

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ
ПЕТРА ВЕЛИКОГО

ИНСТИТУТ ПРИКЛАДНОЙ МАТЕМАТИКИ И МЕХАНИКИ
КАФЕДРА «ПРИКЛАДНАЯ МАТЕМАТИКА»

Отчёт
по лабораторным работам №1-4
по дисциплине
«Математическая статистика»

Выполнил студент:

...

группа: ...

Проверил:

к.ф.-м.н., доцент

Баженов Александр Николаевич

Санкт-Петербург
2020 г.

Содержание

1	Постановка задачи	4
2	Теория	5
2.1	Распределения	5
2.2	Гистограмма	5
2.2.1	Определение	5
2.2.2	Графическое описание	5
2.2.3	Использование	6
2.3	Вариационный ряд	6
2.4	Выборочные числовые характеристики	6
2.4.1	Характеристики положения	6
2.4.2	Характеристики рассеяния	7
2.5	Боксплот Тьюки	7
2.5.1	Определение	7
2.5.2	Описание	7
2.5.3	Построение	7
2.6	Теоретическая вероятность выбросов	8
2.7	Эмпирическая функция распределения	8
2.7.1	Статистический ряд	8
2.7.2	Определение	9
2.7.3	Описание	9
2.8	Оценки плотности вероятности	9
2.8.1	Определение	9
2.8.2	Ядерные оценки	9
3	Реализация	10
4	Результаты	11
4.1	Гистограмма и график плотности распределения	11
4.2	Характеристики положения и рассеяния	12
4.3	Боксплот Тьюки	15
4.4	Доля выбросов	17
4.5	Теоретическая вероятность выбросов	18
4.6	Эмпирическая функция распределения	18
4.7	Ядерные оценки плотности распределения	20
5	Обсуждение	27
5.1	Гистограмма и график плотности распределения	27
5.2	Характеристики положения и рассеяния	27
5.3	Доля и теоретическая вероятность выбросов	27
5.4	Эмпирическая функция и ядерные оценки плотности распределения	27
	Литература	27

Список иллюстраций

1	Нормальное распределение	11
2	Распределение Коши	11
3	Распределение Лапласа	11
4	Распределение Пуассона	12
5	Равномерное распределение	12
6	Нормальное распределение	15
7	Распределение Коши	15
8	Распределение Лапласа	16
9	Распределение Пуассона	16
10	Равномерное распределение	17
11	Нормальное распределение	18
12	Распределение Коши	18
13	Распределение Лапласа	19
14	Распределение Пуассона	19
15	Равномерное распределение	19
16	Нормальное распределение, $n = 20$	20
17	Нормальное распределение, $n = 60$	20
18	Нормальное распределение, $n = 100$	21
19	Распределение Коши, $n = 20$	21
20	Распределение Коши, $n = 60$	22
21	Распределение Коши, $n = 100$	22
22	Распределение Лапласа, $n = 20$	23
23	Распределение Лапласа, $n = 60$	23
24	Распределение Лапласа, $n = 100$	24
25	Распределение Пуассона, $n = 20$	24
26	Распределение Пуассона, $n = 60$	25
27	Распределение Пуассона, $n = 100$	25
28	Равномерное распределение, $n = 20$	26
29	Равномерное распределение, $n = 60$	26
30	Равномерное распределение, $n = 100$	27

Список таблиц

1	Статистический ряд	8
2	Таблица распределения	9
3	Нормальное распределение	12
4	Распределение Коши	13
5	Распределение Лапласа	13
6	Распределение Пуассона	14
7	Равномерное распределение	14
8	Доля выбросов	17
9	Теоретическая вероятность выбросов	18

1 Постановка задачи

Для 5 распределений:

- Нормальное распределение $N(x, 0, 1)$
 - Распределение Коши $C(x, 0, 1)$
 - Распределение Лапласа $L(x, 0, \frac{1}{\sqrt{2}})$
 - Распределение Пуассона $P(k, 10)$
 - Равномерное распределение $U(x, -\sqrt{3}, \sqrt{3})$
1. Сгенерировать выборки размером 10, 50 и 1000 элементов.
Построить на одном рисунке гистограмму и график плотности распределения.
 2. Сгенерировать выборки размером 10, 100 и 1000 элементов.
Для каждой выборки вычислить следующие статистические характеристики положения данных: $\bar{x}, med\ x, z_R, z_Q, z_{tr}$. Повторить такие вычисления 1000 раз для каждой выборки и найти среднее характеристик положения и их квадратов:

$$E(z) = \bar{z} \quad (1)$$

Вычислить оценку дисперсии по формуле:

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

Представить полученные данные в виде таблиц.

3. Сгенерировать выборки размером 20 и 100 элементов.
Построить для них боксплот Тьюки.
Для каждого распределения определить долю выбросов экспериментально (сгенерировав выборку, соответствующую распределению 1000 раз, и вычислив среднюю долю выбросов) и сравнить с результатами, полученными теоретически.
4. Сгенерировать выборки размером 20, 60 и 100 элементов.
Построить на них эмпирические функции распределения и ядерные оценки плотности распределения на отрезке $[-4; 4]$ для непрерывных распределений и на отрезке $[6; 14]$ для распределения Пуассона.

2 Теория

2.1 Распределения

- Нормальное распределение

$$N(x, 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad (3)$$

- Распределение Коши

$$C(x, 0, 1) = \frac{1}{\pi} \frac{1}{x^2 + 1} \quad (4)$$

- Распределение Лапласа

$$L(x, 0, \frac{1}{\sqrt{2}}) = \frac{1}{\sqrt{2}} e^{-\sqrt{2}|x|} \quad (5)$$

- Распределение Пуассона

$$P(k, 10) = \frac{10^k}{k!} e^{-10} \quad (6)$$

- Равномерное распределение

$$U(x, -\sqrt{3}, \sqrt{3}) = \begin{cases} \frac{1}{2\sqrt{3}} & \text{при } |x| \leq \sqrt{3} \\ 0 & \text{при } |x| > \sqrt{3} \end{cases} \quad (7)$$

2.2 Гистограмма

2.2.1 Определение

Гистограмма в математической статистике — это функция, приближающая плотность вероятности некоторого распределения, построенная на основе выборки из него [1].

2.2.2 Графическое описание

Графически гистограмма строится следующим образом. Сначала множество значений, которое может принимать элемент выборки, разбивается на несколько интервалов. Чаще всего эти интервалы берут одинаковыми, но это не является строгим требованием. Эти интервалы откладываются на горизонтальной оси, затем над каждым рисуется прямоугольник. Если все интервалы были одинаковыми, то высота каждого прямоугольника пропорциональна числу элементов выборки, попадающих в соответствующий интервал. Если интервалы разные, то высота прямоугольника выбирается таким образом, чтобы его площадь была пропорциональна числу элементов выборки, которые попали в этот интервал [1].

2.2.3 Использование

Гистограммы применяются в основном для визуализации данных на начальном этапе статистической обработки.

