

1. Параметрические критерии однородности для двух выборок

1.1. Распределение Стьюдента

Пусть ξ, ξ_1, \dots, ξ_n независимые нормально распределенные величины $\mathcal{N}(0, \sigma)$. Случайная величина вида

$$\zeta = \frac{\xi}{\sqrt{\frac{1}{n} \sum_{i=1}^n \xi_i^2}} = \frac{\xi}{\eta}, \quad (1)$$

имеет распределение *Стьюдента с n степенями свободы*.

Не умаляя общности, будем считать $\sigma = 1$. Тогда $\sum_{i=1}^n \xi_i^2$ имеет распределение $\chi^2(n)$ с плотностью

$$\kappa_n(x) = \frac{2^{-\frac{n}{2}}}{\Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}. \quad (2)$$

$$F_{a\xi}(x) = P\{a\xi < x\} = P\left\{\xi < \frac{x}{a}\right\} = F_\xi\left(\frac{x}{a}\right),$$

$$F_{\sqrt{\xi}}(x) = P\left\{\sqrt{\xi} < x\right\} = P\{\xi < x^2\} = F_\xi(x^2), \implies$$

$$f_{a\xi}(x) = \frac{1}{a} f_\xi\left(\frac{x}{a}\right), \quad f_{\sqrt{\xi}}(x) = 2x f_\xi(x^2).$$

Для величины $\frac{1}{n} \sum_{i=1}^n \xi_i^2$ получаем плотность $n\kappa_n(nx)$, а для величины $\eta = \sqrt{\frac{1}{n} \sum_{i=1}^n \xi_i^2}$ плотность $2nx\kappa_n(nx^2)$.

Совместное распределение ξ и η имеет вид:

$$\frac{2n\eta}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} (n\eta^2)^{\frac{n}{2}-1} e^{-\frac{n\eta^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{\xi^2}{2}} = c_n \eta^{n-1} e^{-\frac{\xi^2 + n\eta^2}{2}},$$

$$\text{где } c_n = \frac{\sqrt{2} \left(\frac{n}{2}\right)^{\frac{n}{2}}}{\sqrt{\pi} \Gamma\left(\frac{n}{2}\right)}.$$

$$\begin{aligned} P\{\zeta < x\} &= P\left\{\frac{\xi}{\eta} < x\right\} = P\{\xi < x\eta\} = \\ &= c_n \int \int_{\substack{\eta > 0 \\ \xi < x\eta}} \eta^{n-1} e^{-\frac{\xi^2 + n\eta^2}{2}} d\xi d\eta = c_n \int_{-\infty}^x du \int_0^{\infty} v^{n-1} e^{-\frac{u^2 v^2 + n v^2}{2}} v dv, \end{aligned}$$

используем замену $\xi = uv, \eta = v$ с якобианом преобразования

$$J = \begin{vmatrix} \frac{\partial \xi}{\partial u} & \frac{\partial \xi}{\partial v} \\ \frac{\partial \eta}{\partial u} & \frac{\partial \eta}{\partial v} \end{vmatrix} = \begin{vmatrix} v & u \\ 0 & 1 \end{vmatrix} = v.$$

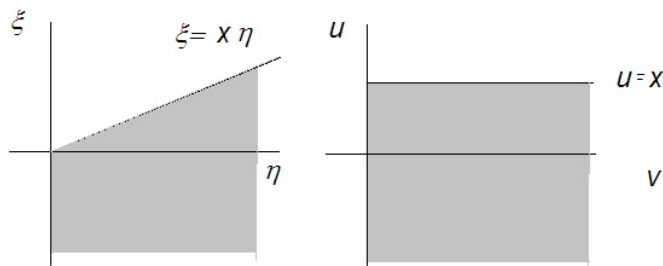


Рис. 1. Замена переменных $\xi = uv$, $\eta = v$ ($\eta > 0, \xi < x\eta$).

