

*STATISTICA*™

**[Парная] корреляция и регрессия**

# Типы статистических задач

Задачи	Инструменты
Описание совокупностей объектов	Анализ одной выборки; расчет параметров распределений (положения, формы); проверка нормальности распределений;
Сравнение параметров	Парные и множественные сравнения средних; сравнение распределений; сравнение частот; t-критерий; тест Манна-Уитни или Краскела-Уоллеса; дисперсионный анализ;
<b>Анализ зависимостей</b>	<b>Установление взаимосвязи между двумя переменными или между многими переменными; установление силы влияния одной или многих переменных на одну результирующую; корреляционный анализ, парная и множественная регрессия, логит-регрессия;</b>
Снижение размерности, ординация	Кластерный, факторный, дискриминантный анализ; анализ соответствий; многомерное шкалирование и др.

# Выбор статистического теста при сравнении распределений (**сравнении центральных тенденций и частот**)

Задача	Количественная шкала, нормальное распределение	Порядковая шкала или отклонение от нормального распределения	Номинальная шкала
Сравнить одну группу с гипотетическим значением	t-тест Стьюдента для одной выборки	Тест Вилкоксона	Тест хи-квадрат
Сравнить две не связанные совокупности	t-тест Стьюдента для не связанных совокупностей	Тест Манна-Уитни	Тест Фишера (тест хи-квадрат)
Сравнить две связанные совокупности	t-тест Стьюдента для связанных совокупностей	Тест Вилкоксона	Тест Мак-Неймера
Сравнить более двух не связанных совокупностей	Однофакторный дисперсионный анализ	Тест Краскела-Уоллиса	Тест хи-квадрат
Сравнить более двух связанных совокупностей	Дисперсионный анализ с повторными	Тест Фридмана	Тест Кохрана

# Задачи оценки взаимосвязи между переменными или

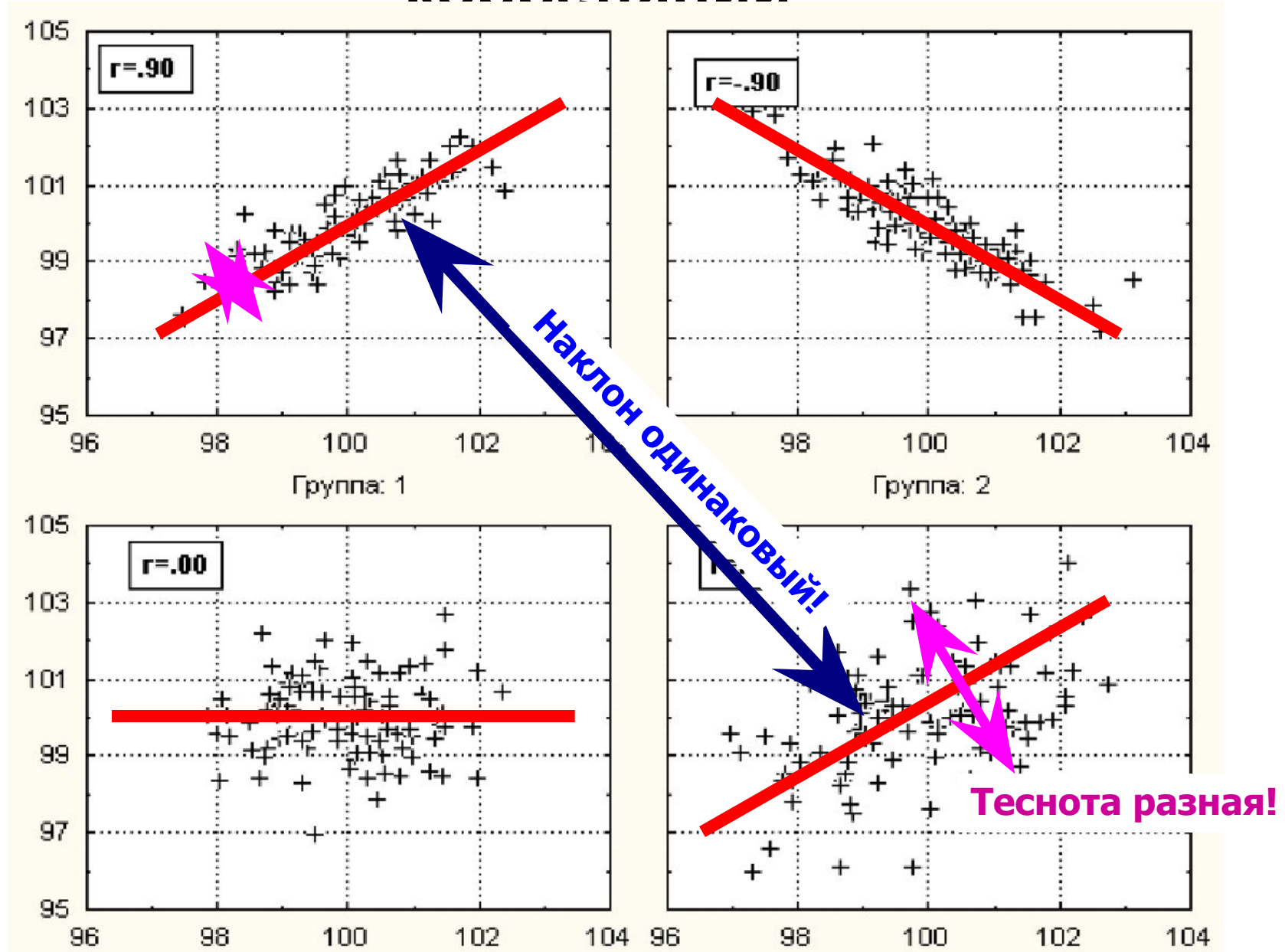
## прогноза

Задача	Количественные нормально распределенные переменные	Количественные ненормально распределенные переменные или ранги	Биноминальные данные (два возможных результата)
Оценить взаимосвязь между двумя	Коэффициент парной корреляции Пирсона	Коэффициенты ранговых корреляций (Спирмена,	Коэффициенты связи
Предсказать изменение одной переменной, если была измерена другая переменная	Простая линейная регрессия или нелинейная регрессия	Непараметрическая (ранговая) регрессия	Простая логистическая регрессия
Предсказать значение, базируясь на нескольких переменных	Множественная линейная (нелинейная) регрессия	Множественная линейная ранговая (медианная) регрессия	Множественная логистическая регрессия

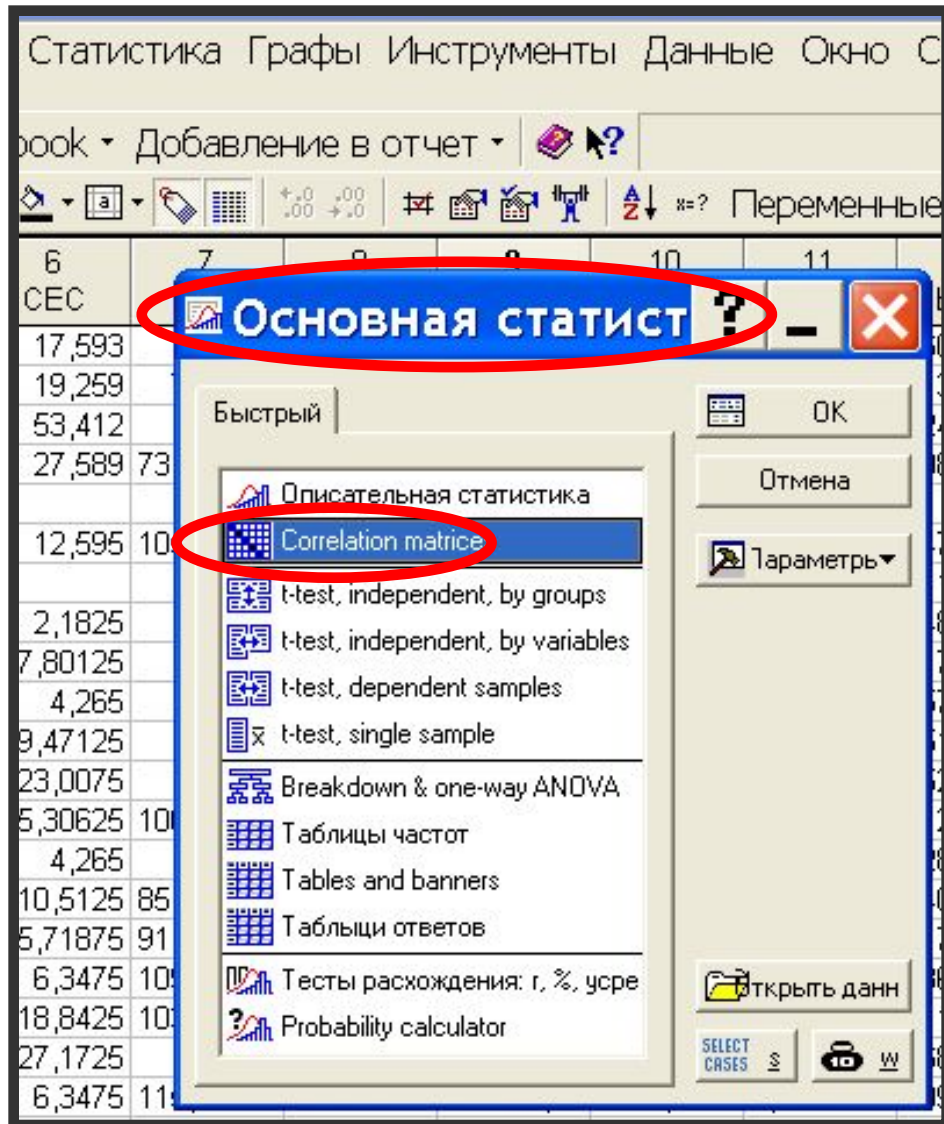
# Корреляция?