Построение гистограмм используется для получения эмпирической оценки плотности распределения случайной величины. Для построения гистограммы наблюдаемый диапазон изменения случайной величины разбивается на несколько интервалов и подсчитывается доля от всех измерений, попавшая в каждый из интервалов. Величина каждой доли, отнесенная к величине интервала, принимается в качестве оценки значения плотности распределения на соответствующем интервале [1].

2.3 Вариационный ряд

Вариационным рядом называется последовательность элементов выборки, расположенных в неубывающем порядке. Одинаковые элементы повторяются [2, с. 409].

Запись вариационного ряда: $x_{(1)}, x_{(2)}, \dots, x_{(n)}$.

Элементы вариационного ряда $x_{(i)}$ ($i = 1, 2, \dots, n$) называются порядковыми статистиками.

2.4 Выборочные числовые характеристики

С помощью выборки образуются её числовые характеристики. Это числовые характеристики дискретной случайной величины X^* , принимающей выборочные значения x_1, x_2, \dots, x_n [2, с. 411].

2.4.1 Характеристики положения

- Выборочное среднее

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (8)$$

- Выборочная медиана

$$\text{med } x = \begin{cases} x_{(l+1)} & \text{при } n = 2l + 1 \\ \frac{x_{(l)} + x_{(l+1)}}{2} & \text{при } n = 2l \end{cases} \quad (9)$$

- Полусумма экстремальных выборочных элементов

$$z_R = \frac{x_{(1)} + x_{(n)}}{2} \quad (10)$$

- Полусумма квартилей

Выборочная квартиль z_p порядка p определяется формулой

$$z_p = \begin{cases} x_{([np]+1)} & \text{при } np \text{ дробном,} \\ x_{(np)} & \text{при } np \text{ целом.} \end{cases} \quad (11)$$

Полусумма квартилей

$$z_Q = \frac{z_{1/4} + z_{3/4}}{2} \quad (12)$$

- Усечённое среднее

$$z_{tr} = \frac{1}{n-2r} \sum_{i=r+1}^{n-r} x_{(i)}, \quad r \approx \frac{n}{4} \quad (13)$$

2.4.2 Характеристики рассеяния

Выборочная дисперсия

$$D = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (14)$$

2.5 Боксплот Тьюки

2.5.1 Определение

Боксплот (англ. box plot) — график, использующийся в описательной статистике, компактно изображающий одномерное распределение вероятностей.

2.5.2 Описание

Такой вид диаграммы в удобной форме показывает медиану, нижний и верхний квартили и выбросы. Несколько таких ящичков можно нарисовать бок о бок, чтобы визуально сравнивать одно распределение с другим; их можно располагать как горизонтально, так и вертикально. Расстояния между различными частями ящичка позволяют определить степень разброса (дисперсии) и асимметрии данных и выявить выбросы [3].

2.5.3 Построение

Границами ящичка служат первый и третий квартили, линия в середине ящичка — медиана. Концы усов — края статистически значимой выборки (без выбросов). Длину «усов» определяют разность первого квартиля и полутора межквартильных расстояний и сумма третьего квартиля и полутора межквартильных расстояний. Формула имеет вид

$$X_1 = Q_1 - \frac{3}{2}(Q_3 - Q_1), \quad X_2 = Q_3 + \frac{3}{2}(Q_3 - Q_1), \quad (15)$$

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 [3].

2.6 Теоретическая вероятность выбросов

Встроенными средствами языка программирования R в среде разработки RStudio можно вычислить теоретические первый и третий квартили распределений (Q_1^T и Q_3^T соответственно). По формуле (15) можно вычислить теоретические нижнюю и верхнюю границы уса (X_1^T и X_2^T соответственно). Выбросами считаются величины x , такие что:

$$\begin{cases} x < X_1^T \\ x > X_2^T \end{cases} \quad (16)$$

Теоретическая вероятность выбросов для непрерывных распределений

$$P_b^T = P(x < X_1^T) + P(x > X_2^T) = F(X_1^T) + (1 - F(X_2^T)), \quad (17)$$

где $F(X) = P(x \leq X)$ - функция распределения.

Теоретическая вероятность выбросов для дискретных распределений

$$P_b^T = P(x < X_1^T) + P(x > X_2^T) = (F(X_1^T) - P(x = X_1^T)) + (1 - F(X_2^T)), \quad (18)$$

где $F(X) = P(x \leq X)$ - функция распределения.

2.7 Эмпирическая функция распределения

2.7.1 Статистический ряд

Статистическим рядом называется последовательность различных элементов выборки z_1, z_2, \dots, z_k , расположенных в возрастающем порядке с указанием частот n_1, n_2, \dots, n_k , с которыми эти элементы содержатся в выборке. Статистический ряд обычно записывается в виде таблицы

z	z_1	z_2	\dots	z_k
n	n_1	n_2	\dots	n_k

Таблица 1: Статистический ряд

2.7.2 Определение

Эмпирической (выборочной) функцией распределения (э. ф. р.) называется относительная частота события $X < x$, полученная по данной выборке:

$$F_n^*(x) = P^*(X < x). \quad (19)$$

2.7.3 Описание

Для получения относительной частоты $P^*(X < x)$ просуммируем в статистическом ряде, построенном по данной выборке, все частоты n_i , для которых элементы z_i статистического ряда меньше x . Тогда $P^*(X < x) = \frac{1}{n} \sum_{z_i < x} n_i$. Получаем

$$F^*(x) = \frac{1}{n} \sum_{z_i < x} n_i. \quad (20)$$

$F^*(x)$ — функция распределения дискретной случайной величины X^* , заданной таблицей распределения

X^*	z_1	z_2	\dots	z_k
P	$\frac{n_1}{n}$	$\frac{n_2}{n}$	\dots	$\frac{n_k}{n}$

Таблица 2: Таблица распределения

Эмпирическая функция распределения является оценкой, т. е. приближённым значением, генеральной функции распределения

$$F_n^*(x) \approx F_X(x). \quad (21)$$

2.8 Оценки плотности вероятности

2.8.1 Определение

Оценкой плотности вероятности $f(x)$ называется функция $\hat{f}(x)$, построенная на основе выборки, приближённо равная $f(x)$

$$\hat{f}(x) \approx f(x). \quad (22)$$

2.8.2 Ядерные оценки

Представим оценку в виде суммы с числом слагаемых, равным объёму выборки:

$$\hat{f}_n(x) = \frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{x-x_i}{h_n}\right). \quad (23)$$

Здесь функция $K(u)$, называемая ядерной (ядром), непрерывна и является плотностью вероятности, x_1, \dots, x_n — элементы выборки, $\{h_n\}$ — любая последовательность положительных чисел, обладающая свойствами

$$h_n \xrightarrow{n \rightarrow \infty} 0; \quad \frac{h_n}{n^{-1}} \xrightarrow{n \rightarrow \infty} \infty. \quad (24)$$

Такие оценки называются непрерывными ядерными [2, с. 421-423].

Замечание. Свойство, означающее сближение оценки с оцениваемой величиной при $n \rightarrow \infty$ в каком-либо смысле, называется состоятельностью оценки.

