

Занимательная статистика

Манга



マンガでわかる

統計学

高橋 信 / 著
トレンド・プロ / マンガ制作



Манга

Занимательная

СТАТИСТИКА

Син Такахаси

Перевод с японского

Захаровой Е. А., Коги Муцуми



Москва
Издательский дом «Додэка-XXI»
2010

УДК 311
ББК 60.6
Т15

Такахаси, Син.

Т15 Занимательная статистика. Манга / Син Такахаси ; пер. с яп. Захаровой Е. А., Коги Муцуми. — М. : Додэка-XXI, 2010. — 224 с. : ил. — (Серия «Образовательная Манга»). — Доп. тит. л. яп. — ISBN 978-5-94120-244-7.
И. Захарова, Е. А., пер.

Если тебя интересует статистика, или тебе просто нужно как-то обработать данные, то «Занимательная статистика» поможет тебе преодолеть чувство, что «ты плохо знаешь математику». Этот иллюстрированный путеводитель легко и непринуждённо проведёт тебя по пути познания статистики. А полученные знания ты сможешь закрепить с помощью упражнений, без которых, как известно, не обходится ни одна книга по математике.

Последуй за всегда невозмутимым Ямамото и ты увидишь, как он научит Руи:

- рассчитать среднее значение, медиану и стандартное отклонение результатов в боулинге;
- построить гистограмму цен на китайскую лапшу рамэн;
- определить вероятность получения проходного балла на экзаменах по математике;
- вычислить коэффициент Крамера, чтобы узнать, как предпочитают признаваться в любви юноши и девушки;
- узнать, как нормируются результаты тестов, когда учителя оценивают успеваемость.

Эти и другие примеры из реальной жизни позволят тебе с лёгкостью усвоить то, что многие находят трудным для понимания.

Если ты хочешь разобраться в статистике, но от обычных учебников статистики у тебя пухнет голова и клонит в сон, или если тебе просто нужно освежить забытые знания, пусть Ямамото-сан и Руи будут твоими гидами.

Книга будет полезна учащимся старших классов средних школ и колледжей, студентам вузов, а также всем, кто интересуется статистикой и хочет, чтобы обучение было лёгким и увлекательным.

УДК 311
ББК 60.6

Все права защищены. Никакая часть этого издания не может быть воспроизведена в любой форме или любыми средствами, электронными или механическими, включая фотографирование, ксерокопирование или иные средства копирования или сохранения информации, без письменного разрешения издательства.

ISBN 978-5-94120-244-7 (рус.)
ISBN 978-4-27406-570-5 (яп.)

© Син Такахаси, Trend-Pro Co., LTD.
© Издательский дом «Додэка-XXI», 2010
© Серия «Образовательная Манга»

Содержание

Предисловие	xiii
Пролог. Любовь и статистика	1
Глава 1. Разберёмся с типами данных	13
1. Количественные и качественные данные	14
2. Примеры качественных данных	20
3. Использование многовариантных ответов на практике	28
Упражнение	29
Ответ	29
Выводы	29
Глава 2. Знакомимся с количественными данными	31
1. Ряды распределения и гистограммы	32
2. Средняя величина	40
3. Медиана	44
4. Стандартное отклонение	48
5. Ряды распределения и величина интервала	54
6. Теория оценивания и описательная статистика	57
Упражнение	57
Ответ	58
Выводы	58
Глава 3. Знакомимся с качественными данными	59
1. Простые статистические таблицы	60
Упражнение	64
Ответ	64
Выводы	64
Глава 4. Нормированное отклонение и рейтинг успеваемости	65
1. Нормирование и нормированное отклонение	66
2. Свойства нормированного отклонения	73
3. Рейтинг успеваемости	74
4. Что такое рейтинг успеваемости?	76
Упражнение	78
Ответ	79
Выводы	80
Глава 5. Вычислим вероятность	81
1. Функция распределения плотности вероятности	82
2. Нормальное распределение	86
3. Стандартное нормальное распределение	89
Пример 1	95
Пример 2	97
4. Распределение хи-квадрат	99

5. Распределение Стьюдента	106
6. Распределение Фишера, или F-распределение	106
7. Распределения и Excel.	107
Упражнение	108
Ответ	108
Выводы	109
Глава 6. Что может связывать две переменные	111
1. Коэффициент линейной корреляции	116
2. Коэффициент корреляции между данными разных типов	121
3. Коэффициент корреляции Крамера.	127
Упражнение	138
Решение.	139
Выводы	142
Глава 7. А что это за проверка гипотезы о независимости?	143
1. Проверка гипотезы	144
2. Проверка гипотезы о независимости.	151
Объяснение.	152
Упражнение.	157
Размышление	158
Вывод.	160
3. Нулевая и альтернативная гипотезы	170
4. Р-значение и порядок проверки	175
5. Проверка гипотезы о независимости и гипотезы об однородности	184
Упражнение.	184
Решение.	185
6. Как выразить словами вывод на основании проверки	187
Упражнение.	188
Ответ	188
Выводы	189
Приложение. Попробуем вычислить с помощью Excel	191
1. Построение таблиц распределения	192
2. Вычисление среднего значения, медианы и стандартного отклонения	195
3. Построение простой статистической таблицы	197
4. Вычисление нормированного отклонения и рейтинга успеваемости	199
4.1. Вычисление нормированного отклонения	199
4.2. Вычисление рейтинга успеваемости	203
5. Вычисление вероятности стандартного нормального распределения	204
6. Вычисление значения χ при распределении хи-квадрат	205
7. Вычисление коэффициента линейной корреляции	207
8. Проверка гипотезы о независимости.	208
Предметный указатель	212

Предисловие

Данная книга — наглядное учебное пособие по статистике, которое, в первую очередь, предназначается тем, кому приходится заниматься анализом различных данных, а также тем, кто пока такой анализ не проводит, но хотел бы знать, что же такое статистика.

Автору также будет весьма приятно, если книга окажется интересной и для тех, кто уже изучал эту дисциплину.

Статистика — одна из областей математики, тесно связанная с жизнью и работой.

Если овладеть всеми премудростями этой науки, то можно, например:

- предусмотреть, сколько коробок жареной лапши будет продано в студенческом киоске, который планируется открыть на университетском празднике;
- оценить вероятность успешной сдачи квалификационного экзамена;
- сравнить вероятность выздоровления, если принимать лекарство X и не принимать это лекарство.

Книга содержит 7 глав. За некоторым исключением, главы книги построены следующим образом:

- манга (комикс);
- объяснение, дополняющее мангу;
- упражнения и ответы;
- выводы.

Книга написана так, что читатель может усвоить материал, прочитав только мангу. А вот чтобы получить более глубокие знания, придётся прочитать и всё остальное.

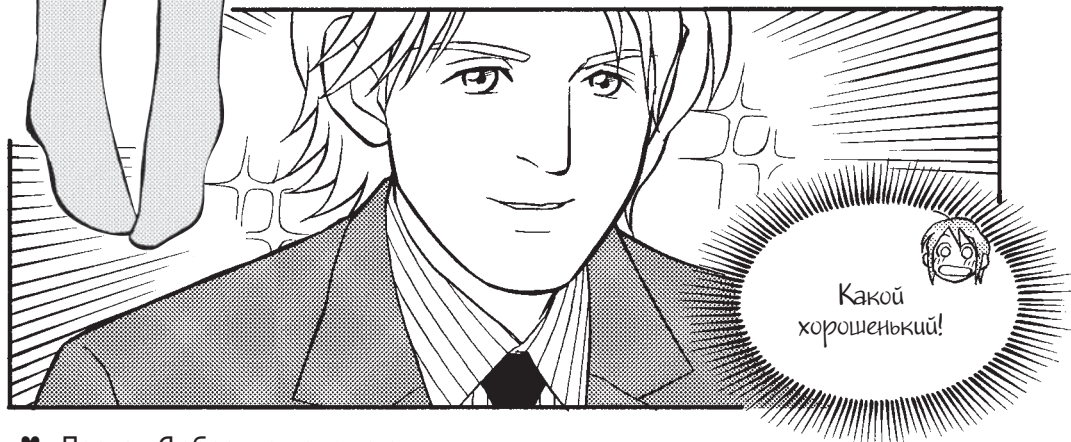
Предел мечтаний автора — читатель, который, перевернув последнюю страницу книги, скажет: «Статистика — это так интересно! Но это ещё и полезно! Да это просто здорово!».

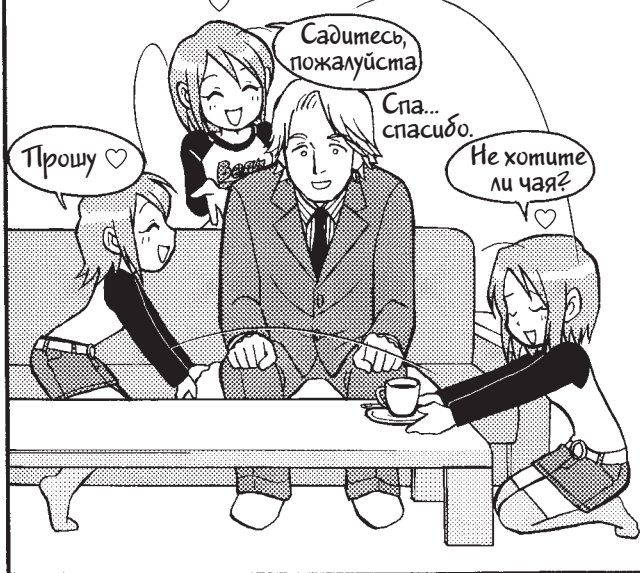
Я бесконечно благодарен всем сотрудникам редакции издательства Ohmsha за предоставленную мне возможность написать эту книгу, а также всем сотрудникам компании Trend-Pro. Я глубоко признателен г-ну Ре Акино, автору сценария, и г-ну Ироха Иноуэ, воплотившему этот сценарий в виде рисунков, за те титанические усилия, которые им пришлось приложить, чтобы на основе моей рукописи создать потрясающий комикс. Я также благодарен г-ну Фумитаке Сакаи (социологический факультет университета Риккё), советами которого я руководствовался во время работы над этой книгой.

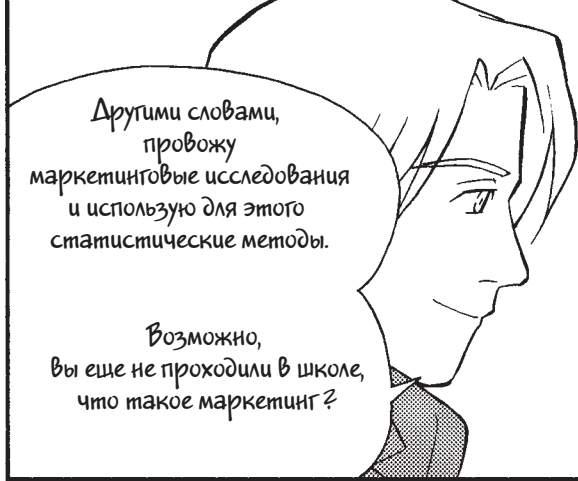
Син Такахаси
Июль, 2004 год

Пролог

**Любовь
и статистика**







Другими словами,
провожу
маркетинговые исследования
и использую для этого
статистические методы.

Возможно,
вы еще не проходили в школе,
что такое маркетинг?

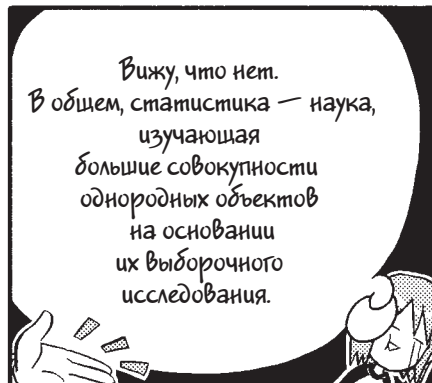


Не помню...
Кажется, нет.



Честная девочка.
Ну, а что такое статистика,
знаешь?

Э-э-э...



Вижу, что нет.
В общем, статистика — наука,
изучающая
большие совокупности
однородных объектов
на основании
их выборочного
исследования.



Что-то
я слишком
затнул.

НЕПОСТИЖНО...

Эй, Руи!
Что с тобой?

А, кстати, ...



... как раз
в сегодняшней газете
есть информация о рейтинге
кабинета министров.



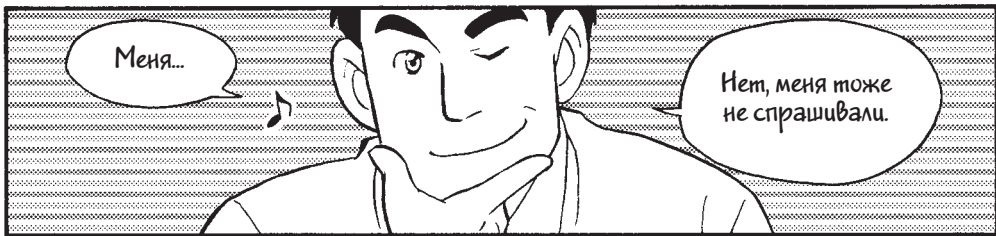
По исследованиям, проведённым газетой "Ведомости", рейтинг кабинета министров среди избирателей 39%.

И что это значит?



Но сотрудники газеты моего мнения не спрашивали.

А как насчет Вас, Такацу-сан?



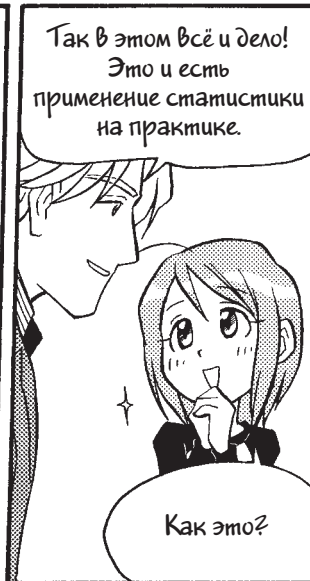
Меня...

Нет, меня тоже не спрашивали.



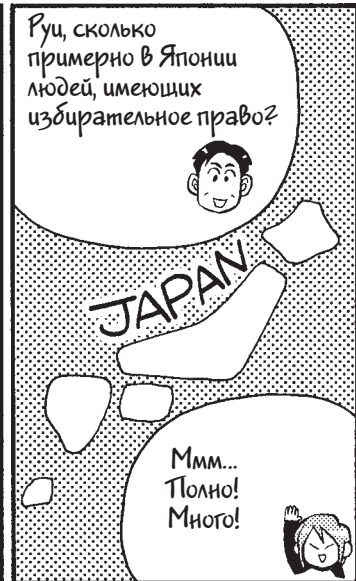
Хм-м... Мнением двоих не поинтересовались, а рейтинг опубликовали. Как такое может быть?

Между тем, избирательного права вас никто не лишал. Более чем странно...



Так в этом всё и дело! Это и есть применение статистики на практике.

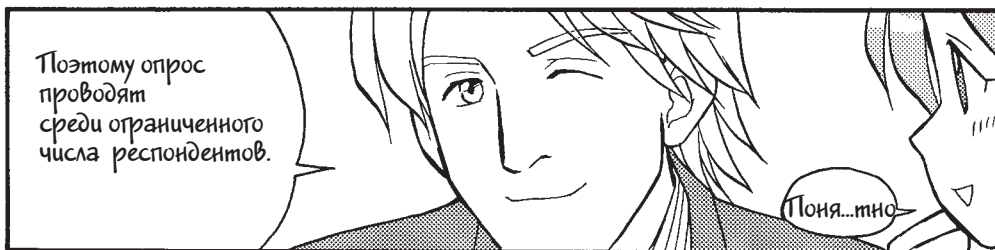
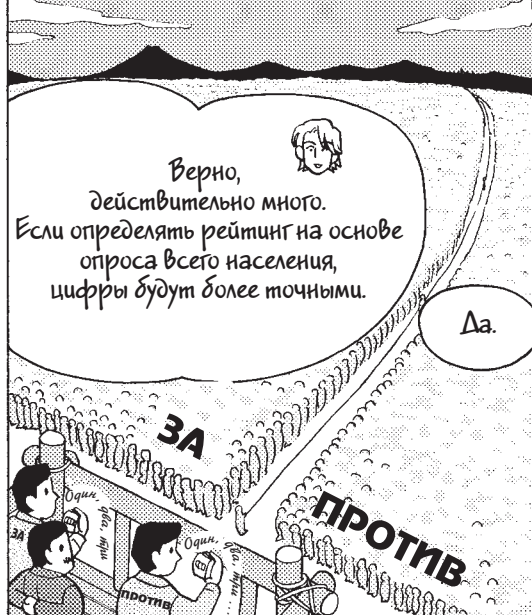
Как это?

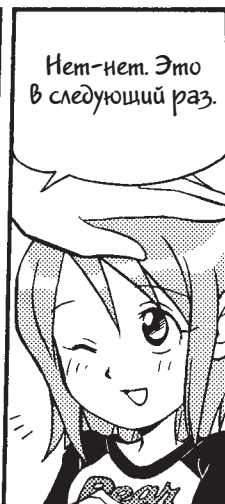


Ну, сколько примерно в Японии людей, имеющих избирательное право?

JAPAN

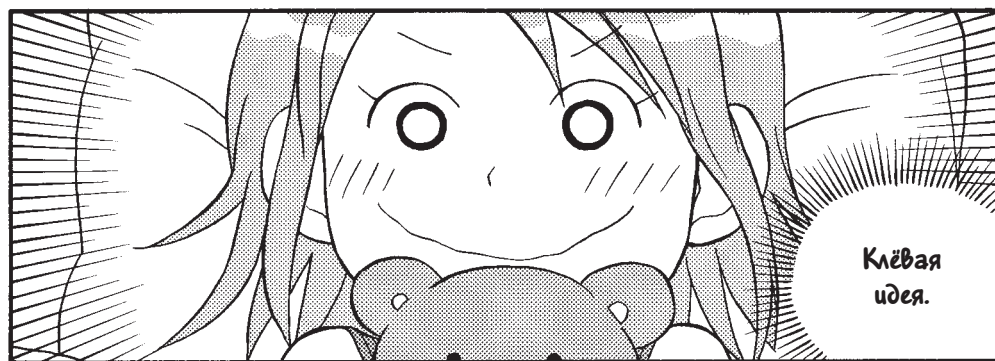
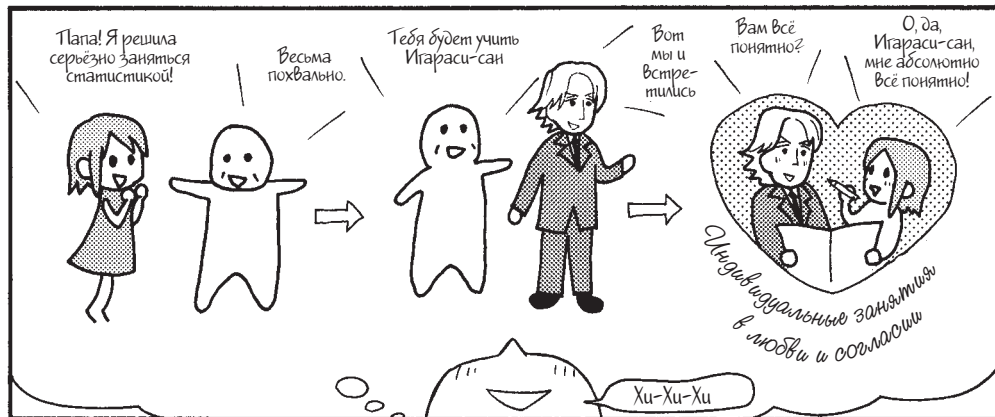
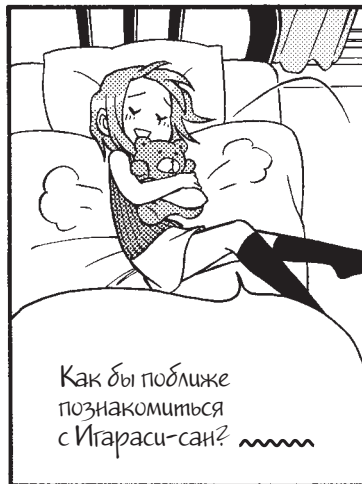
Ммм... Полно! Много!

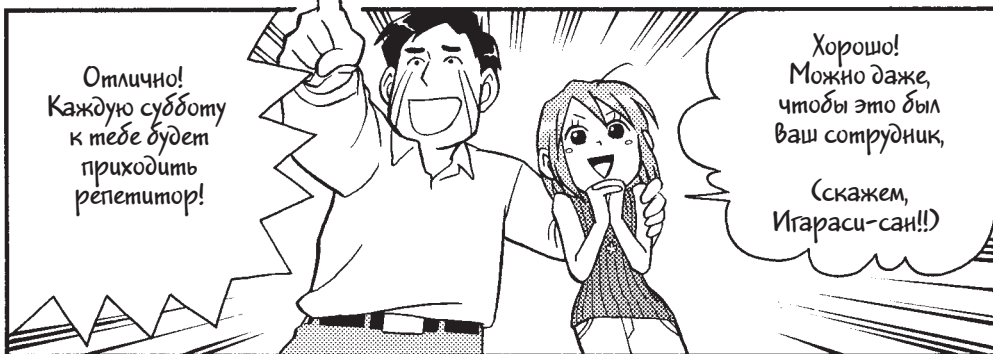
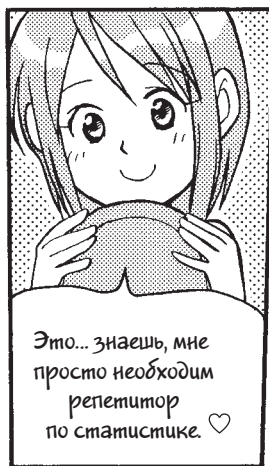


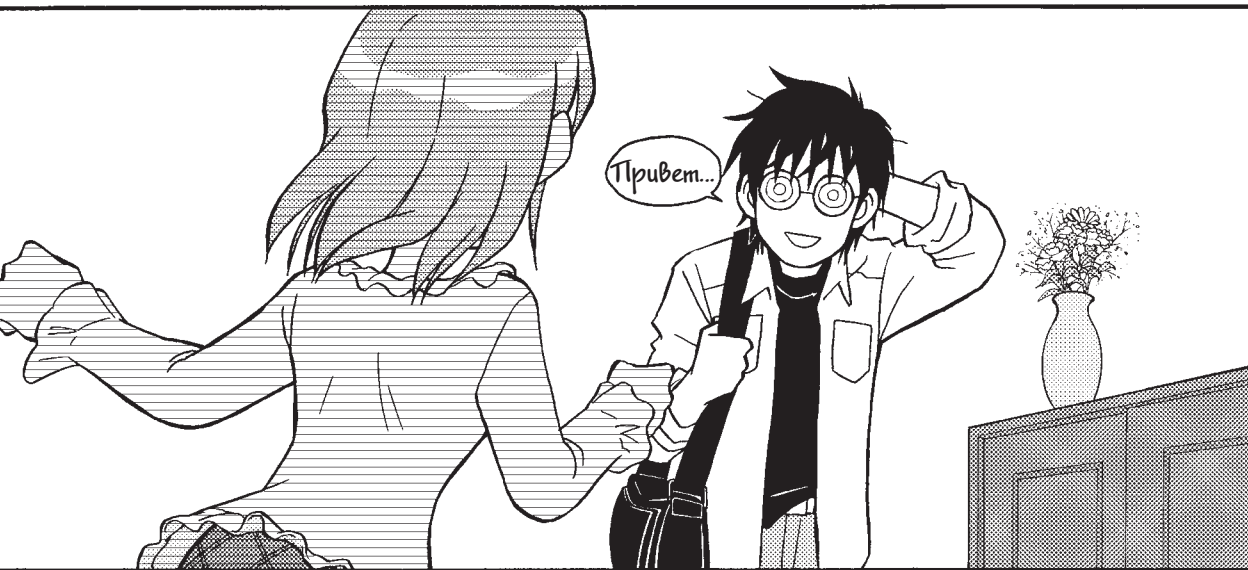
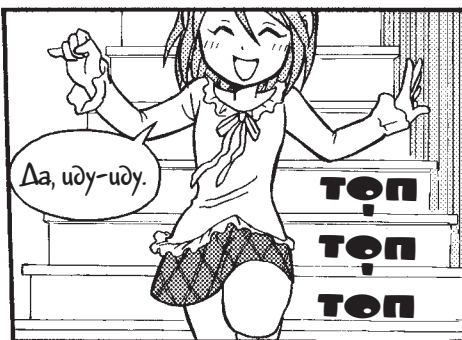


На следующий день

Мечтает







Кто этот парень?!!

Ну, это наш сотрудник, Мамору Ямамото.

Какой-то расстрепанный!

Очень приятно!

Папа, а Игараси-сан?

Причём тут Игараси? Мамору и живёт ближе, и учить будет не хуже.

О, нет!

Ну, я пошёл, счастливо позаниматься.

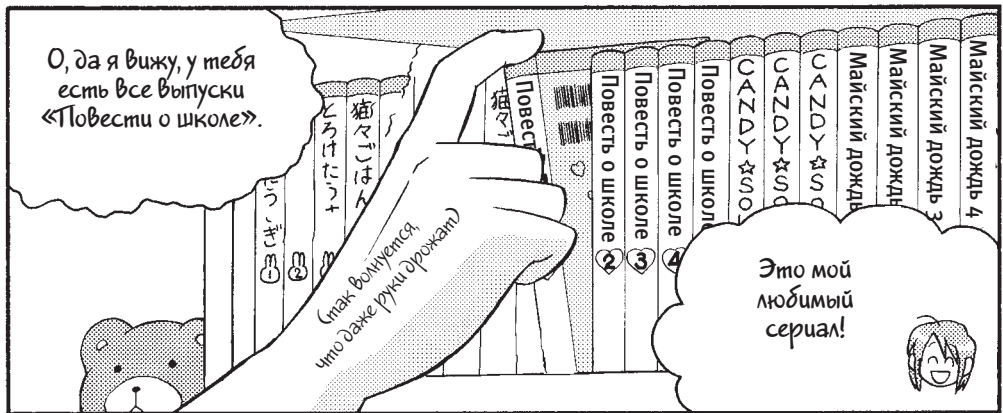
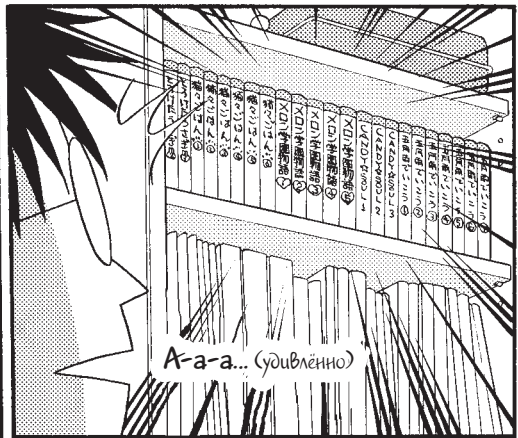
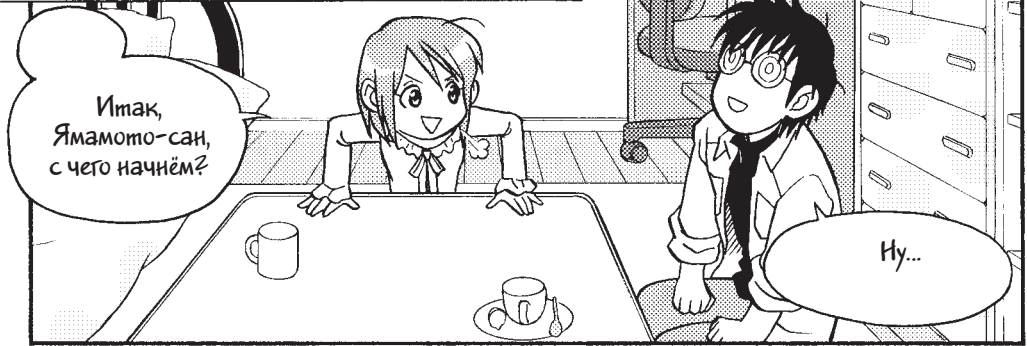
Хе-хе-хе!

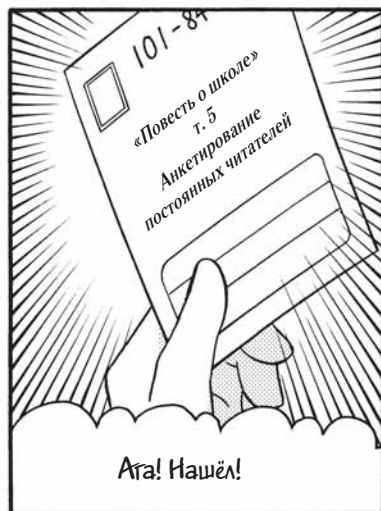


Глава 1

Разберёмся с типами данных

1. Количественные и качественные данные





★ Повесть о школе, т. 5 ★
Анкета постоянных читателей

Вопрос 1. Ваше мнение о 5-м томе «Повести о школе»?

1. Очень интересно
2. Довольно интересно
3. Так себе
4. Скучновато
5. Совершенно неинтересно

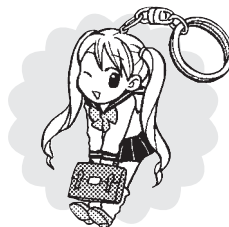
Вопрос 2. Ваш пол?

1. ж 2. м

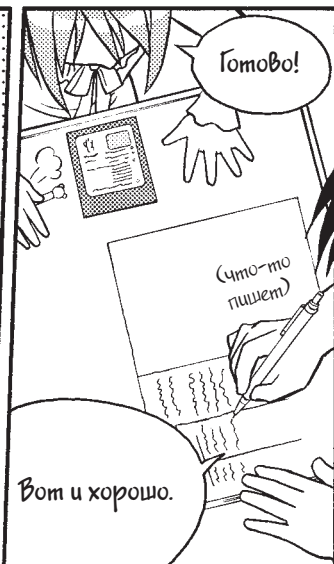
Вопрос 3. Ваш возраст? ___ лет

Вопрос 4. Сколько выпусков журнала вы приобретаете в месяц? ___ шт

Среди участников анкетирования будут разыграны 30 брелоков «Рина»!

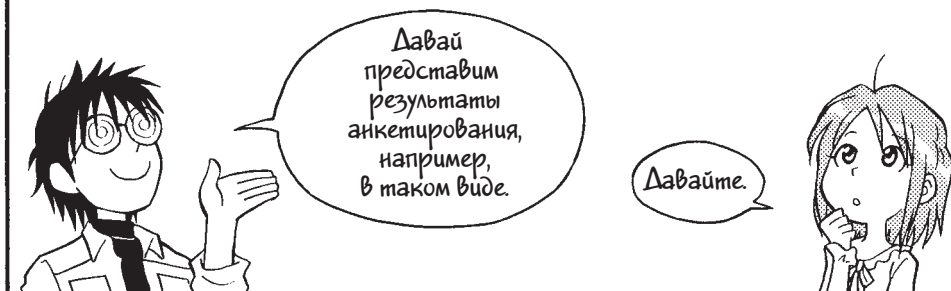


Спасибо за Ваши ответы. Ваше ценное для нас мнение мы учтём в последующих изданиях, а также при разработке будущих проектов.

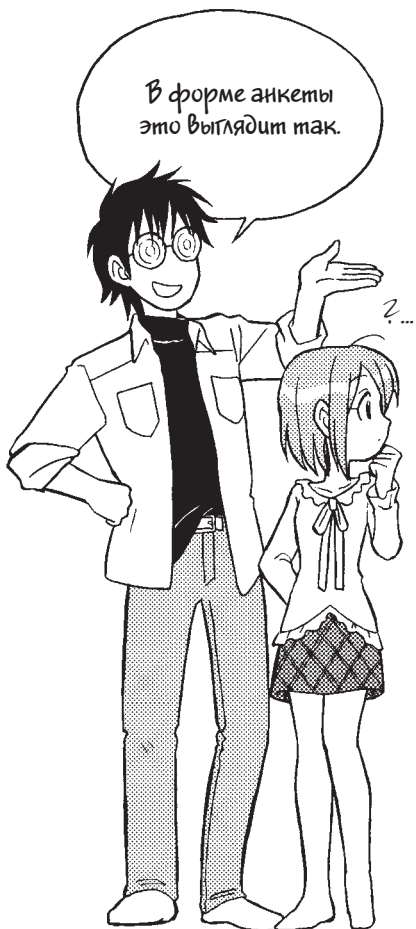


Анкета постоянных читателей

Респондент	Ваше мнение о «Новости о школе»	Пол	Возраст, лет	Кол-во при- обретенных в месяц выпусков, шт.
Рди	очень интересно	ж	17	2
А	довольно интересно	ж	17	1
Б	так себе	м	18	5
В	скудно	м	22	7
Г	довольно интересно	ж	25	4
Д	совершенно неинтересно	м	20	3
Е	очень интересно	ж	16	1
Ж	довольно интересно	ж	17	2
З	так себе	м	18	0
И	так себе	ж	21	3
⋮	⋮	⋮	⋮	⋮







«Повесть о школе», т. 5
Анкета постоянных читателей



Вопрос 1. Ваше мнение о 5-м томе
«Повести о школе»?

- ① Очень интересно
2. Довольно интересно
3. Так-так
4. Качественные данные
5. Совершенно неинтересно

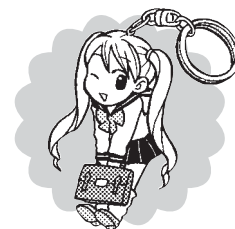
Вопрос 2. Ваш пол?

1. ж 2. м

Вопрос 3. Ваш возраст? 17 лет

Вопрос 4. Сколько журналов журнала
«Рина» приобретаете в месяц? 2 шт

Среди участников анкетирования будут разыграны 30 брелоков «Рина»!



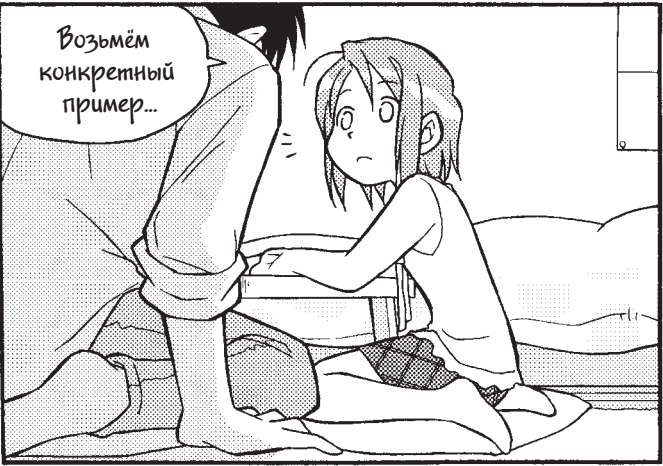
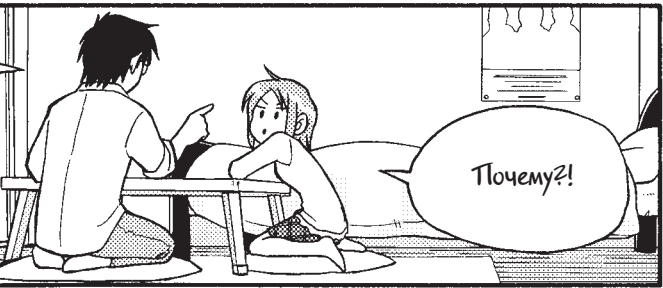
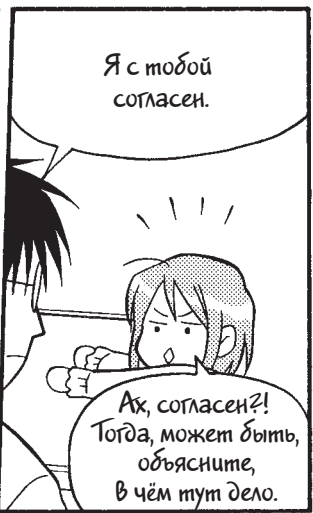
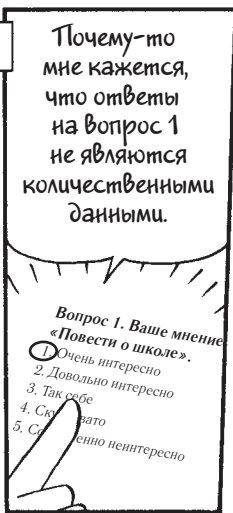
Спасибо за Ваши ответы. Ваше ценное для нас мнение мы учтём в последующих изданиях, а также при разработке будущих проектов.

Качественные данные — это данные, которые нельзя измерить.

Количественные данные — данные, которые можно измерить.

А-а, понятно.

2. Примеры качественных данных



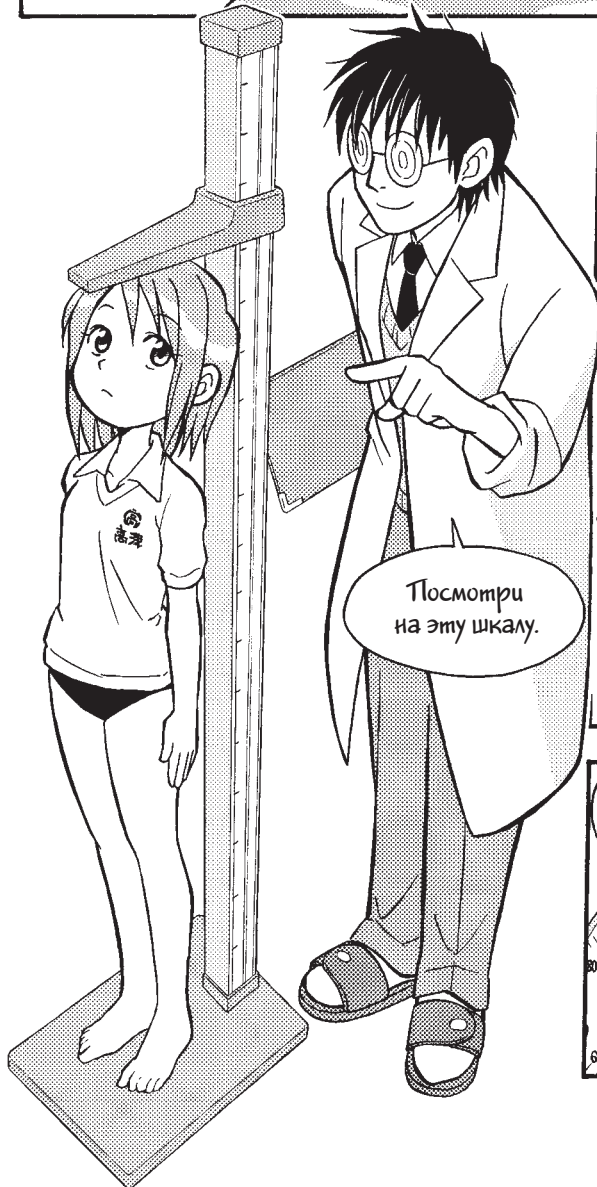


Да как Вы смеете задавать девушке такие вопросы? Это просто возмутительно!



БУМ!

Так, 151 см.



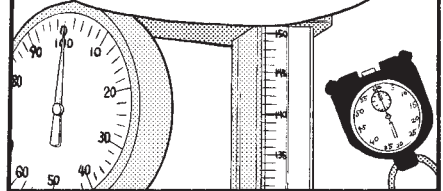
У этой шкалы каждое деление равно 1 см.

