

EXAMINING YOUR DATA

Chapter 2: Multivariate Insurance

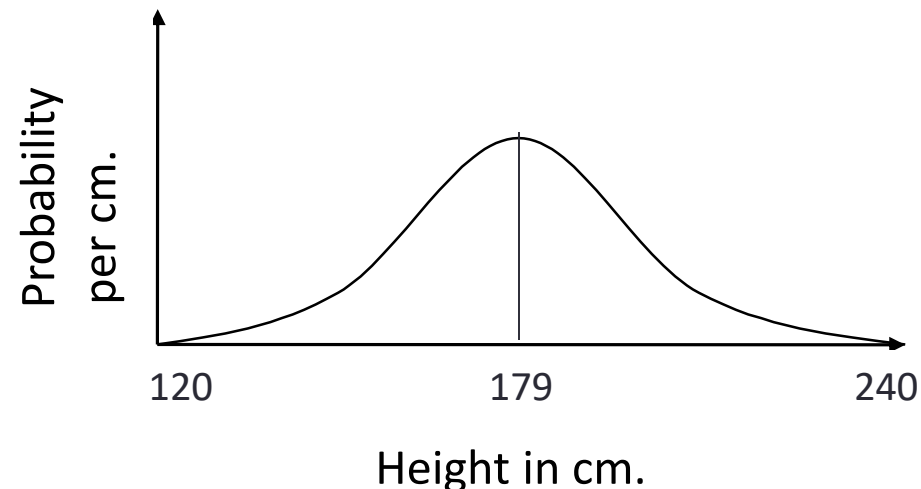
Topics

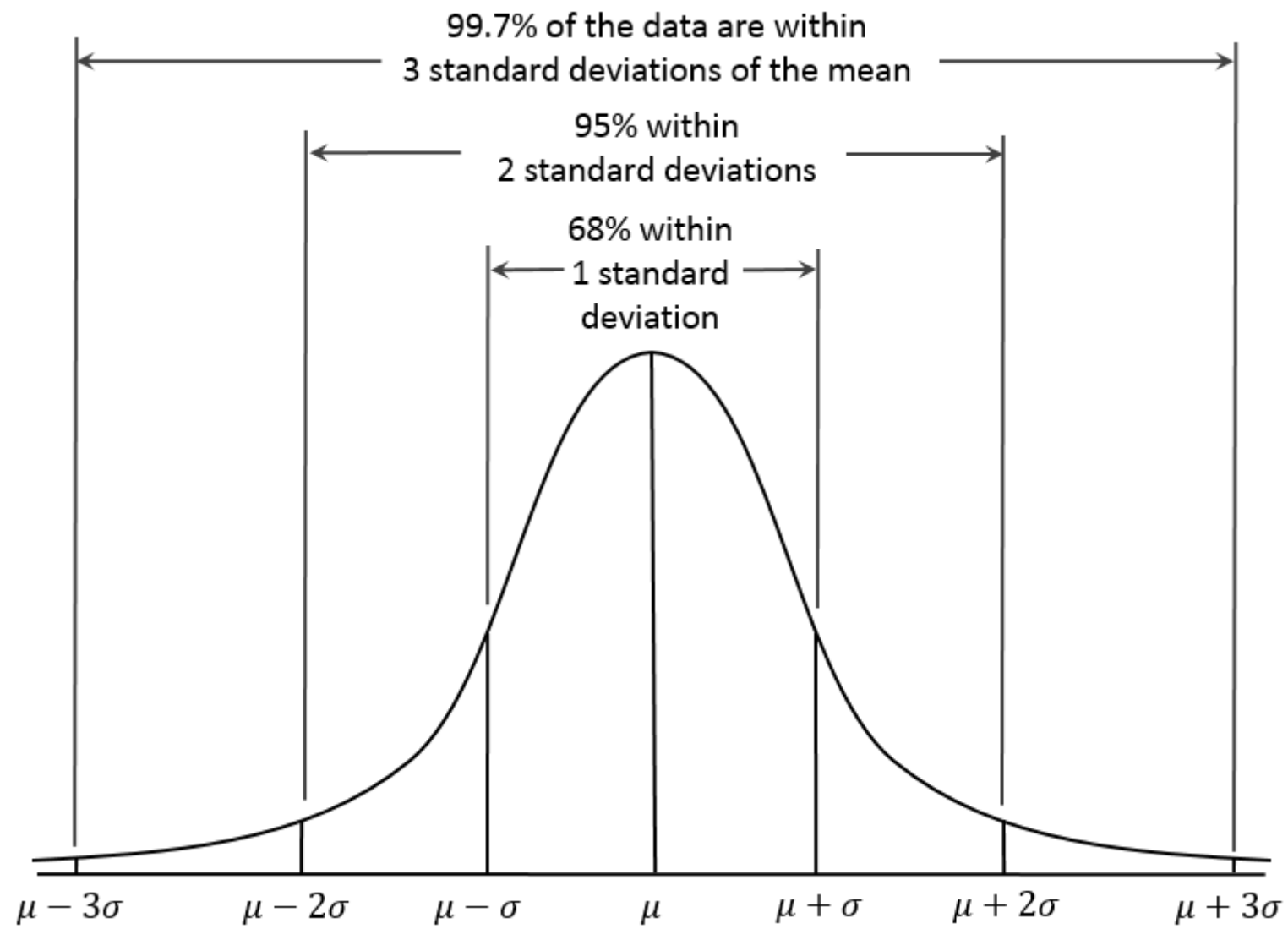
1. The normal distribution and normality.
2. Relationships between variables (correlation).
3. Outliers.
4. Missing data.

1. The Normal Distribution

A probability sample of male, 19 year old Swedes had their height measured. The majority were 179 cm tall, with some shorter and some taller.

The distribution is symmetric around the average (the mean \bar{X}).





Checking Normality

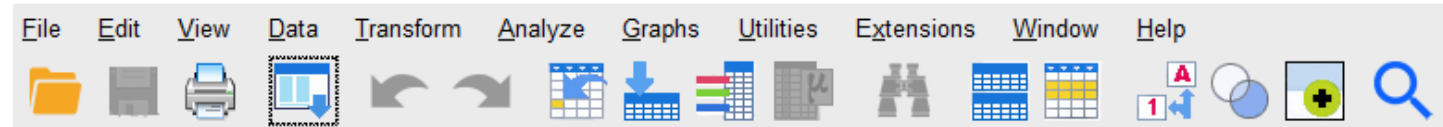
You can examine normality visually or statistically. There is no, "one best way", so it is often good to *consider several*.

- Skewness & Kurtosis
- Histograms: not so good with small samples.
- Statistical tests like Kolmogorov-Smirnov or Shapiro-Wilk ($n < 50$).
- Box plots: Handy for identifying outliers.

How normal does the data need to be?

- It depends on the multivariate technique.
- e.g. Linear regression is quite robust against violations of normality.

