

Chapter 4: The Normal Distribution

TXCL7565/PHSC7565

What This Lecture Covers

- ▶ What is a normal distribution
- ▶ Identifying data that are not normally distributed
- ▶ Proportions of individuals within 1 SD or 2 SD of the mean
- ▶ Skewness and kurtosis
- ▶ Tests for normal distributions

WHAT IS A NORMAL DISTRIBUTION

Normal Distribution

- Many of the things we measure show a characteristic distribution, with the bulk of the data points clustered around the mean and data points become steadily rarer as we move further from the mean.
- When many (independent) random factors contribute to an observed value, the observed values tend to follow a normal distribution.
- Normal Distribution = Gaussian Distribution

Many Random Factors

```
true_value = 10
e1 = runif(100000) - 0.5
e2 = runif(100000) - 0.5
e3 = runif(100000) - 0.5
e4 = runif(100000) - 0.5
e5 = runif(100000) - 0.5

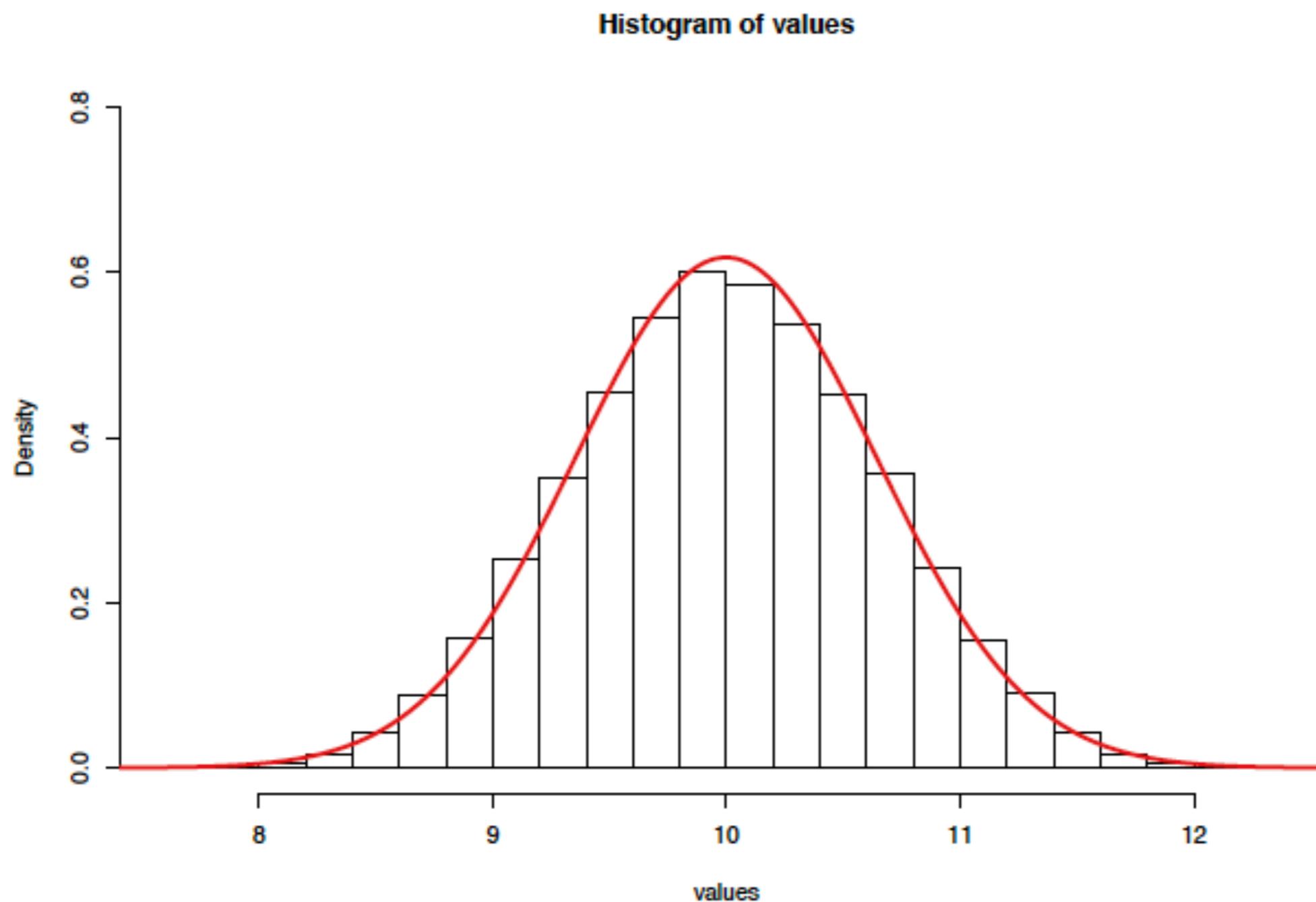
values = true_value + e1 + e2 + e3 + e4 + e5
head(round(values,2))
```

```
## [1] 10.31 9.25 9.61 9.20 9.74 10.60
```

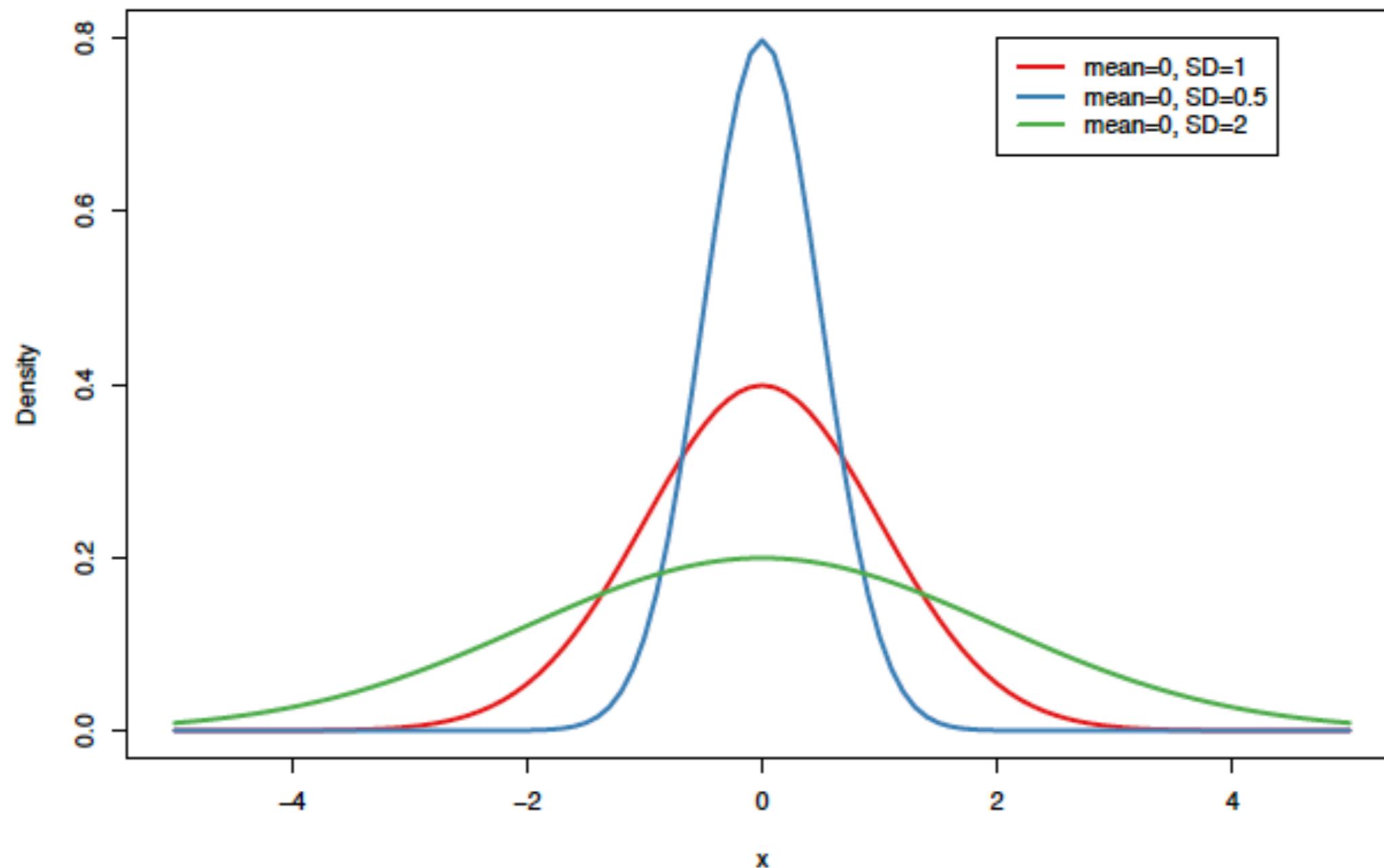
```
tail(round(values,2))
```

```
## [1] 8.57 9.06 10.63 9.87 10.81 9.57
```

Results of many random factors contributing



The Nature of the Normal Distribution



Why Does a Normal Distribution Matter?

Many commonly used statistical tests rely on the assumption that the data have been sampled from a population that follows a normal distribution.

This is often a reasonable assumption.

**IDENTIFYING DATA THAT
ARE NOT NORMALLY
DISTRIBUTED**

How to Spot Non-Normal Data

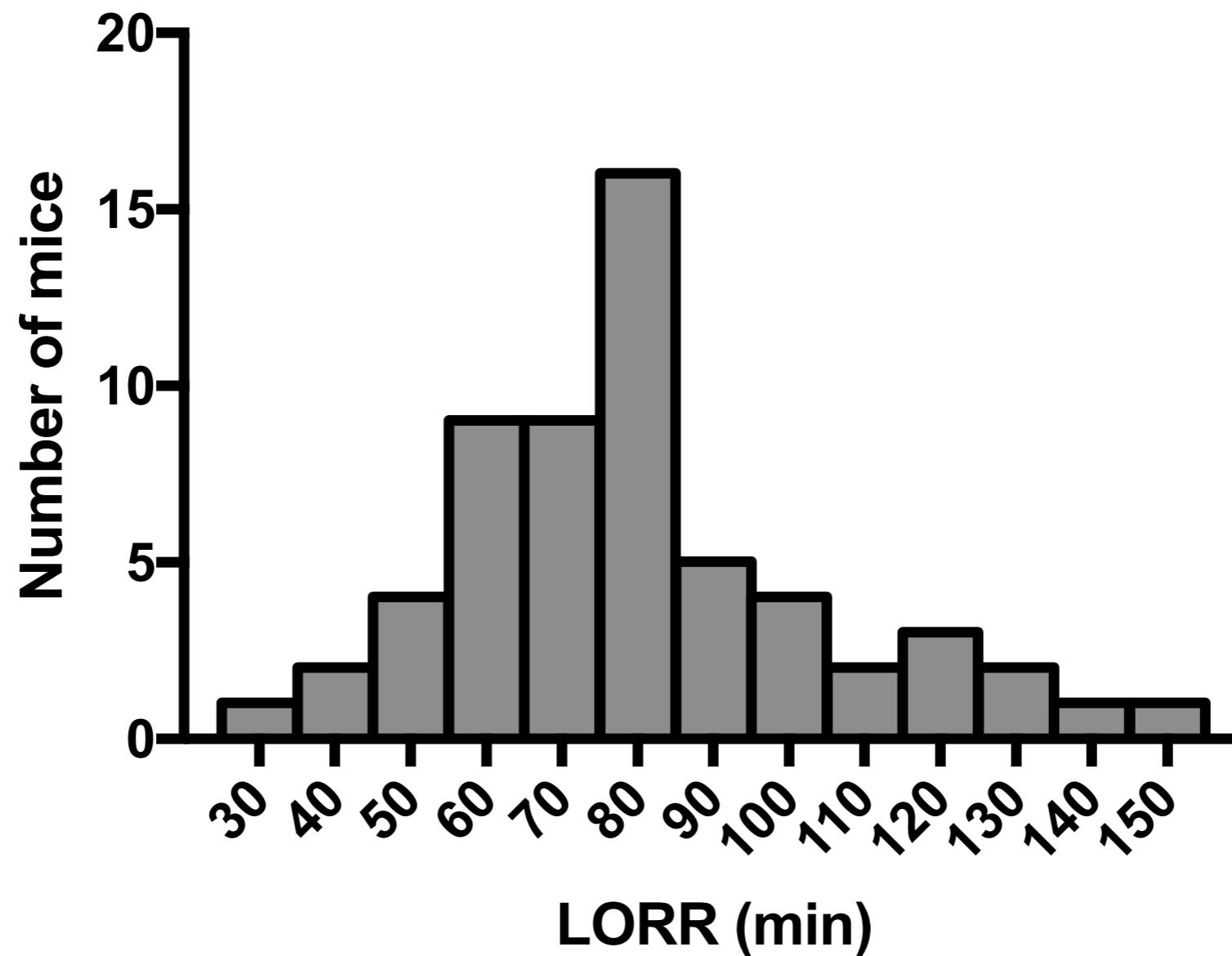
Most often, the easiest/best way to determine if data are normally distributed is to plot them in a histogram

