

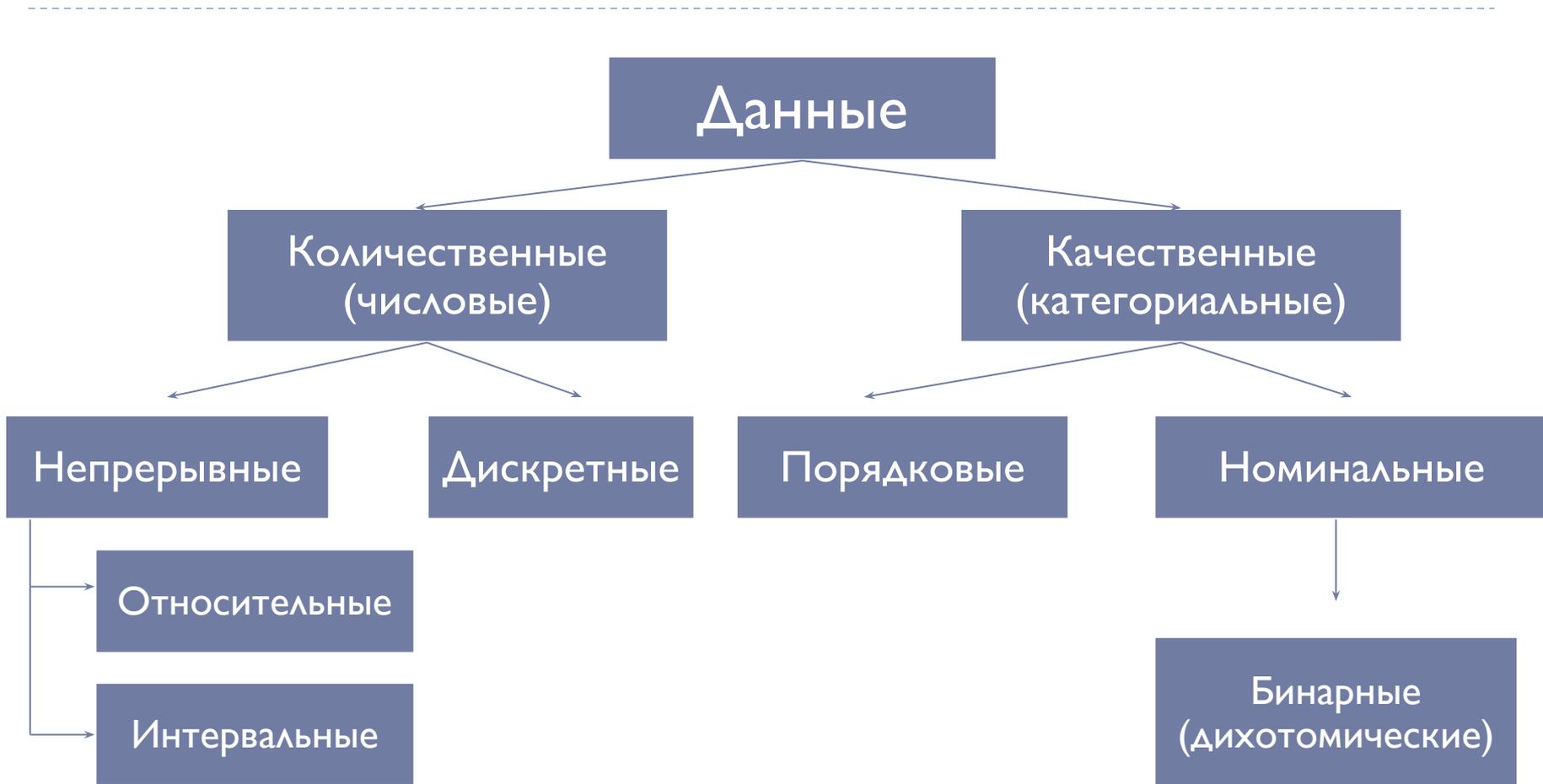
Статистика в медико-биологических исследованиях

К.м.н., доц. Хисамутдинов А.Н.



- Каждое решение врача должно основываться на научных данных
- статистические методы - ключевой, решающий инструмент, который позволяет качественно или количественно доказать, обосновать или опровергнуть новую научную идею и мысль





Количественные (числовые) данные

- **Непрерывные** – данные, которые получают при измерении на непрерывной шкале, т.е. теоретически они могут иметь дробную часть. Примеры: масса тела, рост, артериальное давление.
- **Интервальные данные** – вид непрерывных данных, которые измеряются в абсолютных величинах, имеющих физический смысл (шкала IQ, температура в градусах Цельсия, Фаренгейта)
- **Относительные данные (наличие абсолютной нулевой точки)** – вид непрерывных данных, отражающих долю изменения значения признака по отношению к исходному (или какому-либо другому) значению признака (доза препарата, возраст, абсолютная температура).
- **Дискретные данные** – количественные данные, которые не могут иметь дробную часть (количество детей).

Качественные (категориальные) данные

- **Номинальные (шкалы наименований)** – вид качественных данных, которые отражают условные коды неизмеримых категорий, когда отдельным числам не соответствует никакого эмпирического значения (пол, семейное положение, коды диагноза)
 - **Бинарные (дихотомические) данные** – особо выделяемый вид качественных данных, когда признак имеет два возможных значения (пол, наличие/отсутствие заболевания)
 - **Порядковые** – вид качественных данных, которые отражают условную степень выраженности какого-либо признака (например стадии заболевания, степени сердечной недостаточности)
-

Важнейшие понятия



□ **Генеральная совокупность:**

все множество данных. Пример: если целью исследования является изучение уровня гемоглобина населения Земли, генеральная совокупность – значения уровня гемоглобина в крови каждого жителя земного шара

□ **Выборочная совокупность (выборка):**

часть данных, отобранная из генеральной совокупности

Цель формирования выборки: получить **оценку** некоторого изучаемого параметра генеральной совокупности, не перебирая все данные по всей генеральной совокупности

Описательные статистики

- **Минимум и максимум** – минимальное и максимальное значения переменной в совокупности
- **Размах** – разница между максимальным и минимальным значением (обозначение R)
- **Среднее** – сумма значений переменной, деленное на число значений переменной
- **Дисперсия** – (от англ. variance) и **стандартное (среднеквадратическое) отклонение** (англ. standard deviation) – меры изменчивости переменной
- **Коэффициент вариации** – мера относительного разброса случайной величины; показывает, какую долю среднего значения этой величины составляет ее средний разброс

Описательные статистики (продолжение)

- **Медиана** – разбивает выборку на две равные части. Половина значений переменной лежит ниже медианы, половина - выше
- **Квартили** представляют собой значения, которые делят две половины выборки (разбитые медианой) еще раз пополам
- **Процентили** – величины, которые делят упорядоченные наблюдения на 100 равных частей
- **Мода** представляет собой максимально часто встречающееся значение переменной (наиболее «модное» значение переменной)

| | | | | | | | | | | | | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 165 | 166 | 167 | 168 | 169 | 170 | 171 | 172 | 173 | 174 | 175 | 176 | 177 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

Среднее

- ▶ Пусть имеется переменная X , тогда оценка среднего, или **выборочное среднее**, вычисляется как среднее арифметическое наблюдаемых значений. Выборочное среднее обычно обозначается $\bar{X} = \frac{\sum X}{n}$ (M)
- ▶ Выборочное среднее не устойчиво к выбросам. Пример: среднее чисел: **1,2,3,4,5,6,7** составляет 4, если к ним добавить 62, то среднее будет составлять 16!!! (необходимо избегать использование выборочного среднего, не обработав выбросы в данных)

Дисперсия и стандартное (среднеквадратическое) отклонение

- меры изменчивости переменной
- чем сильнее разбросаны значения переменной относительно среднего, тем больше дисперсия и стандартное отклонение



Дисперсия и стандартное отклонение

- ▶ Выборочная дисперсия переменной X вычисляется по формуле:

$$\bar{\sigma}^2 = \frac{\sum(X - \bar{X})^2}{n}$$

- ▶ Стандартное отклонение равно квадратному корню из выборочной дисперсии:

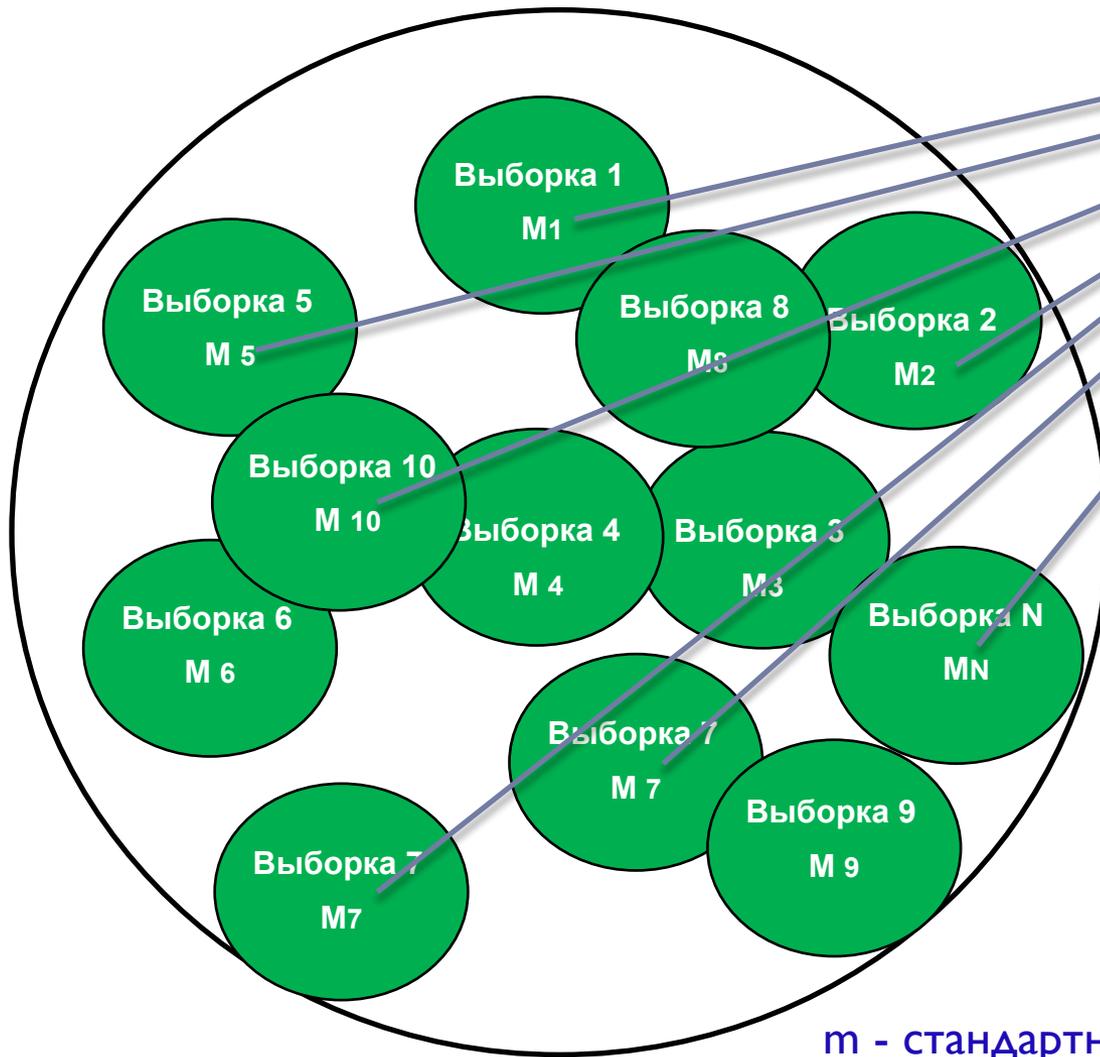
$$\sigma = \sqrt{\bar{\sigma}^2} = \sqrt{\frac{\sum(X - \bar{X})^2}{n}}$$

