

Санкт-Петербургский политехнический университет  
Петра Великого

Институт прикладной математики и механики  
Кафедра «Прикладная математика»

**Отчёт**  
**по лабораторным работам №5-8**  
**по дисциплине**  
**«Математическая статистика»**

Выполнил студент:  
Самутичев Евгений Романович  
группа: 3630102/70201

Проверил:  
к.ф.-м.н., доцент  
Баженов Александр Николаевич

Санкт-Петербург  
2020 г.

# Содержание

<b>1</b>	<b>Постановка задачи</b>	<b>3</b>
<b>2</b>	<b>Теория</b>	<b>4</b>
2.1	Двумерное нормальное распределение . . . . .	4
2.2	Ковариация и коэффициент корреляции . . . . .	4
2.3	Выборочные коэффициенты корреляции . . . . .	4
2.3.1	Пирсона . . . . .	4
2.3.2	Квадрантный . . . . .	4
2.3.3	Спирмена . . . . .	5
2.4	Эллипс равновероятности . . . . .	5
2.5	Простая линейная регрессия . . . . .	6
2.5.1	Критерий наименьших квадратов . . . . .	6
2.5.2	Критерий наименьших модулей . . . . .	6
2.6	Точечное оценивание . . . . .	7
2.6.1	Основные понятия . . . . .	7
2.6.2	Метод максимального правдоподобия . . . . .	7
2.7	Критерий согласия $\chi^2$ . . . . .	7
2.8	Интервальное оценивание . . . . .	8
2.9	Классические оценки . . . . .	9
2.9.1	Для математического ожидания $m$ . . . . .	9
2.9.2	Для среднего квадратичного отклонения $\sigma$ . . . . .	9
2.10	Асимптотически нормальные оценки . . . . .	9
2.10.1	Для математического ожидания $m$ . . . . .	9
2.10.2	Для среднего квадратичного отклонения $\sigma$ . . . . .	10
<b>3</b>	<b>Реализация</b>	<b>11</b>
<b>4</b>	<b>Результаты</b>	<b>12</b>
4.1	Коэффициенты корреляции . . . . .	12
4.2	Эллипсы равновероятности . . . . .	14
4.3	Выборка без выбросов . . . . .	19
4.4	Выборка с выбросами . . . . .	20
4.5	Критерий согласия $\chi^2$ . . . . .	21
4.6	Классические оценки . . . . .	21
4.7	Асимптотически нормальные оценки . . . . .	21
<b>5</b>	<b>Обсуждение</b>	<b>22</b>
5.1	Коэффициенты корреляции . . . . .	22
5.2	Эллипсы равновероятности . . . . .	22
5.3	Линейная регрессия . . . . .	22
5.4	Критерий согласия $\chi^2$ . . . . .	22
<b>6</b>	<b>Приложения</b>	<b>23</b>

## Список иллюстраций

1	$\rho = 0.0, n = 20$ . . . . .	14
2	$\rho = 0.0, n = 60$ . . . . .	14
3	$\rho = 0.0, n = 100$ . . . . .	15

4	$\rho = 0.5, n = 20$	15
5	$\rho = 0.5, n = 60$	16
6	$\rho = 0.5, n = 100$	16
7	$\rho = 0.9, n = 20$	17
8	$\rho = 0.9, n = 60$	17
9	$\rho = 0.9, n = 100$	18
10	Без выбросов	19
11	С выбросами	20

## Список таблиц

1	$\rho = 0$	12
2	$\rho = 0.5$	12
3	$\rho = 0.9$	13
4	Смесь нормальных распределений	13
5	Таблица вычислений $\chi^2$	21

# 1 Постановка задачи

1. Сгенерировать двумерные выборки размера 20, 60, 100 для нормального двумерного распределения  $N(x, y, 0, 0, 1, 1, \rho)$ . Коэффициент корреляции  $\rho$  взять равным 0, 0.5, 0.9. Каждая выборка генерируется 1000 раз и для неё вычисляются: среднее значение, среднее значение квадрата и дисперсия коэффициентов корреляции Пирсона, Спирмена и квадрантного коэффициента корреляции.

Повторить все вычисления для смеси нормальных распределений:

$$f(x, y) = 0.9 \cdot N(x, y, 0, 0, 1, 1, 0.9) + 0.1 \cdot N(x, y, 0, 0, 10, 10, -0.9)$$

. Изобразить сгенерированные точки на плоскости и нарисовать эллипс равновероятности.

2. Найти оценки коэффициентов  $a, b$  линейной регрессии  $y_i = a + bx_i + \varepsilon_i$ , используя 20 точек на отрезке  $[-1.8, 2]$  с равномерным шагом равным 0.2. Ошибку  $\varepsilon_i$  считать нормально распределённой с параметрами  $(0, 1)$ . В качестве эталонной зависимости взять  $y_i = 2 + 2x_i + e_i$ . При построении оценок коэффициентов использовать два критерия: критерий наименьших квадратов и критерий наименьших модулей. Прodelать то же самое для выборки, у которой в значения  $y_1$  и  $y_2$  вносятся возмущения 10 и -10.
3. Сгенерировать выборку объёмом 100 элементов для нормального распределения  $N(0, 1)$ . По сгенерированной выборке оценить параметры  $\mu$  и  $\sigma$  нормального закона методом максимального правдоподобия. В качестве основной гипотезы  $H_0$  будем считать, что сгенерированное распределение имеет вид  $N(\hat{\mu}, \hat{\sigma})$ . Проверить основную гипотезу, используя критерий согласия  $\chi^2$ . В качестве уровня значимости взять  $\alpha = 0.05$ . Привести таблицу вычислений  $\chi^2$ .

**Дополнительное исследование:** для проверки самого критерия, сгенерировать выборки объёма 20, 100 для нормального распределения  $U(-1, 1)$ , после чего проверить их на «нормальность».

4. Для двух выборок размерами 20 и 100 элементов, сгенерированных согласно нормальному закону  $N(0, 1)$ , для параметров положения и масштаба построить асимптотически нормальные интервальные оценки на основе точечных оценок метода максимального правдоподобия и классические интервальные оценки на основе статистик  $\chi^2$  и Стьюдента. В качестве доверительной вероятности взять  $\gamma = 0.95$ .

## 2 Теория

### 2.1 Двумерное нормальное распределение

Двумерная случайная величина  $(X, Y)$  называется *распределенной нормально* (или просто *нормальной*) если её плотность вероятности определена формулой

$$N(x, y, m_1, m_2, \sigma_1, \sigma_2, \rho) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} \left[\frac{(x-m_1)^2}{\sigma_1^2} - 2\rho\frac{(x-m_1)(y-m_2)}{\sigma_1\sigma_2} + \frac{(y-m_2)^2}{\sigma_2^2}\right]\right) \quad (1)$$

Можно показать [1, стр. 133-134] что компоненты  $X, Y$  двумерной нормальной случайной величины также распределены нормально с математическими ожиданиями  $m_X = m_1, m_Y = m_2$  и среднеквадратическими отклонениями  $\sigma_X = \sigma_1, \sigma_Y = \sigma_2$ . В свою очередь параметр  $\rho$  называют *коэффициентом корреляции*. Его значение будет раскрыто далее.