Example Data Set - Sleep Time

One characteristic related to an individual's likelihood of developing an alcohol use disorder is how sensitive they are to alcohol. In the Radcliffe lab, how sensitive a particular mouse is to the sedative effects of alcohol is measured by giving them a large enough dose of alcohol to cause them to 'fall asleep' and then measuring the number of minutes that pass before they wake up.

Sleep time or Loss of Righting Reflex (LORR) is the number of minutes between when the mouse first loses the ability to right themselves when placed on their back to when they can right themselves again.

Histogram of Sleep Time



Characteristics of a Normal Distribution

3 visual characteristics that any true normal distribution will possess:

1. The data are unimodal
2. The distribution is symmetrical
3. The frequencies decline steadily as we move towards higher or lower values, without any sudden, sharp cut-off

Unimodal vs. Polymodal

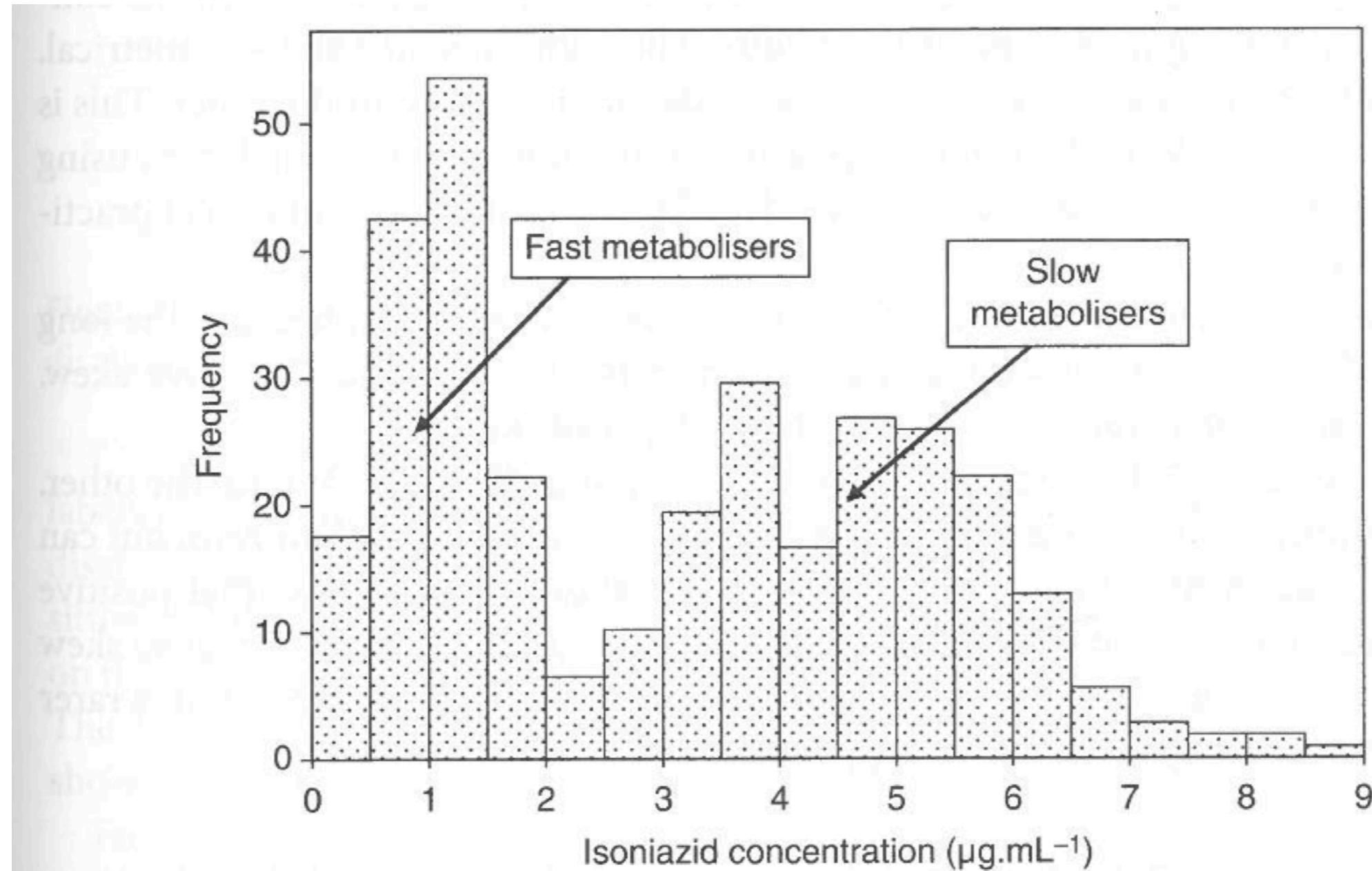


Figure 4.1 Bimodal data. Isoniazid concentrations ($\mu\text{g.mL}^{-1}$) six hours after a standard oral dose

Symmetrical Distribution

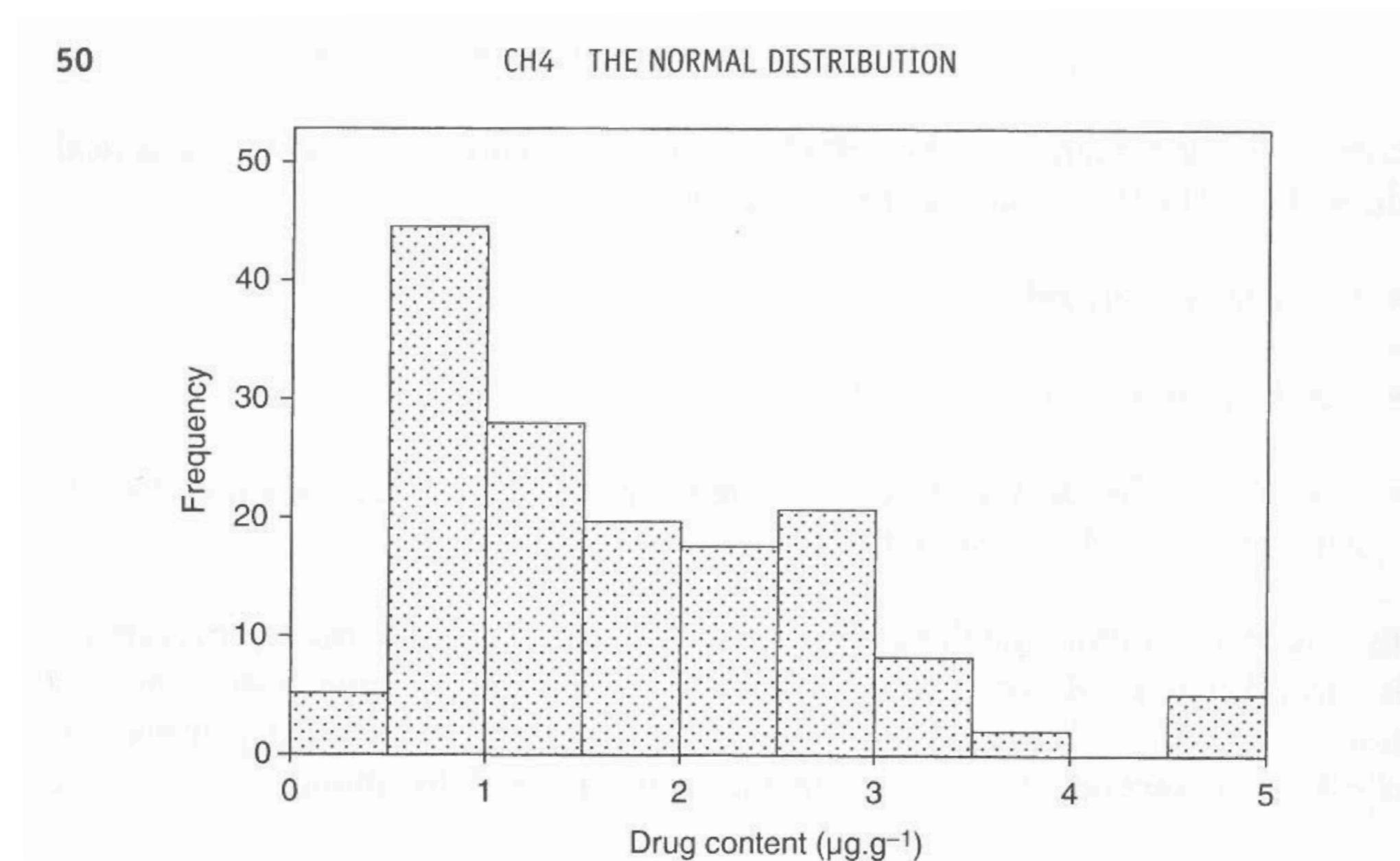


Figure 4.2 Skewed data. Drug content (mg.g^{-1} of dried plant tissue) of a series of individual plants

