

# **Chapter 6:** 95% Confidence Interval for the Mean and Data Transformation

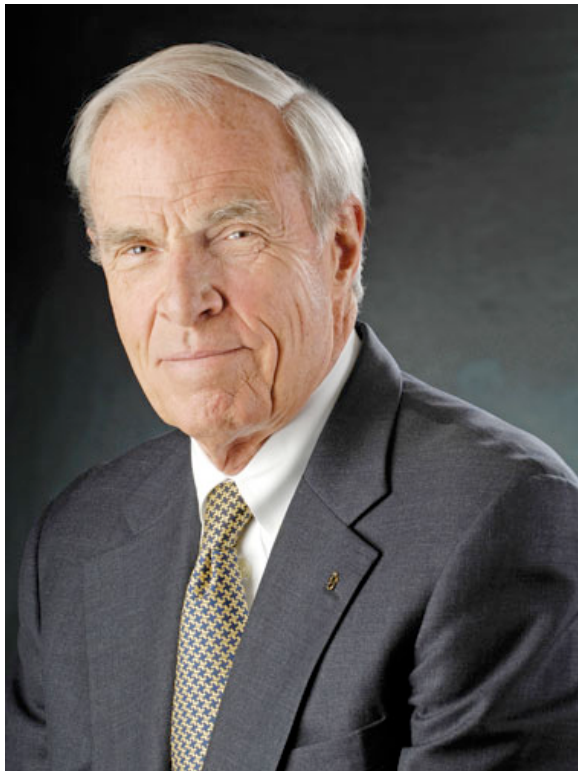
TXCL7565/PHSC7565

# What This Chapter Covers

- What is a confidence interval (CI)?
- What do we mean by '95%' confidence?
- Calculating a CI
- Sensitivity of CI to SD, sample size, and level of confidence
- One-sided CIs
- CI for difference between two means
- Normal distribution and CI

WHAT IS A CONFIDENCE  
INTERVAL?

# Age Guess



Bruce Benson  
President of CU



Don Cheadle  
Actor  
Alumnus of East  
High School



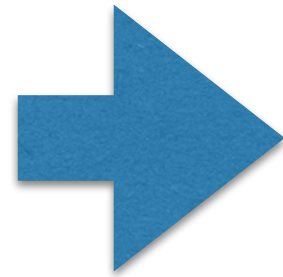
Stranger

# Confidence Interval

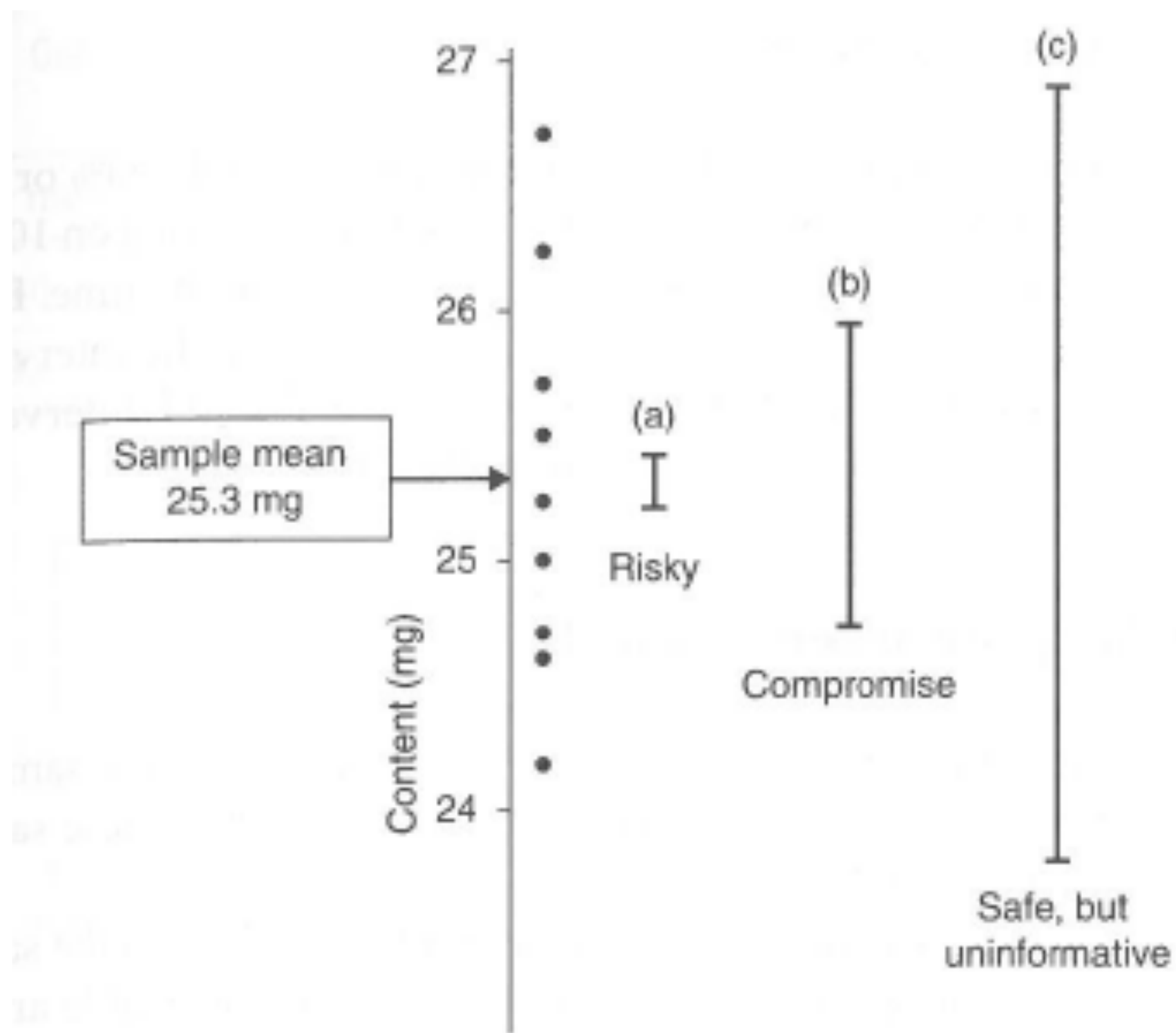
- A mean derived from a sample is unlikely to be a perfect estimate of the population mean.
- Therefore, we normally give a range that the population mean is likely to fall within.
- Since we have no reason to believe there is bias in our sample mean estimate, we create a range by subtracting and adding the same amount to the sample mean estimate.

# Width of the Interval

The **wider**  
the interval



The **more**  
**confident** we  
are that the true  
population  
mean falls  
within the  
interval



# Norm in the Field

- A '95% confidence interval' is the standard confidence interval width.
- It originates from a p-value threshold of 0.05 for significance.
- I'll explain more when we tackle p-values next week.



WHAT DO WE MEAN BY  
'95%' CONFIDENCE?