Стандартная ошибка среднего (ошибка репрезентативности) (Standard Error of Mean, SEM)

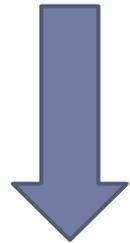
- Представляет собой стандартное отклонение распределения средних отдельных выборок (рассчитанное из средних отдельных выборок)

$$m \text{ или } SEM = \frac{\sigma}{\sqrt{n}}$$

и отображает точность оцененного параметра среднего



\bar{X}



$$\sigma = \sqrt{\frac{\sum(X - \bar{X})^2}{n}}$$

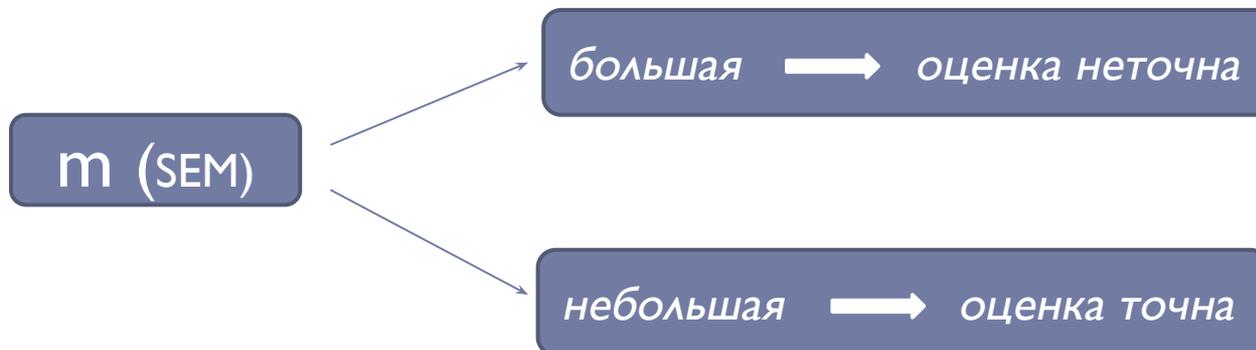
=

SEM, m

m - стандартное отклонение рассчитанное из средних отдельных выборок

Стандартная ошибка (SEM) или стандартное отклонение (σ)?

- несмотря на внешнюю схожесть, параметры SEM и σ используют в разных целях:
 - **стандартное отклонение** отражает разброс значений данных и должно быть указано, если необходимо описать выборку и пояснить изменчивость в наборе данных
 - **стандартная ошибка среднего** отражает точность оцененного параметра среднего



Доверительный интервал для среднего



- ▶ Наряду с выборочным средним, которое является точечной оценкой параметра, часто приводят интервальные оценки.

$$[\bar{X} - (t_{0,05} \times SEM); \bar{X} + (t_{0,05} \times SEM)]$$

$$t_{0,05} = 1,96$$

$$t_{0,01} = 2,58$$

- ▶ Доверительный интервал – диапазон значений, внутри которого находится средняя популяции (с вероятностью 95%, 99%)



Как правильно описать выборочную совокупность?

М Мах Me

М ДИ σ^2

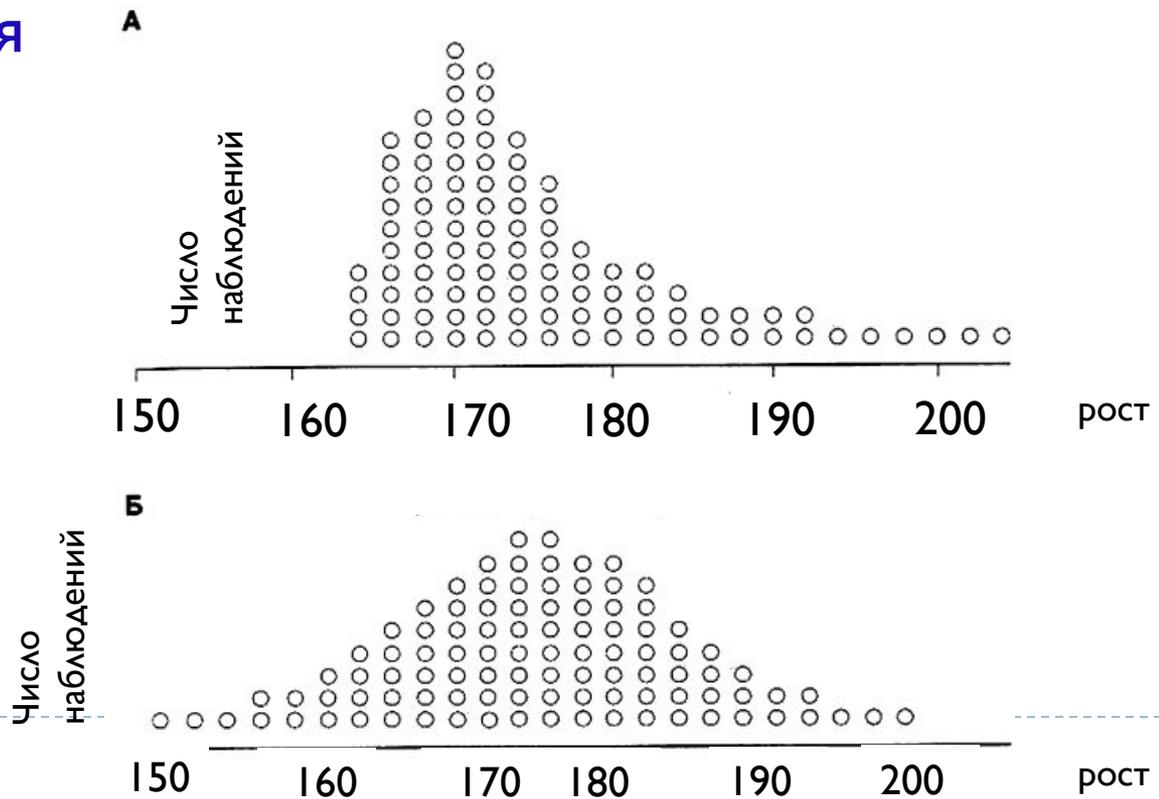
σ m Min Mo

Какие описательные статистики использовать?

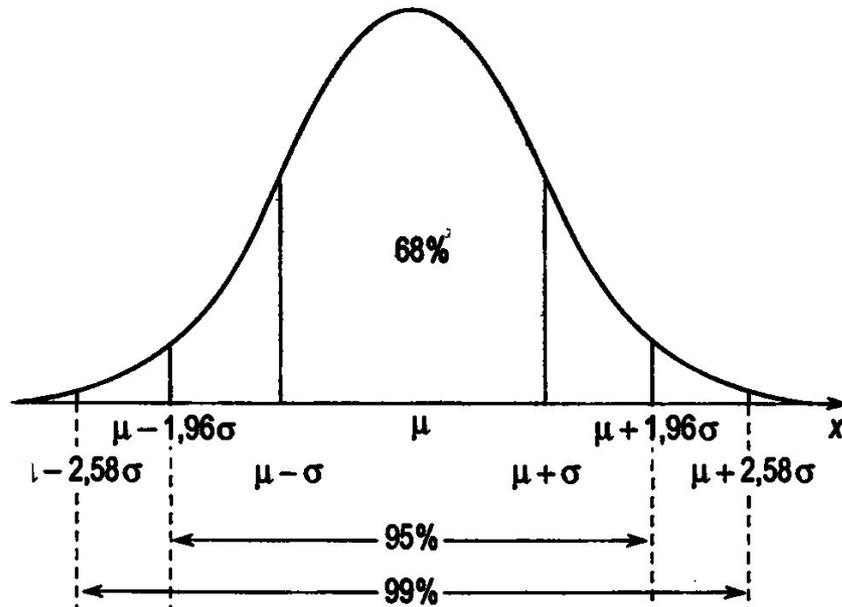


Нормальное распределение

- Для того, чтобы выбрать описательные статистики для совокупности, сначала следует установить, соответствует ли вид распределения значений изучаемого признака закону нормального распределения



Свойства нормального распределения



Нормальное распределение **полностью** определяется средней и стандартным отклонением:
 $\bar{X} \pm \sigma$

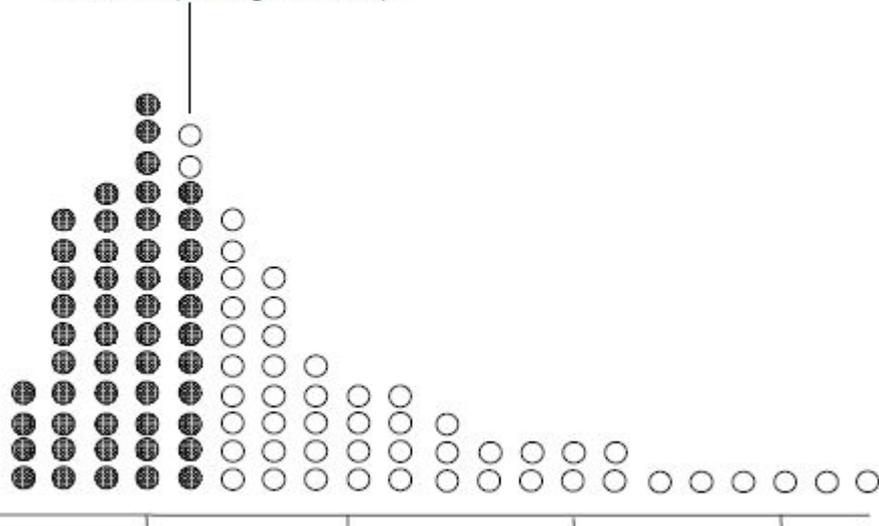
.....и для описания выборочных совокупностей, имеющих нормальное распределение

(и только таких признаков!!!), следует использовать среднее (M) и стандартное отклонение (σ) **в формате $M \pm \sigma$.**

Н.В! Международные научные журналы в качестве описательных статистик нормально распределенных совокупностей используют формат **M**
 (σ)

Если переменная не соответствует закону нормального распределения ...

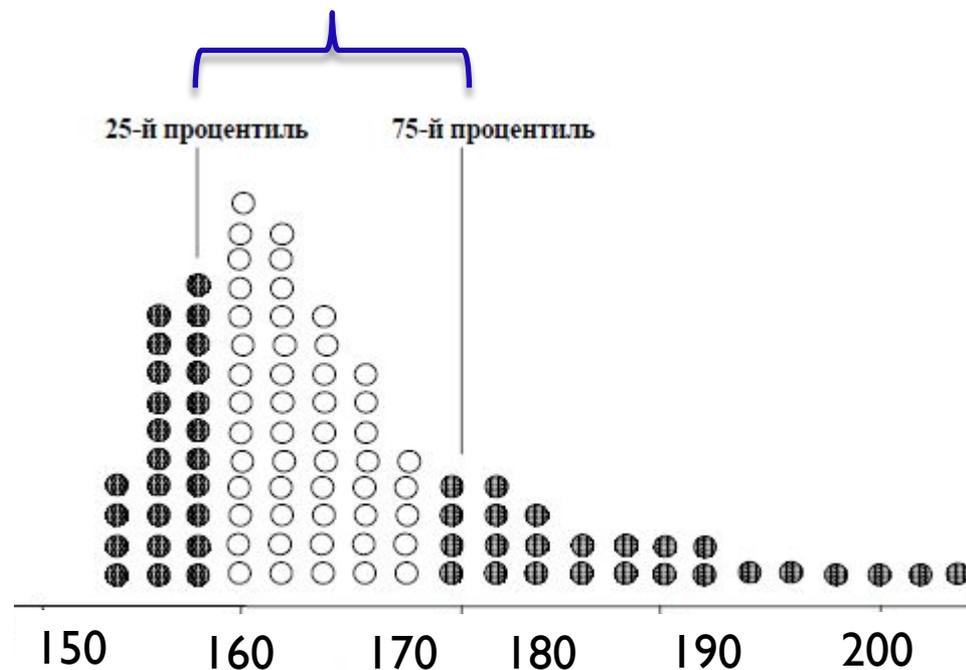
Медиана (50-й процентиль)