Sharp Cut-Offs

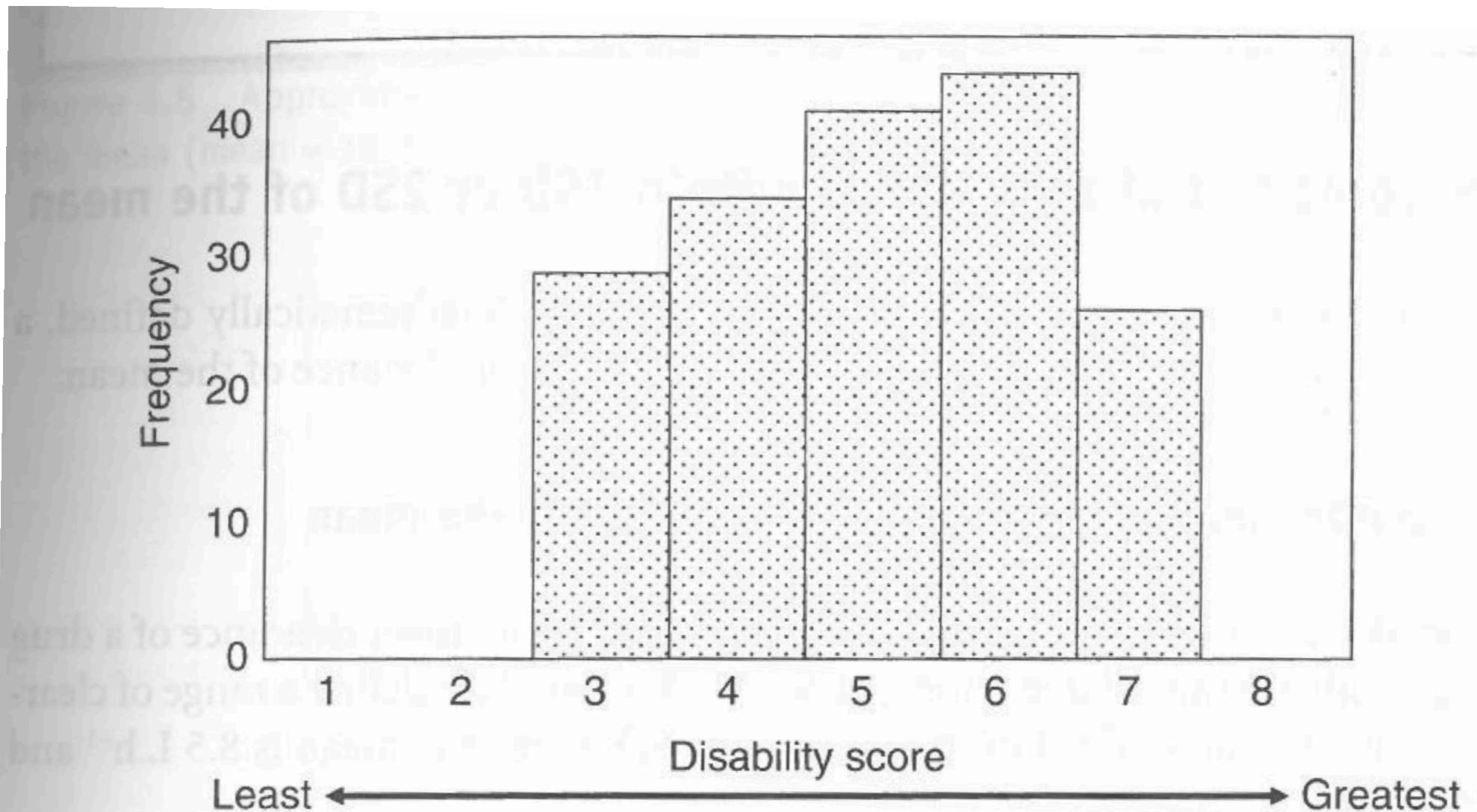
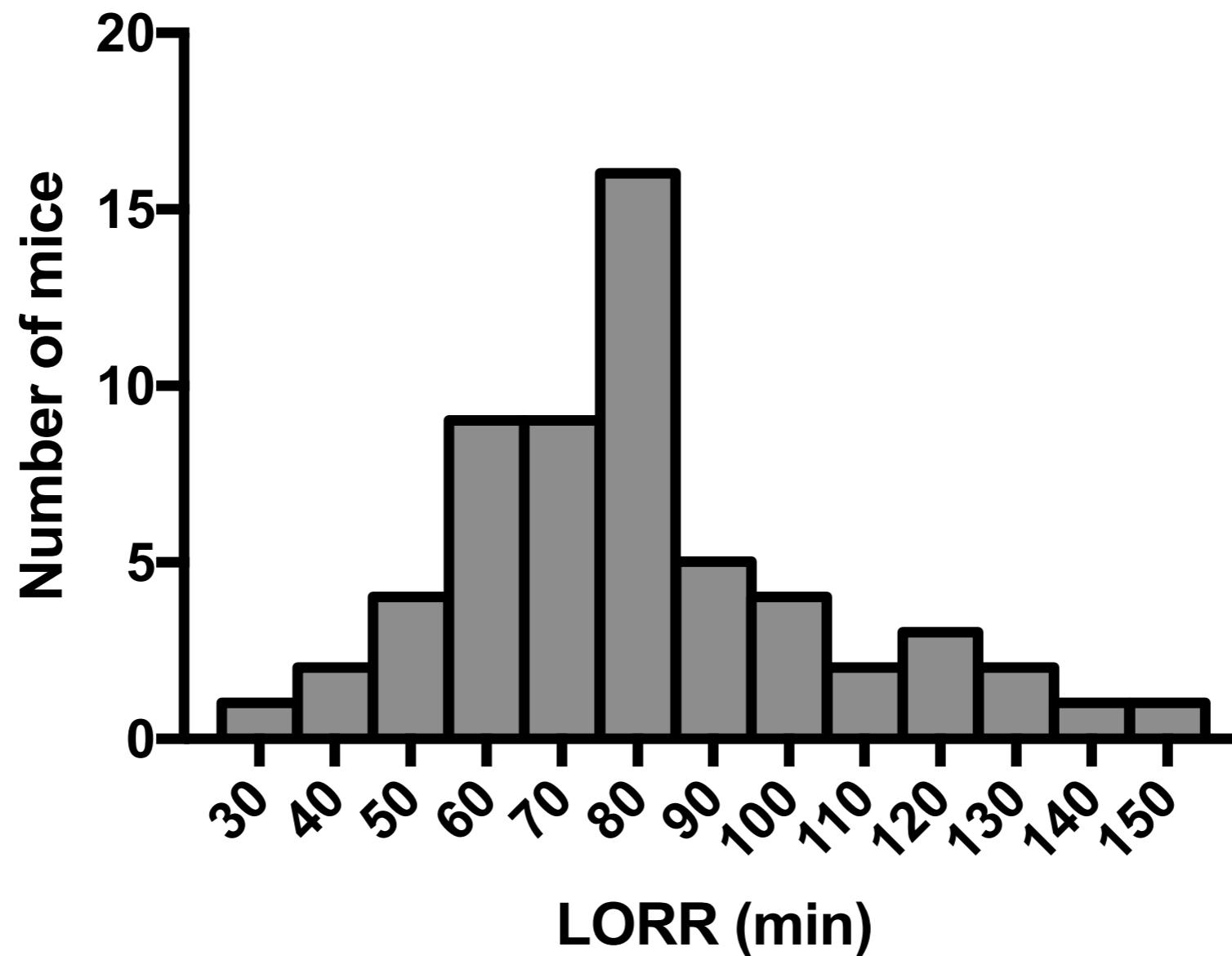
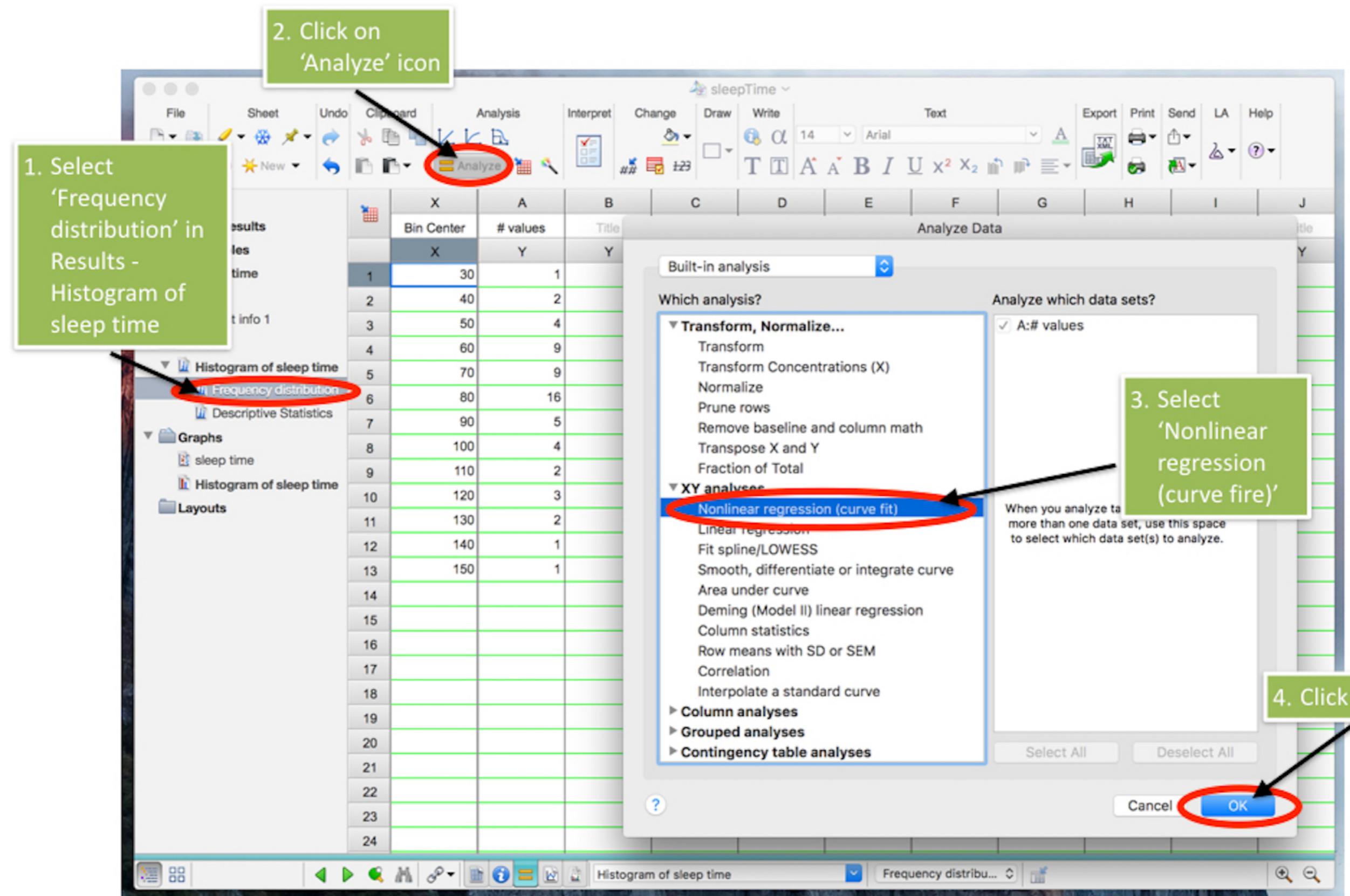


Figure 4.4 Data with sudden cut offs. Patients' self assessment scores for degree of disability due to gout