- Использование коэффициента корреляции позволяет оценить, в какой степени две переменные изменяются совместно – увеличивается ли или уменьшается одна переменная при изменении другой.
- Коэффициент корреляции – мера силы (тесноты и направления) связи между изменчивостью переменных.
- Интерпретация знака коэффициента корреляции – есть вопросы?
- Надежность коэффициент корреляции зависит от его величины и *n*.
- **Никаких причинных интерпретаций коэффициент корреляции сделать не позволяет!**
- **Коэффициент корреляции может быть использован только для прогноза направления (но не величины!) изменения одной переменной в связи с изменением другой**

# Знаки и теснота коэффициента корреляции

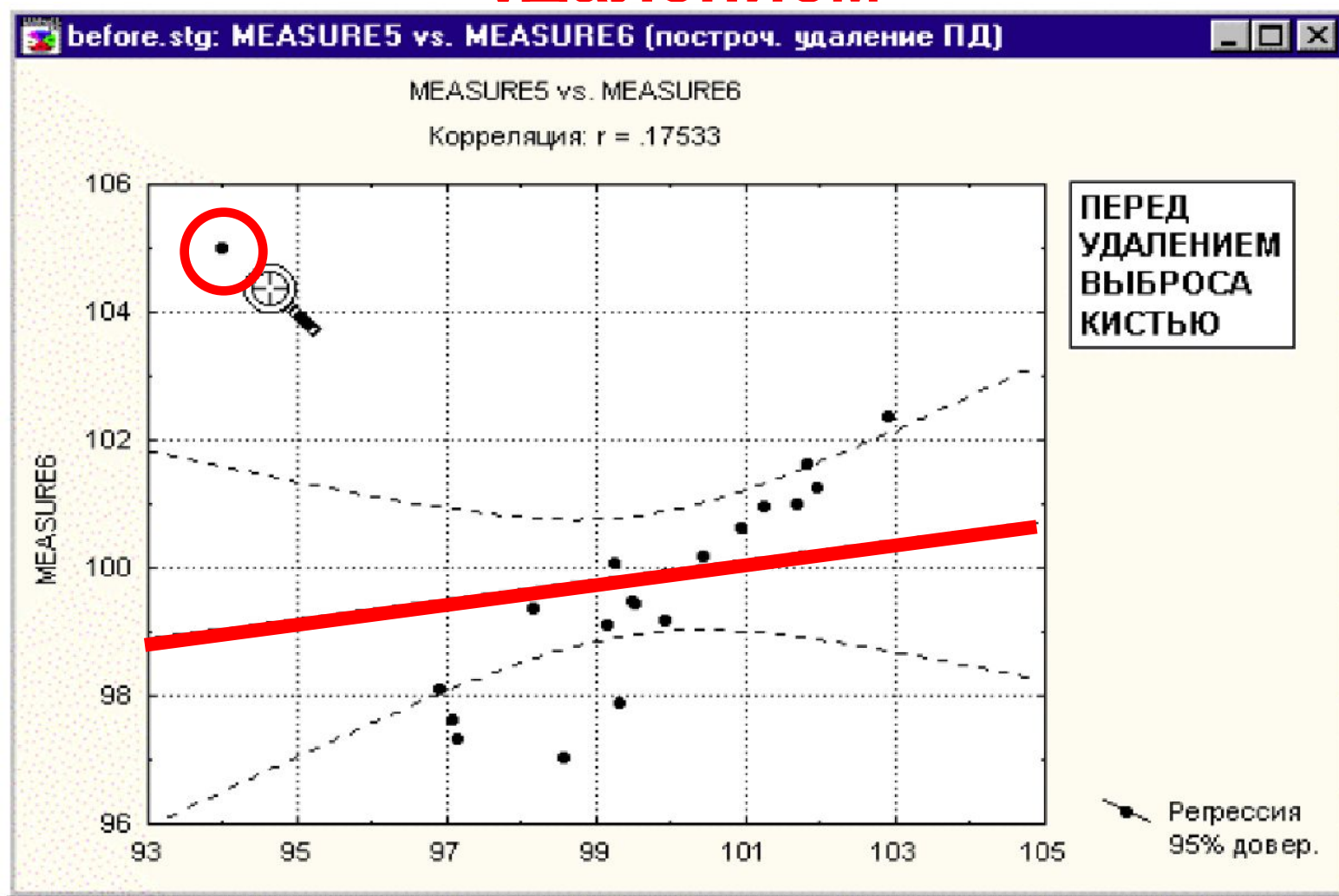


# Техника расчета $r$ Пирсона



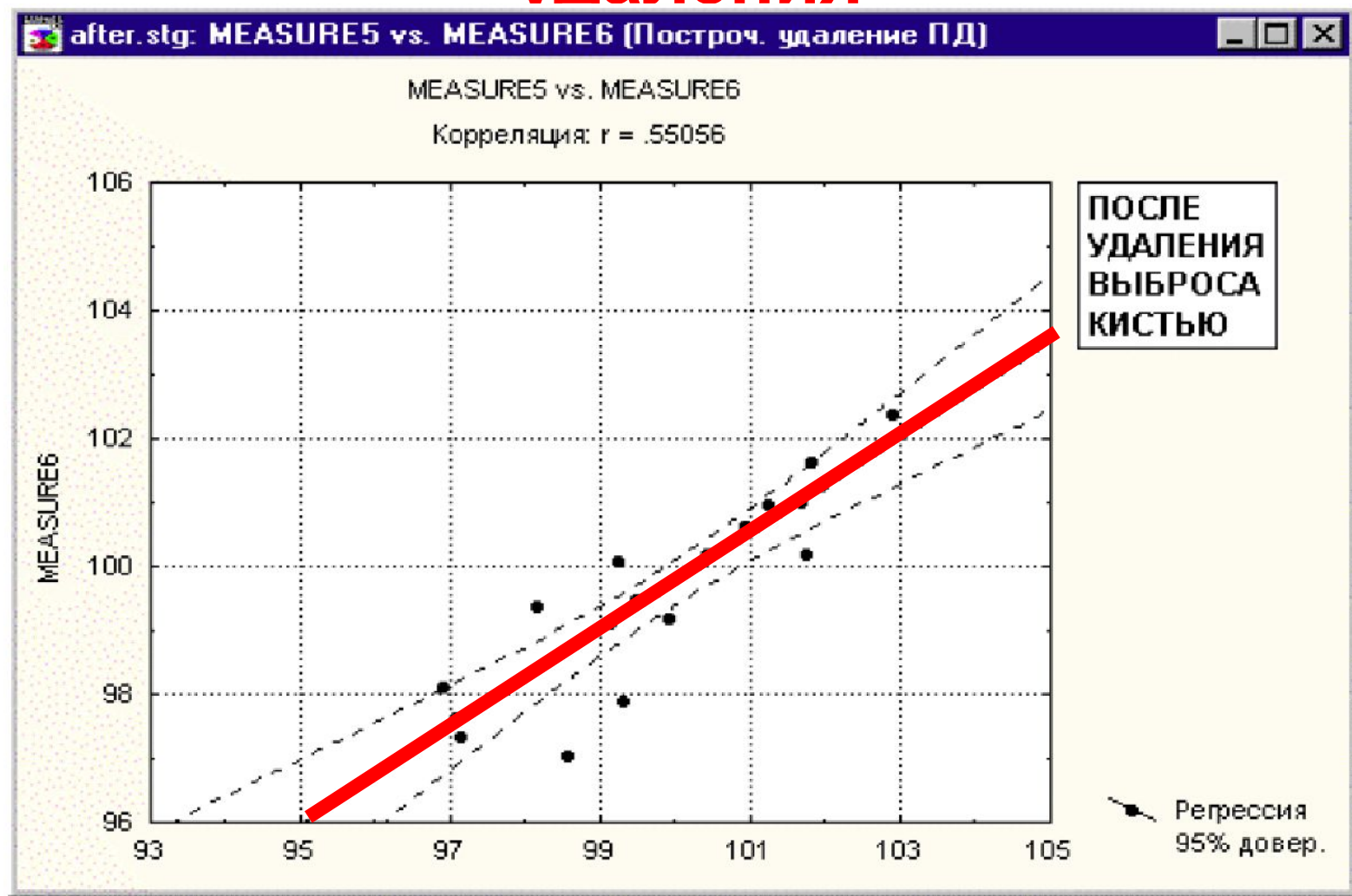
- «Пример\_тм\_токсичность\_преобразования.xls» (Cu\_хлорид \* Cd\_хлорид вместе и по зонам )
- Пары переменных или матрицы;
- Просмотр результатов в разном расширении;
- Иллюстрации;
- Категоризированные зависимости