Если плотность $f(x)$ кусочно-непрерывная, то ядерная оценка плотности является состоятельной при соблюдении условий, накладываемых на параметр сглаживания h_n , а также на ядро $K(u)$.

Гауссово (нормальное) ядро [4, с. 38]

$$K(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}. \quad (25)$$

Правило Сильвермана [4, с. 44]

$$h_n = 1.06 \hat{\sigma} n^{-1/5}, \quad (26)$$

где $\hat{\sigma}$ - выборочное стандартное отклонение.

3 Реализация

Лабораторная работа выполнена с помощью встроенных средств языка программирования R в среде разработки RStudio. Исходный код лабораторной работы приведён в приложении.

4 Результаты

4.1 Гистограмма и график плотности распределения

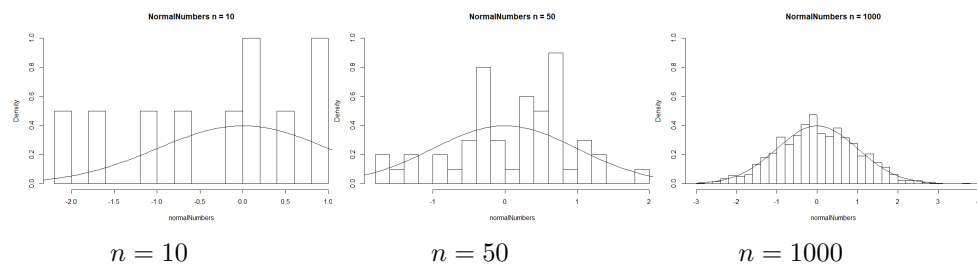


Рис. 1: Нормальное распределение

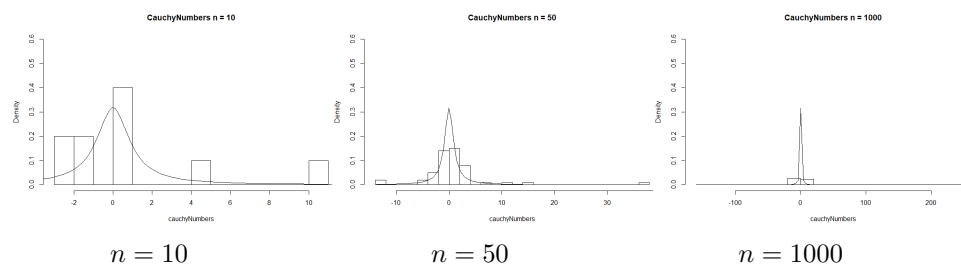


Рис. 2: Распределение Коши

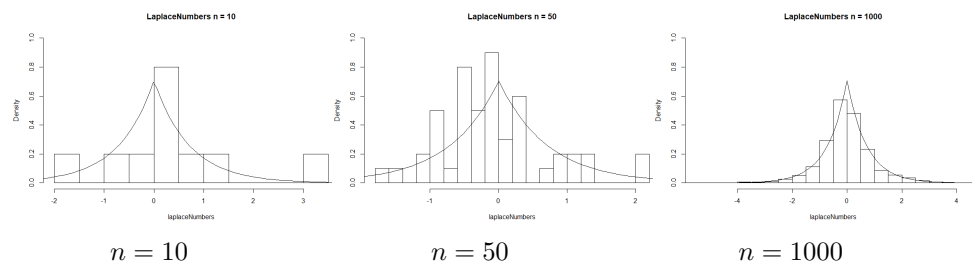


Рис. 3: Распределение Лапласа

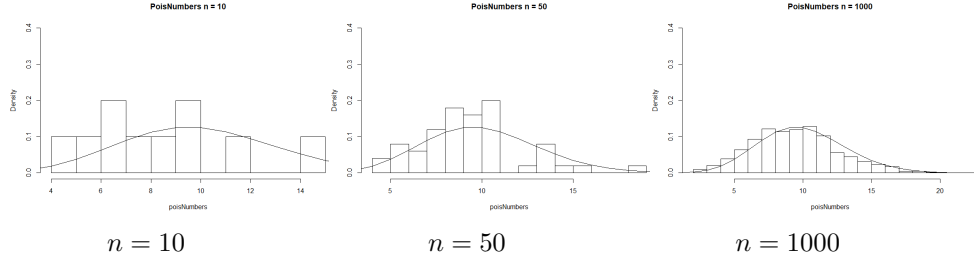


Рис. 4: Распределение Пуассона

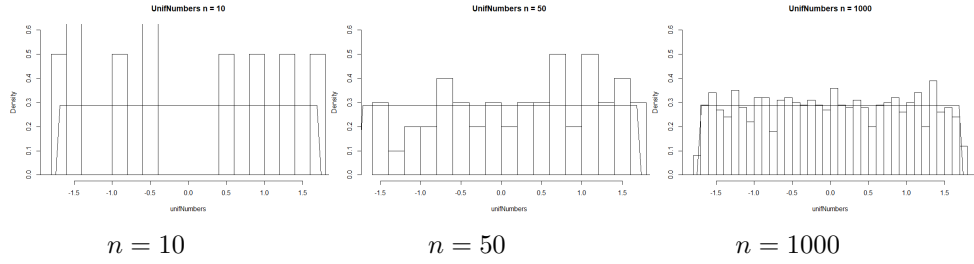


Рис. 5: Равномерное распределение

4.2 Характеристики положения и рассеяния

normal n = 10					
	\bar{x} (8)	$med\ x$ (9)	z_R (10)	z_Q (12)	z_{tr} (13)
$E(z)$ (1)	0.012	0.017	0.014	0.007	0.013
$D(z)$ (2)	0.097	0.136	0.201	0.112	0.121
normal n = 100					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	0.002	-0.002	0.016	-0.014	0.000
$D(z)$	0.010	0.016	0.085	0.013	0.012
normal n = 1000					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	0.000894	0.000890	-0.004452	-0.000313	0.001087
$D(z)$	0.00095	0.00150	0.06169	0.00123	0.00116

Таблица 3: Нормальное распределение

cauchy n = 10					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	0.961	0.015	4.624	0.070	0.022
$D(z)$	857.404	0.301	21094.729	1.202	0.338
cauchy n = 100					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	1.1034	-0.0071	52.9826	-0.0393	-0.0067
$D(z)$	370.392	0.024	899464.699	0.054	0.026
cauchy n = 1000					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	-0.567	-0.002	-303.307	-0.004	-0.001
$D(z)$	1192.6058	0.0026	297518700	0.0054	0.0028

Таблица 4: Распределение Коши

laplace n = 10					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	-0.008	-0.002	-0.021	-0.000	-0.001
$D(z)$	0.098	0.070	0.402	0.099	0.067
laplace n = 100					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	-0.004	-0.001	-0.046	-0.015	-0.002
$D(z)$	0.011	0.006	0.416	0.010	0.007
laplace n = 1000					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	-0.0012	-0.0004	-0.0038	-0.0025	-0.0010
$D(z)$	0.00102	0.00053	0.40565	0.00104	0.00064