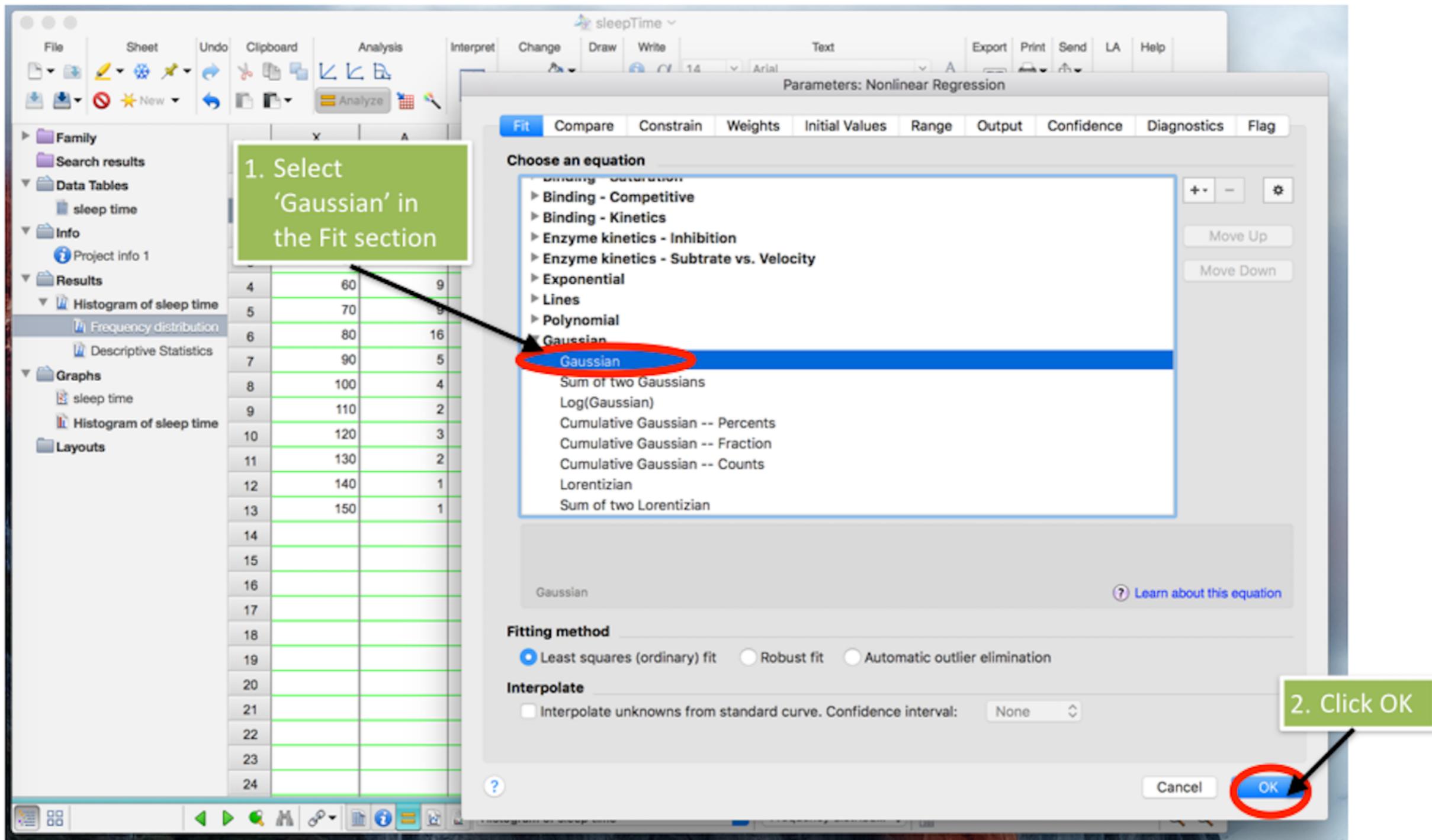
Histogram of Sleep Time



Adding a Normal Distribution Line - GraphPad

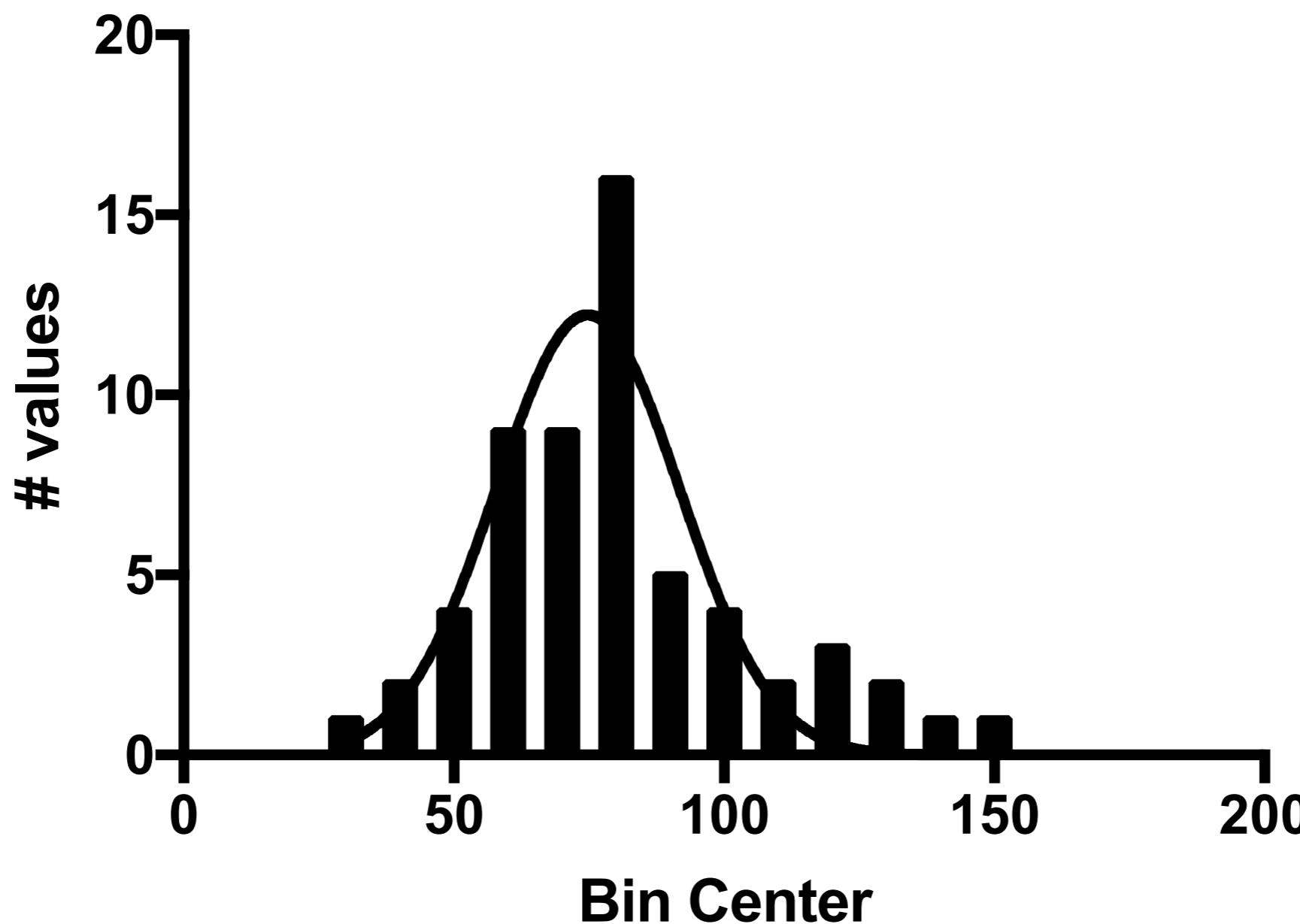


Adding a Normal Distribution Line - GraphPad



Histogram of Sleep Time With Normal Curve

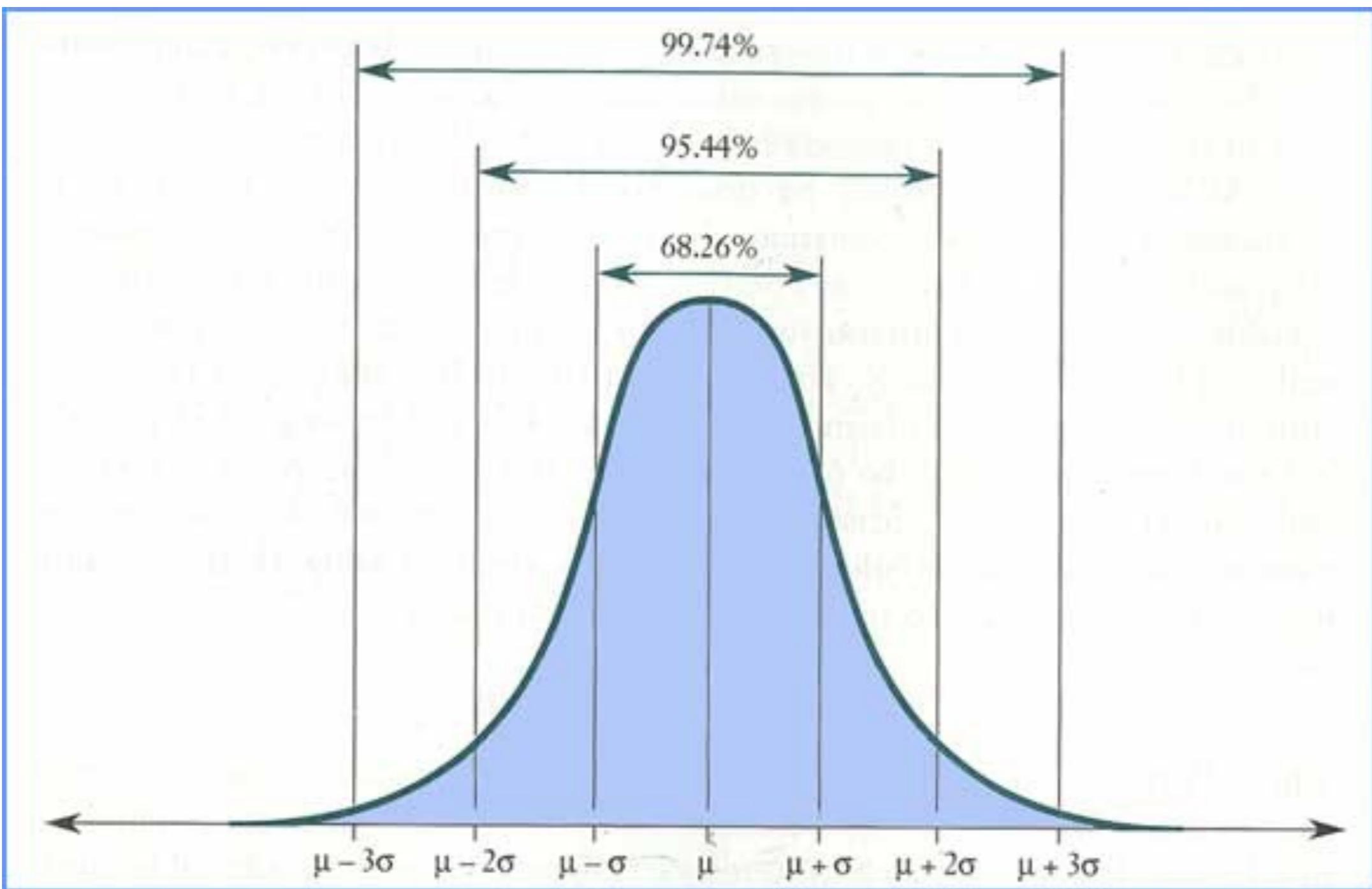
Histogram of sleep time



**PROPORTIONS OF
INDIVIDUALS WITHIN 1 SD
OR 2 SD OF THE MEAN**

Proportions in Normal Distribution

Ideal Normal Distribution In Terms of SD



Sleep Time Example

- Mean = 79.2
- SD = 25.05
- Mean ‘ \pm ’ 1 SD = (54.2, 104.3)
- Proportion of values within 1 SD = 0.73
- Mean ‘ \pm ’ 2 SD = (29.2, 129.3)
- Proportion of values within 2 SD = 0.95

The Standard Normal

- Standard Normal - Normal distribution with mean = 0 and SD = 1
- All normal distributions can be converted to a standard normal distribution using the following formula:

$$Z = \frac{\text{Value} - \text{Mean}}{SD}$$

- Z is the number of SD the value is away from the mean

Calculate the proportion of samples within 1 SD and 2 SD - GraphPad

1. Calculate Z scores for each value using 'Transform'
2. Calculate the absolute value of the z scores using 'Transform'
3. Calculate the cumulative frequency based on the absolute values of the z scores

Calculate Z Scores

The screenshot shows the QI Macros software interface. A data table titled 'sleep time' is open, showing a single column of values from 30.0 to 78.0. The 'Analysis' tab is selected in the ribbon. A yellow box labeled '1.' has an arrow pointing to the 'sleep time' data table. A second yellow box labeled '2.' has an arrow pointing to the 'Analyze' icon in the ribbon. A third yellow box labeled '3.' has an arrow pointing to the 'Transform' option in the 'Which analysis?' dropdown. A fourth yellow box labeled '4.' has an arrow pointing to the 'OK' button in the dialog box.