# 95% Confidence

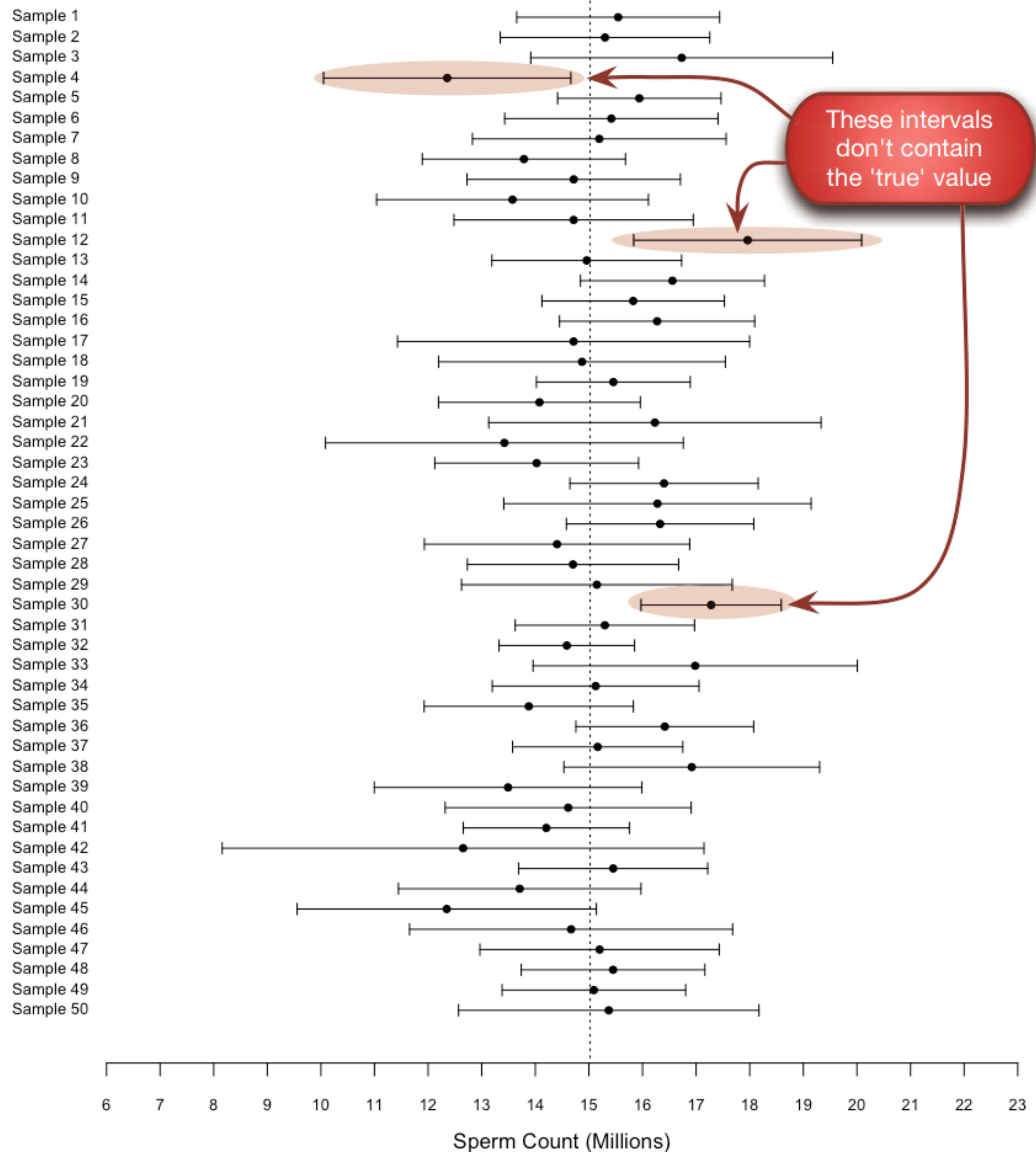
- For a 95% CI, we typically say that we are 95% confident that the interval includes the true population mean
  - CANNOT say that there is a 95% **probability** that the interval includes the true population mean

# Statistical Definition of Level of Confidence

**Level of Confidence** = the proportion of intervals that will cover the true population mean when the intervals are calculated in the same manner and generated from sampling of the same population

**FIGURE 2.9**

The confidence intervals of the sperm counts of Japanese quail (horizontal axis) for 50 different samples (vertical axis)

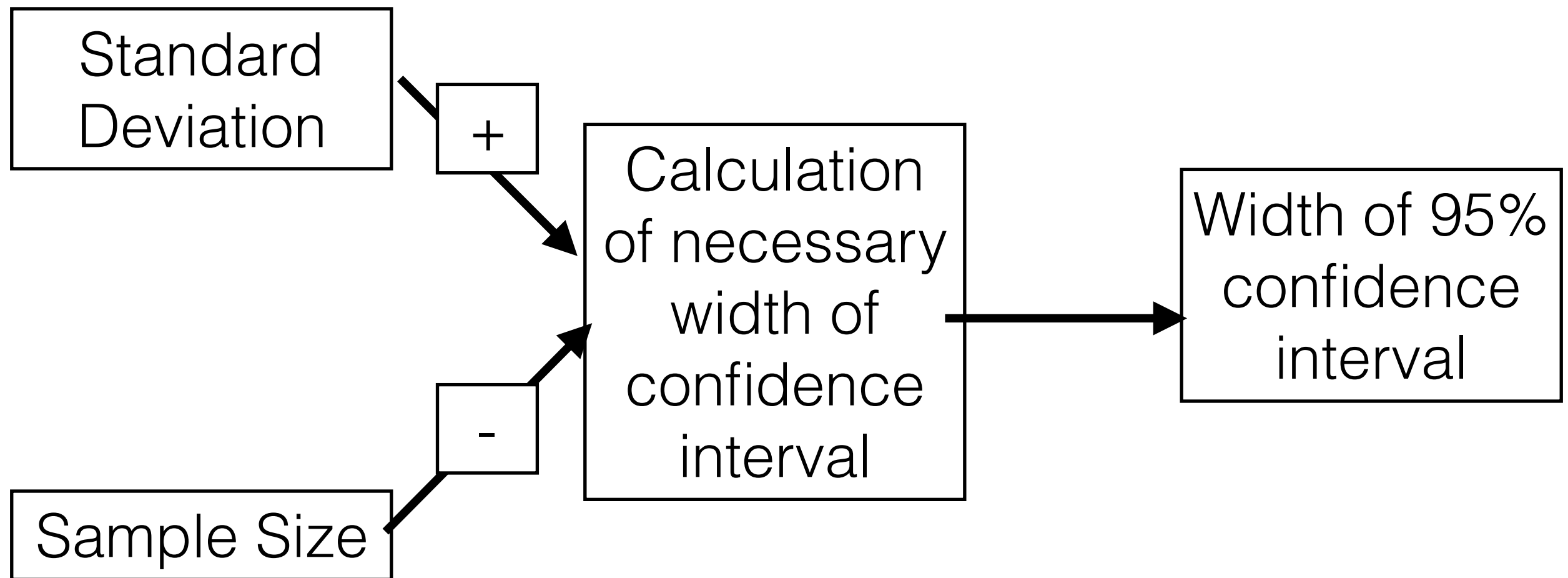


# Calculating a CI

The width of confidence intervals are influenced by:

1. Standard Deviation = higher SD, wider interval
2. Sample Size = smaller sample size, wider interval
3. Level of Confidence = higher level of confidence, wider interval

# Calculating 95% CI



# Formula for CI

$$\bar{x} \pm t_c \frac{s}{\sqrt{n}}$$

$\bar{x}$  = sample mean

$t_c$  = t-statistic

$s$  = sample standard deviation

$n$  = number of observations in sample

1. Click  
'Analyze  
icon'

2. Select  
'Column  
statistics'

3. Click OK

GraphPad PRISM

Create New Analysis

Data to analyze  
Table: imipramine

Type of analysis  
Which analysis?

- Transform, Normalize...
  - Transform
  - Transform Concentrations (X)
  - Normalize
  - Prune rows
  - Remove baseline and column math
  - Transpose X and Y
  - Fraction of Total
- XY analyses
- Column analyses
  - t-tests (and nonparametric tests)
  - One-way ANOVA (and nonparametric)
  - Column statistics
  - Frequency distribution
  - ROC Curve
  - Bland-Altman method comparison
  - Correlation
  - Identify outliers
  - Analyze a stack of P values
- Grouped analyses
- Contingency table analyses
- Survival analyses

Analyze which data sets?  
☒ A:imipramine

When you analyze tables or graphs with more than one data set, use this space to select which data set(s) to analyze.

Select All Deselect All

Cancel OK

	Group A	Group B
	imipramine	Title
	Y	Y
1	24.7	
2	25.8	
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		
21		
22		
23		
24		
25		

imipramine Row 20, Column B



1. Select  
'CI of  
the  
mean'

2. You can  
change the  
level of  
confidence

Parameters: Column Statistics

**Descriptive Statistics**

- ☒ Minimum and maximum
- ☒ Quartiles (Median, 25th and 75th percentile)
- ☐ Percentile 90
- ☒ Mean, SD, SEM
- ☐ Coefficient of variation
- ☐ Geometric mean
- ☐ Skewness and kurtosis
- ☒ Column sum

**Confidence intervals**

- ☒ CI of the mean
- ☐ CI of geometric mean
- ☐ CI of median

**Test if the values come from a Gaussian distribution**

- ☐ D'Agostino-Pearson omnibus normality test (recommended)
- ☐ Shapiro-Wilk normality test
- ☐ Kolmogorov-Smirnov test with Dallal-Wilkinson-Lilliefors P value (not recommended)

**Inferences**

- ☐ One-sample t test. Are column means significantly different than a hypothetical value?
- ☐ Wilcoxon signed-rank test. Compare column medians to a hypothetical value.

