

# Некоторые задачи проверки статистических гипотез и их применения

Выпускная квалификационная работа

Лунев Иван Сергеевич, группа 422

Санкт-Петербургский Государственный Университет  
Прикладная математика и информатика  
Вычислительная стохастика и статистические модели

Научный руководитель: к.ф.-м.н., доцент В. В. Некруткин  
Рецензент: к.ф.-м.н., доцент А. И. Коробейников



Санкт-Петербург  
2018 г.

Две задачи — проверка равенства нулю коэффициента корреляции Пирсона и частного коэффициента корреляции.

Классические критерии — в предположении гауссовости выборок.

Пакеты Statistica, SPSS, R и т. д.



## Проблемы:

- ❶ Когда классические критерии применимы, если распределения не гауссовские?
- ❷ Построение критериев, применимых в общей ситуации.

Если вектор  $(x, y)$  имеет невырожденное двумерное гауссовское распределение, критерий хорошо известен [Robb J. Muirhead, 2009]:

Нулевая гипотеза:  $H_0 : \{\rho = 0\}$ .

Статистика:

$$\mu(\mathbf{X}_n, \mathbf{Y}_n) := \sqrt{n-2} \frac{\hat{\rho}_n}{\sqrt{1 - \hat{\rho}_n^2}},$$

где  $\hat{\rho}_n$  — выборочный коэффициент корреляции.

Распределение статистики:

$$\mathcal{L}_{H_0}(\mu(\mathbf{X}_n, \mathbf{Y}_n)) = t(n-2).$$

## Классический критерий

Получаем критерий с уровнем значимости  $\alpha$ , который отвергает нулевую гипотезу, если  $|\mu(\mathbf{X}_n, \mathbf{Y}_n)| \geq T_{n-2, 1-\alpha/2}$ , где  $T_{m, \gamma}$  —  $\gamma$ -квантиль распределения Стьюдента с  $m$  степенями свободы.

$(\mathbf{X}_n, \mathbf{Y}_n)$  – двумерная независимая повторная выборка, соответствующая вектору  $(x, y)$  с непрерывными координатами, причем четвертые моменты  $x$  и  $y$  конечны. Обозначим

$$x^* = (x - \mathbb{E}x)/\sigma_x, \quad y^* = (y - \mathbb{E}y)/\sigma_y,$$

$$\sigma_{f,\rho}^2 = \mathbb{D} \left( \rho(x^{*2} + y^{*2}) - 2x^*y^* \right) / 4,$$

$$\mu(\mathbf{X}_n, \mathbf{Y}_n) = \frac{\sqrt{n}}{\sigma_{f,\rho}} (\hat{\rho}_n - \rho),$$

$\hat{\rho}_n$  — выборочный коэффициент корреляции.

Утверждение (ход доказательства, например, в [E. Lehman, 1999]).

Если  $|\rho| \neq 1$  и  $n \rightarrow \infty$

$$\mathcal{L}(\mu(\mathbf{X}_n, \mathbf{Y}_n)) \Rightarrow N(0, 1).$$

Для гауссовской выборки  $\sigma_{f,\rho}^2 = (1 - \rho^2)^2$ .

Нулевая гипотеза:  $H_0 : \{\rho = 0\}$ .

Статистика:

$$\mu(\mathbf{X}_n, \mathbf{Y}_n) = \frac{\sqrt{n}}{\widehat{\sigma_f}} \widehat{\rho}_n, \quad \text{где } \widehat{\sigma_f}^2 \text{ — состоятельная оценка } \sigma_{f,0}^2.$$

## Критерий

Пусть  $\alpha \in (0, 1)$ . Критерий отвергает гипотезу  $H_0$ , если

$$|\mu(\mathbf{X}_n, \mathbf{Y}_n)| \geq C_{1-\alpha/2},$$

где  $C_\beta$  — квантиль уровня  $\beta$  распределения  $N(0, 1)$ .

