

*There are three kinds of lies:
lies, damned lies and statistics.*
*Attributed by Mark Twain
to Benjamin Disraeli*

Research Methods and Data Analysis - exercise

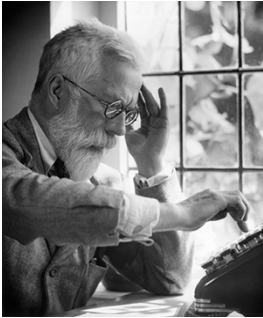
Compiled by
Mart Murdvee

© Mart Murdvee 2000-2013 1

He uses statistics as a drunken man uses a lamppost—for support rather than illumination.
— ANDREW LANG

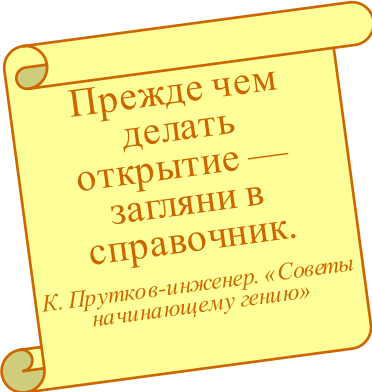
Giving a school man only a little, or very superficial, knowledge of statistics is like putting a razor in the hands of a baby.
— CARTER ALEXANDER

© Mart Murdvee 2000-2013 2



To consult the statistician after an experiment is finished is often merely to ask him to conduct a post mortem examination. He can perhaps say what the experiment died of.
Ronald Fisher

© Mart Murdvee 2000-2013 3




Прежде чем
делать
открытие —
загляни в
справочник.
*К. Прутков-инженер. «Советы
начинающему гению»*

In modern way:
Before You
apply for
Nobel Price,
look into the
Wikipedia!

*Kozma Prutkov – ingénieur.
„Recommendations for
start-up geniuses“*

© Mart Murdvee 2000-2013 4


Every baby knows the
scientific method!



- 1 Make an observation.
- 2 Form a hypothesis.
- 3 Perform the experiment.
- 4 Analyze the data.
- 5 Report your findings.
- 6 Invite others to reproduce the results.

© Mart Murdvee 2000-2013 5

Levels of phenomenon investigation and use of results



1. **Description:**
definitions, facts, measures, structures, conditions, etc.
2. **Explaining:**
regularities, relations, causality, etc.
3. **Prediction:**
using the knowledge about of the regularities predict future changes
4. **Using (manipulation):**
using the knowledge about of the regularities induce desired changes

© Mart Murdvee 2000-2013 6

Statistics



- **Statistics (as scientific discipline)** is the study of the collection, organization, analysis, interpretation, presentation of data and making decisions based on data.
- **Statistics (as practice)** is a set of the methods, procedures, rules and customs for collection, organization, analysis, interpretation, presentation of data and making decisions based on data.

© Mart Murdvee 2000-2013

7

Descriptive and inferential statistics

Descriptive statistics (kirjeldav statistika) – describing data in a generalized / simplified, manageable and understandable format:

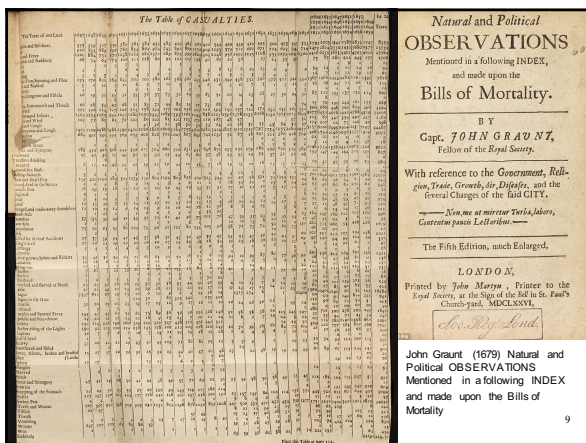
- resumptive statistics (central tendency, variation, distribution, ...);
- presentation of data - graphs, charts, drawings, tables ...