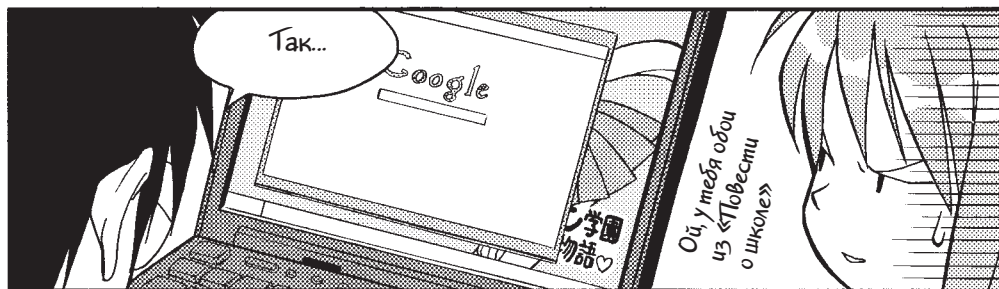
Верно. Поэтому деление, следующее за 151 см, будет 152 см, а следующее — 153 см и так далее.

Посмотри на эту шкалу.

Да...

Это значит, что шкала имеет равные, или одинаковые, интервалы между соседними делениями.





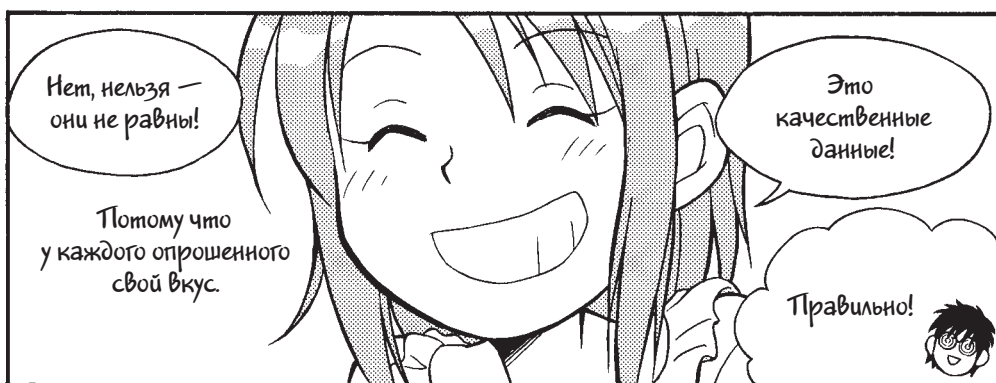
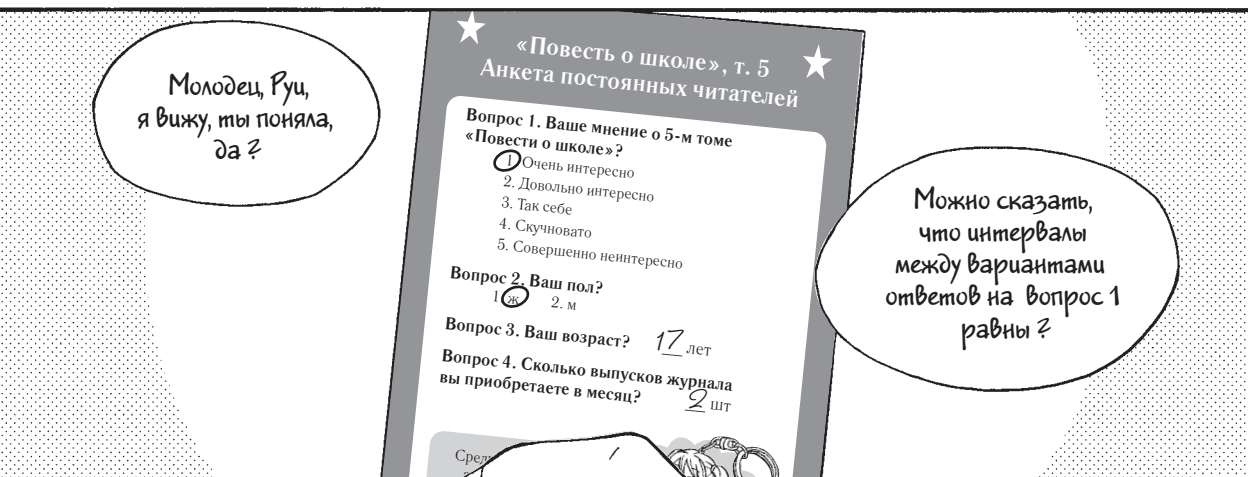


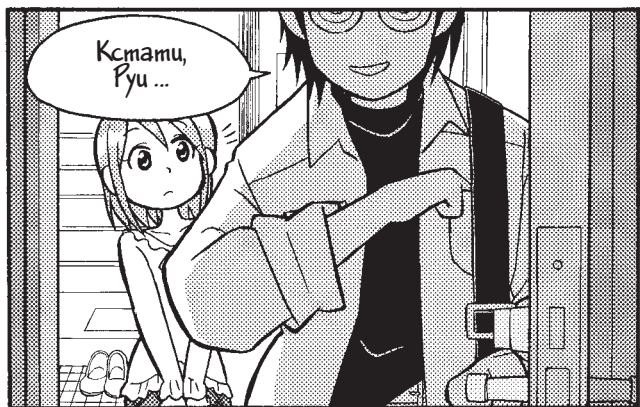
Критерии сложности экзамена Eiken

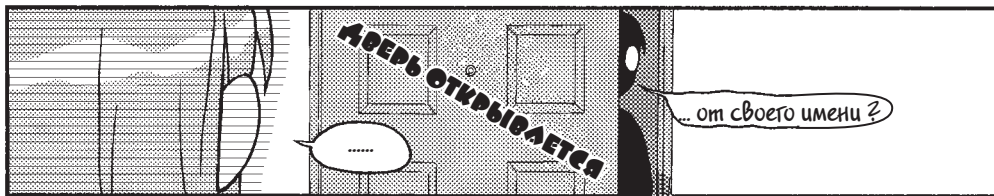
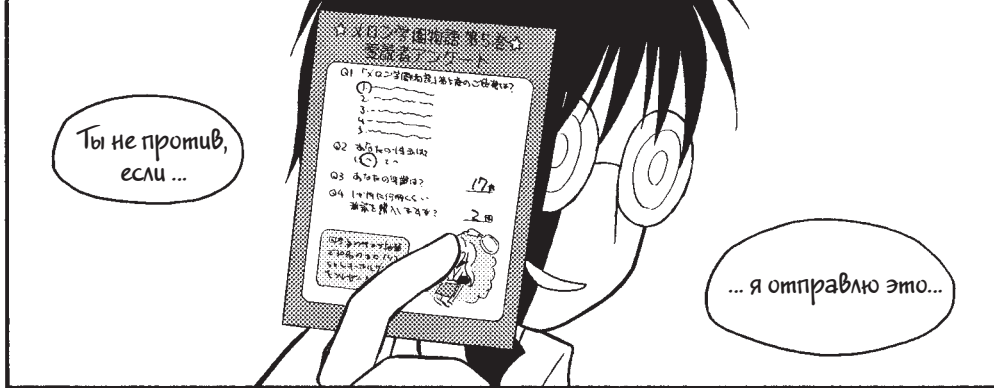
<http://www.eiken.or.jp/>

1-й уровень	2-й уровень	3-й уровень	4-й уровень	5-й уровень
Продвинутый, соответствует программе вуза; словарный запас примерно 10 – 15 тыс. слов	Соответствует программе средней школы; словарный запас примерно 5,1 тыс. слов	Соответствует программе 3-го года обучения средней школы; словарный запас примерно 2,1 тыс. слов	Средний, соответствует программе 2-го года обучения средней школы; словарный запас примерно 1,3 тыс. слов	Начальный, соответствует программе 1-го года обучения средней школы; словарный запас примерно 600 слов









3. Использование многовариантных ответов на практике

Как было показано, вопрос 1 анкеты постоянных читателей относится к качественным данным. Однако на практике, например при опросе потребителей, эти данные часто рассматриваются как количественные. Другими словам, возможны случаи, когда это может выглядеть так:

		балл
очень интересно	⇒	5
довольно интересно	⇒	4
так себе	⇒	3
скучновато	⇒	2
совершенно неинтересно	⇒	1

или так:

		балл
очень интересно	⇒	2
довольно интересно	⇒	1
так себе	⇒	0
скучновато	⇒	-1
совершенно неинтересно	⇒	-2

Существуют мир теории и мир практики, точнее, теоретический мир и реальный. Поэтому одни и те же данные могут рассматриваться и как количественные, и как качественные: всё зависит от того, где они используются — в теории или на практике.

Упражнение

Посмотрите на таблицу:

Респондент	Группа крови	Оценка вкусовых качеств спортивного коктейля X	Комфортная комнатная температура при работающем кондиционере, °C	Лучший результат бега на 100 м, с
А	В (III)	невкусно	25	14,1
Б	А (II)	вкусно	24	12,2
В	AB (IV)	вкусно	25	17,0
Г	О (I)	так себе	27	15,6
Д	А (II)	невкусно	24	18,4
...

Определите, к каким категориям данных относятся графы:

«Группа крови», «Оценка вкусовых качеств спортивного коктейля X»,
«Комфортная комнатная температура при работающем кондиционере»
и «Лучший результат бега на 100 м».

Ответ

«Группа крови» и «Оценка вкусовых качеств спортивного коктейля X» относятся к качественным данным.
«Комфортная комнатная температура при работающем кондиционере» и «Лучший результат бега на 100 м» относятся к количественным данным.

Выводы

Данные делятся на количественные и качественные. Такие категории данных, как, например «Очень интересно», ..., «Совершенно неинтересно» с точки зрения теории относятся к качественным данным. Однако на практике эти же данные могут рассматриваться и как количественные.

Глава 2

Знакомимся с количественными данными

1. Ряды распределения и гистограммы

ВАЖНО

* Вкусный рамэн**.
50 лучших ресторанов

ДВЕРЬ ОТКРЫВАЕТСЯ

Привет,
Руи.

** Рамэн — китайская лапша, блюдо ресторанов быстрого питания.

Очень!
Посмотрела этот журнал и думаю,
в какой ресторан
лучше пойти.

А-а-а...
(удивляется)

Всё так
аппетитно,
правда?

Здравствуйте,
Ямамото-
сан.

Да, никак
ты любишь
рамэн?




Для начала
я свёл в таблицу
цены на рамэн.

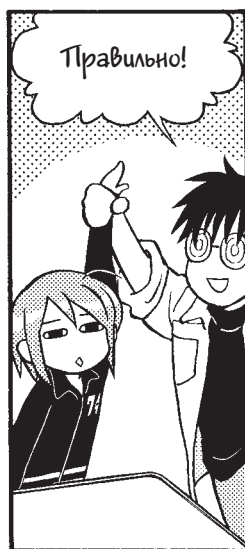
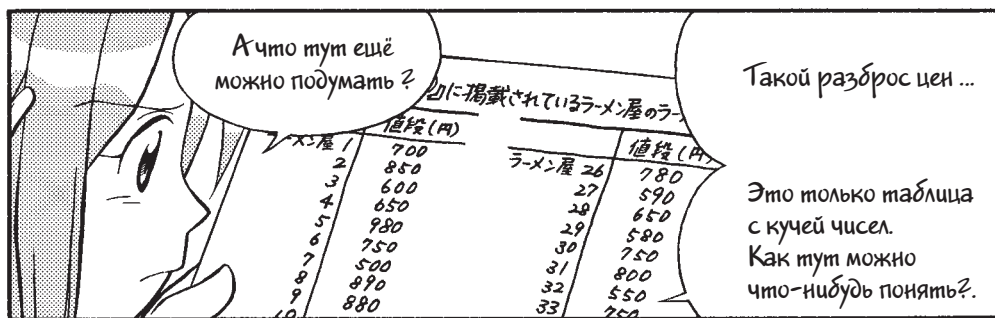


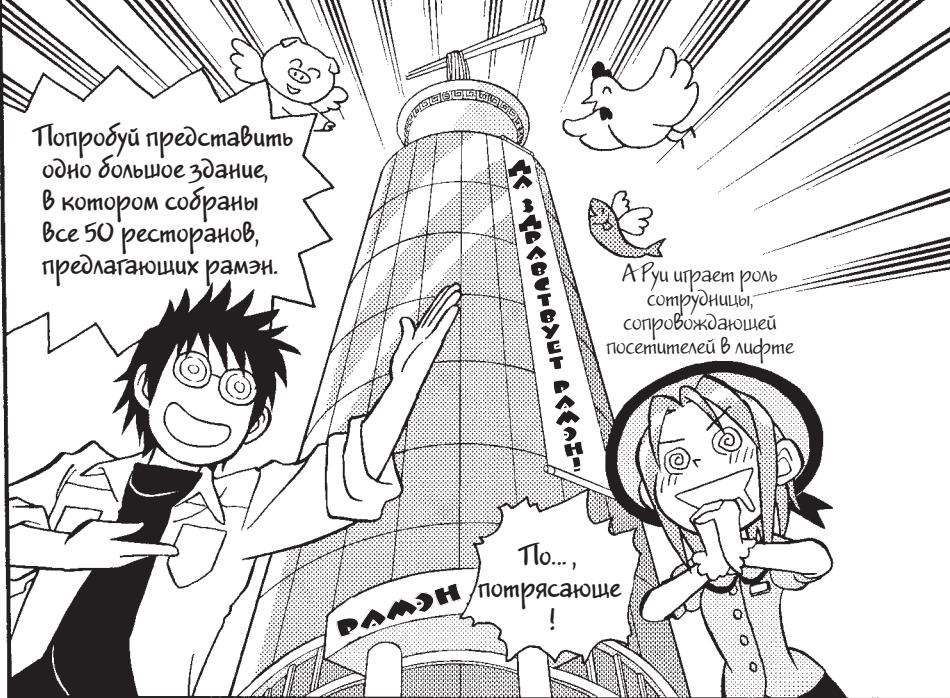
Цены на рамэн в 50 лучших ресторанах
(от журнала «Вкусный рамэн. 50 лучших ресторанов»)

Ресторан	Цена, иены	Ресторан	Цена, иены
1	700	26	780
2	850	27	590
3	600	28	650
4	650	29	580
5	980	30	750
6	750	31	800
7	500	32	550
8	890	33	750
9	880	34	700
10	700	35	600
11	890	36	800
12	720	37	800
13	680	38	880
14	650	39	790
15	790	40	790
16	670	41	780
17	680	42	600
18	900	43	670
19	880	44	680
20	720	45	650
21	850	46	890
22	700	47	930
23	780	48	650
24	850	49	777
25	750	50	700

От обсуждения
ресторанов
главно перешли
к занятию...
Всё-таки
странный тип...









Попробуй представить одно большое здание, в котором собраны все 50 ресторанов, предлагающих рамэн.

А Руи играет роль сотрудницы, сопровождающей посетителей в лифте

Этаж и интервал цен			<p>В каждом ресторане только один вид рамэна ...</p>																	
От	До																			
5 этаж	900—1000	<table border="1"> <tr><td>5</td><td>18</td><td>47</td></tr> </table>	5	18	47	<p>На каждом этаже свой интервал цен на рамэн (от ... до ...)</p> <p>Такое разделение в статистике называется распределением.</p>														
5	18	47																		
4 этаж	800—900	<table border="1"> <tr><td>37</td><td>38</td><td>46</td></tr> <tr><td>2</td><td>8</td><td>9</td><td>11</td><td>19</td><td>21</td><td>24</td><td>31</td><td>36</td></tr> </table>	37	38	46		2	8	9	11	19	21	24	31	36					
37	38	46																		
2	8	9	11	19	21		24	31	36											
3 этаж	700—800	<table border="1"> <tr><td>26</td><td>30</td><td>33</td><td>34</td><td>39</td><td>40</td><td>41</td><td>49</td><td>50</td></tr> <tr><td>1</td><td>6</td><td>10</td><td>12</td><td>15</td><td>20</td><td>22</td><td>23</td><td>25</td></tr> </table>	26	30	33	34	39	40	41	49	50	1	6	10	12	15	20	22	23	25
26	30	33	34	39	40	41	49	50												
1	6	10	12	15	20	22	23	25												
2 этаж	600—700	<table border="1"> <tr><td>43</td><td>44</td><td>45</td><td>48</td></tr> <tr><td>3</td><td>4</td><td>13</td><td>14</td><td>16</td><td>17</td><td>28</td><td>35</td><td>42</td></tr> </table>	43	44	45	48	3	4	13	14	16	17	28	35	42					
43	44	45	48																	
3	4	13	14	16	17	28	35	42												
1 этаж	500—600	<table border="1"> <tr><td>7</td><td>27</td><td>29</td><td>32</td></tr> </table>	7	27	29	32														
7	27	29	32																	



Понятно ...

На каждом этаже есть вывеска, на которой указана средняя цена на рамэн на этом этаже.



На втором этаже цены варьируются от 600 до 700 иен. Значит, средняя цена равна 650 иен!

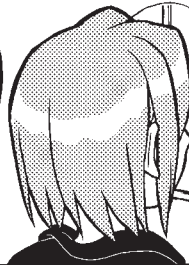
Путеводитель по этажам

Номер этажа и интервал цен	Рестораны, расположенные на каждом этаже	Цифра на вывеске (средина интервала)
5 этаж 900...1000	■ ■ ■ ■	950
4 этаж 800...900	■ ■ ■ ■ ■ ■ ■ ■	850
3 этаж 700...800	■ ■ ■ ■ ■ ■ ■ ■ ■ ■	750
2 этаж 600...700	■ ■ ■ ■ ■ ■ ■ ■ ■ ■ ■ ■	650
1 этаж 500...600	■ ■ ■ ■ ■ ■ ■ ■	550



Это называется серединой интервала

Рестораны распределены по этажам в соответствии с ценами на рамэн. На каждом этаже может быть разное количество ресторанов.



800-900	=====	850
700-800	=====	750
600-700	=====	650
500-600	=====	

(смеется)

Сямамото-сан в статусе сотрудницы, сопровождающей посетителей по этажам)

Действительно.



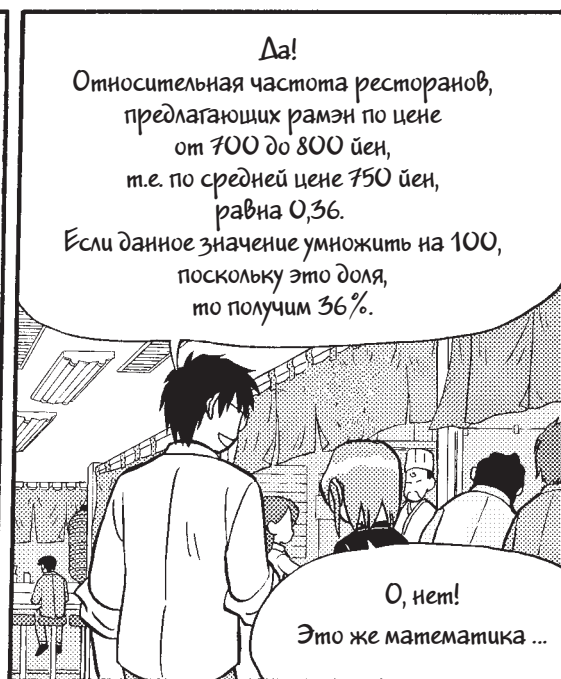
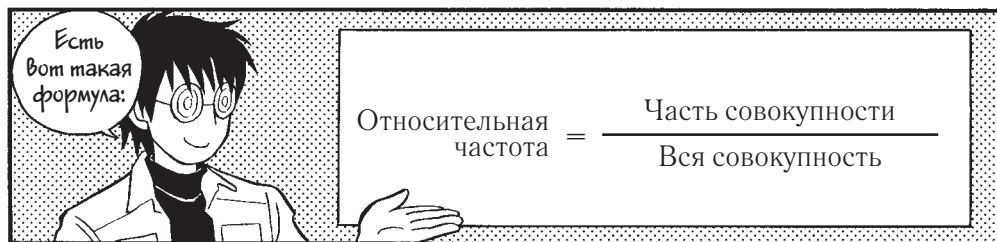
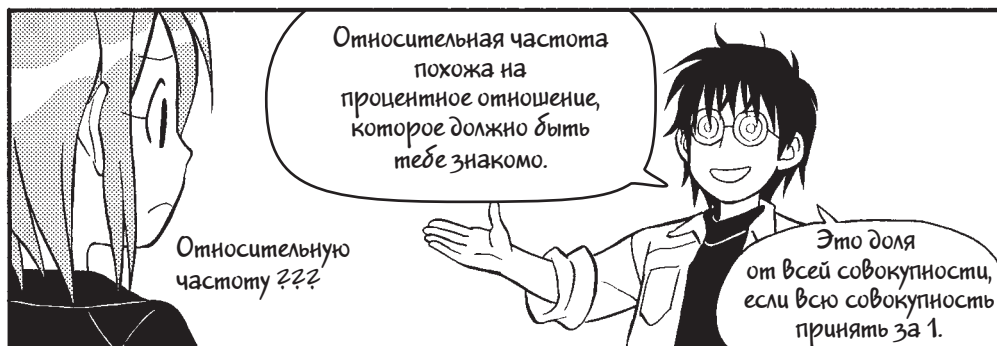
Число ресторанов, расположенных на каждом этаже, называют частотой.

На 1-м этаже — 4,
на 2-ом этаже — 13,
и т.д.



На 3-м этаже больше всего ресторанов — 18!

Теперь вычислим относительную частоту ресторанов на 3-м этаже!





Распределение по цене
50 лучших ресторанов,
предлагающих вкусовой ритэн

Интервал цен	Середина интервала	Кол-во ресторанов, частота	Количество ресторанов, относительная частота
500 ~ 600	550	4	0.08
600 ~ 700	650	13	0.26
700 ~ 800	750	18	0.36
800 ~ 900	850	12	0.24
900 ~ 1000	950	3	0.06
Итого:		50	1.00



... по горизонтали,
(на оси x),
откладывают
переменные.

В нашем случае
это будут
цены на рамэн.

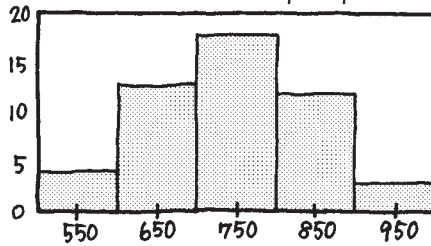
Ширина столбца
равна величине
интервала.

В середине столбца
ставят значение,
равное середине
интервала.

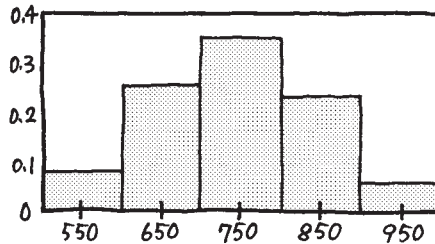


Гистограммы построены на основе
таблицы распределения 50 лучших
ресторанов, предлагающих рамэн

Гистограмма 1.
Частота (количество ресторанов)



Гистограмма 2.
Относительная частота



По
вертикали
(оси y)
отложены:

на верхнем
рисунке —
частота,

на нижнем рисунке —
относительная
частота



Ну, как?

М-мм

Будем
считать,
что с ценами
на рамэн ...

... я
худо-бедно
разобралась.

Вот это
«худо-бедно»
очень важно!
Таблицы (или ряды)
распределения
и гистограммы
помогают лучше
понять данные!

Вот как?..
Понятно!

2. Средняя величина

Мы недавно с девочками из моего класса ходили в боулинг...

Во время перерыва на чай

Удалось сбить хоть одну кеглю?



Что?!

Да я...

Я очень хорошо играю в боулинг!

Я пошутил.

Если все девочки класса, это довольно много, не так ли?

Да, 18 девчонок. Поэтому мы разделились на три команды по 6 человек.

Смотри, вот таблица результатов игры.

(быстро достает таблицу)



Результаты игры в боулинг

Команда А		Команда Б		Команда В	
Игрок	Очки	Игрок	Очки	Игрок	Очки
Руи-Руи	86	Томми	84	Синобу	229
Дэюн	78	Хаси	71	Юки	77
Юми	124	Хана	108	Хитоми	59
Связка	111	Мэй	85	Рисако	95
Токо	90	Канна	90	Май	70
Кэвэ	38	Асати	89	Кэвэ	88



O! Это отличный материал для исследования.

А Руи-Руи — это ты?



Да!
Я набрала 86 очков!

Беглый просмотр результатов позволяет сделать вывод, что у тебя, Руи, был средний результат в команде, да?

И что с того?!



Средний означает результат, который в среднем набрал один человек в каждой команде. Понятно?

Понятно. Это результат, который находится посередине других результатов, набранных игроками команды, так?



Если мой результат окажется выше среднего, вы угостите меня пирожным! ♡

Давай попробуем вычислить среднюю величину.

Может быть ...

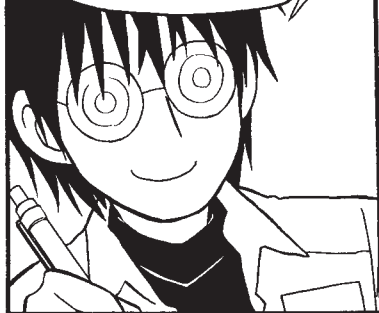


Вы сражались командами. Это значит, что вы сражались за итоговый результат, набранный каждой командой. Верно?



Да, и что?

Средний результат равен общему количеству очков, набранному командой, делённому на число игроков в команде.



Команда А

$$\frac{86+73+124+111+90+38}{6} = \frac{522}{6} = 87$$

Команда Б

$$\frac{84+71+103+85+90+89}{6} = \frac{522}{6} = 87$$

Команда В

$$\frac{229+77+59+95+70+88}{6} = \frac{618}{6} = 103$$

Команда В молодец!!



Значит, среднее количество очков в твоей команде, Руи-Руи, равно 87.

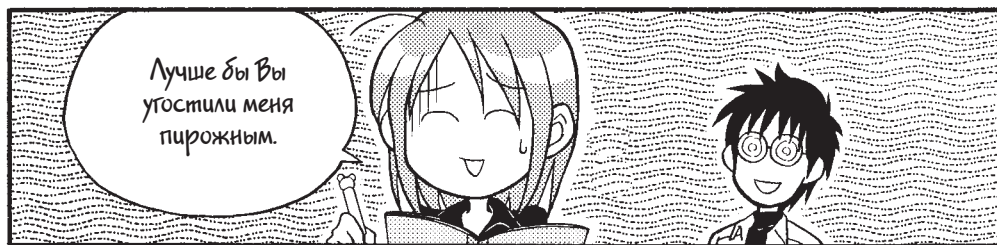


А у тебя было 86, так?

Так Вы угостите меня пирожным?

(злитесь)





3. Медиана

Посмотри
ещё раз
на таблицу
результатов

Ну что
на этот раз?

Давай
посмотрим
на команду В.

Как ты
думаешь,

Результаты игры в боулинг

Команда А		Команда Б		Команда В	
Игрок	Очки	Игрок	Очки	Игрок	Очки
Фуи-Фуи	86	Томи	84	Синобу	229
Дэюн	73	Хаюи	71	Юки	77
Юми	124	Хана	108	Хатоми	59
Свэтика	111	Мэй	85	Рисако	95
Токо	90	Канна	90	Май	70
Каэдэ	38	Асати	89	Козуэ	88

правильно ли
считать
средним
результатом
команды В
103 очка?

Вообще-то не очень.
5 игроков команды
набрали < 100 очков,
а средний результат
> 100 !?

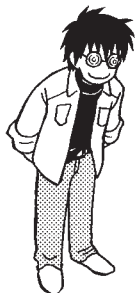
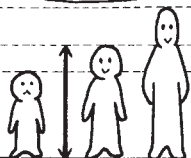
В подобных случаях,
когда имеются
слишком большие
или малые значения ...

... вычисляют
не среднюю величину,
а медиану .

Синобу
просто
моло-
де-е-еи...

Медиану?

Медиана — значение, которое приходится на середину ряда, если расположить данные в порядке возрастания (или убывания).



Для начала попробуем расположить в ряд очки, набранные игроками каждой команды.



Команда А

38 73 86 90 111 124

Команда Б

71 84 85 89 90 103

Команда В

59 70 77 88 95 229

Ряд с нечётным числом элементов

-1041,6 -39,0 **-5,7** 60,4 77,3

медиана

Ряд с чётным числом элементов

-0,4 35,2 **37,8 42,2** 46,1 910,3

Среднее значение этих двух элементов будет медианой

Если ряд состоит из нечётного числа элементов, медианой будет значение, находящееся точно посередине, ...

... а если ряд с чётным числом элементов, как в случае с боулингом, медианой будет среднее значение между 3-м и 4-м элементами.

А теперь попробуем вычислить медиану для команды В.



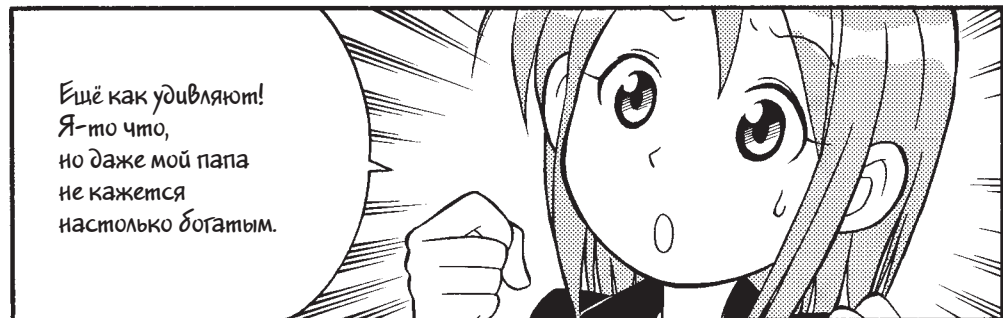
$$\frac{77 + 88}{2} = 82,5$$

Вот что получается.

Правильно!



* Меньше 3500 рублей.





Средние накопления
такие большие
из-за миллионеров.

Не нужно расстраиваться,
если сумма ваших сбережений
намного меньше
средней величины.



Миллионеры...

В подобных случаях,
медиана
гораздо ближе
к размеру сбережений
обычных людей.

Да ты меня,
похоже,
не слушаешь
...



Решено!
Выйду замуж
за богача,
чьи сбережения
намного больше
медианы!

Ты меня
расстраиваешь

4. Стандартное отклонение

Итак, давай рассмотрим результаты ...

... команд А и Б.

Давай.

Нарисуем шкалу...



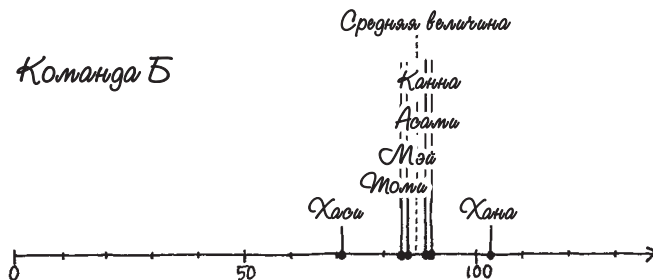
Теперь для каждого игрока отметим значение набранных им очков и напомним его имя.



Команда А

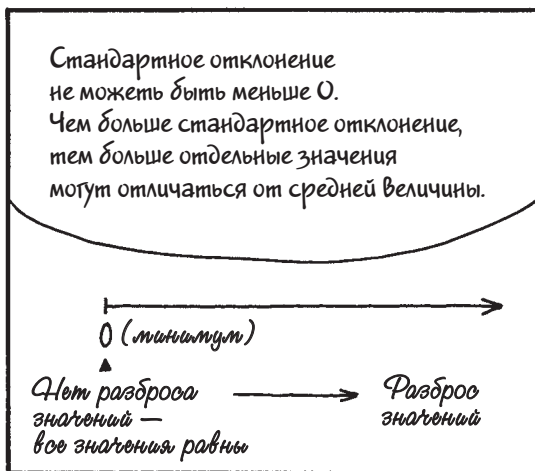
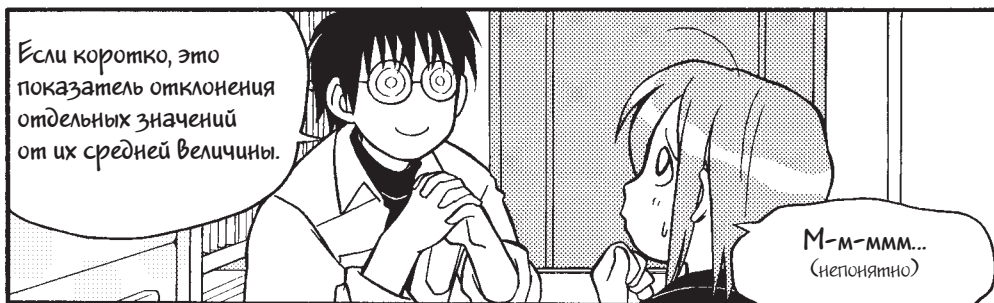


Команда Б



Средняя величина и для команды А, и для команды Б была равна 87,

но ситуация на рисунке (линии на шкале) сильно различается, верно?



Правильно!
А точная формула
имеет вид



И опять математика...



Стандартное
отклонение =

$$\sqrt{\frac{(i\text{-е значение} - \text{среднее значение})^2}{\text{кол-во значений}}}$$

Да не переживай
ты так!
Всего-то и надо —
подставить
в эту формулу
конкретные числа.
Хочешь, попробуем
вместе?

Хочу.



Сначала команда А

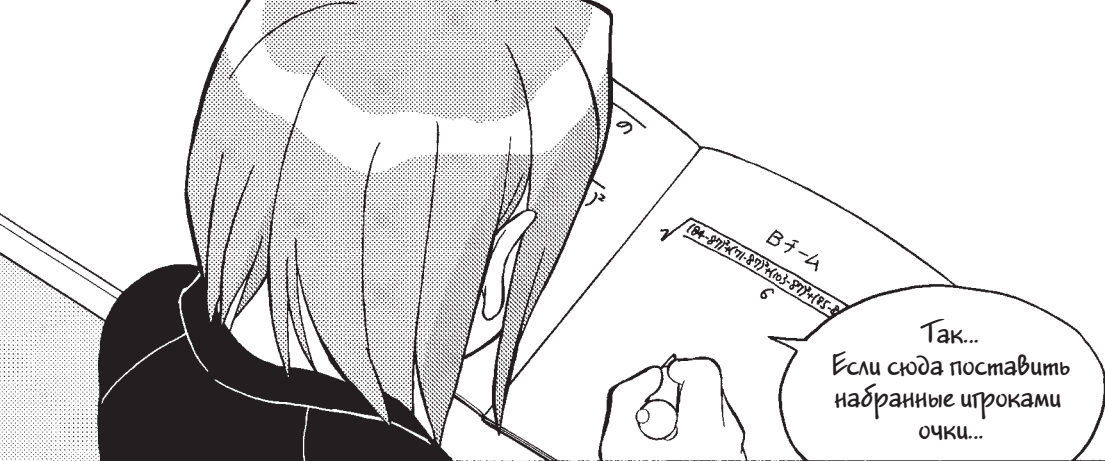
Команда А

$$\begin{aligned} & \sqrt{\frac{(86-87)^2 + (73-87)^2 + (124-87)^2 + (111-87)^2 + (90-87)^2 + (38-87)^2}{6}} = \\ & = \sqrt{\frac{(-1)^2 + (-14)^2 + 37^2 + 24^2 + 3^2 + (-49)^2}{6}} = \\ & = \sqrt{\frac{1 + 196 + 1369 + 576 + 9 + 2401}{6}} = \\ & = \sqrt{\frac{4552}{6}} = \\ & = \sqrt{758.6\dots} = \\ & \approx 27.5 \end{aligned}$$

Так-то
наверно
и я смогу.

Тогда попробуй сама
посчитать стандартное
отклонение
для команды Б.





Готово! Команда Б

$$\sqrt{\frac{(84-87)^2 + (71-87)^2 + (103-87)^2 + (85-87)^2 + (90-87)^2 + (89-87)^2}{6}}$$

$$= \sqrt{\frac{(-3)^2 + (-16)^2 + 16^2 + (-2)^2 + 3^2 + 2^2}{6}}$$

$$= \sqrt{\frac{9 + 256 + 256 + 4 + 9 + 4}{6}}$$

$$= \sqrt{\frac{538}{6}} =$$

$$= \sqrt{89.6\dots} =$$


$$\approx 9.5$$

Правильно!
Видишь,
справилась
же!



Стандартное отклонение:
Команда А — 27,5
Команда Б — 9,5

На самом деле у всех игроков
команды Б похожие результаты.
Стандартное отклонение
здесь меньше, чем в команде А



Я сказал, что формула для стандартного отклонения имеет вид

$$\sqrt{\frac{(\text{i-е значение} - \text{среднее значение})^2}{\text{кол-во значений}}}$$

но есть и другая формула

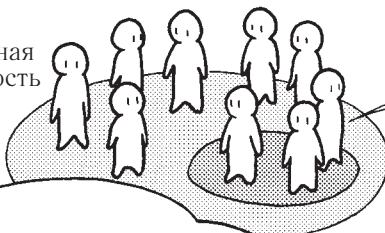
$$\sqrt{\frac{(\text{i-е значение} - \text{среднее значение})^2}{\text{кол-во значений} - 1}}$$



От общего количества значений отнимают 1?



Генеральная совокупность



Выборочная совокупность

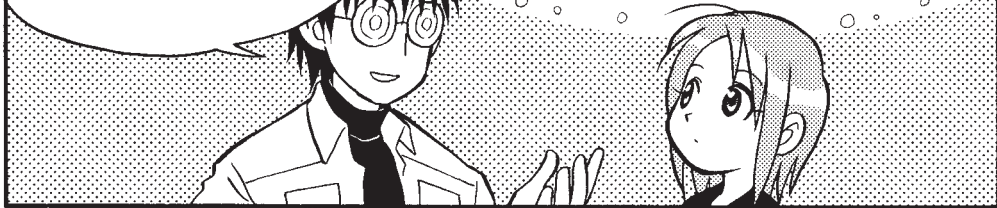
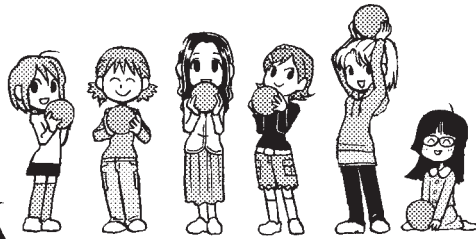
Первая формула используется при вычислении стандартного отклонения генеральной совокупности, ...

... а вторая формула — при вычислении стандартного отклонения в выборочной совокупности.

Генеральная совокупность — вся изучаемая группа людей или объектов, ...

... а выборочная совокупность — это группа людей или объектов, отобранная из генеральной совокупности, так?

Да, верно. Хорошо,
когда есть возможность
получить данные обо всех
объектах совокупности,
как в случае с твоей
командой. Но ...



... обычно сделать
это сложно.

Поэтому почти
всегда используют
вторую формулу.

Вот как...

Ну, на сегодня всё.



Спасибо!



5. Ряды распределения и величина интервала

Возможно, не все читатели до конца разобрались с понятиями «ряды распределения» и «гистограммы». Поэтому автор предлагает еще раз рассмотреть таблицу на **стр. 38**.