# Нарушение «нормальности»: управление выбросами: **перед** **удалением**

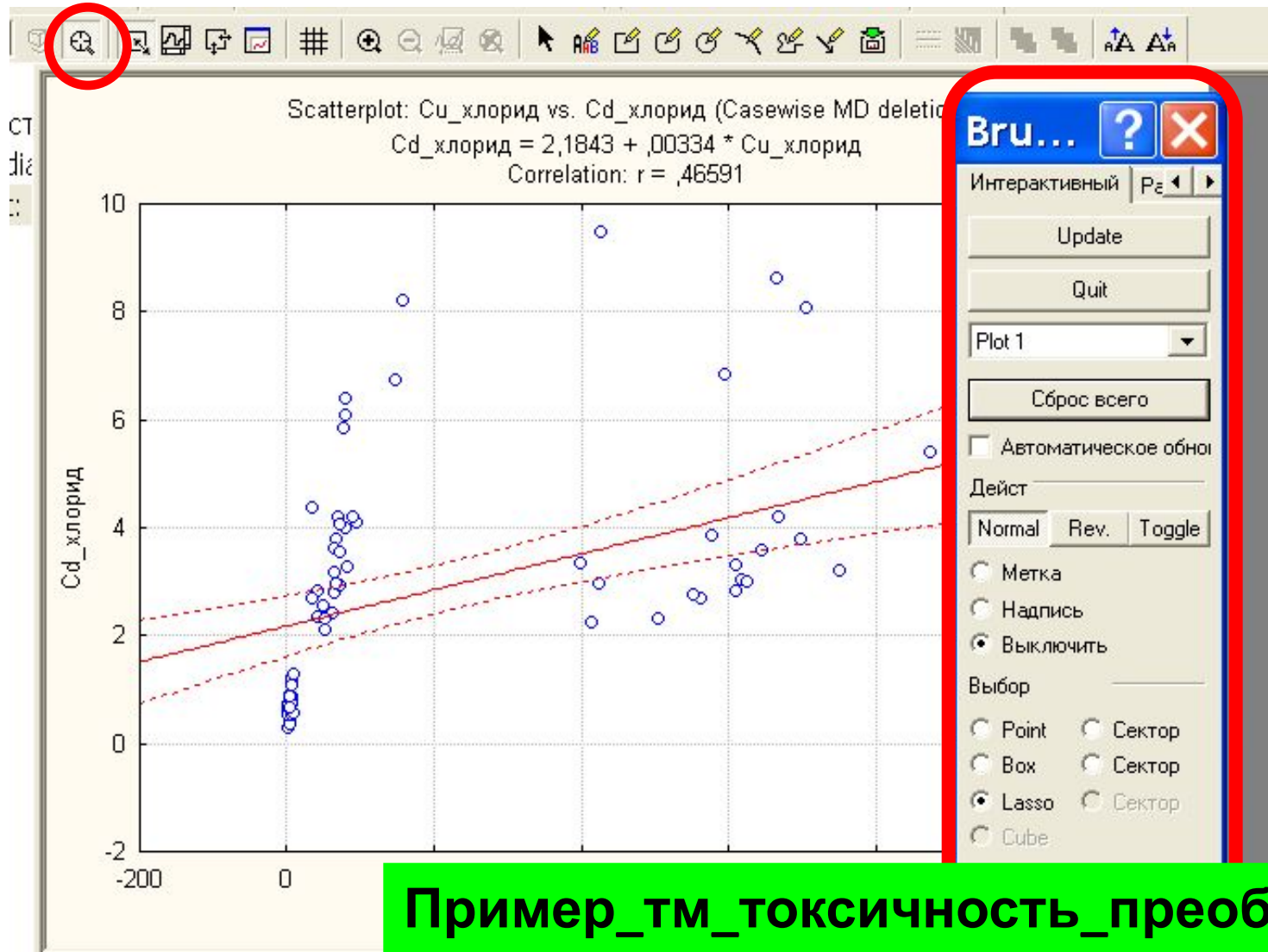




# Нарушение «нормальности»: управление выбросами: **после** **удаления**



# Управление выбросами: инструмент «КИСТЬ»



**Пример\_тм\_токсичность\_преобразован  
ия.xls;**

**Cu\_хлорид \* Cd\_хлорид вместе**

# Управление выбросами: общие правила отсутствуют

## Количественный подход к выбросам.

Некоторые исследователи применяют численные методы удаления выбросов. Например, исключаются значения, которые выходят за границы  $\pm 2$  стандартных отклонений (и даже  $\pm 1.5$  стандартных отклонений) вокруг выборочного среднего. В ряде случаев такая “чистка” данных абсолютно необходима.

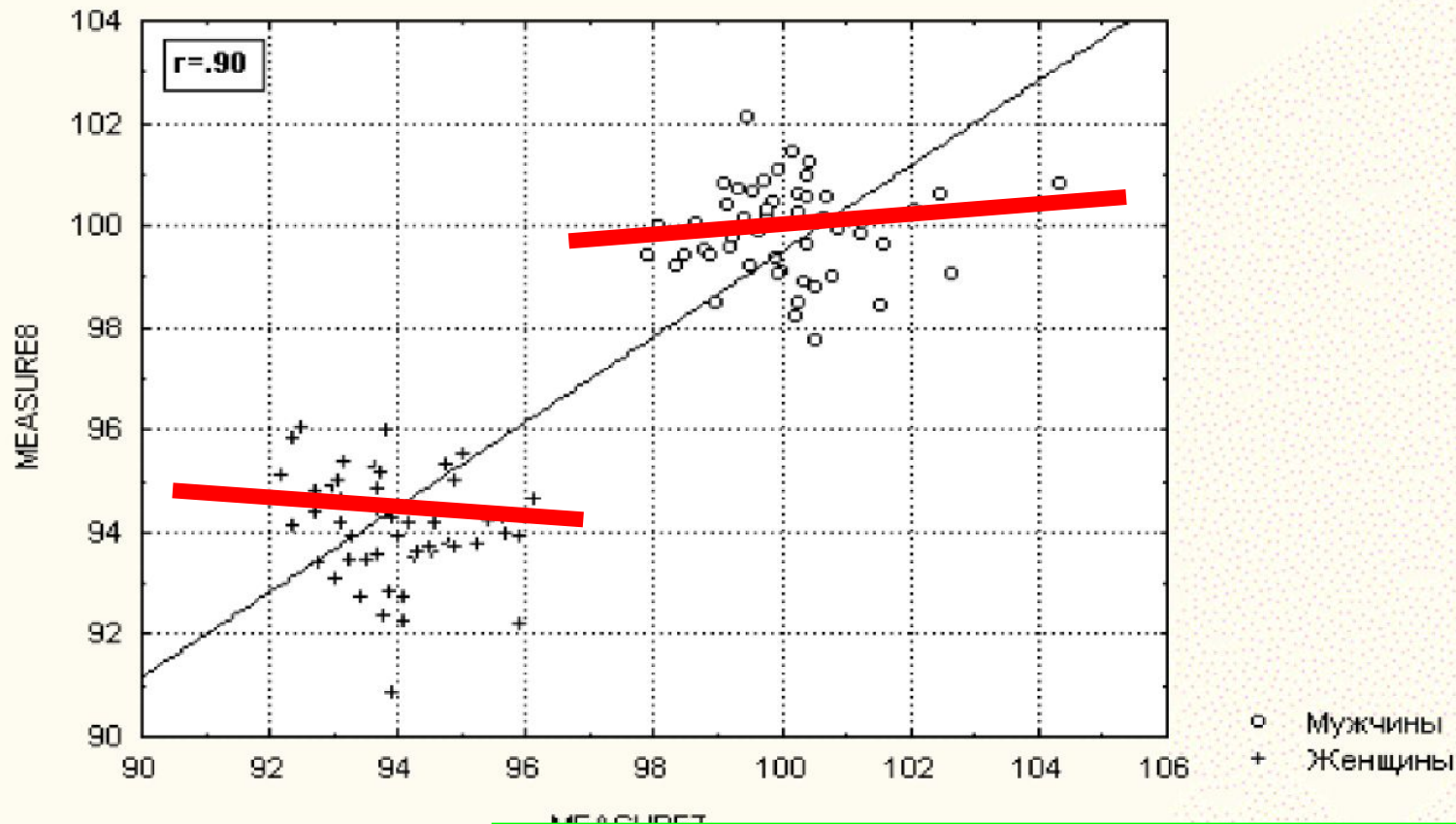
Например, при изучении реакции в когнитивной психологии, даже если почти все значения экспериментальных данных лежат в диапазоне 300-700 *миллисекунд*, то несколько “странных времен реакции” 10-15 *секунд* совершенно меняют общую картину.