Now we get SPSS going for everyone in the class

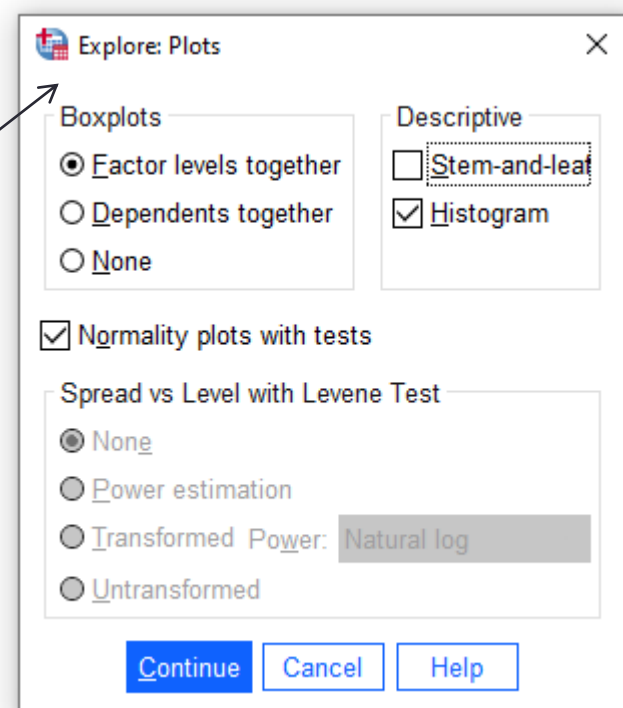
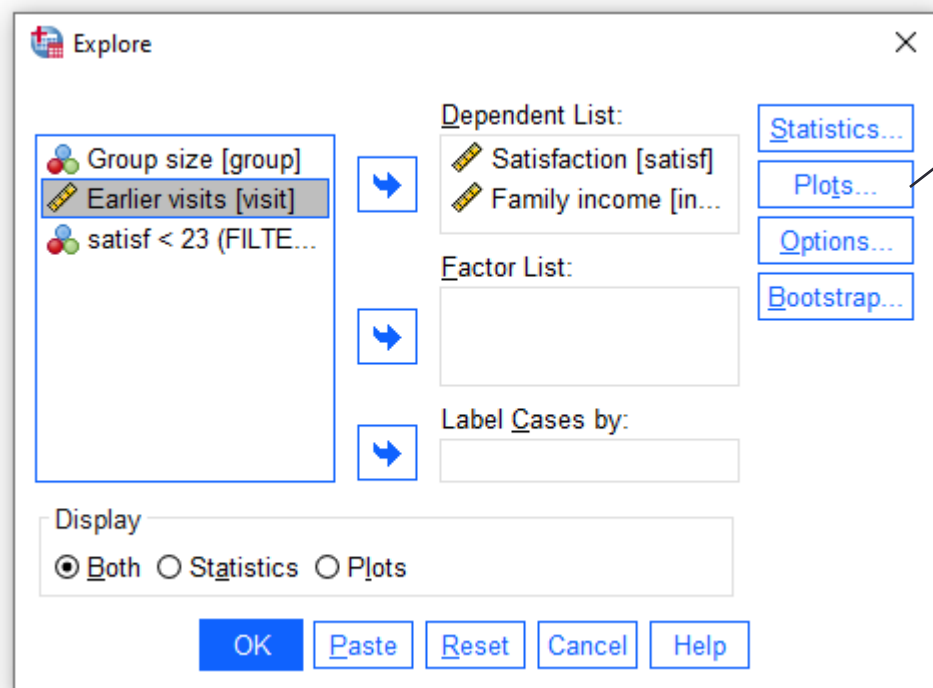


	satisf	group	visit	income	filter_\$	var	var	var	va
1	12	3	2	380	1				
2	9	3	3	320	1				
3	12	2	1	430	1				
4	11	3	3	370	1				
5	8	2	4	270	1				
6	13	4	3	500	1				
7	9	2	4	230	1				
8	9	2	4	260	1				
9	14	2	1	300	1				
10	17	2	1	360	1				
11	11	3	2	290	1				
12	12	3	2	370	1				
13	13	4	2	500	1				
14	13	3	2	370	1				
15	15	3	2	400	1				
16	11	5	3	450	1				
17	11	4	3	400	1				
38	9	2	4	380	1				
39	9	4	3	390	1				
40	8	3	2	360	1				
41									
42									
43									

<

Data View

Variable View



Skewness & Kurtosis

Descriptives			Statistic	Std. Error
Satisfaction	Mean		11,33	,535
	95% Confidence Interval for Mean	Lower Bound	10,24	
		Upper Bound		
	Skewness		1,188	,374
Family income	Kurtosis		2,865	,733
	Mean		369,75	16,316
	95% Confidence Interval for Mean	Lower Bound	336,75	
		Upper Bound	402,75	
	5% Trimmed Mean		365,00	
	Median		370,00	
	Variance		10648,654	
	Std. Deviation		103,192	
	Minimum		200	
	Maximum		700	
	Range		500	
	Interquartile Range		148	
	Skewness		,624	,374
	Kurtosis		1,359	,733

Rule of Thumb:
Less than 1 (absolute value) is normal

Side note: If you ask a statistician who calculates skewness by hand, s/he will say the cutoff is +/- 3. most software standardizes both measures to +/- 1.

Statistical Tests

Rule of Thumb:
Big (> 0.05) is good (normal).