Таблица 2.1. Распределение 50 лучших ресторанов, предлагающих вкусный рамэн (по цене на рамэн)

Интервалы цен, йены	Середина интервала	Количество ресторанов, частота	Кол-во ресторанов, относительная частота
500—600	550	4	0,08
600—700	650	13	0,26
700—800	750	18	0,36
800—900	850	12	0,24
900—1000	950	3	0,06
Всего:		50	1,00

Как следует из таблицы, величина интервала равна 100. Это значение не является стандартом в математике. Просто так захотел Ямамото-сан. Решение о выборе интервала принимает тот, кто анализирует данные.

Не исключено, что среди читателей найдутся и такие, которым не дает покоя вопрос: «Ряды распределения, построенные на основе субъективных решений, неубедительны. Я не смогу показать их другим. Нет ли математического способа определения величины интервала?». Конечно, есть. Покажем, как можно вычислить величину интервала для **Табл. 2.1**.

Шаг 1

Количество интервалов определяется по формуле Стерджесса:

$$\text{Кол-во интервалов} = 1 + \frac{\log_{10} N}{\log_{10} 2},$$

где N — количество значений в совокупности.

$$1 + \frac{\log_{10} 50}{\log_{10} 2} = 1 + 5,6438... = 6,6438... \approx 7$$

Шаг 2

Величина интервала определяется по формуле:

$$\frac{\text{MAX} - \text{MIN}}{\text{Кол-во интервалов}},$$

где MAX — максимальное значение в совокупности,

MIN — минимальное значение в совокупности.

$$\frac{980 - 500}{7} = \frac{480}{7} = 68,5714... \approx 69$$

Ниже приведена таблица распределения 50 лучших ресторанов по цене на рамэн с величиной интервала, рассчитанной по формуле, данной на **Шаге 2**.

Таблица 2.2. Распределение 50 лучших ресторанов, предлагающих вкусный рамэн (по цене на рамэн)

Интервал цен, йены	Середина интервала	Количество ресторанов, частота	Количество ресторанов, относительная частота
500—569	534,5	2	0,04
569—638	603,5	5	0,10
638—707	672,5	15	0,30
707—776	741,5	6	0,12
776—845	810,5	10	0,20
845—914	879,5	10	0,20
914—983	948,5	2	0,04
Итого		50	1,00

Ну и как? Не исключено, что некоторым данная таблица может показаться менее убедительной, чем **Табл. 2.1**.

При этом могут возникнуть такие вопросы, как: «Почему величина интервала равна именно 69 йенам?», «И что это за формула, ну, этого, как его, Стер... ? Да я вообще такого не знаю!» и «Почему интервалы распределены таким непонятным образом?!». Кроме того, среди читателей найдутся и такие, которые не рискнут самостоятельно определить величину интервала.

Случаи, когда распределение непонятно, даже если величина интервала определена математическим способом, встречаются довольно часто. И здесь уместно вспомнить то, о чем шла речь в начале этой главы: таблицы (ряды) распределения позволяют систематизировать данные наблюдения и интуитивно понять общую ситуацию. Следовательно, вполне достаточно выбрать такую величину интервалов, которая будет понятна тем, кто проводит статистический анализ.

6. Теория оценивания и описательная статистика

Объясняя Руи, что такое статистика, Игараси-сан определил её как науку, изучающую «большие совокупности однородных объектов на основании их выборочного исследования». Это не совсем так.

В статистике можно выделить два раздела: теорию оценивания и описательную статистику. В прологе речь шла о теории оценивания. Тогда что же такое описательная статистика? Это набор методов по упорядочиванию данных с целью наиболее простого и ясного восприятия этих данных. Можно считать, что описательная статистика рассматривает выборку как генеральную совокупность.

Возможно, такое определение кому-то покажется абстрактным и сложным для понимания. Поэтому приведём пример. Ямамото-сан вычислял среднее значение очков, набранное игроками команды Руи, и стандартное отклонение, чтобы представить положение в команде Руи в наглядном виде. Именно такая статистика и есть описательная.

Упражнение

Результаты забега на 100 метров приведены в следующей таблице:

Участник забега	Результат бега на 100 м, с
А	16,3
Б	22,4
В	18,5
Г	18,7
Д	20,1

Рассчитайте на основании этих результатов

- средний результат,
- медиану,
- отклонение.

Ответ

$$\text{Средний результат} = \frac{16,3 + 22,4 + 18,5 + 18,7 + 20,1}{5} = \frac{96}{5} = 19,2$$

Медиана: 16,3 18,5 **18,7** 20,1 22,4

Стандартное отклонение =

$$\begin{aligned} &= \sqrt{\frac{(16,3 - 19,2)^2 + (22,4 - 19,2)^2 + (18,5 - 19,2)^2 + (18,7 - 19,2)^2 + (20,1 - 19,2)^2}{5}} = \\ &= \sqrt{\frac{(-2,9)^2 + 3,2^2 + (-0,7)^2 + (-0,5)^2 + 0,9^2}{5}} = \\ &= \sqrt{\frac{20,2}{5}} = \\ &= \sqrt{4,04} = \\ &\approx 2,01 \end{aligned}$$

Выводы

- Чтобы «интуитивно» понять общую ситуацию с данными, строят ряды распределения.
- Величина интервала в рядах распределения определяется по формуле Стерджесса.
- Чтобы «математически» понять ситуацию с данными, вычисляют среднюю величину, медиану и стандартное отклонение.
- Если ряд распределения содержит слишком большие или слишком малые значения, рассчитывают медиану, а не среднюю величину.
- Стандартное отклонение — показатель, отражающий степень разброса (рассеяния) значений.

Глава 3

Знакомимся с качественными данными

1. Простые статистические таблицы

Надеюсь, ты помнишь, что качественные данные — это данные, которые нельзя измерить?

Помню, в общих чертах!

(ставит чашку)



Что?



Ты сегодня в школьной форме...

А, это?



Уже скоро...

... я скажу ей прощай.



Ты что, заканчиваешь школу? А как же ещё год?

Скоро в нашей школе введут новую форму.



(звук франфар)
Вот такую! ♡



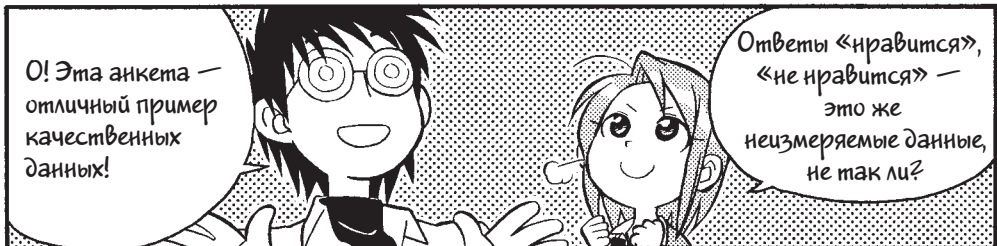
Матроска
в клетку?...
Необычно.

Поэтому-то,
в нашем классе
было проведено
анкетирование.

Анкета. Нравится ли вам новая форма?

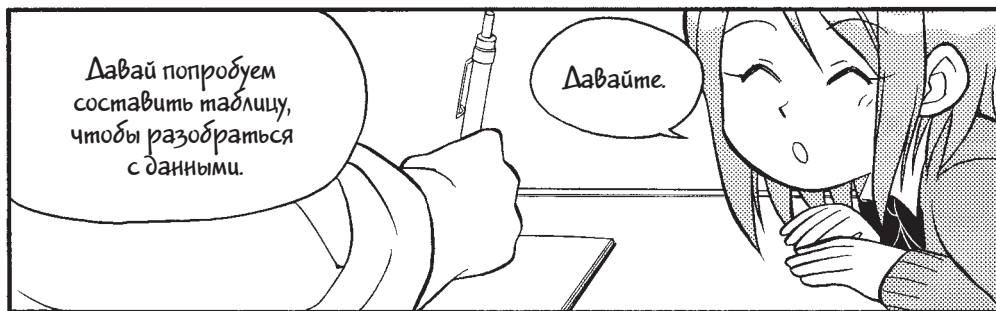
	Новая форма...	Новая форма...	Новая форма...	
1	нравится	16	так себе	
2	так себе	17	нравится	
3	нравится	18	нравится	
4	так себе	19	нравится	
5	не нравится	20	нравится	
6	нравится	21	нравится	
7	нравится	22	нравится	
8	нравится	23	не нравится	
9	нравится	24	так себе	
10	нравится	25	нравится	
11	нравится	26	нравится	
12	нравится	27	не нравится	
13	так себе	28	нравится	
14	нравится	29	нравится	
15	нравится	30	нравится	
			31	так себе
			32	так себе
			33	нравится
			34	не нравится
			35	нравится
			36	нравится
			37	нравится
			38	нравится
			39	так себе
			40	нравится

Вот
результаты.



О! Эта анкета —
отличный пример
качественных
данных!

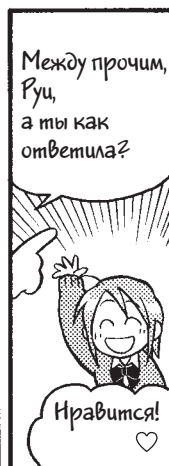
Ответы «нравится»,
«не нравится» —
это же
неизмеряемые данные,
не так ли?



Оценки новой школьной формы

Ответ	Кол-во	%
нравится	28	70
так себе	8	20
не нравится	4	10
Итого:	40	100

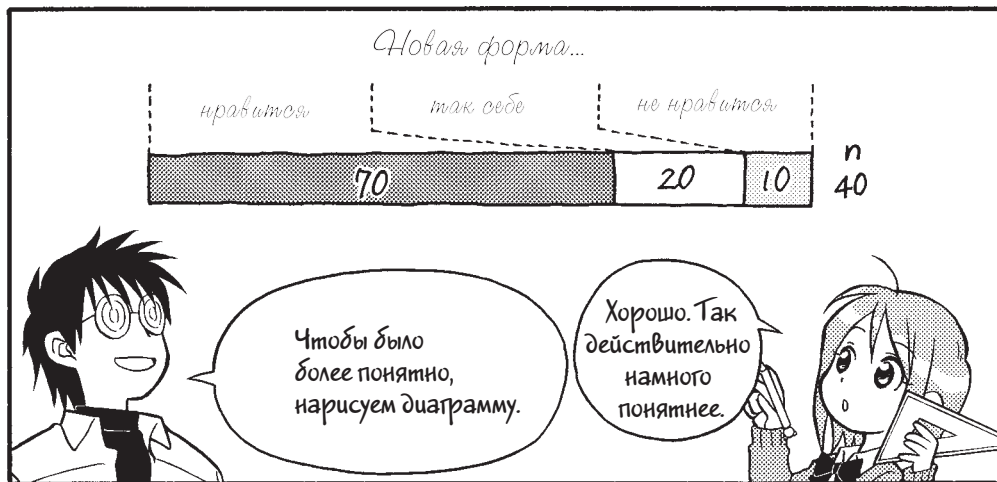
Это простая статистическая таблица.



Значит, процентное соотношение будет таким:

$$\frac{28}{40} \times 100 = \frac{7}{10} \times 100 = 70(\%)$$

OK!



Упражнение

Одна газета провела анкетирование по вопросу, какая из двух политических партий победит на следующих выборах.

Результаты анкетирования приведены в виде следующей таблицы:

Респондент	Победит партия А или Б
1	победит Б
2	победит Б
3	победит Б
4	не знаю
5	победит А
6	победит Б
7	победит А
8	не знаю
9	победит Б
10	победит Б

Используя результаты анкетирования, постройте простую статистическую таблицу.

Ответ

Простая статистическая таблица выглядит так:

Оценка партии	Частота	%
Победит А	2	20
Не знаю	2	20
Победит Б	6	60
Итого	10	100

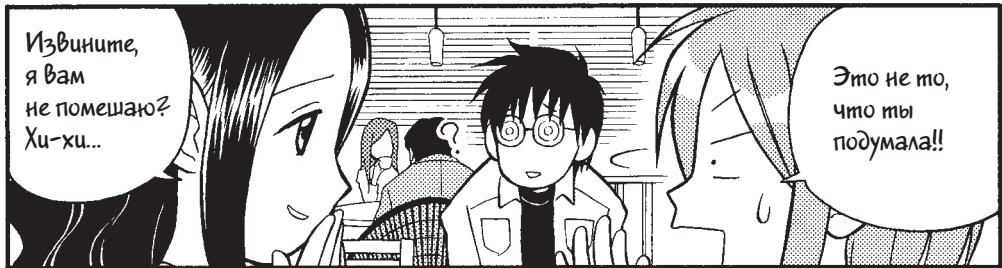
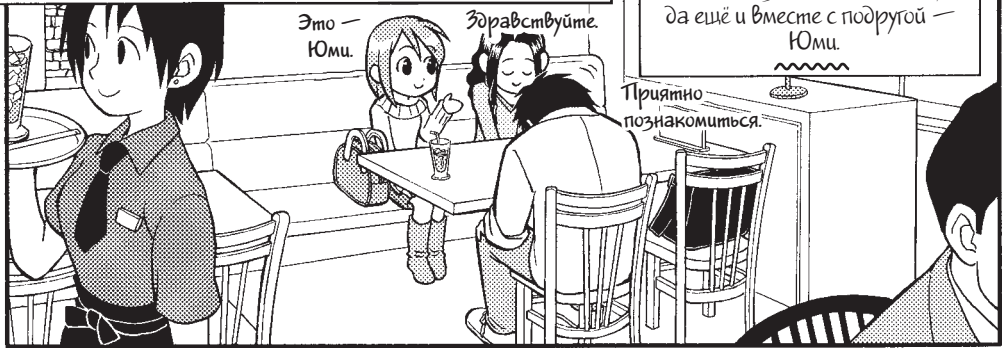
Выводы

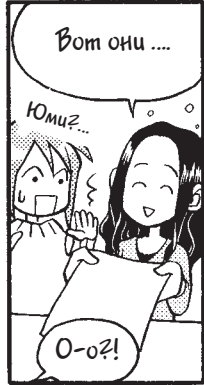
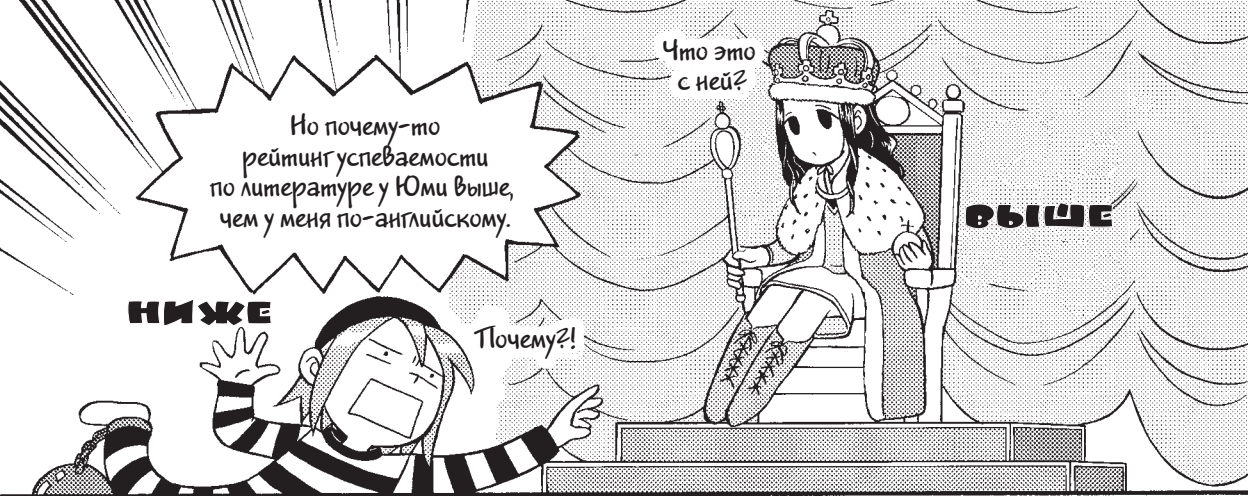
- Чтобы понять общую ситуацию с данными, строят простые статистические таблицы.

Глава 4

Нормированное отклонение и рейтинг успеваемости

1. Нормирование и нормированное отклонение





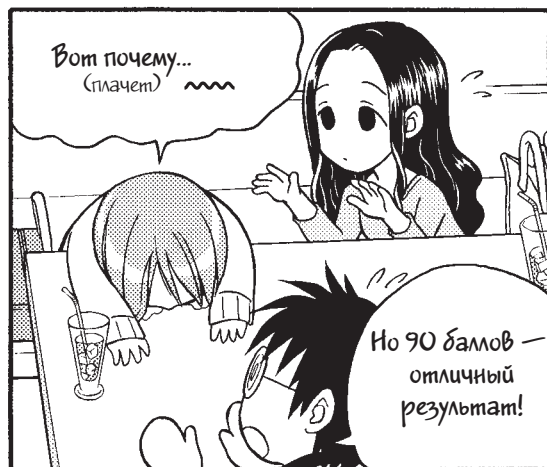
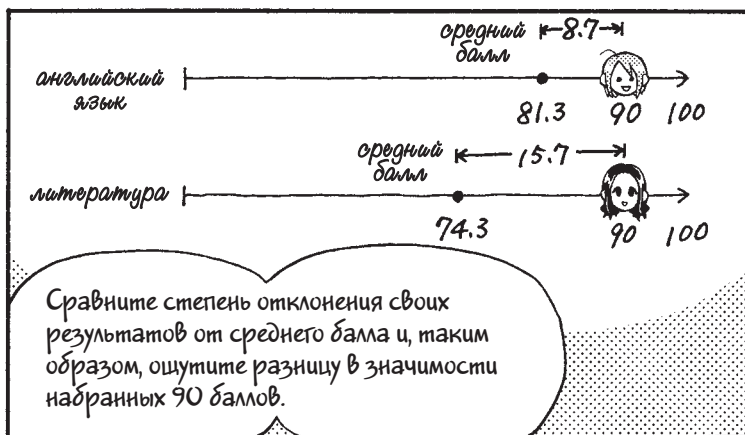
Результаты тестов (максимальное количество баллов – 100)

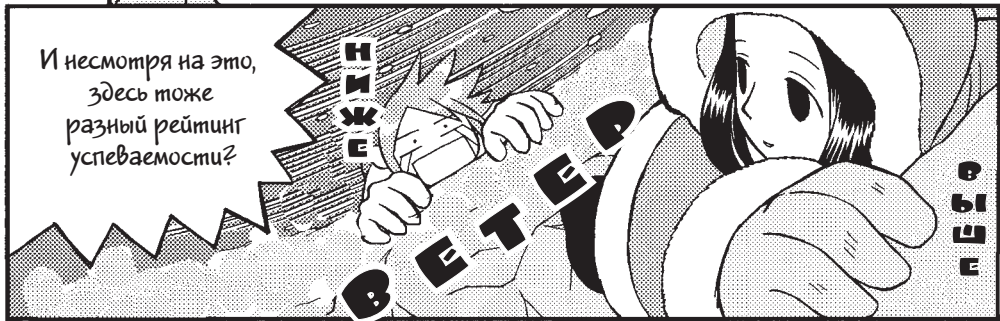
Ученик	Английский язык	Литература	Ученик	Английский язык	Литература
Руи	90	71	З	67	85
Юми	81	90	И	87	93
А	73	79	К	78	89
Б	97	70	Л	85	78
В	85	67	М	96	74
Г	60	66	Н	77	65
Д	74	60	О	100	78
Е	64	83	П	92	53
Ж	72	57	Ф	86	80



Готово!

Средний балл
по английскому = 81.3
по литературе = 74.3





Хотя отклонение от среднего результата одинаковое.

Угу, понятно.

Ученик	История	Биология
Руд	73	59
Юта	61	73
А	14	47
Б	41	38
В	49	63
Г	87	56
Д	69	15
Е	65	53
Ж	36	80
Средний балл	53	53

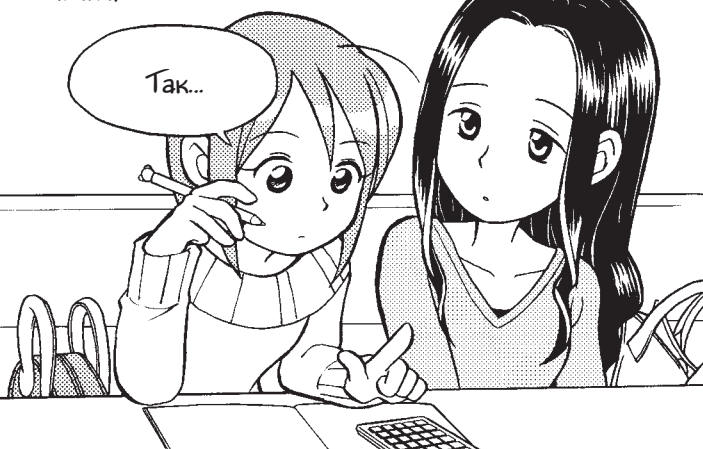
Ученик	История	Биология
З	7	50
И	53	41
К	100	62
Л	57	44
М	45	26
Н	56	91
О	34	35
П	37	53
Ф	70	68
Средний балл	53	53



Формула, по-моему, такая:

$$\sqrt{\frac{(\text{значение} - \text{среднее значение})^2}{\text{кол-во значений}}}$$

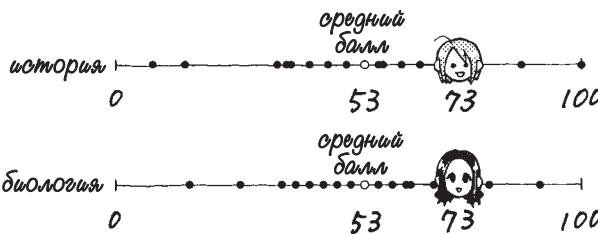
Верно?



Стандартное отклонение равно:
 история = 22,7
 биология = 18,3

Готово!

Чем меньше стандартное отклонение, тем меньше разброс данных, поэтому получается, что результаты теста по биологии у всех более-менее похожи в отличие от результатов теста по истории.



Что это значит...?

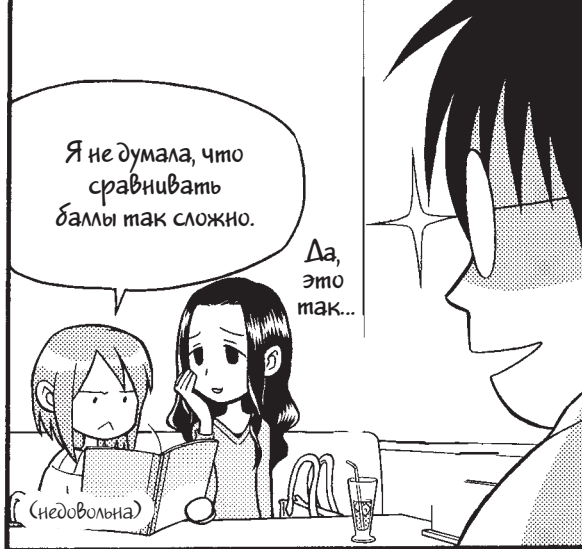
Если бы я готовился к поступлению в институт, я бы поднажал на биологию, так как в тесте по биологии значимость балла выше, чем в тесте по истории.

Какие-то 1-2 бала могут существенно повлиять на конечный результат.

Как ему идёт форма старшеклассника...



Xu-Xu-Xu



Нормирование проводят так:

Нормированное отклонение = (Z-показатель)

$$\frac{\text{Значение} - \text{Среднее значение}}{\text{Стандартное отклонение}}$$



Нормированные отклонения образуют совокупность нормированных значений.



Применительно к тестам нормированное отклонение (соно же Z-показатель) имеет другое название — стандартизованный балл.

Согласны!



Результаты тестов по истории и биологии

Ученик	История	Биология
Руи	73	59
Ю.ми	61	73
А	14	47
Б	41	38
В	49	63
Г	87	56
Д	69	15
Е	65	53
Ж	36	80
З	7	50
И	53	41
К	100	62
Л	57	44
М	45	26
Н	56	91
О	34	35
П	37	53
Ф	70	68
Средний балл	53	53
Станд. отклонение	22.7	18.3

Нормированные отклонения

История	Биология
0.88	0.33
0.35	1.09
-1.71	-0.33
-0.53	-0.82
-0.18	0.55
1.49	0.16
0.70	-2.08
0.53	0
-0.75	1.48
-2.02	-0.16
0	-0.66
2.07	0.49
0.18	-0.49
-0.35	-1.48
0.13	2.08
-0.84	-0.98
-0.70	0
0.75	0.82
0	0
1	1

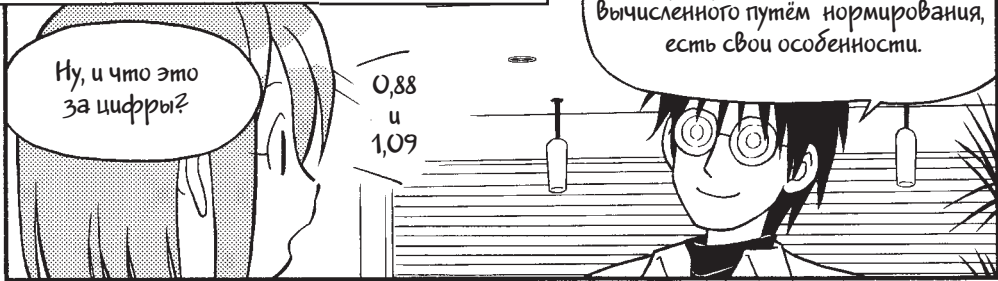
Нормированное отклонение Руи по истории = $\frac{73 - 53}{22.7} = \frac{20}{22.7} = 0.88$

Нормированное отклонение Ю.ми по биологии = $\frac{73 - 53}{18.3} = \frac{20}{18.3} = 1.09$

Да, так.



2. Свойства нормированного отклонения



1. Независимо от максимального количества баллов, среднее значение нормированного отклонения (Z -показателя) всегда равно 0 , а стандартное отклонение нормированных отклонений всегда равно 1 .



Можно сравнивать результаты тестов с максимальным количеством баллов, равным 100 и 200 .

2. В чём бы ни измерялась переменная, среднее значение её нормированных отклонений всегда равно 0 , а стандартное отклонение нормированных отклонений всегда равно 1 .



Можно также сравнивать количество ударов по воротам или угловых в футболе.



3. Рейтинг успеваемости

Рейтинг успеваемости (он же Т-показатель) рассчитывается на основе нормированного отклонения.

О-ооо...
(удивляется)

Формула такая:

$$\text{Рейтинг успеваемости} = \text{Т-показатель} = \\ = \text{Нормированное отклонение} \times 10 + 50$$

И правда, используется значение нормированного отклонения.

Попытаемся вычислить рейтинг успеваемости в ваших тестах.

Рус
(история) $0.88 \times 10 + 50 = 8.8 + 50 = 58.8$

Рус
(биология) $1.09 \times 10 + 50 = 10.9 + 50 = 60.9$

Да-да, такие цифры и были.

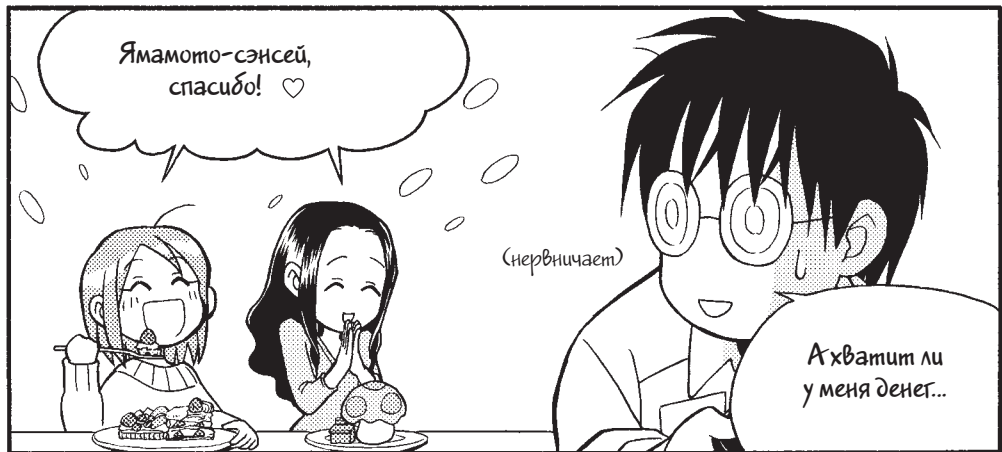
Нормированное отклонение имеет такие свойства:

Нормированное отклонение (Z-показатель):

1. Независимо от максимального количества баллов, среднее значение нормированного отклонения (Z-показателя) всегда = 0, а стандартное отклонение нормированных отклонений всегда = 1.
2. В чём бы ни измерялась переменная, среднее значение её нормированных отклонений всегда = 0, а стандартное отклонение нормированных отклонений всегда = 1.

Рейтинг успеваемости (Т-показатель):

1. Независимо от максимального количества баллов, среднее значение рейтинга успеваемости (Т-показателей) всегда = 50, а квадратичное отклонение рейтингов успеваемости всегда = 10.
2. В чём бы не измерялась переменная, среднее значение рейтинга успеваемости всегда = 50, а стандартное отклонение рейтингов успеваемости всегда = 10.



4. Что такое рейтинг успеваемости?

В общем случае это Т-показатель, который вычисляется по формуле:

$$\begin{aligned} \text{Т-показатель} &= \text{Нормальное отклонение} \times 10 + 50 = \\ &= \frac{\text{Значение} - \text{Среднее значение}}{\text{Стандартное отклонение}} \times 10 + 50. \end{aligned}$$

В классе Руи 40 учеников, и из них 18 девочек. Рейтинги успеваемости были приведены только для двух девочек. Если объектом исследования были бы все ученики класса, среднее значение и стандартное отклонение были бы совсем другими, и, естественно, рейтинги успеваемости Руи и Юми тоже были бы другими. Если бы рассчитывались рейтинги успеваемости всех учеников класса, показатель Руи был бы выше.

Результаты теста всех учеников представлены в Табл. 4.1. Обязательно попробуйте вычислить рейтинги успеваемости. Забегая вперёд, скажу что у Руи по истории он равен 59,1, а у Юми по биологии он равен 56,7.

Теперь представим, что одинаковые тесты проводились в двух классах. Сначала вычислили отдельно средний балл и стандартное отклонение для учеников 1-го класса, и затем на их основе рассчитали рейтинг успеваемости. То же проделали и для 2-го класса. Показатель успеваемости ученика А из 1-го класса был равен 57, и рейтинг успеваемости ученика Б из 2-го класса тоже был равен 57. На первый взгляд кажется, что у них одинаковые знания. Однако, такие необходимые для расчёта рейтинга успеваемости данные, как средний балл и стандартное отклонение, в 1-м и 2-м классах были разными. Следовательно, нельзя сравнивать рейтинги успеваемости этих двух учеников.

И еще один пример. Ученик А в апреле сдавал пробный экзамен на подготовительных курсах. Его рейтинг успеваемости был равен 54. Все лето А продолжал заниматься на других курсах. А в сентябре, чтобы убедиться, что его труды не пропали даром, он решил сдать еще один пробный экзамен, но уже на других курсах. На этот раз рейтинг успеваемости был 62. На первый взгляд кажется, что знания А улучшились. Но он сдавал пробные экзамены на разных курсах, и, следовательно, организаторы проведения экзамена в апреле и в сентябре были разные. Это, в свою очередь, означает, что средний результат и стандартное отклонение, т.е. показатели, необходимые для вычисления рейтинга успеваемости, в апреле и сентябре были разные, и, следовательно, нельзя сравнивать два значения рейтинга успеваемости А.

Ну как? Понятие показатель успеваемости довольно-таки сложное.

Таблица 4.1. Результаты теста по истории и биологии
(всех учеников класса Руи)

Девочка	История	Биология
Руи	73	59
Юми	61	73
А	14	47
Б	41	38
В	49	63
Г	87	56
Д	69	15
Е	65	53
Ж	36	80
З	7	50
И	53	41
К	100	62
Л	57	44
М	45	26
Н	56	91
О	34	35
П	37	53
Р	70	68

Мальчик	История	Биология
а	54	2
б	93	7
в	91	98
г	37	85
д	44	100
е	16	29
ж	12	57
з	44	37
и	4	95
к	17	39
л	66	70
м	53	14
н	14	97
о	73	39
п	6	75
р	22	80
с	69	77
т	95	14
у	16	24
ф	37	91
х	14	36
ц	88	76

Средний балл по всему классу	48,0	54,9
Стандартное отклонение по всему классу	27,5	26,9

Упражнение

В упражнении на стр. 57 приведены результаты бега на 100 метров.

Участник	Результат бега на 100 м, с
А	16,3
Б	22,4
В	18,5
Г	18,7
Д	20,1
Среднее значение	19,2
Стандартное отклонение	2,01

Проверьте с помощью этой таблицы:

1. Равно ли 0 среднее значение нормированных отклонений.
2. Равно ли 1 стандартное отклонение нормированных отклонений.

Ответ

1. Среднее значение нормированных отклонений =

$$= \frac{\left(\frac{16,3 - 19,2}{2,01}\right) + \left(\frac{22,4 - 19,2}{2,01}\right) + \left(\frac{18,5 - 19,2}{2,01}\right) + \left(\frac{18,7 - 19,2}{2,01}\right) + \left(\frac{20,1 - 19,2}{2,01}\right)}{5} =$$

$$= \frac{\left\{ (16,3 - 19,2) + (22,4 - 19,2) + (18,5 - 19,2) + (18,7 - 19,2) + (20,1 - 19,2) \right\}}{2,01 \cdot 5} =$$

упорядочили
числитель

$$= \frac{\left\{ 16,3 + 22,4 + 18,5 + 18,7 + 20,1 - 19,2 - 19,2 - 19,2 - 19,2 - 19,2 \right\}}{2,01 \cdot 5} =$$

$$= \frac{\left\{ 96 - 19,2 \times 5 \right\}}{2,01 \cdot 5} =$$

отдельно индивидуальные значения,
отдельно средние значения (-19,2).

$$= \frac{\left\{ 96 - 96 \right\}}{2,01 \cdot 5} =$$

$$= \frac{0}{5} =$$

$$= 0$$

2. Стандартное отклонение нормированных отклонений =

$$= \sqrt{\frac{\left(\frac{16,3 - 19,2}{2,01} - 0\right)^2 + \left(\frac{22,4 - 19,2}{2,01} - 0\right)^2 + \left(\frac{18,5 - 19,2}{2,01} - 0\right)^2 + \left(\frac{18,7 - 19,2}{2,01} - 0\right)^2 + \left(\frac{20,1 - 19,2}{2,01} - 0\right)^2}{5}} =$$

$$= \sqrt{\frac{\left(\frac{16,3 - 19,2}{2,01}\right)^2 + \left(\frac{22,4 - 19,2}{2,01}\right)^2 + \left(\frac{18,5 - 19,2}{2,01}\right)^2 + \left(\frac{18,7 - 19,2}{2,01}\right)^2 + \left(\frac{20,1 - 19,2}{2,01}\right)^2}{5}} =$$

$$= \sqrt{\frac{\left\{ (16,3 - 19,2)^2 + (22,4 - 19,2)^2 + (18,5 - 19,2)^2 + (18,7 - 19,2)^2 + (20,1 - 19,2)^2 \right\}}{2,01^2 \cdot 5}} =$$

упорядочили
числитель

$$= \sqrt{\frac{1}{2,01^2} \times \left\{ (16,3 - 19,2)^2 + (22,4 - 19,2)^2 + (18,5 - 19,2)^2 + (18,7 - 19,2)^2 + (20,1 - 19,2)^2 \right\}} =$$

$$= \frac{1}{2,01} \times \sqrt{\left\{ (16,3 - 19,2)^2 + (22,4 - 19,2)^2 + (18,5 - 19,2)^2 + (18,7 - 19,2)^2 + (20,1 - 19,2)^2 \right\}} =$$

$$= \frac{1}{\text{Стандартное отклонение}} \times \text{Стандартно отклонение} =$$

$$= 1$$

Внимательно посмотрите на таблицу на стр. 78.

Выводы

- Нормирование (нормировка, стандартизация, нормализация), или Z-преобразование, — преобразование значений, проводимое на основе данных о степени разброса (рассеяния) и отклонения от среднего значения. Нормирование позволяет оценить значимость значений.
- Нормирование позволяет сравнивать различные переменные величины, например:
 - величины, имеющие разный размах (разность между максимальным и минимальным значениями);
 - величины, имеющие разные единицы измерения.
- Нормированные данные — нормированные отклонения отдельных значений.
- Рейтинг успеваемости рассчитывается по формуле T-показателя на основе нормированного отклонения.

Глава 5

Вычислим вероятность

1. Функция распределения плотности вероятности



В статистике иногда говорят:
«вероятность чего-то меньше 0,05»

Сегодня я расскажу,
что нужно знать,
чтобы вычислить
вероятность чего-либо.



Извините!
Вероятность —
это та самая вероятность,
про которую говорят
в прогнозе погоды?

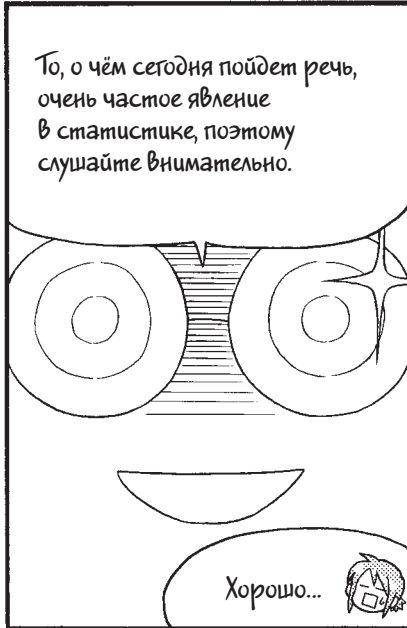
Да.



Содержание
сегодняшнего
занятия
несколько
абстрактное.

Абстрактное?

ДУИ



То, о чём сегодня пойдет речь,
очень частое явление
в статистике, поэтому
слушайте внимательно.

Хорошо...

Результаты теста по английскому языку всех
одинадцатиклассников школ Центрального округа

Уровень	Балл
1	42
2	91
...	...
10421	50
Средний балл	53
Стандартное отклонение	10

Представим,
что все ученики
11-х классов
Центрального
округа ...

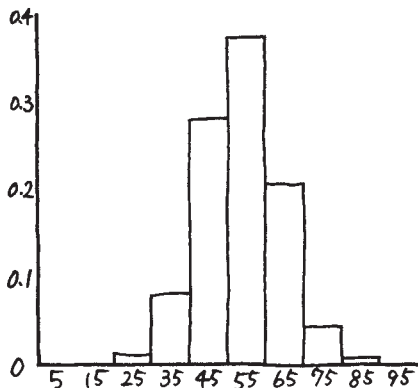
... сдавали экзамен
на подготовительных
курсах.

Вы сегодня
очень хорошо
подготовлены.

Ха-Ха-Ха.
Мы только
начали.

Если эту таблицу
представить
в виде гистограммы,
то получим ...

Гистограмма
«Результаты теста по английскому языку».
Величина интервала равна 10.



О-о-о...
Действительно,
если нарисовать
гистограмму, можно
что-то
понять.

Потому что
наглядно.

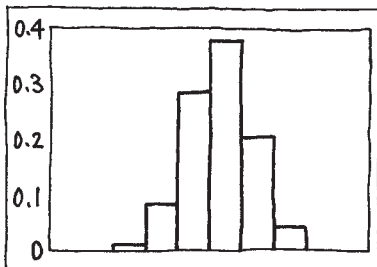
Что произойдет,
если на этой гистограмме
уменьшить величину
интервала?

Что?