1. Select the data table you would like to evaluate

2. Click on the 'Analyze' icon

3. Select 'Transform' under 'Which analysis?'

4. Click 'OK'

	Group A	Group B
	LORR (min)	Title
Y	Y	
1	30.0	
2	36.0	
3	37.0	
4	46.0	
5	50.0	
6	51.0	
7	54.0	
8	55.0	
9	56.0	
10	56.0	
11	57.0	
12	58.0	
13	58.0	
14	60.0	
15	60.0	
16	62.0	
17	65.0	
18	65.0	
19	67.0	
20	69.0	
21	70.0	
22	71.0	
23	72.0	
24	73.0	
25	74.5	
26	75.0	
27	75.5	
28	77.0	
29	78.0	

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

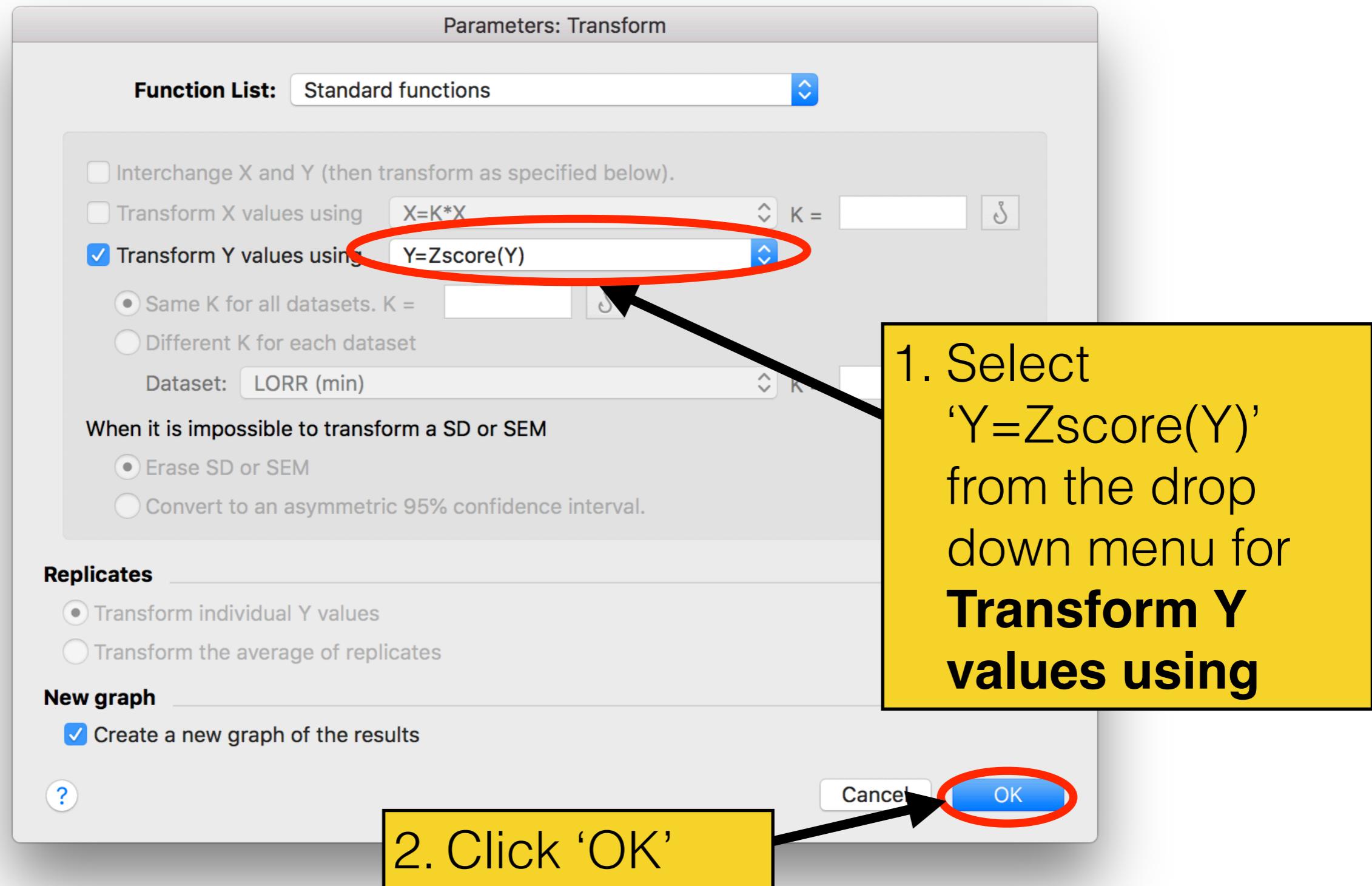
Analyze which data sets?

A:LORR (min)

Select All Deselect All

Cancel OK

Calculate Z scores



Calculate the absolute value of each Z score

The screenshot shows the SigmaPlot software interface with a yellow callout box containing four numbered steps:

1. Select the Transformed data table under the **Results** folder
2. Click on the 'Analyze' Icon
3. Select 'Transform' under 'Which analysis?'
4. Click 'OK'

The software interface includes a toolbar at the top, a left sidebar with project navigation, and a main workspace showing a data table titled "Transform". The data table has columns A and B, with row 1 labeled "LORR (min)" and row 2 labeled "Y" containing the value "-1.966". The "Analyze" icon in the toolbar is circled in red, and the "Transform..." option in the Results folder is also circled in red. The "Transform" option in the "Transform, Normalize..." menu is highlighted and circled in red. The "OK" button in the bottom right corner of the dialog box is also circled in red.

