

Санкт-Петербургский
Политехнический университет Петра Великого

**Отчет по лабораторным работам №1-4
по дисциплине
"Математическая статистика"**

Студент:	Скворцов Владимир Сергеевич
Преподаватель:	Баженов Александр Николаевич
Группа:	5030102/10201

Санкт-Петербург
2024

Содержание

1	Постановка задачи	2
1.1	Описательная статистика	2
1.2	Точечное оценивание характеристик положения и рассеяния	2
2	Теоретическое обоснование	2
2.1	Функции распределения	2
2.2	Характеристики положения и рассеяния	3
3	Описание работы	3
4	Результаты	4
4.1	Гистограммы и графики плотности распределения	4
4.2	Характеристики положения и рассеяния	6
5	Выводы	8
6	Постановка задачи	9
6.1	Боксплот Тьюки	9
6.2	Доверительные интервалы для параметров нормального распределения . . .	9
7	Теоретическое обоснование	9
7.1	Функции распределения	9
7.2	Боксплот Тьюки	10
7.3	Доверительные интервалы для параметров нормального распределения . . .	10
8	Описание работы	10
9	Результаты	11
9.1	Гистограммы и графики плотности распределения	11
9.2	Доверительные интервалы для параметров распределений	13
10	Выводы	14

1 Постановка задачи

1.1 Описательная статистика

Для 5 распределений:

- Нормальное распределение $N(x, 0, 1)$
- распределение Коши $C(x, 0, 1)$
- Распределение Стьюдента $t(x, 0, 3)$ с тремя степенями свободы
- Распределение Пуассона $P(k, 10)$
- Равномерное распределение $U(x, -\sqrt{3}, \sqrt{3})$

Сгенерировать выборки размером 10, 50, 1000 элементов.

Построить на одном рисунке гистограмму и график плотности распределения.

1.2 Точечное оценивание характеристик положения и рассеяния

Сгенерировать выборки размером 10, 50, 1000 элементов.

Для каждой выборки вычислить следующие статистические характеристики положения данных: \bar{x} , $med\ x$, z_Q , z_R , z_{tr} . Повторить такие вычисления 1000 раз для каждой выборки и найти среднее характеристик положения и их квадратов: $E(z) = \bar{z}$. Вычислить оценку дисперсии по формуле $D(z) = \overline{z^2} - \bar{z}^2$.

2 Теоретическое обоснование

2.1 Функции распределения

- Нормальное распределение

$$N(x, 0, 1) = \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} \quad (1)$$

- Распределение Коши

$$C(x, 0, 1) = \frac{1}{\pi} \frac{1}{x^2 + 1} \quad (2)$$

- Распределение Стьюдента $t(x, 0, 3)$ с тремя степенями свободы

$$t(x, 0, 3) = \frac{6\sqrt{3}}{\pi(3 + t^2)^2} \quad (3)$$

- Распределение Пуассона

$$P(k, 10) = \frac{10^k}{k!} e^{-10} \quad (4)$$

- Равномерное распределение

$$U(x, -\sqrt{3}, \sqrt{3}) = \begin{cases} \frac{1}{2\sqrt{3}} & \text{при } |x| \leq \sqrt{3} \\ 0 & \text{при } |x| > \sqrt{3} \end{cases} \quad (5)$$

2.2 Характеристики положения и рассеяния

- Выборочное среднее

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (6)$$

- Выборочная медиана

$$\text{med } x = \begin{cases} x_{(l+1)} & \text{при } n = 2l + 1 \\ \frac{x_{(l)} + x_{(l+1)}}{2} & \text{при } n = 2l \end{cases} \quad (7)$$

- Полусумма экстремальных выборочных элементов

$$z_R = \frac{x_{(1)} + x_{(n)}}{2} \quad (8)$$

- Полусумма квартилей

Выборочная квартиль z_p порядка p определяется формулой

$$z_p = \begin{cases} x_{([np]+1)} & \text{при } np \text{ дробном} \\ x_{(np)} & \text{при } np \text{ целом} \end{cases} \quad (9)$$

Полусумма квартилей

$$z_Q = \frac{z_{1/4} + z_{3/4}}{2} \quad (10)$$

- Усечённое среднее

$$z_{tr} = \frac{1}{n-2r} \sum_{i=r+1}^{n-r} x_{(i)}, \quad r \approx \frac{n}{4} \quad (11)$$

- Среднее характеристики

$$E(z) = \bar{z} \quad (12)$$

- Оценка дисперсии

$$D(z) = \overline{z^2} - \bar{z}^2 \quad (13)$$

3 Описание работы

Лабораторные работы выполнены с использованием Python и его сторонних библиотек `numpy`, `pandas`, `matplotlib`, `seaborn` были построены гистограммы распределений и посчитаны характеристики положения.