### 2.2 Ковариация и коэффициент корреляции

*Ковариацией* двух случайных величин  $X$  и  $Y$  называется величина:

$$K_{XY} = \mathbf{M}[(X - m_X)(Y - m_Y)] \quad (2)$$

В свою очередь *коэффициентом корреляции* называется

$$\rho_{XY} = \frac{K_{XY}}{\sigma_X\sigma_Y} \quad (3)$$

Коэффициент корреляции характеризует зависимость между случайными величинами  $X$  и  $Y$ . Именно его мы задаем в двумерном нормальном распределении как  $\rho$ . Если случайные величины  $X$  и  $Y$  независимы, то  $\rho_{XY} = 0$  т.к. в этом случае очевидно  $K_{XY} = 0$ .

### 2.3 Выборочные коэффициенты корреляции

#### 2.3.1 Пирсона

Пусть по выборке значений  $\{x_i, y_i\}_{i=1}^n$  двумерной случайной величины  $(X, Y)$ . Естественной оценкой для  $\rho_{XY}$  служит *выборочный коэффициент корреляции (Пирсона)*:

$$r = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (4)$$

Важным для приложений свойством является то что при данной оценке гипотеза  $\rho_{XY} \neq 0$  (о наличии зависимости между случайными) величинами может быть принята на уровне значимости 0.05 если выполнено:

$$|r|\sqrt{n-1} > 2.5 \quad (5)$$

это можно найти к примеру в [1, стр. 538]

#### 2.3.2 Квадрантный

*Выборочным квадрантным коэффициентом корреляции* называется величина:

$$r_Q = \frac{(n_1 + n_3) - (n_2 + n_4)}{n} \quad (6)$$

, где  $n_1, n_2, n_3, n_4$  - количества элементов выборки попавших соответственно в I, II, III и IV квадранты декартовой системы координат с центром в  $(\text{med } x, \text{med } y)$  и осями  $x_1 = x - \text{med } x, y_1 = y - \text{med } y$ , где  $\text{med}$  - выборочная медиана.

Формулу (6) можно переписать эквивалентным образом:

$$r_Q = \frac{1}{n} \sum_{i=1}^n \text{sign}(x_i - \text{med } x) \text{sign}(y_i - \text{med } y) \quad (7)$$

Важным свойством этой оценки является робастность. Её мы можем проверить используя схему засорения (смесь нормальных распределений).

### 2.3.3 Спирмена

На практике нередко требуется оценить степень взаимодействия между качественными признаками изучаемого объекта. Качественным называется признак, который нельзя измерить точно, но который позволяет сравнивать изучаемые объекты между собой и располагать их в порядке убывания или возрастания их качества. Для этого объекты выстраиваются в определённом порядке в соответствии с рассматриваемым признаком. Процесс упорядочения называется ранжированием, и каждому члену упорядоченной последовательности объектов присваивается ранг, или порядковый номер.

Например, объекту с наименьшим значением признака присваивается ранг 1, следующему за ним объекту — ранг 2, и т.д. Таким образом, происходит сравнение каждого объекта со всеми объектами изучаемой выборки. Если объект обладает не одним, а двумя качественными признаками — переменными  $X$  и  $Y$ , то для исследования их взаимосвязи используют выборочный коэффициент корреляции между двумя последовательностями рангов этих признаков.

Обозначим ранги, соответствующие значениям переменной  $X$ , через  $u$ , а ранги, соответствующие значениям переменной  $Y$  - через  $v$ . *Выборочный коэффициент ранговой корреляции Спирмена* определяется как выборочный коэффициент корреляции Пирсона между рангами  $u, v$  переменных  $X, Y$ :

$$r_S = \frac{\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(v_i - \bar{v})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})^2 \frac{1}{n} \sum_{i=1}^n (v_i - \bar{v})^2}} \quad (8)$$

## 2.4 Эллипс равновероятности

Рассмотрим выражение для плотности двумерного нормального распределения (1) несколько подробнее, а именно найдем линии уровня или что равносильно проекции сечения графика плотности плоскостями параллельными  $xOy$  на плоскость  $xOy$ :

$$N(x, y, m_1, m_2, \sigma_1, \sigma_2, \rho) = \text{const}$$

, или что равносильно:

$$\frac{(x - m_1)^2}{\sigma_1^2} - 2\rho \frac{(x - m_1)(y - m_2)}{\sigma_1 \sigma_2} + \frac{(y - m_2)^2}{\sigma_2^2} = \text{const} \quad (9)$$

Во всех точках каждого из таких эллипсов плотность двумерного нормального распределения  $N(x, y, m_1, m_2, \sigma_1, \sigma_2, \rho)$  постоянна. Поэтому они и называются *эллипсами равновероятности* [2, стр. 44-45]. Отметим что в предельном случае  $\rho = 1$ :

$$\left( \frac{x - m_1}{\sigma_1} - \frac{y - m_2}{\sigma_2} \right)^2 = \text{const}$$

, такое уравнение задает семейство прямых параллельных прямой:

$$\frac{x - m_1}{\sigma_1} = \frac{y - m_2}{\sigma_2} \quad (10)$$

Аналогично рассматривается предельный случай  $\rho = -1$ .

В данной работе, для выборки построенной по распределению  $N(x, y, m_1, m_2, \sigma_1, \sigma_2, \rho)$  эллипсы равновероятности строились таким образом чтобы покрыть все элементы выборки т.е. в качестве константы, стоящей в правой части уравнения (9) бралась:

$$R = \max_{\{(x_i, y_i)\}_{i=1}^n} \left( \frac{(x_i - m_1)^2}{\sigma_1^2} - 2\rho \frac{(x_i - m_1)(y_i - m_2)}{\sigma_1 \sigma_2} + \frac{(y_i - m_2)^2}{\sigma_2^2} \right) \quad (11)$$

## 2.5 Простая линейная регрессия

Регрессионную модель описания данных называют *простой линейной регрессией*, если

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (12)$$

, где  $\{x_i\}_{i=1}^n$  - значения фактора,  $\{y_i\}_{i=1}^n$  - наблюдаемые значения отклика, а  $\{\varepsilon_i\}_{i=1}^n$  - независимые, нормально распределенные по закону  $N(0, \sigma)$  случайные величины, а  $\beta_0, \beta_1$  - оцениваемые параметры [1, стр. 507]. Для оценки применяются различные методы, в данной работе рассмотрен следующий подход: вводится критерий рассогласования отклика и регрессионной функции, после чего оценки параметров регрессии выводятся из задачи минимизации критерия. Рассмотрим два таких критерия.

### 2.5.1 Критерий наименьших квадратов

Достаточно простые расчетные формулы для оценок получают при выборе критерия в виде суммы квадратов отклонений значений отклика от значений регрессионной функции:

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \rightarrow \min_{\beta_0, \beta_1} \quad (13)$$

Приведем сами расчетные формулы [1, стр. 509]:

$$\hat{\beta}_1 = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - (\bar{x})^2} \quad \hat{\beta}_0 = \bar{y} - \bar{x} \hat{\beta}_1 \quad (14)$$

Важным свойством является несмещенность оценки, однако она чувствительна к выбросам и если нужна робастная оценка, то следует рассмотреть следующий критерий.