1. Select the Transformed data table under the **Results** folder

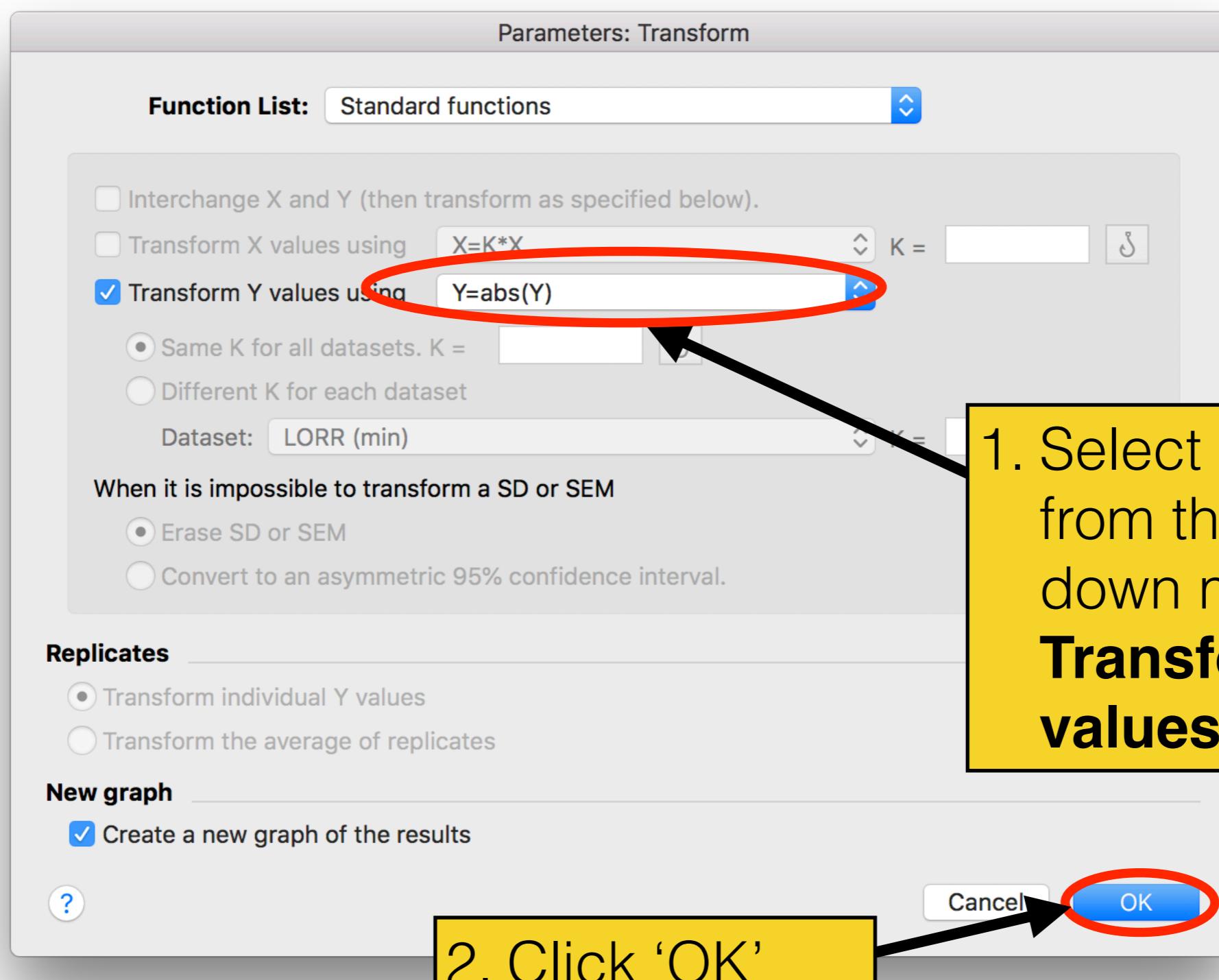
2. Click on the 'Analyze' Icon

3. Select 'Transform' under 'Which analysis?'

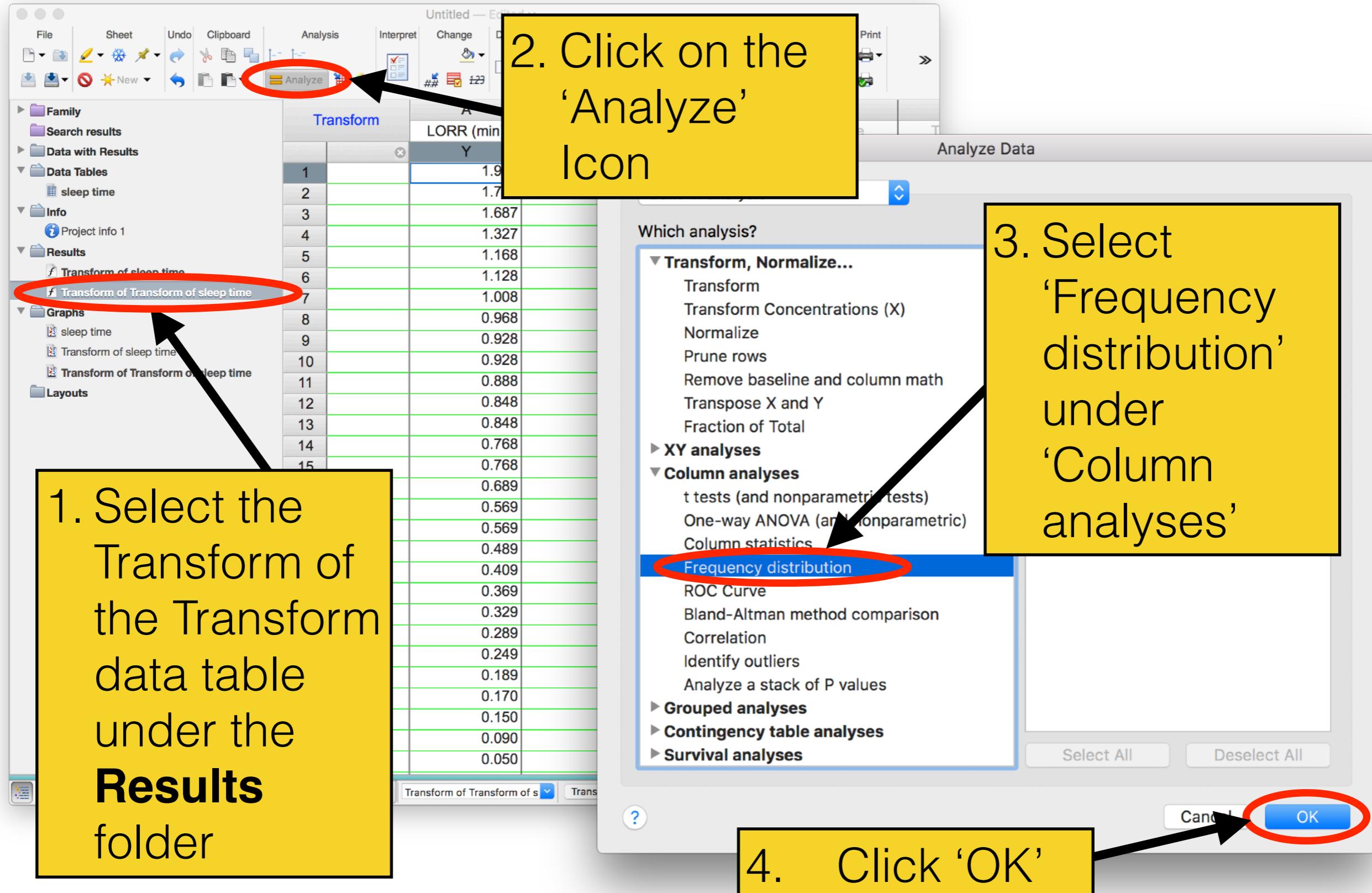
4. Click 'OK'

Untitled — Edit
File Sheet Undo Clipboard Analysis Interpret Change Print
New
Transform A B
A LORR (min) Title
B Y Y
1 -1.966
2 -1.727
3 -1.687
4 -1.327
5 -1.168
6 -1.128
7 -1.008
8 -0.968
9 -0.928
10 -0.928
11 -0.888
12 -0.848
-0.848
-0.768
-0.768
-0.689
-0.569
-0.569
-0.489
-0.409
-0.369
-0.329
-0.289
-0.249
-0.189
-0.170
-0.150
-0.090
-0.050
which analysis?
Analyze which data sets?
 A:LORR (min)
Select All Deselect All
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
OK

Calculate the absolute value of each Z score

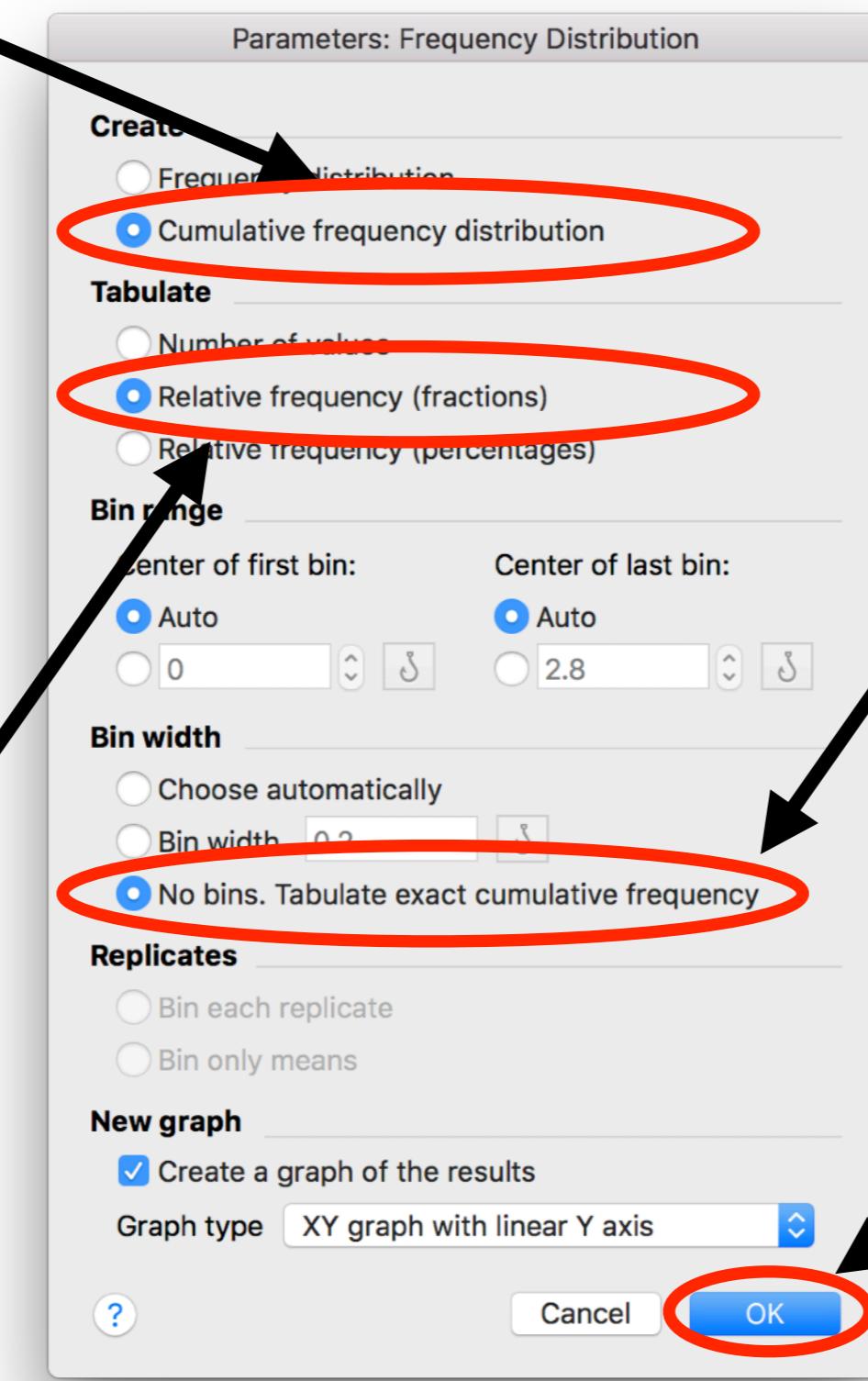


Calculate the cumulative frequency based on the absolute values of the z scores



Calculate the cumulative frequency based on the absolute values of the z scores

1. Select 'Cumulative frequency distribution' under **Create**



2. Select 'Relative frequency (fractions)' under **Tabulate**

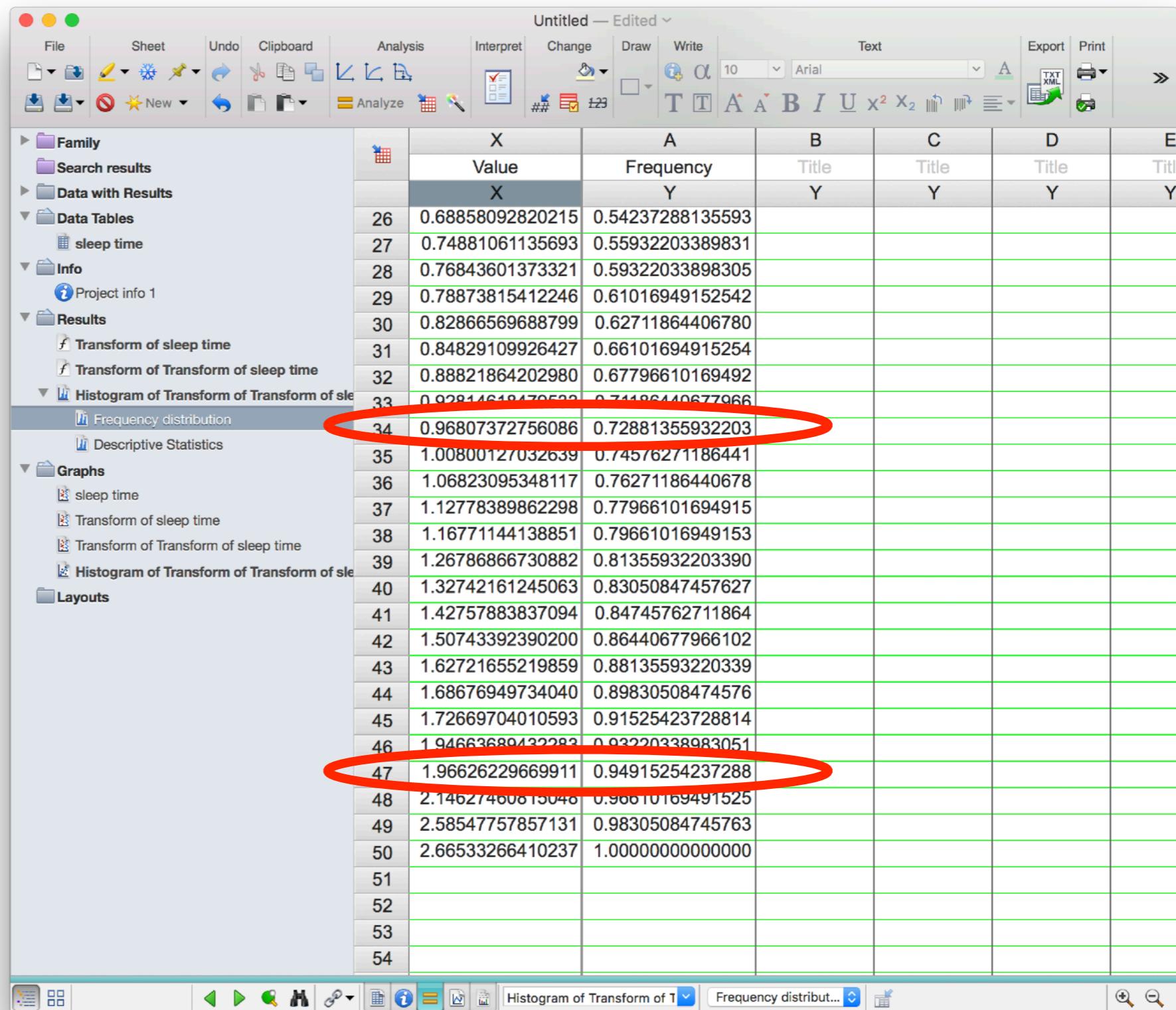
3. Select 'No bins.' under **Bin width**