Далее заменим $\frac{(u^2+n)v^2}{2} = t$, $v = \frac{\sqrt{2t}}{\sqrt{u^2+n}}$, $vdv = \frac{dt}{u^2+n}$.

$$\begin{aligned} P\{\zeta < x\} &= c_n \int_{-\infty}^x du \int_0^{\infty} \left(\frac{\sqrt{2t}}{\sqrt{u^2+n}} \right)^{n-1} e^{-t} \frac{dt}{u^2+n} = \\ &= c_n 2^{\frac{n-1}{2}} \int_{-\infty}^x \frac{du}{(u^2+n)^{\frac{n+1}{2}}} \Gamma\left(\frac{n+1}{2}\right) = \\ &= \frac{\sqrt{2} \left(\frac{n}{2}\right)^{\frac{n}{2}}}{\sqrt{\pi} \Gamma\left(\frac{n}{2}\right)} 2^{\frac{n-1}{2}} \int_{-\infty}^x \frac{du}{(u^2+n)^{\frac{n+1}{2}}} \Gamma\left(\frac{n+1}{2}\right) \end{aligned}$$

Отсюда плотность распределения имеет вид $f_{\zeta}(x, n) =$

$$= \frac{n^{\frac{n}{2} + \frac{1}{2}} \Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\pi n} \Gamma\left(\frac{n}{2}\right)} \frac{1}{(x^2+n)^{\frac{n+1}{2}}} = \frac{1}{\sqrt{\pi n}} \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \left(\frac{x^2}{n} + 1\right)^{-\frac{n+1}{2}}.$$

1.2. Асимптотическое свойство распределения Стюдента

Плотность распределения Стюдента

$$f_{\zeta}(x, n) = \frac{1}{\sqrt{\pi n}} \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \left(\frac{x^2}{n} + 1\right)^{-\frac{n+1}{2}}. \quad (3)$$

можно записать в виде

$$f_{\zeta}(x, n) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)} \frac{1}{\sqrt{2\pi}} \left(\frac{x^2}{n} + 1\right)^{-\frac{n+1}{2}}.$$

Воспользуемся формулой Стирлинга

$$\ln \Gamma(\lambda) = \left(\lambda - \frac{1}{2}\right) \ln \lambda - \lambda + \frac{1}{2} \ln(2\pi) + \frac{1}{12\lambda} + O\left(\frac{1}{\lambda^3}\right).$$

При $p = \frac{n}{2} \rightarrow \infty$ имеют место выражения

$$\begin{aligned} & \ln \Gamma \left(p + \frac{1}{2} \right) - \frac{1}{2} \ln p - \ln \Gamma(p) \approx \\ & \approx \left(p + \frac{1}{2} - \frac{1}{2} \right) \ln \left(p + \frac{1}{2} \right) - p - \frac{1}{2} - \\ & \quad - \frac{1}{2} \ln p - \left(p - \frac{1}{2} \right) \ln p + p = \\ & = p \ln \left(p + \frac{1}{2} \right) - \frac{1}{2} - p \ln p = p \ln \left(1 + \frac{1}{2p} \right) - \frac{1}{2} \approx 0. \end{aligned}$$

При любом фиксированном x имеем

$$-\frac{n+1}{2} \ln \left(1 + \frac{x^2}{n} \right) \xrightarrow{n \rightarrow \infty} -\frac{x^2}{2},$$

откуда

$$f_{\zeta}(x, n) \xrightarrow{n \rightarrow \infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}.$$

1.3. Распределение Фишера

Распределенная по Фишеру случайная величина имеет вид

$$\zeta = \frac{\frac{1}{n} \sum_{i=1}^n \xi_i^2}{\frac{1}{m} \sum_{j=1}^m \eta_j^2} = \frac{m}{n} \zeta_0, \quad (4)$$