Hypothetical value (often 0.0, 1.0 or 100) 0

When a value equals the hypothetical value: Ignore that value entirely, as Prism 5 and earlier versions did

**Calculations**

Subcolumns: Compute the mean of the subcolumns for each row, and then calculate column statistics of those means

**Output**

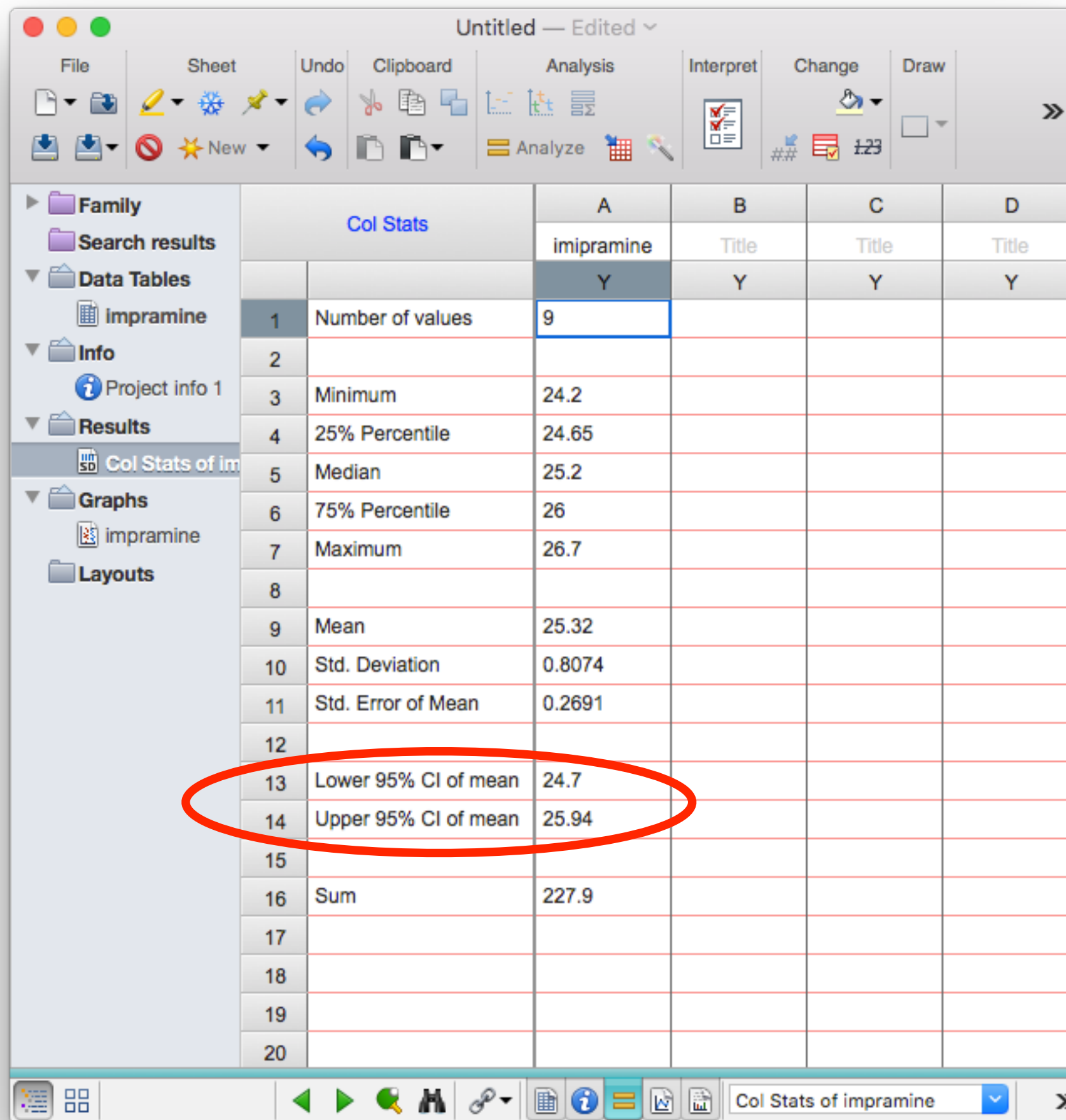
P-value style: GP: 0.1234 (ns), 0.0332 (\*), 0.0021 (\*\*), 0.0002 (\*\*\*), <0.0001 (\*\*\*\*)

Show 4 significant digits.

☐ Make these choices be the default for future analyses.

