Т. 1. М.:

Юнити-Дана, 2001. 656 с. [*Aivazian S. A., Mkhitarian V. S. Applied Statistics and Essentials of Econometrics. Vol. 1. Moscow: Yuniti-Dana, 2001. 656 pp.*]

7. *Лебедев Н. Н. Специальные функции и их приложения. М.: Физ.-мат. лит., 1963. 358 с. [Lebedev N. N. Special functions and their applications. Moscow: Fiz.-Mat. Lit., 1963. 358 pp.]*

Поступила в редакцию 24/VI/2012;
в окончательном варианте — 07/IX/2012.

MSC: 56U78

THE USE OF HERMITE SPECIAL FUNCTIONS FOR INVESTIGATION OF POWER PROPERTIES OF GRABBS STATISTICS

L. K. Shiryaeva

Samara State Economic University,
141, Sovetskoy Armii st., Samara, 443090, Russia.

E-mail: shiryaeva_lk@mail.ru

We consider a normal sample with a single upper outlier. A distribution of studentized form of outlier's deviation from the sample mean is obtained. This distribution uses Hermite special functions with negative integer-valued index. The integral relationships for David's power measures of Grubbs criteria are obtained. We discuss the case, when Grubbs statistic is the likelihood-ratio statistic. We find the maximal deviation of power function for Grubbs criteria from the probability that the contaminant is the outlier and it is identified as discordant. We receive the region of critical values of Grubbs statistic, where the second power measure of David equals to the third and fourth power measures of David. We make calculations of power function for Grubbs criteria in the case of normal samples with a single upper outlier with the right shift. The results of calculations are similar to the theoretically expected facts.

Key words: *outlier, Grubbs statistics, Hermite special function, David's power measures for Grubbs criterion, normal distribution law.*

Original article submitted 24/VI/2012;
revision submitted 07/IX/2012.