«Общий» критерий:

$$\hat{\sigma}_{f,\rho}^2 = \frac{1}{n} \sum_{i=1}^n \left( \hat{\rho}_n \left( \left( \frac{x_i - \bar{x}}{\bar{s}_n(x)} \right)^2 + \left( \frac{y_i - \bar{y}}{\bar{s}_n(y)} \right)^2 \right) - 2 \frac{x_i - \bar{x}}{\bar{s}_n(x)} \frac{y_i - \bar{y}}{\bar{s}_n(y)} \right)^2 / 4$$

— оценка состоятельна для любого  $\rho$  (генеральные моменты заменены на выборочные).

«Модифицированный» критерий:

$$\hat{\sigma}_{f,0}^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 (y_i - \bar{y})^2}{n \bar{s}_n^2(x) \bar{s}_n^2(y)}$$

— оценка состоятельна при  $\rho = 0$ .

## Утверждение

Оба критерия асимптотически точны и состоятельны при  $n \rightarrow \infty$ .

Предположим, что  $\rho = 0$ . Рассмотрим предельную дисперсию

$$\sigma_{f,0}^2 = \mathbb{D}(x^* y^*),$$

где  $x^* = (x - \mathbb{E}x)/\sigma_x$ ,  $y^* = (y - \mathbb{E}y)/\sigma_y$ .

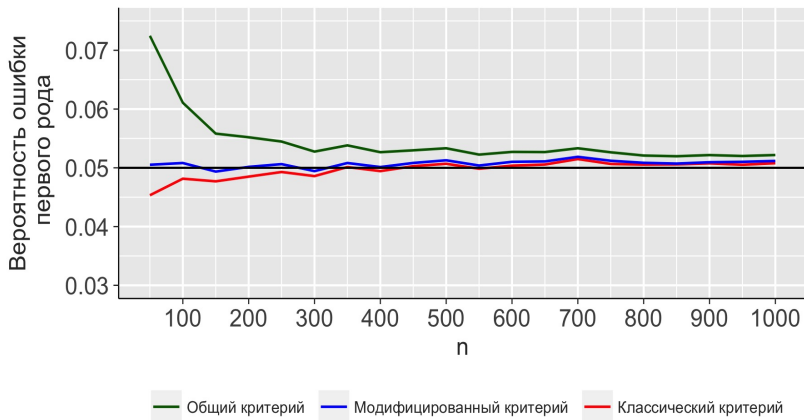
## Утверждение

Классический критерий, примененный к негауссовской выборке, будет асимптотически точным и состоятельным тогда и только тогда, когда  $\sigma_{f,0}^2 = 1$ , что эквивалентно равенству  $\rho(x^2, y^2) = 0$ .

## Следствие

Если случайные величины  $x$  и  $y$  независимы, то  $\sigma_{f,0}^2 = 1$ .

$\rho = 0, \sigma_{f,0}^2 = 1$ : классический критерий асимптотически точен.



**Рис.:** Равномерно распределенные независимые случайные величины. Оценки вероятностей ошибок первого рода для трех критериев при  $\alpha = 0.05, N = 10^5$



$\rho = 0$ ,  $\sigma_{f,0}^2 = 2/3$ : классический критерий асимптотически консервативен.

