

Lecture 6

Empirical Research Methods IN4304

Data preparation methods

IN4304 Empirical Research Methods Spring 2010 Lecture 6

1



TU Delft

Previous lecture

participant-observation and non-participant observation

Does the observer act as a
member of the group?

Sampling strategies

Event sampling, state sampling,
interval sampling, time
sampling

Coding scheme

- Focused
- Objective
- No context-dependent
- Explicitly defined
- Exhaustive
- Mutually exclusive
- Easy to record

Inter-observer agreement (Cohen's Kappa)

$$K = (P_0 - P_e) / (1 - P_e)$$

IN4304 Empirical Research Methods Spring 2010 Lecture 6

2



TU Delft

Today

- Data coding
- Measures of central tendency
- Measures of variability
- Explorative Data Analysis
- Data cleaning
- Normal distribution
- Standardizing data
- Confidence intervals

IN4304 Empirical Research Methods Spring 2010 Lecture 6

3



TU Delft

Learning outcomes of this lecture

After today's lecture you should be able :

- to enter your data into SPSS
- to examine the central tendency of your data set
- to understand the output of some explorative data analysis techniques
- to reflect up on reasons for data cleaning
- to explain the importance of normal distribution

IN4304 Empirical Research Methods Spring 2010 Lecture 6

4



TU Delft

Data coding

- Coding is needed to process the data answers of questionnaire or observation

How much do you like this lecture?

(tick one answer)

Very much	<input type="radio"/>	(code = 1)
Quite a lot	<input type="radio"/>	(code = 2)
Moderately	<input type="radio"/>	(code = 3)
Not very much	<input type="radio"/>	(code = 4)
Not at all	<input type="radio"/>	(code = 5)

Not applicable (code 900)

No answer (missing value) (code = 999)

IN4304 Empirical Research Methods Spring 2010 Lecture 6

5



TU Delft

Entering data into SPSS

- Give variable a clear name and label
- Specify data type (string, number, date etc)
- Specify values code
- Specify missing value code
- Specify level of variable (nominal, ordinal, interval)

SPSS Demo

IN4304 Empirical Research Methods Spring 2010 Lecture 6

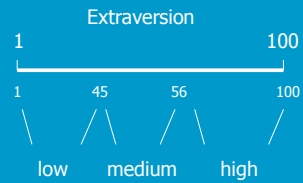
6



TU Delft

Changing data from one level to another

- Changing data from interval to ordinal level
- SPSS recode function



IN4304 Empirical Research Methods Spring 2010 Lecture 6

7



TU Delft

Direction scale

	Sophistication							
Childish	○	○	○	○	○	○	○	Sophisticated
Classy	○	○	○	○	○	○	○	Silly
Novelty	○	○	○	○	○	○	○	Business

Code1 : 1 2 3 4 5 6 7
Code2: -3 -2 -1 0 1 2 3

- Classy-Silly needs to be reversed
- For example:
 - code1 : 8-value
 - Code2 : -1*value
- SPSS command: *Compute*

IN4304 Empirical Research Methods Spring 2010 Lecture 6

8



TU Delft

Measures of central tendency

- Interval level: Mean
- Ordinal level: Median
- Nominal level: Mode

What is the best measure to represent:

1 3 3 3 4 7 8 8 9 10 11 11 12 13 13 15 16 17 17 19 10000
 Mean: 485.71
 Median: 11
 Mode : 3
 5% trimmed mean: 10.47

IN4304 Empirical Research Methods Spring 2010 Lecture 6

9

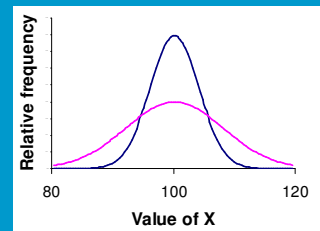


TU Delft

Measures of variability

To express the distribution of the data set.

Mean is 100. The spread is different.



IN4304 Empirical Research Methods Spring 2010 Lecture 6

10



TU Delft

Range

Range of a distribution is the difference between the highest and the lowest value in a set of numbers

$$\text{Range} = \text{Max} - \text{Min}$$

For example

8 9 10 2 4 3 100 12 80

$$\text{Range} = 100 - 2 = 98$$

IN4304 Empirical Research Methods Spring 2010 Lecture 6

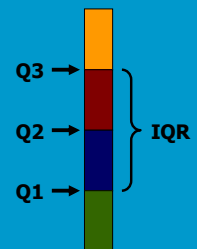
11



TU Delft

Inter Quartile Range

- A quartile is any of the three values which divide the sorted data set into four equal parts
- First quartile (Q1) = cuts off lowest 25% of data = 25th percentile
- Second quartile (Q2) = Median = cuts data set in half = 50th percentile
- Third quartile (Q3) = cuts off highest 25% of data, or lowest 75% = 75th percentile
- Inter Quartile Range (IQR) = Q3-Q1
- The IQR is a more stable statistic than the range, as it is less effected by extreme values.