Tests of Normality

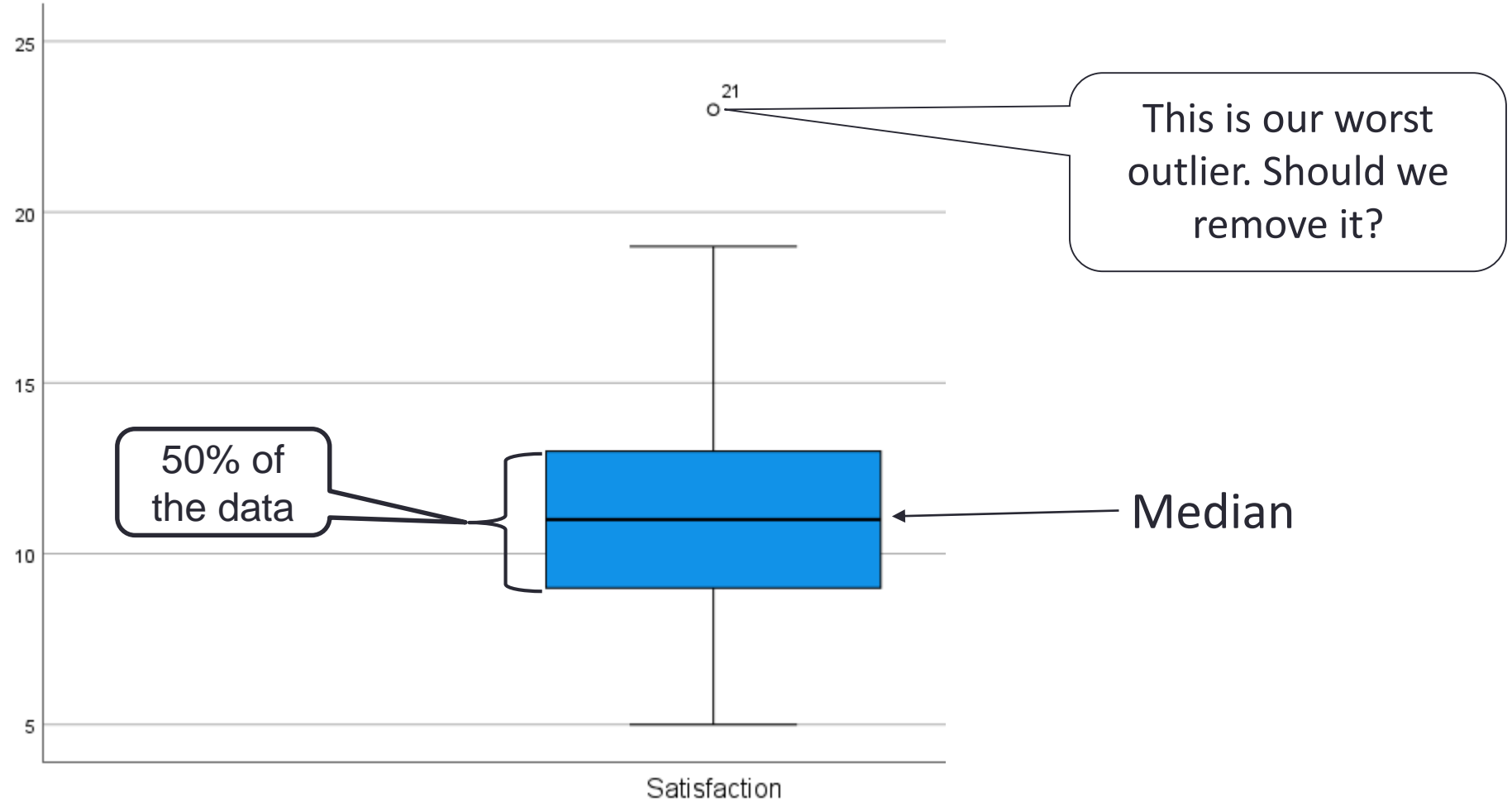
	Kolmogorov-Smirnov ^a			Shapiro-Wilk			
	Statistic	df	Sig.	Statistic	df	Sig.	
Satisfaction	,146	40	,032	,917	40	,006	Abnormal
Family income	,087	40	,200 [*]	,956	40	,122	Normal

*. This is a lower bound of the true significance.