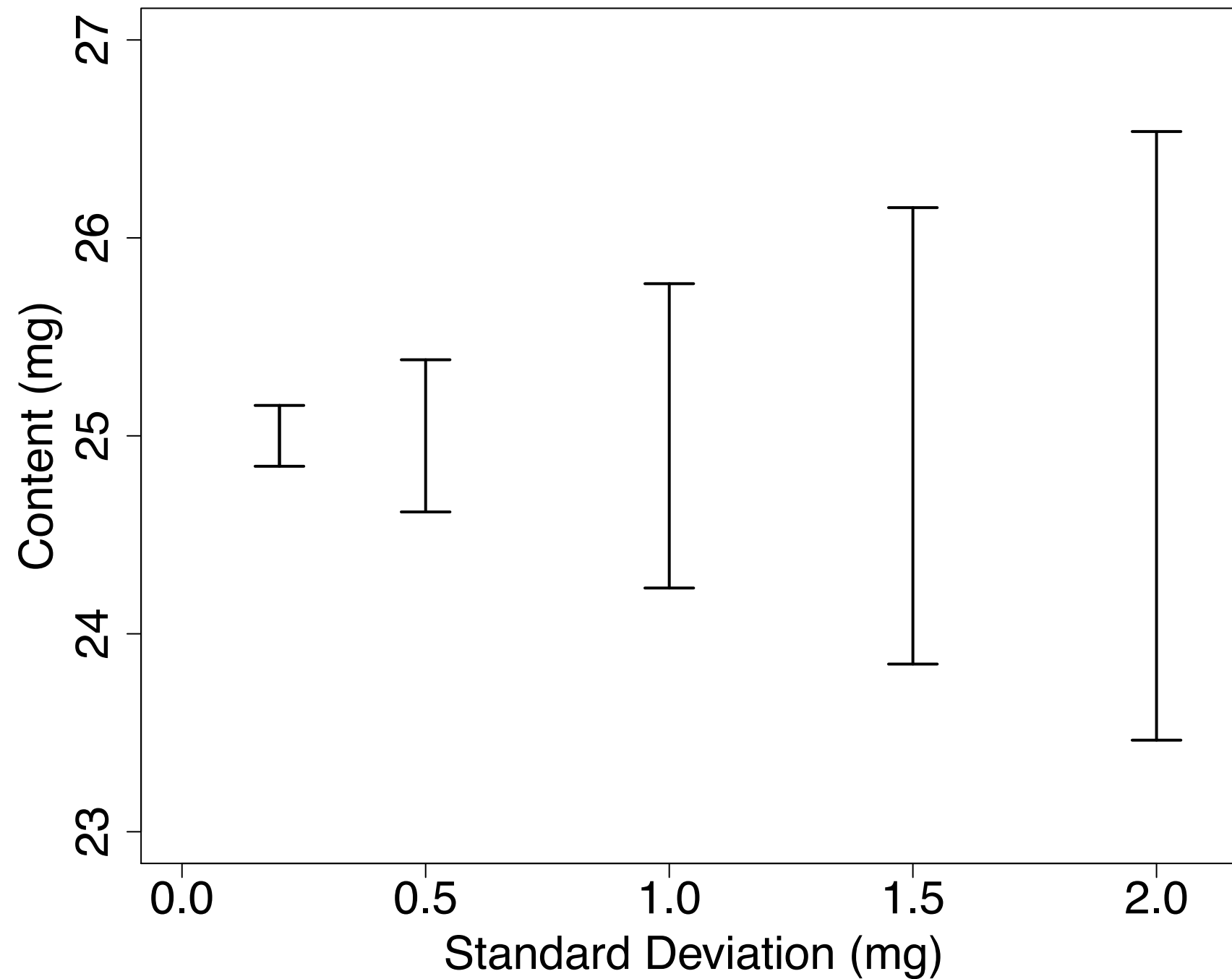
Cancel OK

3. Click OK

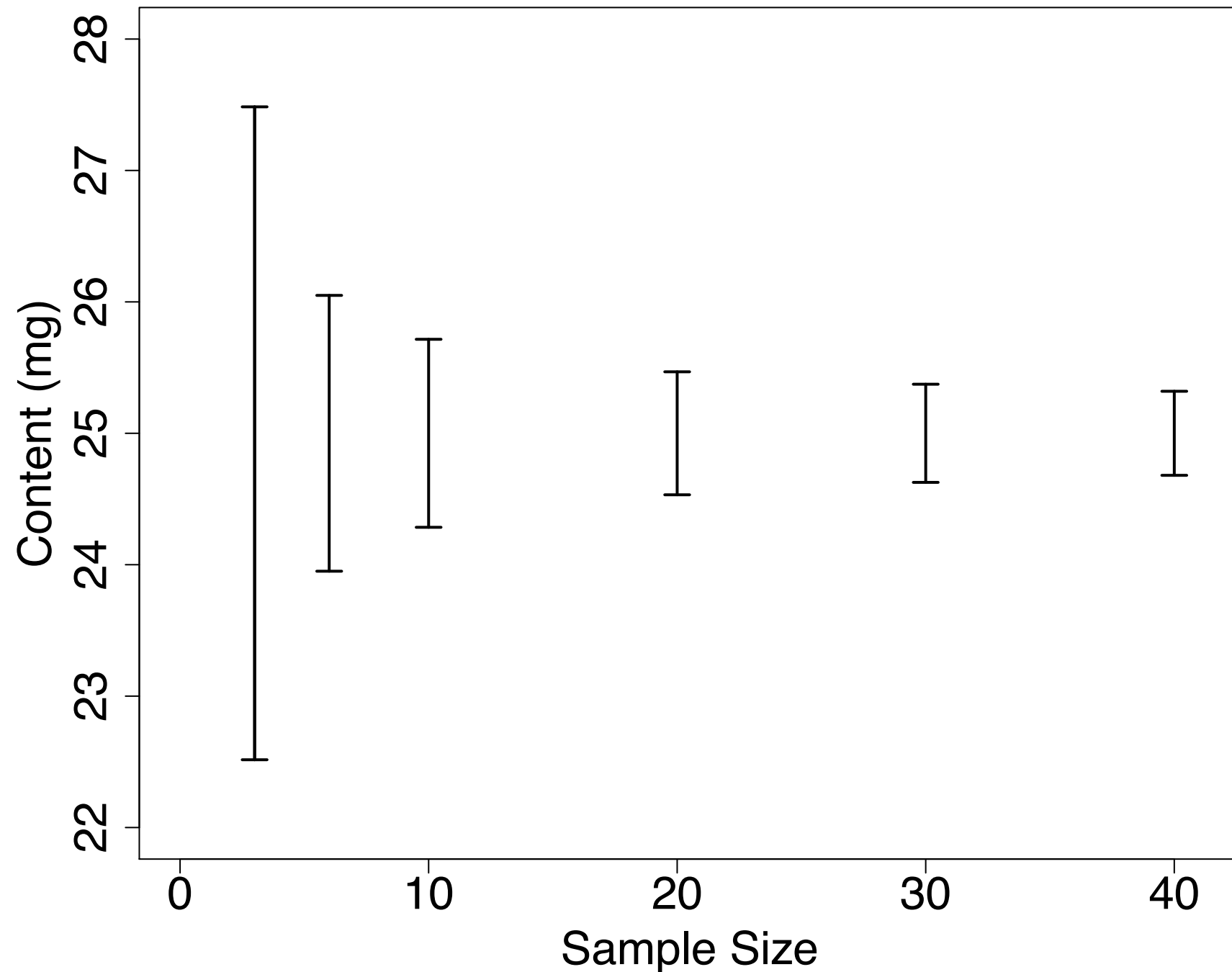


# SENSITIVITY OF CI TO SD, SAMPLE SIZE, AND LEVEL OF CONFIDENCE

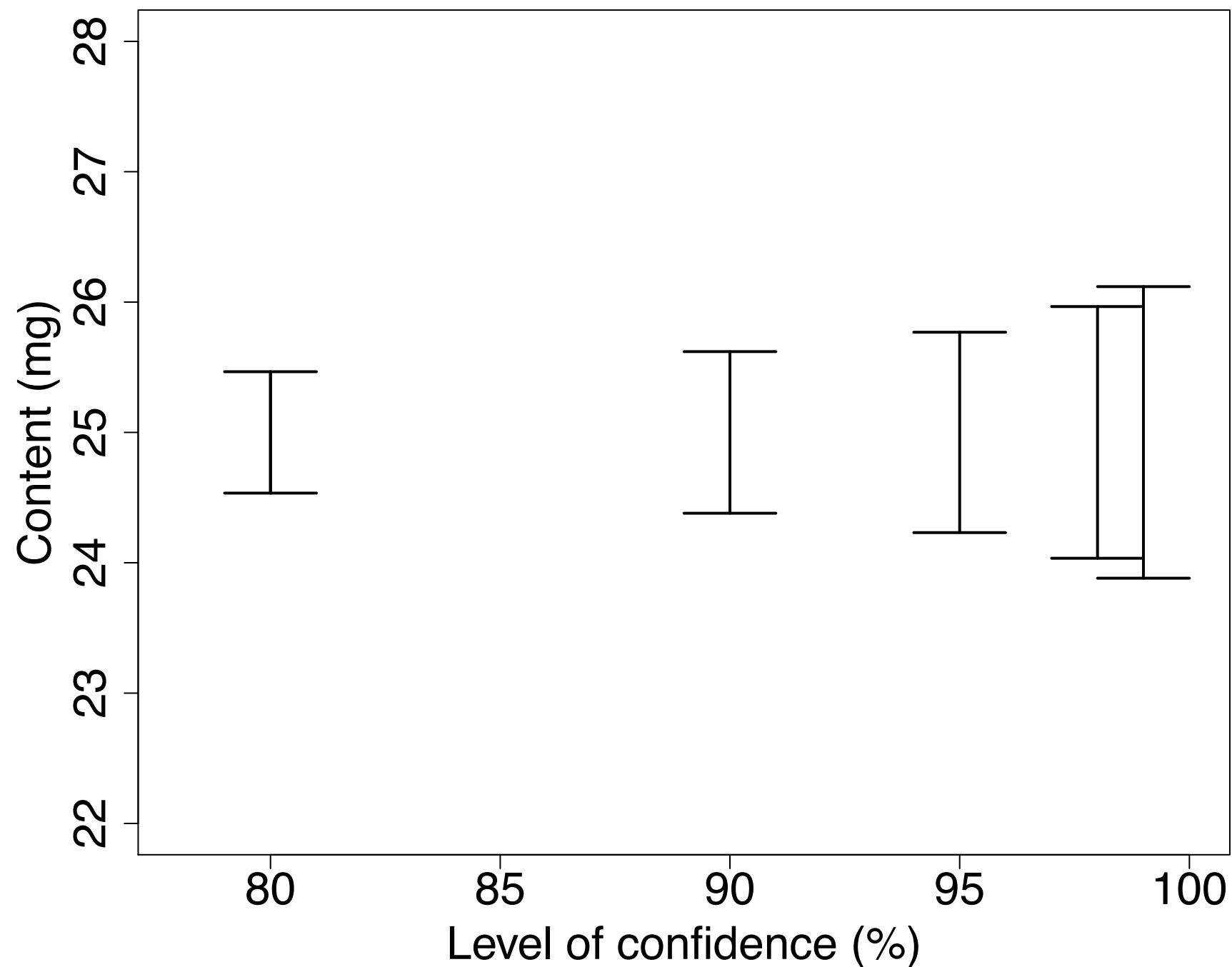
# Sensitivity to SD



# Sensitivity to Sample Size



# Sensitivity to Level of Confidence



# Assessing the reliability of our conclusions

Statistics never draws any absolute conclusions. It will offer an opinion and then back that up with a measure of that conclusion's reliability.

One-sided 95% CIs



# Two-sided / Two-tailed CIs

- In a two-sided CI, we specify both a minimum and maximum value to our range.
- The 'hazard', i.e., 1 - level of confidence, is split between the two ends.