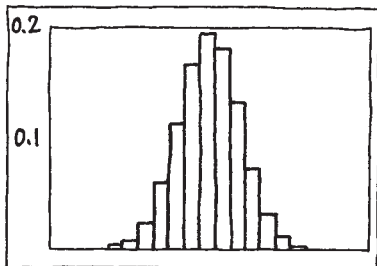
(Амалото-сан
уменьшился)

Величина интервала и гистограмма «Результаты теста по английскому языку»

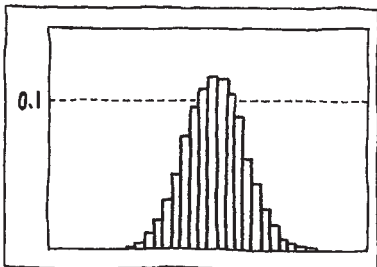
Величина интервала
равна 10



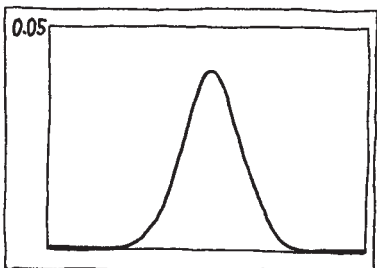
Величина интервала
равна 5



Величина интервала
равна 3

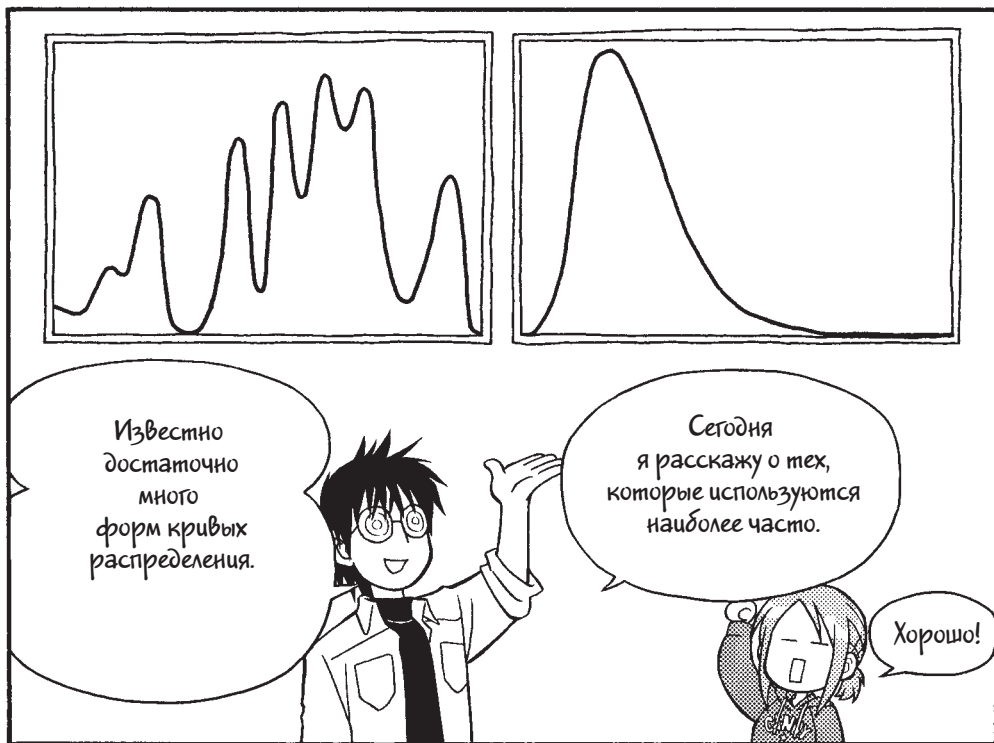


Кривая
распределения



О...
Превращается
в непрерывную
линию!





2. Нормальное распределение

$$f(x) = \frac{1}{\sqrt{2\pi} \times \text{Стандартное отклонение}} e^{-\frac{1}{2} \left(\frac{x - \bar{x}}{\text{Стандартное отклонение}} \right)^2}$$

где \bar{x} — средняя величина x (средняя арифметическая ряда)

Вот.

Что это ?!!

Это — часто встречающаяся в статистике функция распределения вероятности.

А что это за знак «e»?

Символ «e» — математическая константа, основание натурального логарифма; его иногда называют числом Эйлера или числом Непера.
 $e = 2,7182$.

ОСНОВАНИЕ

Ха-Ха-Ха

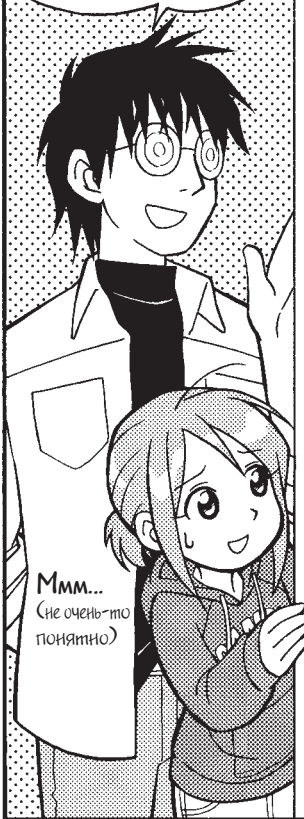
Можешь думать, что это как число π

Ну, тогда ладно...

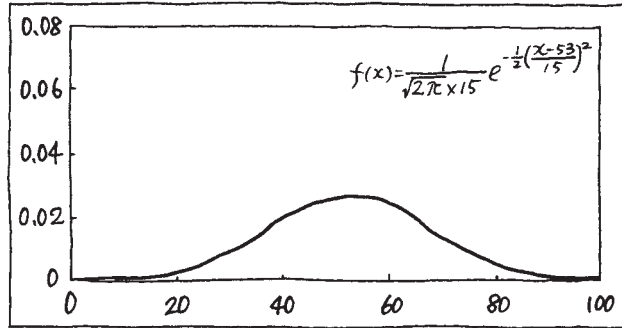
Фу...

График функции распределения вероятности имеет следующие свойства:

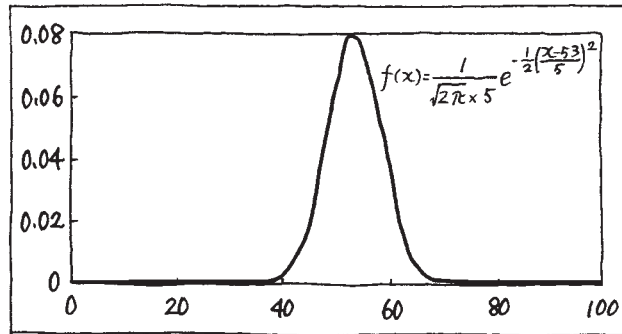
1. Кривая симметрична относительно центра распределения, который находится в точке, соответствующей среднему значению.
2. Функция зависит от среднего значения и от стандартного отклонения.



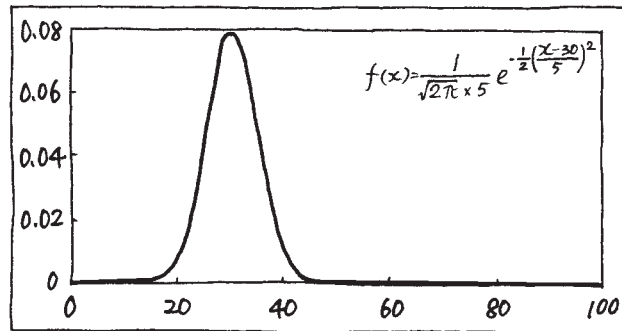
Среднее значение = 53, станд. отклонение = 15



Среднее значение = 53, станд. отклонение = 5



Среднее значение = 30, станд. отклонение = 5

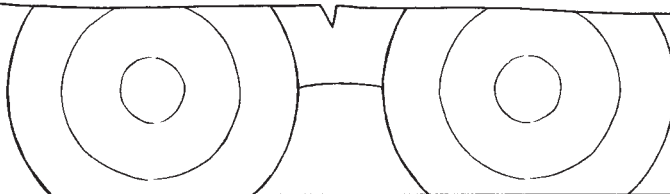


Послушай, существует правило, согласно которому ...



... распределение величины x при определённых значениях средней арифметической ряда (или среднего значения, \bar{x}) и стандартного отклонения называют нормальным распределением, если плотность распределения вероятностей выражается формулой

$$f(x) = \frac{1}{\sqrt{2\pi} \times \text{Стандартное отклонение}} e^{-\frac{1}{2} \left(\frac{x - \bar{x}}{\text{Ст.откл.}} \right)^2}$$



Чего ???!!!



... называют нормальным распределением... !?

Ничего не понимаю...

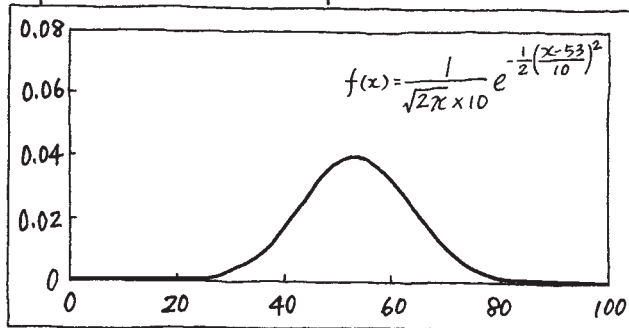


Правило довольно специфическое, поэтому просто запомни его.

Возьмём для примера предыдущий тест. Если кривая функции распределения вероятностей результатов теста по английскому имеет следующий вид...



... то имеет место нормальное распределение при среднем балле 53 и стандартном отклонении 10.





3. Стандартное нормальное распределение



А если плотность распределения вероятностей выражается формулой

$$f(x) = \frac{1}{\sqrt{2\pi} \times \text{Ст.откл.}} e^{-\frac{1}{2} \left(\frac{x - \bar{x}}{\text{Ст.откл.}} \right)^2} = \frac{1}{\sqrt{2\pi} \times 1} e^{-\frac{1}{2} \left(\frac{x - 0}{\text{Ст.откл.}} \right)^2} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} x^2}$$

то в статистике не говорят, что величина x имеет нормальное распределение при значении средней арифметической ряда = 0 и стандартном отклонении = 1.

Принято говорить, что величина x имеет стандартное нормальное распределение.

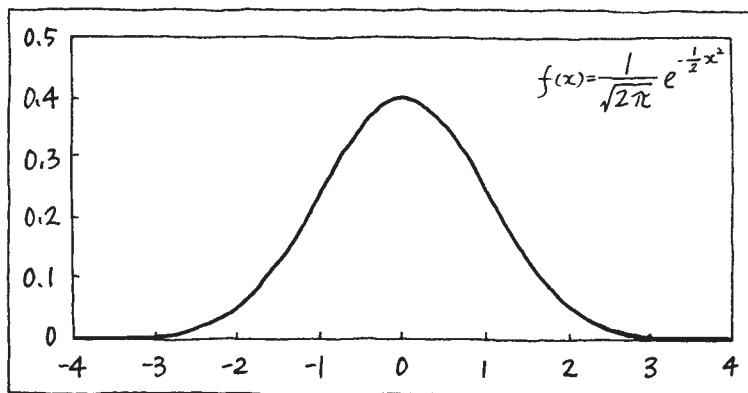


Ученик	Балл	Нормированное отклонение
1	42	-1.1
2	91	3.8
...	⋮	⋮
10421	50	-0.3
среднее значение	53	0
стандартное отклонение	10	1

$$\frac{\text{Балл} - \text{Средний балл}}{\text{Стандартное отклонение}} = \frac{50 - 53}{10} = \frac{-3}{10} = -0.3$$

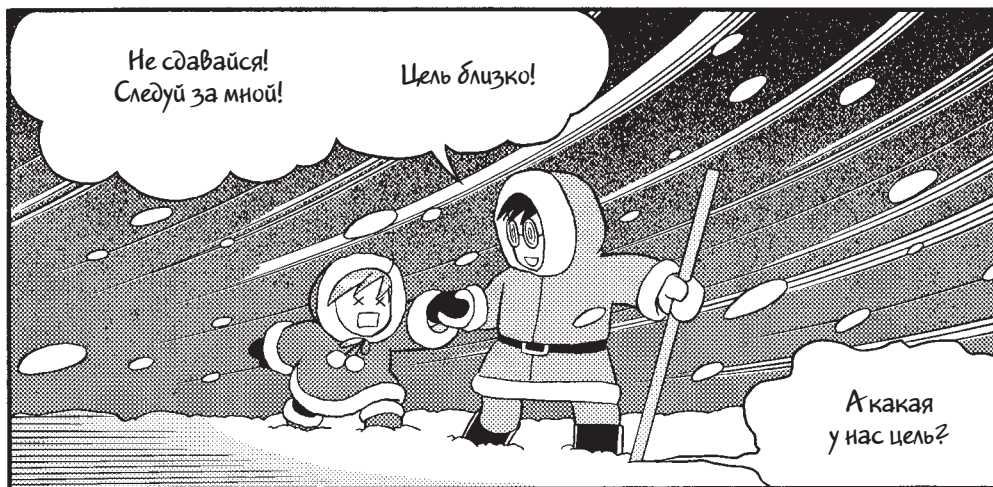
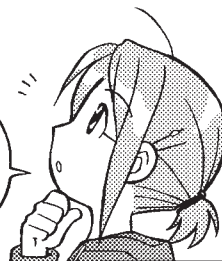
Если это так, то после нормирования результаты теста по английскому ...

Стандартное нормальное распределение



... будут иметь стандартное нормальное распределение.

Понятно!



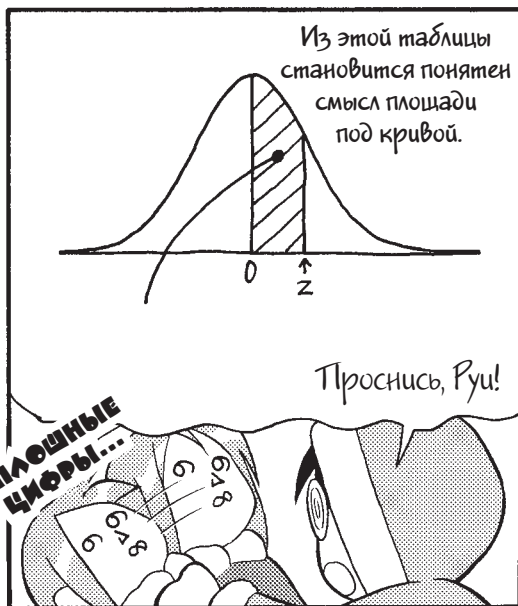
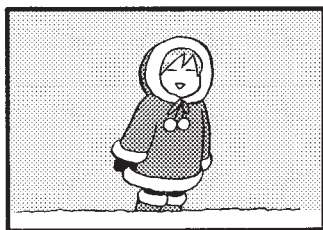
Не сдавайся!
Следуй за мной!

Цель близко!

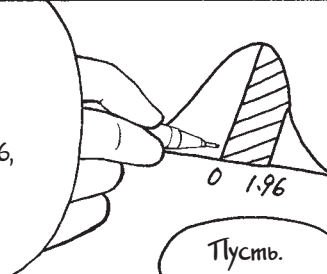
А какая у нас цель?

Таблица стандартного нормального распределения

Z	0,00	0,01	0,02	0,03	0,04	0,05	0,06	0,07	0,08	0,09
0,0	0,0000	0,0040	0,0080	0,0120	0,0160	0,0199	0,0239	0,0279	0,0319	0,0359
0,1	0,0398	0,0438	0,0478	0,0517	0,0557	0,0596	0,0636	0,0675	0,0714	0,0753
0,2	0,0793	0,0832	0,0871	0,0910	0,0948	0,0987	0,1026	0,1064	0,1103	0,1141
:	:	:	:	:	:	:	:	:	:	:
1,8	0,4641	0,4649	0,4656	0,4664	0,4671	0,4678	0,4686	0,4693	0,4699	0,4706
1,9	0,4713	0,4719	0,4726	0,4732	0,4738	0,4744	0,4750	0,4756	0,4761	0,4767
:	:	:	:	:	:	:	:	:	:	:



Пусть Z равно 1,96,



Пусть.

Итак, $Z = 1,96...$

$$Z = 1,9 + 0,06$$

Представим это значение в виде двух чисел:

Разделим десятые и сотые доли, так?

Теперь посмотрим на таблицу,

Z	0,00	0,01	0,02	0,03	0,04	0,05	0,06	0,07	0,08
0,0	0,0000	0,0040	0,0080	0,0120	0,0160	0,0199	0,0239	0,0279	0,0319
0,1	0,0398	0,0438	0,0478	0,0517	0,0557	0,0596	0,0636	0,0675	0,0715
0,2	0,0793	0,0832	0,0871	0,0910	0,0948	0,0987	0,1026	0,1064	0,1103
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
1,8	0,4641	0,4649	0,4656	0,4664	0,4671	0,4678	0,4686	0,4693	0,4700
1,9	0,4713	0,4719	0,4726	0,4732	0,4738	0,4744	0,4750	0,4756	0,4761

Пересечение строки 1,9 и столбца 0,06 ...

даёт 0,4750.



Правильно! Это площадь заштрихованной на графике области при $Z = 1,96$.

Забыл сказать:
и при нормальном распределении,
и при любом другом, площадь области,
ограниченной осью x и всей кривой распределения,
равна 1.

Площадь = 1

Надо же!

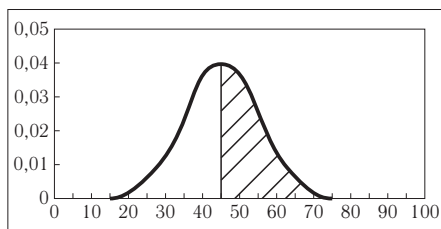


Пример 1

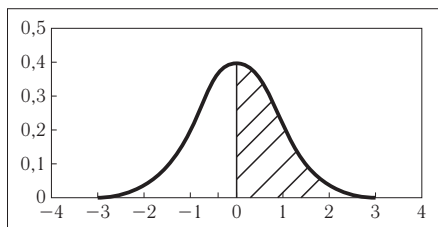


Все десятиклассники школ Восточного округа сдавали тест по математике. Когда им поставили оценки, стало ясно, что распределение результатов теста можно считать нормальным при среднем балле 45, и стандартном отклонении 10. Теперь хорошенько подумай. Следующие пять выводов имеют один и тот же смысл:

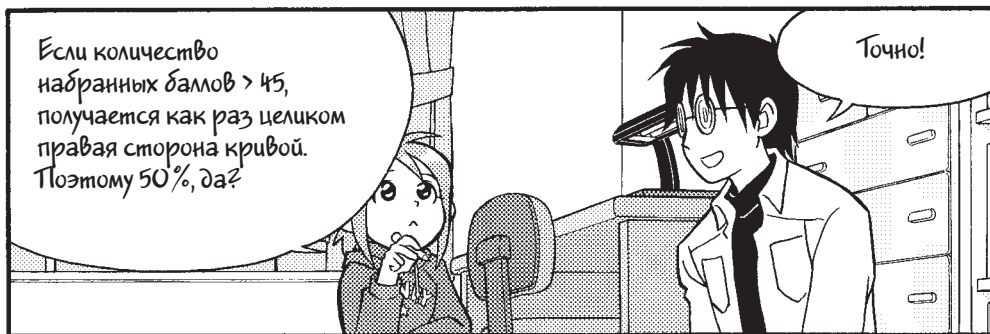
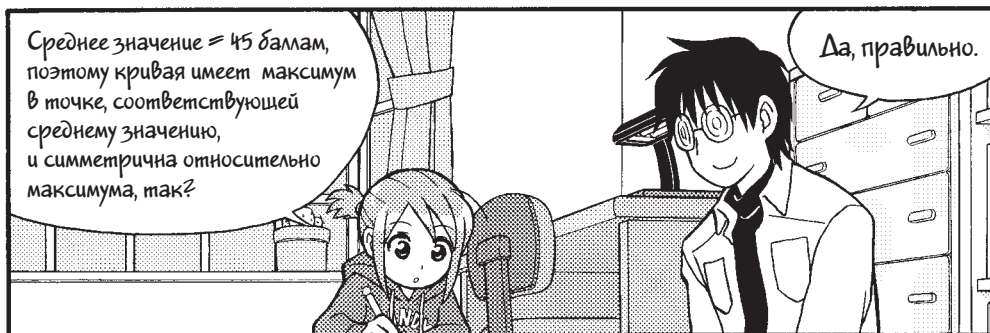
1. При нормальном распределении, когда среднее значение = 45 и стандартное отклонение = 10, площадь заштрихованной области = 0,5.



2. Доля учеников, чей результат был > 45 баллов, составляет 50% от общего числа участников теста.
3. Предположим, что из общего числа учеников произвольно выбрали одного. Вероятность того, что он набрал > 45 баллов, равна 50%.
4. Для стандартного нормального распределения, полученного после нормирования результатов теста, доля учеников с результатом > 0 составляет 50% от общего числа участников теста.



5. Предположим, что из общего числа учеников произвольно выбрали одного. При стандартном нормальном распределении, полученном после нормирования результатов теста по математике, вероятность того, что он набрал положительный балл, составляет 50%.

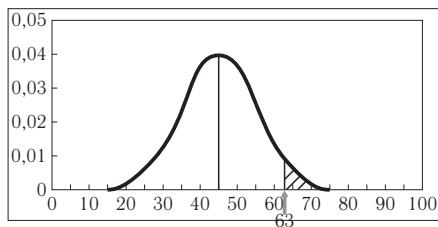


Пример 2

Все десятиклассники школ Восточного округа сдавали тест по математике. Когда им поставили оценки, стало ясно, что распределение результатов теста можно считать нормальным при среднем балле = 45 и стандартном отклонении = 10. Теперь напряги мозги. Как и в предыдущем примере, все следующие пять выводов имеют один и тот же смысл. Но на этот раз ты для начала прочти вывод 4.



1. При нормальном распределении, когда среднее значение = 45 и стандартное отклонение = 10, площадь заштрихованной на графике области = $0,5 - 0,4641 = 0,0359$.

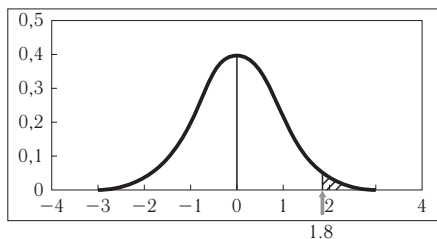


2. Доля учеников, чей результат > 63 баллов, = $0,5 - 0,4641 = 0,0359$ или 3,59% от общего числа сдававших экзамен.
3. Предположим, что из общего числа учеников был произвольно выбран один. Вероятность того, что он набрал > 63 баллов, = $0,5 - 0,4641 = 0,0359$ или 3,59%.
4. При нормальном распределении доля учеников с нормированным отклонением

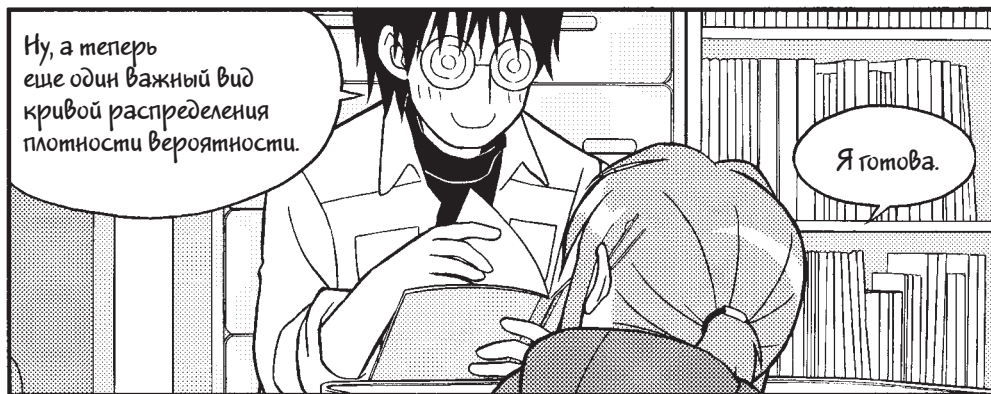
$$\text{балла} > 1,8 = \frac{18}{10} = \frac{63 - 45}{10} = \frac{\text{Значение} - \text{Среднее значение}}{\text{Стандартное отклонение}}$$

составляет 3,59% ($0,5 - 0,4641 = 0,0359$)

(см. Таблицу стандартного нормального распределения).



5. Предположим, что оценки учеников после нормирования распределены по стандартному нормальному закону. Вероятность того, что нормированное отклонение произвольно выбранного ученика > 1,8 равна 3,59%.



4. Распределение хи-квадрат



Если функция распределения вероятностей выражается формулой

$$f(x) = \begin{cases} \frac{1}{2^{\frac{n}{2}} \times \int_0^{\infty} x^{\frac{n}{2}-1} e^{-x} dx} \times x^{\frac{n}{2}-1} \times e^{-\frac{x}{2}} & \text{при } x > 0, \\ 0 & \text{при } x \leq 0, \end{cases}$$

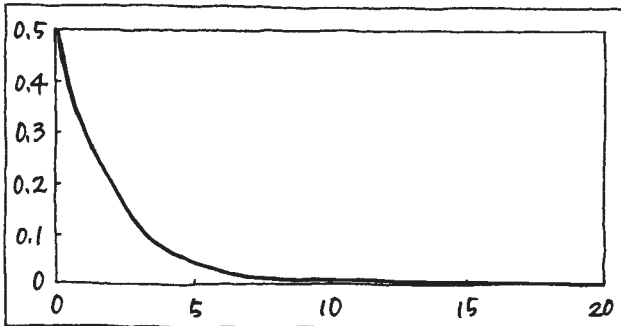
то в статистике говорят, что величина x имеет распределение хи-квадрат с числом степеней свободы n .



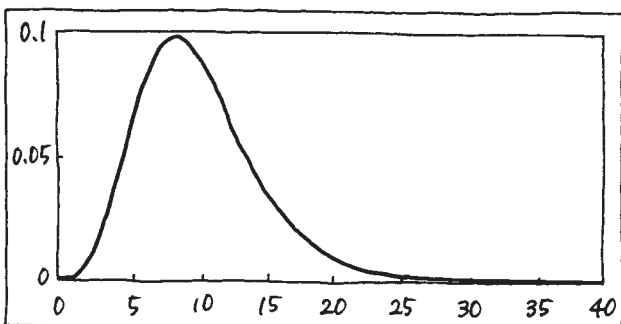
Спасите!



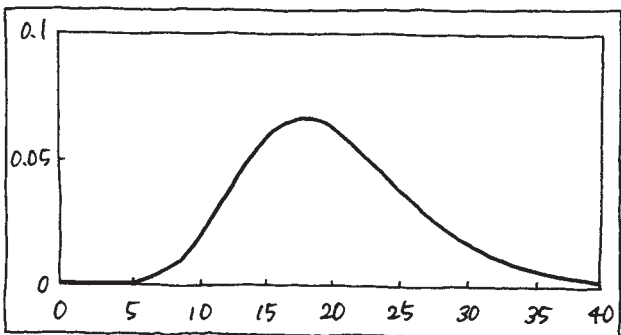
Число степеней свободы равно 2



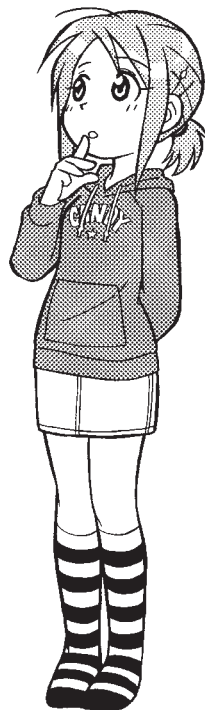
Число степеней свободы равно 10

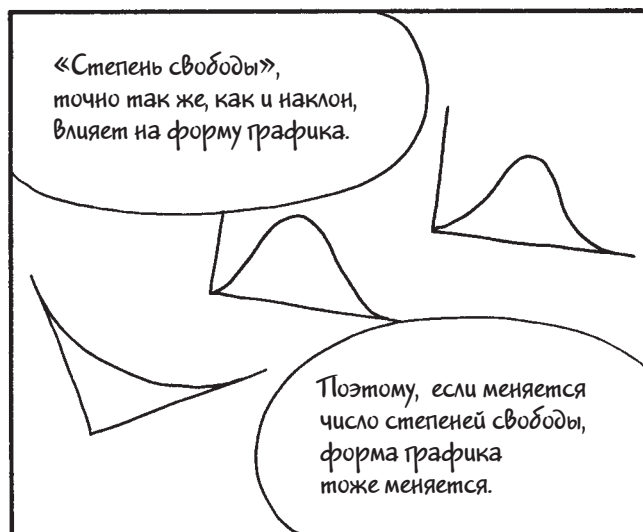


Число степеней свободы равно 20



В зависимости от числа степеней свободы форма графика совершенно меняется.

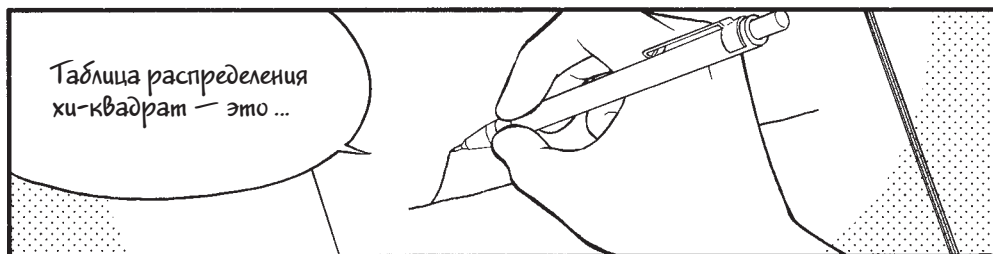




Так же, как существует
таблица стандартного
нормального распределения,
есть и таблица
для распределения
хи-квадрат.



Таблица распределения
хи-квадрат — это ...



... таблица, в которой указывается
значение χ^2 (см. ось x на графике),
соответствующее значению
вероятности (которая, как мы знаем,
равна площади и доле)
заштрихованной области P .

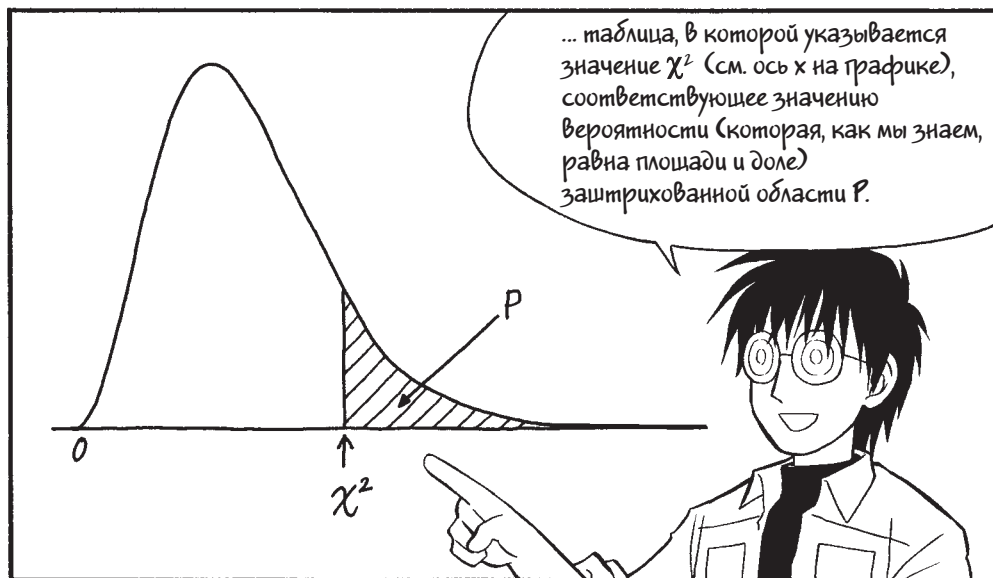
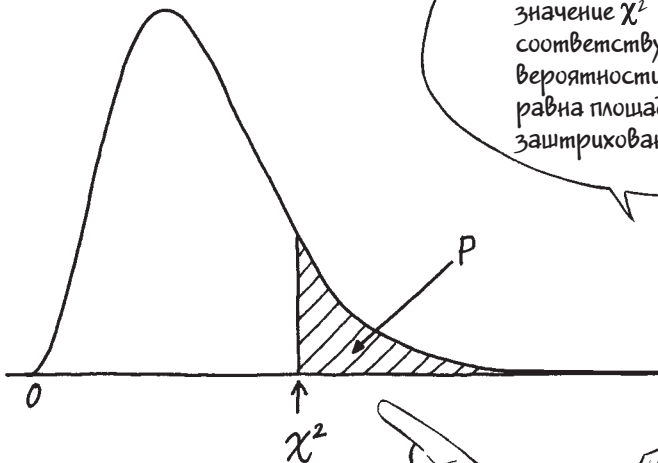




Таблица распределения хи-квадрат

Р Степень свободы	0.995	0.99	0.975	0.95	0.05	0.025	0.01	0.005
1	0.000039	0.0002	0.0010	0.0039	3.8415	5.0239	6.6349	7.8794
2	0.0100	0.0201	0.0506	0.1026	5.9915	7.3778	9.2104	10.5965
3	0.0717	0.1148	0.2158	0.3518	7.8147	9.3484	11.3449	12.8381
4	0.2070	0.2971	0.4844	0.7107	9.4877	11.1433	13.2767	14.8602
5	0.4118	0.5543	0.8312	1.1455	11.0705	12.8325	15.0863	16.7496
6	0.6757	0.8721	1.2373	1.6354	12.5916	14.4494	16.8119	18.5475
7	0.9893	1.2390	1.6899	2.1673	14.0671	16.0128	18.4753	20.2777
8	1.3444	1.6465	2.1797	2.7326	15.5073	17.5345	20.0902	21.9549
9	1.7349	2.0879	2.7004	3.3251	16.9190	19.0228	21.6660	23.5893
10	2.1558	2.5582	3.2470	3.9403	18.3070	20.4832	23.2093	25.1881
:	:	:	:	:	:	:	:	:

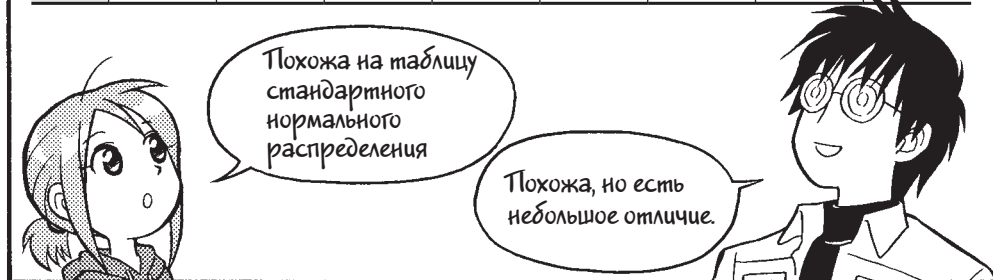


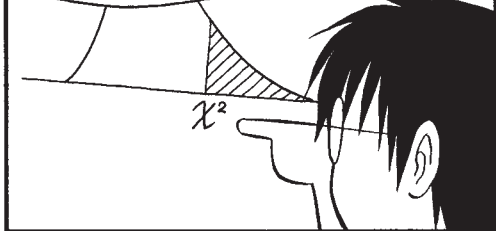
Таблица стандартного нормального распределения позволяет по значению координаты x (в пределах заштрихованной области) найти соответствующую вероятность.

Вероятность равна площади, или доле. Так?

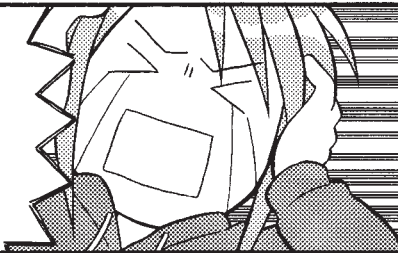


А таблица распределения хи-квадрат позволяет по вероятности найти соответствующую координату на оси x .

Вот это значение!



В голове у меня полная каша!!!



Ну подожди, не нервничай.

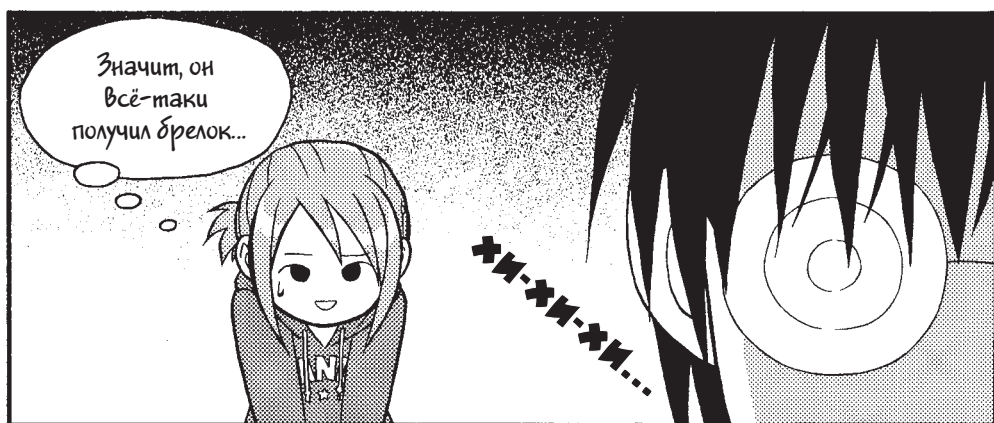
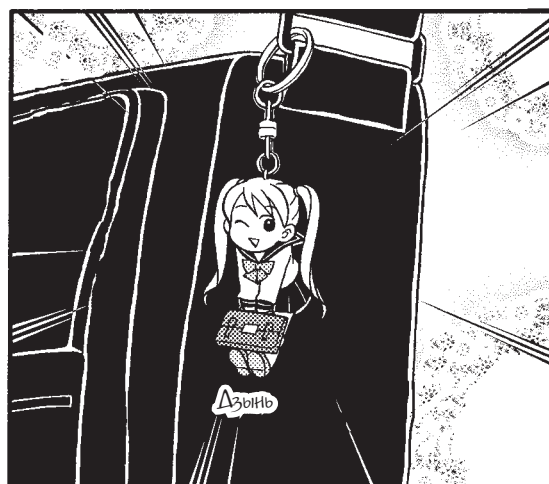
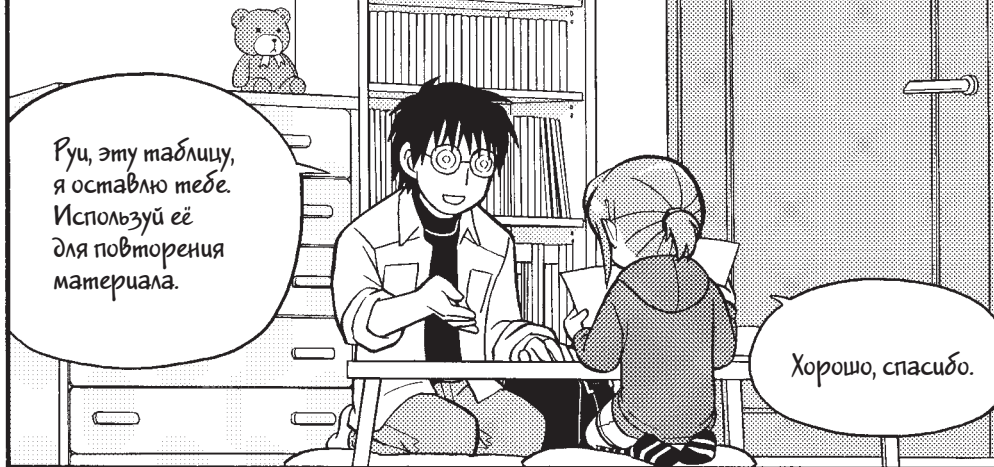


Давай посмотрим, какое будет значение, если число степеней свободы = 1 и $P = 0,05$.