К сожалению, определение выбросов субъективно, и решение должно приниматься индивидуально в каждом эксперименте (с учетом особенностей эксперимента и/или “сложившейся практики” в данной области). Следует заметить, что в некоторых случаях относительная частота выбросов к численности групп может быть исследована и разумно проинтерпретирована с точки зрения самой организации эксперимента.

# Осторожно:

## корреляция в неоднородных группах !

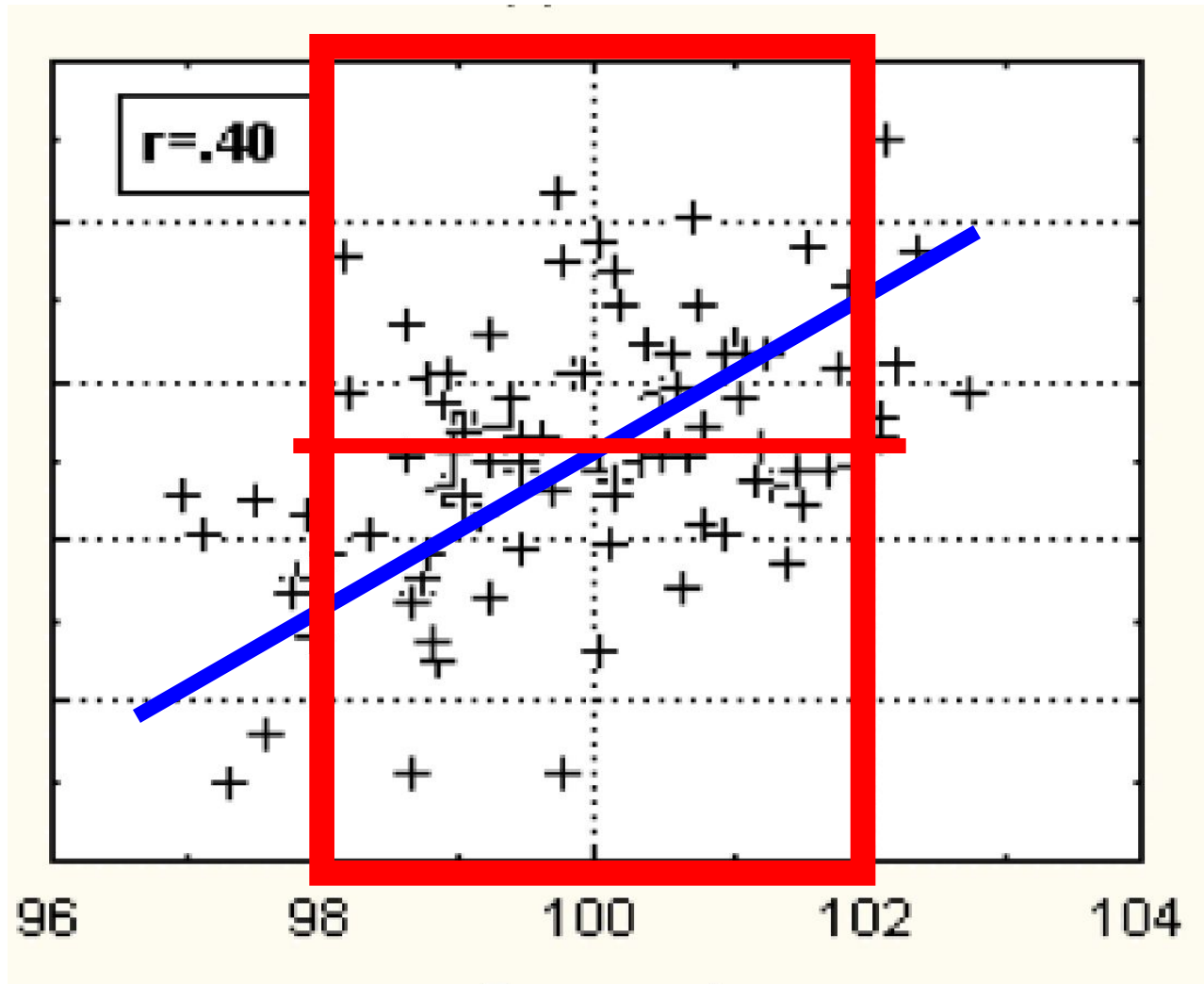
Диаграмма рассеяния (CORRS.STA.11п\*400н)



Пример\_тм\_токсичность\_преобразования.xls;

Сд\_хлорид \* Сд\_хлорид вместе и по

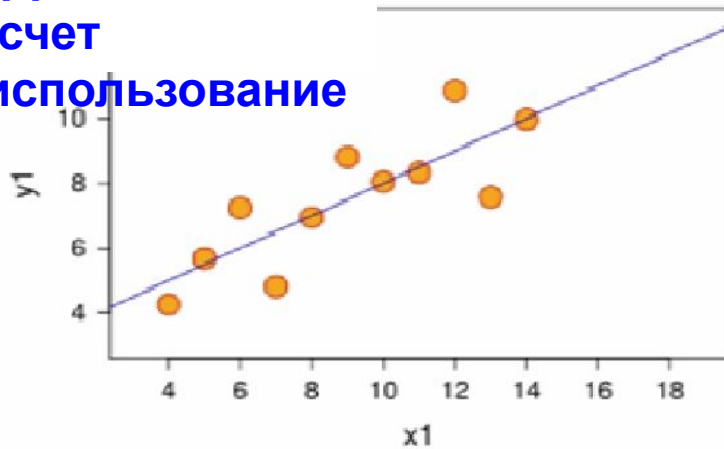
# Условие продуктивного использования коэффициентов корреляции: достаточная дисперсия данных



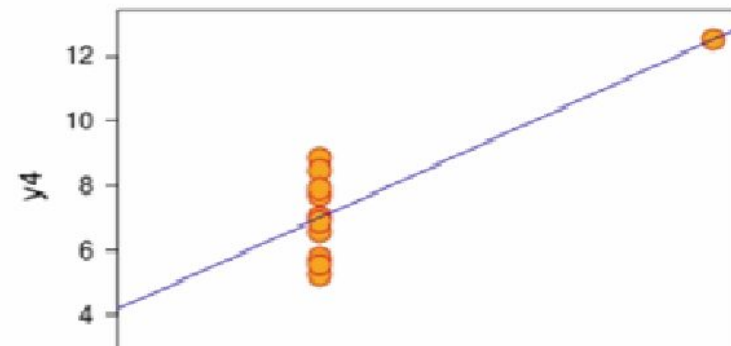
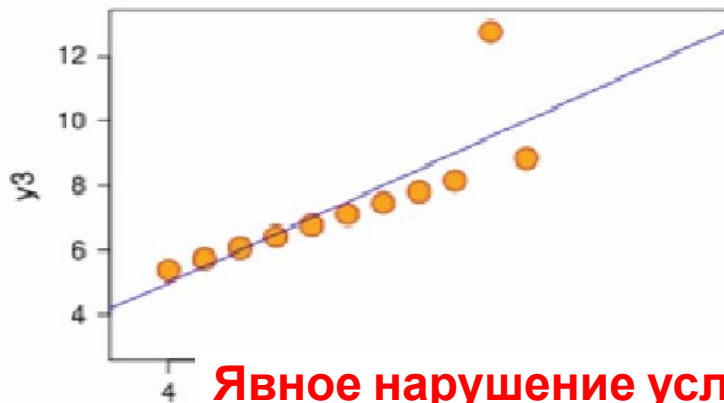
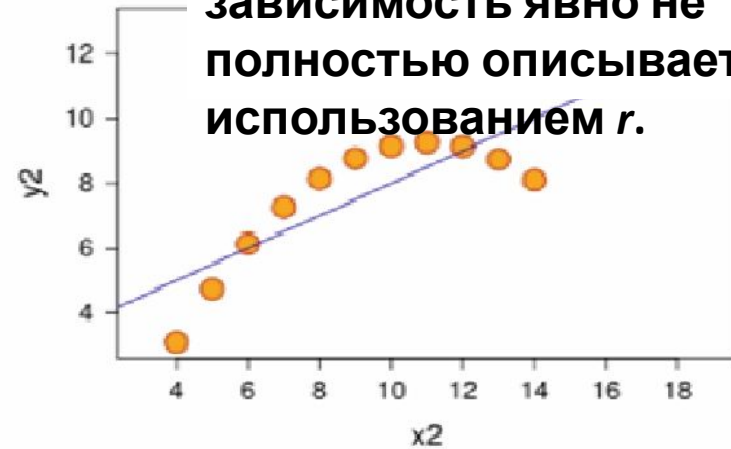
# Правильный/неправильный расчет и интерпретация

$r$  (во всех случаях  $r=0,816$  и  $P$  одинаковая)

Корректный расчет и использование



Так делать можно, но зависимость явно не полностью описывается с использованием  $r$ .