1. The true population mean is no *less* than some stated figure (2.5% chance this is false)
2. The true population mean is no *greater* than some stated figure (2.5% chance this is false)

# One-sided 95% CIs

- Sometimes we are only interested in whether the mean is above a certain level

The true population mean is no *less* than some stated figure (5% chance this is false)

OR

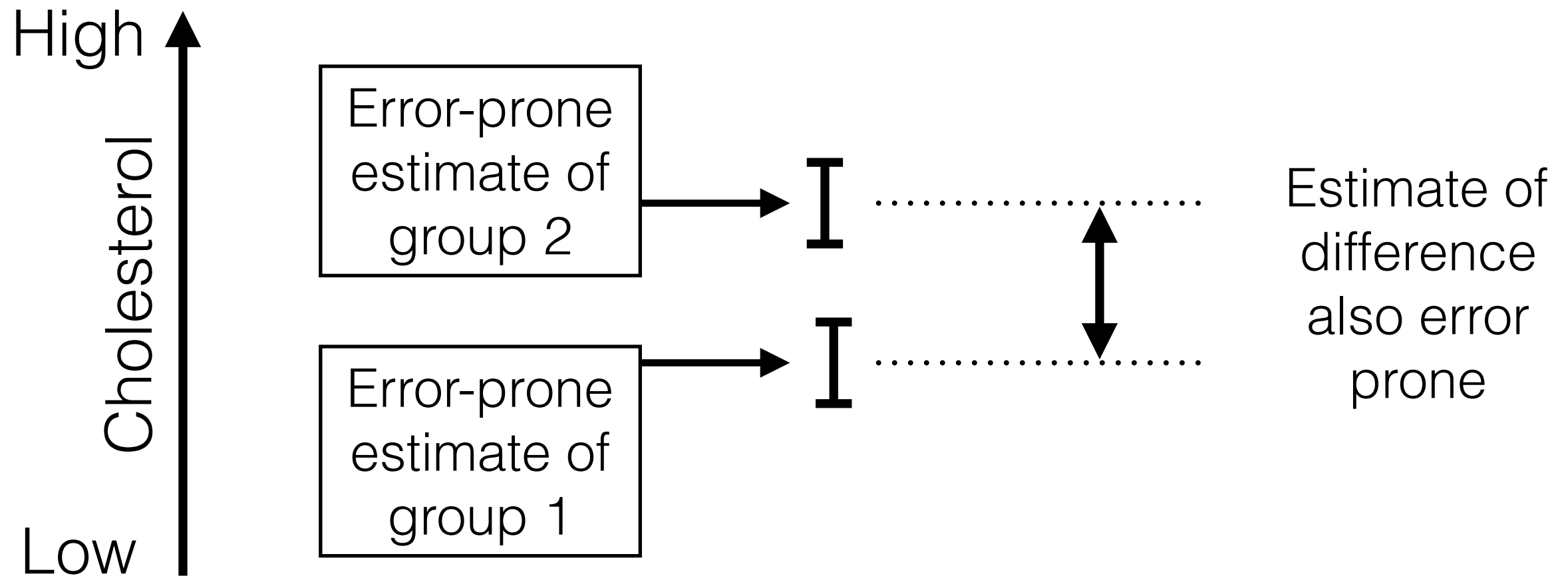
The true population mean is no *greater* than some stated figure (5% chance this is false)

# One-sided CIs in GraphPad

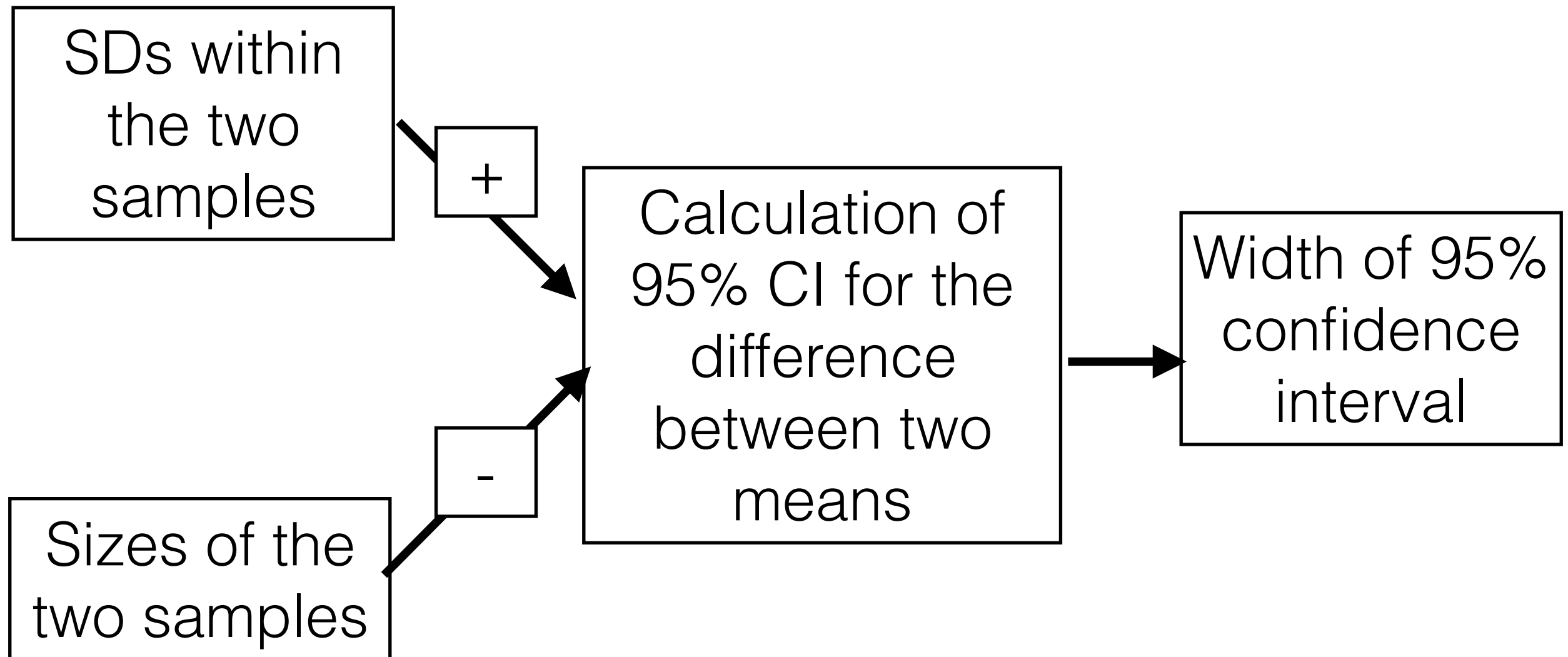
- GraphPad doesn't offer a one-sided confidence interval options
- Solution - Change the level of confidence by doubling the risk, i.e., to get a 95% one-sided confidence interval, ask for a 90% two-sided confidence interval

95% CI FOR THE  
DIFFERENCE BETWEEN  
TWO TREATMENTS

# 95% CI for Difference Between 2 Means



# 95% CI of Difference Between 2 Means



# Formula for 95% CI of Difference Between 2 Means

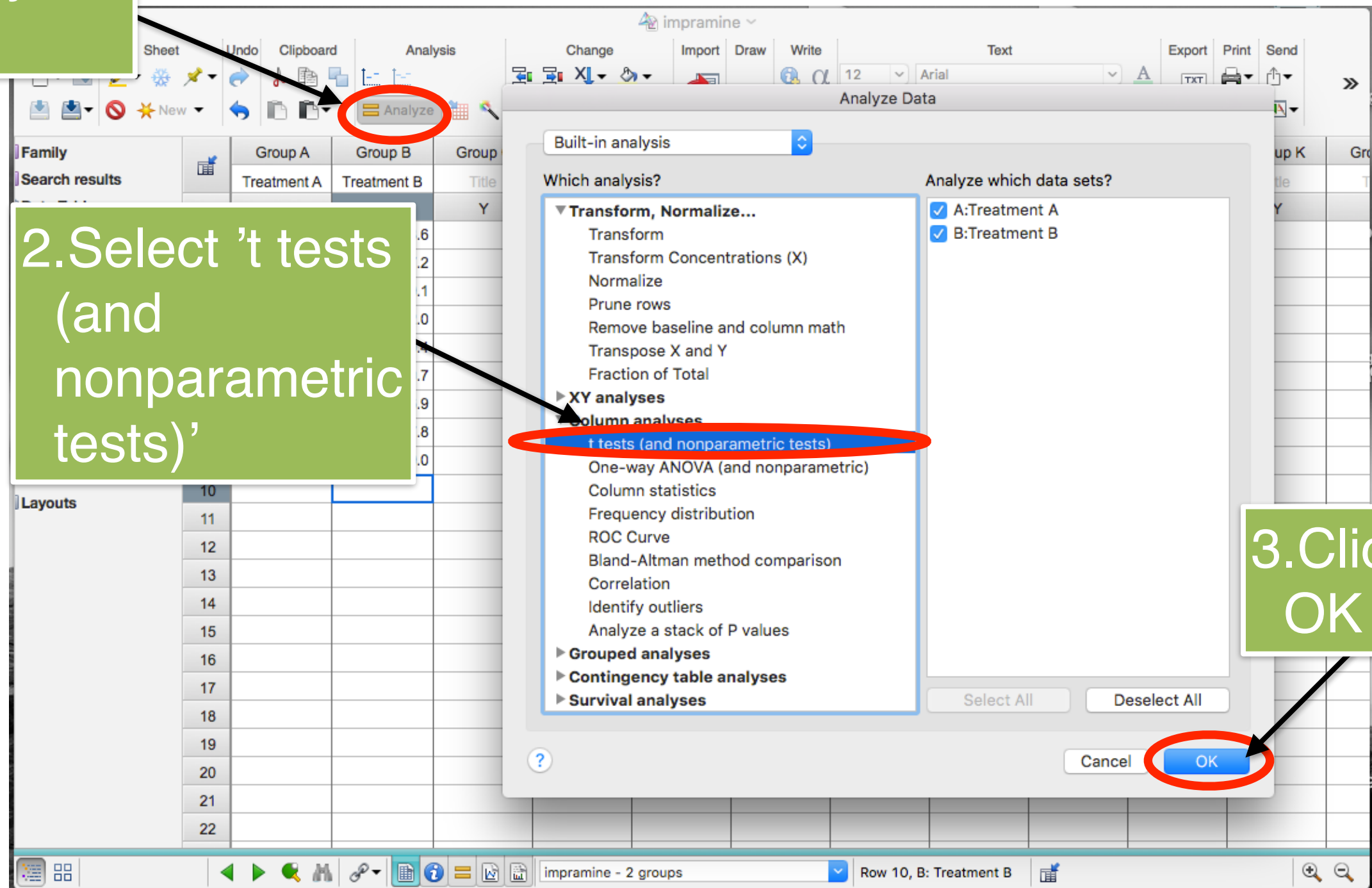
To come in later chapters...