0.99	0.975	0.95	0.05	0.025	0.01
0.002	0.0010	0.0039	3.8415	5.0239	6.6349
0.201	0.0506	0.1026	7.15	7.3778	9.2104
0.148	0.2158	0.3518	8.147	9.3484	11.3449
0.2971	0.4844	0.717	9.4877	11.1433	12.8381
0.5543	0.8312	1.14	10.595	12.8381	14.4494
0.7	1.026	1.385	11.8307	14.4494	16.0137
0.9893	0.8721	1.2373	12.8381	16.0137	17.5345
1.3444	1.2390	1.6899	14.4494	17.5345	19.0228
		2.16	16.0137	19.0228	
		2.71	17.5345		
		3.84	19.0228		

Значение, находящееся на пересечении строки «1», и столбца «0.05», ...

... будет 3,8415



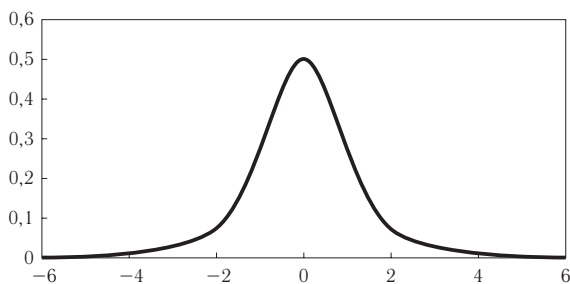
5. Распределение Стьюдента

В статистике часто используется и такая формула распределения вероятностей:

$$f(x) = \frac{\int_0^{\infty} x^{\frac{n+1}{2}-1} e^{-x} dx}{\sqrt{n\pi} \int_0^{\infty} x^{\frac{n}{2}-1} e^{-x} dx} \times \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}$$

где n — число степеней свободы. Если плотность распределения вероятностей можно выразить с помощью этой формулы, в статистике это означает, что величина x имеет распределение Стьюдента с числом степеней свободы n .

Распределение Стьюдента с числом степеней свободы 5



6. Распределение Фишера, или F-распределение

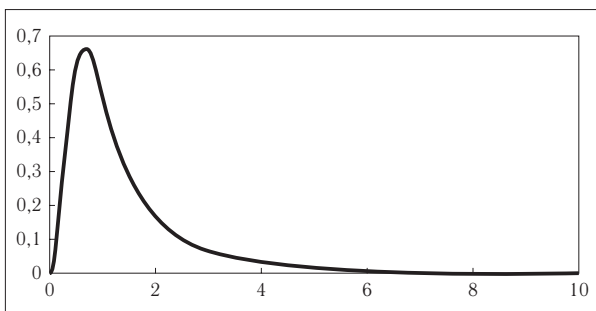
Не менее часто используется и такая функция распределения плотности вероятности:

$$f(x) = \begin{cases} \frac{\left(\int_0^{\infty} x^{\frac{n+m}{2}-1} e^{-x} dx\right) \times n^{\frac{n}{2}} \times m^{\frac{m}{2}}}{\left(\int_0^{\infty} x^{\frac{n}{2}-1} e^{-x} dx\right) \times \left(\int_0^{\infty} x^{\frac{m}{2}-1} e^{-x} dx\right)} \times \frac{x^{\frac{n}{2}-1}}{(nx+m)^{\frac{n+m}{2}}} & \text{при } x > 0, \\ 0 & \text{при } x \leq 0, \end{cases}$$

где n и m — число степеней свободы величины x .

Если формула плотности распределения вероятности имеет такой вид, в статистике это означает, что x имеет F-распределение с числом степеней свободы n и m .

Случай, когда степень свободы $n = 10$, а $m = 5$



7. Распределения и Excel

До начала 90-х вычисление вероятности и значения x было настолько сложным и трудоёмким, что его можно было выполнить только с помощью таблиц стандартного нормального распределения и распределения хи-квадрат. Однако, по мере развития компьютерных технологий, необходимость в этих таблицах отпала, поскольку величины, которые указаны в таблицах, можно получить с помощью программы Excel.

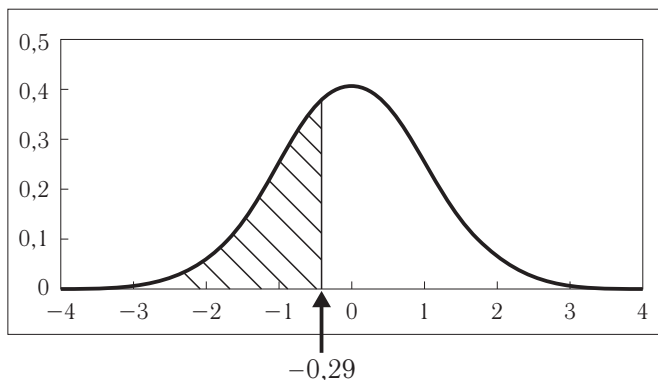
Таблица 5.1.

Распределение	Функция		Значение функции
	русский	английский	
Нормальное распределение*	НОРМРАСП	NORMDIST	Вероятность для заданного x
	НОРМОБР	NORMINV	Значение x для заданной вероятности
Стандартное нормальное распределение	НОРМСТРАСП	NORMSDIST	Вероятность для заданного x
	НОРМСТОБР	NORMSINV	Значение x для заданной вероятности
Распределение хи-квадрат	ХИ2РАСП	CHIDIST	Вероятность для заданного x
	ХИ2ОБР	CHIINV	Значение x для заданной вероятности
Распределение Стьюдента	СТЮДРАСП	TDIST	Вероятность для заданного x
	СТЮДРАСПОБР	TINV	Значение x для заданной вероятности
Распределение Фишера	ФРАСП	FDIST	Вероятность для заданного x
	ФРАСПОБР	FINV	Значение x для заданной вероятности

* Функция распределения вероятностей (речь идёт о нормальном распределении) зависит от таких параметров, как среднее значение и стандартное отклонение. Поэтому таблицу нормального распределения, даже при большом желании, невозможно составить. Но с помощью Excel можно сформировать таблицу, соответствующую таблице нормального распределения, что очень удобно.

Упражнение

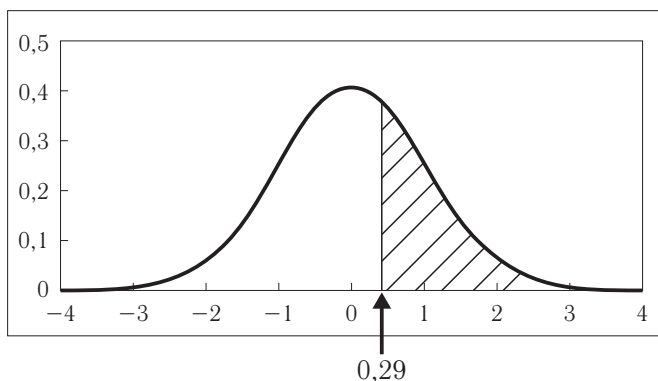
1. Используя таблицу стандартного нормального распределения на стр. 93, вычислите вероятность для заштрихованной на графике области.



2. Используя таблицу распределения хи-квадрат на стр. 103, вычислите значение χ^2 , если число степеней свободы равно 2, а $P = 0,05$.

Ответ

1. Искомая вероятность равна вероятности заштрихованной области.



Вероятность при $z = 0,29 = 0,2 + 0,09$, как следует из таблицы стандартного нормального распределения, равна $0,1141$. Следовательно, искомая вероятность будет равна: $0,5 - 0,1141 = 0,3859$.

2. Как следует из таблицы распределения хи-квадрат, значение величины χ^2 , которое нужно найти, равно $5,9915$.

Выводы

- К широко используемым видам функции распределения вероятностей можно отнести:
 - нормальное распределение;
 - стандартное нормальное распределение;
 - распределение хи-квадрат;
 - распределение Стьюдента;
 - распределение Фишера.
- Площадь области, ограниченной осью x и кривой распределения вероятностей, равна 1.
- Площадь области, ограниченной осью x и кривой распределения вероятностей, тождественна и доле, и вероятности.
- С помощью таблиц распределений или программы Excel можно вычислить:
 - вероятность при определенном значении x ;
 - значение x при определенном значении вероятности.

Глава 6

Что может связывать две переменные

(кто-то идет
шаркая ногами)

На самом деле ...

... хорошо
иногда
провести урок
на улице

Что это у Вас
за вид?
Ботинки
разные!!

А-а-а,
действи-
тельно.

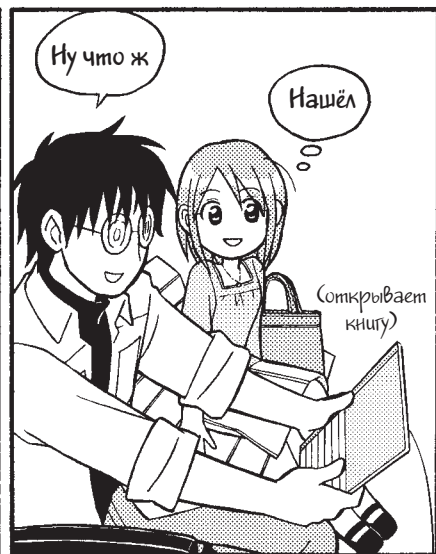
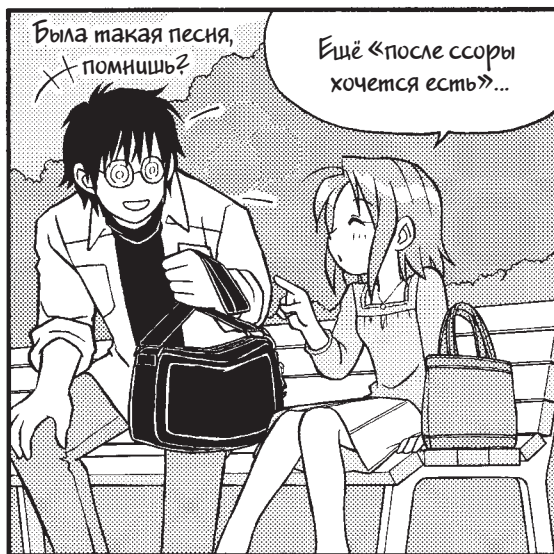
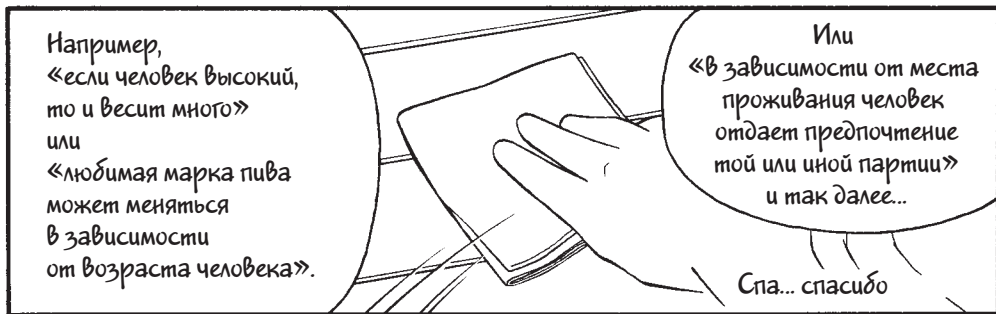
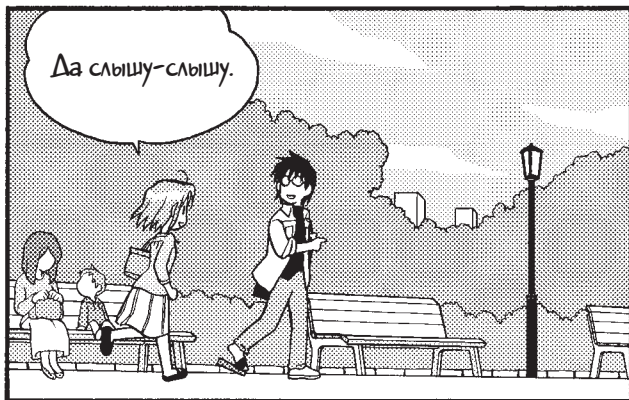
Ах, если бы
это был
Игараси-сан,
было бы так
здорово...

Игараси-сан...

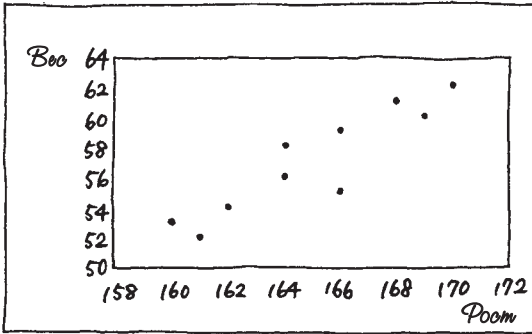
Игараси-сан...
Э-э-э?

Чёрт...
Я уже стала
забывать, как
он выглядит...

Сегодня
я тебе расскажу
про связь между
двумя переменными

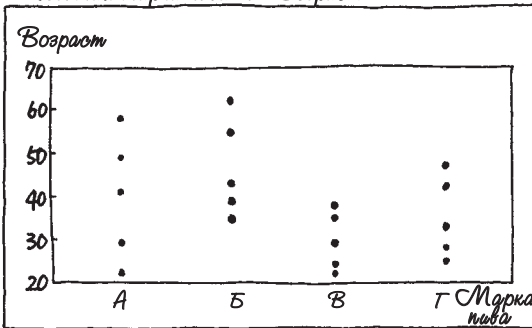


Точечная диаграмма «Рост» и «Вес»



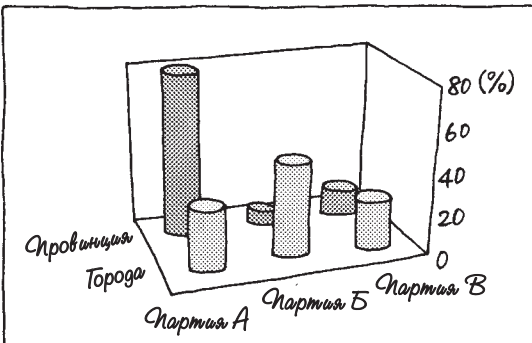
← Количественные данные
и количественные данные

Точечная диаграмма
«Любимая марка пива» и «Возраст»



← Количественные данные
и качественные данные

Столбчатая диаграмма
«Место проживания» и «Поддерживаемая партия»

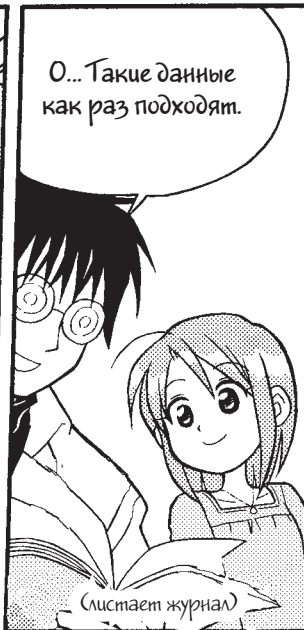
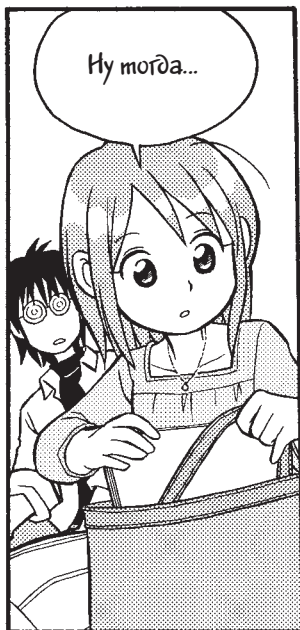
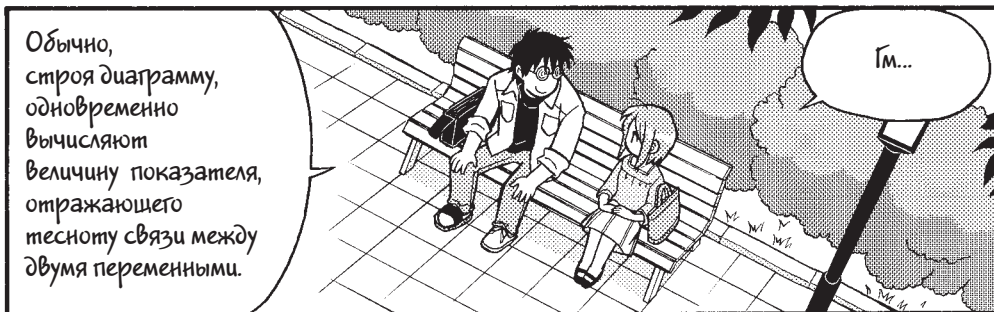


← Качественные
данные
и качественные
данные

Когда
построишь диаграмму,
сразу ясно, есть или нет
взаимосвязь между
двумя переменными.

Ага





1. Коэффициент линейной корреляции

Ой, есть анкета:
«расходы на косметику»
и «расходы на одежду».

Количественные
и количественные
данные



Расходы на косметику и
Расходы на одежду (в месяц)

Ответы 10 женщин в возрасте от 20 до 29 лет

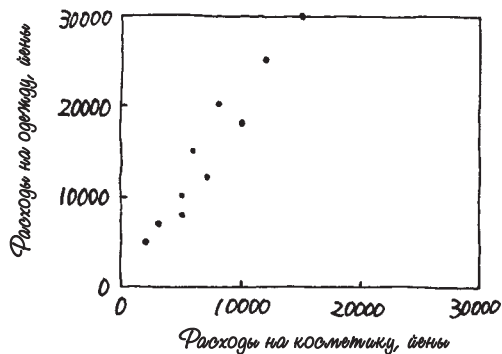
Респондент	Расходы на косметику, йены	Расходы на одежду, йены
А	3000	7000
Б	5000	8000
В	12000	25000
Г	2000	5000
Д	7000	12000
Е	15000	30000
Ж	5000	10000
З	6000	15000
И	8000	20000
К	10000	18000

Для начала
построим
диаграмму.

Ага.

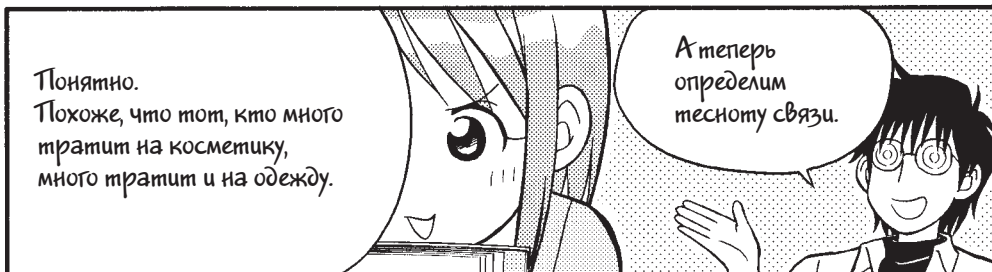


Точечная диаграмма



Понятно.
Похоже, что тот, кто много
тратит на косметику,
много тратит и на одежду.

А теперь
определим
тесноту связи.



Тип данных	Показатель	Значение	Формула
Количеств. и количеств.	Коэффициент линейной корреляции	-1 ... 1	$\frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$
Количеств. и качеств.	Корреляционное отношение	0 ... 1	$\frac{\text{Межгрупповая дисперсия}}{\text{Внутригрупповая дисперсия} + \text{Межгрупповая дисперсия}}$
Качеств. и качеств.	Коэффициент корреляции Крамера	0 ... 1	$\sqrt{\frac{\chi_0^2}{n(\min\{\text{кол-во строк, кол-во столбцов}\} - 1)}}$

*«Корреляционное отношение» см. на стр. 121, «Коэффициент корреляции Крамера» см. стр. 127.



В зависимости от того, какие у нас данные, различается и показатель.



Вот как...



Расходы на косметику и расходы на одежду линейно зависимы.

$$\frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$$

количественные данные и количественные данные, поэтому ... не хочется

... давай, не будем торопиться и вычислим всё не спеша.

Поехали!

Спасибо!

Порядок вычисления коэффициента линейной корреляции (для определения тесноты связи между «расходами в месяц на косметику» и «расходами в месяц на одежду»)

Респондент	Расходы на косметику, йены	Расходы на одежду, йены					
	x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(y - \bar{y})^2$	$(x - \bar{x})(y - \bar{y})$
А	3000	7000	-4300	-8000	18490000	64000000	34400000
Б	5000	8000	-2300	-7000	5290000	49000000	16100000
В	12000	25000	4700	10000	22090000	100000000	47000000
Г	2000	5000	-5300	-10000	28090000	100000000	53000000
Д	7000	12000	-300	-3000	90000	9000000	900000
Е	15000	30000	7700	15000	59290000	225000000	115500000
Ж	5000	10000	-2300	-5000	5290000	25000000	11500000
З	6000	15000	-1300	0	1690000	0	0
И	8000	20000	700	5000	490000	25000000	3500000
К	10000	18000	2700	3000	7290000	9000000	8100000
Сумма	73000	150000	0	0	148100000	606000000	290000000
Ср. знач.	7300	15000			\downarrow S_{xx}	\downarrow S_{yy}	\downarrow S_{xy}

Итак,
подставим числа.

$$\frac{S_{xy}}{\sqrt{S_{xx} \times S_{yy}}} = \frac{290000000}{\sqrt{148100000 \times 606000000}} = 0.9680$$



Если воспользоваться компьютером, всё будет намного быстрее.

Коэффициент линейной корреляции = 0,9680.



Когда между двумя переменными существует тесная связь, коэффициент корреляции приближается к ± 1 , а когда связь слабая, он приближается к 0.



Гм...



Этот результат довольно близок к 1, т.е. между расходами на косметику и расходами на одежду существует тесная связь, так?



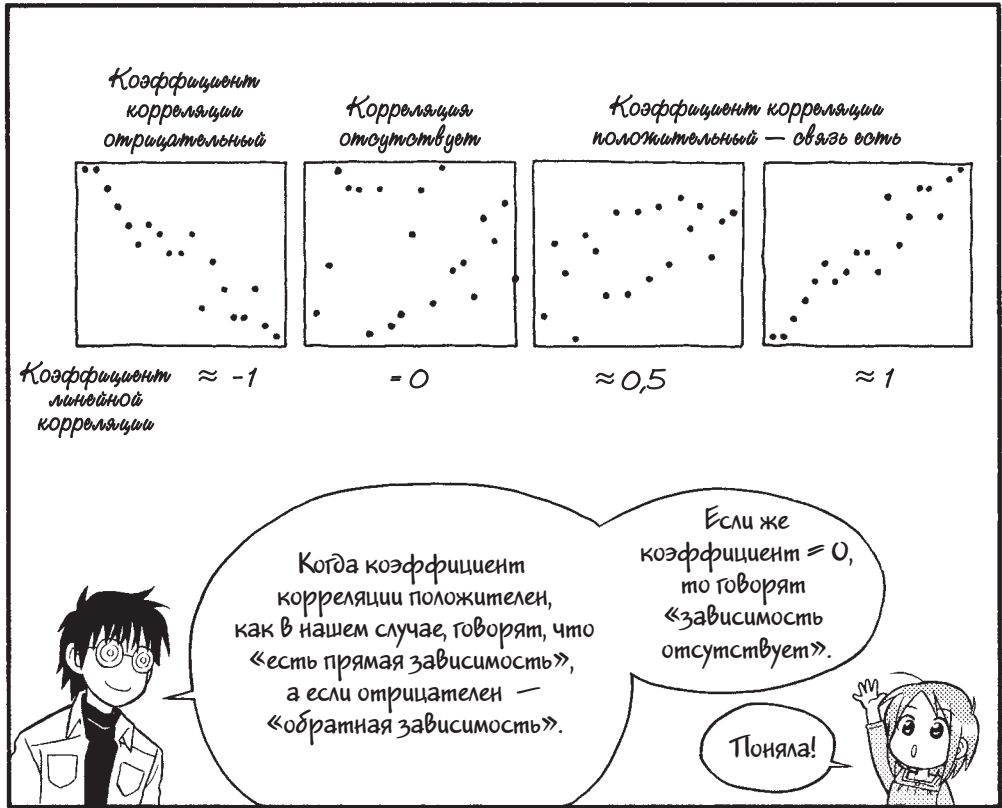
Примерно так.

А когда коэффициент корреляции приближается к -1 ?



В том случае, когда увеличение расходов на косметику ведёт к уменьшению расходов на одежду.





Критерии величины коэффициента линейной корреляции

Абсолютная величина коэфф-та лин. корреляции	Вывод о степени взаимосвязи	Вывод о наличии взаимосвязи
1,0—0,9	⇒ Очень тесная	Есть
0,9—0,7	⇒ Достаточно тесная	
0,7—0,5	⇒ Слабая	
< 0,5	⇒ Очень слабая	Нет



Это для справки.

Вот как!

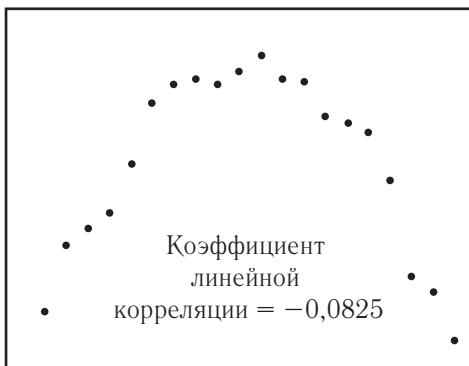


Примечание

До этого я говорил, что коэффициент линейной корреляции — это показатель тесноты связи между количественными данными. Строго говоря, это не совсем так. Коэффициент корреляции еще показывает, является эта зависимость линейной или нет.



Пример нелинейной зависимости



Например, эта диаграмма демонстрирует очевидную связь между двумя переменными. Но поскольку зависимость нелинейная, коэффициент линейной корреляции практически равен нулю.

2. Коэффициент корреляции между данными разных типов

Теперь, перейдём к следующей теме! Здесь есть данные о «возрасте» и «любимом бренде одежды».

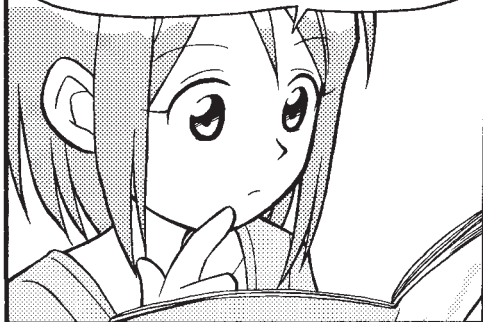
Это количественные и качественные данные.



Возраст, и Любимый бренд одежды

Респондент	Возраст	Бренд
А	27	Benetton
Б	33	Zara
В	16	O'STIN
Г	29	Бапари
Д	32	Zara
Е	23	Benetton
Ж	25	Zara
З	28	Benetton
И	22	O'STIN
К	18	O'STIN
Л	26	Zara
М	26	Benetton
Н	15	O'STIN
О	29	Zara
П	26	O'STIN

Для случая, когда одни данные количественные, а другие качественные, вычисляется коорреляционное отношение, которое может принимать значения от 0 до 1.



А в этом случае приближение значения к 1 тоже свидетельствует о наличии тесной связи между данными?



«Любимый бренд одежды» и «Возраст»

Попробуем-ка упорядочить эту таблицу.

Да-да, попробуем...

	Benetton	Zana	O'STICH	
	23	25	15	
	26	26	16	
	27	29	18	
	28	32	22	
		33	26	
			29	
Сумма	104	145	126	375
Среднее значение	26	29	21	25

Диаграмма
«Любимый бренд одежды» и «Возраст»

Следующим этапом будет построение диаграммы.

Ух ты! Похоже, что есть взаимосвязь!

Ну тогда вычислим корреляционное отношение.

Согласна!

Вычисление корреляционного отношения предусматривает выполнение нижеследующих шагов.



Шаг 1

Вычислить суммы стандартных отклонений для каждого столбца таблицы:

	(значение — среднее значение в столбце Benetton) ²	(значение — среднее значение в столбце Zara) ²	(значение — среднее значение в столбце O'STIN) ²
	$(23 - 26)^2 = (-3)^2 = 9$	$(25 - 29)^2 = (-4)^2 = 16$	$(15 - 21)^2 = (-6)^2 = 36$
	$(26 - 26)^2 = 0^2 = 0$	$(26 - 29)^2 = (-3)^2 = 9$	$(16 - 21)^2 = (-5)^2 = 25$
	$(27 - 26)^2 = 1^2 = 1$	$(29 - 29)^2 = 0^2 = 0$	$(18 - 21)^2 = (-3)^2 = 9$
	$(28 - 26)^2 = 2^2 = 4$	$(32 - 29)^2 = 3^2 = 9$	$(22 - 21)^2 = 1^2 = 1$
		$(33 - 29)^2 = 4^2 = 16$	$(26 - 21)^2 = 5^2 = 25$
			$(29 - 21)^2 = 8^2 = 64$
Сумма	14	50	160
	↓ S_{BB}	↓ S_{ZZ}	↓ S_{OO}

Шаг 2

Вычислить внутригрупповую дисперсию, показывающую как сильно отличаются данные в каждой категории, как сумму $S_{TT} + S_{CC} + S_{BB}$:

$$S_{BB} + S_{ZZ} + S_{OO} = 14 + 50 + 160 = 224$$

Шаг 3

Рассчитать межгрупповую дисперсию — меру разброса данных между категориями. Для этого надо найти стандартные отклонения средних значений от общего среднего значения, умножить их на количество данных в соответствующей графе и вычислить сумму полученных произведений, т. е.:

$$\begin{aligned} & \text{Кол-во значений Benetton} \times (\text{Ср. знач. Benetton} - \text{общее ср. знач.})^2 + \\ & + \text{Кол-во значений Zara} \times (\text{Ср. знач. Zara} - \text{общее ср. знач.})^2 + \\ & + \text{Кол-во значений O'STIN} \times (\text{Ср. знач. O'STIN} - \text{общее ср. знач.})^2 \end{aligned}$$

$$\begin{aligned} & 4 \times (26 - 25)^2 + 5 \times (29 - 25)^2 + 6 \times (21 - 25)^2 = \\ & = 4 \times 1 + 5 \times 16 + 6 \times 16 = \\ & = 4 + 80 + 96 = \\ & = 180 \end{aligned}$$

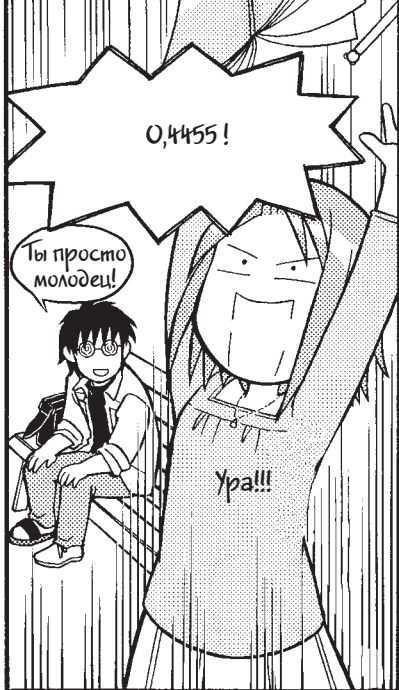
Шаг 4

Вычислить корреляционное отношение по формуле:

$$\frac{180}{224 + 180} = \frac{180}{404} = 0,4455$$

Значение корреляционного отношения между «возрастом» и «любимым брендом одежды» равно ...

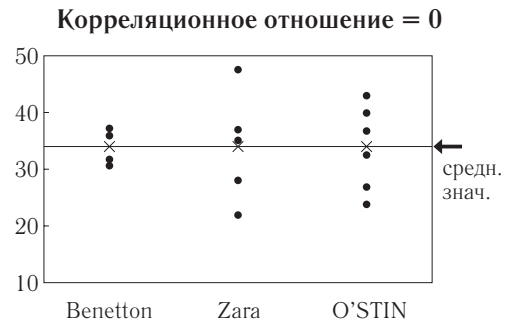




Как я уже говорил, корреляционное отношение может принимать значения от 0 до 1. Чем теснее взаимосвязь между двумя переменными, тем значение ближе к 1, чем слабее взаимосвязь, тем ближе к 0.



Диаграмма «Любимый бренд одежды» и «Возраст»



Корреляционное отношение = 1

⇕

Данные внутри каждой группы одинаковы

⇕

Внутригрупповая дисперсия = 0

Корреляционное отношение = 0

⇕

Средние значения всех групп одинаковы

⇕

Межгрупповая дисперсия = 0



К сожалению, в статистике нет какого-то определённого значения корреляционного отношения, выше которого переменные считаются тесно связанными. Для справки приведём таблицу соответствия корреляционного отношения и степени взаимосвязи.

Критерии величины корреляционного отношения

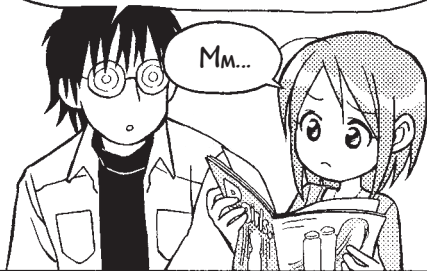
Корреляционное отношение	⇒	Вывод о степени взаимосвязи	Вывод о наличии взаимосвязи
1,0—0,8	⇒	Очень тесная	Есть
0,8—0,5	⇒	Достаточно тесная	
0,5—0,25	⇒	Слабая	Нет
< 0,25	⇒	Очень слабая	

В нашем примере корреляционное отношение было равно 0,4455: значит связь слабая!



3. Коэффициент корреляции Крамера

Дальше...
Было бы здорово привести пример, который мог бы пояснить связь между качественными данными.



Мм...

Ой! Как насчёт вот этих данных?

Опросили 300 школьников!
Какой способ признания в любви вы предпочитаете?

Опросили 300 школьников на тему «какой способ признания в любви вы предпочитаете?»

Ну-ка, ну-ка...
Вот как... Три способа
признания в любви:
«по телефону»,
«по SMS»
и «при встрече».
Эти данные можно
использовать.



Да!

Однако женский
журнал проводит
странное анкетирование ...



Не
твоего
ума
дело!

Таблица взаимной сопряженности «пола» и «способа признания в любви»

Пол респондента	Способ признания в любви, люди			Итого
	по телефону	по SMS	при встрече	
женский	34	61	53	148
мужской	38	40	74	152
Итого:	72	101	127	300

Это значит, что из 152 опрошенных молодых людей 74 хотели бы, чтобы им признавались в любви при встрече.

Таблица взаимной сопряженности в %

Пол респондента	Способ признания в любви, %			Итого
	по телефону	по SMS	при встрече	
женский	23	41	36	100
мужской	25	26	49	100
Итого:	24	34	42	100

Это значит, что 49% опрошенных молодых людей хотели бы, чтобы им признавались в любви при встрече.

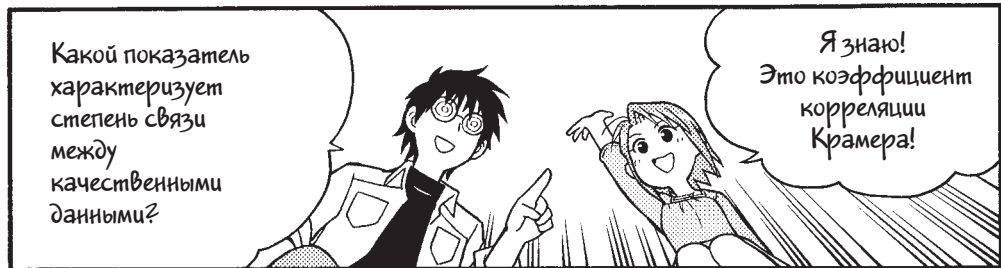
Такая
таблица,
где две
переменные
перекре-
щаются,
называется
таблицей
взаимной
сопряжен-
ности.



Вот как...
Девушки предпочитают,
чтобы им признавались
в любви «по SMS», ...



... а юноши
предпочитают, чтобы им
признавались в любви
«при встрече».



Вычисление коэффициента корреляции Крамера предусматривает выполнение следующих шагов:



Шаг 1

Построить таблицу взаимной сопряжённости. Величины в клетках, обведённых жирной чертой, называются **эмпирическими частотами**.

Пол респондента	Способ признания в любви			Итого
	по телефону	по SMS	при встрече	
женский	34	61	53	148
мужской	38	40	74	152
Итого:	72	101	127	300

Шаг 2

Выполнить вычисления, как показано в следующей таблице. Величины в обведённых клетках называются **теоретическими частотами**.

Пол респондента	Способ признания в любви			Итого
	по телефону	по SMS	при встрече	
женский	$\frac{148 \times 72}{300}$	$\frac{148 \times 101}{300}$	$\frac{148 \times 127}{300}$	148
мужской	$\frac{152 \times 72}{300}$	$\frac{152 \times 101}{300}$	$\frac{152 \times 127}{300}$	152
Итого:	72	101	127	300

$$\frac{\text{Число опрошенных юношей} \times \text{Число ответов «при встрече»}}{\text{Общее число опрошенных}}$$

Если между «полом» и «способом признания в любви» полностью отсутствует какая-либо связь, то отношение эмпирических частот в строках «женский» и «мужской» будет таким-же, как отношение значений в строке «Итого» (см. Шаг 1).

Теоретические частоты, вычисленные на Шаге 2, отражают, каково было бы число соответствующих респондентов в случае полного отсутствия какой-либо связи между «полом» и «способом признания в любви».



Шаг 3

Вычислить:

$$\frac{(\text{Эмпирическая частота} - \text{Теоретическая частота})^2}{\text{Теоретическая частота}}$$

Пол респондента	Способ признания в любви			Итого
	по телефону	по SMS	при встрече	
женский	$\frac{\left(34 - \frac{148 \times 72}{300}\right)^2}{\frac{148 \times 72}{300}}$	$\frac{\left(61 - \frac{148 \times 101}{300}\right)^2}{\frac{148 \times 101}{300}}$	$\frac{\left(53 - \frac{148 \times 127}{300}\right)^2}{\frac{148 \times 127}{300}}$	148
мужской	$\frac{\left(38 - \frac{152 \times 72}{300}\right)^2}{\frac{152 \times 72}{300}}$	$\frac{\left(40 - \frac{152 \times 101}{300}\right)^2}{\frac{152 \times 101}{300}}$	$\frac{\left(74 - \frac{152 \times 127}{300}\right)^2}{\frac{152 \times 127}{300}}$	152
Итого:	72	101	127	300



Чем больше разница между эмпирическими и теоретическими частотами, т.е. чем теснее связь между «полом» и «способом признания в любви», тем больше значения величин в клетках таблицы.

Шаг 4

Вычислить сумму величин в клетках таблицы на **Шаге 3**.

Иначе говоря, вычислить значение критерия согласия Пирсона (χ_0^2).

$$\begin{aligned}\chi_0^2 &= \frac{\left(34 - \frac{148 \times 72}{300}\right)^2}{\frac{148 \times 72}{300}} + \frac{\left(61 - \frac{148 \times 101}{300}\right)^2}{\frac{148 \times 101}{300}} + \frac{\left(53 - \frac{148 \times 127}{300}\right)^2}{\frac{148 \times 127}{300}} + \\ &+ \frac{\left(38 - \frac{152 \times 72}{300}\right)^2}{\frac{152 \times 72}{300}} + \frac{\left(40 - \frac{152 \times 101}{300}\right)^2}{\frac{152 \times 101}{300}} + \frac{\left(74 - \frac{152 \times 127}{300}\right)^2}{\frac{152 \times 127}{300}} = \\ &= 8,0091\end{aligned}$$

Согласно Шагу 3, чем больше разница между эмпирическими и теоретическими частотами, т.е. чем теснее связь между «полом» и «способом признания», тем больше значение критерия согласия Пирсона (χ_0^2).