Inferential statistics (järeluslik statistika) - the probability calculations to draw conclusions about the population, to determine the accuracy and reliability of the findings and hypotheses tested:

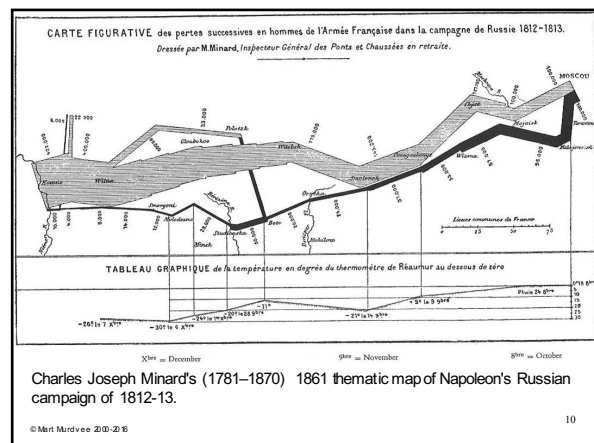
- reliability of statistics (confidence interval...);
- probability / error of differences or correlations (statistical significance, p-value).

© Mart Murdvee 2000-2013

8



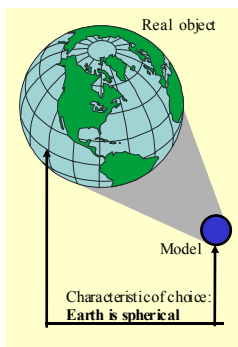
9



© Mart Murdvee 2000-2013

10

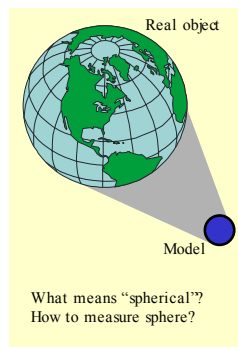
Model



- The model is a description of real world object or phenomenon, which describes the major structures and/or properties of explored object(s).
- Models are inherently incomplete - model includes only some aspects (characteristics) of the real-world object.
- To build the model, some assumptions about the essential and important structures and relations of real-world objects or events must be made.

11

Operationalization

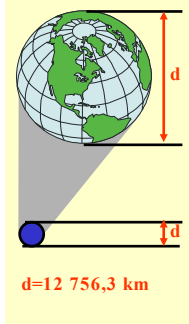


is the process of defining a fuzzy concept so as to make the concept clearly distinguishable or measurable and to understand it in terms of empirical observations.

- First step: defining of the concept
- Operationalization specifies methods for measuring

12

Types of Data



Ratio data – values can be ranked, interval size is known and is same for the whole range of scale. Scale has absolute zero point. For example - height, weight, age, salary, distance.

Interval data – values can be ranked, interval size is known and is same for the whole interval of scale. For example - IQ, SAT score...

Ordinal data – values can be ranked, but the interval is unknown. For example - education, satisfaction, military rank...

Nominal data – values can not be ranked and used directly in analysis. For example - gender, nationality, place of residence... Nominal data can be used only for sample grouping.

Dichotomic nominal data (such as gender) can be used directly for the statistical analysis (eg, correlation).

© Mart Murdvee 2000-2013

13



Using Excel

© Mart Murdvee 2000-2013

14

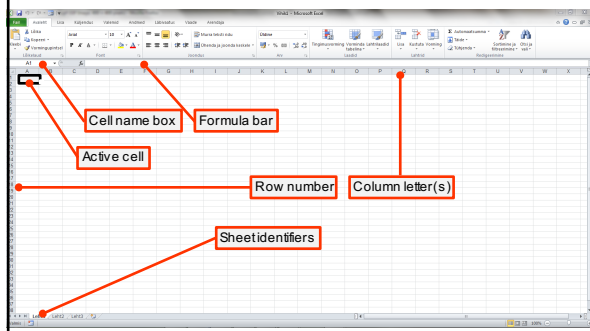
Some useful Windows shortcuts

Press this key	To do this
Ctrl+C (or Ctrl+Insert)	Copy the selected item
Ctrl+V (or Shift+Insert)	Paste the selected item
Ctrl+Z	Undo an action!!!
Ctrl+Y	Redo an action
Ctrl+A	Select all items (in a document)
Esc	Cancel the current task
Left Alt+Shift	Switch the input language when multiple input languages are enabled