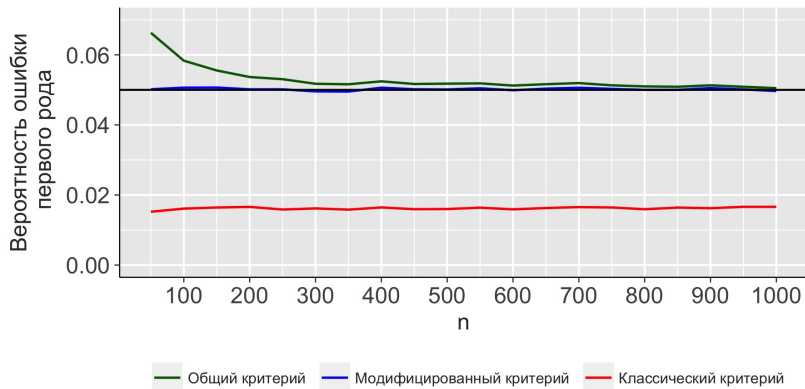
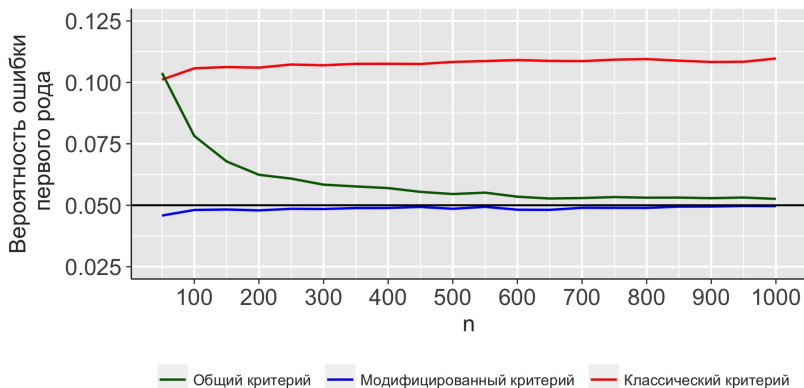


Рис.: Равномерное распределение в круге. Оценки вероятностей ошибок первого рода для трех критериев при  $\alpha = 0.05$ ,  $N = 10^5$

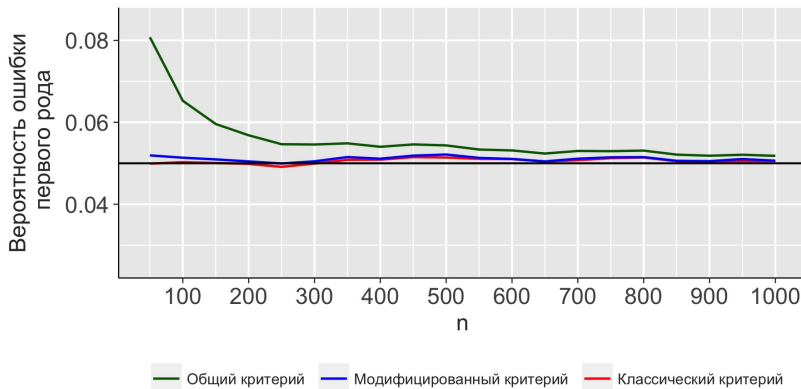
$\rho = 0$ ,  $\sigma_{f,0}^2 = 1.5$ : классический критерий асимптотически либерален.



**Рис.:** Круговая симметрия с  $r^2 \in G(1/2, 1)$ . Оценки вероятностей ошибок первого рода для трех критериев при  $\alpha = 0.05$ ,  $N = 10^5$

*Равномерная смесь равномерных распределений на трех окружностях с радиусами 1,  $\sqrt{2}$  и  $\sqrt{6 + \sqrt{39}}$ .*

Зависимость,  $\rho = 0$ ,  $\sigma_{f,0}^2 = 1$ : классический критерий асимптот. точен.



**Рис.:** Смесь равномерных распределений на трех окружностях. Оценки вероятностей ошибок первого рода для трех критериев при  $\alpha = 0.05$ ;  $N = 10^5$ .

Случайные величины  $x, y, z$  с конечными вторыми моментами,  $|\rho_{x,z}| < 1$ ,  $|\rho_{y,z}| < 1$ .

Линейные регрессии  $x$  на  $z$  и  $y$  на  $z$  с остатками  $\epsilon_1$  и  $\epsilon_2$ .

$$\rho_{x,y|z} = \rho(\epsilon_1, \epsilon_2) = \frac{\rho_{x,y} - \rho_{x,z}\rho_{y,z}}{\sqrt{1 - \rho_{x,z}^2}\sqrt{1 - \rho_{y,z}^2}} :$$

частный коэффициент корреляции между  $x, y$  при исключении влияния  $z$ .