Шаг 5

Величина коэффициента корреляции Крамера вычисляется по формуле:

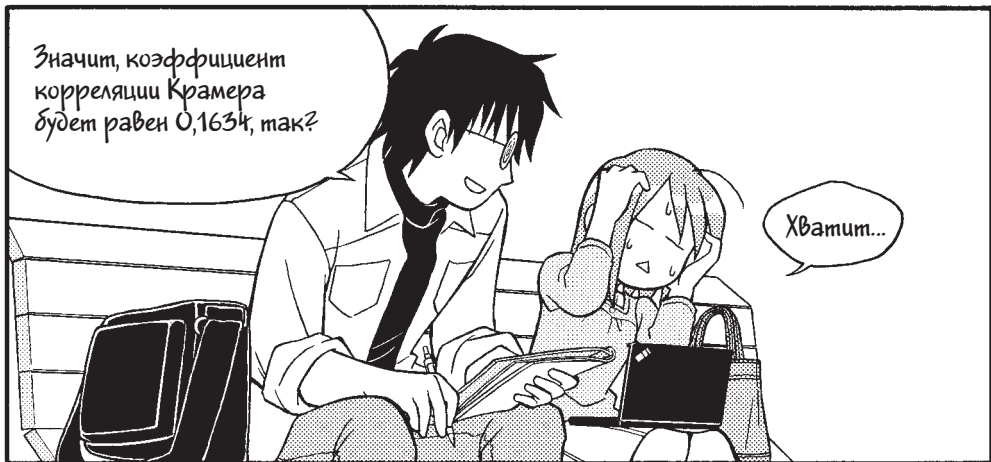
$$\sqrt{\frac{\chi_0^2}{n \times (\min \{ \text{кол-во строк в таблице}; \text{кол-во столбцов в таблице} \} - 1)}}$$

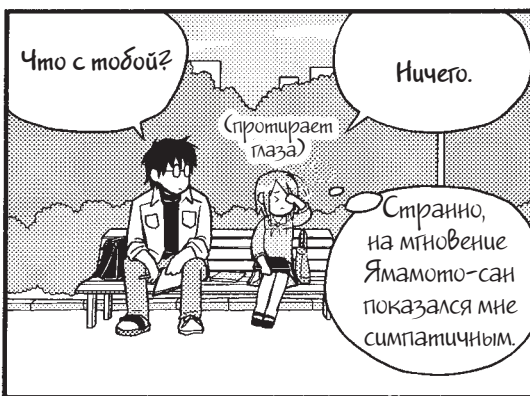
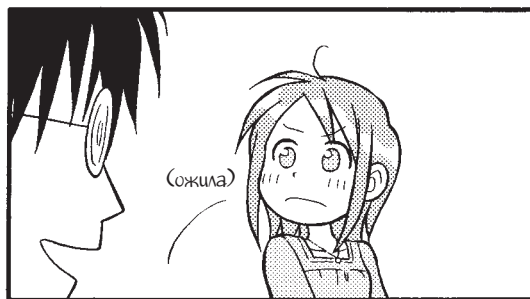
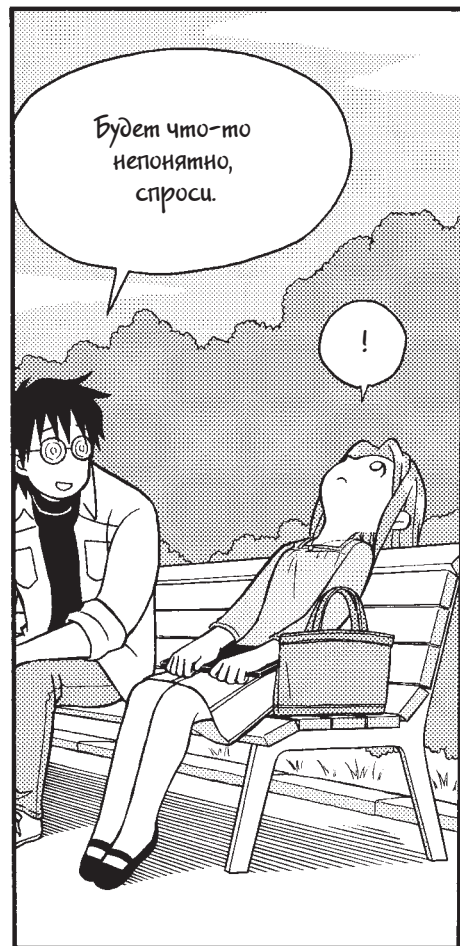
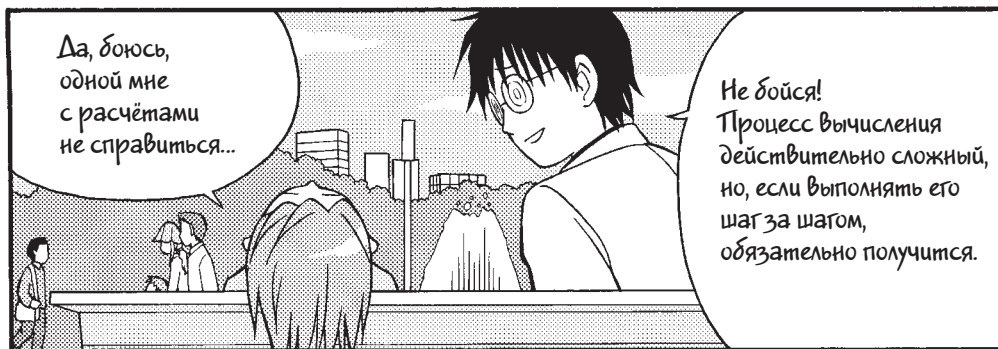
где

n — общее число единиц совокупности,

$\min \{a; b\}$ — из величин a и b надо взять меньшую.

$$\sqrt{\frac{8,0091}{300 \times (\min \{2, 3\} - 1)}} = \sqrt{\frac{8,0091}{300 \times (2 - 1)}} = \sqrt{\frac{8,0091}{300}} = 0,1634$$





Как я уже говорил, коэффициент корреляции Крамера может принимать значение от 0 до 1. Чем теснее связаны две переменные, тем ближе коэффициент Крамера к 1, а чем слабее связь, тем ближе к 0.



Таблицы взаимной сопряжённости «Пола» и «Способа признания в любви»

Величина коэффициента Крамера равна 1

Пол респондента	Способ признания в любви, %			Итого
	по телефону	по SMS	при встрече	
женский	17	83	0	100
мужской	0	0	100	100

Величина коэффициента Крамера равна 1



Предпочтения девушек и юношей совершенно различаются.

Величина коэффициента Крамера равна 0

Пол респондента	Способ признания в любви, %			Итого
	по телефону	по SMS	при встрече	
женский	17	48	35	100
мужской	17	48	35	100

Величина коэффициента Крамера равна 0



Предпочтения девушек и юношей одинаковы

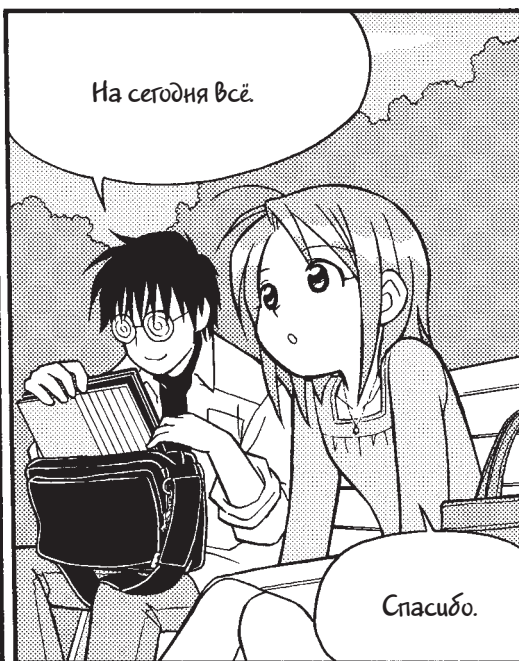


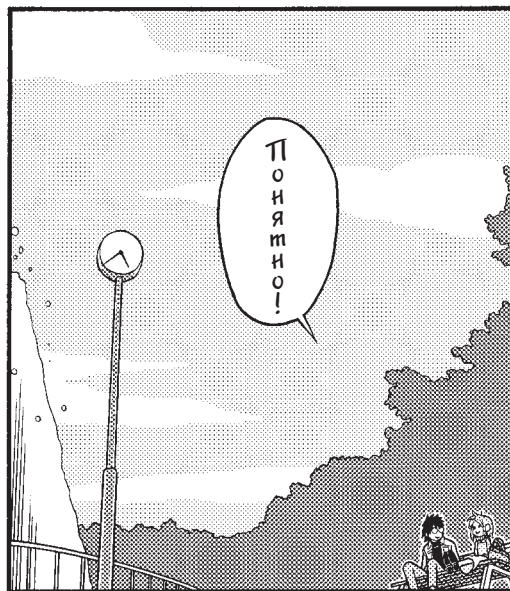
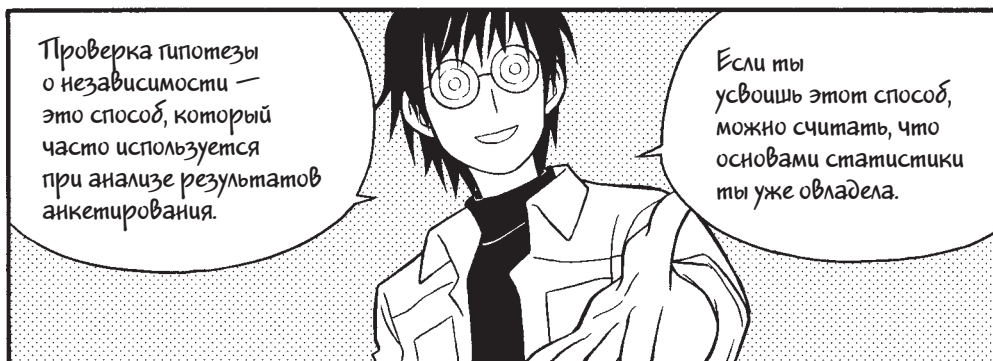
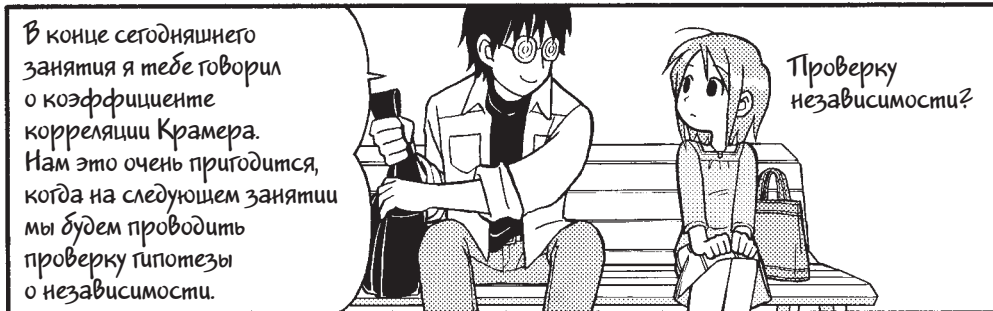
К сожалению, в статистике нет кого-то определённого значения коэффициента корреляции Крамера, выше которого переменные считаются тесно связанными.

Для справки приведём таблицу соответствия коэффициента Крамера и степени взаимосвязи.

Критерии величины коэффициента корреляции Крамера

Коэффициент Крамера		Вывод о степени взаимосвязи	Вывод о наличии взаимосвязи
1,0—0,8	⇒	Очень тесная	Есть
0,8—0,5	⇒	Достаточно тесная	
0,5—0,25	⇒	Слабая	Нет
< 0,25	⇒	Очень слабая	





Упражнение

В последнее время положение дел в компании X, владеющей сетью семейных ресторанов, трудно назвать благополучным. Чтобы узнать пожелания клиентов, компания провела анкетирование, объектами которого стали произвольно выбранные жители Японии в возрасте старше 20 лет.

Результаты анкетирования представлены в виде следующей таблицы:

Респондент	Какую кухню Вы обычно заказываете в семейном ресторане?	Если после обеда будет предлагаться бесплатный напиток, Вы предпочтёте чай или кофе?
1	китайская	кофе
2	европейская	кофе
...
250	японская	чай

На основании этой таблицы был составлен один из вариантов таблицы взаимной сопряжённости:

Обычно заказываемая кухня	Предпочитаемый напиток		Итого
	Кофе	Чай	
Японская	43	33	76
Европейская	51	53	104
Китайская	29	41	70
Итого:	123	127	250

Вычислите коэффициент корреляции Крамера, чтобы установить взаимосвязь между «обычно заказываемой кухней» и «предпочитаемым напитком».

Решение**Шаг 1**

Составим таблицу взаимной сопряжённости.

Обычно заказываемая кухня	Предпочитаемый напиток		Итого
	Кофе	Чай	
Японская	43	33	76
Европейская	51	53	104
Китайская	29	41	70
Итого:	123	127	250

Шаг 2

Вычислим теоретические частоты.

Обычно заказываемая кухня	Предпочитаемый напиток		Итого
	Кофе	Чай	
Японская	$\frac{76 \times 123}{250}$	$\frac{76 \times 127}{250}$	76
Европейская	$\frac{104 \times 123}{250}$	$\frac{104 \times 127}{250}$	104
Китайская	$\frac{70 \times 123}{250}$	$\frac{70 \times 127}{250}$	70
Итого:	123	127	250

Шаг 3

Вычислим:

$$\frac{\text{Эмпирическая частота} - \text{Теоретическая частота}}{\text{Теоретическая частота}}$$

Обычно заказываемая кухня	Предпочитаемый напиток		Итого
	Кофе	Чай	
Японская	$\left(43 - \frac{76 \times 123}{250}\right)^2$ $\frac{76 \times 123}{250}$	$\left(33 - \frac{76 \times 127}{250}\right)^2$ $\frac{76 \times 127}{250}$	76
Европейская	$\left(51 - \frac{104 \times 123}{250}\right)^2$ $\frac{104 \times 123}{250}$	$\left(53 - \frac{104 \times 127}{250}\right)^2$ $\frac{104 \times 127}{250}$	104
Китайская	$\left(29 - \frac{70 \times 123}{250}\right)^2$ $\frac{70 \times 123}{250}$	$\left(41 - \frac{70 \times 127}{250}\right)^2$ $\frac{70 \times 127}{250}$	70
Итого:	123	127	250

Шаг 4

Вычислим сумму величин в обведённых жирной чертой клетках таблицы на **Шаге 3**.

Иначе говоря,

вычислим значение критерия согласия Пирсона.

$$\begin{aligned}\chi_0^2 &= \frac{\left(43 - \frac{76 \times 123}{250}\right)^2}{\frac{76 \times 123}{250}} + \frac{\left(33 - \frac{76 \times 127}{250}\right)^2}{\frac{76 \times 127}{250}} \\ &+ \frac{\left(51 - \frac{104 \times 123}{250}\right)^2}{\frac{104 \times 123}{250}} + \frac{\left(53 - \frac{104 \times 127}{250}\right)^2}{\frac{104 \times 127}{250}} \\ &+ \frac{\left(29 - \frac{70 \times 123}{250}\right)^2}{\frac{70 \times 123}{250}} + \frac{\left(41 - \frac{70 \times 127}{250}\right)^2}{\frac{70 \times 127}{250}} \\ &= 3,3483\end{aligned}$$

Шаг 5

Вычисляем коэффициент корреляции Крамера по формуле:

$$\begin{aligned}&\sqrt{\frac{\chi_0^2}{n \times (\min\{\text{кол-во строк в таблице; кол-во столбцов в таблице}\} - 1)}} = \\ &= \sqrt{\frac{3,3483}{250 \times (\min\{3, 2\} - 1)}} = \sqrt{\frac{3,3483}{250 \times (2 - 1)}} = \sqrt{\frac{3,3483}{250}} = 0,1157\end{aligned}$$

Выводы

- Коэффициент линейной корреляции является показателем тесноты связи между двумя количественными данными.
- Корреляционное отношение является показателем тесноты связи между количественными и качественными данными.
- Коэффициент корреляции Крамера является показателем тесноты связи между качественными данными.
- Перечисленные коэффициенты корреляции имеют следующие свойства:

Вид функции распределения вероятностей	Значение коэффициента корреляции		Отсутствие какой-либо связи между двумя переменными величинами	Переменные величины максимально тесно связаны
	max	min		
Коэффициент линейной корреляции	-1	1	0	-1 или 1
Корреляционное отношение	0	1	0	1
Коэффициент корреляции Крамера	0	1	0	1

- В статистике не существует какого-то определённого значения указанных коэффициентов корреляции, выше которого переменные считаются тесно связанными.

Глава 7

**А что это
за проверка
гипотезы
о независимости?**

1. Проверка гипотезы

Итак,
сегодняшнее
занятие..

Эй!
Почему бы не посмотреть
на меня!!

(обижается)

Ох, извини.
Это новая форма?

Да! Это образец,
но Вам я покажу.

Ну как!

(юбка
шелестит)

(крутится)

О-о...
Тебе идет.

Спасибо.

Ну, что
мы будем
сегодня
проходить?

На прошлом уроке мы изучали коэффициент корреляции Крамера, помнишь?

Опросили 300 школьников!

Какой способ признания в любви вы предпочитаете?

А... про признания в любви, да?

Коэффициент корреляции Крамера в том примере был равен 0,1634. Мы сделали вывод, что «связь очень слабая», так?

Да, так и было.

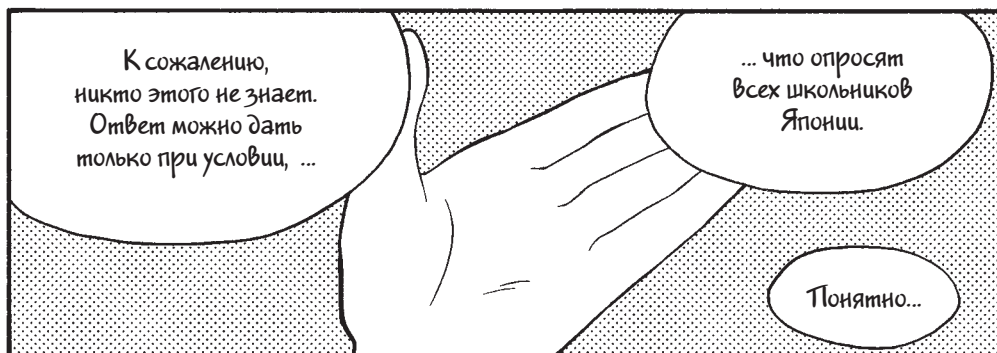
Ну, теперь хорошенько подумай.

Результат этого анкетирования основан всего лишь на данных, ...

... полученных от 300 школьников, произвольно выбранных из всех школьников Японии.

Если бы были выбраны другие школьники, величина коэффициента Крамера была бы, скорее всего, другая.

Если хорошо подумать, наверное вы правы.



... зная, что коэффициент корреляции Крамера, вычисленный на основе ответов, полученных от 300 школьников, произвольно выбранных из генеральной совокупности, равен 0,1634, считаем, что это значение соответствует коэффициенту корреляции Крамера для всей генеральной совокупности.



Нам придётся сделать подобное субъективное заключение.

Довольно туманное...



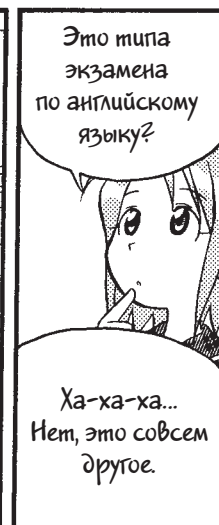
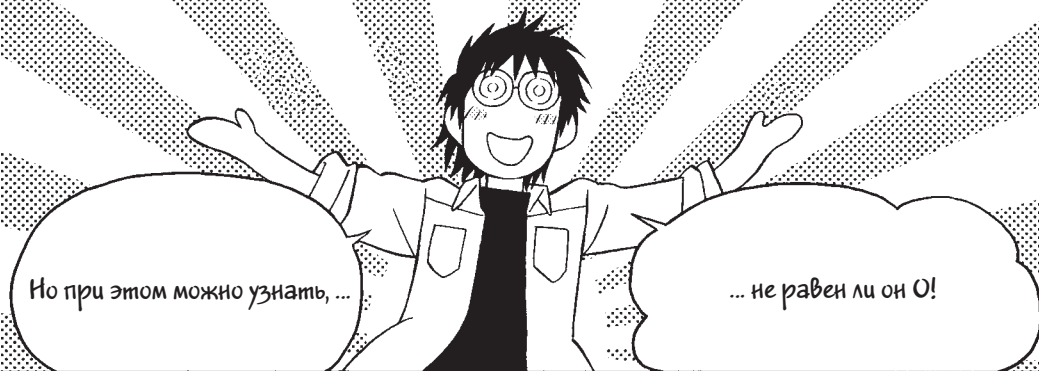
Но ведь наверняка есть какие-то статистические методы?



Нет. К сожалению, даже если очень хорошо владеть статистикой, невозможно узнать точный коэффициент корреляции Крамера для генеральной совокупности.



Что?
Правда?



Проверка гипотезы о независимости — это один из способов анализа, называемых в статистике проверкой.

Проверка

проверка гипотезы о независимости

проверка гипотезы о равенстве средних величин (двух) совокупностей

проверка гипотезы об отсутствии корреляции

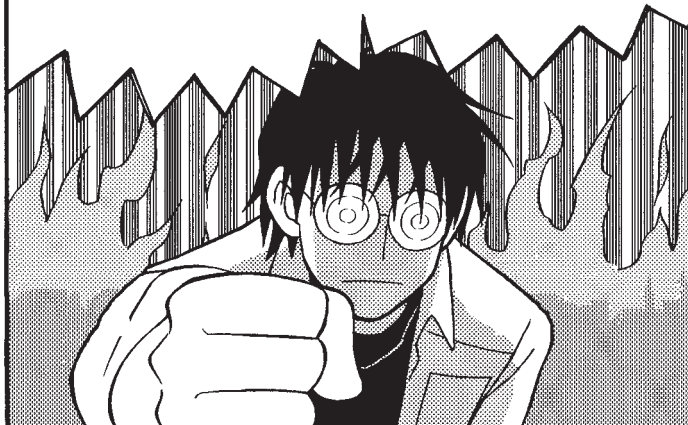
проверка корреляционного отношения

проверка гипотезы о равенстве долей в (двух) совокупностях

В первую очередь я тебе объясню, что такое проверка.

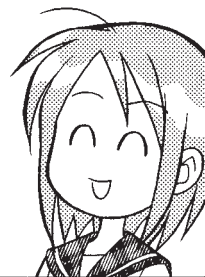
Ладно.

Проверка гипотезы — один из способов анализа, позволяющий на основе данных выборочной совокупности сделать вывод о справедливости гипотезы, которую выдвинул исследователь о генеральной совокупности.



Правильнее было бы сказать — это статистическая проверка гипотез.

А, вот это мне более понятно.



Есть разные виды проверяемых гипотез.

Примеры проверяемых гипотез	
Гипотезы	Примеры применения
О независимости	Проверяют, равняется ли нулю коэффициент корреляции Крамера между полом и способом признания в любви для генеральной совокупности.
О корреляционном отношении	Проверяют, равняется ли нулю корреляционное отношение между любимым брендом одежды и возрастом для генеральной совокупности.
Об отсутствии корреляции	Проверяют, равняется ли нулю коэффициент линейной корреляции между расходами в месяц на косметику и расходами в месяц на одежду для генеральной совокупности.
О равенстве средних величин (двух) совокупностей	Проверяют, получают ли школьницы Токио и школьницы Осаки одну и ту же или разную сумму на карманные расходы. (Будьте внимательны: предполагаются две генеральные совокупности).
О равенстве долей в (двух) совокупностях	Проверяют, различается ли рейтинг кабинета министров среди избирателей, проживающих в городах, и избирателей из сельской местности. (Будьте внимательны: предполагаются две генеральные совокупности).



Порядок статистической проверки гипотез	
Шаг 1.	Определить генеральную совокупность.
Шаг 2.	Сформулировать нулевую и альтернативную гипотезы.
Шаг 3.	Выбрать вид статистической проверки гипотезы.
Шаг 4.	Определить уровень значимости.
Шаг 5.	Вычислить фактическое значение выбранного статистического критерия на основе данных выборочной совокупности.
Шаг 6.	Проверить, входит ли вычисленное на Шаге 5 значение статистического критерия в критическую область.
Шаг 7.	Если фактическое значение выбранного статистического критерия (Шаг 6) входит в критическую область, делают вывод о том, что верна альтернативная гипотеза. В противном случае полагают, что нет оснований считать нулевую гипотезу ошибочной.



2. Проверка гипотезы о независимости

Объясню суть проверки гипотезы о независимости. Сегодня — это главная тема.



Проверка гипотезы о независимости — один из способов анализа, проводимый с целью выяснить, не равен ли нулю коэффициент корреляции Крамера для генеральной совокупности.

Понятно.

Другими словами, это способ анализа, предполагающий наличие взаимосвязи между двумя переменными в таблице взаимной сопряжённости.

Поня-я-тно...

	34	61	53	148
38	40	74	152	
72	101	127		

Проверка гипотезы о независимости называется также проверкой критерия согласия Пирсона χ^2 .

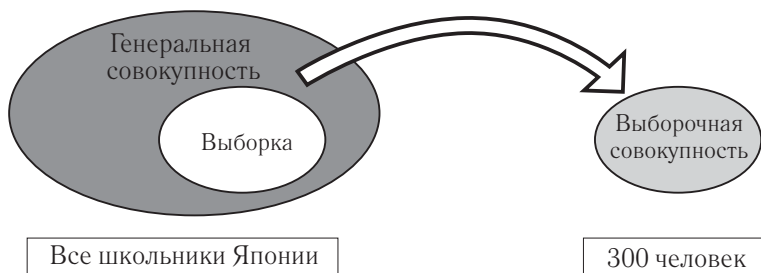
Опять этот χ^2 !
Ужасно...



Перед тем, как объяснить проверку гипотезы о независимости на конкретном примере, расскажу о том, что лежит в основе такой проверки. Предположим, что провели следующее исследование (хотя в действительности это невозможно):

Шаг 1

Из генеральной совокупности, которой являются «школьники Японии», произвольно выбрали 300 человек.



Шаг 2

Среди выбранных 300 школьников провели анкетирование (см. стр. 127), а затем рассчитали величину критерия согласия Пирсона.

Шаг 3

Выбранных 300 школьников «вернули» в генеральную совокупность.

Шаг 3

Шаг 1 — Шаг 3 повторяли множество раз.

Функцией распределения критерия согласия Пирсона, полученной в этом исследовании, является распределение хи-квадрат с числом степеней свободы, равным 2, при условии, что коэффициент корреляции Крамера для генеральной совокупности «все школьники Японии» = 0. Другими словами, если коэффициент корреляции Крамера для генеральной совокупности «все школьники Японии» = 0, то критерий согласия Пирсона (χ_0^2) имеет хи-квадрат-распределение с числом степеней свободы, равным 2.

1. Способ вычисления критерия согласия Пирсона (χ_0^2) см. на стр. 130—133.
2. О распределении хи-квадрат с числом степеней свободы 2 см. на стр. 100.

Описанное выше исследование было проведено на самом деле, но с учётом следующих ограничений:



- Поскольку невозможно провести исследование среди всех школьников, живущих в Японии, считается, что генеральная совокупность «все школьники Японии» — это совокупность, состоящая из 10 000 человек, как указано в **Табл. 7.1**.
- Было сделано предположение, что коэффициент корреляции Крамера для совокупности «все школьники Японии» = 0. Другими словами, предположили, что ответы юношей и девушек на вопрос «Какой способ признания в любви вы предпочитаете?» (по телефону : по SMS : при встрече) одинаковы (см. стр. 135). Предположили также, что на основе **Табл. 7.1** была создана таблица взаимной сопряжённости.
- Исследования (Шаг 1 — Шаг 3) провели 20 тысяч раз.

Таблица 7.1. Какой способ признания в любви вы предпочитаете?
(все школьники Японии)

Респондент	Пол респондента	Способ признания в любви
1	ж	при встрече
2	ж	по телефону
...
10000	м	по SMS

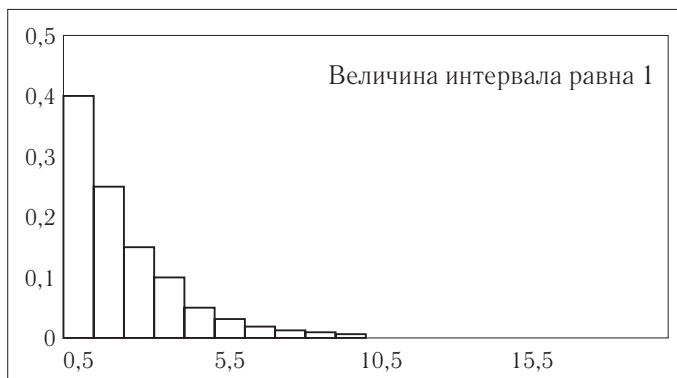
Таблица 7.2. Взаимная сопряжённость пола и способа признания в любви

Пол респондента	Способ признания в любви			Итого
	по телефону	по SMS	при встрече	
женский	400	1600	2000	4000
мужской	600	2400	3000	6000
Итого:	1000	4000	5000	10000

Таблица 7.3. Результаты исследования

Номер исследования	Критерий согласия Пирсона (χ_0^2)
1	0,8598
2	0,7557
...	...
20000	2,7953

Рис. 7.1. Гистограмма на основе данных Табл. 7.3



Действительно, Рис. 7.1 очень похож на график функции распределения вероятностей, когда число степеней свободы равно 2 (см. стр. 100). Нет сомнений, что величина критерия согласия Пирсона (χ_0^2) имеет хи-квадрат-распределение с числом степеней свободы 2. Это не имеет непосредственного отношения к исследованию, но я скажу вам одну очень важную вещь: число степеней свободы, равное 2, получается из выражения:

$$(2 - 1) \times (3 - 1) = 1 \times 2 = 2$$

↑
два варианта:
девушки, юноши

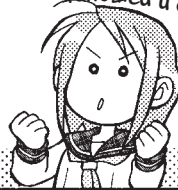
↑
три варианта:
по телефону, по SMS, при встрече



Почему такой странный расчет? Ответ на этот вопрос выходит за рамки данной книги. Поэтому не беспокойтесь, если вы не до конца поняли этот способ расчёта.

Допустим, что коэффициент корреляции Крамера для совокупности «все школьники Японии» = 0. Другими словами, связи между полом и способом признания в любви не существует.

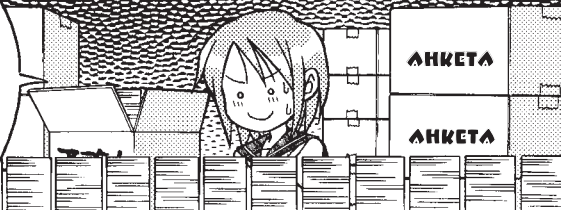
Предположим, что пропорции предпочтений юношей и девушек равны.



Возьмём анкеты 300 человек, выбранных из генеральной совокупности «все школьники Японии»...



... затем следующие 300 человек, затем следующие, ... и так много раз.



Если вычислить величину критерия согласия Пирсона (χ^2), то окажется, ...

Найдём сумму клеток таблицы, вычисленных по формуле:
$$\frac{(\text{эмпирическая частота} - \text{теоретическая частота})^2}{\text{теоретическая частота}}$$



... что величина (χ^2) имеет распределение хи-квадрат с числом степеней свободы, равным 2!

Наконец-то,







Упражнение

Женский журнал «P-girls» решил провести опрос среди школьников по вопросу: «Какой способ признания в любви Вы предпочитаете?». Для этого корреспондент из всех школьников Японии произвольно выбрал 300 человек и провёл анкетирование. Результаты анкетирования представлены в виде следующей таблицы:

Респондент	Способ признания в любви	Возраст респондента	Пол респондента
1	при встрече	17	ж
2	по телефону	15	ж
...
300	по SMS	18	м

Таблица взаимной сопряжённости «пола» и «способа признания в любви»

Пол респондента	Способ признания в любви			Итого
	по телефону	по SMS	при встрече	
женский	34	61	53	148
мужской	38	40	74	152
Итого:	72	101	127	300

Выясните путём проверки гипотезы о независимости, больше ли 0 коэффициент корреляции Крамера между полом и способом признания в любви для генеральной совокупности «все школьники Японии». Другими словами, надо выяснить, есть ли взаимосвязь между полом и способом признания в любви. Предположим, что уровень значимости (объясню позже) равен 0,05.



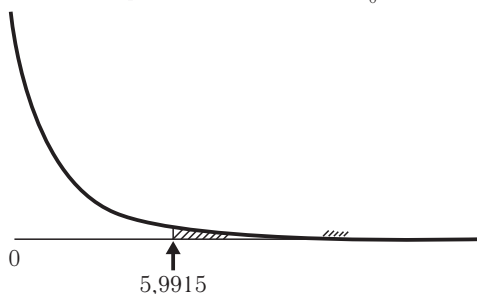


Размышление

Как уже было сказано (стр. 152—154), если коэффициент корреляции Крамера для генеральной совокупности «все школьники Японии» = 0, критерий согласия Пирсона (χ_0^2) имеет распределение хи-квадрат с числом степеней свободы, равным 2. Следовательно, если коэффициент корреляции Крамера = 0 для генеральной совокупности «все школьники Японии», вероятность того, что величина, полученная на основе данных, поступивших от случайно отобранных 300 школьников, будет, например, $> 5,9915$, равна 0,05. (см. Таблицу распределения хи-квадрат на стр. 103).



Рис. 7.2. Вероятность того, что $\chi_0^2 > 5,9915$.



Величина уже была вычислена и равняется 8,0091 (см. стр. 132).

Ну, как, не кажется ли вам, что величина слишком большая, хотя и вычислена на основе ответов, полученных от 300 человек, случайно выбранных из генеральной совокупности. Если поразмыслить, учитывая комментарий, данный на стр. 132, не естественно ли, что величина коэффициента корреляции Крамера для генеральной совокупности «все школьники Японии» > 0 ?

Не только в этом примере, но всегда при проверке гипотезы о независимости следует придерживаться следующей последовательности действий:

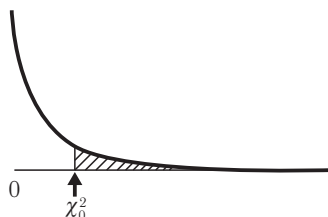
- 1) Сначала делают предположение, что коэффициент корреляции Крамера для генеральной совокупности = 0.
- 2) Затем рассчитывают величину χ_0^2 для выборочной совокупности.
- 3) Если χ_0^2 очень большой, делают вывод, что коэффициент корреляции Крамера для генеральной совокупности > 0 .



Это надо запомнить.

Добавлю кое-что к пункту 3. Чем больше величина χ_0^2 , тем меньше вероятность, т.е. площадь заштрихованной области на **Рис. 7.3**.

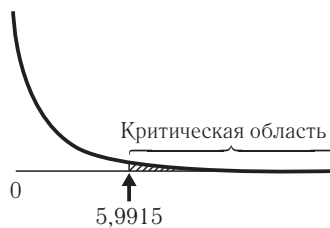
Рис. 7.3. Вероятность при определённом значении



Если при проверке гипотезы о независимости оказывается, что значение вероятности заштрихованной области меньше так называемого уровня значимости, делают вывод, что величина коэффициента корреляции Крамера для генеральной совокупности > 0 . Обычно в статистических исследованиях используют такой уровень значимости, как 0,05 или 0,01; выбор конкретного значения — прерогатива исследователя. Предположим, что выбрали значение уровня значимости, равное 0,05. В действительности, уровень значимости — это заштрихованная область на графике, показаном на **Рис. 7.3**.

При этом область, указанная на **Рис. 7.4**, называется критической областью.

Рис. 7.4. Критическая область при значении уровня значимости, равном 0,05





Вывод

Порядок выполнения проверки

Шаг 1

Определяем генеральную совокупность.

Генеральная совокупность — это:

Все школьники Японии



В этом упражнении изначально определено, что генеральной совокупностью являются «все школьники Японии». Поэтому нет необходимости выполнять Шаг 1.



Пример.

При проведении проверки гипотезы о равенстве долей в (двух) совокупностях (см. стр. 149) предполагалось, что генеральными совокупностями являются «избиратели, проживающие в городах и в сельской местности». При этом необходимо уточнить, что подразумевается под городом — Токио и Осака? Или это столицы префектур? Решение этой проблемы остаётся за исследователем, поскольку это его прерогатива — определить, что будет представлять собой генеральная совокупность при проверке. Если чётко не выделить генеральную совокупность, при проведении любой проверки можно оказаться в ситуации: «Ой! Что же я хотел исследовать?!». Подобная ситуация наблюдается довольно часто, поэтому следует быть предельно внимательными.

Шаг 2

Выстраиваем нулевую и альтернативную гипотезы.

Нулевая гипотеза:

Коэффициент корреляции Крамера для генеральной совокупности $= 0$.
«Пол» и «способ признания в любви» не связаны.

Альтернативная гипотеза:

Коэффициент корреляции Крамера для генеральной совокупности > 0 .
«Пол» и «способ признания в любви» связаны.



О нулевой и альтернативной гипотезах расскажу позже.

Шаг 3

Выбирают вид гипотезы для статистической проверки.

Проведём проверку гипотезы о независимости.



В этом примере изначально решено было проводить проверку гипотезы о независимости. Поэтому нет необходимости выполнять Шаг 3. В действительности, когда проводят проверку, исследователь сам должен выбрать вид гипотезы с учётом целей анализа.

Шаг 4

Определяют уровень значимости.

Пусть уровень значимости будет равен 0,05.



Поскольку будет использоваться уровень значимости, равный 0,05, нет необходимости выполнять Шаг 4. В действительности при проведении проверки исследователь должен сам выбрать уровень значимости. Обычно это 0,05 или 0,01. Уровень значимости обозначается α (альфа).

Шаг 5

Вычисляют фактическое значение выбранного статистического критерия на основе данных выборочной совокупности.

Я собираюсь провести проверку гипотезы о независимости. Поэтому в качестве статистического критерия будет выступать критерий согласия Пирсона χ_0^2 . Для данного примера значение критерия Пирсона χ_0^2 уже вычислено и равно 8,0091 (см. стр. 132).



Выбранный статистический критерий — это формула, которая преобразует данные выборочной совокупности в одну величину. В зависимости от вида проверяемых гипотез выбираются разные критерии. В случае, когда проводят проверку гипотезы о независимости, критерием является величина, указанная выше, а в случае проведения проверки об отсутствии корреляции (см. стр. 149) критерием будет величина, рассчитываемая по формуле:

$$\frac{\text{линейный коэффициент корреляции}^2 \times \sqrt{\text{общее число единиц совокупности} - 2}}{\sqrt{1 - \text{линейный коэффициент корреляции}^2}}$$

Шаг 6

Выясняют, входит ли фактическое значение выбранного статистического критерия, вычисленное на Шаге 5, в критическую область.