IN4304 Empirical Research Methods Spring 2010 Lecture 6

12



TU Delft

Standard Deviation (SD)

Standard Deviation measures the spread from the mean. You can think of it as the mean deviation from the mean

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

Population mean: μ
 Sample mean: \bar{x}
 Population SD: σ
 Sample SD: s

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

IN4304 Empirical Research Methods Spring 2010 Lecture 6

13



Explorative Data Analysis (SDA)

Examine your data before you apply statistical analysis on it.

Use techniques such as

- Box plot
- stem and leaf plot
- histogram

IN4304 Empirical Research Methods Spring 2010 Lecture 6

14



SPSS – explorative command

Stemplot (stem-and-leaf plot)

Number of times a week that a website was visited Stem-and-Leaf Plot

Frequency	Stem & Leaf	Weekno	Hits
11	0 . 00012233447	43	3
2	1 . 22	41	5
3	2 . 057	44	9
7	3 . 1357888	42	12
6	4 . 234569	37	25
1	5 . 0	36	28
5	6 . 25689	38	34
4	7 . 2359	39	35
2	8 . 69	35	44
		40	46
		14	72
		22	121

Stem width: 100
 Each leaf: 1 case(s)

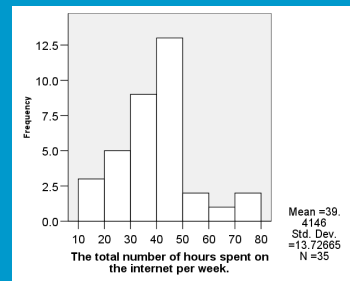
IN4304 Empirical Research Methods Spring 2010 Lecture 6

15



Histogram

- Gives an overview of shape of the distribution
- SPSS – graph - histogram

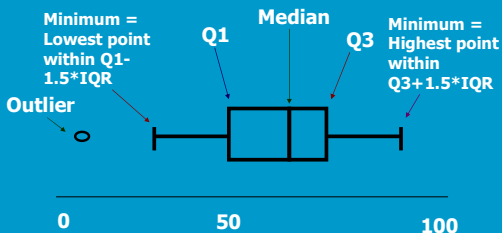


IN4304 Empirical Research Methods Spring 2010 Lecture 6

16



Box plot



SPSS – Explorative analysis or graph command

IN4304 Empirical Research Methods Spring 2010 Lecture 6

17



Outliers

- A case whose value is very different from most others
- Can bias statistics such as mean
- What should you do with outliers?
 - Exclude them?
 - Investigate them?
 - Include them in the analysis?

IN4304 Empirical Research Methods Spring 2010 Lecture 6

18



Data cleaning

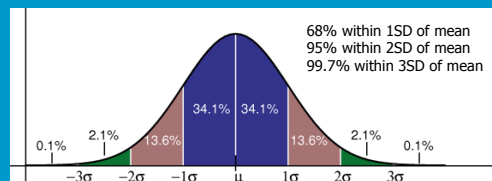
- Reason to remove cases
 - Observation was distorted, e.g. did not understand task, questionnaire, purposely messing up your experiment..., error in measuring instrument
 - Outliers, but only statistical reason is questionable
 - Check your data for unlikely data points

IN4304 Empirical Research Methods Spring 2010 Lecture 6

19



Normal distribution



Why normal distribution important?

- Gives a good model of data for some real world data sets (e.g. large populations)
- Good approximation of results of random outcomes (e.g. throwing a dice many times)
- Large number of inferential statistics are based on normal distributions (Moore and McCabe, 2006)

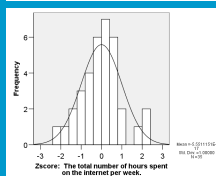
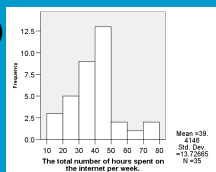
IN4304 Empirical Research Methods Spring 2010 Lecture 6

20



Standardizing data (1)

- Transforming data into z-score, makes it possible to compare data from different unit of measurement



$$z = \frac{x - \mu}{\sigma}$$

$$z = \frac{x - \bar{x}}{s}$$

IN4304 Empirical Research Methods Spring 2010 Lecture 6

21



Standardizing data (2)

Repeated measure design – focus on relative difference

- z-scores to overcome individual use of questionnaire scale e.g.
 - People that only use the extremes
 - People that never use the extremes
 - People that only give high marks

For each individual calculate mean and sd across the his/her questions

Calculate z-score for value an individual gave to question based on this persons mean and sd of his/her questions

Unlikely

- extremely
- quite
- slightly
- neither
- slightly
- quite
- extremely

Likely

IN4304 Empirical Research Methods Spring 2010 Lecture 6

22



Standardizing data (3) – Example

Participant 1 answered:

Q1=1 Q2=3 Q3=7 Q4=7

z-score would be:

Q1= -1.17 Q2=-0.50 Q3=0.83 Q4=0.83

Participant 2 answered:

Q1=4 Q2=5 Q3=6 Q4=6

z-score would be:

Q1=-1.31 Q2=-0.26 Q3=0.78 Q4=0.78

Unlikely

- extremely
- quite
- slightly
- neither
- slightly
- quite
- extremely

Likely

IN4304 Empirical Research Methods Spring 2010 Lecture 6

23

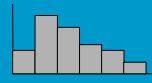


Medium, mean and distribution shape



Symmetrical
(median = mean)

Notice that for the sample
2 8 110 2
The median = 5
The mean = 30.5



Positively skewed or skewed to the right (median < mean)



Negatively skewed or skewed to the left (median > mean)

IN4304 Empirical Research Methods Spring 2010 Lecture 6

24



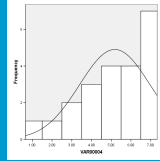
Skewness

Measure

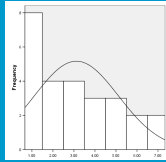
Skewness = 0 --> symmetrical distribution

Skewness > 0 --> positively skewed

Skewness < 0 --> negatively skewed



Ceiling effects



Floor effect

IN4304 Empirical Research Methods Spring 2010 Lecture 6

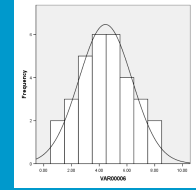
25



TU Delft

Check for normality (1)

Look at the histogram and the projected normal curve



Rule of the thumb for serious deviation from normality

- skewness (or kurtosis) statistic > 2*Std Error (i.e. sd/\sqrt{n})
- Difference between mean and median > 0.5*SD (Coolican, 2004)

IN4304 Empirical Research Methods Spring 2010 Lecture 6

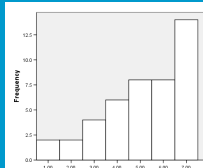
26



TU Delft

Check for normality (2)

- Kolmogorov-Smirnov Normality test



VAR00004	Kolmogorov-Smirnov ^b		
	Statistic	df	Sig.
	.178	44	.001

a. Lilliefors Significance Correction

Remember:

- it is easier to show deviation from normality with large samples
- It does not tell about the size of deviation

Probability that this sample was obtained from normal distributed population

IN4304 Empirical Research Methods Spring 2010 Lecture 6

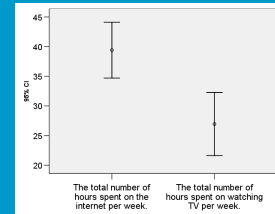
27



TU Delft

Confidence intervals

- If we assume the sample is taken from a normal distribution it is possible to make a **prediction** of the interval in which the population mean might be with e.g. 95% certainty.



IN4304 Empirical Research Methods Spring 2010 Lecture 6

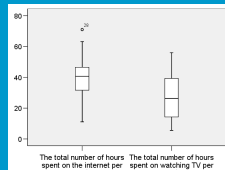
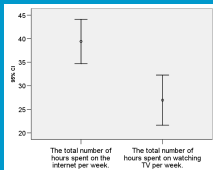
28



TU Delft

Class Questions

- What is the difference between box plot and confidence interval?



IN4304 Empirical Research Methods Spring 2010 Lecture 6

29



TU Delft

Summary

Data Entry

Code data, direction scales, clear name, missing data, level

Central tendency

Mean, median, mode, IQR, SD

Explorative data analysis

- Boxplot
- Histogram
- Stem and leaf plot

Normal distribution

1. Gives a good model of data for some real world data sets (e.g. large populations)
2. Good approximation of results of random outcomes (e.g. throwing a dice many times)
3. Large number of inferential statistics are based on normal distributions

Check for Normality if you want to use parametric statistics

IN4304 Empirical Research Methods Spring 2010 Lecture 6

30



TU Delft

This week in practicum

Explorative Data Analysis

- Box-plot
- Scatter plot
- Stem-and-leaf plot

IN4304 Empirical Research Methods Spring 2010 Lecture 6

31



TU Delft

Next time

Quantitative Data analysis I – Differences

- Hypothesis testing
- Chi-square test, t -test
- (M)ANOVA
- Mann-Whitney U-test
- (Robson ch. 13)

IN4304 Empirical Research Methods Spring 2010 Lecture 6

32



TU Delft

References

- Coolican, H., (2004). *Research methods and statistics in psychology* (4th ed). London, UK: Hodder Arnold.
- Moore D.S, and McCabe, G.P. (2006) *Statistiek in de praktijk; theorieboek* (5 ed) Den Haag: The Netherlands, Academic Service.

IN4304 Empirical Research Methods Spring 2010 Lecture 6

33



TU Delft