# 95% Confidence Interval for Difference Between 2 Means - GraphPad

1. Click  
'Analyze  
icon

2. Select 't tests  
(and  
nonparametric  
tests)'

3. Click  
OK





Parameters: t Tests (and Nonparametric Tests)

Experimental Design Options

Experimental design

☒ Unpaired

☐ Paired

	Group A	Group B
	Control	Treated
	Y	Y
1		
2		
3		
4		
5		

Assume Gaussian distribution?

☒ Yes. Use parametric test.

☐ No. Use nonparametric test.

Choose test

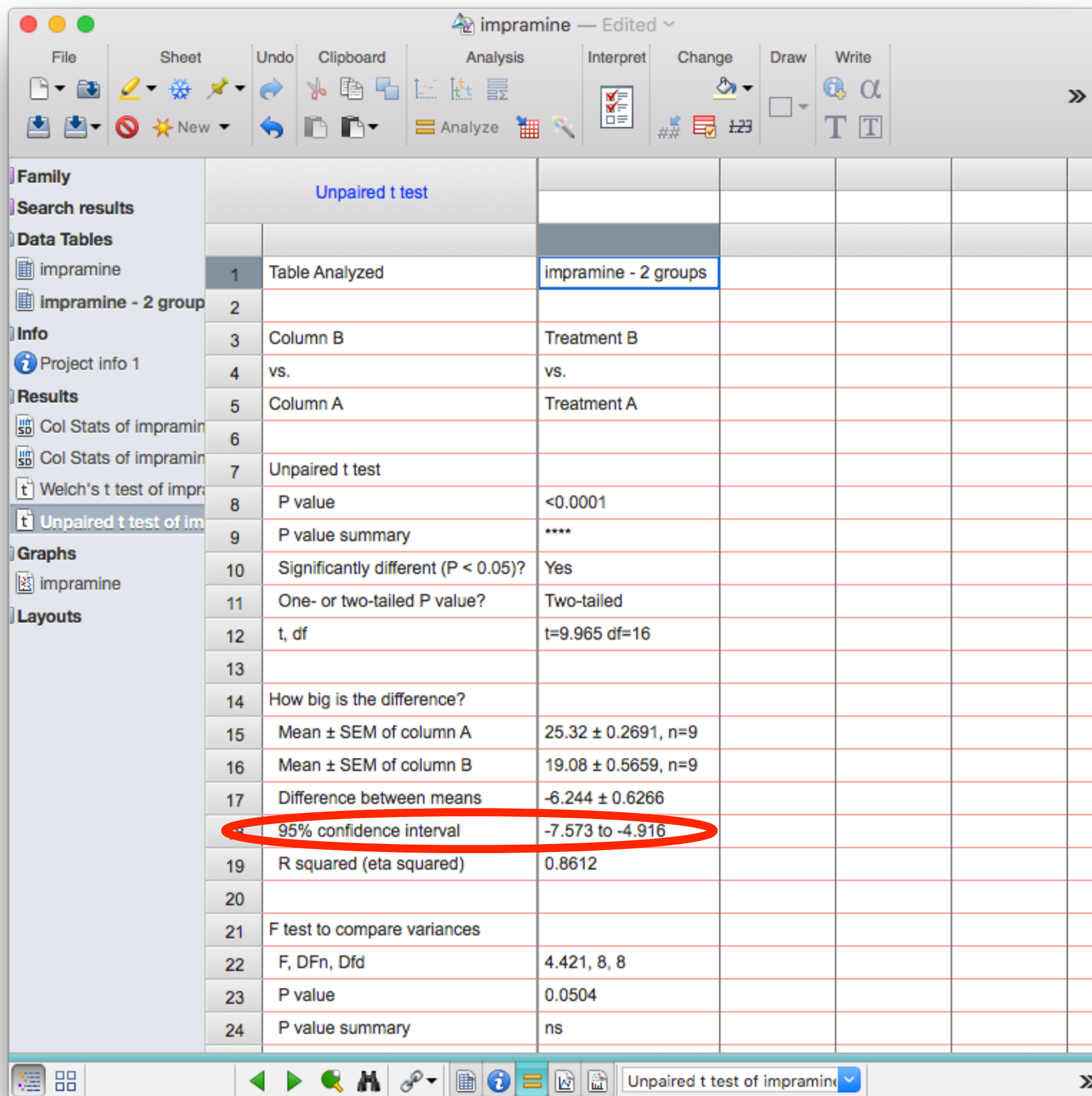
☒ Unpaired t test. Assume both populations have the same SD