где ξ_i, η_j независимы, одинаково распределены $\mathcal{N}(0, \sigma)$. Не умаляя общности, будем считать $\sigma^2 = 1$. Тогда $\xi = \sum_{i=1}^n \xi_i^2$ и $\eta = \sum_{j=1}^m \eta_j^2$ имеют распределение (2) хи-квадрат.

$$P\{\zeta_0 < x\} = P\left\{ \frac{\xi}{\eta} < x \right\} = P\{\xi < x\eta\}.$$

Совместное распределение величин ξ и η имеет вид

$$\begin{aligned} \kappa_n(\xi) \kappa_m(\eta) &= \frac{2^{-\frac{n}{2}}}{\Gamma(\frac{n}{2})} \xi^{\frac{n}{2}-1} e^{-\frac{\xi}{2}} \frac{2^{-\frac{m}{2}}}{\Gamma(\frac{m}{2})} \eta^{\frac{m}{2}-1} e^{-\frac{\eta}{2}} = \\ &= a_{nm} \xi^{\frac{n}{2}-1} \eta^{\frac{m}{2}-1} e^{-\frac{\xi+\eta}{2}}, \\ & \text{где } a_{nm} = \frac{2^{-\frac{m+n}{2}}}{\Gamma(\frac{n}{2}) \Gamma(\frac{m}{2})}. \end{aligned}$$

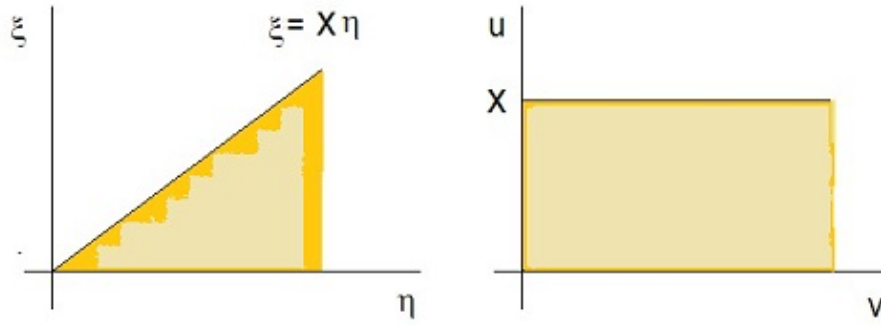


Рис. 2. Замена переменных $\xi = uv$, $\eta = v$.

$$P\{\zeta_0 < x\} = a_{nm} \int_{\eta > 0} \int_{0 < \xi < x\eta} \xi^{\frac{n}{2}-1} \eta^{\frac{m}{2}-1} e^{-\frac{\xi+\eta}{2}} d\xi d\eta. \quad (5)$$

Используем замену $\xi = uv$, $\eta = v$ с якобианом преобразования

$$J = \begin{vmatrix} \frac{\partial \xi}{\partial u} & \frac{\partial \xi}{\partial v} \\ \frac{\partial \eta}{\partial u} & \frac{\partial \eta}{\partial v} \end{vmatrix} = \begin{vmatrix} v & u \\ 0 & 1 \end{vmatrix} = v.$$

Область интегрирования из $\eta > 0$, $0 < \xi < x\eta$ перейдет в область $v > 0$, $0 < u < x$ (рис.2). Следовательно,

$$\begin{aligned} P\{\zeta_0 < x\} &= a_{nm} \int_0^x u^{\frac{n}{2}-1} du \int_0^\infty v^{\frac{n+m}{2}-2} e^{-\frac{uv+v}{2}} v dv = \\ &= a_{nm} \int_0^x u^{\frac{n}{2}-1} du \int_0^\infty v^{\frac{n+m}{2}-1} e^{-\frac{(u+1)v}{2}} dv. \end{aligned}$$

Заменим переменные $t = \frac{u+1}{2}v$, $v = \frac{2}{u+1}t$, $dv = \frac{2}{u+1}dt$.