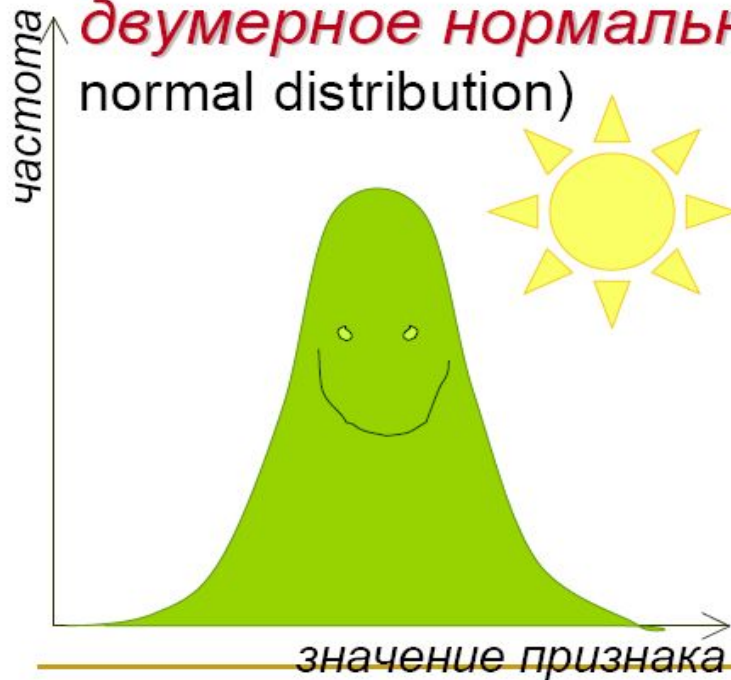


Явное нарушение условий использования  $r$ : «выбросы» и отклонение от нормального распределения

Требование к выборке для тестирования гипотезы о коэффициенте корреляции Пирсона:

Для каждого  $X$  значения  $Y$  должны быть распределены нормально, и для каждого  $Y$  все  $X$  должны иметь нормальное распределение -

**двумерное нормальное распределение** (bivariate normal distribution)



# Оперирование «пропущенными значениями» при расчете корреляционных матриц

The image shows a software interface with two overlapping dialog boxes. The background is a data table with columns of numerical values.

**Dialog Box 1: "Select the variables fo..."**

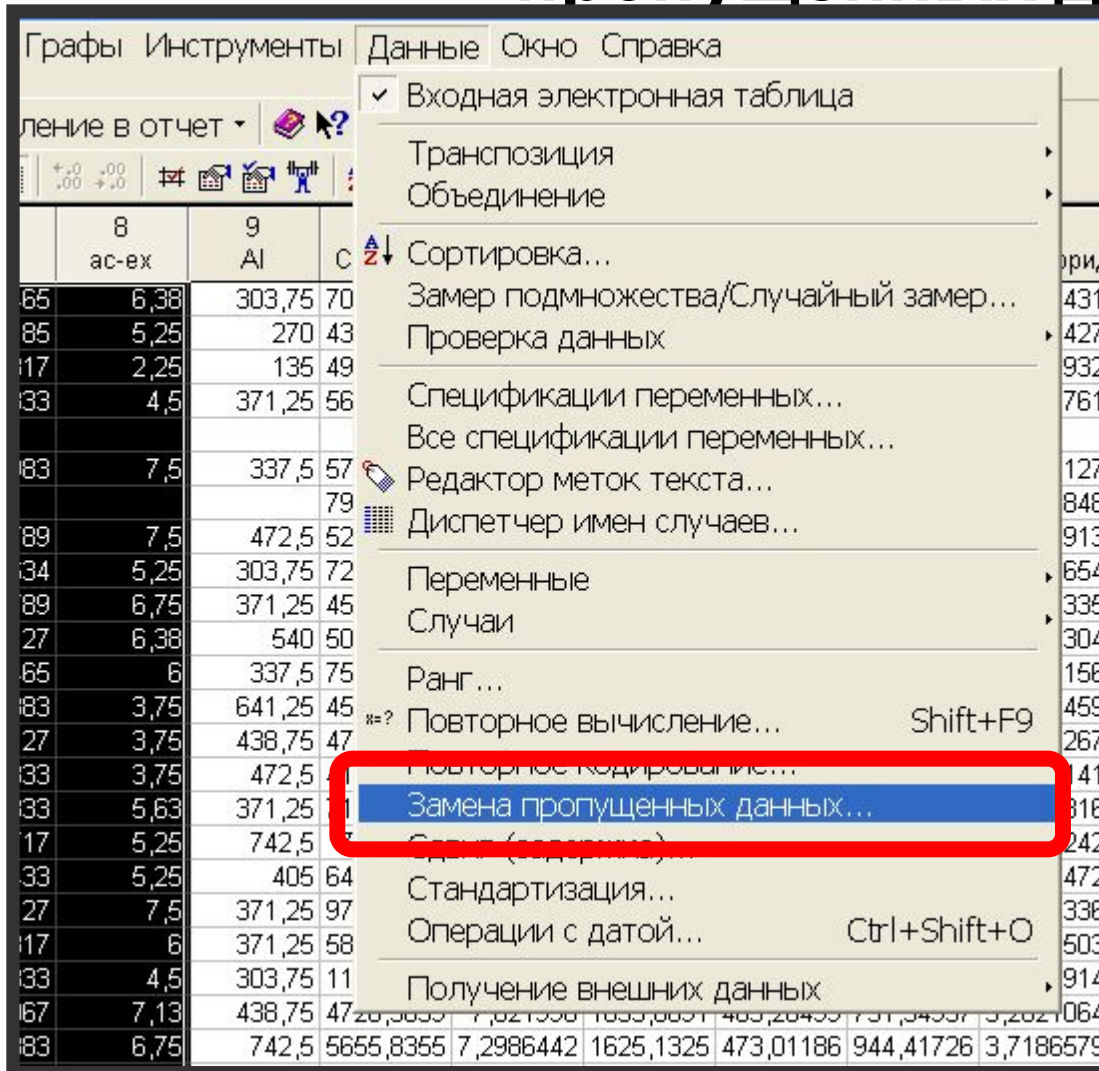
- Buttons: ? (Help), X (Close)
- List of variables:
  - 1-Зона
  - 2-Код1
  - 3-Код2
  - 4-Код2\_ проверка
  - 5-pH
  - 6-СЕС
  - 7-гидр
  - 8-ас-ех
  - 9-Аl
  - 10-Сu\_кисл
  - 11-Cd\_кисл
  - 12-Pb\_кисл
  - 13-Zn\_кисл
  - 14-Cu\_хлорид
  - 15-Cd\_хлорид
  - 16-Pb\_хлорид
  - 17-Zn\_хлорид
  - 18-ТОКСИЧНОСТ
  - 19-ТОКСИЧНОСТ
  - 20-длина1
- Buttons: ОК, Отмена
- Buttons: <, |||, >
- Text: "Выбор переменных:"
- Text: "5-12"
- Buttons: "Выбор всего", "Spread", "Zoom"

**Dialog Box 2: "Произведения и корреляции"**

- Buttons: ? (Help), - (Minus), X (Close)
- Text: "Summary"
- Buttons: "Отмена", "Параметры"
- Buttons: "Matrix", "Matrix"
- Buttons: "SELECT CASES", "W", "W"
- Text: "Взвешенные моменты"
- Text: "DF ="
- Buttons: "W-1", "N-1"
- Text: "Удаление MD"
- Buttons: "По случай", "Попарны"



# Замена пропущенных значений средними: ВОЗМОЖНОСТЬ МИНИМИЗИРОВАТЬ УЩЕРБ ОТ ПРОПУЩЕННЫХ ДАННЫХ



Пример\_тм\_токсич  
ность\_преобразов  
ания.xls;

Операции с  
переменными:

«СЕС»