Ссылка на GitHub репозиторий: <https://github.com/vladimir-skvortsov/spbstu-mathematical-statistics>

4 Результаты

4.1 Гистограммы и графики плотности распределения

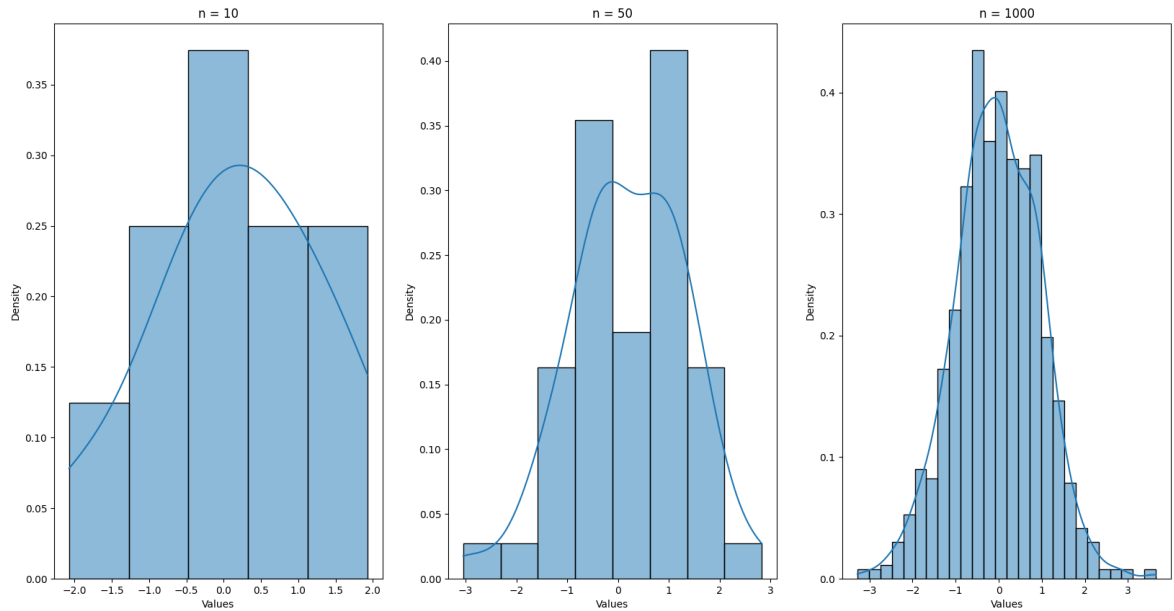


Рис. 1: Нормальное распределение (14)

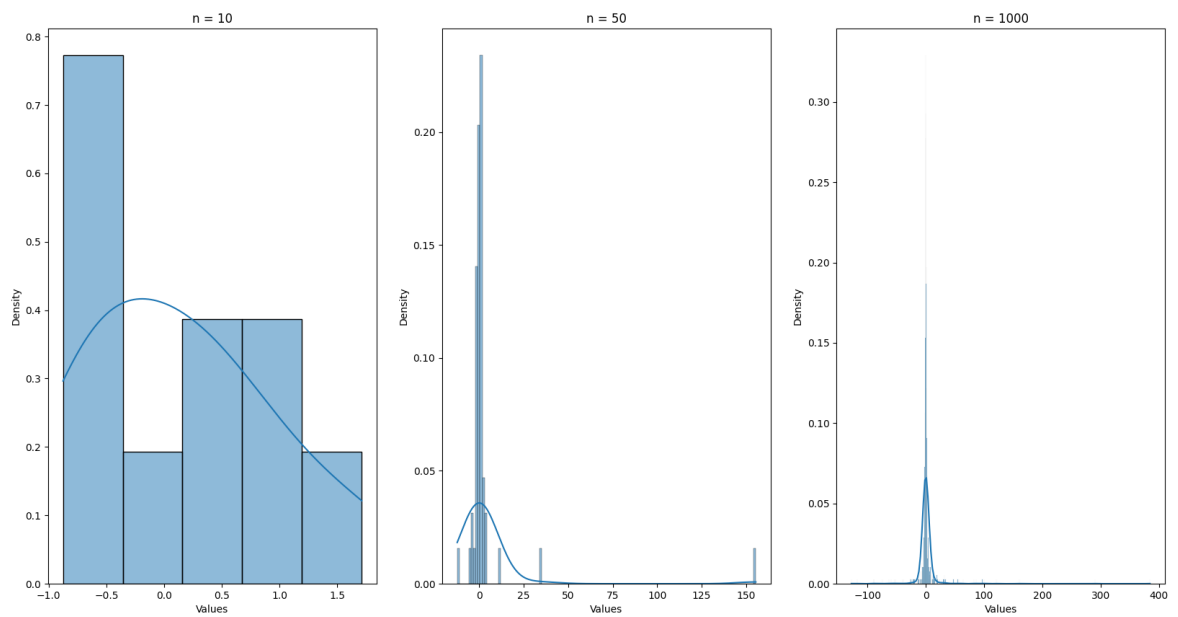


Рис. 2: Распределение Коши (15)