50% наблюдений

25-й процентиль

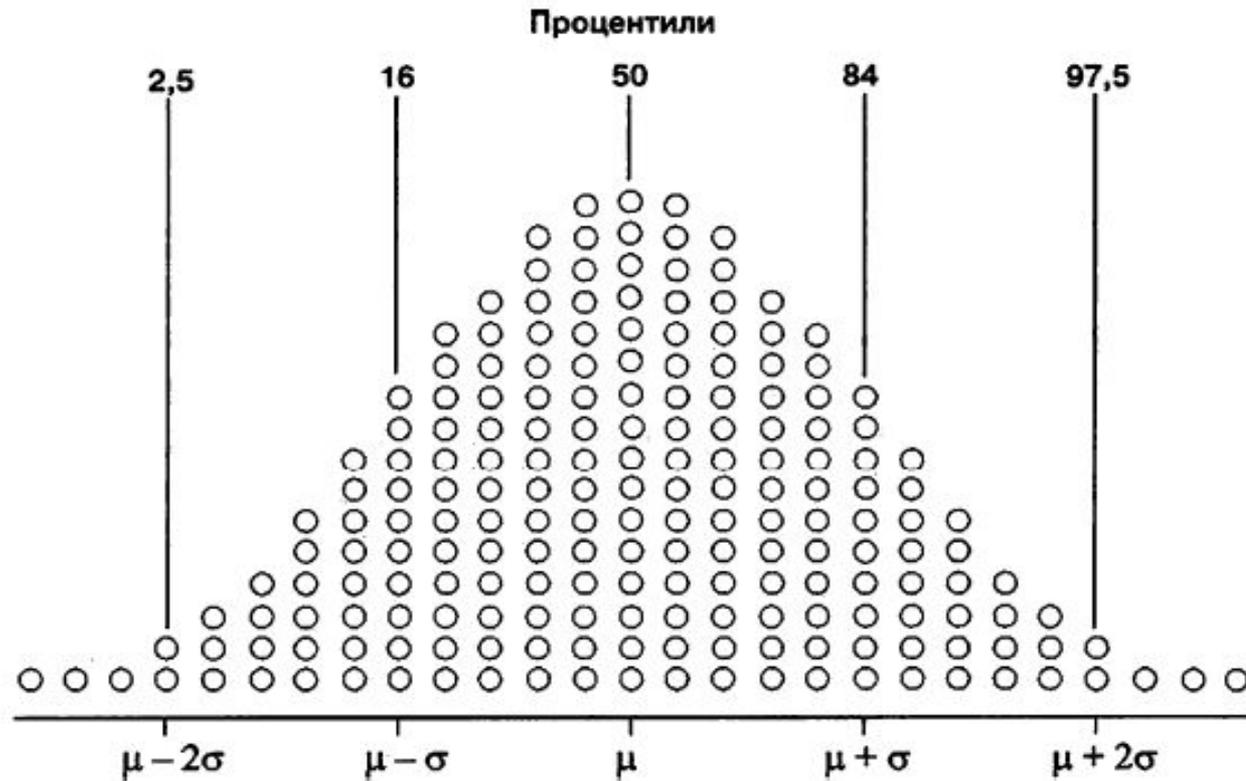
75-й процентиль



□ ...совокупность описывается:

Me [квартиль 1; квартиль 3]

Свойства нормального распределения



Среднее и стандартное отклонение

- Среднее и медиана нормального распределения равны

Важно! Отличия в описательном анализе различных типов данных

- Количественные данные + нормальное распределение:

среднее \pm стандартное отклонение

- Количественные данные + распределение, отличное от нормального:

медиана [1 и 3 квартили]

- Порядковая / номинальная шкала:

таблица частот

Важно! В медико-биологических исследованиях:

- нормальное распределение $\approx 20\%$
- распределение, отличное от нормального $\approx 80\%$

Важно!

- Возможности обработки переменных, относящихся к номинальной шкале очень ограничены: возможен только частотный анализ таких переменных и в некоторых ситуациях для дихотомических переменных - рассчитать ранговую корреляцию, а рассчитать среднее значение для переменной «Семейное положение», совершенно бессмысленно.
- Как правило, переменные, относящиеся к номинальной шкале часто используются для группировки, с помощью которых совокупная выборка разбивается по категориям этих переменных. В частичных выборках проводятся одинаковые статистические тесты, результаты которых затем сравниваются друг с другом.

Важно!

- Переменные с порядковой шкалой, кроме частотного анализа, допускают также вычисление определенных статистических характеристик, таких как медианы. Если должна быть установлена связь (корреляция) с другими переменными такого рода, для этой цели можно использовать коэффициент ранговой корреляции

Точность представления описательных статистик количественных данных

- Принято приводить оценки параметров (M , σ , m , M_e ...) с той же точностью, с которой были представлены исходные данные
- Пример: если АД измерялось с точностью до разряда единицы, то следует приводить параметры , **не в виде $145,36 \pm 27,458$ мм.рт.ст**, а в виде **145 ± 27 мм.рт.ст.**

Этапы анализа данных



Этапы анализа данных

- Планирование исследования
- Сбор информации и формирование базы данных
- Чистка данных
- Описательный и визуальный анализ
- Группировка
- Вычисление статистик для групп
- Нахождение связей и зависимостей
- Построение математических уравнений для прогноза
- Верификация (кросс проверка) уравнений для прогноза



Начало и конец. Кто неправильно застегнул первую пуговицу, уже не застегнется как следует

Формирование базы данных

Переменные



Наблюдения



| | A | B | C | D | E | F | G | H |
|----|------------|-----|---------|------------|------------|-----------|------------|--------|
| 1 | | Пол | Возраст | Эритроциты | Тромбоциты | Лейкоциты | Гемоглобин | Группа |
| 2 | пациент 1 | жен | 64 | 3,4 | 198 | 1,1 | 61 | I |
| 3 | пациент 2 | муж | 39 | 3,9 | 141 | 7,9 | 130 | II |
| 4 | пациент 3 | жен | 45 | 3,9 | 141 | 6,3 | 106 | I |
| 5 | пациент 4 | жен | 33 | 4,0 | 175 | 7,9 | 137 | II |
| 6 | пациент 5 | муж | 63 | 3,9 | 141 | 7,1 | 106 | II |
| 7 | пациент 6 | жен | 57 | 3,3 | 209 | 1,1 | 65 | I |
| 8 | пациент 7 | жен | 56 | 4,1 | 175 | 7,1 | 133 | II |
| 9 | пациент 8 | жен | 83 | 3,8 | 73 | 6,1 | 106 | I |
| 10 | пациент 9 | жен | 50 | 3,7 | 186 | 6,6 | 116 | I |
| 11 | пациент 10 | жен | 54 | 3,3 | 232 | 0,4 | 61 | I |
| 12 | пациент 11 | жен | 55 | 3,8 | 164 | 6,3 | 102 | I |
| 13 | пациент 12 | муж | 39 | 3,8 | 130 | 7,1 | 109 | I |
| 14 | пациент 13 | муж | 41 | 3,9 | 164 | 7,2 | 123 | II |
| 15 | пациент 14 | жен | 38 | 3,6 | 107 | 5,1 | 92 | I |
| 16 | пациент 15 | жен | 50 | 3,9 | 164 | 7,7 | 116 | II |
| 17 | пациент 16 | жен | 52 | 3,7 | 130 | 7,1 | 120 | II |
| 18 | пациент 17 | жен | 20 | 4,0 | 186 | 8,4 | 133 | II |
| 19 | пациент 18 | жен | 50 | 3,5 | 198 | 1,5 | 72 | I |
| 20 | пациент 19 | жен | 57 | 3,7 | 141 | 6,1 | 99 | I |
| 21 | пациент 20 | жен | 57 | 3,8 | 209 | 7,6 | 133 | II |
| 22 | пациент 21 | жен | 81 | 4,3 | 255 | 9,7 | 130 | II |
| 23 | пациент 22 | муж | 76 | 3,9 | 198 | 6,4 | 109 | I |
| 24 | пациент 23 | жен | 70 | 4,0 | 198 | 8,0 | 140 | II |

Чистка данных

- Обработка пропусков
- Поиск некорректных показателей
- Поиск выбросов
- Удаление повторных наблюдений
- Верификация текстовых меток
- Проверка диапазонов

Пример: исследование препаратов, влияющих на

| | A | B | C | D | E | F | G | H |
|----|------------|---------|---------|------------|------------|-----------|------------|--------|
| 1 | | Пол | Возраст | Эритроциты | Тромбоциты | Лейкоциты | Гемоглобин | Группа |
| 2 | пациент 1 | жен | 64 | 3,4 | 198 | 1,1 | 61 | I |
| 3 | пациент 2 | муж | 39 | 3,9 | 141 | 7,9 | 130 | II |
| 4 | пациент 3 | ж | 45 | 3,9 | 141 | 6,3 | 106 | I |
| 5 | пациент 4 | жен | 33 | 4,0 | 175 | 7,9 | 137 | II |
| 6 | пациент 5 | муж | 63 | 3,9 | 141 | 7,1 | 106 | II |
| 7 | пациент 6 | жен | 57 | 3,3 | 2090 | 1,1 | 65 | I |
| 8 | пациент 7 | жен | 56 | 4,1 | 175 | 7,1 | 133 | II |
| 9 | пациент 8 | же | 83 | 3,8 | 73 | 6,1 | 106 | I |
| 10 | пациент 9 | жен | 500 | 0,3,7 | 186 | 6,6 | 116 | I |
| 11 | пациент 10 | жен | 54 | 3,3 | 232 | 0,4 | 61 | I |
| 12 | пациент 11 | жен | 55 | 3,8 | 164 | 6,3 | 102 | I |
| 13 | пациент 12 | мужчина | 39 | 3,8 | 130 | 7,1 | 109 | I |
| 14 | пациент 13 | муж | 41 | 3,9 | 164 | 7,2 | 123 | II |
| 15 | пациент 14 | жен | 38 | 3,6 | 107 | 5,1 | 92 | I |
| 16 | пациент 15 | жен | 50 | 3,9 | 164 | 7,7 | 116 | II |
| 17 | пациент 16 | жен | 52 | 3,7 | 130 | 7,1 | 120 | II |
| 18 | пациент 17 | жен | 20 | 4,0 | 186 | 8,4 | 133 | II |
| 19 | пациент 18 | жен | 50 | 3,5 | 198 | 1,5 | 72 | I |
| 20 | пациент 19 | жен | 57 | 3,7 | 141 | 6,1 | 99 | I |

Визуальный анализ

...сначала данные нужно увидеть...