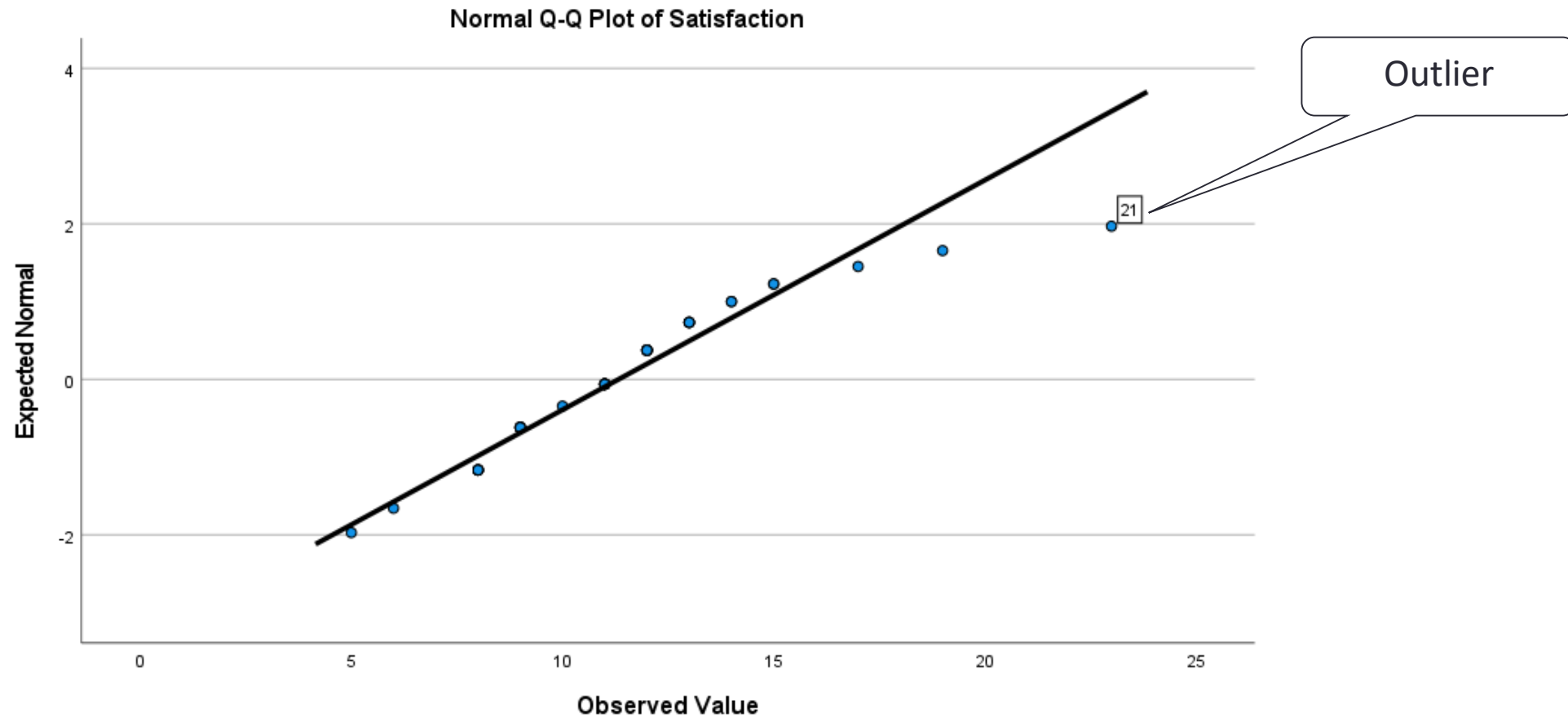
a. Lilliefors Significance Correction

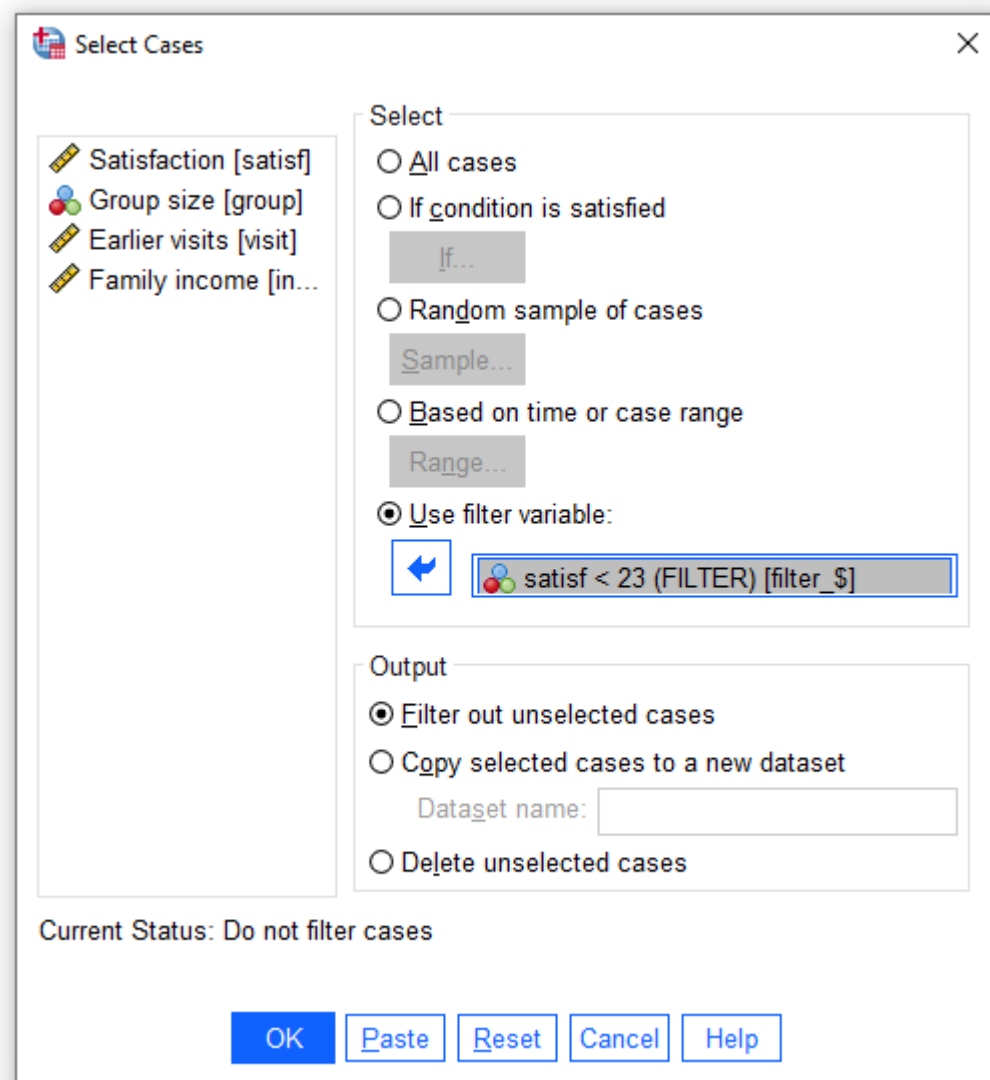
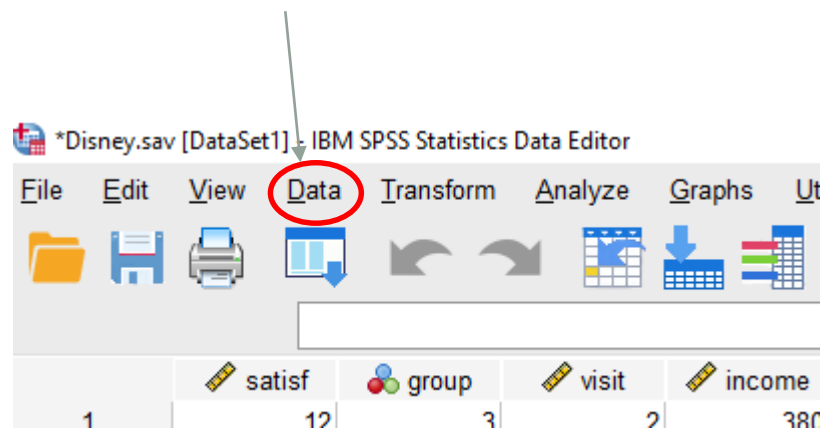
We look here because $n < 50$.

Boxplot



Normal Q-Q Plot





Without the Outlier (21)

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Satisfaction	,121	39	,158	,965	39	,267
Family income	,110	39	,200 [*]	,972	39	,430

*. This is a lower bound of the true significance.

a. Lilliefors Significance Correction

Now both are normally distributed.

Rules of Thumb!

- Consider which technique you are using before worrying too much about normality.
- Remove outliers.
- Consider transforming the data.
- With larger sample sizes normality is less of an issue.

2. Relationships Between Variables

In, for example, regression we want strong linear relationships between the independent variables and the dependent variable, but weak linear relationships between independent variables.

- Examine scatterplots.
- Examine bivariate correlations.

Disney.sav [DataSet1] - IBM SPSS Statistics Data Editor

File Edit View Data Transform Analyze **Graphs** Utilities Extensions Window Help

Chart Builder

Variables: Chart prev.