☐ Unpaired t test with Welch's correction. Do not assume equal SDs

1.Click OK

Cancel

OK



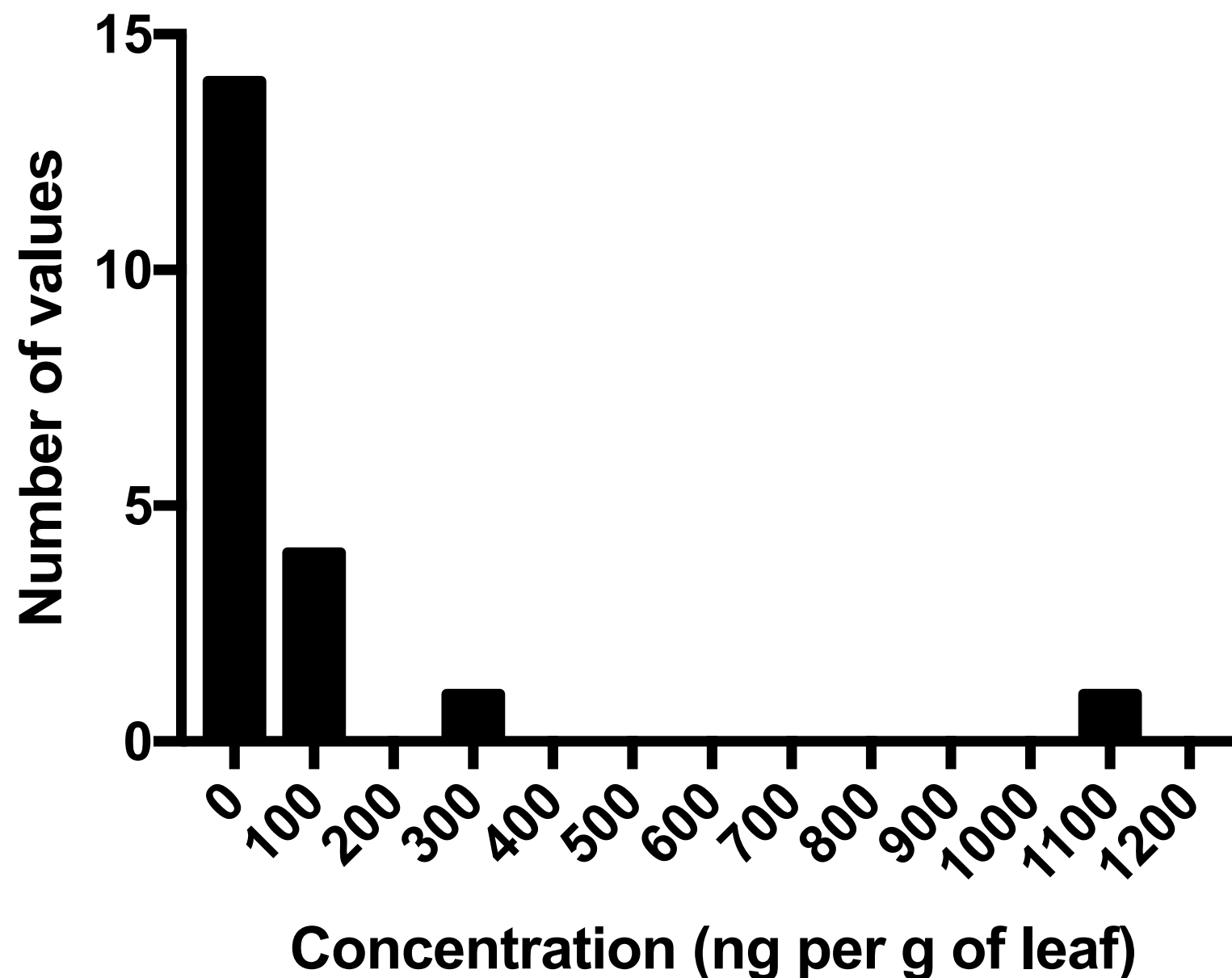
# NORMAL DISTRIBUTION AND CI

# Normal Distribution and CI

- The methods presented for calculating CI are based on the assumption that the data have a normal distribution.
- However, the CI is pretty robust and only performs poorly when the data are grossly non-normal.