$$\begin{aligned} a_{nm} \int_0^x u^{\frac{n}{2}-1} du \int_0^\infty t^{\frac{n+m}{2}-1} \left(\frac{2}{u+1}\right)^{\frac{n+m}{2}-1} \frac{2}{u+1} e^{-t} dt &= \\ = \frac{2^{-\frac{m+n}{2}}}{\Gamma(\frac{n}{2})\Gamma(\frac{m}{2})} \int_0^x u^{\frac{n}{2}-1} \left(\frac{2}{u+1}\right)^{\frac{n+m}{2}} du \Gamma\left(\frac{n+m}{2}\right) &= \\ = \frac{\Gamma\left(\frac{n+m}{2}\right)}{\Gamma(\frac{n}{2})\Gamma(\frac{m}{2})} \int_0^x \frac{u^{\frac{n}{2}-1}}{(u+1)^{\frac{n+m}{2}}} du. \end{aligned}$$

Таким образом, используя $C_{nm} = \frac{\Gamma\left(\frac{n+m}{2}\right)}{\Gamma(\frac{n}{2})\Gamma(\frac{m}{2})}$, получаем

$$\begin{aligned} f_0(x) &= C_{nm} \frac{x^{\frac{n}{2}-1}}{(x+1)^{\frac{n+m}{2}}}, \\ f(x) &= \frac{n}{m} f_0\left(\frac{n}{m}x\right) = \frac{n}{m} C_{nm} \frac{\left(\frac{n}{m}x\right)^{\frac{n}{2}-1}}{\left(\frac{n}{m}x+1\right)^{\frac{n+m}{2}}}. \end{aligned}$$

1.4. Совместное распределение \bar{x} и m_2

Лемма 1. (Фишера) Пусть x_1, \dots, x_n — независимые и нормально распределенные $\mathcal{N}(0, \sigma)$ компоненты вектора X . В результате ортогонального преобразования $Y = CX$ имеем также независимые и нормально распределенные $\mathcal{N}(0, \sigma)$ компоненты y_1, \dots, y_n . Тогда квадратичная форма

$$Q(X) = \sum_{i=1}^n x_i^2 - y_1^2 - \dots - y_p^2$$

независима от y_1, \dots, y_p и имеет плотность $\frac{1}{\sigma^2} \kappa_{n-p} \left(\frac{x}{\sigma^2} \right)$, где $\kappa_n(x)$ — плотность распределения хи-квадрат с n степенями свободы (2).

Теорема 1. Пусть $x_1, \dots, x_n \sim \mathbf{N}(\mu, \sigma)$, \bar{x} и m_2 — выборочные среднее и второй центральный момент. Тогда

1. \bar{x} и m_2 независимы;
2. $\bar{x} \sim \mathbf{N}(\mu; \frac{\sigma}{\sqrt{n}})$;
3. статистика $\frac{nm_2}{\sigma^2} \sim \chi^2(\mathbf{n} - 1)$ имеет распределение хи-квадрат с $(n - 1)$ степенью свободы.

Без ограничения общности будем считать $\mu = 0$.

$$nm_2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2.$$

Выражение $n\bar{x}^2 = \left(\frac{x_1}{\sqrt{n}} + \dots + \frac{x_n}{\sqrt{n}} \right)^2$ есть квадрат линейной формы $c_1x_1 + \dots + c_nx_n$, для которой $c_1^2 + \dots + c_n^2 = 1$. Поэтому применим лемму 1, положив $p = 1$ и $y_1 = \sqrt{n}\bar{x}$.

1.5. Статистики Стьюдента и Фишера

Пусть x_1, \dots, x_n — выборка из нормальной совокупности со средним μ и дисперсией σ^2 . Тогда отношение

$$\tau = \frac{\bar{x} - \mu}{\sqrt{m_2}} \sqrt{n-1} \quad (6)$$

имеет распределение Стьюдента с $n-1$ степенью свободы.