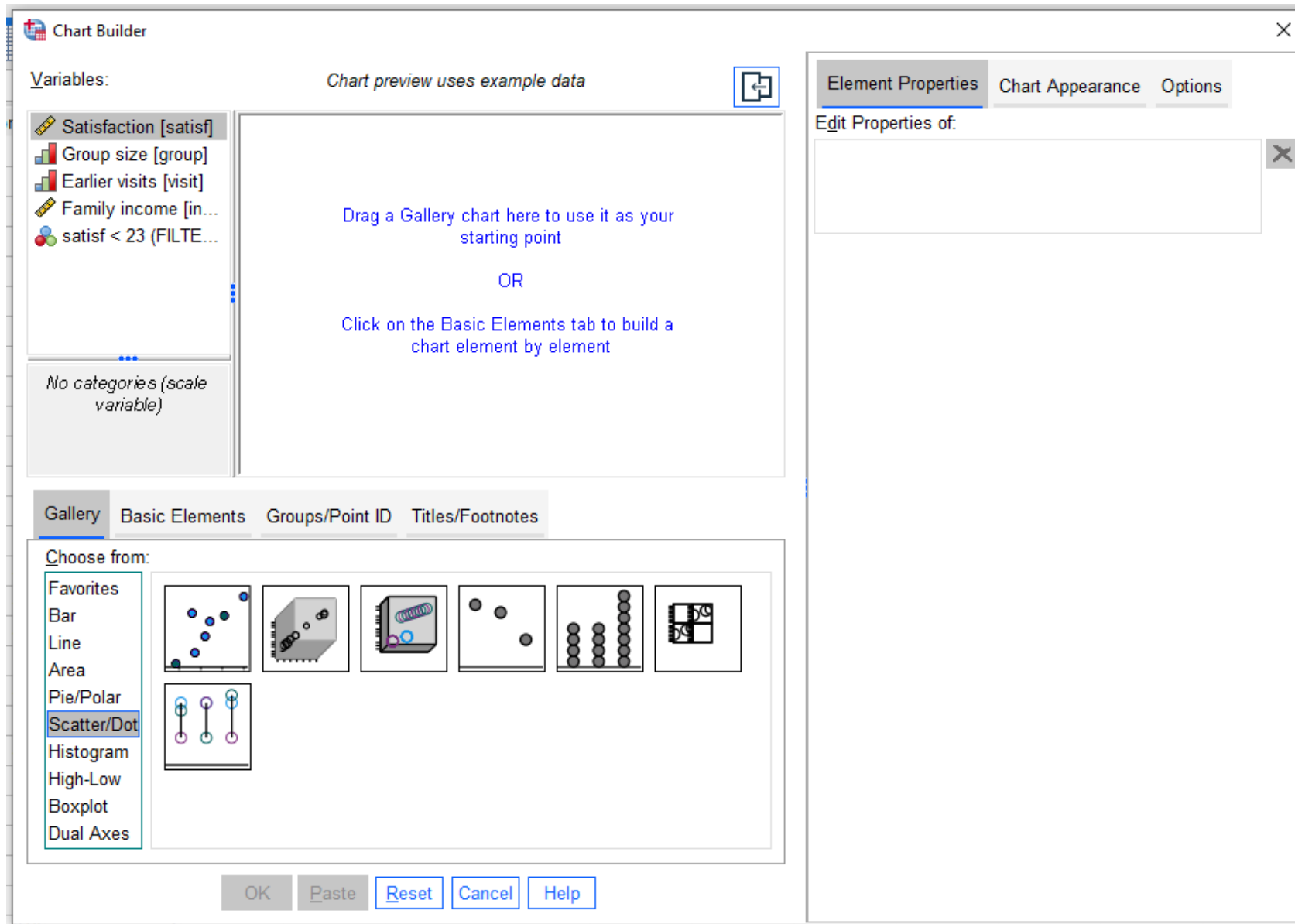
Drag a C

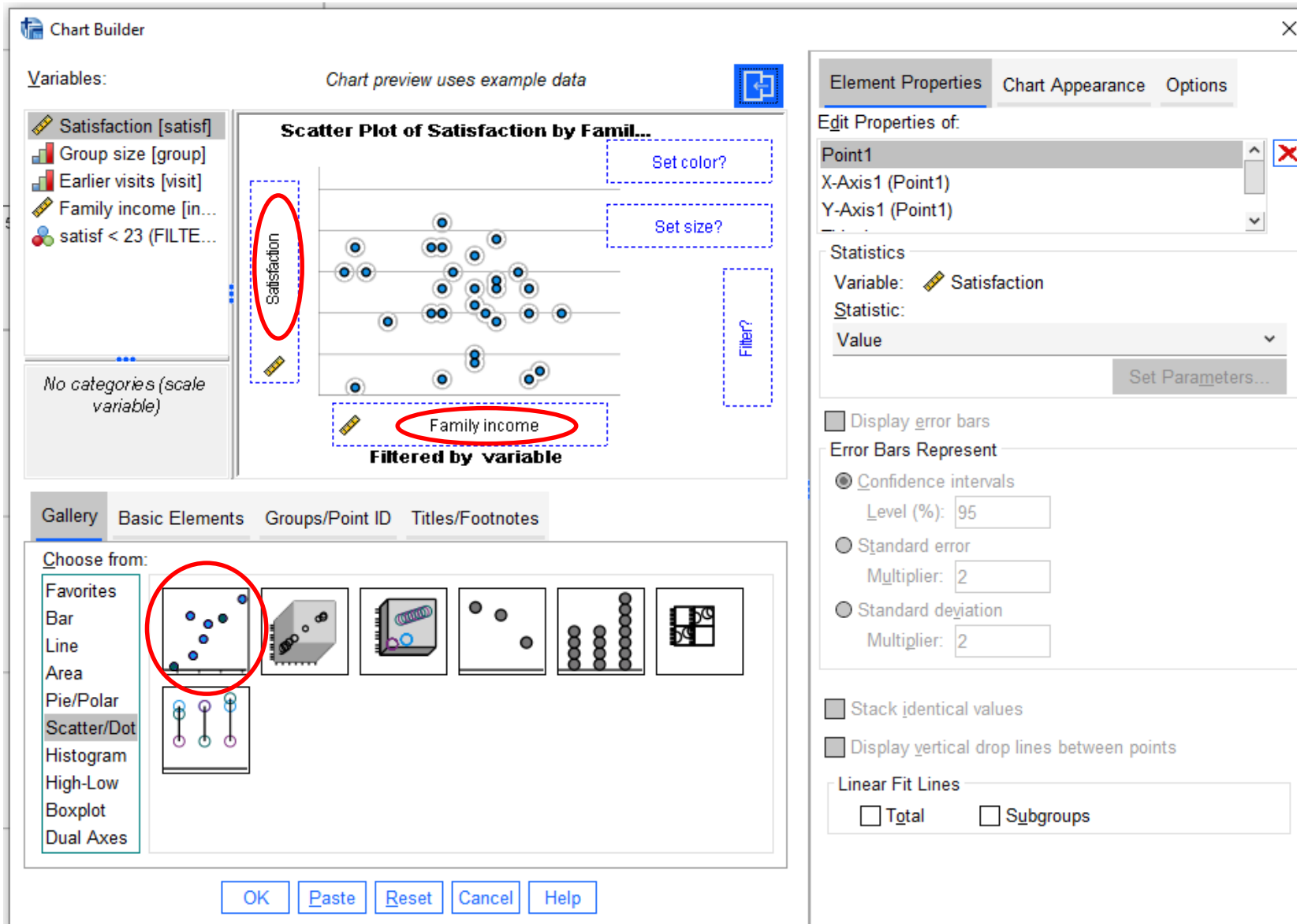
Click on

	satisf	group	visit	income
1	12	3	2	
2	9	3	3	
3	12	2	1	
4	11	3	3	
5	8	2	4	
6	13	4	3	
7	9	2	4	
8	9	2	4	
9	14	2	1	
10	17	2	4	