### 2.5.2 Критерий наименьших модулей

В отличие от задач метода наименьших квадратов, для этого критерия минимизацию на практике проводят численно, решая:

$$M(\beta_0, \beta_1) = \sum_{i=1}^n |y_i - \beta_0 - \beta_1 x_i| \rightarrow \min_{\beta_0, \beta_1} \quad (15)$$

В данной работе был использован метод Нелдера-Мида [3], применимый к негладким функциям (в том числе к  $M(\beta_0, \beta_1)$ ). Подробнее см. реализация.

## 2.6 Точечное оценивание

### 2.6.1 Основные понятия

Пусть имеется выборка  $\{x_i\}_{i=1}^n$  из генеральной совокупности с плотностью распределения  $f(x, \theta)$ . Предполагается что функциональный вид зависимости задан с точностью до неизвестного параметра  $\theta$ . Требуется по выборке наблюдений  $\{x_i\}_{i=1}^n$  определить число  $\hat{\theta}_n$  которое можно принять за значение параметра  $\theta$ . Точечной оценкой неизвестного параметра  $\theta$  распределения называется борелевская функция наблюдений  $\hat{\theta}_n = \hat{\theta}_n(x_1, \dots, x_n)$ , приближенно равная  $\theta$ . Следует заметить что параметр может быть векторным, к примеру  $\theta = (\mu, \sigma)$  для нормального распределения.

### 2.6.2 Метод максимального правдоподобия

Рассмотрим один общий метод построения точечных оценок. Для начала введем важное понятие, функцией правдоподобия (ФП) называется совместная плотность вероятности распределения  $n$  независимых с.в.  $x_1, \dots, x_n$ :

$$L(x_1, \dots, x_n, \theta) = f(x_1, \theta) f(x_2, \theta) \dots f(x_n, \theta) \quad (16)$$

Оценкой максимального правдоподобия (о. м. п.) будем называть такое значение  $\hat{\theta}_{\text{мп}}$  из множества допустимых значений параметра  $\theta$ , для которого ФП принимает наибольшее значение при заданных  $x_1, \dots, x_n$ :

$$\hat{\theta}_{\text{мп}} = \arg \max_{\theta} L(x_1, \dots, x_n, \theta) \quad (17)$$

Легко обобщается на случай векторного параметра  $\theta = (\theta_1, \dots, \theta_m)$ :

$$\hat{\theta}_{\text{мп}} = \arg \max_{\theta_1, \dots, \theta_m} L(x_1, \dots, x_n, \theta_1, \dots, \theta_m) \quad (18)$$

Известно [1, стр. 444] что о. м. п. нормального распределения являются выборочное среднее и выборочная дисперсия:

$$\hat{\mu}_{\text{мп}} = \bar{x} \quad \hat{\sigma}_{\text{мп}} = \sqrt{s^2} \quad (19)$$

## 2.7 Критерий согласия $\chi^2$

Для проверки гипотезы о законе распределения применяются различные критерии согласия. В данной работе рассматривается наиболее обоснованный и наиболее часто используемый в практике - критерий  $\chi^2$  [1, стр. 482]. И так, выдвинута гипотеза  $H_0$  о генеральном законе распределения с функцией распределения  $F(x)$ . Под конкурирующей гипотезой  $H_1$  понимается гипотеза о справедливости одного из конкурирующих распределений.

Разобьем множество значений изучаемой случайной величины  $X$  на  $k$  непересекающихся подмножеств  $\Delta_1, \dots, \Delta_k$  и пусть  $p_i = \mathbf{P}(X \in \Delta_k)$ . Если множество значений представляет вещественную ось, то подмножества имеют вид:

$$\Delta_i = (a_{i-1}, a_i], i = 2, \dots, k-1 \quad \Delta_1 = (-\infty, a_1] \quad \Delta_k = (a_{k-1}, +\infty) \quad (20)$$

Пусть  $n_1, \dots, n_k$  - частоты попадания выборочных элементов в подмножества  $\Delta_1, \dots, \Delta_k$  соответственно. В случае справедливости гипотезы относительные частоты  $\frac{n_i}{n}$  должны быть близки к  $p_i$  при  $i = 1, \dots, k$ . Поэтому за меру отклонения было предположено (К. Пирсоном) [1, стр. 483] выбрать значение

$$\chi_B^2 = \sum_{i=1}^k \frac{n}{p_i} \left( \frac{n_i}{n} - p_i \right)^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} \quad (21)$$



Существует **теорема**: статистика критерия  $\chi^2$  асимптотически распределена по закону  $\chi^2$  с  $k - 1$  степенями свободы. На основе этой теоремы формируется правило проверки гипотезы о законе распределения по методу  $\chi^2$ : можно принять гипотезу  $H_0$  на уровне значимости  $\alpha$  если  $\chi_B^2 < \chi_{1-\alpha}^2$ , в противном случае она отвергается.

В данной работе  $k$  и длины  $\Delta_1, \dots, \Delta_k$  выбирались по правилам, которые обычно используют при построении гистограмм [4]. Правило Райса для числа интервалов:

$$k = \lceil 1.72\sqrt[3]{n} \rceil \quad (22)$$

и правило Фридмана-Дайсона для ширины (считаем все интервалы кроме крайних одинаковой ширины)

$$a_i = \text{med } N(\hat{\mu}, \hat{\sigma}) + \left(i - \frac{k-1}{2}\right) h, \text{ где } h = 2 \frac{\text{IQR}(x_1, \dots, x_n)}{\sqrt[3]{n}}, i = 2, \dots, k-1 \quad (23)$$

, где  $\text{IQR}(x_1, \dots, x_n)$  - выборочная интерквартильная широта,  $\text{med } N(\hat{\mu}, \hat{\sigma})$  - медиана гипотетического распределения (т.к. предполагается что именно в окрестности медианы будет большая часть элементов выборки).

## 2.8 Интервальное оценивание

*Интервальной оценкой* (или *доверительным интервалом*) числовой характеристики или параметра распределения  $\theta$  генеральной совокупности с доверительной вероятностью  $\gamma$  называется интервал  $(\theta_1, \theta_2)$ , границы которого являются случайными функциями:  $\theta_1 = \theta_1(x_1, \dots, x_n)$ ,  $\theta_2 = \theta_2(x_1, \dots, x_n)$ , который покрывает  $\theta$  с вероятностью  $\gamma$ :

$$\mathbf{P}(\theta_1 < \theta < \theta_2) = \gamma \quad (24)$$

Часто вместо доверительной вероятности  $\gamma$  рассматривается *уровень значимости*  $\alpha = 1 - \gamma$ . Важной характеристикой данной интервальной оценки является половина длины доверительного интервала, она называется *точностью* интервального оценивания

$$\Delta = \frac{\theta_2 - \theta_1}{2} \quad (25)$$

Рассмотрим общий метод построения интервальных оценок [1, стр. 456- 457]. Пусть известна статистика  $Y(\hat{\theta}, \theta)$ , содержащая оцениваемый параметр  $\theta$  и его точечную оценку  $\hat{\theta}$  со следующими свойствами:

- Функция распределения  $F_Y(x)$  известна и не зависит от  $\theta$
- Функция  $Y(\hat{\theta}, \theta)$  непрерывна и строго монотонна (для определенности - строго возрастает) по  $\theta$

которые мы будем проверять при построении интервальных оценок нормального распределения. Зададим уровень значимости  $\alpha$  и будем строить доверительный интервал так чтобы  $(-\infty, \alpha_1), (\alpha_2, +\infty)$  покрывали  $\theta$  с вероятностью  $\frac{\alpha}{2}$ .