Значение критерия согласия Пирсона χ_0^2 , который в данном примере является статистическим критерием, равняется 8,0091.

Так как уровень значимости равен 0,05, из Таблицы хи-квадрат-распределения (стр. 103) получаем, что критическая область больше, чем 5,9915.

Как следует из рисунка, значение выбранного статистического критерия входит в критическую область.



Критическая область меняется в зависимости от уровня значимости α . Если бы уровень значимости был бы равен 0,01, а не 0,05, то критическая область, как следует из Таблицы распределения хи-квадрат на стр. 103, была бы больше, чем 9,2104.

Шаг 7

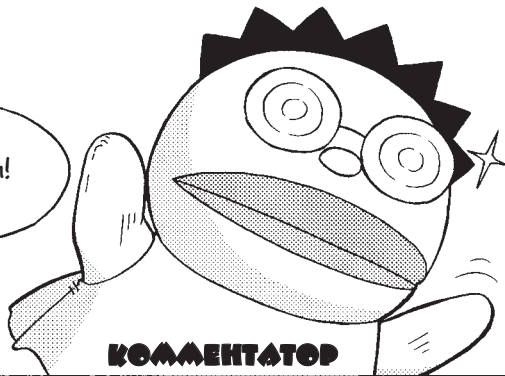
Если значение статистического критерия входит в критическую область (Шаг 6), делают вывод, что «верна альтернативная гипотеза». В противном случае вывод таков: «нельзя утверждать, что нулевая гипотеза ошибочна».

Значение выбранного статистического критерия входит в критическую область. Следовательно, верна альтернативная гипотеза — Величина коэффициента корреляции Крамера для генеральной совокупности > 0 . «Пол» и «способ признания в любви» связаны!



Даже если величина статистического критерия входит в критическую область, нельзя на основе проверки делать вывод, что «альтернативная гипотеза **абсолютно** верна». Можно сделать лишь такой вывод: «хотелось бы утверждать, что альтернативная гипотеза **абсолютно** верна, но существует вероятность ($\alpha \times 100\%$) того, что верна нулевая гипотеза».

Вот так вот!



КОММЕНТАТОР

Понятно...



Однако
меня беспокоит
Шаг 7.

?

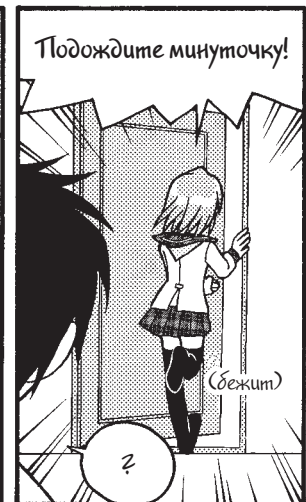


К сожалению, так сказать нельзя.
Можно только сделать вывод:
«нельзя утверждать,
что нулевая гипотеза
ошибочна».

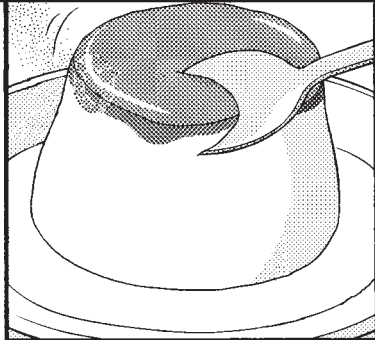
Вот как?





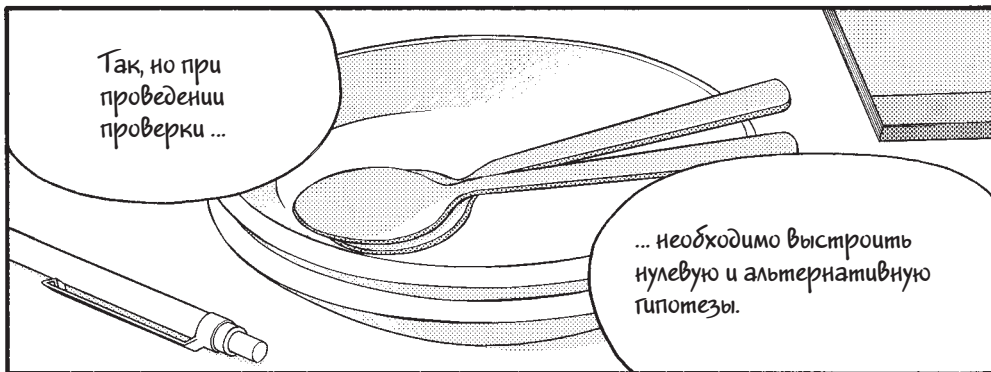


3. Нулевая и альтернативная гипотезы



Знаешь, если бы не твой пример, я бы так и не вспомнила, что у нас в холодильнике был пудинг.

Хорошо, что его никто не украл.



Так, но при проведении проверки ...

... необходимо выстроить нулевую и альтернативную гипотезы.



Что это за нулевая и альтернативная гипотезы?

Ты же сказал, что попозже объяснишь, а я до сих пор так ничего про них и не услышала..?



На самом деле, довольно сложно объяснить в двух словах, что такое нулевая и альтернативная гипотезы.

Вот как?



Примеры проверяемых гипотез	
Гипотезы	Примеры применения
О независимости	Проверяют, равняется ли нулю коэффициент корреляции Крамера между полом и способом признания в любви для генеральной совокупности.
О корреляционном отношении	Проверяют, равняется ли нулю корреляционное отношение между любимым брэндом одежды и возрастом для генеральной совокупности.
Об отсутствии корреляции	Проверяют, равняется ли нулю коэффициент линейной корреляции между расходами в месяц на косметику и расходами в месяц на одежду для генеральной совокупности.
О равенстве средних величин (двух) совокупностей	Проверяют, получают ли школьницы Токио и школьницы Осаки одну и ту же или разную сумму на карманные расходы. (Будьте внимательны: предполагаются две генеральные совокупности).
О равенстве долей в (двух) совокупностях	Проверяют, различается ли рейтинг кабинета министров среди избирателей, проживающих в городах, и избирателей из сельской местности. (Будьте внимательны: предполагаются две генеральные совокупности).



Проверка гипотезы о независимости

Нулевая гипотеза	Коэффициент корреляции Крамера между «полом» и «способом признания в любви» для генеральной совокупности $= 0$.
Альтернативная гипотеза	Коэффициент корреляции Крамера между «полом» и «способом признания в любви» для генеральной совокупности > 0 .

Проверка гипотезы о корреляционном отношении

Нулевая гипотеза	Корреляционное отношение между «любимым брендом одежды» и «возрастом» для генеральной совокупности $= 0$.
Альтернативная гипотеза	Корреляционное отношение между «любимым брендом одежды» и «возрастом» для генеральной совокупности > 0 .

Проверка гипотезы об отсутствии корреляции

Нулевая гипотеза	Коэффициент линейной корреляции между «расходами в месяц на косметику» и «расходами в месяц на одежду» для генеральной совокупности $= 0$.
Альтернативная гипотеза	Коэффициент линейной корреляции между «расходами в месяц на косметику» и «расходами в месяц на одежду» для генеральной совокупности $\neq 0$. или Коэффициент линейной корреляции между «расходами в месяц на косметику» и «расходами в месяц на одежду» для генеральной совокупности > 0 . или Коэффициент линейной корреляции между «расходами в месяц на косметику» и «расходами в месяц на одежду» для генеральной совокупности < 0 .

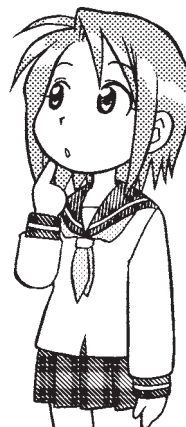
Проверка гипотезы о равенстве средних величин (двух) совокупностей

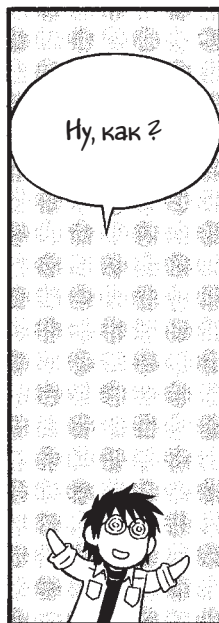
Нулевая гипотеза	«Суммы, получаемые на карманные расходы» школьниками Токио и школьниками Осаки, одинаковы.
Альтернативные гипотезы	«Суммы, получаемые на карманные расходы» школьниками Токио и школьниками Осаки, различны.
	«Сумма, получаемая на карманные расходы» школьниками Осаки, больше, чем сумма, получаемая школьниками Токио.
	«Сумма, получаемая на карманные расходы» школьниками Осаки меньше, чем сумма, получаемая школьниками Токио.

Проверка гипотезы о равенстве долей в (двух) совокупностях

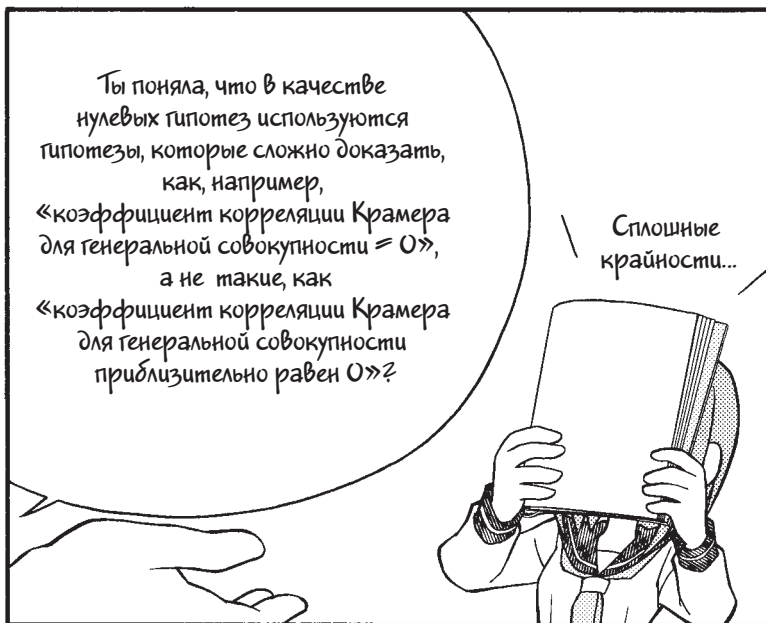
Нулевая гипотеза	Рейтинг кабинета министров среди избирателей, проживающих в городах и в сельской местности, одинаков.
Альтернативные гипотезы	Рейтинг кабинета министров среди избирателей, проживающих в городах и в сельской местности, различен.
	Рейтинг кабинет министров среди избирателей, проживающих в городах, выше, чем среди избирателей, проживающих в сельской местности.
	Рейтинг кабинет министров среди избирателей, проживающих в городах, ниже, чем среди избирателей, проживающих в сельской местности.

Понятно...



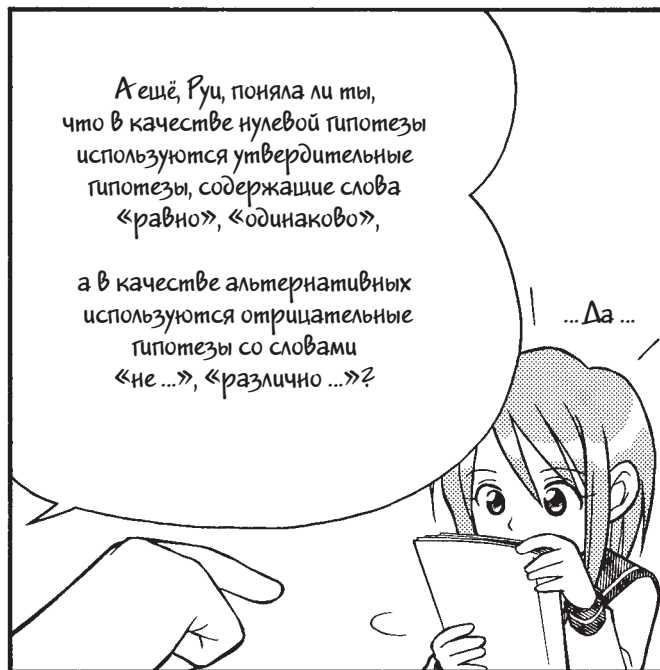


Ну, как ?



Ты поняла, что в качестве нулевых гипотез используются гипотезы, которые сложно доказать, как, например, «коэффициент корреляции Крамера для генеральной совокупности $\neq 0$ », а не такие, как «коэффициент корреляции Крамера для генеральной совокупности приблизительно равен 0»?

Сплошные крайности...



А ещё, Ру, поняла ли ты, что в качестве нулевой гипотезы используются утвердительные гипотезы, содержащие слова «равно», «одинаково», а в качестве альтернативных используются отрицательные гипотезы со словами «не...», «различно...»?

... Да ...

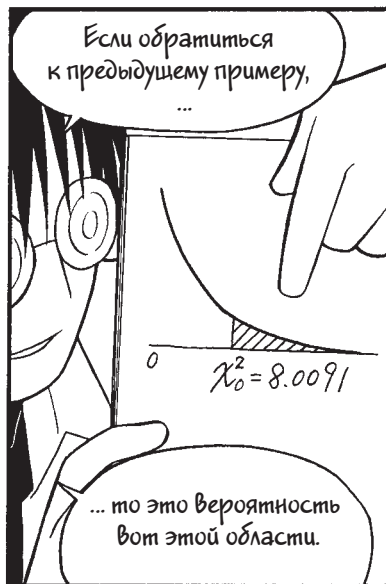
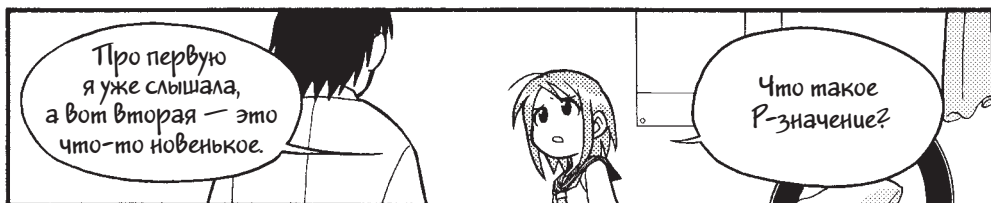


В качестве нулевых используют утвердительные гипотезы, которые сложно доказать, а в качестве альтернативных — отрицательные гипотезы..

Мм...
Вот как

Достаточно, если ты поймёшь только это.

4. Р-значение и порядок проверки





Шаг 6р

Выясняют, меньше ли P -значение, соответствующее величине выбранного статистического критерия, вычисленного на **Шаге 5**, чем уровень значимости.

Уровень значимости равен 0,05.
Поскольку величина критерия согласия Пирсона равна 8,0091, P -значение равно 0,0182 и, следовательно, $0,0182 < 0,05$.
Другими словами, P -значение меньше, чем выбранный статистический критерий.



Как я уже говорил, P -значение можно вычислить, используя функции Excel (но способ вычисления зависит от вида гипотезы).
Например, в Excel можно вычислить величину P -значения для проверки гипотезы о независимости.
Подробное объяснение смотрите на стр. 208.

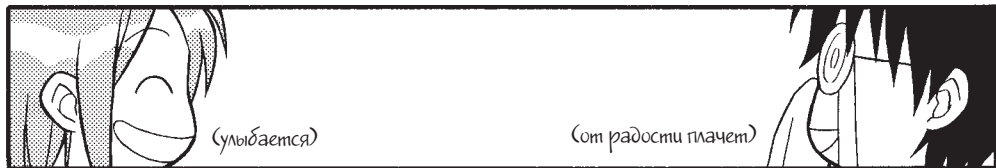
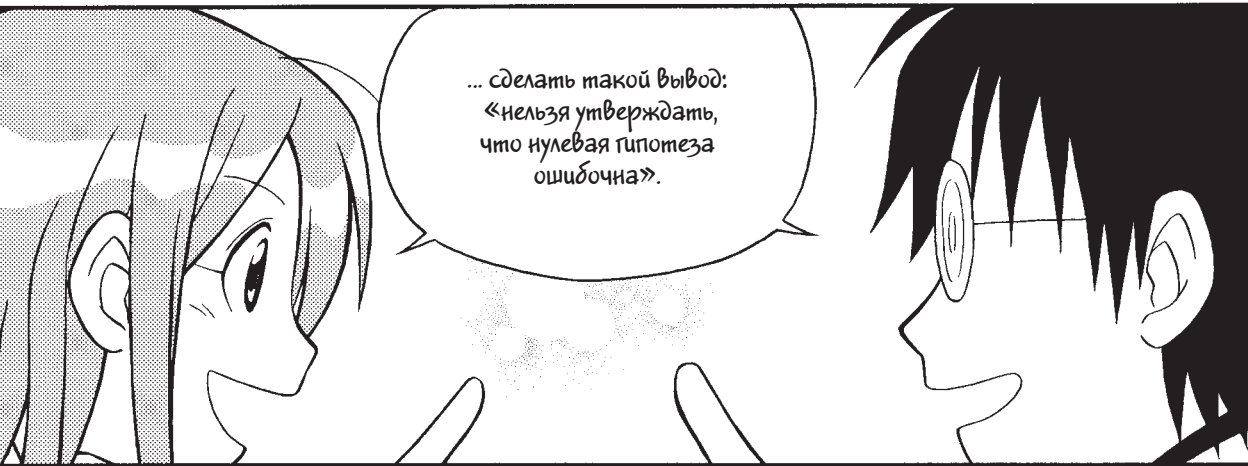
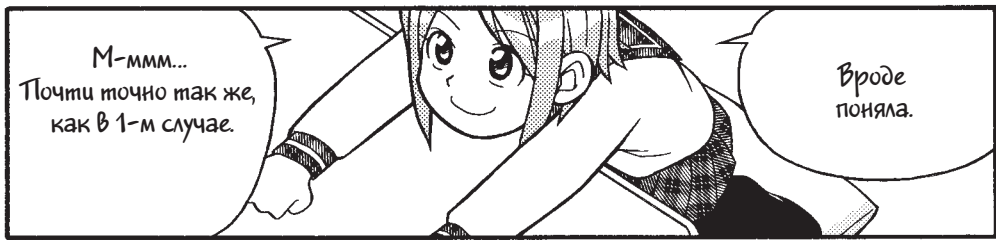
Шаг 7р

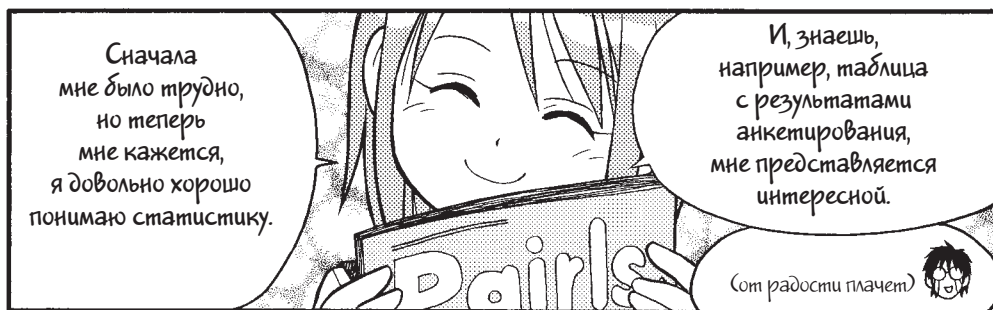
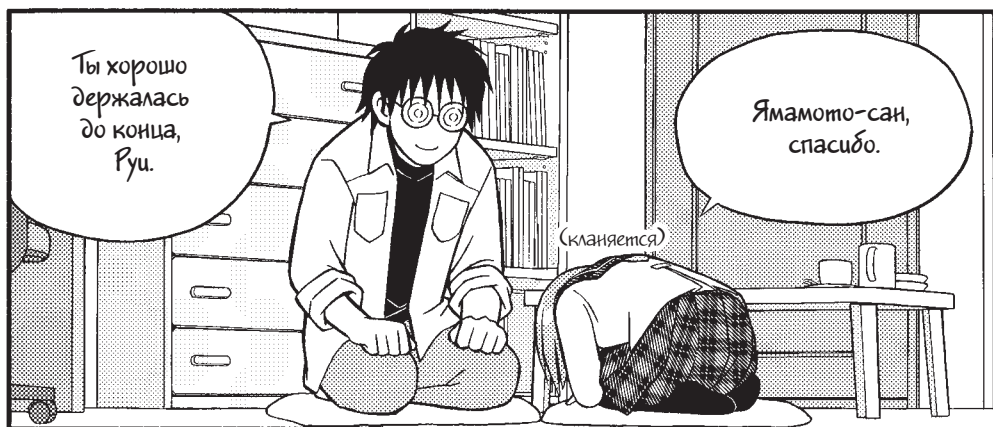
Если P -значение меньше, чем уровень значимости (см. *Шаг 6р*), делают вывод, что «верна альтернативная гипотеза». В противном случае вывод такой: «нельзя утверждать, что нулевая гипотеза ошибочна».

P -значение оказалось меньше уровня значимости. Следовательно, верна альтернативная гипотеза: «величина коэффициента корреляции Крамера для генеральной совокупности > 0 », т.е. «пол» и «способ признания в любви» связаны!



Даже когда P -значение оказывается меньше уровня значимости, нельзя на основе проверки делать вывод, что «альтернативная гипотеза абсолютно верна». Можно сделать лишь такой вывод: «хотелось бы утверждать, что альтернативная гипотеза абсолютно верна, но существует вероятность ($\alpha \times 100\%$) того, что верна нулевая гипотеза».

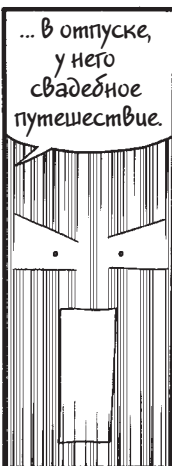






Решено.
Я пойду
на встречу
с Игараси-саном.

Ой, а
Игараси
сейчас ...



... в отпуске,
у него
свадебное
путешествие.



Что????!!!!!!
Он женился?



Зачем же я столько
времени и сил угрохала
на эту статистику...?!

Он женат....

Что?
Разве твой
интерес был
неискренним?



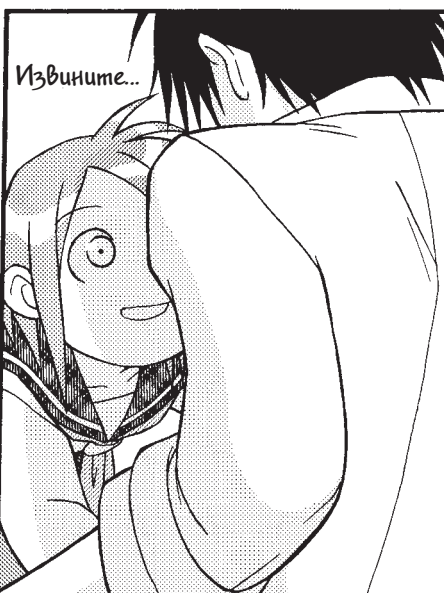
Оставь меня в покое!!

(бежит
в сторону
Ямамото)

Ой,
осторожно!



БУМ



Извините...



Все в порядке?



Я и не знала, что Ямамото-сан такой симпатичный!



Ямамото-сан,
научите меня
ещё чему-нибудь!

И их занятия продолжились ...
... а может и нет.

5. Проверка гипотезы о независимости и гипотезы об однородности

Проверка гипотезы об однородности очень похожа на проверку гипотезы о независимости. Ниже приводится пример такой проверки. Попробуйте разобраться, в чём заключается разница между этими двумя проверками.

Упражнение

Женский журнал «P-girls» решил провести опрос школьников на тему «Какой способ признания в любви Вы предпочитаете:

- по телефону;
- по SMS;
- при встрече»?

Журнал выдвинул такую гипотезу:

Гипотеза:

пропорции ответов «по телефону» : «по SMS» : «при встрече» зависят от пола респондентов.

Чтобы выяснить, правильна ли эта гипотеза, журналист произвольно выбрал определённое число юношей и девушек из «всех школьников Японии» и провёл анкетирование. Результаты этого анкетирования приведены в следующей таблице:

Респондент	Способ признания в любви	Возраст респондента	Пол респондента
1 ...	при встрече	17	ж
148	по SMS	16	ж
149 ...	по телефону	15	м
300	по SMS	18	м

Таблица взаимной сопряжённости «пола» и «способа признания в любви»

Пол респондента	Способ признания в любви			Итого
	по телефону	по SMS	при встрече	
женский	34	61	53	148
мужской	38	40	74	152
Итого:	72	101	127	300

Проверьте правильность сформулированной гипотезы путём проверки гипотезы об однородности. Предположим, что уровень значимости равен 0,05.

Решение

1	Определим генеральную совокупность	Генеральных совокупностей будет две: «все школьницы Японии» и «все школьники Японии».
2	Сформулируем основную и противоположную гипотезы	Нулевая гипотеза: «пропорции предпочтительных способов признания в любви — по телефону : по SMS : при встрече — у школьниц и школьников одинаковы». Альтернативная гипотеза: «пропорции предпочтительных способов признания в любви — по телефону : по SMS : при встрече — у школьниц и школьников различны».
3	Выбираем вид гипотезы для статистической проверки	Проведём проверку гипотезы об однородности.
4	Определим уровень значимости	Пусть уровень значимости равен 0,05.
5	Вычислим фактическое значение выбранного статистического критерия на основе данных выборочной совокупности	В этом упражнении будет проведена проверка гипотезы об однородности. Следовательно, статистическим критерием будет критерий согласия Пирсона. Величина была уже вычислена и равна 8,0091 (см. стр. 132). Если нулевая гипотеза верна, критерий согласия Пирсона имеет распределение хи-квадрат с числом степеней свободы, равным $(2 - 1) \times (3 - 1) = 1 \times 2 = 2$.
6	Проверим, входит ли значение вычисленного на Шаге 5 статистического критерия в критическую область	Величина критерия согласия Пирсона, который является статистическим критерием, равна 8,0091. Так как уровень значимости = 0,05, критическая область, как следует из таблицы распределения хи-квадрат на стр. 103, больше 5,9915. Это значит, что значение статистического критерия входит в критическую область.
7	Если значение статистического критерия входит в критическую область (Шаг 6), делают вывод: «верна альтернативная гипотеза». Если нет, — «нельзя утверждать, что нулевая гипотеза ошибочна»	Значение статистического критерия входит в критическую область. Следовательно, верна альтернативная гипотеза: «пропорции предпочтительных способов признания в любви — по телефону : по SMS : при встрече — у школьниц и школьников различны».

Как вам? И упражнение, и ответ как две капли воды похожи на проверку гипотезы о независимости. Давайте уточним, чем же отличается проверка гипотезы о независимости от проверки гипотезы об однородности.

Во-первых, определённые нами генеральные совокупности разные. В проверке гипотезы о независимости генеральная совокупность только одна — «все школьники Японии», а в проверке гипотезы об однородности генеральных совокупностей две: «все школьницы Японии» и «все школьники Японии».

Во-вторых, нулевые и альтернативные гипотезы разные. При проверке гипотезы о независимости были сформулированы следующие гипотезы:

Нулевая гипотеза	Коэффициент корреляции Крамера для генеральной совокупности = 0 и, значит, «пол» и «способ признания в любви» не связаны.
Альтернативная гипотеза	Коэффициент корреляции Крамера для генеральной совокупности > 0 и, значит «пол» и «способ признания в любви» связаны.

а при проверке гипотезы об однородности были сформулированы гипотезы:

Нулевая гипотеза	Пропорции предпочтительных способов признания в любви — по телефону : по SMS : при встрече — у школьниц и школьников одинаковы.
Альтернативная гипотеза	Пропорции предпочтительных способов признания в любви — по телефону : по SMS : при встрече — у школьниц и школьников различны.

К тому же, гипотезы формулировались в разные моменты: в случае проверки гипотезы о независимости — после сбора данных, а в случае проверки гипотезы об однородности — до сбора данных. Однако, несмотря на очевидные различия на практике часто бывает так: собираются проводить проверку гипотезы о независимости, а на деле проводят проверку гипотезы об однородности или наоборот. **Будьте внимательны!**

6. Как выразить словами вывод на основании проверки

Вывод, полученный на основании проверки, формулируется следующим образом:

Если величина статистического критерия входит в критическую область, делают вывод о том, что верна альтернативная гипотеза. В противном случае вывод таков: «нельзя утверждать, что нулевая гипотеза ошибочна».

На самом деле такие выражения для формулирования вывода не используются. Различные выражения, которые действительно используются для формулирования вывода, сделанного на основе результатов статистической проверки, приведены в Табл. 7.4.

Таблица 7.4. Выражения, используемые для формулирования выводов на основе проверки

Случаи, когда величина статистического критерия входит в критическую область	Случаи, когда величина статистического критерия не входит в критическую область
<ul style="list-style-type: none">• верна альтернативная гипотеза• альтернативная гипотеза значима• нулевая гипотеза отвергается	<ul style="list-style-type: none">• нельзя утверждать, что нулевая гипотеза ошибочна• нельзя отвергнуть нулевую гипотезу• нулевая гипотеза отклонена• нельзя утверждать, что нулевая гипотеза верна• нулевая гипотеза принимается

Выражения «гипотеза значима» и «гипотеза отвергнута» используются достаточно часто. Вместе с тем я специально использовал формулировки, которые, как правило, не используются. Объясняю почему. Я заметил, что среди начинающих изучать статистическую проверку гипотез, есть такие, которые часто говорят «гипотеза значима», причём не очень хорошо понимая, какой смысл имеет данное выражение. Очевидно, они используют это выражение только потому, что уверены в величине P -значения и получили значение статистического критерия. Другими словами, они проводят статистическую проверку гипотез, не сформулировав чётко нулевую и альтернативную гипотезы. И, как мне кажется, генеральная совокупность также чётко не определена. Я раньше полагал, что не следует делать замечания начинающим постигать премудрости статистической проверки гипотез. Однако без чётко выстроенных гипотез невозможно сделать какой-либо вывод.

В связи с этим я использую такие выражения, как «верна альтернативная гипотеза» и «нельзя утверждать, что основная гипотеза ошибочна». Это позволит читателю понять и усвоить, что такое нулевая и альтернативная гипотезы.

Упражнение

Таблица взаимной сопряжённости (см. стр. 138).

Обычно заказываемая кухня	Предпочитаемый напиток		Итого
	Кофе	Чай	
Японская	43	33	76
Европейская	51	53	104
Китайская	29	41	70
Итого:	123	127	250

Выясните путём проверки гипотезы о независимости, больше ли 0 коэффициент корреляции Крамера между видом обычно заказываемой кухни и предпочитаемым напитком для генеральной совокупности «жители Японии старше 20 лет». Другими словами, есть ли взаимосвязь между видом обычно заказываемой кухни и предпочитаемым напитком? Пусть уровень значимости равен 0,01.

Ответ

1	Определим генеральную совокупность	Генеральной совокупностью будут «жители Японии старше 20 лет»
2	Сформулируем нулевую и альтернативную гипотезы	Основная гипотеза: «вид обычно заказываемой кухни и предпочитаемый напиток не связаны». Альтернативная гипотеза: «вид обычно заказываемой кухни и предпочитаемый напиток взаимосвязаны».
3	Выбираем вид гипотезы	Проведём проверку гипотезы о независимости.
4	Определим уровень значимости	Пусть уровень значимости равен 0,01.
5	Вычислим фактическое значение выбранного статистического критерия на основе данных выборочной совокупности	В этом упражнении будет проведена проверка гипотезы о независимости. Следовательно, статистическим критерием будет являться критерий согласия Пирсона. Величина была вычислена ранее и равна 3,34839 (см. стр. 141)
6	Проверим, входит ли значение вычисленного на Шаге 5 статистического критерия в критическую область	Величина критерия согласия Пирсона, являющегося статистическим критерием, равна 3,3483. Так как уровень значимости $\alpha = 0,01$, критическая область, как следует из таблицы распределения хи-квадрат на стр. 103, больше 9,2104. Это значит, что значение статистического критерия не входит в критическую область.
7	Если значение статистического критерия входит в критическую область (Шаг 6), делают вывод: «верна альтернативная гипотеза». Если нет — «нельзя утверждать, что нулевая гипотеза ошибочна».	Значение статистического критерия не входит в критическую область. Следовательно, нельзя утверждать, что нулевая гипотеза: «вид обычно заказываемой кухни и предпочитаемый напиток не связаны» ошибочна.

Выводы

- Проверка — один из способов анализа, который позволяет установить, правильна ли гипотеза, сделанная исследователем о генеральной совокупности на основе данных выборочной совокупности.
- Такая проверка называется статистической проверкой гипотез.
- Статистический критерий — формула, с помощью которой преобразуются данные выборочной совокупности.
- Обычно используются такие значения уровня значимости, как 0,05 или 0,01.
- Критическая область — область значений, соответствующая определённому уровню значимости.
- Проверка гипотезы о независимости — один из способов анализа, который позволяет выяснить, не равен ли 0 коэффициент корреляции Крамера для генеральной совокупности. Можно также сказать, что такой анализ позволяет выяснить, есть ли связь между двумя переменными в таблице взаимной сопряжённости.
- Если коэффициент корреляции Крамера для генеральной совокупности равен 0, величина критерия согласия Пирсона имеет распределение хи-квадрат.
- Р-значение в случае проверки гипотезы о независимости при условии, что нулевая гипотеза верна, — вероятность того, что критерий согласия Пирсона χ_0^2 больше или равен наблюдаемой величины.
- Сделать вывод о проверке можно на основании того, что:
 - 1) значение статистического критерия входит в критическую область, либо
 - 2) Р-значение меньше уровня значимости.
- Будь то проверка гипотезы о независимости или любая другая проверка, порядок проведения анализа один и тот же, при этом он предусматривает выполнение следующих действий:

Шаг 1	Определить генеральную совокупность.
Шаг 2	Сформулировать нулевую и альтернативную гипотезы.
Шаг 3	Выбрать вид статистической проверки гипотезы.
Шаг 4	Определить уровень значимости.
Шаг 5	Вычислить фактическое значение выбранного статистического критерия на основе данных выборочной совокупности.
Шаг 6	Проверить, входит ли вычисленное на Шаге 5 значение статистического критерия в критическую область.
Шаг 7	Если Р-значение меньше, чем уровень значимости (см. Шаг 6р), то делают вывод о том, что верна альтернативная гипотеза. Если нет, то вывод «нельзя утверждать, что нулевая гипотеза ошибочна».
Шаг 6р	Проверить, меньше ли Р-значение, соответствующее величине выбранного статистического критерия, вычисленного на Шаге 5, чем уровень значимости.

Приложение

**Попробуем
вычислить
с помощью
Excel**

В этой главе объясняется, как с помощью программы Excel:

1. Построить таблицу (ряд) распределения.
2. Вычислить среднее значение, медиану и стандартное отклонение.
3. Составить простую статистическую таблицу.
4. Вычислить нормированное отклонение и рассчитать T-показатель.
5. Вычислить вероятность стандартного нормального распределения.
6. Вычислить значение x при распределении хи-квадрат.
7. Вычислить коэффициент линейной корреляции.
8. Проверить гипотезу о независимости.

Файлы Excel можно загрузить по адресу <http://www.dodeca.ru/books/33081.php>. Читателю, не имеющему опыта работы в Excel, рекомендуется сначала попробовать вычислить среднее значение, медиану и стандартное отклонение (см. стр. 195).

1. Построение таблиц распределения

Используются данные со стр. 33.

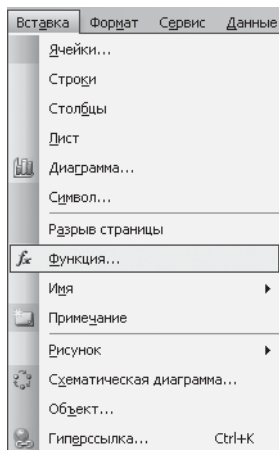


Выберите ячейку J3.

	A	B	C	D	E	F	G	H	I	J
1	Ресторан	Цена, йены		Ресторан	Цена, йены					
2	ресторан 1	700		ресторан 26	780		от	до	меньше	частота
3	ресторан 2	850		ресторан 27	590		500	600	599	
4	ресторан 3	600		ресторан 28	650		600	700	699	
5	ресторан 4	650		ресторан 29	580		700	800	799	
6	ресторан 5	980		ресторан 30	750		800	900	899	
7	ресторан 6	750		ресторан 31	800		900	1000	999	
8	ресторан 7	500		ресторан 32	550					
9	ресторан 8	890		ресторан 33	750					
10	ресторан 9	880		ресторан 34	700					
11	ресторан 10	700		ресторан 35	600					
12	ресторан 11	890		ресторан 36	800					
13	ресторан 12	720		ресторан 37	800					
14	ресторан 13	680		ресторан 38	880					
15	ресторан 14	650		ресторан 39	790					
16	ресторан 15	790		ресторан 40	790					
17	ресторан 16	670		ресторан 41	780					
18	ресторан 17	680		ресторан 42	600					
19	ресторан 18	900		ресторан 43	670					
20	ресторан 19	880		ресторан 44	680					
21	ресторан 20	720		ресторан 45	650					
22	ресторан 21	850		ресторан 46	890					
23	ресторан 22	700		ресторан 47	930					
24	ресторан 23	780		ресторан 48	650					
25	ресторан 24	850		ресторан 49	777					
26	ресторан 25	750		ресторан 50	700					

Шаг 2

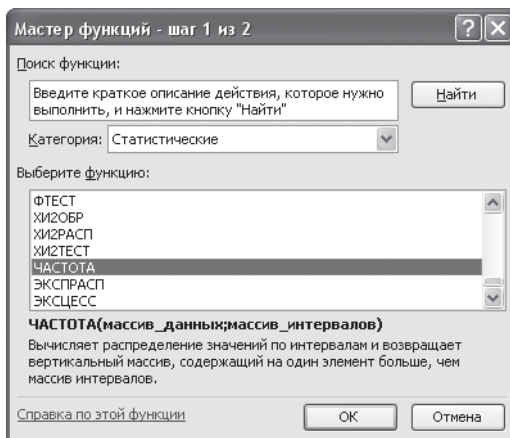
Выберите **Вставка** ▶ **Функция**.



Шаг 3

Выберите **Статистические** в строке **Категория**, а затем **ЧАСТОТА*** в графе **Выберите функцию**.

* FREQUENCY в английской версии Excel.



Шаг 4 Выделите нижеуказанную область и нажмите кнопку **OK**.

	A	B	C	D	E	F	G	H	I	J
1	Ресторан	Цена, йены		Ресторан	Цена, йены					
2	ресторан 1	700		ресторан 26	780					
3	ресторан 2	850		ресторан 27	590		от	до	меньше	частота
4	ресторан 3	600		ресторан 28	650		500	600	599	{3;13:17}
5	ресторан 4	650		ресторан 29	580		600	700	699	
6	ресторан 5	980		ресторан 30	750		700	800	799	
7	ресторан 6	750		ресторан 31	800		800	900	899	
8	ресторан 7	500		ресторан 32	550		900	1000	999	
9	рест	Аргументы функции								
10	рест	ЧАСТОТА								
11	рестс	Массив_данных B2:E26 = {700;0;"ресторан 2								
12	рестс	Массив_интервалов I3:I7 = {599;699;799;899;9								
13	рестс	= {4;13;18;12;3;0}								
14	рестс	Вычисляет распределение значений по интервалам и возвращает вертикальный массив, содержащий на один элемент больше, чем массив интервалов.								
15	рестс	Массив_интервалов массив интервалов или ссылка на интервалы, в которых группируются значения из массива данных.								
16	рестс	Справка по этой функции								
17	рестс	Значение: 4								
18	рестс	<input type="button" value="OK"/> <input type="button" value="Отмена"/>								
19	рестс									
20	рестс									
21	рестс									
22	рестс									
23	рестс									
24	рестс									
25	рестс									
26	ресторан 25	750		ресторан 30	700					

Шаг 5 Выделите ячейки, начиная с ячейки J3 и до ячейки J7.