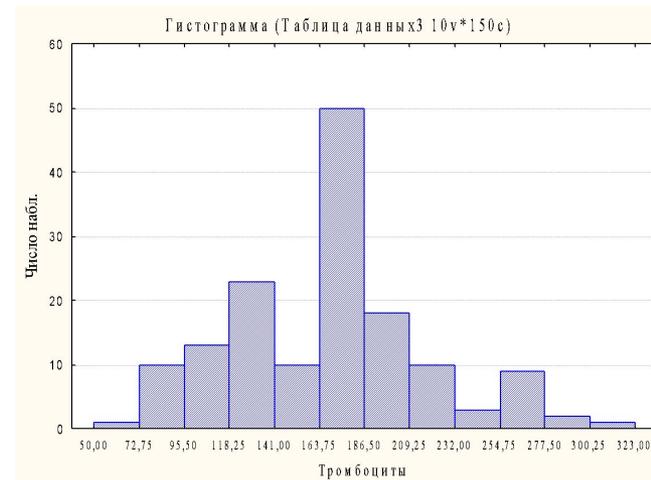
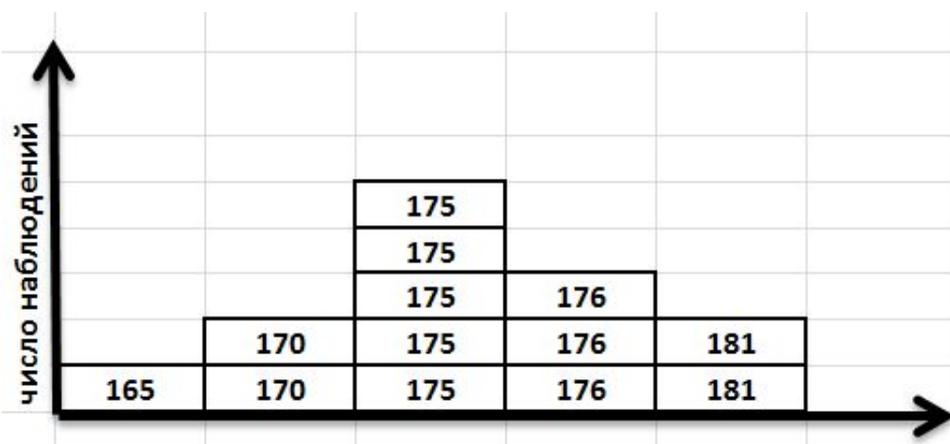
Типы графиков, наиболее часто используемые при статистическом анализе

- Гистограмма
- График средних с ошибками
- Диаграмма размаха
- Диаграмма рассеяния

Гистограмма (frequency plot, histogram, bar chart)

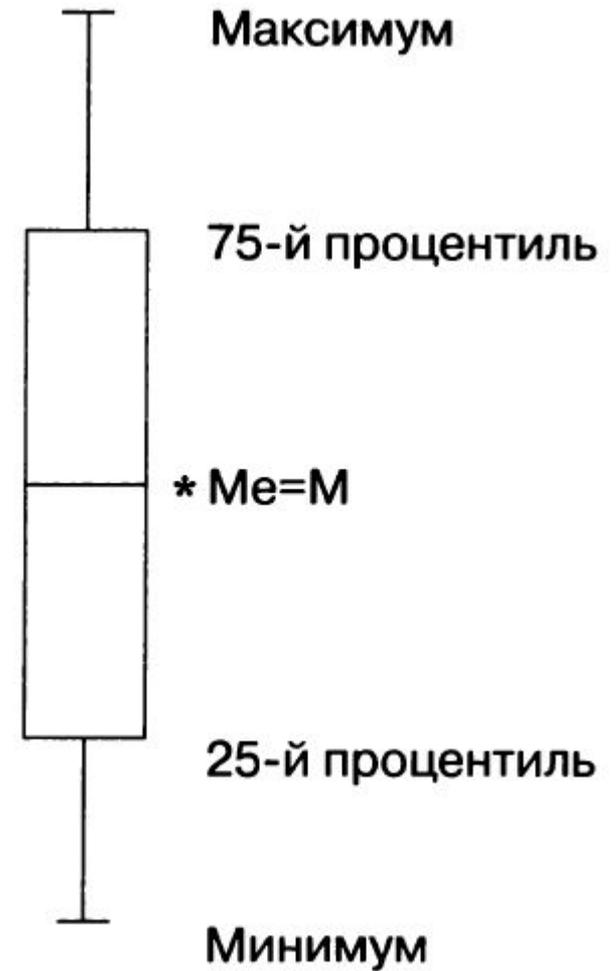
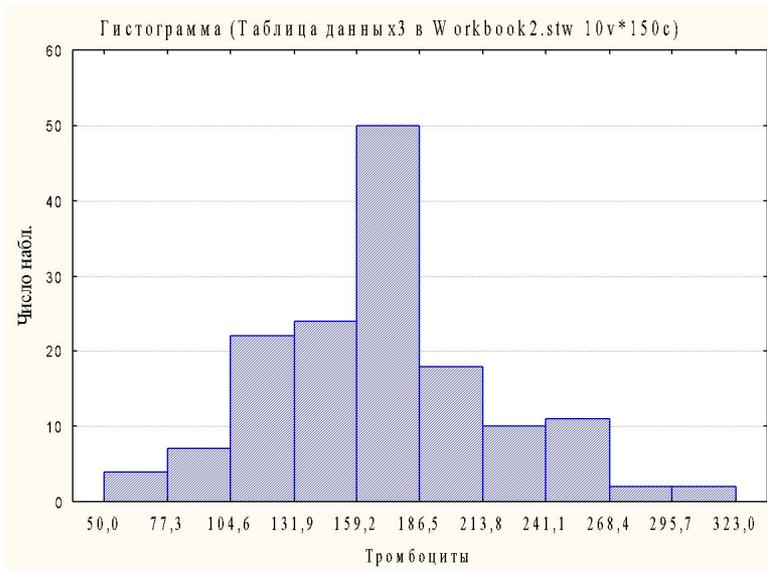
| | | | | | | | | | | | | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 165 | 170 | 170 | 175 | 175 | 175 | 175 | 175 | 176 | 176 | 176 | 181 | 181 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

| | | | | | |
|-----------------------|-----|-----|-----|-----|-----|
| значение признака, см | 165 | 170 | 172 | 176 | 181 |
| число наблюдений | 1 | 2 | 5 | 3 | 2 |



□ Визуальный анализ распределения признака

Диаграмма размаха



□ ...в описательной статистике

при оценке статистической значимости различий

График средних с ошибками

График средних (Таблица данных3 в Workbook2.stw 10v *150с)

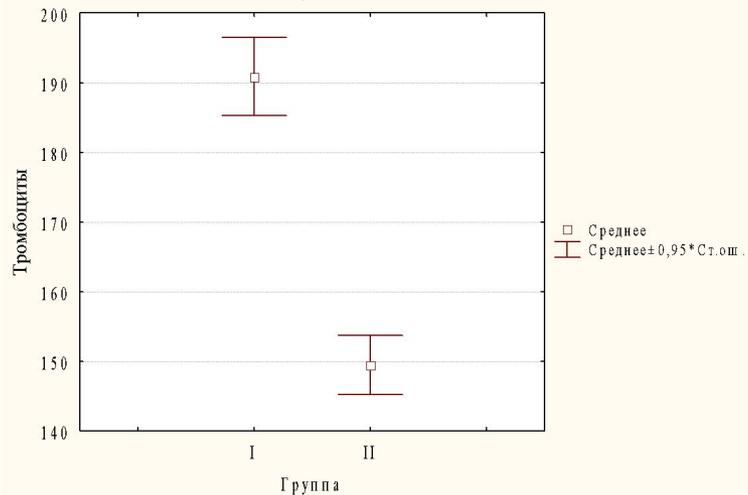


График средних (Таблица данных3 в Workbook2.stw 10v *150с)

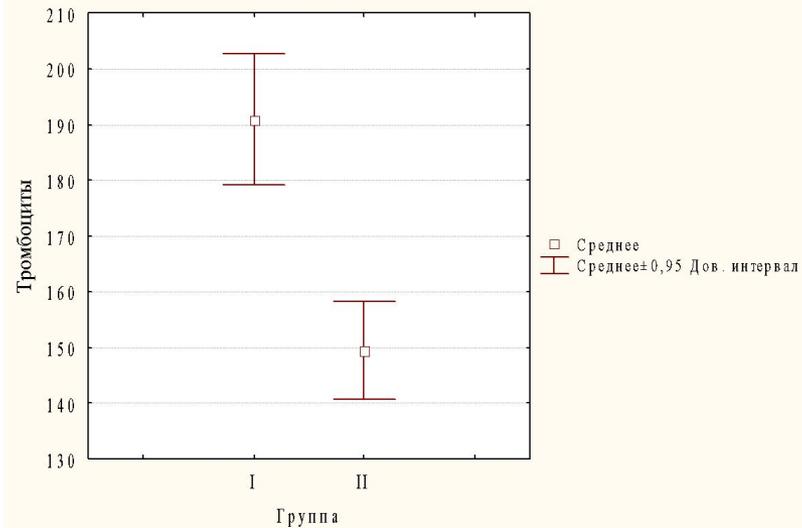


График средних (Таблица данных3 в Workbook2.stw 12v *150с)

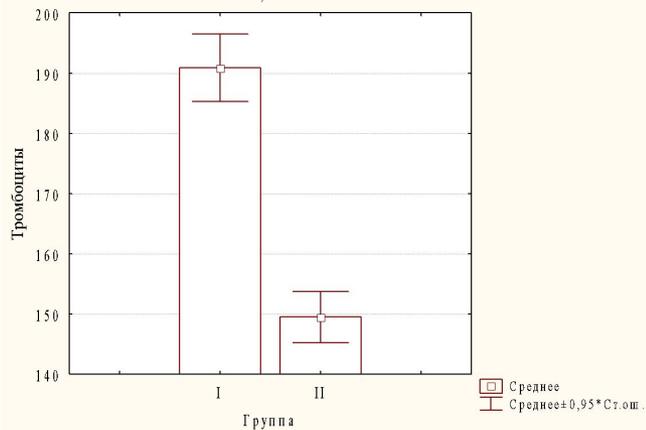


Диаграмма рассеяния

Диаграмма рассеяния (Таблица данных3 в Workbook2.stw 10v*150с)

$$\text{Возраст} = 39,7642 - 0,02 * x$$

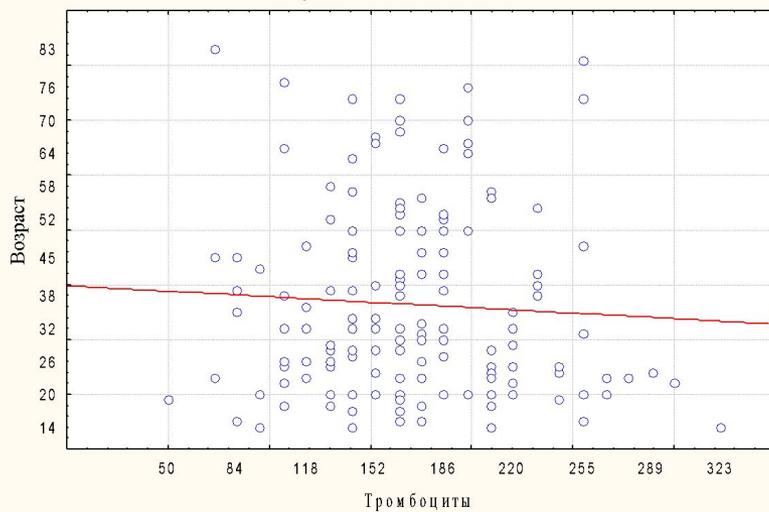


Диаграмма рассеяния (Таблица данных3 в Workbook2.stw 12v*150с)

$$\text{Вес} = -35,7946 + 0,5916 * x$$

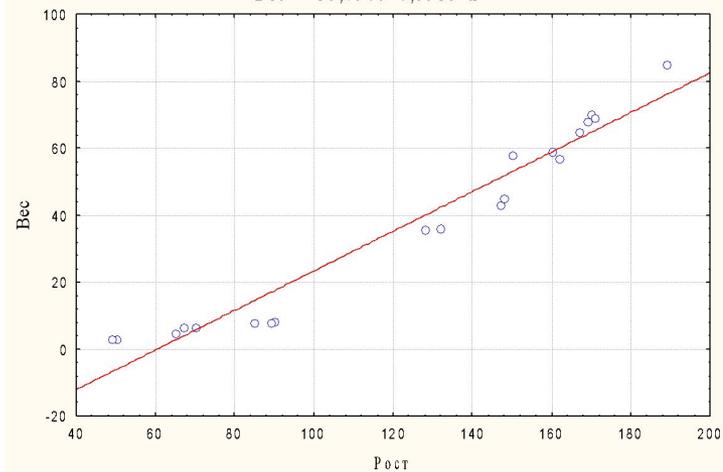
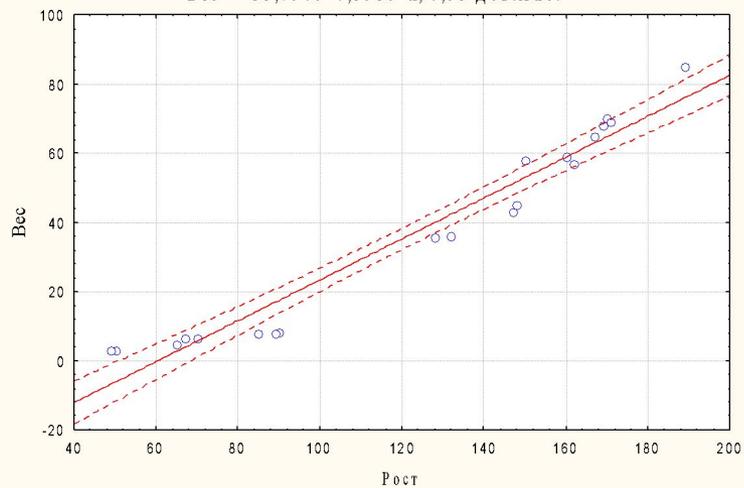