Пусть  $y_{\alpha/2}, y_{1-\alpha/2}$  - квантили распределения  $Y$  соотв. порядков, тогда

$$\begin{aligned} \mathbf{P}\left(y_{\alpha/2} < Y(\hat{\theta}, \theta) < y_{1-\alpha/2}\right) &= F_Y(y_{1-\alpha/2}) - F_Y(y_{\alpha/2}) = \\ &= 1 - \alpha/2 - \alpha/2 = 1 - \alpha = \gamma \end{aligned} \quad (26)$$

Т.к.  $Y(\hat{\theta}, \theta)$  – строго возрастает по  $\theta$ , то у неё есть обратная функция  $Y^{-1}(y)$  относительно  $\theta$  и она также строго возрастает, а значит:

$$\begin{aligned} y_{\alpha/2} &< Y(\hat{\theta}, \theta) < y_{1-\alpha/2} \\ Y^{-1}(y_{\alpha/2}) &< \theta < Y^{-1}(y_{1-\alpha/2}) \end{aligned} \quad (27)$$

итого  $\theta_1 = Y^{-1}(y_{\alpha/2})$  и  $\theta_2 = Y^{-1}(y_{1-\alpha/2})$  – мы построили границы интервала. Применим это для построения интервальных оценок нормального распределения по выборке  $(x_1, \dots, x_n)$ .

## 2.9 Классические оценки

### 2.9.1 Для математического ожидания $m$

Доказано что случайная величина  $T = \sqrt{n-1} \cdot \frac{\bar{x}-m}{s}$  называемая статистикой Стьюдента, распределена по закону Стьюдента с  $n-1$  степенями свободы, применяя с некоторыми деталями [1, стр. 457 – 458] выкладки, получаем оценки границ интервала:

$$\begin{aligned} m_1 &= \bar{x} - \frac{xt_{1-\alpha/2}(n-1)}{\sqrt{n-1}} \\ m_2 &= \bar{x} + \frac{xt_{1-\alpha/2}(n-1)}{\sqrt{n-1}} \end{aligned} \quad (28)$$

, где  $t_{1-\alpha/2}(n-1)$  – квантиль порядка  $1-\alpha/2$  распределения Стьюдента с  $n-1$  степенями свободы.

### 2.9.2 Для среднего квадратичного отклонения $\sigma$

Доказано что случайная величина  $ns^2/\sigma^2$  распределена по закону  $\chi^2$  с  $n-1$  степенями свободы. Применяя общий метод построения интервальных оценок получаем оценки границ интервала:

$$\begin{aligned} \sigma_1 &= \frac{s\sqrt{n}}{\sqrt{\chi_{1-\alpha/2}^2(n-1)}} \\ \sigma_2 &= \frac{s\sqrt{n}}{\sqrt{\chi_{\alpha/2}^2(n-1)}} \end{aligned} \quad (29)$$

, где  $\chi_{1-\alpha/2}^2(n-1), \chi_{\alpha/2}^2(n-1)$  – квантили соотв. порядков  $\chi^2$ -распределения с  $n-1$  степенями свободы.

## 2.10 Асимптотически нормальные оценки

### 2.10.1 Для математического ожидания $m$

В силу центральной предельной теоремы центрированная и нормированная случайная величина  $\sqrt{n}(\bar{x}-m)/\sigma$  распределена приблизительно нормально с параметрами 0 и 1. Исходя из этого [1, стр. 460] получаем оценку:

$$\begin{aligned} m_1 &= \bar{x} - \frac{su_{1-\alpha/2}}{\sqrt{n}} \\ m_2 &= \bar{x} + \frac{su_{1-\alpha/2}}{\sqrt{n}} \end{aligned} \quad (30)$$

, где  $u_{1-\alpha/2}$  – квантиль нормального распределения  $N(0, 1)$  порядка  $1-\alpha/2$

### 2.10.2 Для среднего квадратичного отклонения $\sigma$

Аналогично, в силу центральной предельной теоремы центрированная и нормированная случайная величина  $(s^2 - \mathbf{M}s^2)/\sqrt{\mathbf{D}s^2}$  при большом объеме выборки  $n$  распределена приблизительно нормально с параметрами 0 и 1. Исходя из этого [1, стр. 461] получаем оценку:

$$\begin{aligned}\sigma_1 &= s \left( 1 + u_{1-\alpha/2} \sqrt{(e+2)/n} \right)^{-1/2} \\ \sigma_2 &= s \left( 1 - u_{1-\alpha/2} \sqrt{(e+2)/n} \right)^{-1/2}\end{aligned}\tag{31}$$

, где  $e$  - выборочный эксцесс, определяемый как

$$e = \frac{m_4}{s^4} - 3\tag{32}$$

, где  $m_4$  - четвертый выборочный центральный момент, определяемый как

$$m_4 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4\tag{33}$$

### 3 Реализация

Работа выполнена с использованием языка **Python** в интегрированной среде разработки **PyCharm**, были задействованы библиотеки:

- **NumPy** - векторизация вычислений, работа с массивами данных, вычисление выборочных характеристик
- **SciPy** - модуль **stats** для генерации данных по распределениям и эталонной зависимости, вычисления коэффициентов корреляции, оценок МНК, оценки методом максимального правдоподобия, модуль **optimize** для метода Нелдера-Мида
- **Matplotlib** - построение эллипсов рассеяния, построение графиков

Исходный код лабораторных работ приведен в приложении.

## 4 Результаты

### 4.1 Коэффициенты корреляции

$n = 20$	$r(4)$	$r_S(8)$	$r_Q(7)$
$E(z)$	0.0	0.0	0.0
$E(z^2)$	0.1	0.1	0.1
$D(z)$	0.053556	0.053729	0.054264
$n = 60$	$r$	$r_S$	$r_Q$
$E(z)$	0.0	0.0	0.0
$E(z^2)$	0.02	0.02	0.02
$D(z)$	0.016997	0.017097	0.018564
$n = 100$	$r$	$r_S$	$r_Q$
$E(z)$	0.0	0.0	0.0
$E(z^2)$	0.01	0.01	0.01
$D(z)$	0.010095	0.010177	0.010416

Таблица 1:  $\rho = 0$

$n = 20$	$r$	$r_S$	$r_Q$
$E(z)$	0.5	0.5	0.3
$E(z^2)$	0.3	0.2	0.2
$D(z)$	0.03309	0.036427	0.046232
$n = 60$	$r$	$r_S$	$r_Q$
$E(z)$	0.49	0.47	0.33
$E(z^2)$	0.25	0.23	0.12
$D(z)$	0.009512	0.010584	0.014097
$n = 100$	$r$	$r_S$	$r_Q$
$E(z)$	0.5	0.48	0.33
$E(z^2)$	0.25	0.23	0.12
$D(z)$	0.006043	0.00652	0.008687