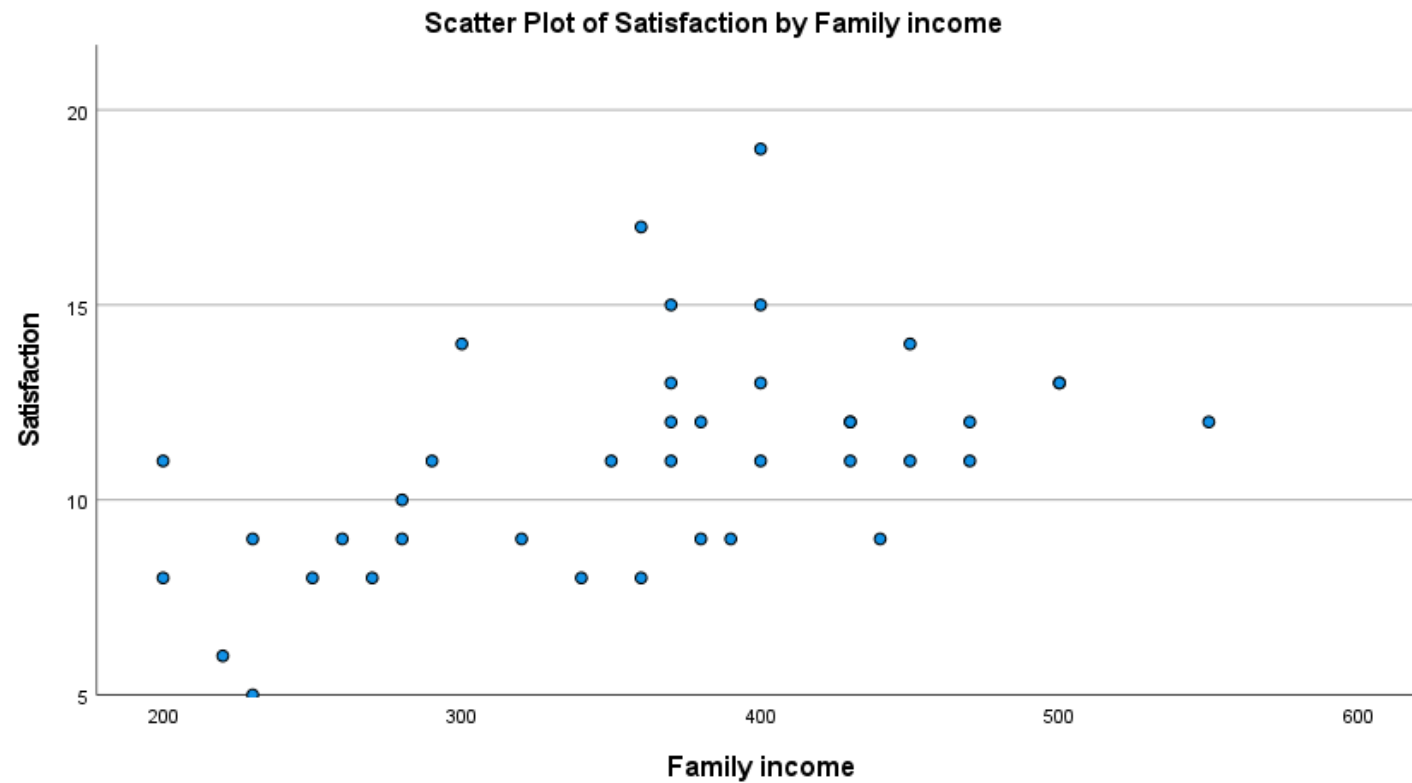
Satisfaction [satisf]
Group size [group]
Earlier visits [visit]
Family income [in...]
satisf < 23 (FILTE...)

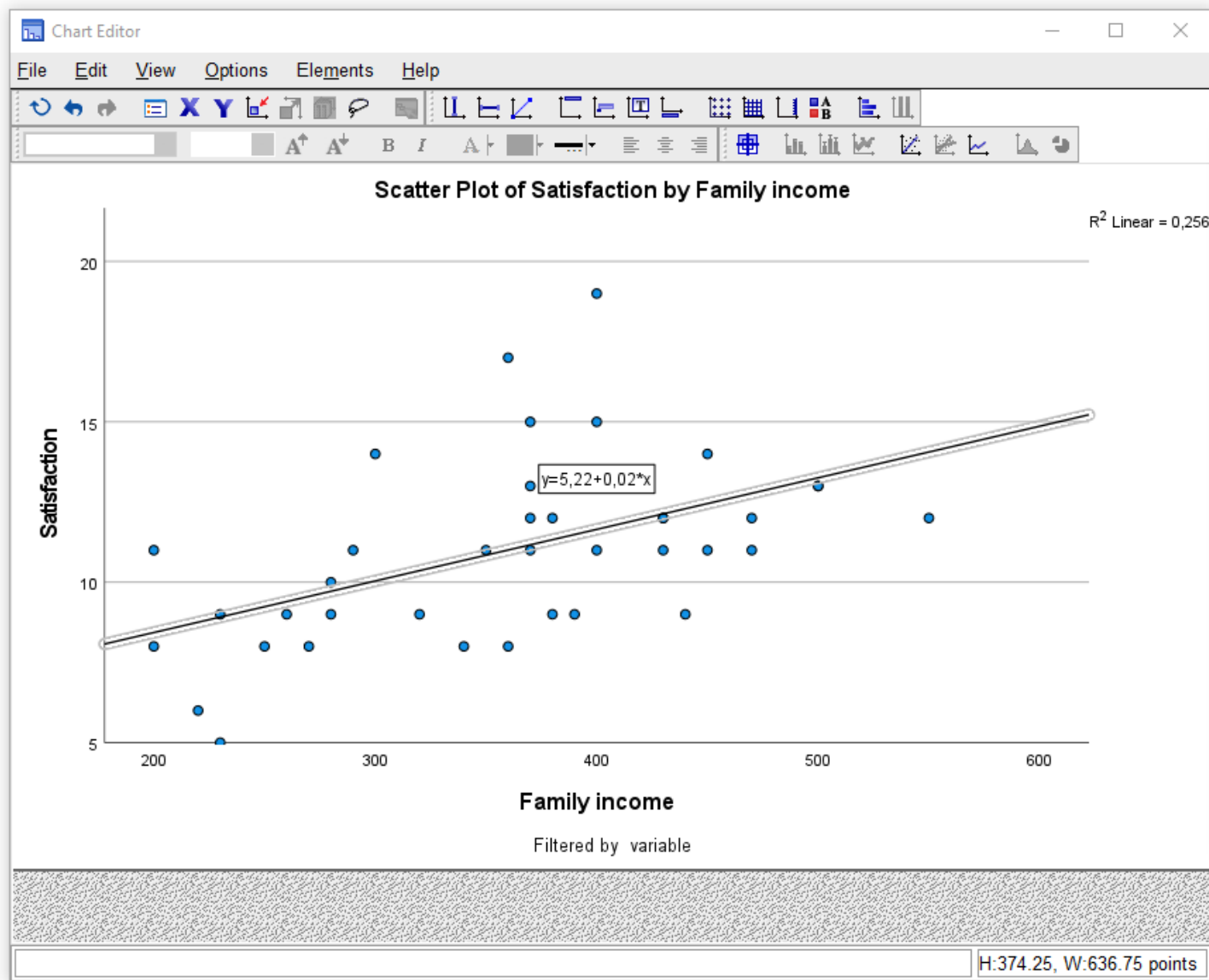
No categories (scale variable)