Диаграмма рассеяния (Таблица данных3 в Workbook2.stw 12v*150с)

$$\text{Вес} = -35,7946 + 0,5916 * x; 0,95 \text{ Дов.Инт.}$$



Статистический анализ

Выдвижение и проверка гипотез

- ▣ Нулевая гипотеза (H_0) представляет собой утверждение, в котором исследователь констатирует факт отсутствия каких-либо отличий, либо влияний на исходные данные
- ▶ Исследователю необходимо сформулировать нулевую гипотезу так, чтобы отказ от нее приводил к желательному заключению
- ▶ Альтернативная гипотеза (H_1) предназначена для определения согласованности данных с нулевой гипотезой и опровергает ее

Статистическая гипотеза подтверждается
или отклоняется с помощью ...



Статистические критерии: выбор

- Строгое математическое правило, по которому принимается или отвергается та или иная статистическая гипотеза с известным уровнем значимости
- **Параметрические критерии** – группа статистических критериев, которые включают в расчет параметры вероятностного распределения признака (средние и дисперсии) и предполагают нормальность распределения
- **Непараметрические методы** разработаны для тех ситуаций, когда исследователь ничего не знает о параметрах исследуемой популяции, непараметрические методы не основываются на оценке параметров (таких как среднее или стандартное отклонение) при описании выборочного распределения интересующей величины.
- Непараметрические методы позволяют обрабатывать данные "низкого качества" из выборок малого объема с переменными, про распределение которых мало что или вообще ничего не известно.

Расчет величины статистического критерия

- Выбрать соответствующие формулы для расчета статистических критериев
- Принять решение о нулевой гипотезе: либо она отвергается, либо принимается. Принятие гипотезы не означает, что она является единственно верной.

v-число степеней свободы

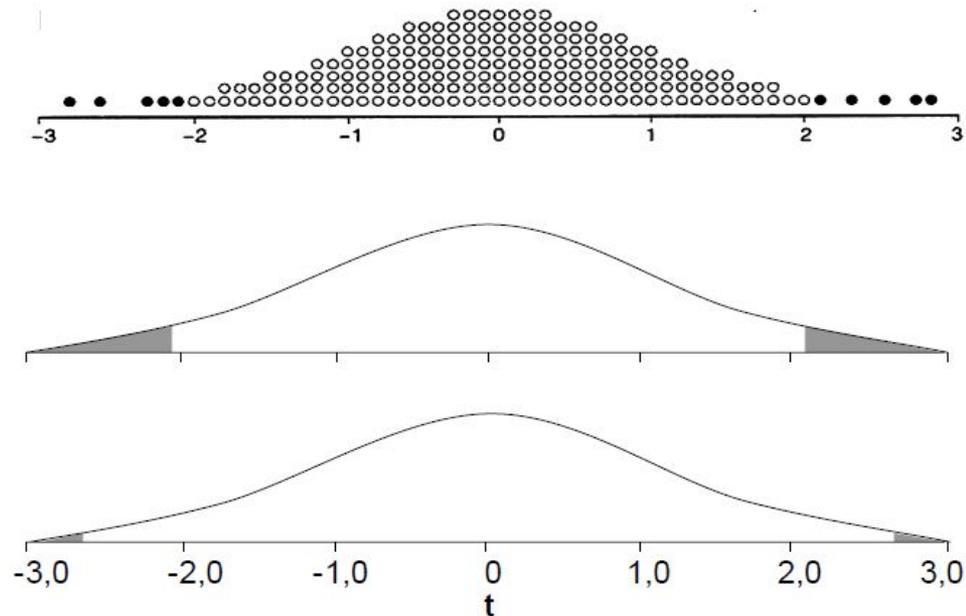
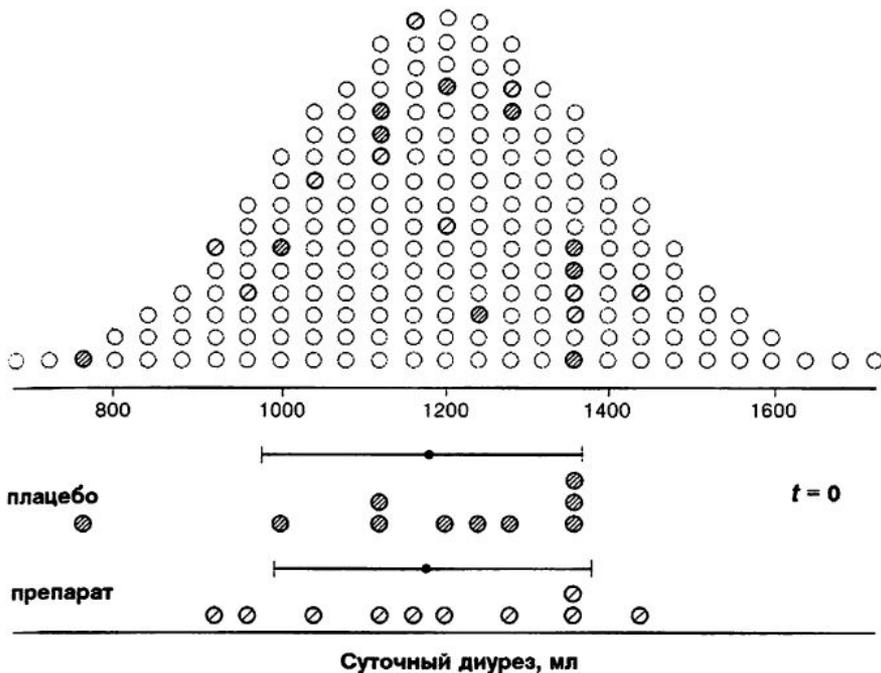
Таблица 4.1. Критические значения t (двусторонний вариант)

| v | Уровень значимости α | | | | | | | | |
|----|-----------------------------|-------|-------|--------|--------|--------|---------|---------|---------|
| | 0,5 | 0,2 | 0,1 | 0,05 | 0,02 | 0,01 | 0,005 | 0,002 | 0,001 |
| 1 | 1,000 | 3,078 | 6,314 | 12,706 | 31,821 | 63,656 | 127,321 | 318,289 | 636,578 |
| 2 | 0,816 | 1,886 | 2,920 | 4,303 | 6,965 | 9,925 | 14,089 | 22,328 | 31,600 |
| 3 | 0,765 | 1,638 | 2,353 | 3,182 | 4,541 | 5,841 | 7,453 | 10,214 | 12,924 |
| 4 | 0,741 | 1,533 | 2,132 | 2,776 | 3,747 | 4,604 | 5,598 | 7,173 | 8,610 |
| 5 | 0,727 | 1,476 | 2,015 | 2,571 | 3,365 | 4,032 | 4,773 | 5,894 | 6,869 |
| 6 | 0,718 | 1,440 | 1,943 | 2,447 | 3,143 | 3,707 | 4,317 | 5,208 | 5,959 |
| 7 | 0,711 | 1,415 | 1,895 | 2,365 | 2,998 | 3,499 | 4,029 | 4,785 | 5,408 |
| 8 | 0,706 | 1,397 | 1,860 | 2,306 | 2,896 | 3,355 | 3,833 | 4,501 | 5,041 |
| 9 | 0,703 | 1,383 | 1,833 | 2,262 | 2,821 | 3,250 | 3,690 | 4,297 | 4,781 |
| 10 | 0,700 | 1,372 | 1,812 | 2,228 | 2,764 | 3,169 | 3,581 | 4,144 | 4,587 |
| 11 | 0,697 | 1,363 | 1,796 | 2,201 | 2,718 | 3,106 | 3,497 | 4,025 | 4,437 |
| 12 | 0,695 | 1,356 | 1,782 | 2,179 | 2,681 | 3,055 | 3,428 | 3,930 | 4,318 |
| 13 | 0,694 | 1,350 | 1,771 | 2,160 | 2,650 | 3,012 | 3,372 | 3,852 | 4,221 |
| 14 | 0,692 | 1,345 | 1,761 | 2,145 | 2,624 | 2,977 | 3,326 | 3,787 | 4,140 |
| 15 | 0,691 | 1,341 | 1,753 | 2,131 | 2,602 | 2,947 | 3,286 | 3,733 | 4,073 |
| 16 | 0,690 | 1,337 | 1,746 | 2,120 | 2,583 | 2,921 | 3,252 | 3,686 | 4,015 |
| 17 | 0,689 | 1,333 | 1,740 | 2,110 | 2,567 | 2,898 | 3,222 | 3,646 | 3,965 |
| 18 | 0,688 | 1,330 | 1,734 | 2,101 | 2,552 | 2,878 | 3,197 | 3,610 | 3,922 |
| 19 | 0,688 | 1,328 | 1,729 | 2,093 | 2,539 | 2,861 | 3,174 | 3,579 | 3,883 |
| 20 | 0,687 | 1,325 | 1,725 | 2,086 | 2,528 | 2,845 | 3,153 | 3,552 | 3,850 |
| 21 | 0,686 | 1,323 | 1,721 | 2,080 | 2,518 | 2,831 | 3,135 | 3,527 | 3,819 |
| 22 | 0,686 | 1,321 | 1,717 | 2,074 | 2,508 | 2,819 | 3,119 | 3,505 | 3,792 |
| 23 | 0,685 | 1,319 | 1,714 | 2,069 | 2,500 | 2,807 | 3,104 | 3,485 | 3,768 |
| 24 | 0,685 | 1,318 | 1,711 | 2,064 | 2,492 | 2,797 | 3,091 | 3,467 | 3,745 |
| 25 | 0,684 | 1,316 | 1,708 | 2,060 | 2,485 | 2,787 | 3,078 | 3,450 | 3,725 |
| 26 | 0,684 | 1,315 | 1,706 | 2,056 | 2,479 | 2,779 | 3,067 | 3,435 | 3,707 |
| 27 | 0,684 | 1,314 | 1,703 | 2,052 | 2,473 | 2,771 | 3,057 | 3,421 | 3,689 |
| 28 | 0,683 | 1,313 | 1,701 | 2,048 | 2,467 | 2,763 | 3,047 | 3,408 | 3,674 |
| 29 | 0,683 | 1,311 | 1,699 | 2,045 | 2,462 | 2,756 | 3,038 | 3,396 | 3,660 |
| 30 | 0,683 | 1,310 | 1,697 | 2,042 | 2,457 | 2,750 | 3,030 | 3,385 | 3,646 |

- Для равных выборок: $v=2(n-1)$
- Для произвольных выборок: $v=n_1+n_2-2$

Статистический уровень значимости (p-уровень)

- вероятность ошибочного отклонения нулевой гипотезы
- вероятность справедливости нулевой гипотезы
- p-уровень, равный 0.05 рассматривается как приемлемая граница уровня ошибки
- если $p > 0,05$, то нет достаточных оснований, чтобы отвергнуть H_0



Важно!