Таблица 2:  $\rho = 0.5$

$n = 20$	$r$	$r_S$	$r_Q$
$E(z)$	0.891	0.86	0.7
$E(z^2)$	0.797	0.75	0.52
$D(z)$	0.002823	0.005097	0.02752
$n = 60$	$r$	$r_S$	$r_Q$
$E(z)$	0.899	0.882	0.71
$E(z^2)$	0.809	0.78	0.51
$D(z)$	0.000641	0.001065	0.008925
$n = 100$	$r$	$r_S$	$r_Q$
$E(z)$	0.898	0.885	0.71
$E(z^2)$	0.807	0.784	0.5
$D(z)$	0.000417	0.000637	0.004951

Таблица 3:  $\rho = 0.9$

$n = 20$	$r$	$r_S$	$r_Q$
$E(z)$	-0.0	0.5	0.5
$E(z^2)$	1.0	0.3	0.3
$D(z)$	0.448101	0.078135	0.039344
$n = 60$	$r$	$r_S$	$r_Q$
$E(z)$	-0.6	0.48	0.56
$E(z^2)$	0.5	0.26	0.33
$D(z)$	0.079885	0.027009	0.01148
$n = 100$	$r$	$r_S$	$r_Q$
$E(z)$	-0.7	0.47	0.56
$E(z^2)$	0.51	0.24	0.33
$D(z)$	0.029483	0.015814	0.006452