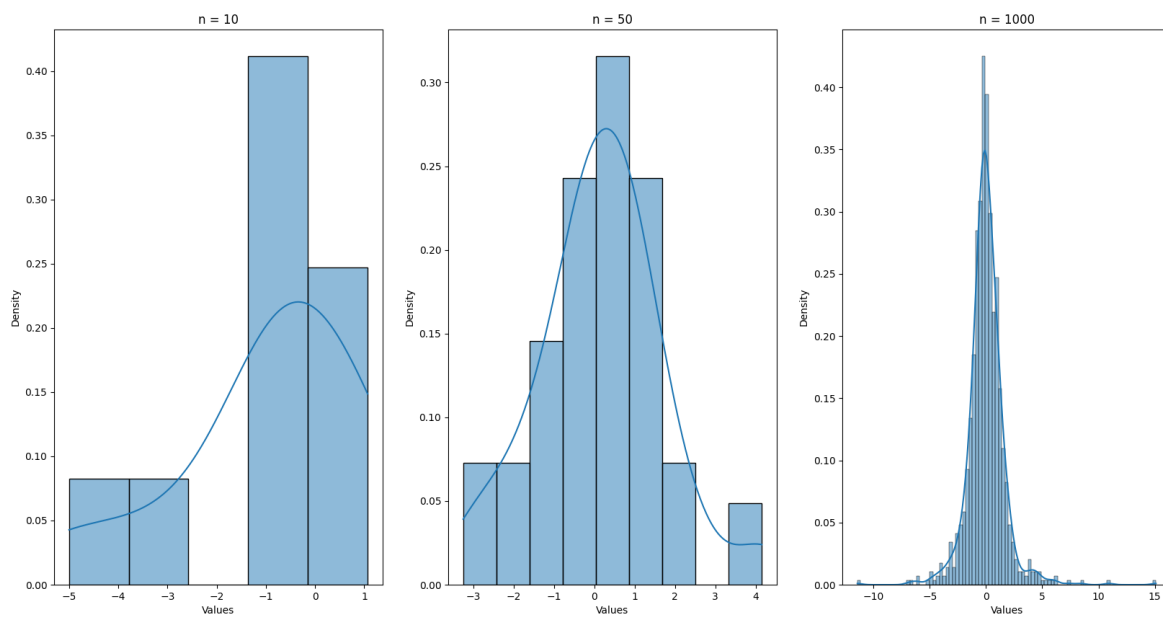


Рис. 3: Распределение Стьюдента (16)

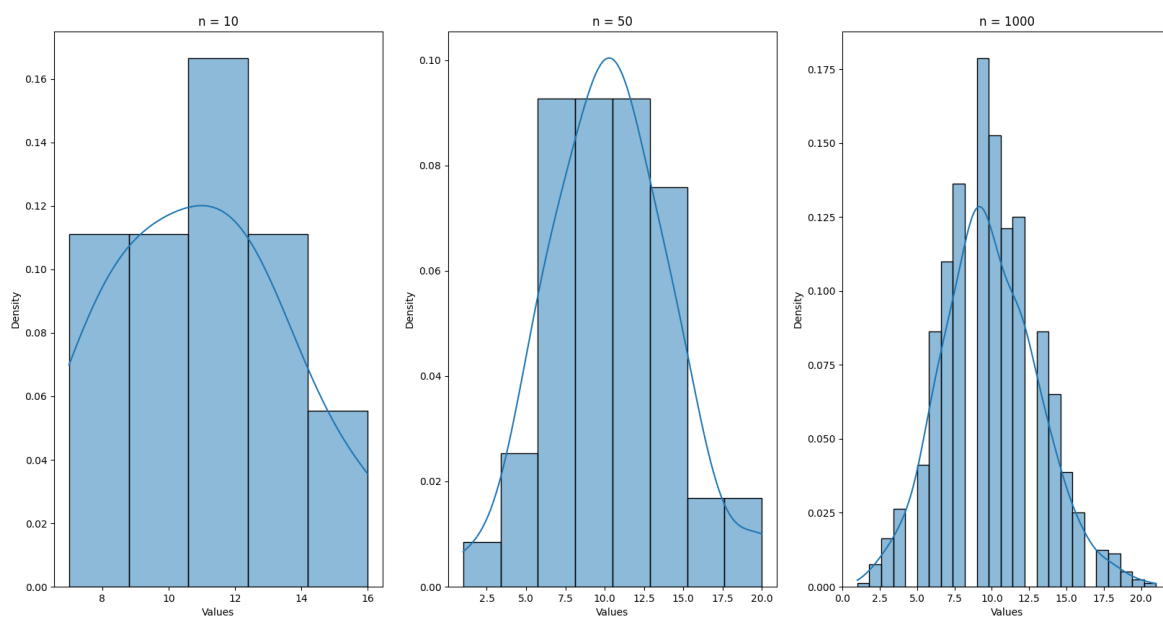


Рис. 4: Распределение Пуассона (17)