Таблица 5: Распределение Лапласа

pois n = 10					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	10.017	9.858	10.329	9.952	9.872
$D(z)$	1.087	1.607	1.923	1.351	1.387
pois n = 100					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	9.986	9.828	10.924	9.854	9.843
$D(z)$	0.095	0.223	0.987	0.142	0.115
pois n = 1000					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	9.999	9.996	11.636	9.995	9.857
$D(z)$	0.010	0.004	0.596	0.003	0.011

Таблица 6: Распределение Пуассона

unif n = 10					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	0.0076	0.0084	0.0079	0.0063	0.0113
$D(z)$	0.099	0.221	0.046	0.134	0.188
unif n = 100					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	0.004	0.008	0.000	-0.014	0.007
$D(z)$	0.009	0.028	0.001	0.014	0.018
unif n = 1000					
	\bar{x}	$med\ x$	z_R	z_Q	z_{tr}
$E(z)$	-0.0012	-0.0019	0.0001	-0.0033	-0.0020
$D(z)$	0.0009	0.0027	0.0000	0.0015	0.0018

Таблица 7: Равномерное распределение

4.3 Боксплот Тьюки

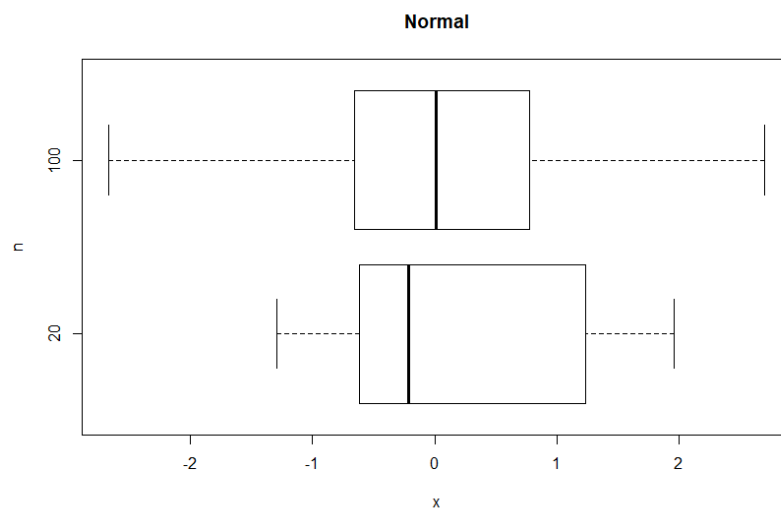


Рис. 6: Нормальное распределение

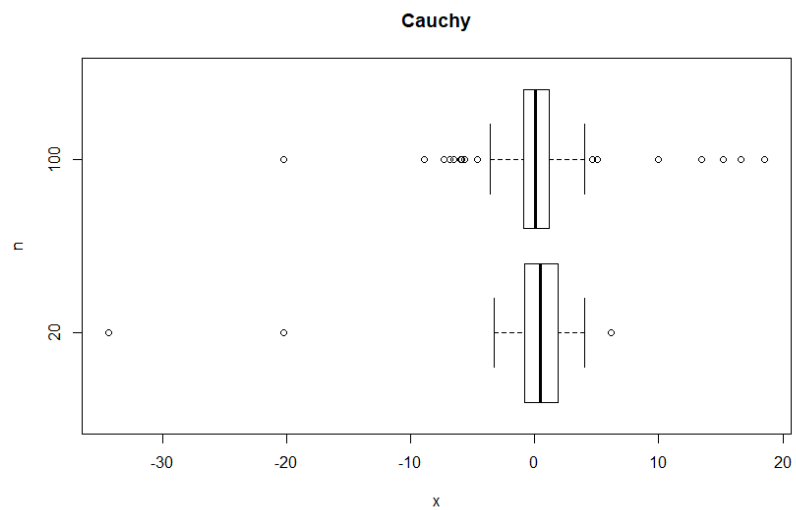


Рис. 7: Распределение Коши

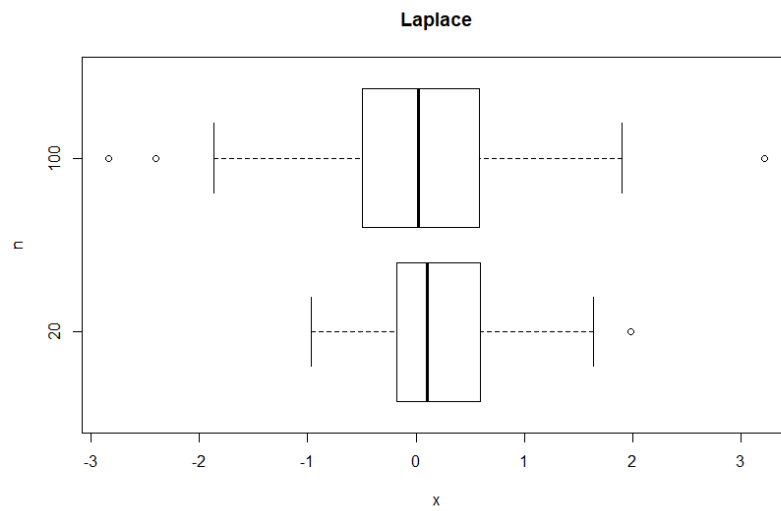


Рис. 8: Распределение Лапласа

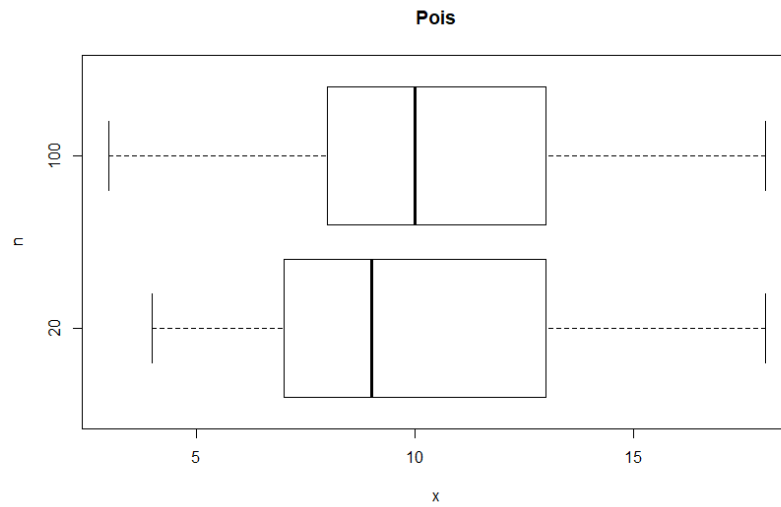


Рис. 9: Распределение Пуассона

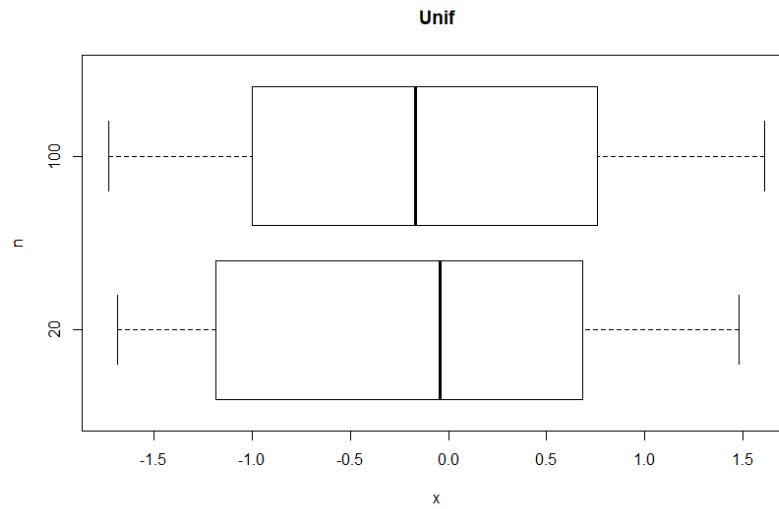


Рис. 10: Равномерное распределение

4.4 Доля выбросов

Выборка	Доля выбросов
normal n = 20	0.02
normal n = 100	0.01
cauchy n = 20	0.15
cauchy n = 100	0.16
laplace n = 20	0.07
laplace n = 100	0.06
pois n = 20	0.02
pois n = 100	0.01
unif n = 20	0
unif n = 100	0

Таблица 8: Доля выбросов

4.5 Теоретическая вероятность выбросов

Распределение	Q_1^T	Q_3^T	X_1^T (15)	X_2^T (15)	P_B^T (17), (18)
Нормальное распределение	-0.674	0.674	-2.698	2.698	0.007
Распределение Коши	-1	1	-4	4	0.156
Распределение Лапласа	-0.490	0.490	-1.961	1.961	0.063
Распределение Пуассона	8	12	2	18	0.008
Равномерное распределение	-0.866	0.866	-3.464	3.464	0