Таблица 4: Смесь нормальных распределений

4.2 Эллипсы равновероятности

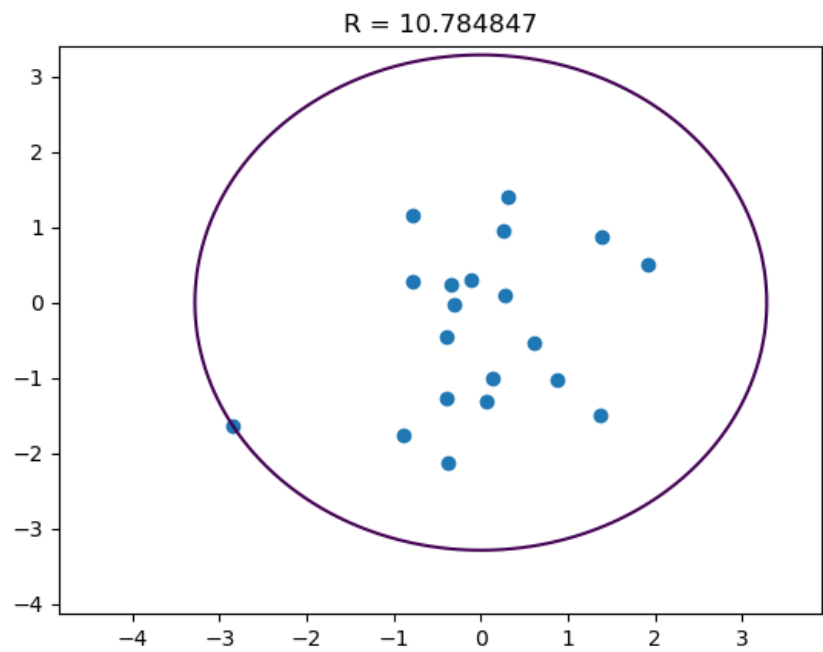


Рис. 1:  $\rho = 0.0, n = 20$

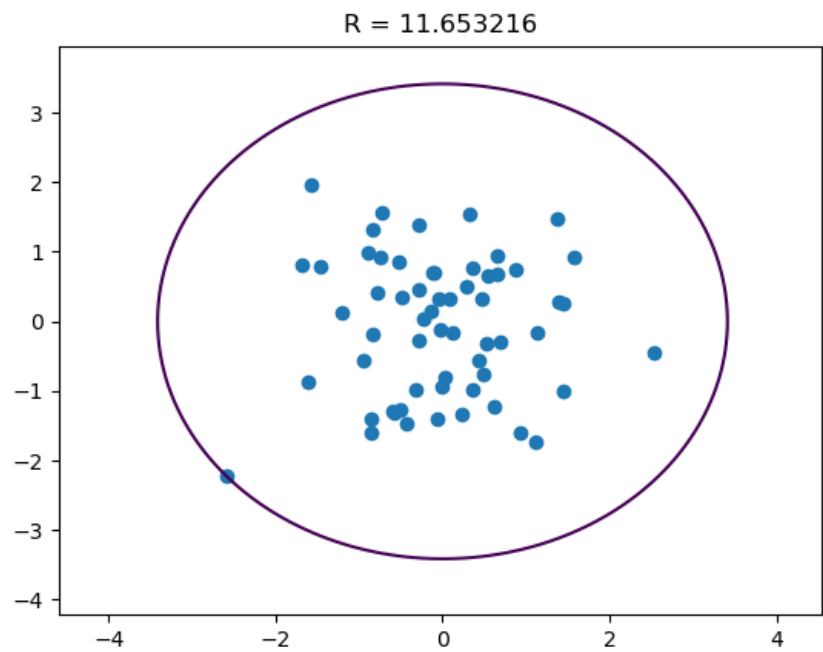


Рис. 2:  $\rho = 0.0, n = 60$

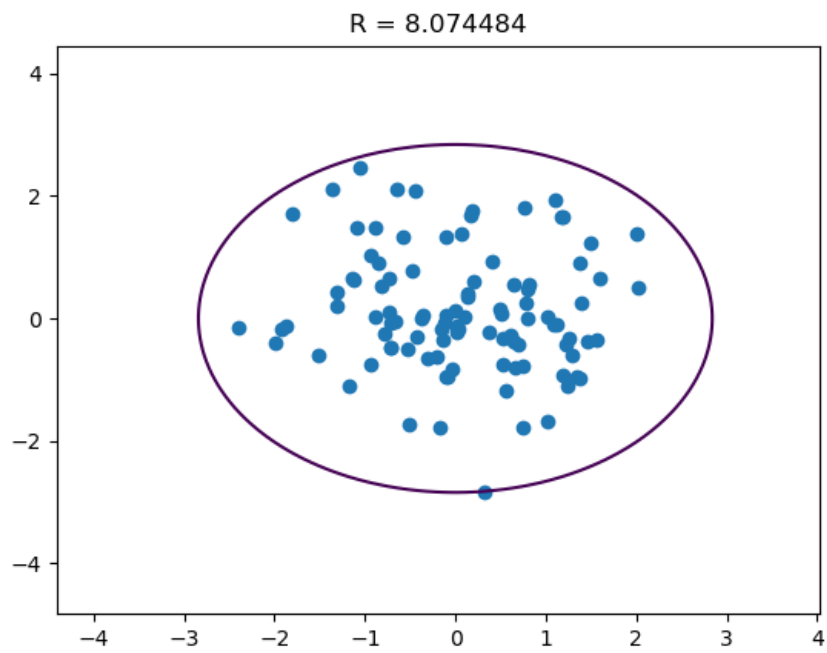


Рис. 3:  $\rho = 0.0, n = 100$

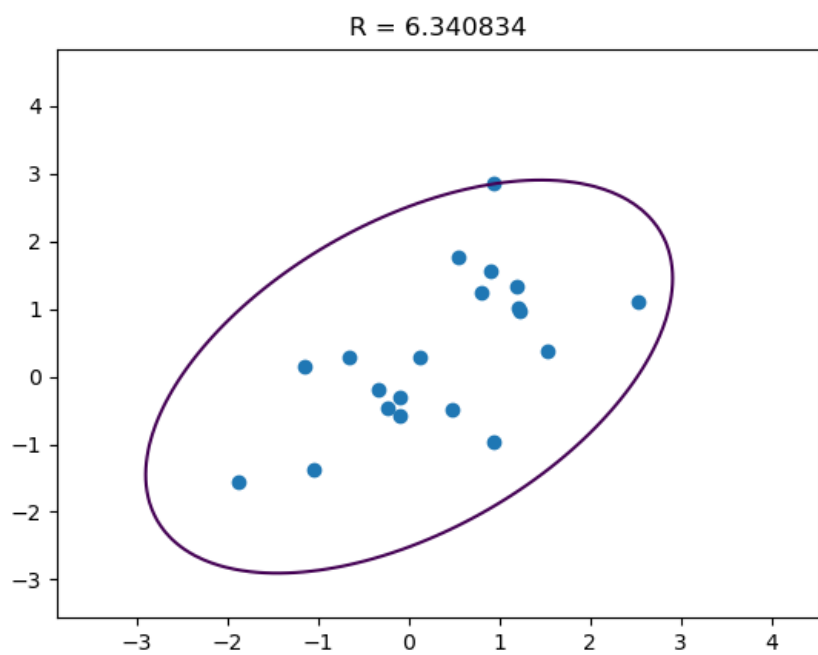


Рис. 4:  $\rho = 0.5, n = 20$



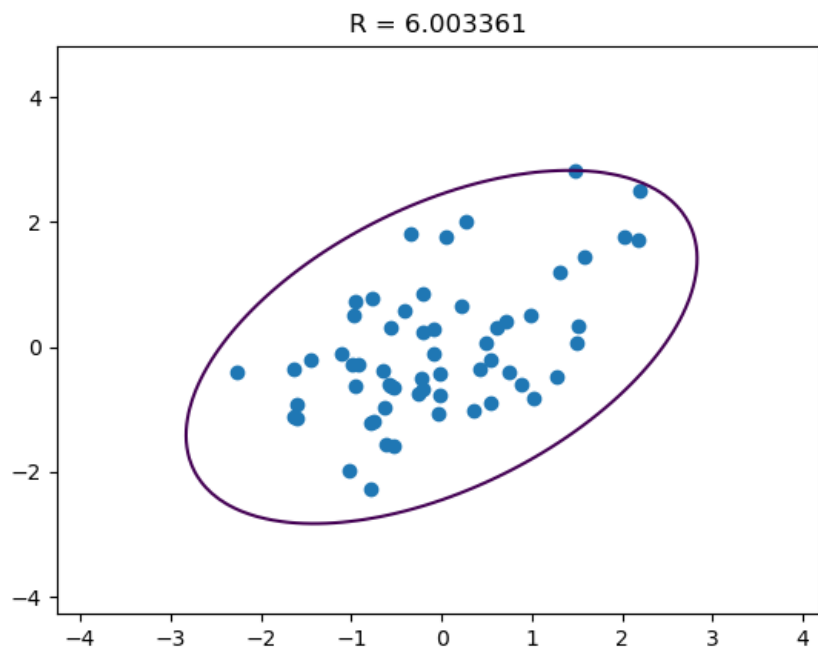


Рис. 5:  $\rho = 0.5, n = 60$

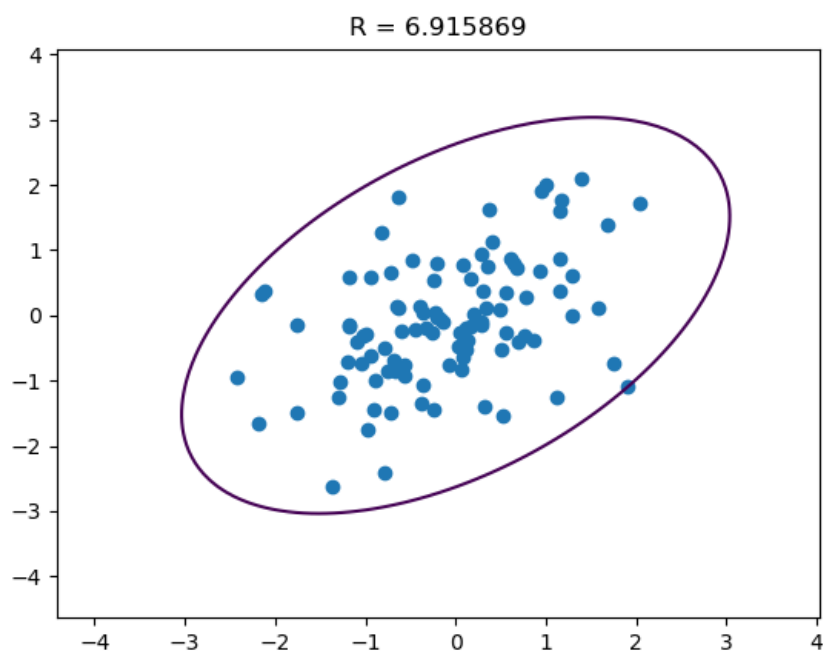


Рис. 6:  $\rho = 0.5, n = 100$

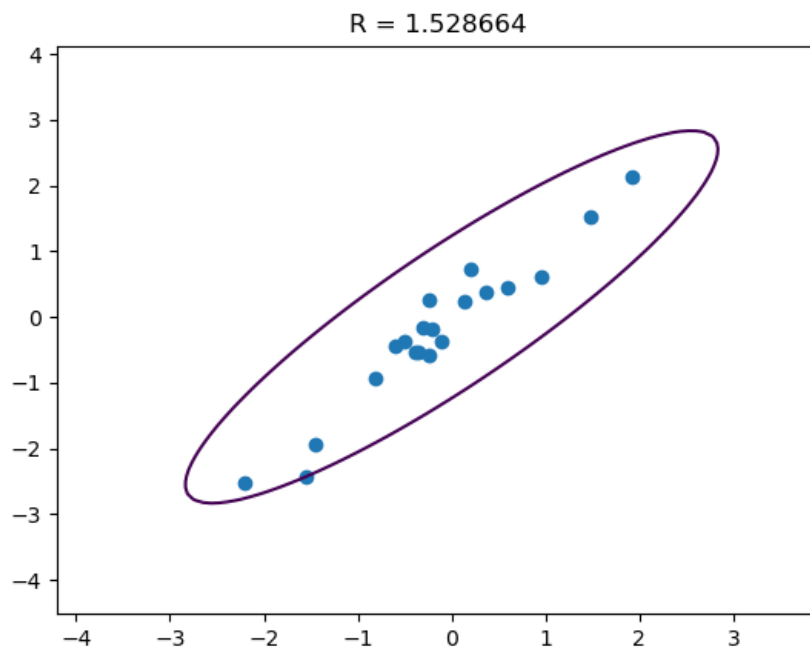


Рис. 7:  $\rho = 0.9, n = 20$

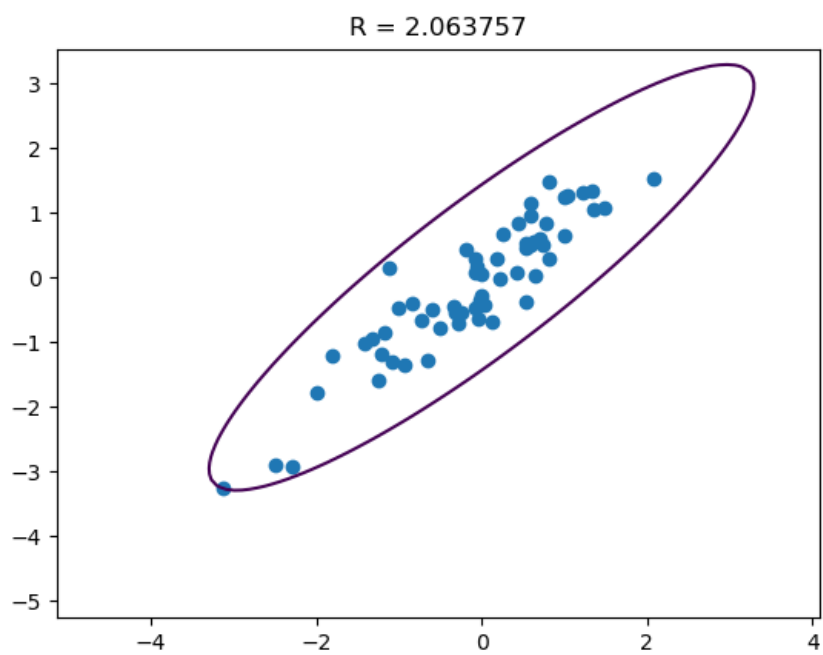


Рис. 8:  $\rho = 0.9, n = 60$