необходимо указывать:

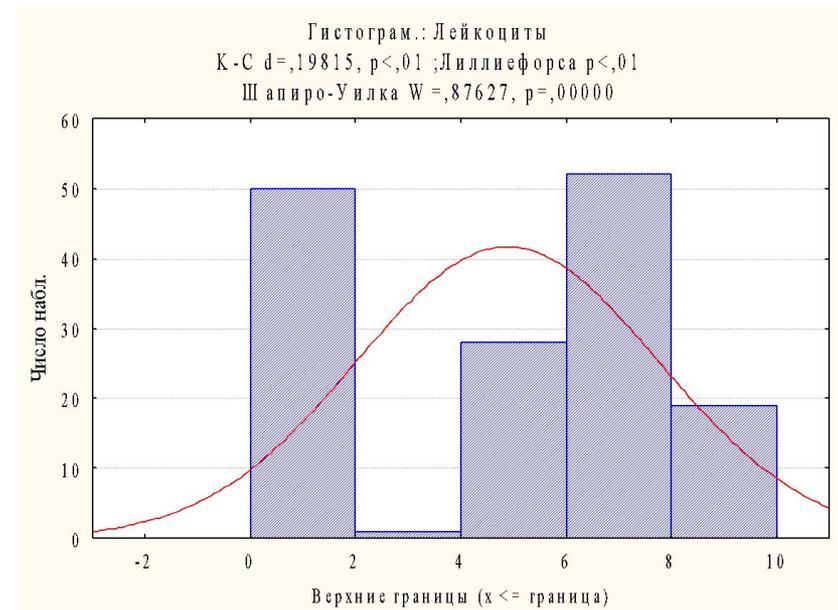
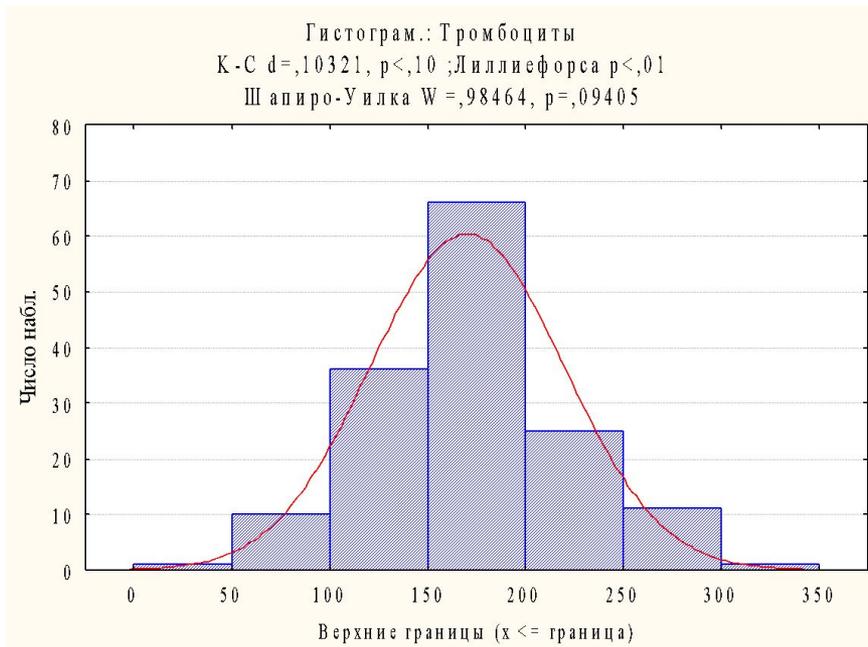
- название и значение статистического критерия
- действительный p -уровень (до $p > 0,001$)

Проверка распределения на нормальность

- Гистограмма (визуальная проверка)
- Применение критериев (статистическая проверка)
 - Критерий Колмогорова-Смирнова
 - Критерий Лиллиефорса
 - Критерий Шапиро-Уилка

- H_0 : распределение нормальное
- H_1 : распределение отличается от нормального

□ Если W статистика значима, то гипотеза о нормальном распределении значений переменной отвергается. Т.е., если $p \leq 0,05$, то переменная имеет распределение отличное от нормального (ненормальное распределение)



Корреляционный анализ

Корреляционный анализ

- **Параметрический корреляционный анализ** Пирсона – для исследования взаимосвязи нормально распределенных количественных признаков
- **Непараметрические методы** корреляционного анализа Спирмена, Кендалла, гамма – для исследования взаимосвязи:
 - Количественных признаков независимо от вида их распределения
 - Количественного и качественного порядкового признака
 - Двух порядковых признаков
- Категориальные – таблица сопряженностей

Корреляционный анализ

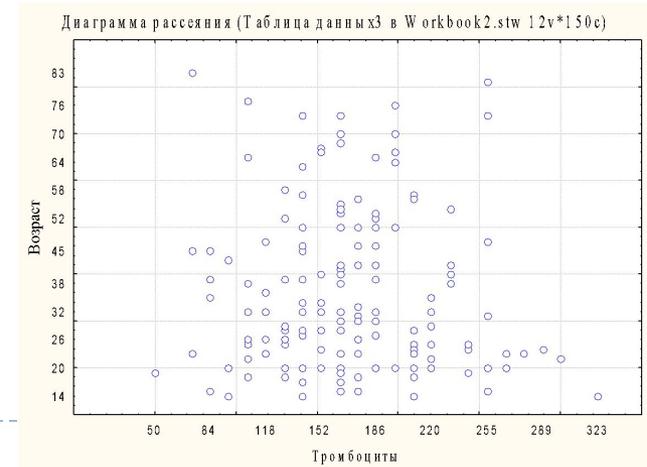
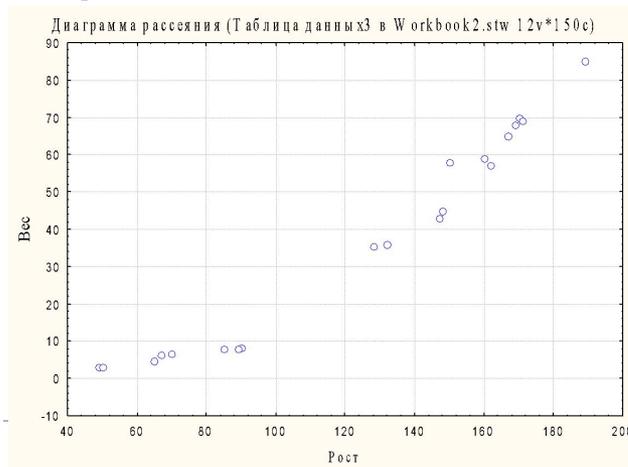
□

- ▶ Correlation – взаимосвязь
- ▶ Мера зависимости между двумя переменными (направление и сила связи)
- ▶ Пример: корреляция между ростом и массой тела
- ▶ Коэффициент корреляции Пирсона определяется:

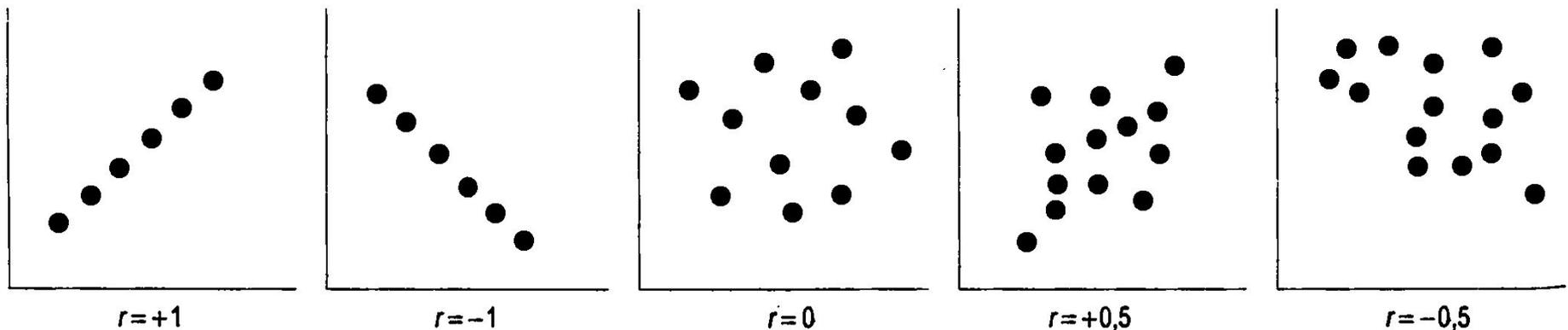
$$r = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sqrt{\sum(X - \bar{X})^2 \sum(Y - \bar{Y})^2}}$$

Корреляционный анализ: коэффициент корреляции

- Значения от -1 до $+1$
- Знак коэффициента показывает направление связи
- Чем больше коэффициент корреляции по абсолютной величине, тем сильнее коррелированы переменные
- Корреляция между x и y не означает соотношение причины и следствия
- безразмерен



- Принята (условно) следующая классификация силы корреляции в зависимости от значения коэффициента корреляции r .
- $|r| \leq 0,25$ - слабая корреляция
- $0,25 < |r| < 0,75$ - умеренная корреляция
- $|r| \geq 0,75$ – сильная корреляция



Корреляционный анализ

Когда не следует рассчитывать коэффициент корреляции?

- Нелинейное соотношение между переменными
- Есть аномальные значения (выбросы)
- Данные содержат подгруппы, для которых средние уровни наблюдений различны

Расчет коэффициента корреляции

▣ Гипотезы:

- ▶ H_0 - связи между переменными нет
- ▶ H_1 - связь между переменными есть

- ▶ Для обоснования использования коэффициента корреляции Пирсона, необходимо провести проверку на нормальность распределения переменной

Параметрический корреляционный анализ Пирсона – для исследования взаимосвязи нормально распределенных количественных признаков

Непараметрические методы корреляционного анализа
Спирмена, Кендалла, гамма – для исследования взаимосвязи:

- Количественных признаков независимо от вида их распределения
- Количественного и качественного порядкового признака
- Двух порядковых признаков

Подходы к сравнению двух групп по количественному признаку:

- с использованием доверительных интервалов
(ответ на вопрос: насколько велики различия совокупностей?)
- путем проверки статистических гипотез
(ответ на вопрос: в какой степени можно быть уверенным, что различия между совокупностями действительно существуют?)