Таблица 9: Теоретическая вероятность выбросов

4.6 Эмпирическая функция распределения

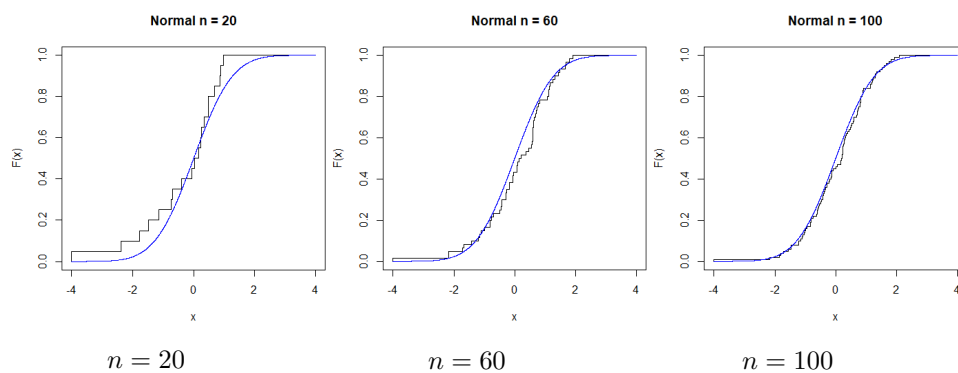


Рис. 11: Нормальное распределение

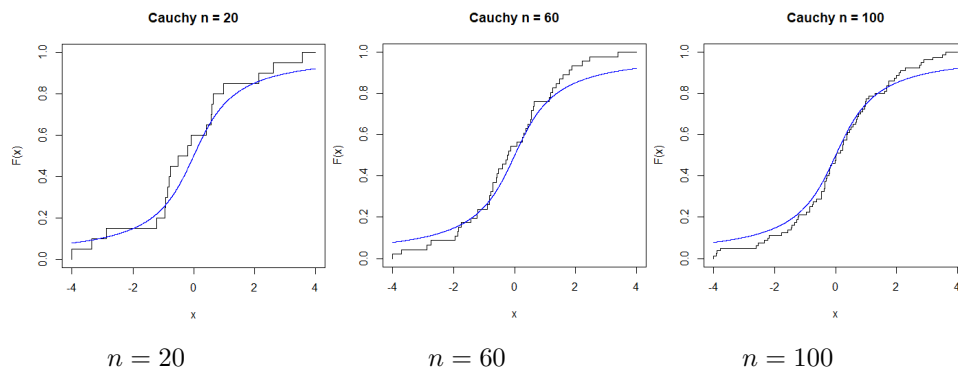


Рис. 12: Распределение Коши

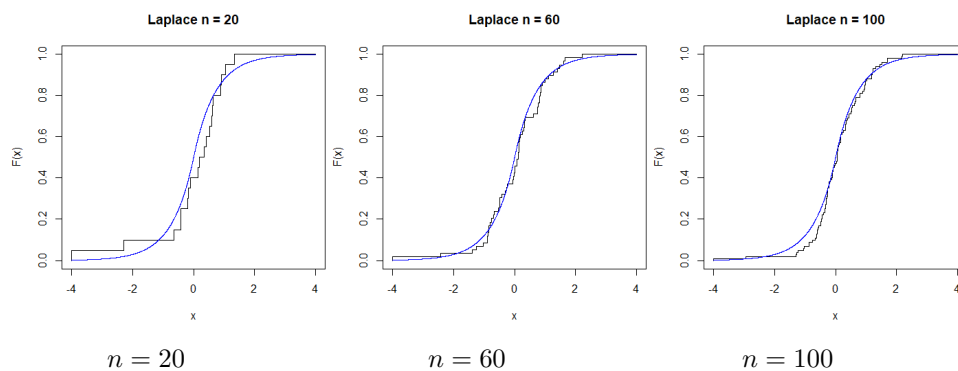


Рис. 13: Распределение Лапласа

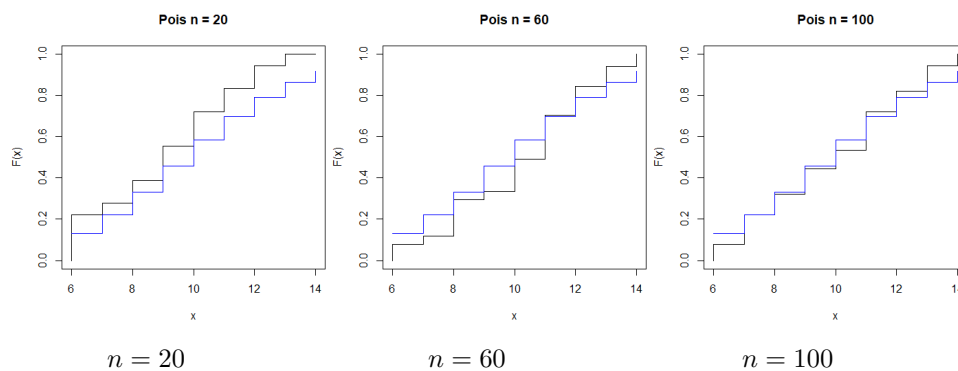


Рис. 14: Распределение Пуассона

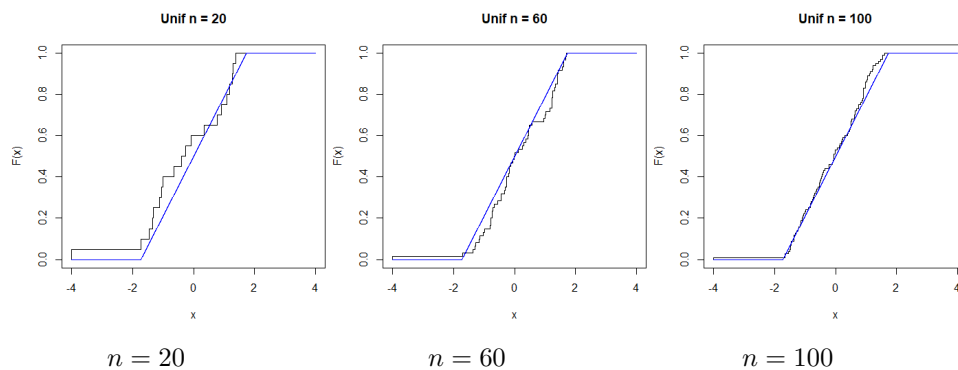


Рис. 15: Равномерное распределение

4.7 Ядерные оценки плотности распределения

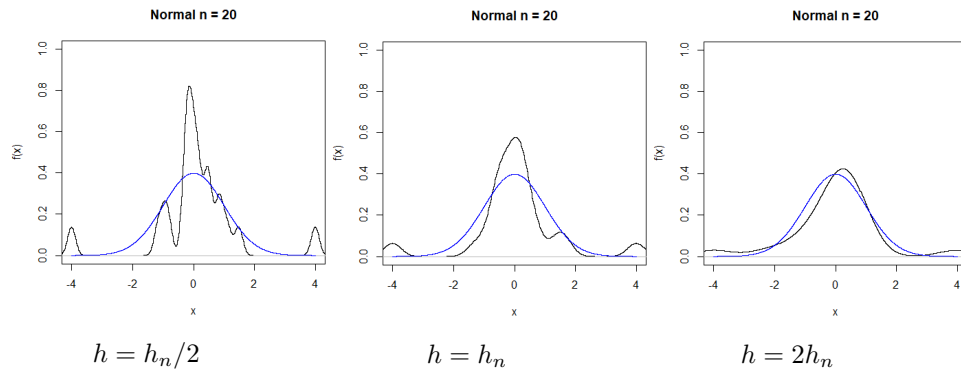


Рис. 16: Нормальное распределение, $n = 20$