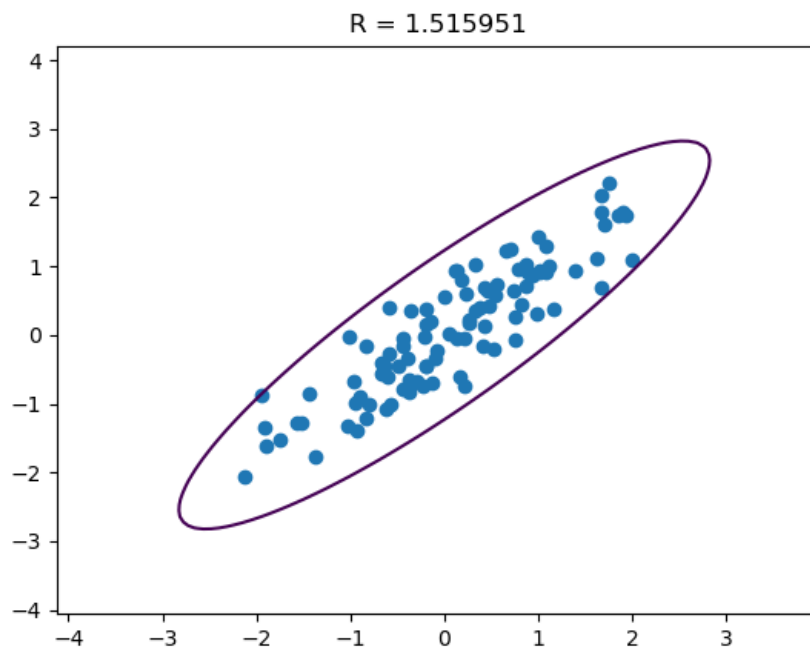


Рис. 9:  $\rho = 0.9, n = 100$

### 4.3 Выборка без выбросов

- Критерий наименьших квадратов:

$$\hat{\beta}_0 = 2.47 \quad \hat{\beta}_1 = 1.95 \quad Q(13) = 13.9637 \quad M(15) = 13.9182$$

- Критерий наименьших модулей:

$$\hat{\beta}_0 = 2.49 \quad \hat{\beta}_1 = 1.68 \quad Q = 15.9356 \quad M = 13.3737$$

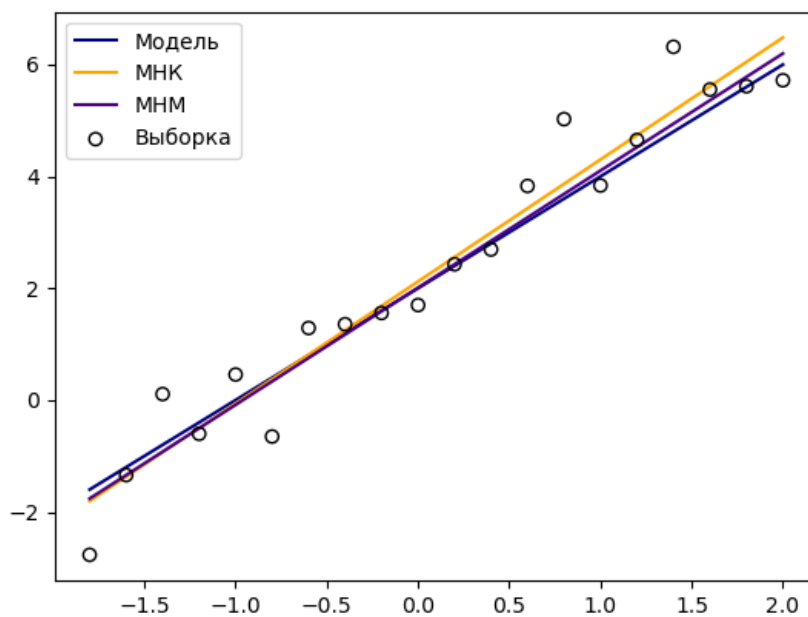


Рис. 10: Без выбросов

## 4.4 Выборка с выбросами

- Критерий наименьших квадратов:

$$\hat{\beta}_0 = 2.61 \quad \hat{\beta}_1 = 0.52 \quad Q = 154.2302 \quad M = 37.381$$

- Критерий наименьших модулей:

$$\hat{\beta}_0 = 2.67 \quad \hat{\beta}_1 = 1.35 \quad Q = 172.7536 \quad M = 29.9906$$