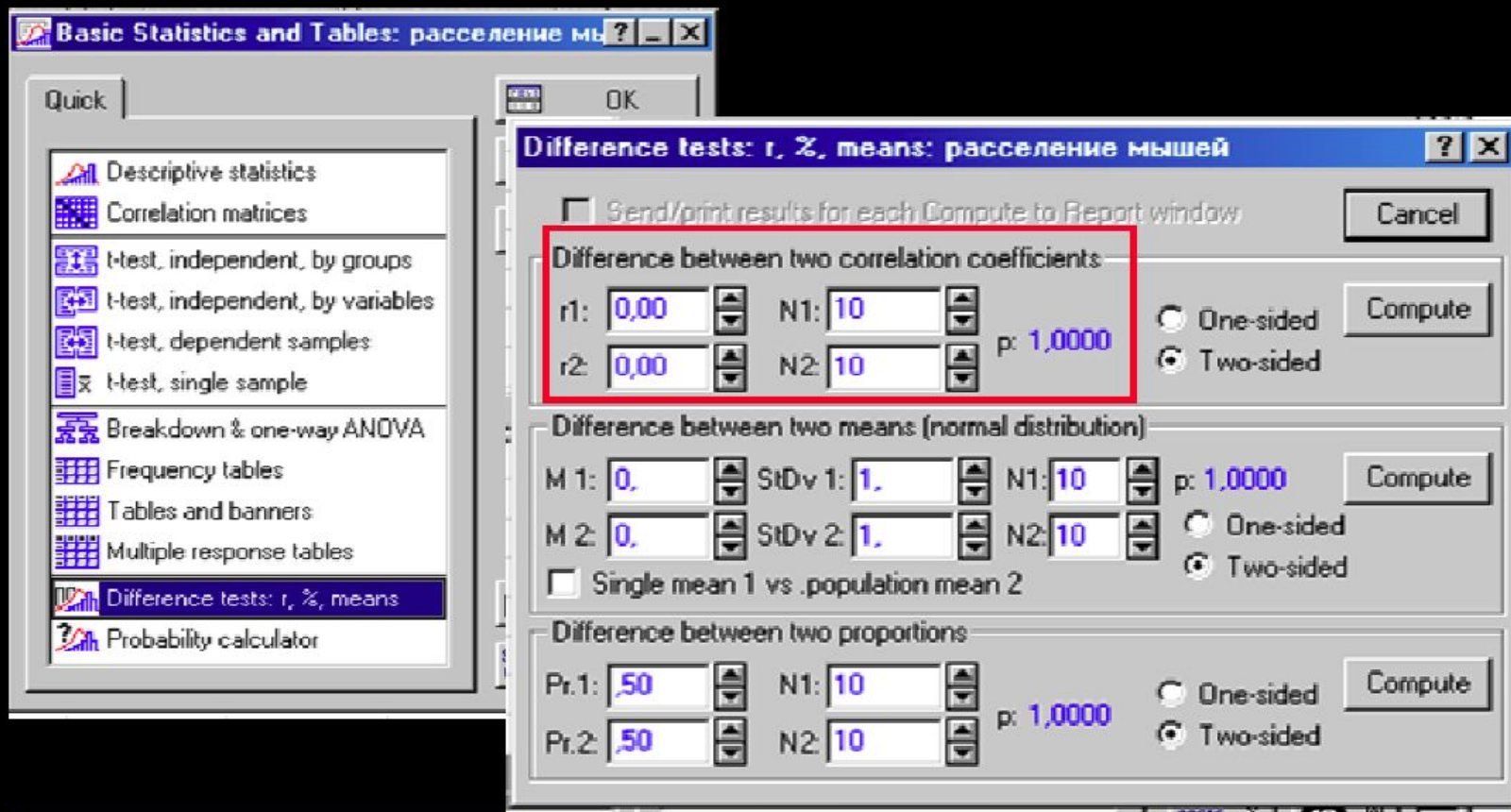
«Hidr»

«ac-ех»

«AI»

С учетом «зон»!

Можно сравнить два коэффициента корреляции от двух выборок

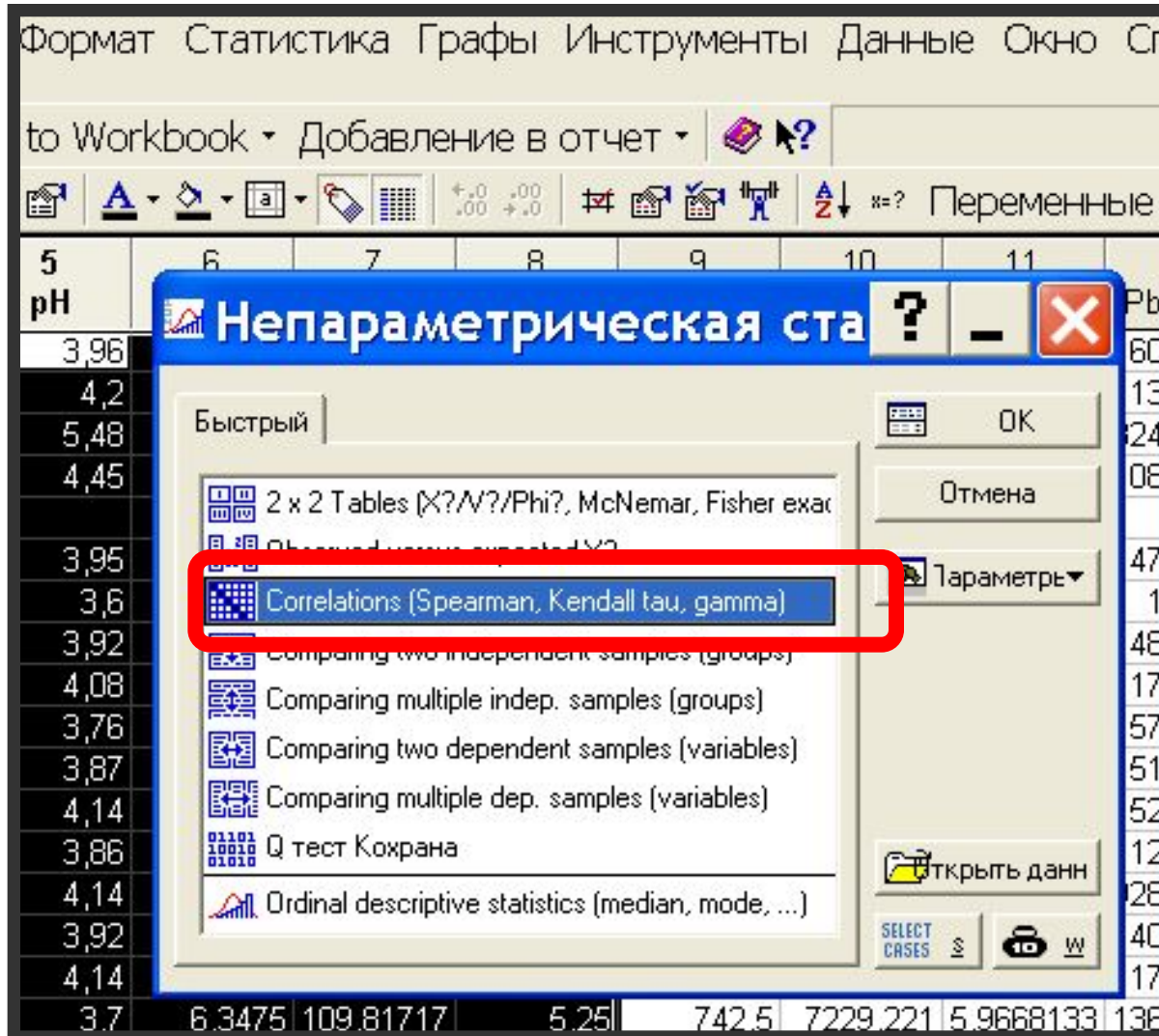


Для двумерного нормального распределения

# Непараметрическая корреляция

Задача	Количественные нормально распределенные переменные	Количественные ненормально распределенные переменные или ранги
Оценить взаимосвязь между двумя переменными	Коэффициент парной корреляции Пирсона	Коэффициенты ранговых корреляций (Спирмена, Кендалла)

# Коэффициент корреляции Спирмена



аналог  
коэффициента  
Пирсона;

подходит для  
расчета  
корреляционных  
матриц;

Размер выборки:  
>10.