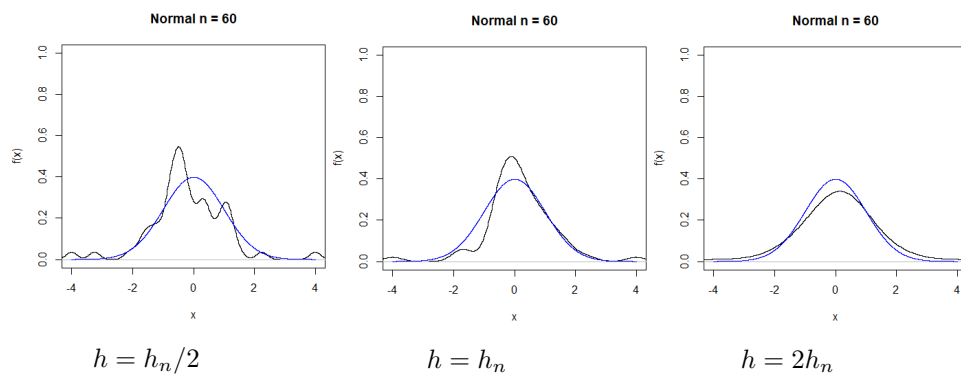


Рис. 17: Нормальное распределение, $n = 60$

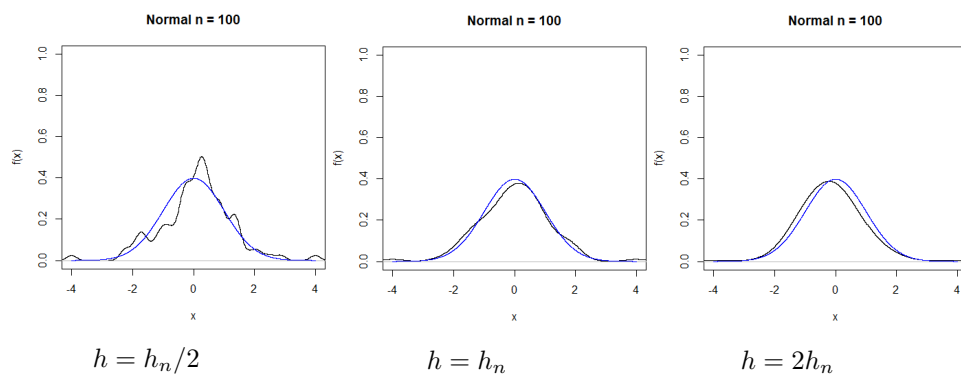


Рис. 18: Нормальное распределение, $n = 100$

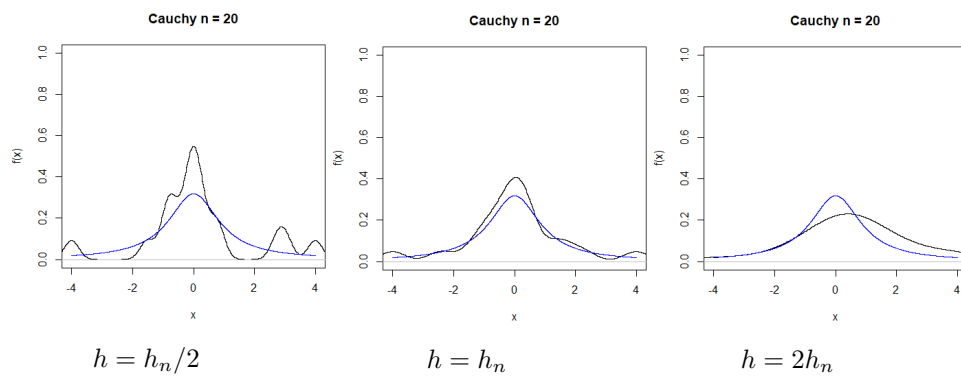


Рис. 19: Распределение Коши, $n = 20$

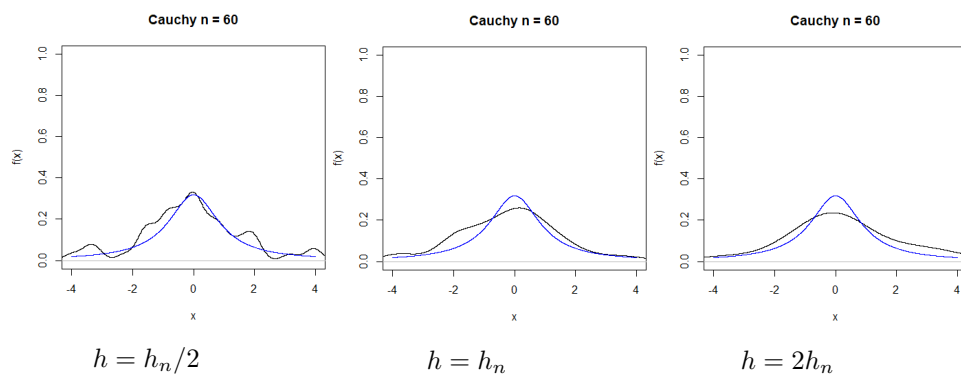


Рис. 20: Распределение Коши, $n = 60$

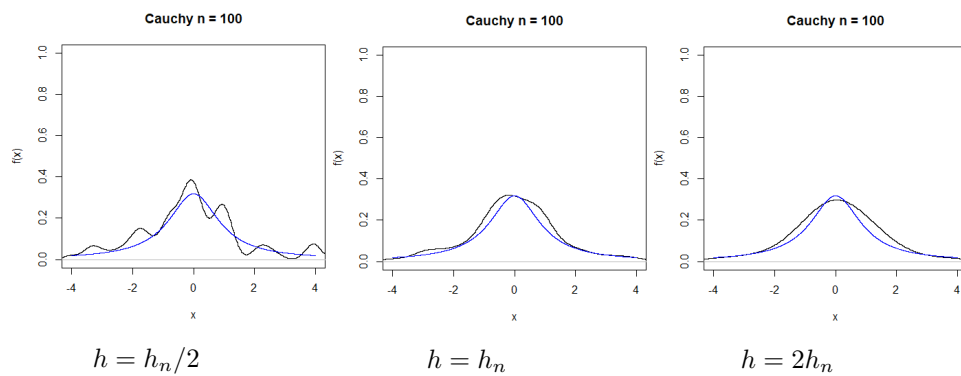


Рис. 21: Распределение Коши, $n = 100$

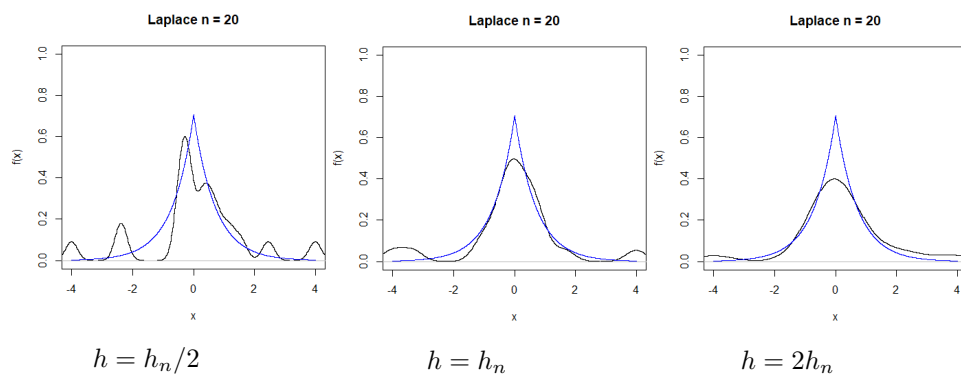


Рис. 22: Распределение Лапласа, $n = 20$

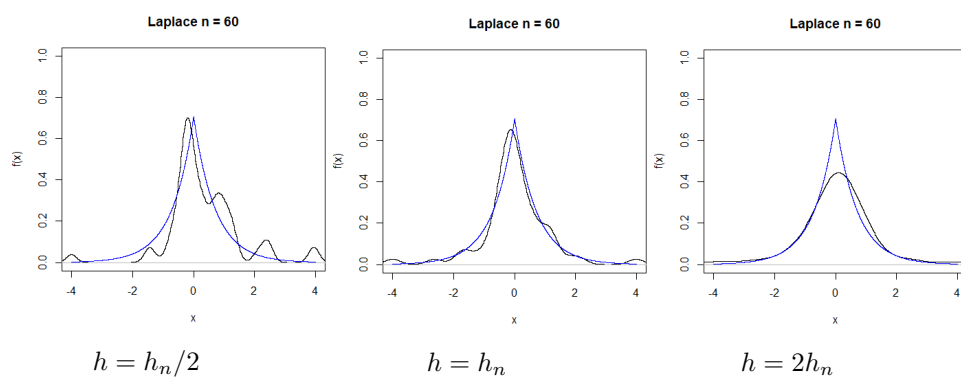


Рис. 23: Распределение Лапласа, $n = 60$