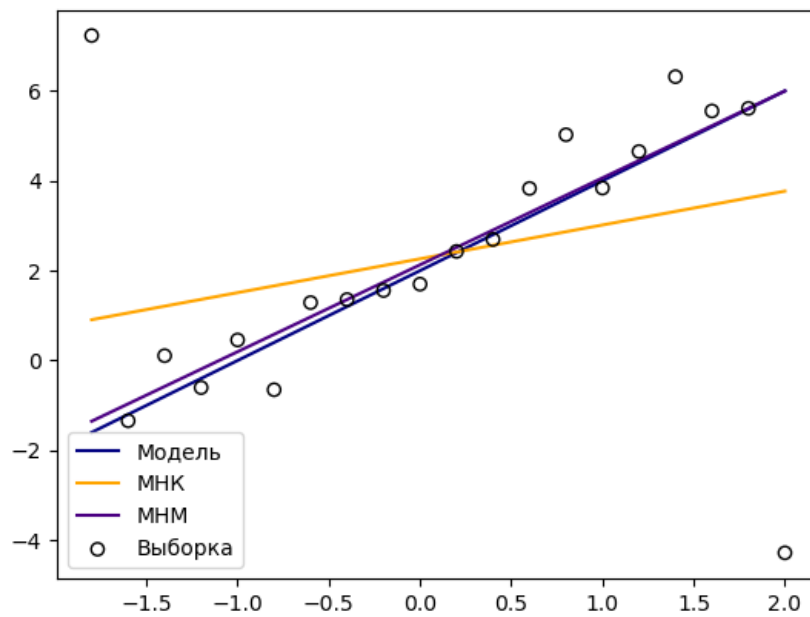


Рис. 11: С выбросами

## 4.5 Критерий согласия $\chi^2$

Оценки:

$$\hat{\mu} = 0.03 \quad \hat{\sigma} = 1.01 \quad (34)$$

Число промежутков:  $k = \lceil 1.72 \cdot \sqrt[3]{100} \rceil = 8$

Таблица вычислений  $\chi^2$ :

$i$	$\Delta_i$	$n_i$	$p_i$	$n_i - np_i$
1	$(-\infty, -1.79]$	2	0.0366	-1.6609
2	$(-1.79, -1.27]$	7	0.0637	0.627
3	$(-1.27, -0.75]$	12	0.121	-0.0972
4	$(-0.75, -0.23]$	22	0.1777	4.2306
5	$(-0.23, 0.29]$	22	0.202	1.8009
6	$(0.29, 0.81]$	15	0.1777	-2.7694
7	$(0.81, 1.33]$	8	0.121	-4.0972
8	$(1.33, +\infty)$	12	0.1003	1.9661

Таблица 5: Таблица вычислений  $\chi^2$

При  $\alpha = 0.05$ :  $\chi_{1-\alpha}^2(k-1) \approx 14.0671$ , а вычисленное  $\chi_B^2 = 4.1883$ , видно что  $\chi_B^2 < \chi_{1-\alpha}^2(k-1)$

В результате доп. исследования, было получено что при  $n = 20$  критерий дает вывод что генеральное распределение является нормальным  $N(0.024, 0.59)$ , в результате вычислений  $\chi_B^2 = 4.8612 < 4.8784 = \chi_{1-\alpha}^2$ , а при  $n = 100$  уже  $\chi_B^2 = 19.2086 \geq 8.3834 = \chi_{1-\alpha}^2$  т.е. установлено что генеральное распределение не является нормальным (и это соответствует тому что оно задано как равномерное)

## 4.6 Классические оценки

	$m(28)$	$\sigma(29)$
$n = 20$	$-0.6 < m < 0.27$	$0.71 < \sigma < 1.36$
$n = 100$	$-0.04 < m < 0.34$	$0.84 < \sigma < 1.12$

## 4.7 Асимптотически нормальные оценки

	$m(30)$	$\sigma(31)$
$n = 20$	$-0.56 < m < 0.23$	$0.71 < \sigma < 1.46$
$n = 100$	$-0.03 < m < 0.34$	$0.84 < \sigma < 1.13$

## 5 Обсуждение

### 5.1 Коэффициенты корреляции

Для начала воспользуемся (5) для анализа экспериментов по которым были получены таблицы 1, 2. Выясним можно ли принять гипотезу о зависимости между случайными величинами на уровне значимости  $\alpha = 0.05$  для  $n = 100$  по коэффициенту Пирсона.

$$0\sqrt{100-1} \leq 2.5, 4.98 \approx 0.5\sqrt{100-1} > .2.5$$

В эксперименте 1 эту гипотезу принять нельзя, а в эксперименте 2 можно. При этом в эксперименте 1 с.в. заведомо независимы, а в эксперименте 2 зависимы, так что все согласуется с теорией.

Из таблиц 1, 2 и 3 видно что  $r, r_S$  являются состоятельными оценками  $\rho_{XY}$  т.к. они все ближе к нему с ростом  $n$ .

Из таблицы 4 видим что  $r_Q$  устойчивая к выбросам (робастная) оценка. Квадрантный коэффициент корреляции показывает лучшие результаты в устойчивости.

### 5.2 Эллипсы равновероятности

Видно что чем ближе  $\rho$  к 1, тем эллипс равновероятности становится все больше похож на прямую, заданную как (10). Т.е. наглядно показано как между с.в.  $X$  и  $Y$  возникает линейная зависимость.

### 5.3 Линейная регрессия

Из графиков видно, что оценка по критерию наименьших модулей значительно лучше приближает эталонную зависимость при наличии выбросов и это согласуется с теорией т.к. она является робастной. В тоже время, критерий наименьших квадратов дает более точное приближение в отсутствие выбросов и, к тому же, проще для вычислений. Полученные значения  $M, Q$  упорядочены как и ожидалось, для оценки МНК значение  $Q$  меньше, чем для любой другой, аналогично для оценки МНМ и значения  $M$

### 5.4 Критерий согласия $\chi^2$

Согласно результатам эксперимента, заданное по оценкам (34) распределение  $N(\hat{\mu}, \hat{\sigma})$  является генеральным законом по которому построена выборка с уровнем значимости 0.05. Теоретически это обосновывается тем что оценки максимального правдоподобия состоятельны. Было установлено что при небольших объемах выборки уверенности в полученных результатах нет, ведь статистика критерия  $\chi^2$  лишь асимптотически распределена по закону  $\chi^2(k-1)$  т.е.  $n$  предполагается достаточно большим.

### 5.5 Интервальное оценивание

Полученные интервальные оценки говорят о том что с вероятностью 0.95 значения  $m = 0$  и  $\sigma = 1$  лежат в соответствующих интервалах. По постановке эксперимента, интервалы действительно накрывают истинные значения параметров. Следует заметить что при большом объеме  $n$  выборки - асимптотические оценки практически совпадают с классическими.

## 6 Приложения

1. Исходный код лабораторной 5 <https://github.com/zhenyatos/statlabs/tree/master/Lab5>
2. Исходный код лабораторной 6 <https://github.com/zhenyatos/statlabs/tree/master/Lab6>
3. Исходный код лабораторной 7 <https://github.com/zhenyatos/statlabs/tree/master/Lab7>
4. Исходный код лабораторной 8 <https://github.com/zhenyatos/statlabs/tree/master/Lab8>

## Список литературы

- [1] **Вероятностные разделы математики.** Учебник для бакалавров технических направлений. // Под ред. Максимова Ю.Д. - СПб «Иван Федоров», 2001. - 592 с., илл
- [2] Вентцель Е.С. *Теория вероятностей: Учеб. для вузов.* — 6-е изд. стер. — М.: Высш. шк., 1999.— 576 с.
- [3] Метод Нелдера — Мида // Википедия. [2019—2019]. Дата обновления: 11.09.2019. URL: <https://ru.wikipedia.org/?oldid=102111276> (дата обращения: 11.09.2019).
- [4] Wikipedia contributors. (2020, March 19). Histogram. In Wikipedia, The Free Encyclopedia. Retrieved 18:27, May 14, 2020, from <https://en.wikipedia.org/w/index.php?title=Histogram&oldid=946321806>