Действительно, в (1) полагая $\xi_1 = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$ и $\sum_{i=1}^n \xi_i^2 = \frac{nm_2}{\sigma^2}$, получим

$$\tau = \frac{\xi_1}{\sqrt{\frac{1}{n-1} \sum_{k=2}^n \xi_k^2}} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}} \sqrt{\frac{1}{n-1} \frac{nm_2}{\sigma^2}}} = \frac{\bar{x} - \mu}{\sqrt{m_2}} \sqrt{n-1}.$$

Так как $m_2 = \frac{n-1}{n} S^2$, $\sqrt{m_2} = \frac{\sqrt{n-1}}{\sqrt{n}} S$,

$$\tau = \frac{\bar{x} - \mu}{S} \sqrt{n} \sim \mathbf{T}(n-1). \quad (7)$$

Статистика Фишера для проверки гипотезы о равенстве дисперсий строится с учетом структуры распределения Фишера (4) как функция от двух независимых нормально распределенных выборок: $x_1, \dots, x_n \sim \mathcal{N}(\mu_1; \sigma)$ и $y_1, \dots, y_m \sim \mathcal{N}(\mu_2; \sigma)$

$$F = \frac{S_x^2}{S_y^2} \sim \mathbf{F}(n; m), \quad (8)$$

где S_x^2 и S_y^2 — несмещенные оценки их дисперсий соответствующих выборок. Действительно, так как $(n-1)S^2/\sigma^2 \sim \chi^2(n-1)$, то в соответствии с (4)

$$F = \frac{\frac{1}{n-1}(n-1)S_x^2/\sigma^2}{\frac{1}{m-1}(m-1)S_y^2/\sigma^2} = \frac{S_x^2}{S_y^2} \sim \mathbf{F}(n; m).$$

1.6. Критерий Стьюдента для зависимых выборок

Часто встречаются данные, в которых за одними и теми же индивидами наблюдают в разные моменты времени, например, у $n = 67$ больных с воспалением легких логарифм показателя скорости оседания эритроцитов (СОЭ) в начале заболевания в среднем равен 3.82 со стандартным отклонением 0.13, а при выписке из стационара 3.28 со стандартным отклонением 0.19. Требуется выяснить, можно ли считать уменьшение СОЭ значимым, или это уменьшение можно объяснить случайностью.

Итак, имеется парная выборка (x_i, y_i) , $i = 1, \dots, n$, зависимых наблюдений объема n над случайными величинами ξ_1 и ξ_2 с нормальным распределением $\mathcal{N}(\mu_1, \sigma_1)$ и $\mathcal{N}(\mu_2, \sigma_2)$, относительно которых нет предположения о независимости — скорее всего, чем меньше был уровень СОЭ в начале заболевания, тем меньше он при выздоровлении. Рассмотрим разность $\eta = \xi_1 - \xi_2$ и соответствующую выборку $z_i = x_i - y_i$ с выборочным средним \bar{z} и несмещенной оценкой дисперсии S_z^2 . Если справедлива нулевая гипотеза $H_0 : \mu_1 = \mu_2$, то для случайной выборки $z_i = x_i - y_i$ должна быть справедлива гипотеза о равенстве нулю ее математического ожидания.

$$t = \frac{\bar{z}}{S_z} \sqrt{n} \sim \mathbf{T}(n - 1) \text{ при справедливости } H_0 : \mu_1 = \mu_2.$$

. Например, $\bar{z} = 0.54$, $S = 0.19$, $t = \frac{\bar{z}-0}{S} \sqrt{n} = \frac{0.54}{0.19} \sqrt{66} = 23.45 = t_*$, двусторонний доверительный уровень вероятности равен

$$p = P\{|t| > |t_*|\} = \text{СТЮДРАСП}(23.45; 66; 2) = 10^{-33} < 0.05,$$

поэтому гипотеза $H_0 : \mu_1 = \mu_2$ отвергается,¹ и уменьшение уровня СОЭ в процессе выздоровления нельзя объяснить случайностью.