$(X_n, Y_n, Z_n)$  – трехмерная повторная выборка объема  $n$ .

Нулевая гипотеза

$$H_0 : \{\rho_{x,y|z} = 0\}.$$

$(x, y, z)$  — невырожденное гауссовское распределение: **классический критерий** (например, [ван дер Варден, 1960]).

Те же вопросы:

- 1 Когда классический критерий применим, если распределение не гауссовское?
- 2 Построение критериев, применимых в общей ситуации.

## Ответ

Так как частный коэффициент сводится к коэффициенту корреляции Пирсона между остатками  $\epsilon_1$  и  $\epsilon_2$ , то все результаты про коэффициент корреляции Пирсона переносятся на частный коэффициент корреляции.

$(\mathbf{X}_n, \mathbf{Y}_n, \mathbf{Z}_n)$  – независимая повторная выборка из распределения вектора  $(x, y, z)$  с непрерывными распределениями  $x, y, z$  и конечными четвертыми моментами. Пусть  $x^*, y^*$  и  $z^*$  – стандартизированные  $(x, y, z)$  и  $\rho_{x,z}^2 \neq 1$ ,  $\rho_{y,z}^2 \neq 1$ .

$$\sigma_{xy|z}^2 = \frac{1}{(1 - \rho_{x,z}^2)(1 - \rho_{y,z}^2)} \mathbb{D} \left( x^* y^* + z^{*2} \rho_{x,z} \rho_{y,z} - x^* z^* \rho_{y,z} - y^* z^* \rho_{x,z} \right),$$

где  $\hat{\rho}_{x,y|z}$  – выборочный коэффициент частной корреляции.

## Утверждение

Если  $\rho_{x,y|z} = 0$ , тогда

$$\mathcal{L}(\sqrt{n} \hat{\rho}_{x,y|z}) \Rightarrow N(0, \sigma_{x,y|z}^2).$$

## Замечание

Если вектор  $(x, y, z)$  имеет гауссовское распределение, то  $\sigma_{xy|z}^2 = 1$ .

Пусть  $\hat{\rho}_{x,z}$  и  $\hat{\rho}_{y,z}$  — выборочные коэффициенты корреляции между  $x, z$  и  $y, z$ . Кроме того, положим

$$\hat{u}_i^* = \frac{\frac{x_i - \bar{x}}{\bar{s}_n(x)} - \hat{\rho}_{x,z} \frac{z_i - \bar{z}}{\bar{s}_n(z)}}{\sqrt{1 - \hat{\rho}_{x,z}^2}}, \quad \hat{v}_i^* = \frac{\frac{y_i - \bar{y}}{\bar{s}_n(y)} - \hat{\rho}_{y,z} \frac{z_i - \bar{z}}{\bar{s}_n(z)}}{\sqrt{1 - \hat{\rho}_{y,z}^2}}.$$

Тогда при  $\rho_{xy|z} = 0$

$$\hat{\sigma}_{x,y|z}^2 = \frac{1}{n} \sum_{i=1}^n (\hat{u}_i^* \hat{v}_i^*)^2$$

— состоятельная оценка  $\sigma_{x,y|z}^2$ .

### Утверждение

Критерий, порожденный оценкой  $\hat{\sigma}_{x,y|z}^2$ , является асимптотически точным и состоятельным при  $n \rightarrow \infty$ .

Таким образом, в работе были получены следующие результаты.

## ❶ Для коэффициента корреляции Пирсона

- Предложены два асимптотического критерия для проверки гипотезы равенства нулю коэффициента корреляции Пирсона, доказаны их асимптотические точность и состоятельность;
- С помощью вычислительных экспериментов проведено сравнение этих критериев;
- Проанализированы свойства классического (гауссовского) критерия при его применении к негауссовской выборке.

## ❷ Для частного коэффициента корреляции

- Все теоретические результаты работы, относящиеся к коэффициенту корреляции Пирсона, перенесены на частный коэффициент корреляции.