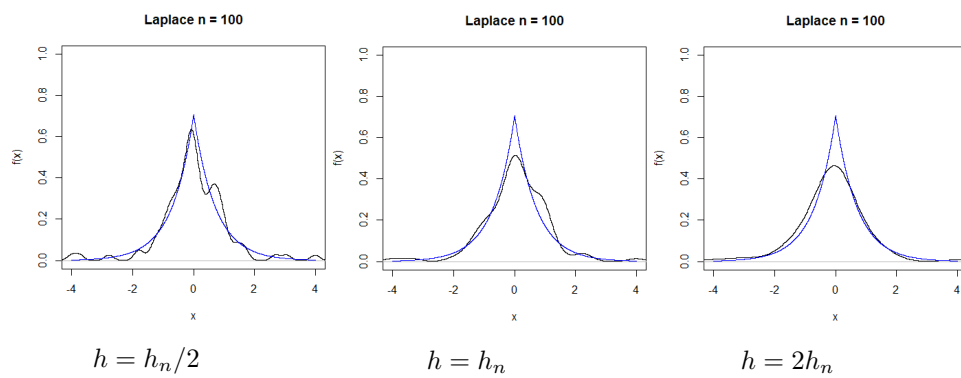


Рис. 24: Распределение Лапласа, $n = 100$

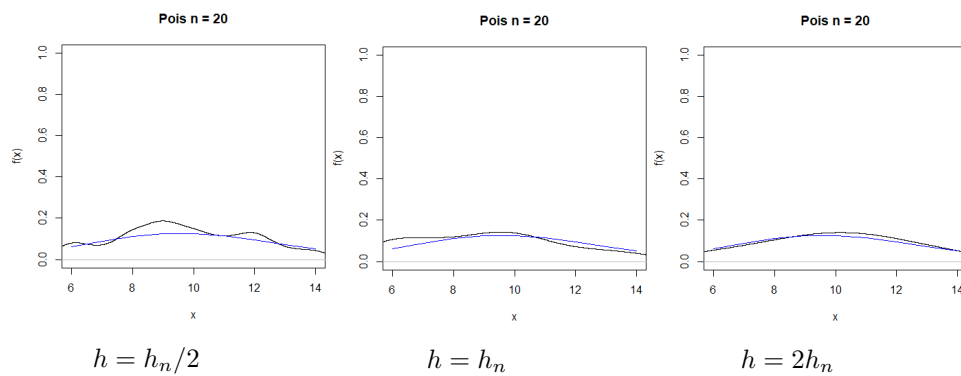


Рис. 25: Распределение Пуассона, $n = 20$

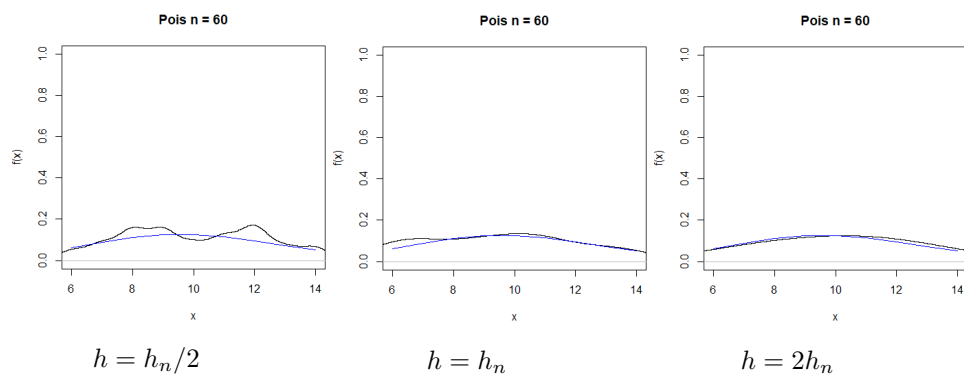


Рис. 26: Распределение Пуассона, $n = 60$

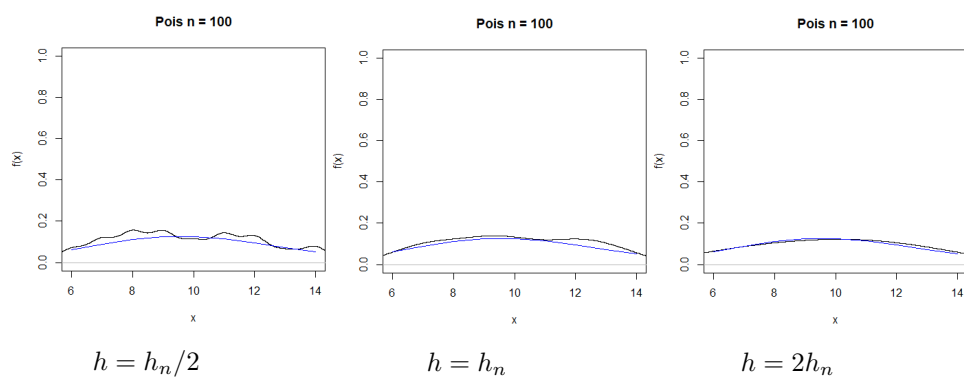


Рис. 27: Распределение Пуассона, $n = 100$

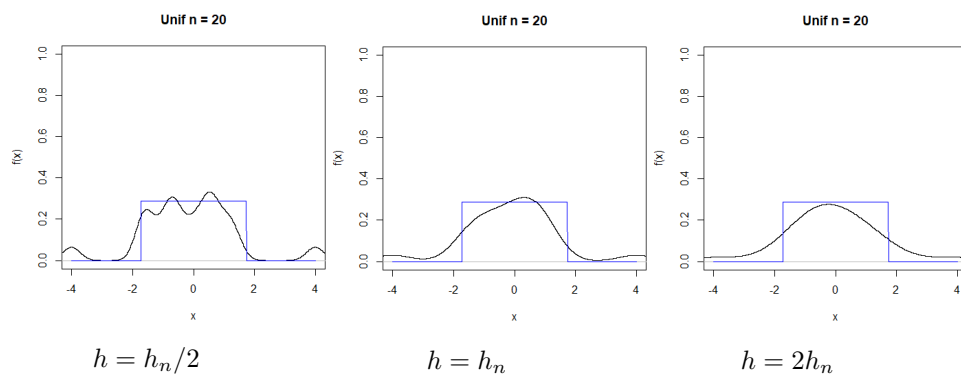


Рис. 28: Равномерное распределение, $n = 20$

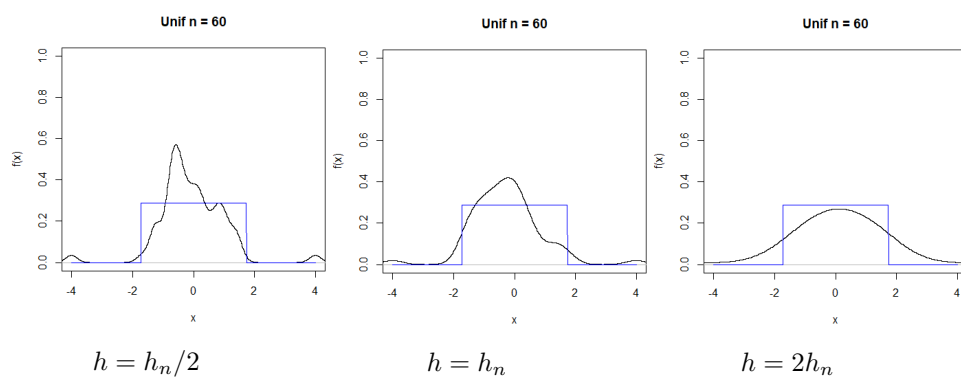


Рис. 29: Равномерное распределение, $n = 60$

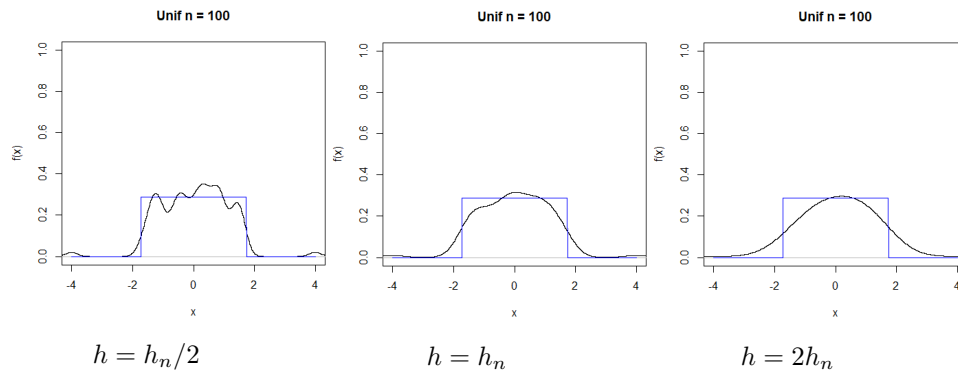


Рис. 30: Равномерное распределение, $n = 100$

5 Обсуждение

5.1 Гистограмма и график плотности распределения

5.2 Характеристики положения и рассеяния

5.3 Доля и теоретическая вероятность выбросов

5.4 Эмпирическая функция и ядерные оценки плотности распределения

Литература

- [1] Histogram. URL: <https://en.wikipedia.org/wiki/Histogram>
- [2] Вероятностные разделы математики. Учебник для бакалавров технических направлений. //Под ред. Максимова Ю.Д. — Спб.: «Иван Федоров», 2001. — 592 с., илл.
- [3] Box plot. URL: https://en.wikipedia.org/wiki/Box_plot
- [4] Анатольев, Станислав (2009) «Непараметрическая регрессия», Квантиль, №7, стр. 37-52.