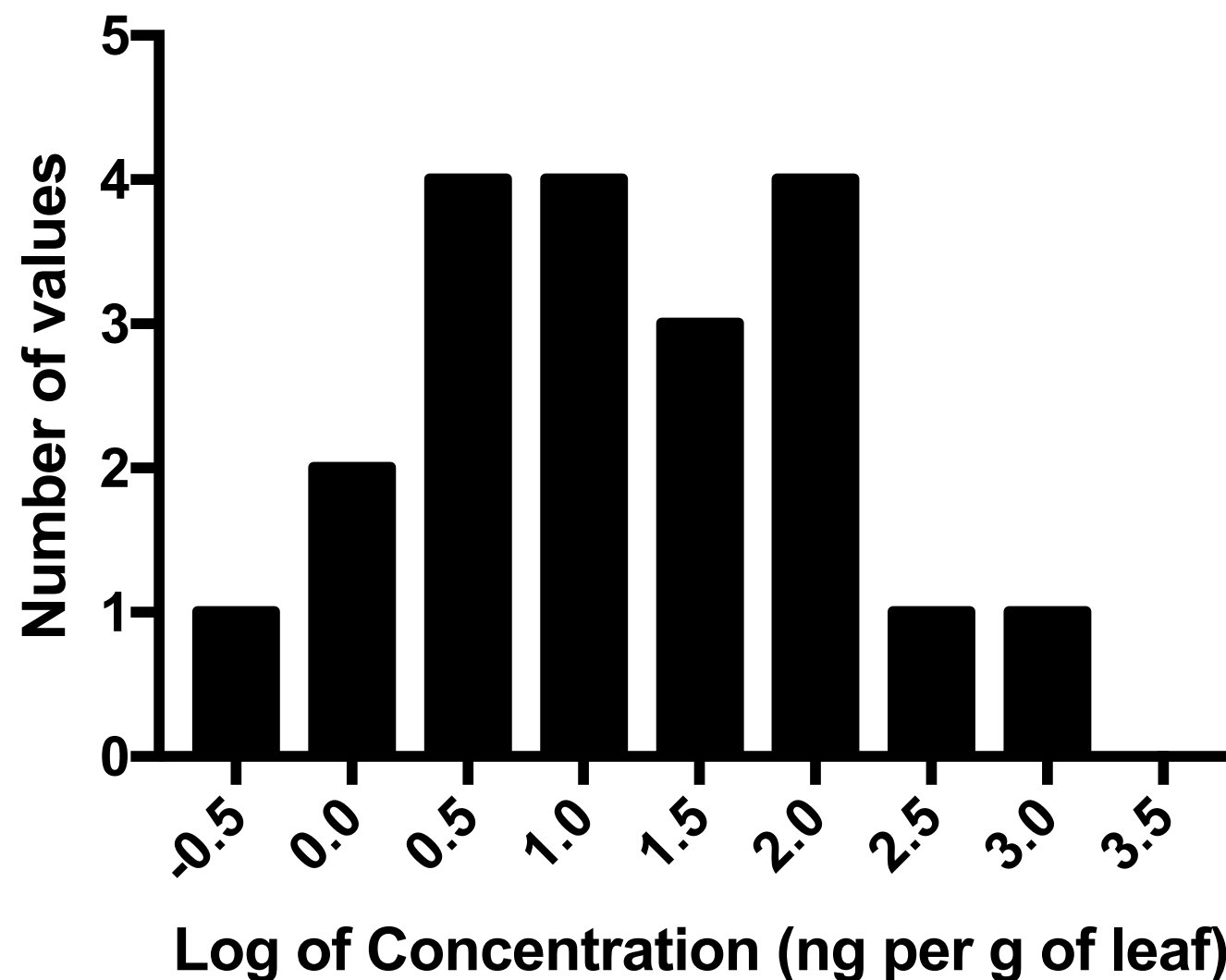
# Example of Non-Normal Data

**Histogram of pesticide residue**



# Common Solution to Positive Skewness - Log Transform

**Histogram of Transform of pesticide residue**

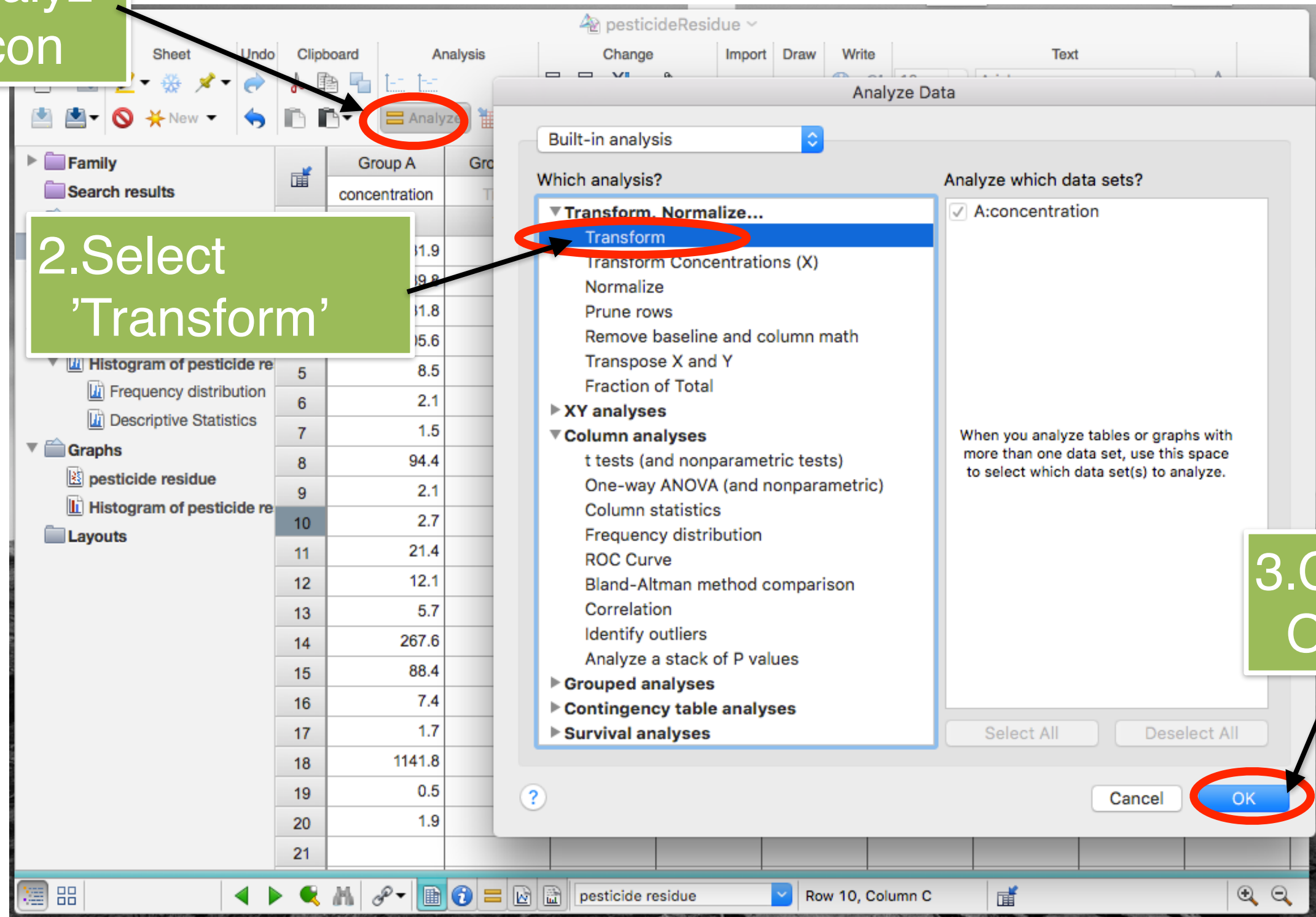


# Log Transform Data in GraphPad

1. Click  
'Analyze'  
icon

2. Select  
'Transform'

3. Click  
OK



1. Select  
'Y=Log(Y)'  
transformation

Parameters: Transform

Function List: Standard functions

☐ Interchange X and Y (then transform as specified below).

☐ Transform X values using  $X=K*X$  K =

☒ Transform Y values using  $Y=\text{Log}(Y)$

☒ Same K for all datasets. K =

☐ Different K for each dataset

Dataset: concentration K =

When it is impossible to transform a SD or SEM

☒ Erase SD or SEM

☐ Convert to an asymmetric 95% confidence interval.

**Replicates**

☒ Transform individual Y values

☐ Transform the average of replicates

**New graph**

☒ Create a new graph of the results

2. Click  
OK



# Geometric Mean

**To calculate a geometric mean:**

1. Log transform data
2. Calculate mean of log transformed data
3. Take the antilog of the mean of the log transformed data

$$\text{geometric mean} = 10^{\left[\frac{1}{n} \sum_{i=1}^n \log a_i\right]}$$

# Calculate 95% Based on Transformed Data

1. Log transform the data
2. Calculate the 95% CI based on the log transformed data
3. Take the antilog of the lower and upper bounds of the 95% CI calculated on the log transformed data

**Log Transform Gives Asymmetrical Limits**

# GraphPad will calculate CI for the geometric mean automatically

Parameters: Column Statistics

**Descriptive Statistics**

- ☒ Minimum and maximum
- ☒ Quartiles (Median, 25th and 75th percentile)
- ☐ Percentile
- ☒ Mean, SD, SEM
- ☐ Coefficient of variation
- ☒ Geometric mean
- ☐ Skewness and kurtosis
- ☒ Column sum

**Confidence intervals**

- ☒ CI of the mean
- ☒ CI of geometric mean
- ☐ CI of median

Confidence level:

**Test if the values come from a Gaussian distribution**

- ☐ D'Agostino-Pearson omnibus normality test (recommended)
- ☐ Shapiro-Wilk normality test
- ☐ Kolmogorov-Smirnov test with Dallal-Wilkinson-Lilliefors P value (not recommended)

**Inferences**

- ☐ One-sample t test. Are column means significantly different than a hypothetical value?
- ☐ Wilcoxon signed-rank test. Compare column medians to a hypothetical value.

Hypothetical value (often 0.0, 1.0 or 100)

When a value equals the hypothetical value:

**Calculations**

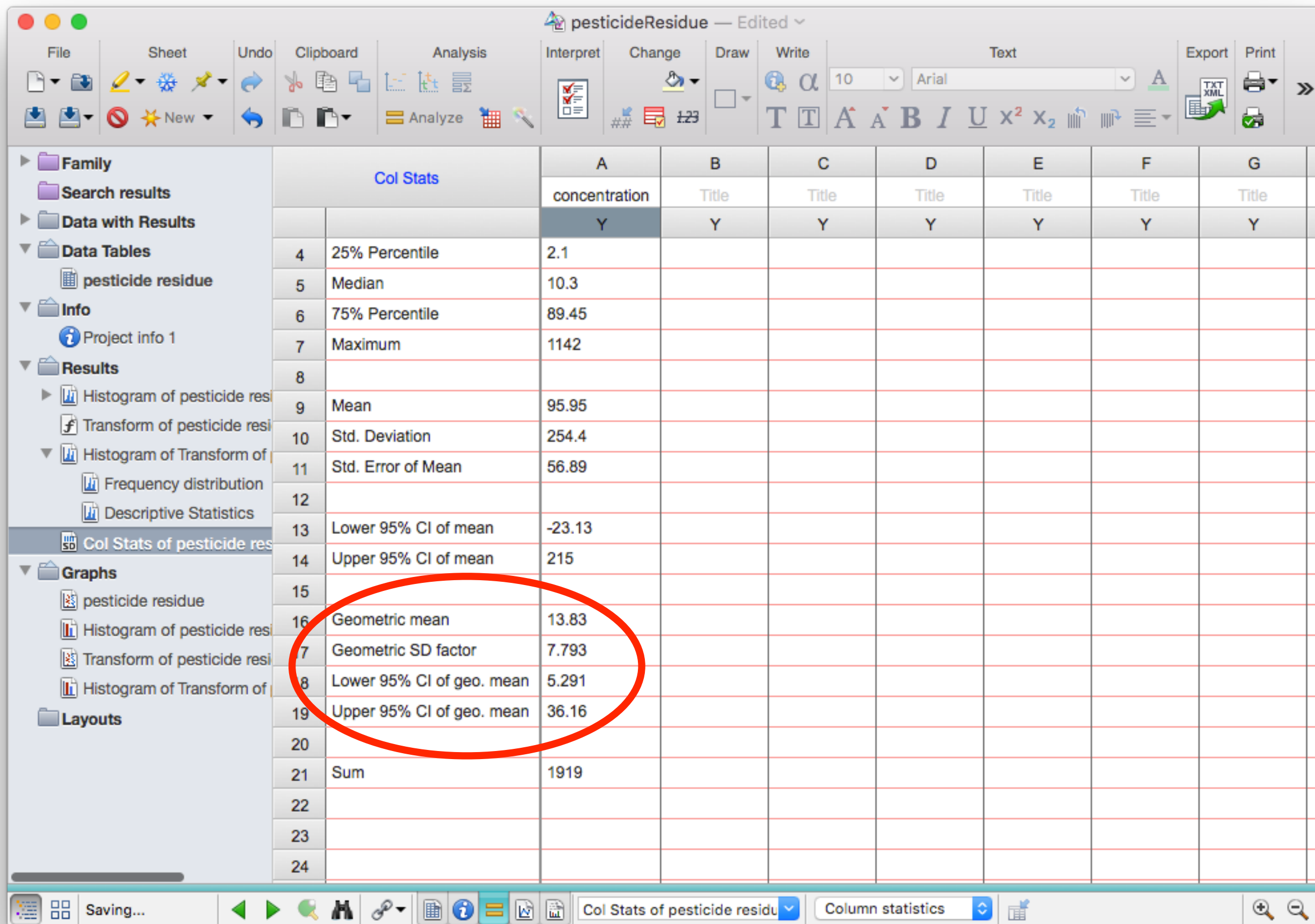
Subcolumns:

**Output**

P-value style:

Show  significant digits.

☐ Make these choices be the default for future analyses.



# Other transformations

- **Log transformation with an added constant** - if the data contains zeros or negative values, it is impossible to log transform the data. The usual solution is to add a fixed value to each point prior to log transforming the data.
- **Square-root transformation** - Data that represent counts tend to be positively skewed. Transform the data by taking the square root.

# What did we learn?

- Confidence intervals give a sense of how accurately we estimated the population mean from our sample mean.
- Level of confidence refers to the number of intervals generated in the same manner that would cover the true mean.
- Larger SD results in larger CI, smaller sample size results in larger CI, higher level of confidence results in larger CI
- 95% CI for the difference between two means takes into account the error in estimating both means.
- Deviations from normality can cause CI to be incorrect
- Data transformation can alleviate the problems of non-normality