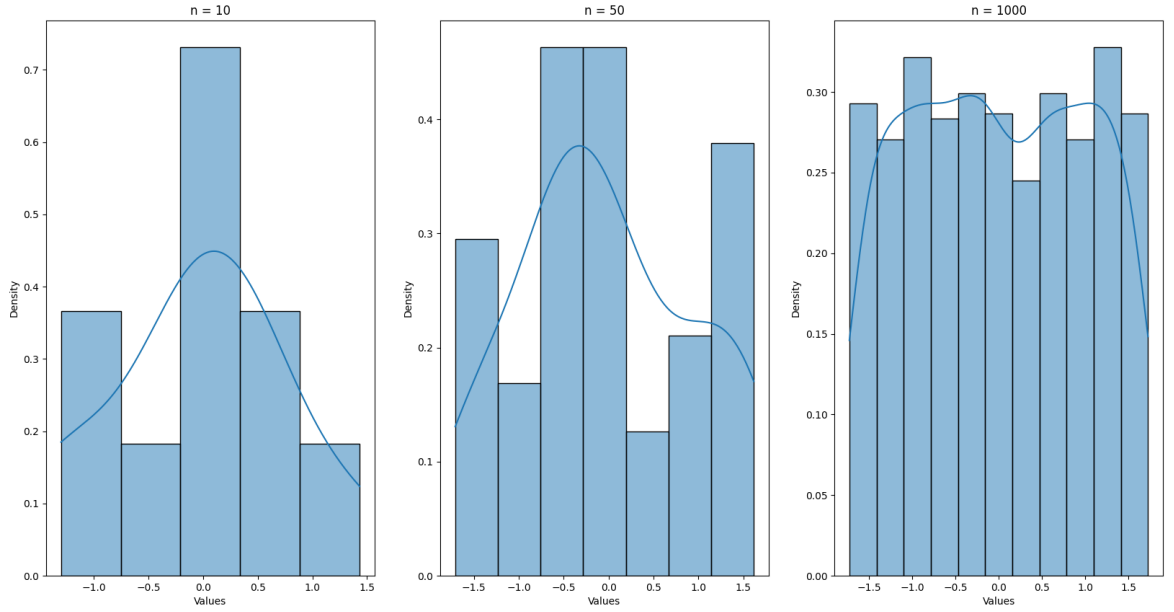


Рис. 5: Равномерное распределение (18)

4.2 Характеристики положения и рассеяния

n = 10					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	-0.017466	-0.019283	-0.019494	-0.014486	-0.007937
$D(z)$ (13)	0.100879	0.142707	0.187775	0.115437	0.160836
n = 50					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	-0.007937	0.100879	0.142707	0.187775	0.115437
$D(z)$ (13)	0.009941	0.015535	0.095586	0.012392	0.020005
n = 1000					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	0.000038	-0.001779	-0.002971	0.001002	-0.000085
$D(z)$ (13)	0.000985	0.001682	0.061385	0.001243	0.001939

Таблица 1: Нормальное распределение

n = 10					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	-4.724165	-0.015986	-23.612109	-0.015176	-8.310631
$D(z)$ (13)	11477.749749	0.337081	286469.541418	1.163577	31698.450396
n = 50					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	0.781733	0.012225	37.029997	0.008637	0.857304
$D(z)$ (13)	431.900044	0.025323	1060046.375320	0.055008	167.707860
n = 1000					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	-0.336134	-0.001532	-129.057477	-0.001540	-0.049715
$D(z)$ (13)	240.553988	0.002310	50362265.313181	0.004735	174.261104

Таблица 2: Распределение Коши

n = 10					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	0.016265	0.004667	0.040925	0.014315	0.000750
$D(z)$ (13)	0.259126	0.183832	1.659319	0.184564	0.431912
n = 50					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	-0.002158	-0.001389	0.021238	0.003592	-0.016753
$D(z)$ (13)	0.026907	0.019051	9.893951	0.018478	0.052782
n = 1000					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	0.000335	-0.000238	-0.054818	0.000162	0.000679
$D(z)$ (13)	0.002898	0.001903	32.527888	0.001944	0.005656

Таблица 3: Распределение Стьюдента

n = 10					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	10.002500	9.874000	10.294500	9.917625	9.937000
$D(z)$ (13)	1.081944	1.477624	2.018020	1.283917	1.699309
n = 50					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	10.013690	9.855500	10.896000	9.945125	10.013560
$D(z)$ (13)	0.095748	0.197370	0.957184	0.139817	0.204837
n = 1000					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	10.005702	9.997000	11.627000	9.994000	10.004912
$D(z)$ (13)	0.010137	0.002991	0.634371	0.002964	0.020719

Таблица 4: Распределение Пуассона

n = 10					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	-0.005450	-0.006939	-0.005412	-0.007901	-0.015610
$D(z)$ (13)	0.104110	0.240206	0.044016	0.144291	0.172234
n = 50					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	-0.001915	-0.006312	-0.001349	0.001960	-0.004766
$D(z)$ (13)	0.010019	0.029723	0.000599	0.014276	0.018935
n = 1000					
	\bar{x} (6)	$med\ x$ (7)	z_R (8)	z_Q (10)	z_{tr} (11)
$E(z)$ (12)	0.000470	0.000924	-0.000133	-0.000355	-0.000387
$D(z)$ (13)	0.001014	0.003127	0.000005	0.001469	0.001887

Таблица 5: Равномерное распределение

5 Выводы

В процессе выполнения лабораторной работы был проведен анализ пяти уникальных распределений: нормальное, Коши, Стьюдента, Пуассона и равномерное. Были сгенерированы выборки разных объемов для каждого из них - 10, 50 и 1000 элементов. Были созданы гистограммы каждого распределения и нанесены на них графики плотности соответствующих распределений, что облегчило наглядное сопоставление формы распределения выборок с их теоретическими аналогами. Были также рассчитаны разные показатели положения и рассеяния для каждой выборки, включая выборочную среднюю величину, медиану, полусумму крайних элементов выборки, полусумму квартилей и усеченное среднее. Использовалась стандартная формула для оценки дисперсии.