-
- При описании результатов исследования рекомендуется представлять результаты применения обоих подходов



Доверительный интервал для разности средних

- Расчет объединенной оценки дисперсии

$$s^2 = \frac{1}{2} (s_{\text{Д}}^2 + s_{\text{П}}^2)$$

- Расчет стандартной ошибки разности средних

$$s_{\bar{X}_{\text{Д}} - \bar{X}_{\text{П}}} = \sqrt{\frac{s^2}{n_{\text{Д}}} + \frac{s^2}{n_{\text{П}}}}$$

- Расчета доверительного интервала разности средних

$$(\bar{X}_1 - \bar{X}_2) - t_{0,05} s_{\bar{X}_1 - \bar{X}_2} < \mu_1 - \mu_2 < (\bar{X}_1 - \bar{X}_2) + t_{0,05} s_{\bar{X}_1 - \bar{X}_2}$$

- Не должен содержать «0»
-



Сравнимые группы:

- **независимые (несвязанные)**

если набор объектов исследования (участников) в каждую из групп осуществляется независимо от того, какие объекты исследования (участники) включены в другую группу (рандомизация)

- **зависимые (связанные)**

динамические исследования, когда изучаются одни и те же объекты в разные моменты времени



Независимые выборки



Независимые выборки

Проверка статистической гипотезы

Расчет ДИ

Параметрические
методы

Непараметрические
методы

Нормальное
распределение

Любое распределение

t – критерий для
независимых
выборок

- U критерий Манна-Уитни
- критерий Вальда - Вольфовица
- критерий Колмогорова-Смирнова

Параметрический метод

t критерий для независимых выборок



t – критерий (t-test, Student's t-test)

Алгоритм действий



**Вильям ГОССЕТ -
"СТЮДЕНТ"**
(W. S. Gosset – "Student"
1876-1937)

известный английский статистик

- Зависимые или независимые наблюдения?
- Чему равен α -уровень критерия?
- Являются ли распределения переменных нормальными?
- Равны ли дисперсии в группах?
- Какой вывод можно сделать?
- Какая надежность вывода?



t критерий для независимых выборок: ***соблюдение условий***

□ Классический вариант:

- *значения признаков в каждой из сравниваемых групп должны иметь нормальное распределение (должна проводиться проверка распределения признака на соответствие закону нормального распределения)*
- *дисперсии распределений признаков в сравниваемых группах равны (может быть проверено с помощью критерия Левена)*

□ Модифицированный вариант:

- *при невыполнении условия равенства дисперсий – расчет t-критерия с отдельными оценками дисперсий*
-



Выдвижение и проверка гипотез

- Представляет собой стандартное отклонение распределения средних отдельных выборок (рассчитанное из средних отдельных выборок)

$$\mathbf{m} \text{ или SEM} = \frac{\sigma}{\sqrt{n}}$$

и отображает точность оцененного параметра среднего

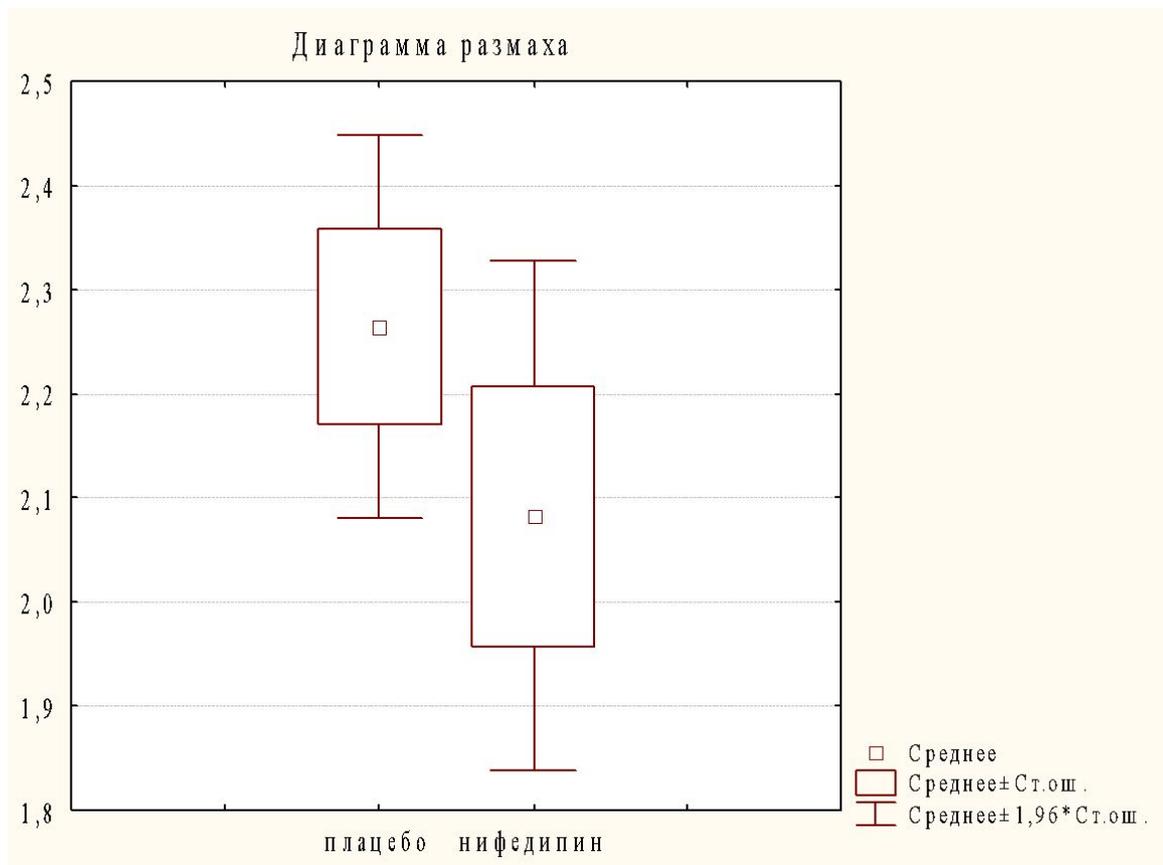
t-критерий для независимых выборок

□

\bar{X}

Пример: исследование препаратов, влияющих на диаметр коронарных сосудов

- ▣ Метод визуализации: диаграмма размаха
- ▣ Статистический метод: T-критерий для независимых выборок



Представление результатов:

- Число объектов исследования в каждой из групп
 - Средние и СКО изучаемого признака для каждой из групп
 - Результаты применения критериев для оценки нормальности распределения и равенства дисперсий в случае, если используется классический критерий Стьюдента
 - Результаты применения критерия для оценки нормальности распределения и указание модифицированного критерия Стьюдента для групп с различными дисперсиями
 - Диаграммы размаха
-



Непараметрические методы



Когда используются методы непараметрической статистики

Ответ: когда распределение данных отличается от нормального

Преимущество:

- критерии непараметрической статистики не содержат никаких предположений относительно распределения данных
 - отсутствие больших выборок
 - шкала измерений может быть порядковой

 - **Недостаток:** низкая мощность
-



Если условия применимости t критериев
не выполнены...

Непараметрические критерии (*non-parametric tests*)

- критерий Вальда-Вольфовица
(*Wald-Wolfowitz runs test*)
 - критерий Колмогорова Смирнова
(*Kolmogorov-Smirnov test*)
 - критерий Манна-Уитни (U-критерий)
(*Mann-Uitney U-test*)
- независимые
выборки
-
- критерий знаков
(*sign test*)
 - критерий Вилкоксона
(*Wilcoxon signed-rank test*)
- зависимые выборки
-



Критерий серий Вальда-Вольфовица

- *непараметрическая альтернатива t критерия для независимых выборок*
- Значения сравниваемых групп выстраиваются в единую последовательность по рангу. Производится подсчет количества смен группирующего признака, с помощью которого можно найти количество непрерывных последовательностей (число смен + 1) - серий.
- Если нет различия между группами, то число и длина серий, относящихся к одной и той же группе, будут примерно одинаковыми. В противном случае две группы отличаются друг от друга.



Двухвыборочный критерий Колмогорова-Смирнова

- *непараметрическая альтернатива t критерия для независимых выборок*
- Критерий основан на максимуме абсолютного значения разности эмпирических функций первой и второй выборки, также чувствителен к различию общей формы распределений двух выборок (в частности, различие в дисперсии, асимметрии и т.д.).



U критерий Манна-Уитни

- *непараметрическая альтернатива t критерия для независимых выборок*
- U критерий вычисляется, как сумма индикаторов попарного сравнения элементов первой выборки с элементами второй выборки.
- U критерий - наиболее мощная (чувствительная) непараметрическая альтернатива t -критерия для независимых выборок.



Представление результатов

- Число объектов исследования для каждой из групп
- Медианы и границы интерквартильного отрезка для каждой из групп
- Точное значение критерия и p - уровень
- Диаграммы размаха



Зависимые (связанные) выборки



t критерий для зависимых выборок

- проверить, различаются ли средние значения количественного признака до и после лечения, если известно, что в обоих случаях распределение признака является нормальным*



Представление результатов

- Число объектов исследования в каждой из выборок
- Аргументированная информация о выполнении условий применимости метода
- Средние значения изучаемого признака и СКО для каждой из групп
- Точное значение критерия и p уровень
- Диаграмма размаха



Если условия применимости t критериев не выполнены...

- критерий знаков
(*sign test*)

- критерий Вилкоксона
(*Wilcoxon signed-rank test*)



зависимые выборки



Критерий знаков

- *непараметрическая альтернатива t критерия для зависимых выборок*
 - **Критерий основан на следующих простых соображениях:** подсчитывает, сколько раз определенное значение первой переменной (A) больше соответствующего значения переменной (B), иными словами, определяется количество положительных разностей между значениями переменной (A) и значениями переменной (B).
 - **N.B!** Учитываются только знаки разностей (а не их значения).
-



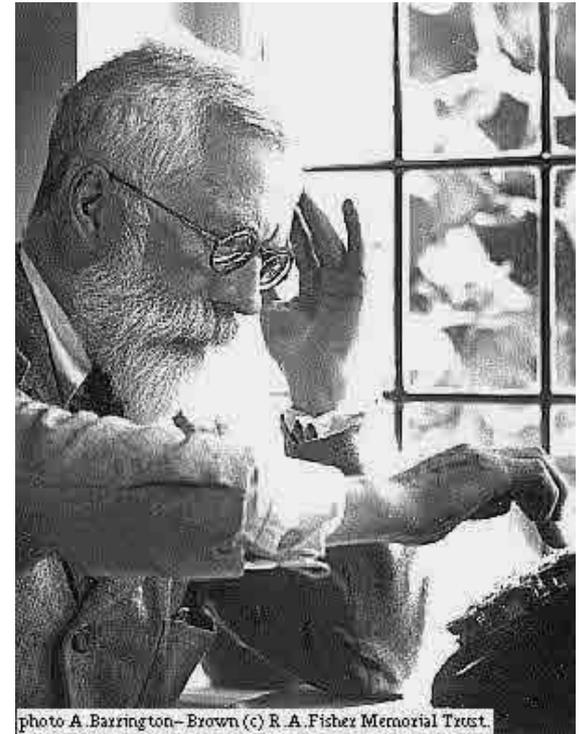
W критерий знаковых рангов Вилкоксона

- *непараметрическая альтернатива t критерия для зависимых выборок*
 - Критерий принимает во внимание не только знаки разностей, но и их величину.
 - Более мощный критерий (по сравнению с критерием знаков). Если предположения параметрического t -критерия для зависимых выборок (интервальная шкала) выполнены, то критерий имеет почти такую же мощность, как и t -критерий.
-



Дисперсионный анализ

- ▣ *ANOVA – analysis of variance*
(1920 г. Рональд Фишер,
английский статистик и генетик)