# Линейная (парная) регрессия

**Задача:** предсказать значение одной переменной на основании другой на основе аппроксимации – линии.

**Переменные:** зависимая ( $Y$ ) и независимая ( $X$ ).

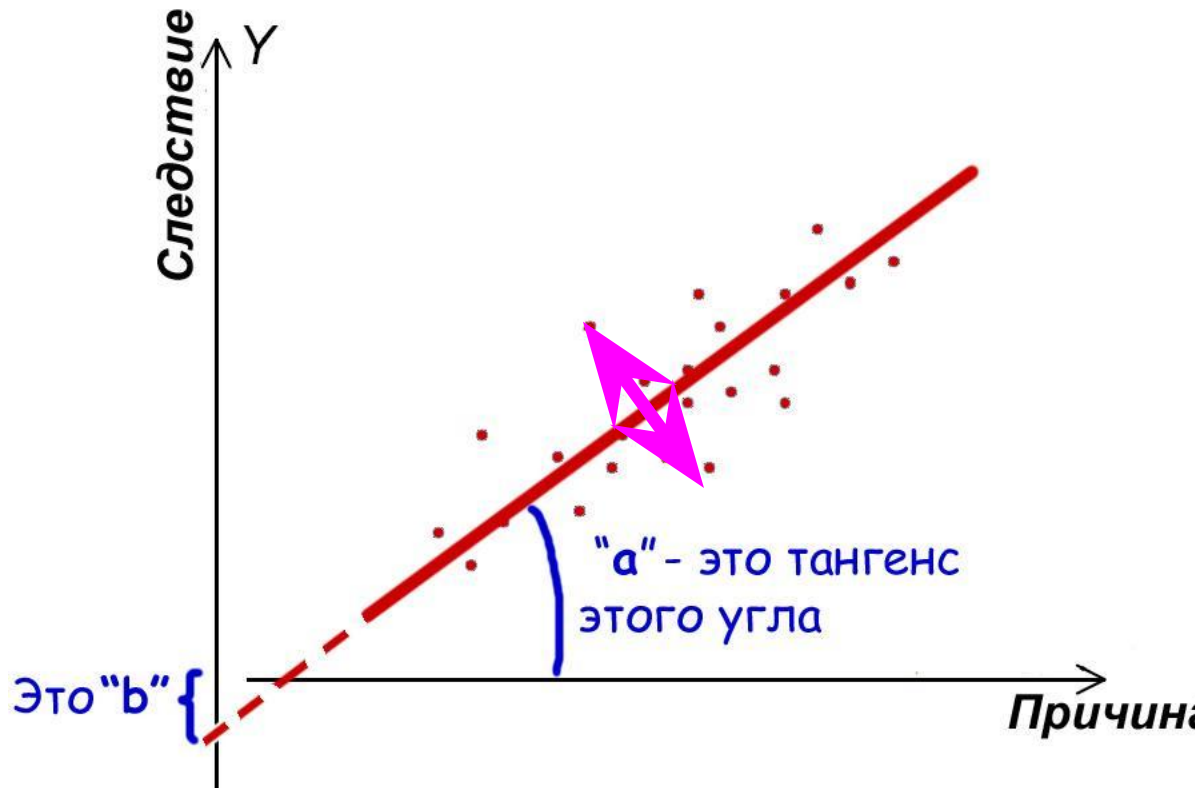
**Предположения:**

- линейная зависимость между переменными;
- независимость измерений отдельных  $X$  и  $Y$  от других измерений  $X$  и  $Y$ ;
- двумерное нормальное распределение и нормальное распределение «остатков», т.е. разностей между наблюдаемыми и предсказываемыми величинами  $Y$ .
- 

**Интерпретация** (при правильной постановке вопроса и правильном расчете): причинная и объясняющая.

**Формальное выражение:**  $Y = aX + b$ .

$$Y = aX + b$$

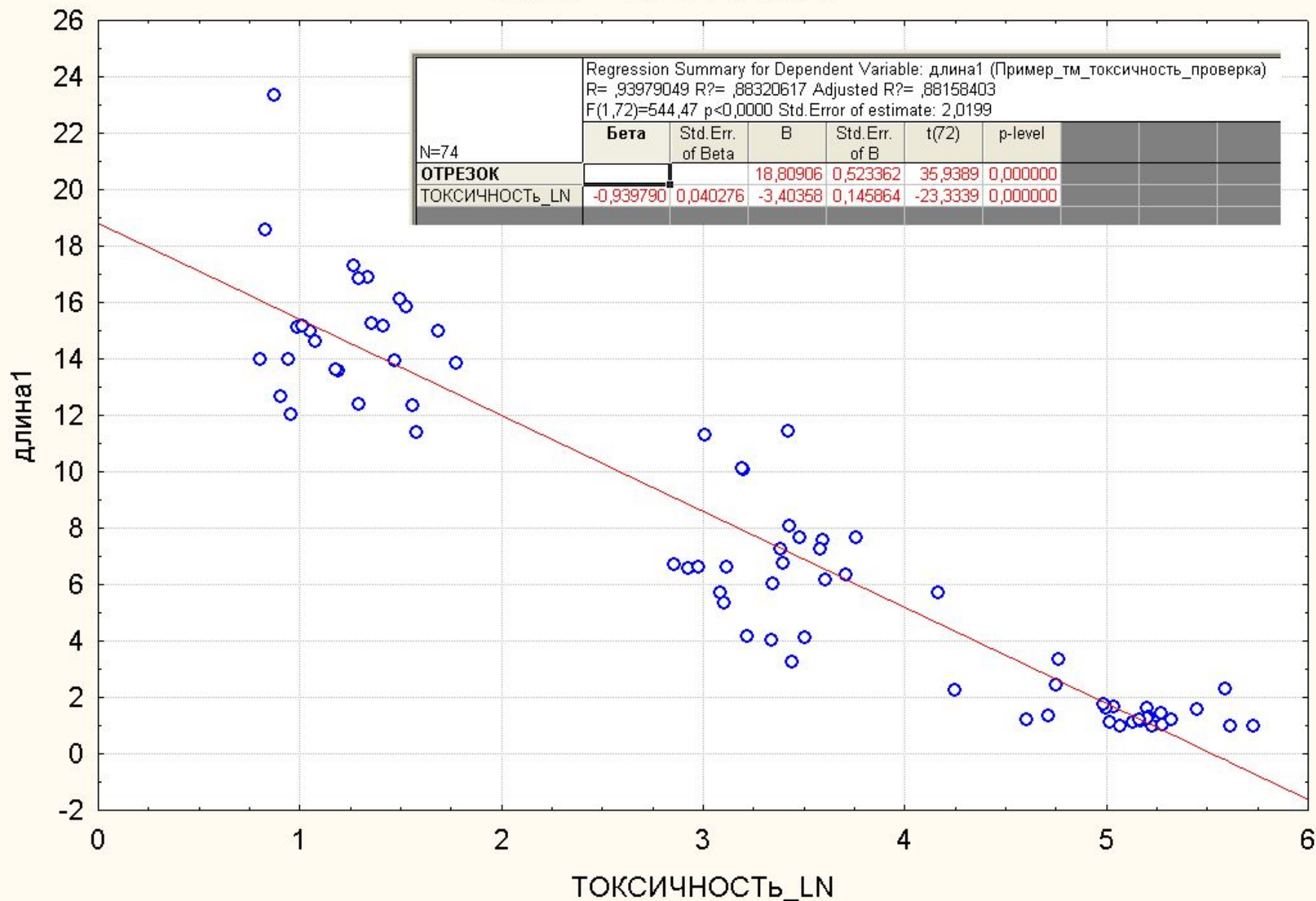


Изменчивость данных возле линии регрессии характеризует параметр  $R^2$  – простой квадрат коэффициента корреляции Пирсона (в случае линейной регрессии).