На основании полученных данных были сделаны следующие выводы:

1. В случае нормального распределения, оценки показателей положения и рассеяния становятся ближе к их теоретическим значениям по мере увеличения размера выборки.
2. Для распределения Коши показатели положения и рассеяния менее стабильны и могут сильно отличаться от теоретических даже при больших размерах выборки.
3. Распределение Стьюдента при небольших размерах выборки также демонстрирует определенную нестабильность оценок, однако с увеличением размера выборки результаты становятся более точными.
4. Для распределения Пуассона и равномерного распределения, оценки показателей положения и рассеяния кажутся стабильными при любом объеме выборки.
5. В общем, выборочное среднее является наиболее чувствительным к экстремальным значениям по сравнению с медианой, особенно в меньших выборках. Однако с увеличением размера выборки, влияние этих экстремальных значений на среднее значение уменьшается. В то же время, медиана обычно более устойчива к выбросам и мало варьирует с изменением размера выборки.
6. Медиана является чувствительной к типу распределения: в нормальном и распределении Стьюдента медиана равна среднему, в распределении Коши она дает надежные, устойчивые к выбросам оценки, в Пуассоновском приближается к среднему, и в равномерном равна половине суммы минимального и максимального значений.

6 Постановка задачи

6.1 Боксплот Тьюки

Сгенерировать выборки размером 20 и 100 элементов. Построить для них боксплот Тьюки.

6.2 Доверительные интервалы для параметров нормального распределения

Сгенерировать выборки размером 20 и 100 элементов. Вычислить параметры положения и рассеяния:

- для нормального распределения,
- для произвольного распределения.

7 Теоретическое обоснование

7.1 Функции распределения

- Нормальное распределение

$$N(x, 0, 1) = \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} \quad (14)$$

- Распределение Коши

$$C(x, 0, 1) = \frac{1}{\pi} \frac{1}{x^2 + 1} \quad (15)$$

- Распределение Стьюдента $t(x, 0, 3)$ с тремя степенями свободы

$$t(x, 0, 3) = \frac{6\sqrt{3}}{\pi(3 + t^2)^2} \quad (16)$$

- Распределение Пуассона

$$P(k, 10) = \frac{10^k}{k!} e^{-10} \quad (17)$$

- Равномерное распределение

$$U(x, -\sqrt{3}, \sqrt{3}) = \begin{cases} \frac{1}{2\sqrt{3}}, & |x| \leq \sqrt{3} \\ 0, & |x| > \sqrt{3} \end{cases} \quad (18)$$

7.2 Боксплот Тьюки

Боксплот (англ. box plot) — график, использующихся в описательной статистике, компактно изображающий одномерное распределение вероятностей. Такой вид диаграммы в удобной форме показывает медиану, нижний и верхний квартили и выбросы. Границами ящика служат первый и третий квартили, линия в середине ящика — медиана. Концы усов — края статистически значимой выборки (без выброса). Длину «усов» определяют разность первого квартиля и полутора межквартильных расстояний и сумма третьего квартиля и полутора межквартильных расстояний. Формула имеет вид

$$X_1 = Q_1 - \frac{3}{2}(Q_3 - Q_1), \quad X_2 = Q_3 + \frac{3}{2}(Q_3 - Q_1), \quad (19)$$

где X_1 — нижняя граница уса, X_2 — верхняя граница уса, Q_1 — первый квартиль, Q_3 — третий квартиль. Данные, выходящие за границы усов (выбросы), отображаются на графике в виде маленьких кружков. Выбросами считаются величины, такие что:

$$\begin{cases} x < X_1^T \\ x > X_2^T \end{cases} \quad (20)$$

7.3 Доверительные интервалы для параметров нормального распределения

Пусть $F_T(x)$ — функция распределения Стюдента с $n - 1$ степенями свободы. Полагая, что $2F_T(x) - 1 = 1 - \alpha$, где α — выбранный уровень значимости. Тогда $F_T(x) = 1 - \alpha/2$. Пусть $st_{1-\alpha/2}(n - 1)$ — квантиль распределения Стюдента с $n - 1$ степенями свободы и порядка $1 - \alpha/2$. Тогда получаем

$$P\left(\bar{x} - \frac{st_{1-\alpha/2}(n - 1)}{\sqrt{n - 1}} < m < \bar{x} + \frac{st_{1-\alpha/2}(n - 1)}{\sqrt{n - 1}}\right) = 1 - \alpha, \quad (21)$$

что и даст доверительный интервал для m с доверительной вероятностью $\gamma = 1 - \alpha$ для нормального распределения.

Случайная величина $n \frac{s^2}{\sigma^2}$ распределена по закону χ^2 с $n - 1$ степенями свободы. Тогда

$$P\left(\bar{x} - \frac{st_{1-\alpha/2}(n - 1)}{\sqrt{n - 1}} < m < \bar{x} + \frac{st_{1-\alpha/2}(n - 1)}{\sqrt{n - 1}}\right) = 1 - \alpha, \quad (22)$$

8 Описание работы

Лабораторные работы выполнены с использованием Python и его сторонних библиотек: `numpy`, `pandas`, `matplotlib`, `seaborn`.