Scatterplot





Properties

Chart Size Lines **Fit Line** Variables

☐ Display Spikes ☐ Suppress intercept

Fit Method

☐ Mean of Y ☐ Quadratic

☒ Linear ☐ Cubic

☐ Lgess

% of points to fit: 50

Kernel: Epanechnikov

Confidence Intervals

☒ None

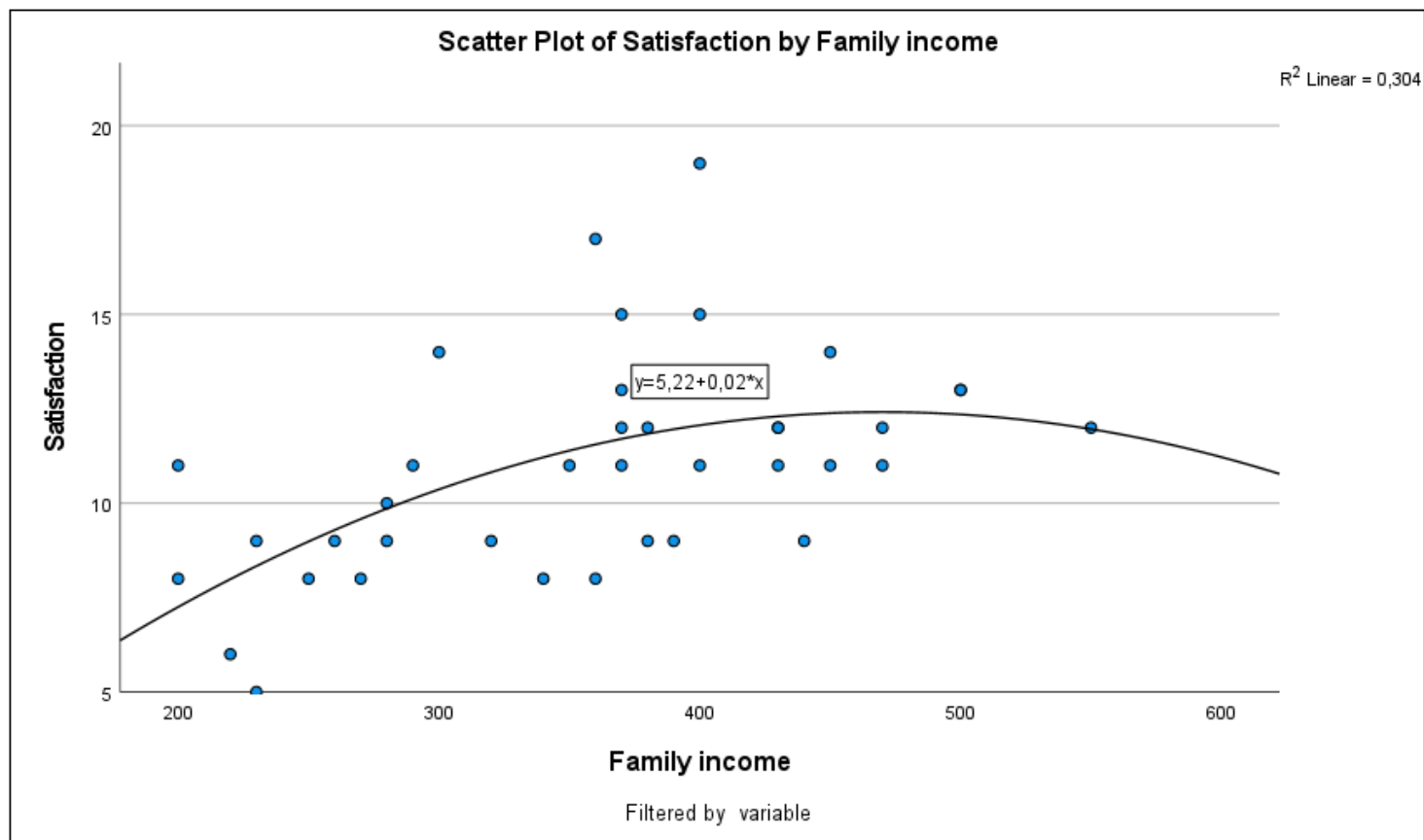
☐ Mean

☐ Individual

%: 95

☒ Attach label to line

Apply Close Help



Properties

Chart Size Lines **Fit Line** Variables

☐ Display Spikes ☐ Suppress intercept

Fit Method

☐ Mean of Y ☒ Quadratic

☐ Linear ☐ Cubic

☐ Lgess

% of points to fit: 50

Kernel: Epanechnikov

Confidence Intervals

☒ None

☐ Mean

☐ Individual

%. 95

☒ Attach label to line

Apply Cancel Help

The thick line represents the range of the true data. The thin line shows values beyond the data.

A. Schmuck et al.

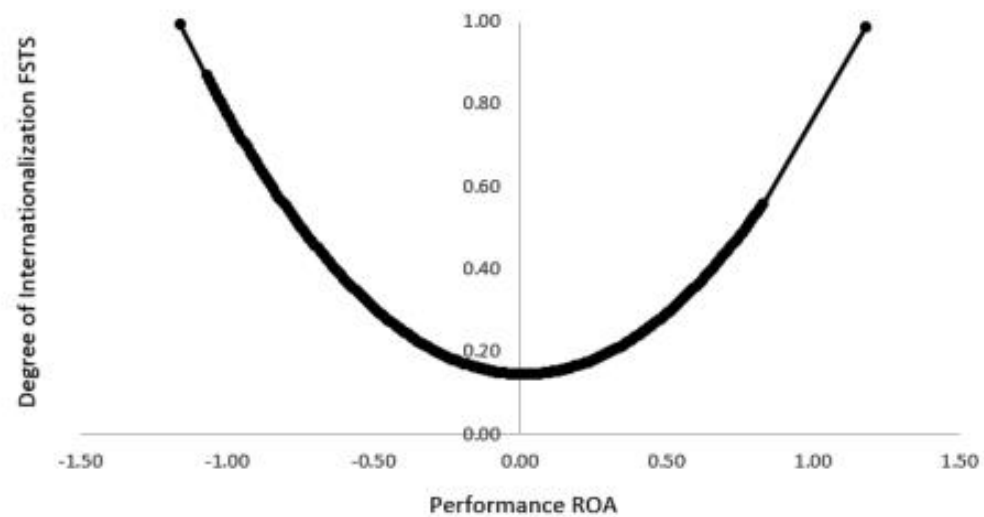
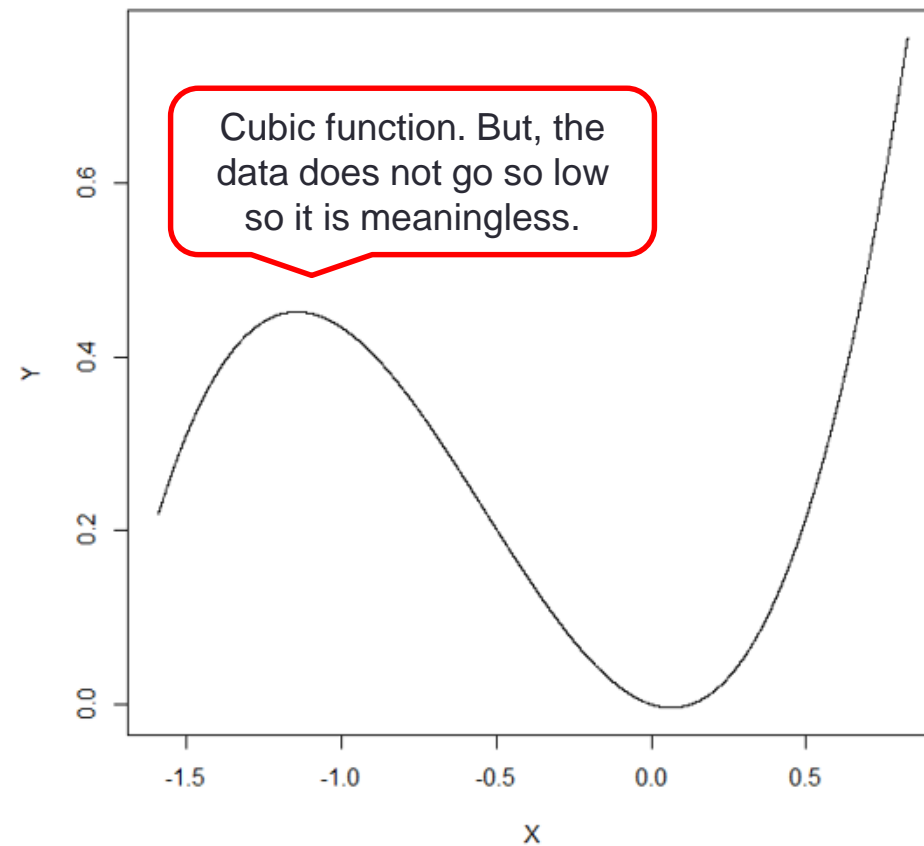