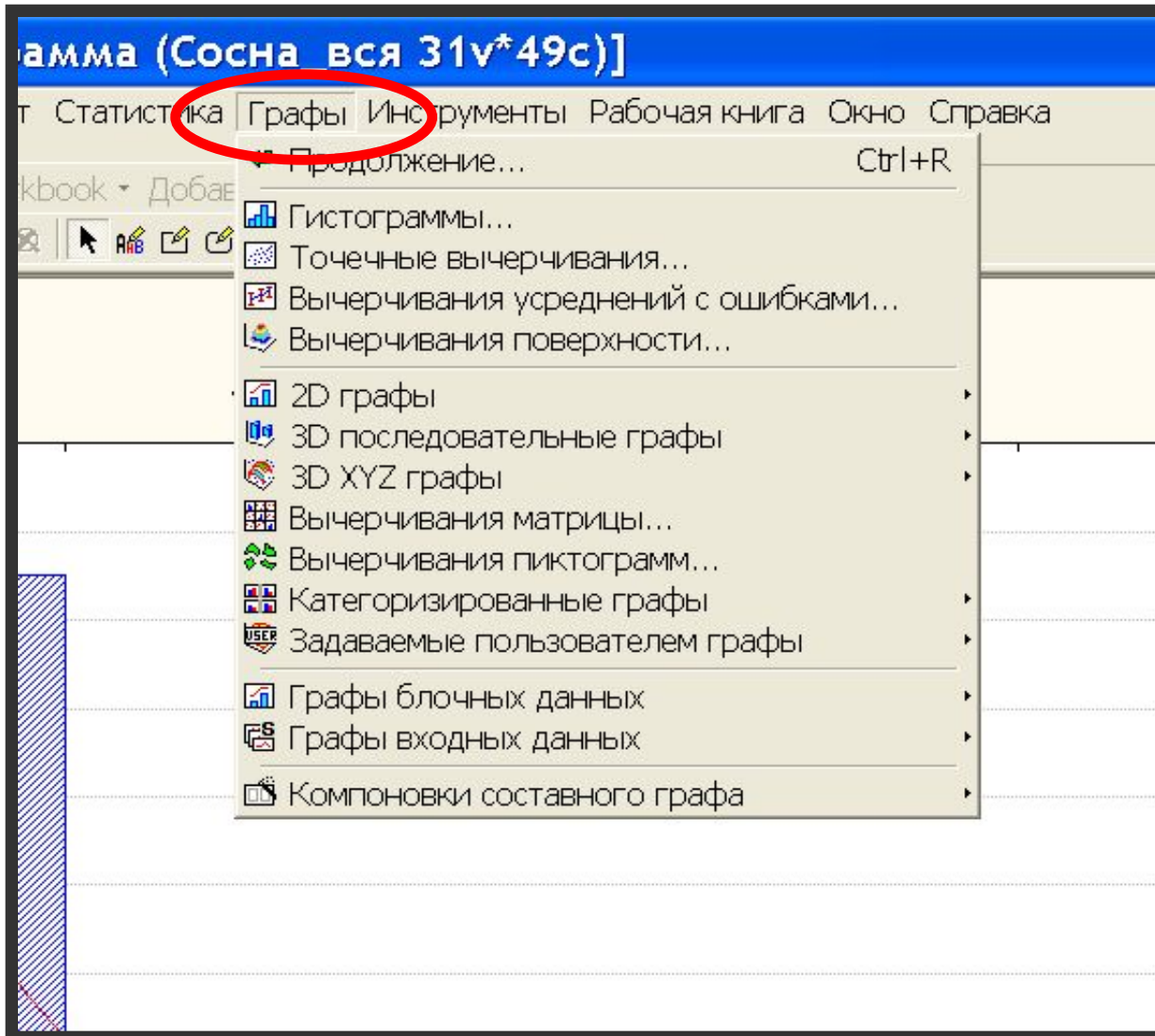
**НО!!!**  
Показатель  $R^2$  приемлем и для нелинейных и для множественных зависимостей. Интерпретируется он как.....?

# Scatterplot (Пример\_тм\_токсичность\_проверка 31v\*75c)

$$\text{длина1} = 18,8091 - 3,4036 * x$$



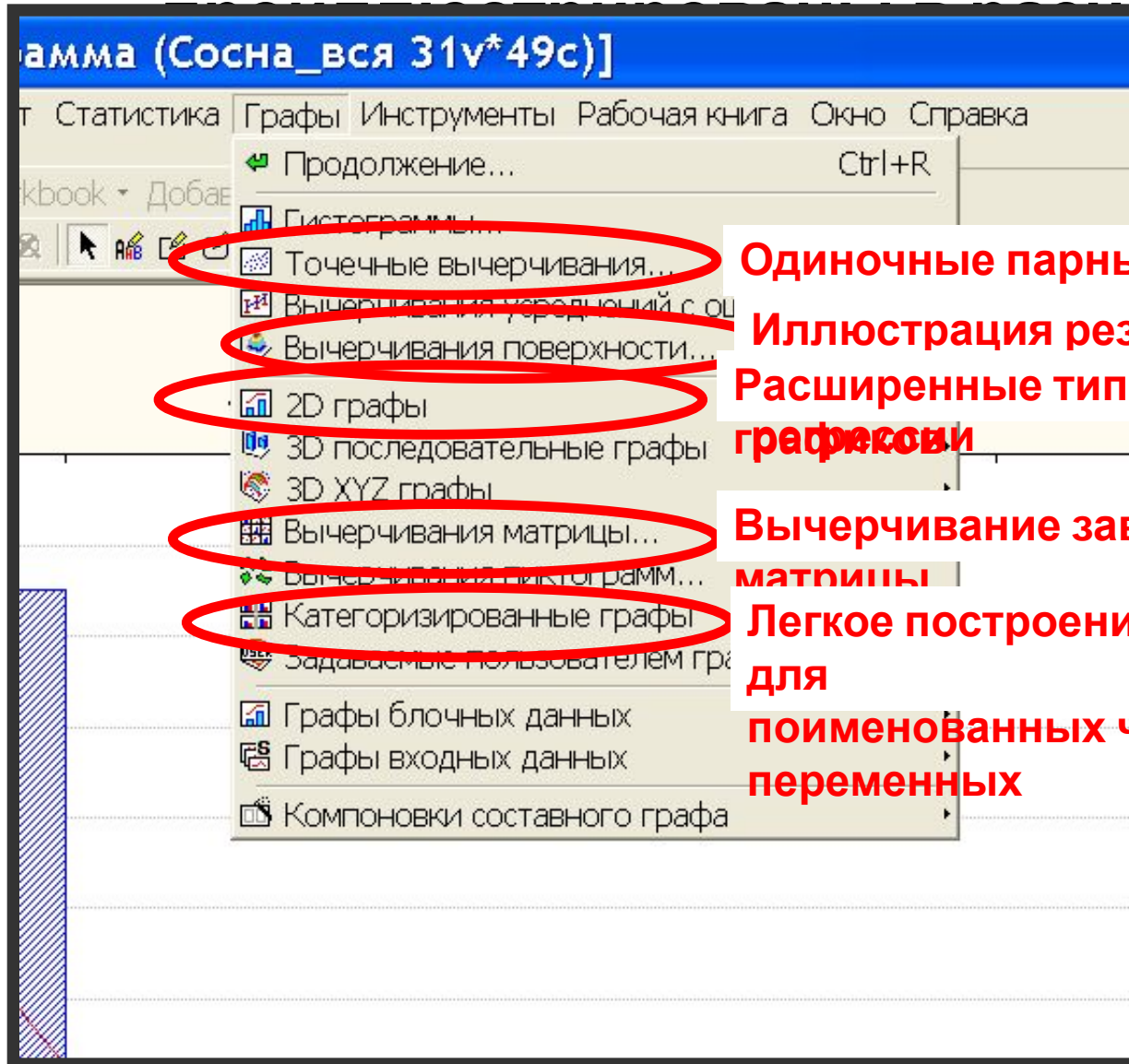
**Пункт меню «Графы»: ШИРОЧАЙШИЕ  
возможность построения диаграмм как без  
расчета статистик,  
так и с расчетом таковых**





# Корреляционные и регрессионные зависимости могут быть

ых пунктах:



Одиночные парные

Иллюстрация результатов

Расширенные типы

графов

Вычерчивание зависимостей в виде

матрицы

Легкое построение зависимостей

для

поименованных частей

переменных

# STATISTICA™

## Том II: ГРАФИКА

1. Введение .....	2001	12. Блочные статистические графики .....	2761
2. Примеры .....	2015	13. Размещение нескольких графиков .....	2769
3. Опции, общие для всех графиков .....	2111	14. Пустые графические окна ..	2781
4. Основные типы графиков..	2291	15-19. Пользовательские графики:	
5. Быстрые статистические графики .....	2421	15. Двумерные графики .....	2791
6-11. Статистические графики:		16. Трехмерные последовательные графики.....	2811
6. Двумерные графики .....	2457	17. XYZ графики.....	2825
7. Трехмерные последовательные графики.....	2567	18. Матричные графики .....	2839
8. XYZ графики.....	2603	19. Пиктографики .....	2851
9. Матричные графики .....	2663	20. Статистические графики пользователя .....	2861
10. Пиктографики .....	2679	21. Примечания .....	2875
11. Категоризованные графики.....	2693	Литература.....	2901

# Настройка вида графиков: Инструменты → Параметры → Графы

The screenshot shows a software interface with a menu bar and a data table. The 'Инструменты' menu is open, and the 'Параметры...' option is selected. The 'Параметры' dialog box is open, showing the 'Графы 1' tab. The dialog contains a grid of line styles and colors for 10 series, and a section for line properties.

Series	Line Style	Color
1	Solid blue	Red
2	Dashed red	Green
3	Dotted green	Blue
4	Long dashed purple	Yellow
5	Short dashed yellow	Magenta
6	Dash-dot cyan	Cyan
7	Long dash-short dash red	Dark Blue
8	Short dash-long dash green	Brown
9	Long dash-short dash blue	Olive
10	Long dash-short dash black	Dark Red

Line properties:  
Ширина: 0 points  
Размер: 5 points  
Размер: 10 points