Ссылка на GitHub репозиторий: <https://github.com/vladimir-skvortsov/spbstu-mathematical-statistics>

9 Результаты

9.1 Гистограммы и графики плотности распределения

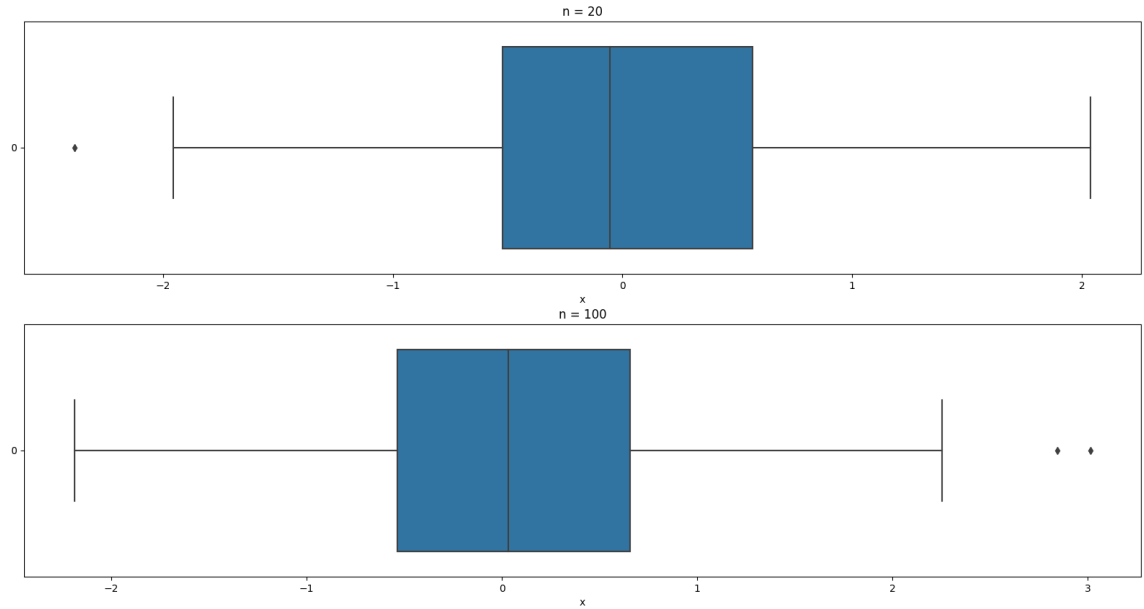


Рис. 6: Нормальное распределение (14)

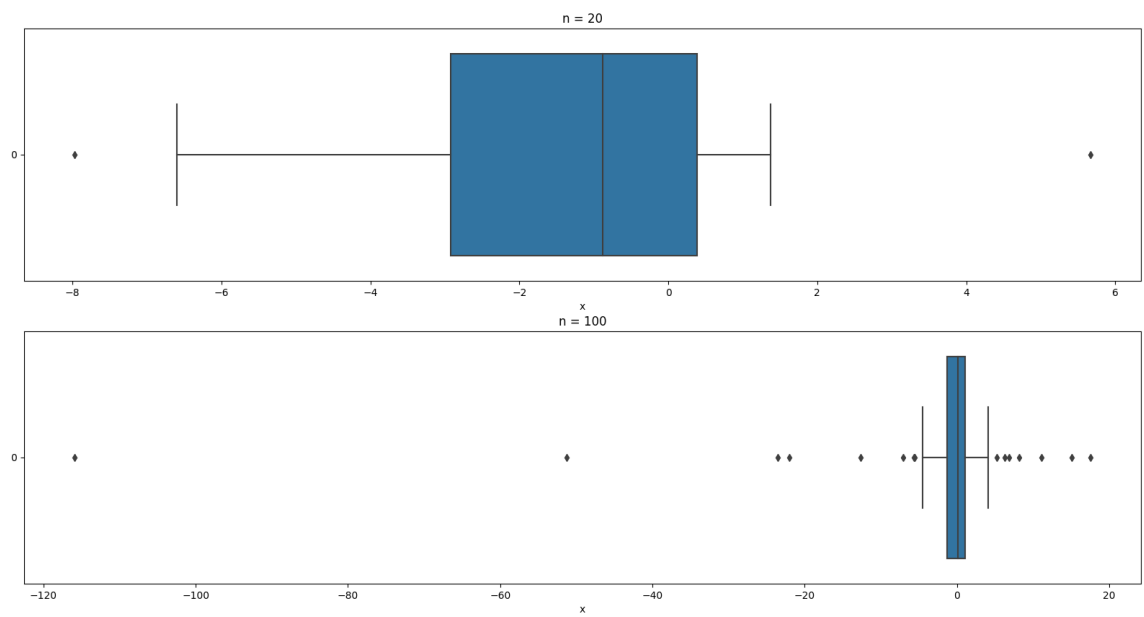


Рис. 7: Распределение Коши (15)