¹ Значимость критерия в R считается как $p = 2 * (1 - pt(t_*, n - 1))$ или используется общая функция $t.test$ с опцией $paired = TRUE$.

1.7. Распределение разности выборочных средних

Пусть имеется две независимые выборки x_1, \dots, x_{n_1} и y_1, \dots, y_{n_2} , о которых известно, что они имеют нормальное распределение со средними μ_1 и μ_2 и одинаковой дисперсией σ^2 . Выборочные средние обозначим через \bar{x} и \bar{y} , а несмещенные оценки дисперсии через S_1^2 и S_2^2 . Тогда статистика

$$T = \frac{(\bar{x} - \bar{y} - (\mu_1 - \mu_2))\sqrt{n_1 + n_2 - 2}}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}\sqrt{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}} \sim \mathbf{T}(n_1 + n_2 - 2) \quad (9)$$

имеет распределение Стьюдента с числом степеней свободы $n_1 + n_2 - 2$.

Во-первых, из независимости выборок выборочные средние \bar{x} и \bar{y} также независимы, а по теореме Фишера (т.1) нормально распределены соответственно с параметрами $\mathcal{N}(\mu_1, \sigma/\sqrt{n_1})$ и $\mathcal{N}(\mu_2, \sigma/\sqrt{n_2})$. Следовательно, их разность имеет нормальное распределение с параметрами

$$\bar{x} - \bar{y} \sim \mathcal{N}\left(\mu_1 - \mu_2, \sigma\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}\right),$$

откуда получаем

$$\xi = \frac{\bar{x} - \bar{y} - (\mu_1 - \mu_2)}{\sigma\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim \mathcal{N}(0; 1). \quad (10)$$

С другой стороны, также по теореме 1

$$\begin{aligned} (n_1 - 1)S_1^2/\sigma^2 &\sim \chi^2(\mathbf{n}_1 - \mathbf{1}), \\ (n_2 - 1)S_2^2/\sigma^2 &\sim \chi^2(\mathbf{n}_2 - \mathbf{1}). \end{aligned}$$

Из аддитивности распределения χ^2 их сумма вида

$$\eta = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{\sigma^2} \sim \chi^2(\mathbf{n}_1 + \mathbf{n}_2 - \mathbf{2}) \quad (11)$$

имеет распределение хи-квадрат с числом степеней свободы, равным сумме степеней свободы слагаемых. Рассмотрим случайную величину вида

$$T = \frac{\xi}{\sqrt{\frac{1}{n_1 + n_2 - 2}\eta}},$$

которая, согласно структуре (1), имеет распределение Стьюдента с числом степеней свободы, равным $n_1 + n_2 - 2$. Подставив соответствующие ξ из (10) и η из (11), получим вид статистики из (9).

При известных дисперсиях σ_1^2 и σ_2^2 (возможно неодинаковых), учитывая нормальность средних

$$\bar{x} \sim \mathcal{N}(\mu_1, \sigma_1/\sqrt{n_1}), \quad \bar{y} \sim \mathcal{N}(\mu_2, \sigma_2/\sqrt{n_2}),$$

вычисляем дисперсию разности

$$\mathbf{D}(\bar{x} - \bar{y}) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2},$$

откуда получаем

$$Z = \frac{\bar{x} - \bar{y} - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim \mathcal{N}(0; 1). \quad (12)$$

В случае неизвестных неодинаковых дисперсий $\sigma_1^2 \neq \sigma_2^2$ статистика T имеет приближенно распределение Стьюдента с нецелым числом степеней свободы.

$$T = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \sim \mathbf{T}(k), \quad \text{где } k = \frac{(S_1^2/n_1 + S_2^2/n_2)^2}{\frac{(S_1^2/n_1)^2}{n_1-1} + \frac{(S_2^2/n_2)^2}{n_2-1}}. \quad (13)$$