Fig. 1 Regression plot for predicted degree of internationalization



Bivariate Correlations

Variables:

Satisfaction [satisf]
Group size [group]
Earlier visits [visit]
Family income [income]
satisf < 23 (FILTER) [filter_\$]

Options...
Style...
Bootstrap...
Confidence interval...

Correlation Coefficients
☒ Pearson ☐ Kendall's tau-b ☐ Spearman

Test of Significance
☒ Two-tailed ☐ One-tailed

☒ Flag significant correlations ☐ Show only the lower triangle ☒ Show diagonal

OK Paste Reset Cancel Help

Parametric:
When people just say,
"correlation", they mean
Pearson correlation

Non-parametric
(e.g. ordinal)

Correlations

		Correlations			
		Satisfaction	Group size	Earlier visits	Family income
Satisfaction	Pearson Correlation	1	,282	-,573**	,506**
	Sig. (2-tailed)		,081	<,001	,001
	N	39	39	39	39
Group size	Pearson Correlation	,282	1	-,132	,764**
	Sig. (2-tailed)	,081		,423	<,001
	N	39	39	39	39
Earlier visits	Pearson Correlation	-,573**	-,132	1	-,345*
	Sig. (2-tailed)	<,001	,423		,031
	N	39	39	39	39
Family income	Pearson Correlation	,506**	,764**	-,345*	1
	Sig. (2-tailed)	,001	<,001	,031	
	N	39	39	39	39

** . Correlation is significant at the 0.01 level (2-tailed).

* . Correlation is significant at the 0.05 level (2-tailed).

Correlations — when considering regression

		Correlations			
		Y	Group size	# of earlier visits	Family income
		Satisfaction			
Satisfaction	Pearson Correlation Sig. (2-tailed) N				
Group size	Pearson Correlation Sig. (2-tailed) N				
X_1					
# of earlier visits	Pearson Correlation Sig. (2-tailed) N				
X_2					
Family income	Pearson Correlation Sig. (2-tailed) N				
X_3					

Want high correlations

Want low correlations (under 1. 91)

Rule of Thumb!

- Correlations between independent variables should not exceed 0.9. If they do you most likely have problems with multicollinearity.
- Consider the substantive meaning of correlations. Small correlations, even if they are significant, are meaningless. By small I mean below |.3|.
- Cohen (1988) established practical guidelines for interpreting statistically significant correlations.

Table 8.14 Correlation strength

Correlation strength	Absolute value of correlation coefficient
Small	0.10–0.29
Medium	0.30–0.49
Large	0.50–1.00


3. Outliers

- There is no rule of thumb. You simply have to look at them and decide if they are an aberration or if they are a natural part of the data.
- Examine how they influence your results.

Article

Best-Practice Recommendations for Defining, Identifying, and Handling Outliers

**Herman Aguinis¹, Ryan K. Gottfredson¹,
and Harry Joo¹**

Organizational Research Methods
16(2) 270-301
© The Author(s) 2013
Reprints and permission:
sagepub.com/journalsPermissions.nav
DOI: 10.1177/1094428112470848
orm.sagepub.com


4. Missing Data

Two questions:

1. Why is it missing?

- You need to determine if the missing data is missing at random, or if there is some sort of pattern as to why it is missing.

2. How much is missing?

- You have to consider how much data is missing. If not so much, just run your analyses with the default "listwise deletion".
- Pallant often suggests pairwise deletion – this way you don't lose so much data.
- If a lot is missing, you may opt to replace missing values.
- You may drop variables with a lot of missing values.

Four Steps

1. Type of missing data: ignorable – not ignorable?
 - Ignorable if substantively missing at random.
2. Extent of missing data: If not ignorable, is there so much missing as to bias results?
 - 10% cutoff for no problem.
3. Diagnose randomness (many months & large brain).
4. What to do?
 - Listwise deletion (complete case approach).
 - Pairwise deletion.
 - Imputation.

Rule of Thumb!

In SPSS you can test for differences between missing and non-missing data across all variables.

If there are significant patterns – hmm.

- Try to avoid missing data.
- If the % of missing values is low, ignore the patterns.
 - Around 10%.
- Are they variables you really need?
- Are the missing values concentrated in a few cases? Consider deleting them.

Use the “Banking” data

Missing Value Analysis

Quantitative Variables:

- Education [educ...]
- Gender [gender]
- Satisfaction (X1)...
- Trust (X2) [trust]
- Commitment (X3)

Categorical Variables:

Maximum Categories: 25

Case Labels:

Use All Variables

Patterns...

Descriptives...

Estimation

- ☐ Listwise
- ☐ Pairwise
- ☐ EM
- ☐ Regression

Variables...

EM...

Regression...

OK Paste Reset Cancel Help

Missing Value Analysis: Patterns

Display

- ☐ Tabulated cases, grouped by missing value patterns
Omit patterns with less than 1 % of cases
- ☒ Sort variables by missing value pattern
- ☐ Cases with missing values, sorted by missing value patterns
- ☒ Sort variables by missing value pattern
- ☒ All cases, optionally sorted by selected variable

Variables

Missing Patterns for:

- education
- gender
- satisfac
- trust
- commit
- loyalty
- satisf2

Additional Information for:

Sort by:

Sort Order

- ☒ Ascending
- ☐ Descending

Continue Cancel Help

Missing data

↓ Delete a bad variable?

Data Patterns (all cases)									
Case	#Missing	% Missing	Missing and Extreme Value Patterns						
			education	gender	satisfac	trust	commit	loyalty	satisf2
1	0	,0							
2	0	,0							
3	0	,0							
4	0	,0							
5	0	,0							
226	1	14,3						A	
227	0	,0							
228	3	42,9				A	A	A	
229	0	,0							
230	0	,0							
231	1	14,3				A			
232	0	,0							
233	0	,0							
234	5	71,4			A	A	A	A	A
235	1	14,3						A	

Delete a bad case?