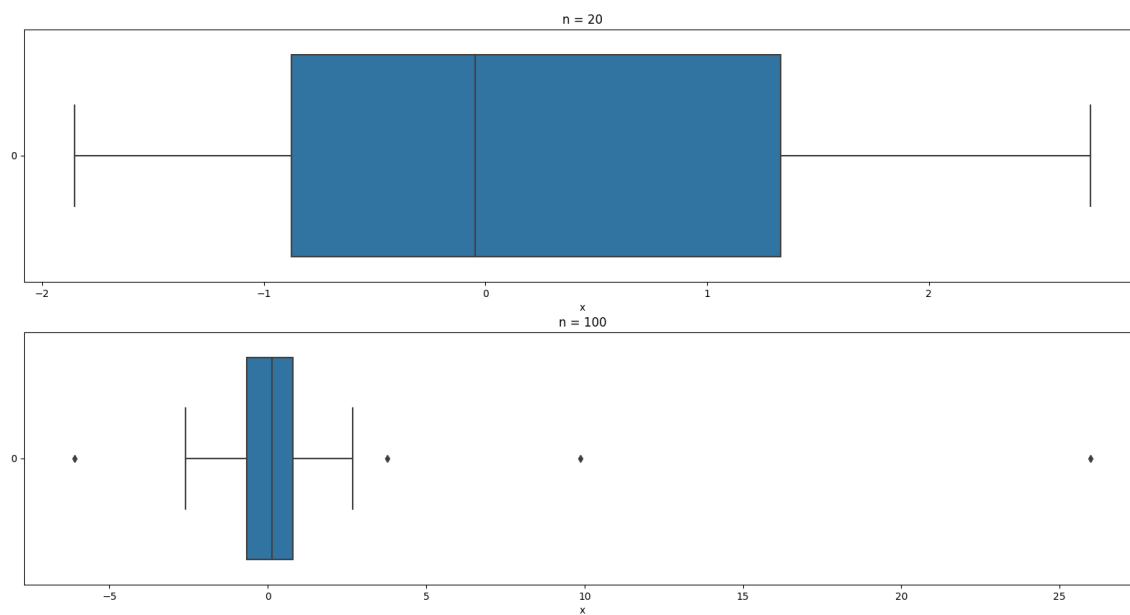


Рис. 8: Распределение Стьюдента (16)

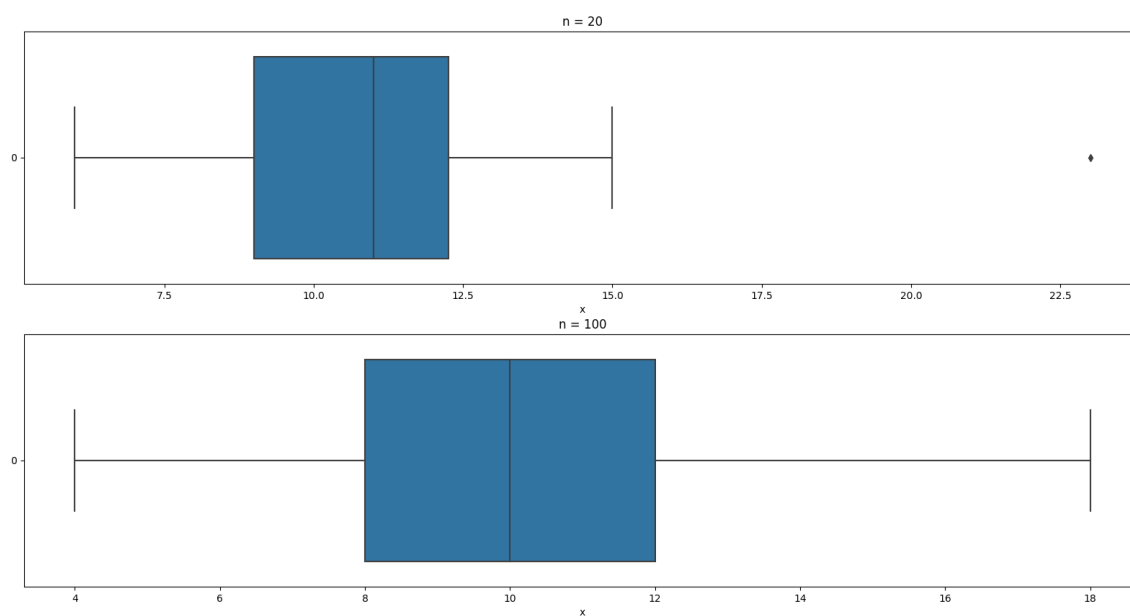


Рис. 9: Распределение Пуассона (17)