Общее назначение

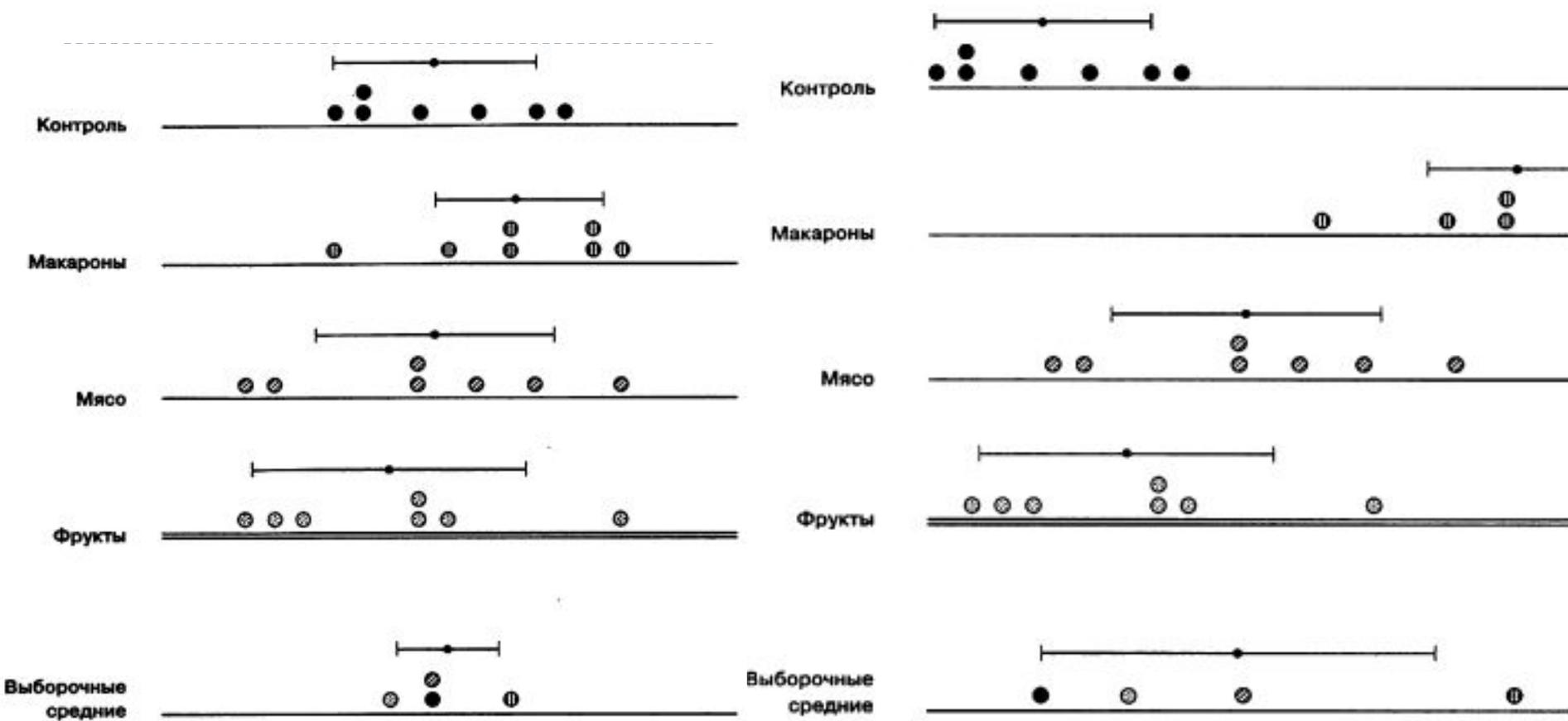
- Сравнение **средних** в **нескольких** группах
- Сравнение групп проводится с помощью оценки межгрупповой и внутригрупповой **дисперсий**. Отсюда термин «дисперсионный анализ»



Дисперсионный анализ

- Для оценки различий, необходимо сравнить разброс выборочных средних с разбросом значений внутри каждой из групп





$F = \text{межгрупповая дисперсия} / \text{внутригрупповая дисперсия}$
(разброс выборочных средних) / (разброс внутри групп)



Дисперсионный анализ

$F = \text{межгрупповая дисперсия} / \text{внутригрупповая дисперсия}$

или

$$F = \frac{\text{Дисперсия совокупности, оцененная по выборочным средним}}{\text{Дисперсия совокупности, оцененная по выборочным дисперсиям}}$$

- Числитель и знаменатель соотношения – оценки одной и той же величины – дисперсии совокупности

- Если верна нулевая гипотеза, то как внутригрупповая, так и межгрупповая дисперсии служат оценками одной и той же дисперсии и должны быть приближенно равны



Проверяемая гипотеза

- Нулевая гипотеза: различий между группами нет
- При истинности нулевой гипотезы, оценка дисперсии, связанной с внутригрупповой изменчивостью, должна быть близкой к оценке межгрупповой дисперсии
- При ложности – значимо отличаться



Дисперсионный анализ - этапы

- Проверка нормальности
- Проверка равенства дисперсий
- ANOVA
- Апостериорные сравнения групп



Методы множественного сравнения

Если ДА показал наличие значимых различий между средними значениями выборок

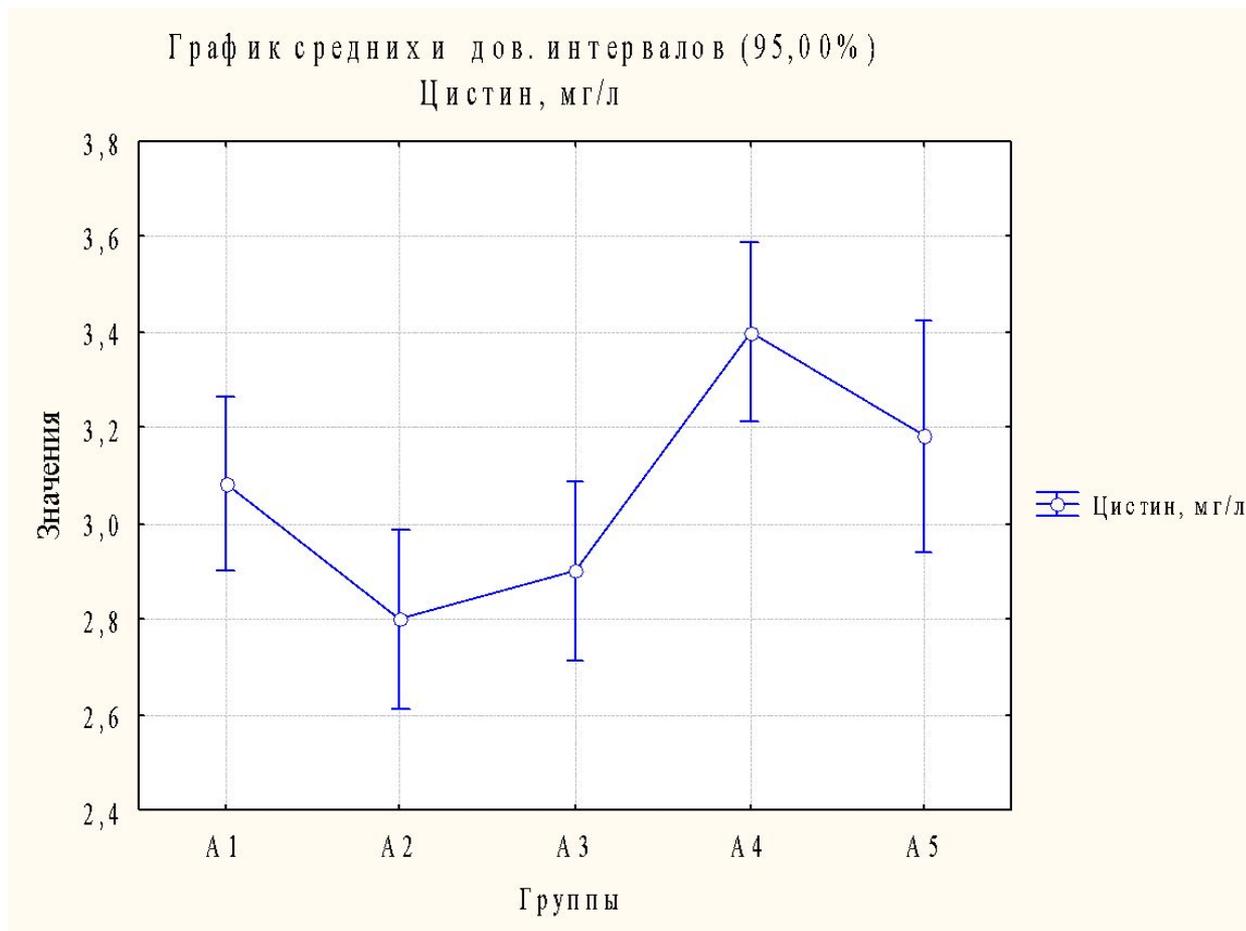


апостериорные сравнения с использованием:

- Поправки Бонферрони
- Критерия Фишера (наименьшей значимой разности)
- Критерия Шеффе
- Критерия Тьюки
- Критериев размахов Ньюмана-Кеулса и Дункана



Графическое представление результатов



Представление результатов

- Число объектов исследования в каждой из выборок
- Аргументированная информация о выполнении условий применимости метода
- Средние значения изучаемого признака и СКО для каждой из групп
- Точное значение критерия и p -уровень
- Диаграмма размаха



т.е. данные отвечают на вопрос о том, между какими именно группами различие статистически значимо!

□ Выход: апостериорные сравнения



Окончательный результат

| Группы | Крит. НЗР; перем.: Цистин, мг/л (Таблица данных2 в Д Отмечены различия, значимые на уровне $p < ,05000$ | | | | |
|--------|--|----------|----------|----------|----------|
| | {1} | {2} | {3} | {4} | {5} |
| | M=3,0833 | M=2,8000 | M=2,9000 | M=3,4000 | M=3,1833 |
| A1 {1} | | 0,015735 | 0,106063 | 0,007739 | 0,369157 |
| A2 {2} | 0,015735 | | 0,369157 | 0,000011 | 0,001740 |
| A3 {3} | 0,106063 | 0,369157 | | 0,000113 | 0,015735 |
| A4 {4} | 0,007739 | 0,000011 | 0,000113 | | 0,058627 |
| A5 {5} | 0,369157 | 0,001740 | 0,015735 | 0,058627 | |



Расчет поправки Бонферрони

- $p = 1 - (1 - 0,05)^k$,
- или $p = 0,05 \times k$,
- где k – число сравнений.
- Например, при сравнении 4 групп необходимо сделать 6 сравнений: $\alpha = \alpha^* / 6$, при $\alpha = 0,05$ расчет имеет следующий вид: $0,05 / 6 = 0,008.$, т.е. « p » должно быть меньше 0,008.



Дисперсионный анализ повторных измерений



Дисперсионный анализ - этапы

- Проверка нормальности
- Проверка равенства дисперсий
- ANOVA
- Апостериорные сравнения групп



Различия между несколькими несвязанными группами – непараметрический H-критерий Краскела-Уоллиса

- Обобщение критерия Манна-Уитни для трех и более независимых выборок
 - Критерий базируется на общей ранговой последовательности значений всех выборок и не требует предположения о нормальности распределения
 - Анализируемый признак должен быть количественным или порядковым
-



т.е. данные не отвечают на вопрос о том, между какими именно группами различие статистически значимо!

- Выход: апостериорные сравнения с использованием непараметрического теста Манна-Уитни, применяя поправку Бонферрони при оценке значения p



Расчет поправки Бонферрони

- $p = 1 - (1 - 0,05)^k$,
- или $p = 0,05 \times k$,
- где k – число сравнений.
- Например, при сравнении 4 групп необходимо сделать 6 сравнений: $\alpha = \alpha^* / 6$, при $\alpha = 0,05$ расчет имеет следующий вид: $0,05 / 6 = 0,008.$, т.е. « p » должно быть меньше 0,008.

При 6 сравнениях
вероятность «ошибки»
составит 30% !!!!



Использованная литература

- Гланц, С. Медико-биологическая статистика / С. Гланц; пер. англ. — М.: Практика, 1998. — 459 с.
 - Петри, А. Наглядная медицинская статистика / А. Петри, К. Сэбин; пер. с англ. под ред. В. П. Леонова. — 2-е изд., перераб. и доп. — М.: ГЭОТАР-Медиа, 2009. — 168 с.
 - Т.А.Ланг. Как описывать статистику в медицине. Аннотированное руководство для авторов, редакторов и рецензентов / Т.А.Ланг, М.Сесик; пер. с англ. под ред. В.П. Леонова. — М.: Практическая медицина, 2011. — 480 с.
 - О.Ю.Реброва, Статистический анализ медицинских данных., 2002. — 312с.
 - Материалы интернет ресурса: statsoft.ru
-



Стентон Гланц

Медико-биологическая
СТАТИСТИКА

Перевод с английского
доктора физ.-мат. наук
Ю. А. Данилова
под редакцией
Н. Е. Бузикашвили
и Д. В. Самойлова



п р а к т и к а
Москва 1999

А. Петри, К. Сэбин

**НАГЛЯДНАЯ
СТАТИСТИКА
В МЕДИЦИНЕ**

Перевод с английского



Москва
Издательский дом
ГЭОТАР-МЕД
2003

Как описывать
статистику
в **медицине**

Руководство для авторов,
редакторов и рецензентов

Т. А. Ланг
М. Сесик

Перевод с английского
под редакцией В. П. Леонова

практическая медицина