	G	H	I	J
	от	до	меньше	частота
	500	600	599	4
	600	700	699	
	700	800	799	
	800	900	899	
	900	1000	999	

Шаг 6 Щелкните мышью на области в строке формул.

id	Вставка	Формат	Сервис	Данные
[Иконки]				
fx =ЧАСТОТА(B2:E26;I3:I7)				
B	C	D	E	

Шаг 7 Нажмите комбинацию клавиш **Shift + Ctrl** и, удерживая ее в нажатом состоянии, нажмите клавишу **Enter**.

Шаг 8 Вычисление закончено!

G	H	I	J
от	до	меньше	частота
500	600	599	4
600	700	699	13
700	800	799	18
800	900	899	12
900	1000	999	3

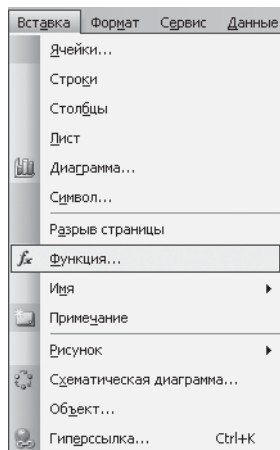
2. Вычисление среднего значения, медианы и стандартного отклонения

Используются данные со стр. 41.

Шаг 1 Выберите ячейку В10.

	A	B	C
1		Команда А	
2	Руи-Руи	86	
3	Дэюн	73	
4	Юми	124	
5	Сизука	111	
6	Токо	90	
7	Каэдэ	38	
8			
9			
10	среднее		
11	медиана		
12	отклонение		

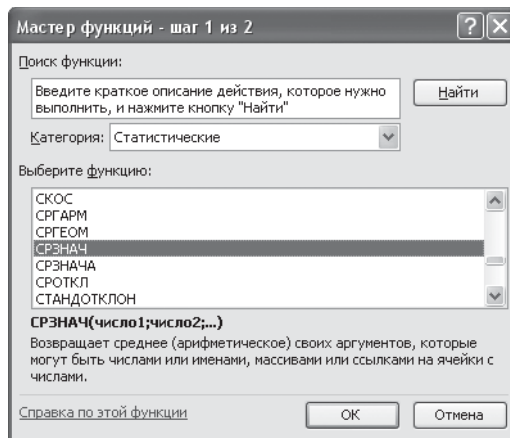
Шаг 2 Выберите Вставка ► Функция.



Шаг 3

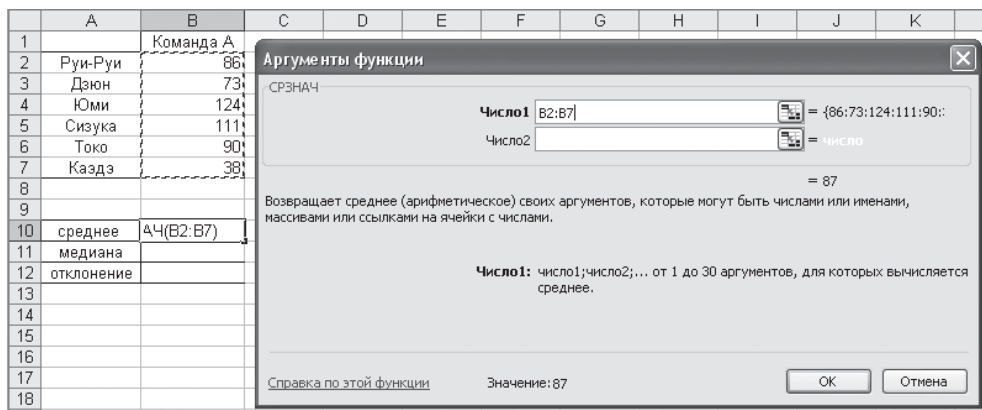
Выберите **Статистические** в строке **Категория**, и **СРЗНАЧ*** в графе **Выберите функцию**.

* AVERAGE в английской версии Excel.



Шаг 4

Выделите нижеуказанную область и нажмите кнопку **OK**.



Шаг 5

Вычисление закончено!

	А	В
1		Команда А
2	Руи-Руи	86
3	Дзюн	73
4	Юми	124
5	Сизука	111
6	Токо	90
7	Каздз	38
8		
9		
10	среднее	87
11	медиана	
12	отклонение	

Шаг 6

Вычислите медиану и стандартное отклонение, выполняя последовательно Шаг 1 — Шаг 5. При вычислении медианы выберите **МЕДИАНА*** в графе **Выберите функцию**, а при вычислении стандартного отклонения — **СТАНДОТКЛОНП****.

* MEDIAN в английской версии Excel.

** STDEVP в английской версии Excel.

3. Построение простой статистической таблицы

Используются данные со стр. 61.

Шаг 1

Выберите ячейку F20.

	A	B	C	D	E	F	G	H
1		Новая форма...			Новая форма...			Новая форма...
2	1	нравится		16	так себе		31	так себе
3	2	так себе		17	нравится		32	так себе
4	3	нравится		18	нравится		33	нравится
5	4	так себе		19	нравится		34	не нравится
6	5	не нравится		20	нравится		35	нравится
7	6	нравится		21	нравится		36	нравится
8	7	нравится		22	нравится		37	нравится
9	8	нравится		23	не нравится		38	нравится
10	9	нравится		24	так себе		39	так себе
11	10	нравится		25	нравится		40	нравится
12	11	нравится		26	нравится			
13	12	нравится		27	не нравится			
14	13	так себе		28	нравится			
15	14	нравится		29	нравится			
16	15	нравится		30	нравится			
17								
18								
19						частота		
20					нравится			
21					так себе			
22					не нравится			
23								

Шаг 2

Выберите пункт **Вставка ▶ Функция**.

Шаг 3

Выберите **Статистические** в строке **Категория**, а затем **СЧЕТЕСЛИ*** в графе **Выберите функцию**.

* COUNTIF в английской версии Excel.

- Шаг 4** Выделите указанную ниже область, напишите **нравится** в графе **Критерий** и нажмите кнопку **ОК**.

	A	B	C	D	E	F	G	H	I
1		Новая форма...			Новая форма...			Новая форма...	
2									
3									
4									
5									
6									
7									
8									
9									
10									
11									
12									
13									
14									
15									
16									
17									
18									
19									
20									
21									
22									
23									

Аргументы функции

СЧЁТЕСЛИ

Диапазон: A2:H16 = {1;"нравится";0};16;

Критерий: нравится =

= 0

Подсчитывает количество непустых ячеек в диапазоне, удовлетворяющих заданному условию.

Критерий: условие в форме числа, выражения или текста, который определяет, какие ячейки надо подсчитывать.

Справка по этой функции Значение: 0

	A	B	C	D	E	F	G	H
19								
20								
21								
22								
23								

- Шаг 5**
- Вычисление закончено!

	A	B	C	D	E	F	G	H
1		Новая форма...			Новая форма...			Новая форма...
2	1	нравится		16	так себе		31	так себе
3	2	так себе		17	нравится		32	так себе
4	3	нравится		18	нравится		33	нравится
5	4	так себе		19	нравится		34	не нравится
6	5	не нравится		20	нравится		35	нравится
7	6	нравится		21	нравится		36	нравится
8	7	нравится		22	нравится		37	нравится
9	8	нравится		23	не нравится		38	нравится
10	9	нравится		24	так себе		39	так себе
11	10	нравится		25	нравится		40	нравится
12	11	нравится		26	нравится			
13	12	нравится		27	не нравится			
14	13	так себе		28	нравится			
15	14	нравится		29	нравится			
16	15	нравится		30	нравится			
17								
18								
19								
20								
21								
22								
23								

- Шаг 6** Выполняя последовательно Шаг 1 — Шаг 5, подсчитайте частоту ответов «так себе» и «не нравится».

4. Вычисление нормированного отклонения и рейтинга успеваемости

Используются данные со стр. 72.

Функция для расчета нормированного отклонения в Excel есть, а функции для расчета рейтинга успеваемости (Т-показателя) нет.

Несмотря на это, можно довольно быстро вычислить рейтинг успеваемости, используя результат вычисления нормированного отклонения.

Поэтому в данной книге считается, что рейтинг успеваемости можно рассчитать в Excel.

4.1. Вычисление нормированного отклонения

Шаг 1 Выберите ячейку E2.

		История		Нормированное отклонение	Рейтинг успеваемости
1					
2	Руи	73	Руи		
3	Юми	61	Юми		
4	А	14	А		
5	Б	41	Б		
6	В	49	В		
7	Г	87	Г		
8	Д	69	Д		
9	Е	65	Е		
10	Ж	36	Ж		
11	З	7	З		
12	И	53	И		
13	К	100	К		
14	Л	57	Л		
15	М	45	М		
16	Н	56	Н		
17	О	34	О		
18	П	37	П		
19	Р	70	Р		
20	Среднее	53			
21	Отклонение	22.7			

Шаг 2 Выберите **Вставка** ▶ **Функция**.

Шаг 3 Выберите **Статистические** в строке **Категория**, а затем **НОРМАЛИЗАЦИЯ*** в графе **Выберите функцию**.

* STANDARDIZE в английской версии Excel.

Шаг 4 Выберите ячейку B2.

	A	B	C	D	E	F	G	H	I	J
1		История			Нормированное отклонение	Рейтинг успеваемости				
2	Руи	73		Руи	НОРМАЛИЗАЦИЯ(B2)					
3	Юми	61		Юми						
4	А	14								
5	Б	41								
6	В	49								
7	Г	87								
8	Д	69								
9	Е	65								
10	Ж	36								
11	З	7								
12	И	53								
13	К	100								
14	Л	57								
15	М	45								
16	Н	56								
17	О	34								
18	П	37								
19	Р	70								
20	Среднее	53								
21	Отклонение	22.7								

Аргументы функции

НОРМАЛИЗАЦИЯ

x B2 = 73

Среднее = число

Стандартное_откл = число

=

Возвращает нормализованное значение.

X нормализуемое значение.

Справка по этой функции Значение:

Шаг 5 Выберите ячейку B20 в графе **Среднее** и нажмите клавишу F4. Убедитесь, что B20 в графе **Среднее** превратилось в **\$B\$20**.

Аргументы функции

НОРМАЛИЗАЦИЯ

x B2 = 73

Среднее \$B\$20 = 53

Стандартное_откл = число

=

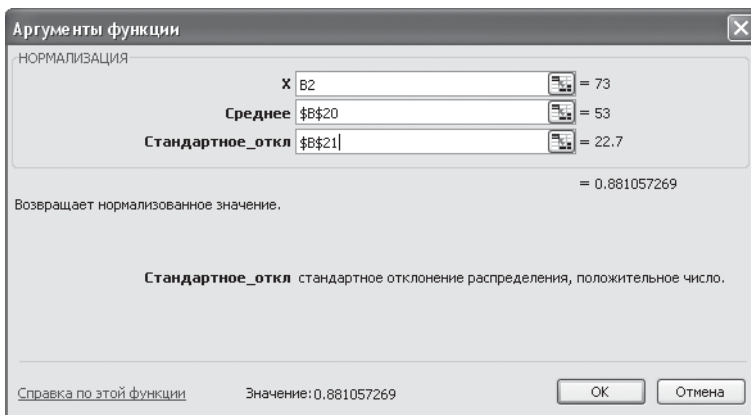
Возвращает нормализованное значение.

Среднее арифметическое среднее распределения.

Справка по этой функции Значение:

Шаг 6

Выберите ячейку **B21** в графе **Стандартное_откл** и нажмите клавишу **F4**. Убедитесь, что **B21** в графе **Стандартное_откл** превратилось в **\$B\$21** и нажмите кнопку **ОК**.



Шаг 7

Убедитесь, что было вычислено нормированное отклонение Руи.

	A	B	C	D	E	F
1		История			Нормированное отклонение	Рейтинг успеваемости
2	Руи	73		Руи	0.88	
3	Юми	61		Юми		
4	А	14		А		
5	Б	41		Б		
6	В	49		В		
7	Г	87		Г		
8	Д	69		Д		
9	Е	65		Е		
10	Ж	36		Ж		
11	З	7		З		
12	И	53		И		
13	К	100		К		
14	Л	57		Л		
15	М	45		М		
16	Н	56		Н		
17	О	34		О		
18	П	37		П		
19	Р	70		Р		
20	Среднее	53				
21	Отклонение	22.7				

Шаг 8

Подведите указатель мыши к правому углу ячейки **E2**, и как только он превратится в чёрный крестик, нажмите клавишу мыши и, удерживая её в нажатом состоянии, растяните область до ячейки **E19**; отпустите клавишу мыши.

D	E	F
	Нормированное отклонение	Рейтинг успеваемости
Руи	0.88	
Юми		
А		
Б		
В		
Г		
Д		
Е		
Ж		
З		
И		
К		
Л		
М		
Н		
О		
П		
Р		

Шаг 9

Вычисление нормированного отклонения закончено!

D	E	F
	Нормированное отклонение	Рейтинг успеваемости
Руи	0.88	
Юми	0.35	
А	-1.71	
Б	-0.53	
В	-0.18	
Г	1.49	
Д	0.70	
Е	0.53	
Ж	-0.75	
З	-2.02	
И	0.00	
К	2.07	
Л	0.18	
М	-0.35	
Н	0.13	
О	-0.84	
П	-0.70	
Р	0.75	

4.2. Вычисление рейтинга успеваемости

Шаг 10

Выберите ячейку **F2** и «напишите» точно так же, как пишете текст в Word, выражение $=E2*10+50$, а затем нажмите клавишу **Enter**.

D	E	F
	Нормированное отклонение	Рейтинг успеваемости
Руи	0.88	$=E2*10+50$
Юми	0.35	
А	-1.71	
Б	-0.53	
В	-0.18	
Г	1.49	
Д	0.70	
Е	0.53	
Ж	-0.75	
З	-2.02	
И	0.00	
К	2.07	
Л	0.18	
М	-0.35	
Н	0.13	
О	-0.84	
П	-0.70	
Р	0.75	

Шаг 11

Повторите Шаг 8.

Шаг 12

Вычисление рейтинга успеваемости закончено!

D	E	F
	Нормированное отклонение	Рейтинг успеваемости
Руи	0.88	58.79
Юми	0.35	53.52
А	-1.71	32.85
Б	-0.53	44.72
В	-0.18	48.24
Г	1.49	64.95
Д	0.70	57.04
Е	0.53	55.28
Ж	-0.75	42.53
З	-2.02	29.77
И	0.00	50.00
К	2.07	70.67
Л	0.18	51.76
М	-0.35	46.48
Н	0.13	51.32
О	-0.84	41.65
П	-0.70	42.96
Р	0.75	57.47

5. Вычисление вероятности стандартного нормального распределения

Используются данные со стр. 93.

Шаг 1 Выберите ячейку **B2**.

	A	B
1	z	1.96
2	промежуточное значение	
3	площадь (= доля = вероятность)	

Шаг 2 Выберите **Вставка** ▶ **Функция**.

Шаг 3 Выберите **Статистические** в строке **Категория**, а затем **НОРМСТРАСП*** в графе **Выберите функцию**.

* NORMDIST в английской версии Excel.

Шаг 4 Выберите ячейку **B1** и нажмите кнопку **ОК**.

	A	B	C	D	E	F	G
1	z	1.96					
2	промежуточное значение	СП(B1)					
3	площадь (= доля = вероятность)						

Аргументы функции

НОРМСТРАСП

z: B1 = 1.96

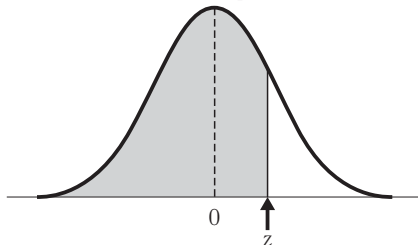
= 0.975002105

Возвращает стандартное нормальное интегральное распределение.

z значение, для которого строится распределение.

Справка по этой функции Значение: 0,975002105

Шаг 5 На самом деле функция **НОРМСТРАСП** предназначена для вычисления вероятности, указанной ниже на рисунке.



Поэтому «напишите» в ячейке В3 точно так же, как пишете в Word, выражение [=B2-0.5].

	A	B
1	z	1.96
2	промежуточное значение	0.975002
3	площадь (= доля = вероятность)	=B2-0.5

Шаг 6 Вычисление закончено!

	A	B
1	z	1.96
2	промежуточное значение	0.975002
3	площадь (= доля = вероятность)	0.475002

6. Вычисление значения χ при распределении хи-квадрат

Используются данные со стр. 104.

Шаг 1 Выберите ячейку В3.

	A	B
1	P	0.05
2	степень свободы	1
3	хи-квадрат	

Шаг 2 Выберите **Вставка** ▶ **Функция**.

Шаг 3 Выберите **Статистические** в строке **Категория**, а затем **ХИ2ОБР*** в графе **Выберите функцию**.

* CHIINV в английской версии Excel.

Шаг 4 Выберите ячейки B1 и B2 и нажмите кнопку **ОК**.

	A	B	C	D	E	F	G	H	I	
1	Р	0.05								
2	степень свободы	1								
3	хи-квадрат	(B1;B2)								
4										
5	Аргументы функции									
6	:ХИ2ОБР									
7	Вероятность		B1	= 0.05						
8	Степени_свободы		B2	= 1						
9										
10	= 3.841459149									
11	Возвращает значение обратное к односторонней вероятности распределения хи-квадрат.									
12										
13										
14										
15	Степени_свободы число степеней свободы - число от 1 до 10 ¹⁰ , исключая 10 ¹⁰ .									
16										
17										
18										
19										
20	Справка по этой функции		Значение: 3.841459149		ОК		Отмена			
21										

Шаг 5 Вычисление закончено!

	A	B
1	Р	0.05
2	степень свободы	1
3	хи-квадрат	3.841459
4		

7. Вычисление коэффициента линейной корреляции

Используются данные со стр. 116

Шаг 1 Выберите ячейку **B14**.

	A	B	C
1		Расходы на косметику, йены	Расходы на одежду, йены
2	A	3000	7000
3	B	5000	8000
4	B	12000	25000
5	Г	2000	5000
6	Д	7000	12000
7	Е	15000	30000
8	Ж	5000	10000
9	З	6000	15000
10	И	8000	20000
11	К	10000	18000
12			
13			
14	Кoeffициент линейной корреляции		
15			

Шаг 2 Выберите **Вставка** ▶ **Функция**.

Шаг 3 Выберите **Статистические** в строке **Категория**, а затем **KORPEЛ*** в графе **Выберите функцию**.

* CORREL в английской версии Excel.

Шаг 4 Выделите указанную ниже область и нажмите кнопку **OK**.

	A	B	C	D	E	F	G	H	I	J	K
1		Расходы на косметику, йены	Расходы на одежду, йены								
2		3000	7000								
3	A	5000									
4	B	12000									
5	Г	2000									
6	Д	7000									
7	Е	15000									
8	Ж	5000									
9	З	6000									
10	И	8000									
11	К	10000									
12											
13											
14	Кoeffициент линейной корреляции	:B11;C2:C11)									
15											
16											
17											
18											
19											

Аргументы функции

KORPEЛ

Массив1 {B2:B11} = {3000;5000;12000;2000;7000;15000;5000;6000;8000;10000}

Массив2 {C2:C11} = {7000;8000;25000;5000;12000;30000;10000;15000;20000;18000}

= 0.968019613

Возвращает коэффициент корреляции между двумя множествами данных.

Массив2 второй диапазон значений. Значениями могут быть числа, имена, массивы или ссылки с именами.

Справка по этой функции Значение: 0.968019613 **OK** **Отмена**

Шаг 5 Вычисление закончено!

	A	B	C
1		Расходы на косметику, йены	Расходы на одежду, йены
2	A	3000	7000
3	B	5000	8000
4	B	12000	25000
5	Г	2000	5000
6	Д	7000	12000
7	Е	15000	30000
8	Ж	5000	10000
9	З	6000	15000
10	И	8000	20000
11	К	10000	18000
12			
13			
14	Кoeffициент линейной корреляции	0.968019613	

8. Проверка гипотезы о независимости

Используются данные со стр. 157.

Шаг 1 Выберите ячейку B8.

	A	B	C	D	E
1		по телефону	по SMS	при встрече	Итого
2	ж	34	61	53	148
3	м	38	40	74	152
4	Итого	72	101	127	300
5					
6					
7		по телефону	по SMS	при встрече	
8	ж				
9	м				
10					
11					
12	P-значение				

Шаг 2 Напишите в ячейке B8 точно так же, как пишете в Word, выражение: $=E2*B4/E4$ и пока не нажимайте клавишу Enter.

	A	B	C	D	E
1		по телефону	по SMS	при встрече	Итого
2	ж	34	61	53	148
3	м	38	40	74	152
4	Итого	72	101	127	300
5					
6					
7		по телефону	по SMS	при встрече	
8	ж	$=E2*B4/E4$			
9	м				
10					
11					
12	P-значение				

Шаг 3

Подведите курсор к надписи E2 в ячейке B8 и трижды нажмите кнопку F4. Убедитесь, что E2 превратилось в \$E2, и пока не нажимайте клавишу Enter.

	A	B	C	D	E
1		по телефону	по SMS	при встрече	Итого
2	ж	34	61	53	148
3	м	38	40	74	152
4	Итого	72	101	127	300
5					
6					
7		по телефону	по SMS	при встрече	
8	ж	=E2			
9	м				
10					
11					
12	Р-значение				

Шаг 4

Подведите курсор к надписи B4 в ячейке B8, дважды нажмите клавишу F4 и убедитесь, что B4 превратилось в B\$4. Затем подведите курсор к надписи E4 в ячейке B8, нажмите клавишу F4 и проверьте, превратилось ли E4 в \$E\$4.

	A	B	C	D	E
1		по телефону	по SMS	при встрече	Итого
2	ж	34	61	53	148
3	м	38	40	74	152
4	Итого	72	101	127	300
5					
6					
7		по телефону	по SMS	при встрече	
8	ж	=E2*B\$4/\$E\$4			
9	м				
10					
11					
12	Р-значение				

После этого нажмите клавишу Enter.

Шаг 5

Выберите ячейку B8, подведите указатель мыши к правому углу ячейки B8 и как только указатель примет вид черного крестика нажмите клавишу мыши. Удерживая клавишу мыши в нажатом состоянии, растяните область до ячейки D8, а затем отпустите клавишу мыши.

	A	B	C	D	E
1		по телефону	по SMS	при встрече	Итого
2	ж	34	61	53	148
3	м	38	40	74	152
4	Итого	72	101	127	300
5					
6					
7		по телефону	по SMS	при встрече	
8	ж	35,52			
9	м				
10					
11					
12	Р-значение				

- Шаг 6** Выделите ячейки B8–D8, подведите указатель мыши к правому углу ячейки D8 и как только он превратится в черный крестик нажмите клавишу мыши. Удерживая клавишу мыши в нажатом состоянии, растяните область до ячейки D9, а затем отпустите клавишу мыши.

	A	B	C	D	E
1		по телефону	по SMS	при встрече	Итого
2	ж	34	61	53	148
3	м	38	40	74	152
4	Итого	72	101	127	300
5					
6					
7		по телефону	по SMS	при встрече	
8	ж	35,52	49,82666667	62,65333333	
9	м				
10					
11					
12	R-значение				

- Шаг 7** Выберите ячейку B12. После этого выберите **Вставка** ▶ **Функция**, а затем **Статистические** в строке **Категория** и **ХИ2ТЕСТ*** в графе **Выберите функцию**.

* СНИТЕСТ в английской версии Excel.

	A	B	C	D	E	F	G	H
1		по телефону	по SMS	при встрече	Итого			
2	ж	34						
3	м	38						
4	Итого	72						
5								
6								
7		по телефону						
8	ж	35,52						
9	м	36,48						
10								
11								
12	R-значение	=						
13								
14								
15								
16								
17								
18								
19								
20								
21								
22								

Мастер функций - шаг 1 из 2

Поиск функции:

Введите краткое описание действия, которое нужно выполнить, и нажмите кнопку "Найти"

Найти

Категория: **Статистические**

Выберите функцию:

- ФИШЕРОБР
- ФТЕСТ
- ХИ2ОБР
- ХИ2РАСП
- ХИ2ТЕСТ**
- ЧАСТОТА
- ЭКСПРАСП

ХИ2ТЕСТ(фактический_интервал;ожидаемый_интервал)

Возвращает тест на независимость: значение распределения Хи-квадрат для статистического распределения и соответствующего числа степеней свободы.

Справка по этой функции

OK Отмена

Шаг 8 Выделите указанную ниже область и нажмите кнопку **OK**.

	A	B	C	D	E	F	G	H	I	J	K
1		по телефону	по SMS	при встрече	Итого						
2	ж	34									
3	м	38									
4	Итого	72									
5											
6											
7		по телефону									
8	ж	35.52									
9	м	36.48									
10											
11											
12	P-значение	D3;B8:D9)									
13											
14											
15											
16											
17											
18											
19											

Аргументы функции x

ХИ2ТЕСТ

Фактический_интервал B2:D3 = {34;61;53;38;40;74}

Ожидаемый_интервал B8:D9 = {35.52;49.82666667}

= 0.01823258

Возвращает тест на независимость: значение распределения Хи-квадрат для статистического распределения и соответствующего числа степеней свободы.

Ожидаемый_интервал диапазон, содержащий отношение произведений итогов по строкам и столбцам к общему итогу.

Справка по этой функции Значение: 0.01823258

Шаг 9 Вычисление закончено!
 Проверьте, соответствует ли полученная величина P-значению на стр. 177).

	A	B	C	D	E
1		по телефону	по SMS	при встрече	Итого
2	ж	34	61	53	148
3	м	38	40	74	152
4	Итого	72	101	127	300
5					
6					
7		по телефону	по SMS	при встрече	
8	ж	35.52	49.82666667	62.65333333	
9	м	36.48	51.17333333	64.34666667	
10					
11					
12	P-значение	0.01823258			

Предметный указатель

α

α, 163, 166

A

AVERAGE, функция, 196

C

CHIDIST, функция, 107

CHINV, функция, 107, 205-206

CHITEST, функция, 210-211

CORREL, функция, 207

COUNTIF, функция, 197-198

E

Eiken, экзамен, 23-25

F

FDIST, функция, 107

FINV, функция, 107

FREQUENCY, функция, 193-194

FRACП, функция, 107

FRACПОБР, функция, 107

F-распределение, 106-107

M

Microsoft Excel, см. *вычисления в Excel, функции Excel*

N

NORMDIST, функция, 107

NORMINV, функция, 107

NORMSDIST, функция, 107, 204

NORMSINV, функция, 107

P

P-значение

альтернативная гипотеза, 175-179

гипотеза о независимости, 175

нулевая гипотеза, 175-179

проверка гипотезы, 163, 175-179, 189

S

STANDARDIZE, функция, 199-201

T

TDIST, функция, 107

TINV, функция, 107

T-показатель, см. *рейтинг успеваемости*

V

V Крамера, см. *Крамера коэффициент корреляции*

Z

Z-показатель, см. *нормированное отклонение*

Z-преобразование, см. *нормирование*

A

альтернативная гипотеза

P-значение, 175-179

коэффициент корреляции Крамера, 186

обзор, 170-174

определение, 174

примеры, 161, 171-173

проверка гипотезы о равенстве долей в совокупностях, 173

анкетирование, 4-6

качественные данные, 60-64

ограничения, 4-7

проверка независимости, 137, 208-211

таблицы распределения, 62-64

B

вероятность, 81-109

F-распределение, 106-107

нормальное распределение, 86-89

определение, 82

распределение и Excel, 107-109

распределение Стюдента, 106

распределение хи-квадрат, 99-105, 205-206

результаты теста, 83-84

соответствующая, 104

стандартное нормальное распределение, 89-98, 204-205

степень свободы, 99-108

выборка, см. *выборочная совокупность*

выборочная совокупность, 6, 7, 52, 57

вычисления в Excel, 191-211

коэффициент корреляции, 206-207

медиана, 195-196

нормированное отклонение, 74-80, 199-203

проверка гипотезы о независимости, 208-211

простая статистическая таблица, 197-198

распределение, 107-109

распределение хи-квадрат, 205-206

рейтинг успеваемости, 199-202

среднее значение, 195-196

стандартное нормальное распределение, 204-205

стандартное отклонение, 195-196

таблицы распределения, 192-195

G

генеральная совокупность

выборка, 52

коэффициент корреляции Крамера, 145-150, 157, 186

определение, 6

проверка гипотезы, 149, 186

стандартное отклонение, 52

гистограммы
величина интервала, 84, 85
достоинства, 83
обзор, 38-39
переменные, 39
примеры, 39, 83, 84, 154
функция плотности вероятности, 83-84

Д

данные
измеряемые, см. *количественные данные*
качественные, см. *качественные данные*
количественные, см. *количественные данные*
неизмеряемые, см. *качественные данные*
разброс, 49, 58, 69, 70, 80
диаграмма распределения
ежемесячные расходы, 116-120
корреляционное отношение, 122, 126
примеры, 114, 116
диаграмма цены, 33-39
диаграммы
корреляционное отношение, 126
построение, 33-39
расходы, 116-120
столбиковая, 114
теснота связи, 115
точечная, см. *точечная диаграмма*
доли совокупностей, 149

З

зависимость
корреляционное отношение, 117, 121-127
линейная, 120
нелинейная, 120
относительная частота, 36-37, 39
переменные, 112-115
степень, 115, 116-120
загрузка файлов Excel, 192
значение
Р-значение, 163, 175-179, 189
медиана, 44-47
значимость гипотезы, 187

И

интервал, 39, 54-57, 84
величина, 39, 54-57, 84
формула Стерджесса, 55

К

качественные данные, 14-29
корреляционное отношение, 121
обзор, 14-19
определение, 19
показатели, 117
примеры, 20, 23-26
результат исследования, 60-64

создание таблиц, 60-64
столбиковая диаграмма, 114
точечная диаграмма, 114
количественные данные, 14-29
гистограммы, 38-39, 54, 58
корреляционное отношение, 121
медиана, 44-47
обзор, 31-58
описательная статистика, 57-58
определение, 19
показатели, 117
примеры, 21-23, 26
средняя величина, 40-43
стандартное отклонение, 48-53, 70-79
таблицы распределения, 32-39, 54-56, 58
теория оценивания, 57-58
точечная диаграмма, 114
КОРРЕЛ, функция, 207
корреляционное отношение, 117, 121-127, 207
корреляция, 115, 119
коэффициент корреляции Крамера, см. Крамера
коэффициент корреляции
коэффициент линейной корреляции, 116-120, 206-207
коэффициент независимости, см. *Крамера коэффициент корреляции*
Крамера коэффициент корреляции, 127-138
вычисление в Excel, 207
альтернативная гипотеза, 186
вычисление, 130-135, 141
критерии величины, 136
нулевая гипотеза, 168, 186
показатель тесноты связи, 117, 129
примеры, 127-136
пропорции предпочтений, 155
точность, 147
критерий согласия Пирсона, 132, 152-155, 158
символ, 103
критическая область, 159, 165-167, 187

Л

линейная зависимость, 120

М

медиана
вычисление в Excel, 195-196
определение, 45
применимость, 44
примеры, 45-47
межгрупповая дисперсия, 117, 124, 126
многовариантные ответы, 28

Н

наклон графика, 101
неизмеряемые данные, см. *качественные данные*
нелинейная зависимость, 120

Непера число, 86
нормализация, см. *нормирование*
НОРМАЛИЗАЦИЯ, функция, 199-201
нормальное распределение, 86-91
нормирование, 71-72
нормированное отклонение, 65-80, 73, 199-202
нормировка, см. *нормирование*
НОРМОБР, функция, 107
НОРМРАСП, функция, 107
НОРМСТОБР, функция, 107
НОРМСТРАСП, функция, 107, 204
нулевая гипотеза
Р-значение, 175-179
для проверки корреляционного отношения, 172
для проверки независимости, 172
для проверки отсутствия корреляции, 172
для проверки равенства долей в совокупностях, 173
для проверки равенства средних величин совокупностей, 173
коэффициент Крамера, 168, 186
нельзя отвергнуть, 150, 167, 178, 179, 187
обзор, 170-174
описание, 174
примеры, 167-174
трудность доказательства, 174

О

обратная зависимость, 119
описательная статистика, 57-58
опрос общественного мнения, 4-6
ось у, 39
ось x, 39, 102, 107, 109, 125
отсутствие корреляции, 119

П

переменные, 111-142
гистограмма, 39
зависимость, 112-115
корреляционное отношение, 121-127
коэффициент Крамера, 127-138, 141-142
коэффициент линейной корреляции, 116-120
степень связи, 115, 116-120
Пирсона критерий согласия, см. *критерий согласия Пирсона*
показатели
количественные данные, 117
коэффициент корреляции Крамера, 117, 129
коэффициент линейной корреляции, 120
проверка гипотез, 143-189
Р-значение, 163, 175-179, 189
альтернативная гипотеза, см. *альтернативная гипотеза*
виды, 149, 171
выводы, 187

критерий согласия Пирсона, 151-160
нулевая гипотеза, см. *нулевая гипотеза*
о корреляционном отношении, 149, 171, 172
о независимости, 149, 171
о равенстве долей в совокупностях, 149, 171, 173
о равенстве средних величин совокупностей, 149, 171, 173
об однородности, 184-186
об отсутствии корреляции, 149, 171, 172
обзор, 144-150
определение, 149
порядок проведения, 150, 175-179
примеры, 149, 168-174
проверка критерия согласия Пирсона, 151-169
проверка гипотезы о независимости, 151-161
Р-значение, 175
применимость, 137, 149
примеры, 149, 171, 184-186
сравнение с проверкой однородности, 186
хи-квадрат, 151-169
проверка статистических гипотез, см. *проверка гипотез*
прогноз погоды, 82
процентное отношение, 5, 37, 62, 64
прямая зависимость, 119

Р

разброс данных, 49, 58, 69, 70, 80
распределение
F, 106-107
вычисление в Excel, 107-109
нормальное, 86-91
стандартное нормальное, 89-98, 204-205
Стьюдента, 106
хи-квадрат, см. *хи-квадрат распределение*
рассеяние данных, см. *разброс данных*
результаты тестов
нормальное распределение, 86-89
стандартное нормальное распределение, 89-98
функция плотности вероятности, 83-84
рейтинг успеваемости, 74-80, 199-203
ряды распределения, см. *таблицы распределения*

С

свободы степень, 99-108
середина интервала, 36-39, 54, 56
совокупность, см. *генеральная совокупность*
среднее значение, см. *средняя*
средние накопления, 46-47
средняя
арифметическая, 43, 73, 74
вычисление в Excel, 195-196
гармоническая, 43
геометрическая, 43
нормальное распределение, 87-89

определение, 43
примеры, 40-44
стандартное нормальное распределение, 89-90
СРЗНАЧ, функция, 196
стандартизация, см. *нормирование*
стандартное нормальное распределение, 89-98,
204-205
стандартное отклонение, 48-53, 70-79
вычисление в Excel, 195-196
количественные данные, 48-53, 70-79
нормальное распределение, 87-91
совокупность, 52
стандартное нормальное распределение, 89-90
статистика
описательная, 57-58
определение, 4
теория оценивания, 4-6
степень свободы, 99-108
степень тесноты связи, 115, 116-120
Стерджесса формула, 55
Стьюдента распределение, 106
СТБЮДРАСП, функция, 107
СТБЮДРАСПРОБР, функция, 107
СЧЕТЕСЛИ, функция, 197-198

Т
таблица сопряжённости, 128, 130, 135, 151, 153,
197-198
таблицы распределения, 54-56, 192-195
качественные данные, 60-64
нормальное распределение, 107
распределение хи-квадрат, 102-105, 205-206
стандартное нормальное распределение, 92-93,
104, 108
таблицы сопряжённости, 128, 130, 135, 151, 153
частота, 54-56
теоретическая частота, 130, 131
теория оценивания, 57-58
типы данных, 13-29, 117
точечная диаграмма, 116, 119, 120

У
уровень значимости (α), 159, 163

Ф
функции Excel
AVERAGE, 196
CHIDIST, 107
CHINV, 107, 205-206
CHITEST, 210-211
CORREL, 207
COUNTIF, 197-198

FDIST, 107
FINV, 107
FREQUENCY, 193-194
FRАСП, 107
FRАСПОБР, 107
NORMDIST, 107
NORMINV, 107
NORMSDIST, 107, 204
NORMSINV, 107
STANDARDIZE, 199-201
TDIST, 107
TINV, 107
КОРРЕЛ, 207
НОРМАЛИЗАЦИЯ, 199-201
НОРМОБР, 107
НОРМРАСП, 107
НОРМСТОБР, 107
НОРМСТРАСП, 107, 204
СРЗНАЧ, 196
СТБЮДРАСП, 107
СТБЮДРАСПРОБР, 107
СЧЕТЕСЛИ, 197-198
ХИ2РАСП, 107
ХИ2РАСПОБР, 107, 205-206
ХИ2ТЕСТ, 210-211
ЧАСТОТА, 193-194
функция распределения плотности вероятности, 82-
85, 99, 107, 109

Х
ХИ2РАСП, функция, 107
ХИ2РАСПОБР, функция, 107, 205-206
ХИ2ТЕСТ, функция, 210-211
хи-квадрат, распределение, 99-105
вычисление, 130-133
вычисление χ , 205-206
описание, 99
примеры, 99-105, 152
степень свободы, 99-108

Ч
частота
описание, 36
относительная, 36-37, 39
таблицы распределения, 32-39
теоретическая, 130, 131
эмпирическая, 130, 131
ЧАСТОТА, функция, 193-194

Э
Эйлера число, 86
эмпирическая частота, 130, 131

Книги Издательского дома «Додэка-XXI» можно заказать в торговом-издательском холдинге «АЛЬЯНС-КНИГА» наложенным платежом, выслать открытку или письмо по почтовому адресу: **123242, Москва, а/я 20** или по электронному адресу: **orders@alians-kniga.ru**.

При оформлении заказа следует указать адрес (полностью), по которому должны быть высланы книги; фамилию, имя и отчество получателя. Желательно также указать свой телефон и электронный адрес.

Эти книги вы можете заказать и в Интернет-магазине: **www.alians-kniga.ru**.

Оптовые закупки: тел. **(495) 258-91-94, 258-91-95**; электронный адрес
books@alians-kniga.ru.

Син Такахаси

Занимательная статистика Манга

Подписано в печать 15.12.2009. Формат 70х90/16. Бумага офсетная.
Гарнитура «LiteraturnayaC», «JacobC». Печать офсетная.
Объем 14,0 п. л. Усл. п. л. 16,3.
Тираж 1500 экз. Код ОНМ03.

Издательский дом «Додэка-XXI»

105318 Москва, а/я 70
Тел./факс: (495) 366-04-56, 366-11-55
E-mail: red@dodeca.ru
Web-сайт издательства: www.dodeca.ru
Интернет-магазин: www.alians-kniga.ru

Отпечатано с готовых диапозитивов в ОАО «Щербинская типография»
117623, Москва, ул. Типографская, д. 10.