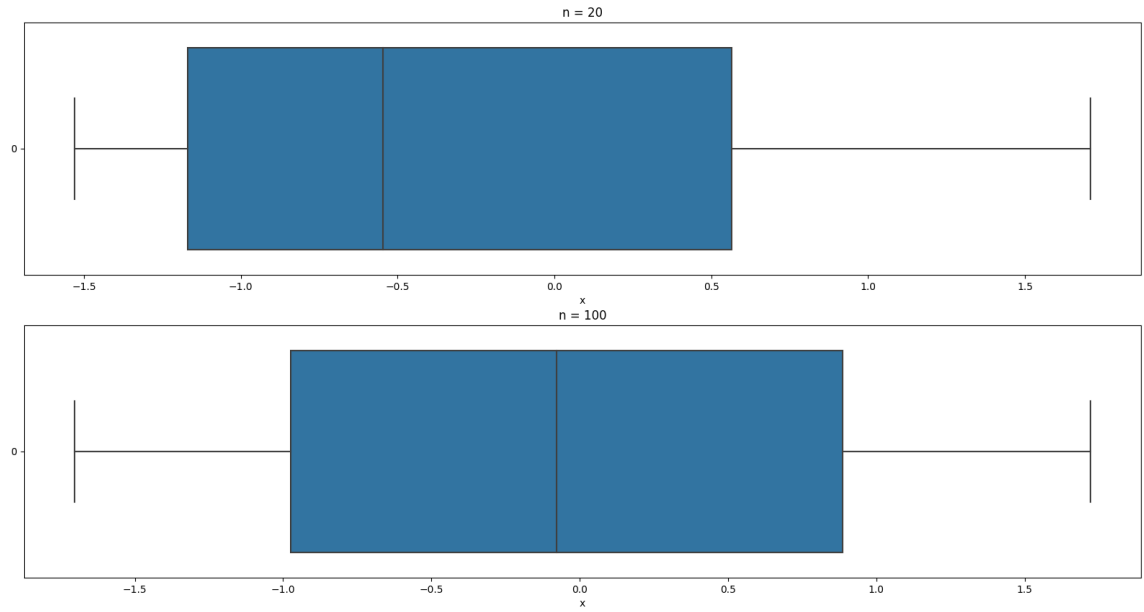


Рис. 10: Равномерное распределение (18)

9.2 Доверительные интервалы для параметров распределений

n = 20	m	σ
	$-0.43 < m < 0.37$	$0.66 < \sigma < 1.25$
n = 100	m	σ
	$-0.12 < m < 0.24$	$0.81 < \sigma < 1.07$

Таблица 6: Доверительные интервалы для параметров нормального распределения (14)

n = 20	m	σ
	$0.11 < m < 0.97$	$0.29 < \sigma < 0.33$
n = 100	m	σ
	$0.30 < m < 0.67$	$0.28 < \sigma < 0.33$

Таблица 7: Доверительные интервалы для параметров произвольного распределения. Асимптотический подход

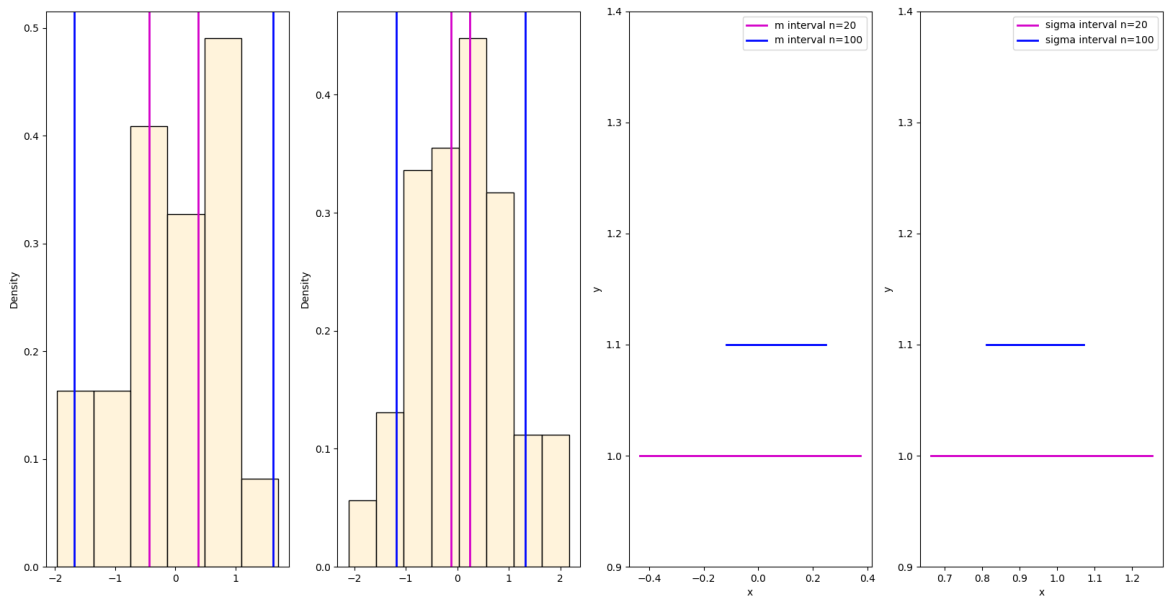


Рис. 11: Гистограммы и оценки для параметров нормального распределения

10 Выводы

По результатам выполнения лабораторной работы были сгенерированы выборки размером 20 и 100 элементов и построены для них боксплоты Тьюки.

Боксплот позволяет наглядно представить основные характеристики выборки - медиану, квартили, межквартильный размах и выбросы. На основе построенных графиков можно увидеть разницу в распределении данных для двух выборок. Для выборки размером в 100 элементов представленные метрики имеют более проработанный вид, ведь с увеличением размера выборки улучшается точность оценок параметров распределения.

Также в ходе выполнения лабораторной работы были сгенерированы две выборки размерами 20 и 100 элементов для нормального и произвольного распределения. Затем для каждой из них были вычислены параметры распределения: среднее значение и дисперсия.

Результаты, представленные графически, демонстрируют, что количество элементов в выборке влияет на точность оценок параметров. Более большое количество наблюдений (т.е. 100 элементов) приводит к более точным и стабильным оценкам среднего и дисперсии, как для нормального, так и для произвольного распределения. Для выборки с меньшим количеством элементов (20 элементов) оценки могут сильно варьироваться в зависимости от конкретной выборки, что также наглядно отображено на графиках.

Лабораторная работа иллюстрирует важнейший статистический принцип: точность статистической оценки увеличивается с ростом объема выборки. Результаты этого исследования подчеркивают значимость использования достаточно больших выборок для надежного анализа данных.