4. Click 'OK'

Calculate the proportion of samples within 1 SD and 2 SD - GraphPad

73% of the data points are less than 1 SD from the mean

95% of the data points are less than 2 SD from the mean

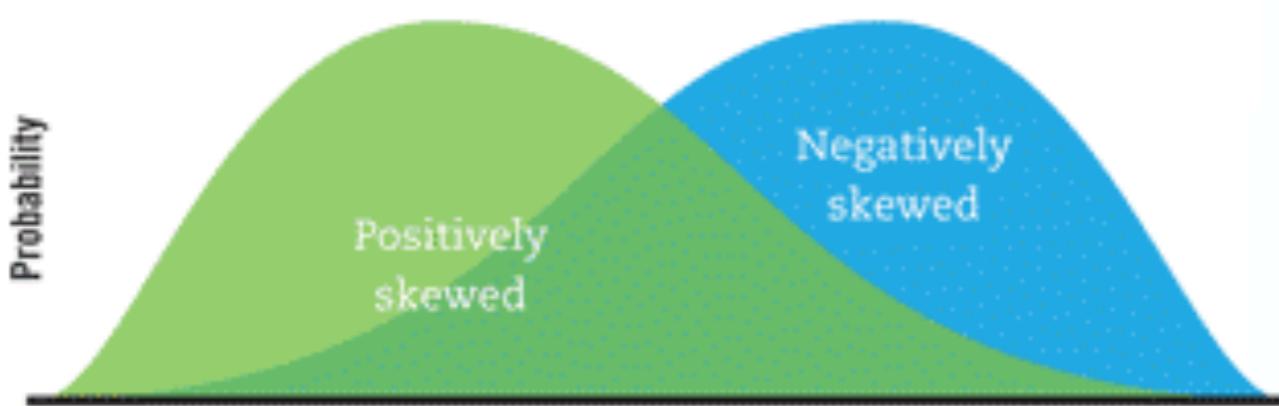


The Normal Distribution Does Not Define Normal Limits

- Normal distribution doe not equate to ‘normal’ range.
- Defining the normal limits of a clinical measurement is not straightforward and requires clinical thinking, not just statistics.

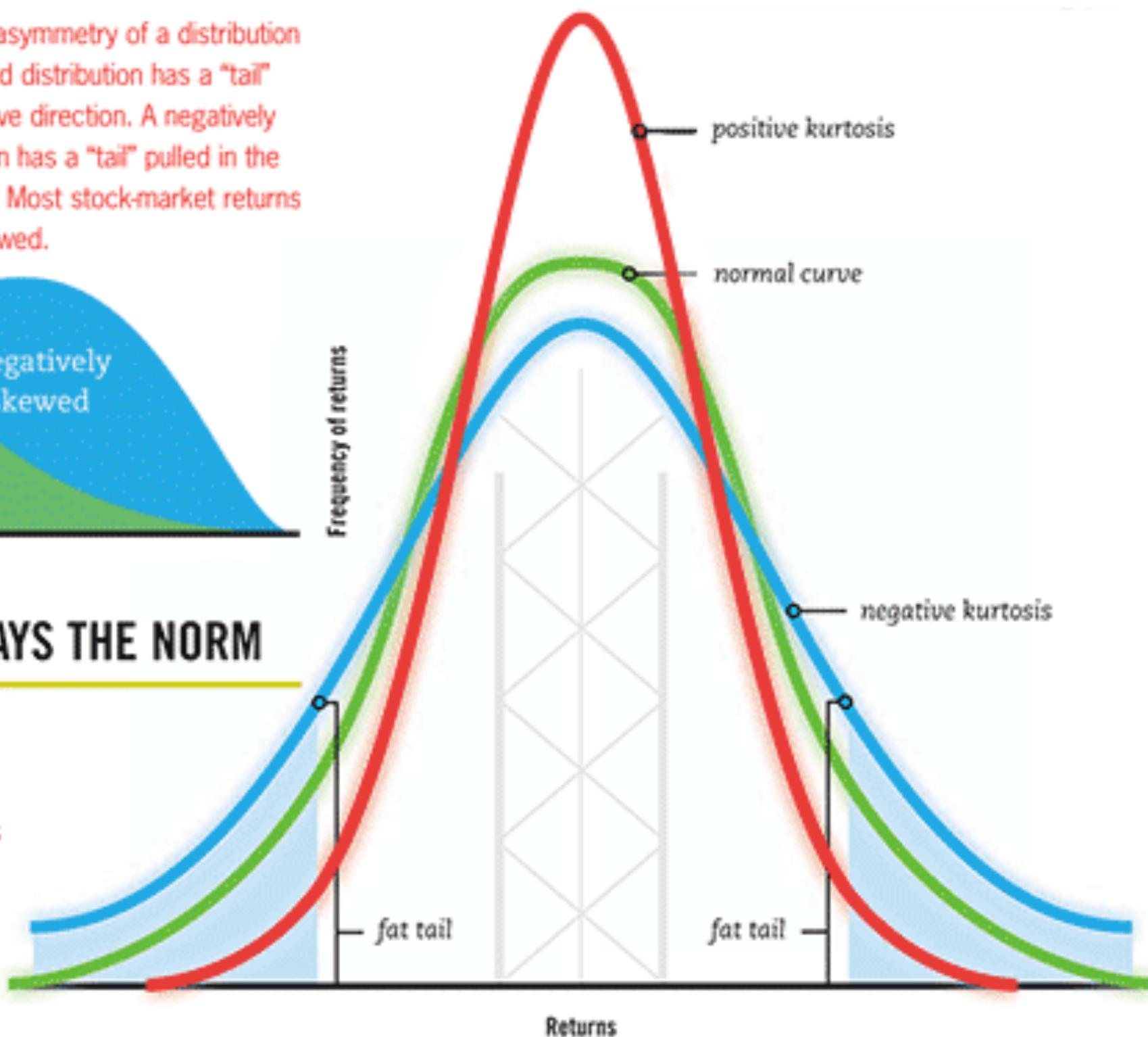
SKEWNESS AND KURTOSIS

Skewness is the asymmetry of a distribution
A positively skewed distribution has a "tail" pulled in the positive direction. A negatively skewed distribution has a "tail" pulled in the negative direction. Most stock-market returns are negatively skewed.



NORMAL NOT ALWAYS THE NORM

Kurtosis refers to how peaked the curve is:
steeper means positive kurtosis and flatter means negative kurtosis. Fat tails occur when there are more outsize returns on the downside or upside, or both, than the normal curve suggests.



Skewness

Skewness can be quantified numerically by a number that ranges from negative infinity to positive infinity, with a value of 0 indicating that no skewness is present

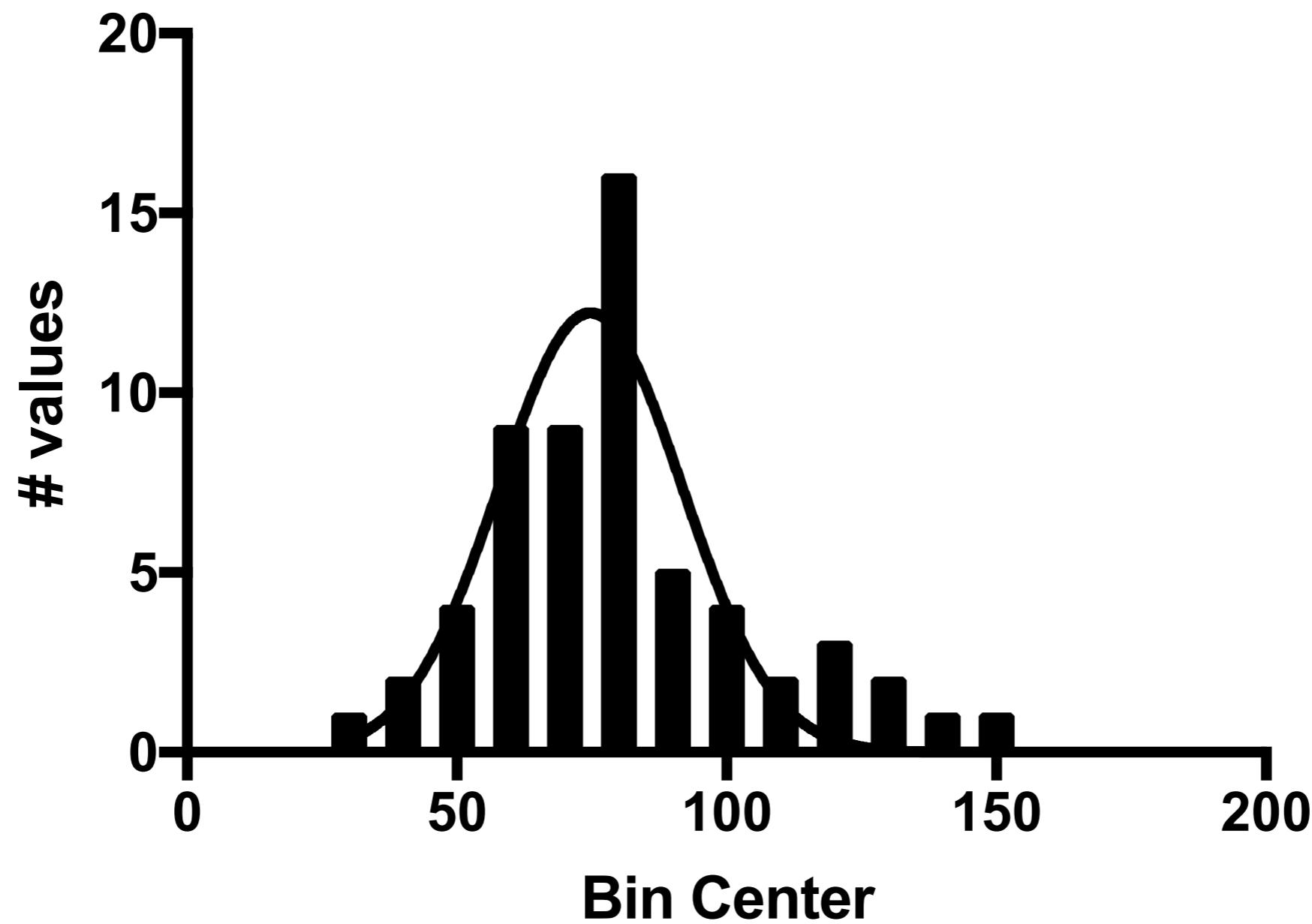
Kurtosis

Like skewness, **kurtosis** can also be quantified numerically with a range of negative infinity to positive infinity, with a value of 0 indicating that no kurtosis is present.

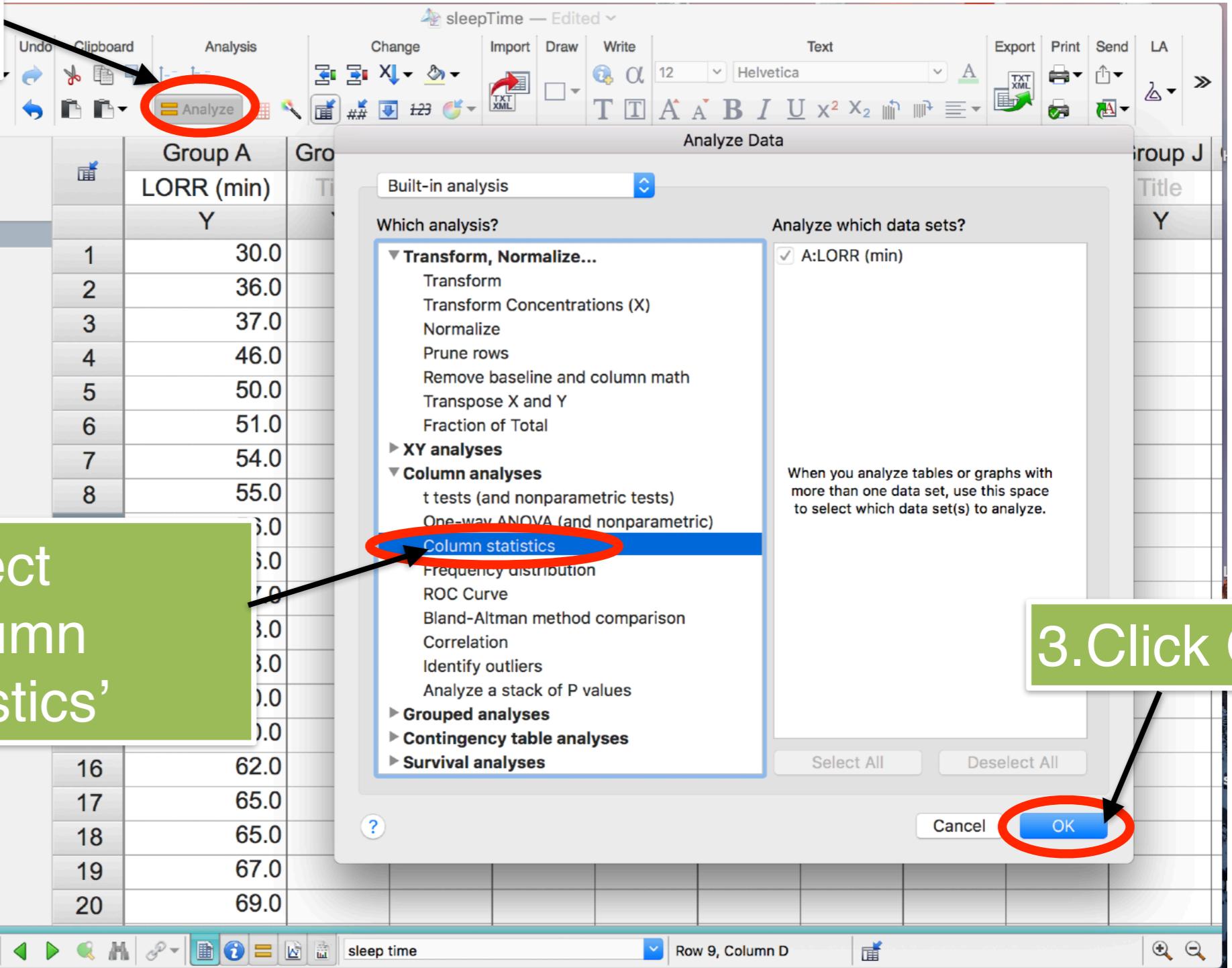
- A *negative* kurtosis value indicates that the ‘shoulders’ of the distribution are too high and wide
- A *positive* kurtosis value indicates that the ‘shoulders’ of the distribution are too low and narrow.

Skewness and Kurtosis in Sleep Time

Skewness = 0.70 (positive skew); kurtosis = 0.61 (narrow peak)



1. Click
'Analyze'
icon



2. Select
'Column
statistics'

3. Click OK

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

Confidence level: 95%

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

1. Click 'Skewness and kurtosis'

2. Click OK

sleepTime — Edited

File | Sheet | Undo | Clipboard | Analysis | Interpret | Change | Draw | Write | Text | Export | Print | Send | LA | Help | GraphPad PRISM®

Col Stats

	A	B	C	D	E	F
	LORR (min)	Title	Title	Title	Title	Title
	Y	Y	Y	Y	Y	Y
1	Number of values	59				
2						
3	Minimum	30				
4	25% Percentile	60				
5	Median	78				
6	75% Percentile	89				
7	Maximum	146				
8						
9	Mean	79.25				
10	Std. Deviation	25.05				
11	Std. Error of Mean	3.261				
12						
13	Lower 95% CI of mean	72.72				
14	Upper 95% CI of mean	85.77				
15						
16	Geometric mean	75.41				
17	Geometric SD factor	1.381				
18						
19	Skewness	0.6964				
20	Kurtosis	0.6086				
21						
22	Sum	4676				
23						
24						



Col Stats of sleep time

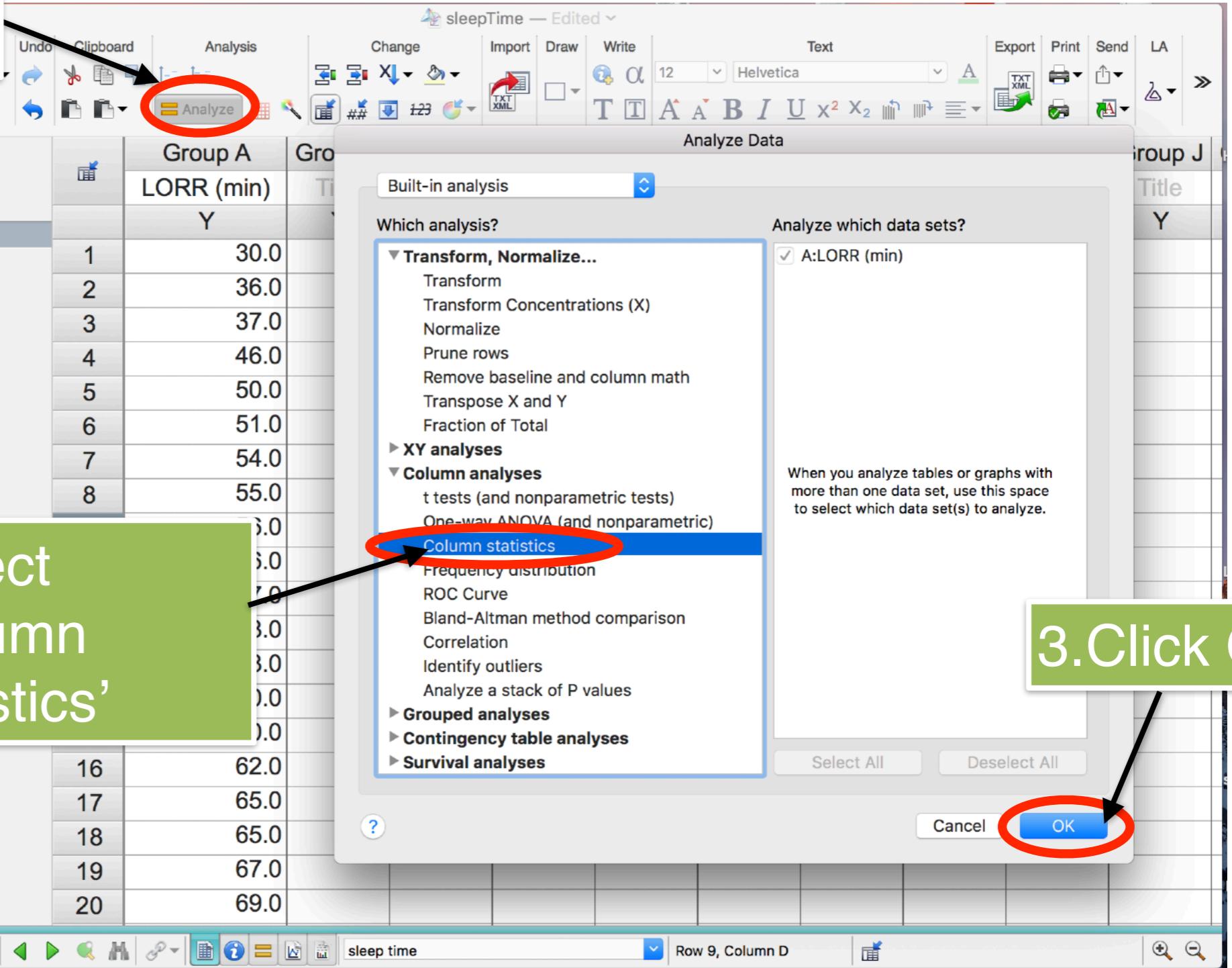
Column statistics

TESTS FOR NORMAL DISTRIBUTIONS

D'Agostino-Pearson Omnibus Normality Test

- Recommended in GraphPad
- Test for both skewness and kurtosis
- A significant p-value ($p<0.05$) indicates that the distribution is significantly different from normal, i.e., the data are non-normal.
- CAVEATS:
 - With a small sample size even large departures from normality may not be detected
 - With a large sample size even small departures from normality will be detected

1.Click
'Analyze
icon'



Parameters: Column Statistics

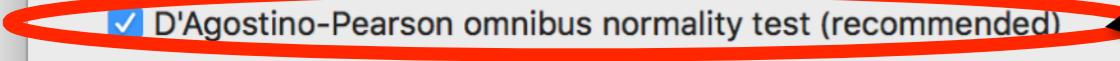
Descriptive Statistics

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
Confidence level: 95%

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

1. Click
'D'Agostino-
Pearson
omnibus
normality test'

2. Click OK

sleepTime — Edited

File | Sheet | Undo | Clipboard | Analysis | Interpret | Change | Draw | Write | Text | Export | Print | Send | LA | Help | GraphPad PRISM®

Col Stats

	A	B	C	D	E
	LORR (min)	Title	Title	Title	Title
	Y	Y	Y	Y	Y
1	Number of values	59			
2					
3	Minimum	30			
4	25% Percentile	60			
5	Median	78			
6	75% Percentile	89			
7	Maximum	146			
8					
9	Mean	79.25			
10	Std. Deviation	25.05			
11	Std. Error of Mean	3.261			
12					
13	Lower 95% CI of mean	72.72			
14	Upper 95% CI of mean	85.77			
15					
16	Sum	4676			
17					
18	D'Agostino & Pearson normality test				
19	K2	5.95			
20	P value	0.0510			
21	Passed normality test (alpha=0.05)?	Yes			
22	P value summary	ns			
23					
24					



What did we learn?

- The normal distribution often results from many random factors creating variability and therefore, is common in a research setting.
- 3 main characteristics of a normal distribution: 1) the data are unimodal, 2) the distribution is symmetric, and 3) the frequencies decline steadily as we move towards higher and lower values, without any sudden sharp cut-off.
- If the data are normally distributed, 68% of observations will fall within 1 SD of the mean and 95% of observations will fall within 2 SD.
- Skewness is a quantitative measure of the lack of symmetry in the data distribution and kurtosis is a quantitative measure of how closely the peak, shoulders, and tails of the distribution match a normal distribution.
- Use tests of non-normality with caution.