© Mart Murdvee 2000-2013

15

Excel 2010 screen parts



© Mart Murdvee 2000-2013

16

Excel: Managing screen and data

Widening a column

Place your cursor over the column headings (A, B, C etc.) You can drag the edges of the column heading (when the cursor becomes \leftrightarrow) to widen the column. If you double-click the edge of the column heading, the column widens to fit the widest piece of data in it.

Inserting a row or column

Select row or column. Right-click and add. Rows are inserted above your cursor. Columns are inserted to the left of your cursor.

Copying and moving data

Select the cells containing the data with the mouse. Right-click and Copy (Ctrl+C) to copy the data. Move your cursor to where you want the data to start. Paste (Ctrl+V) or select paste mode. Going to another location and pressing Ctrl+V will copy the data there as well.

The last time you insert the data you can also use Enter instead of Ctrl+V.

Clearing the data from a row or column

Select the row(s) or column(s) by clicking the numbers to the left and the letters above respectively. Press the Delete key.

Deleting a row or column

Select the row(s) or column(s). Right-click and select Delete.

Keeping columns and rows visible as you scroll

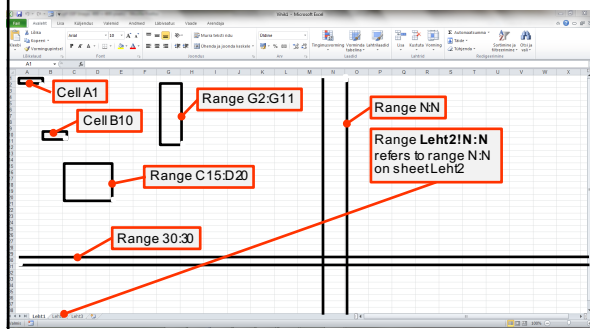
This is useful for keeping column and row labels etc. visible as you scroll around the data.

To "freeze" the top row(s) and left row(s), select the cell where the scrolling starts. Go to View | Freeze Panes.

© Mart Murdvee 2000-2013

17

Excel: cell address, or "cell reference"



© Mart Murdvee 2000-2013

18

Excel: Simple calculations

- Cells containing numbers for calculations should not contain anything else!
- Click in the cell where you want the answer to appear.
- Insert an = and then enter numbers or click cells and insert operands:
 - + for addition [=2+3]
 - for subtraction [=2-1]
 - * for multiplication [=4*5]
 - / for division [=4/2]
- There is a quick way to add — highlight the data to be added up and click the AutoSum icon (Σ) and select cell range.

or use
equations
with cell
references
→

	A	B	C
1	2	3	=A1+B1
2	2	1	=A2-B2
3	4	5	=A3*B3
4	4	2	=A4/B4

© Mart Murdvee 2000-2013

19

Relative and absolute references

- References to cells in Excel are by default "relative references". This means if you copy the cells somewhere else, any formulae in them will still work, but references are relative.

	A	B	C	D
1	2	3	=A1*B1	
2	5	6	=A2*B2	
3	6	7	=B3*C3	

copied cell

pasted cell

- Using \$ sign, the references can be „frozen“:
 - \$ before column letter „freezes“ column
 - \$ before row number „freezes“ row

	A	B	C	D
1	2	3	=A\$1*B1	
2	5	6	=A\$1*B2	
3	6	7	=A\$1*B3	

copied cell

pasted cell

© Mart Murdvee 2000-2013

20

Series

	A	B	C
1	1	Q1	esmaspäev
2	2	Q2	teisipäev
3	3	Q3	kolmapäev
4	4	Q4	neljapäev
5	5	Q5	reede
6	6	Q6	laupäev
7	7	Q7	pühapäev
8	8	Q8	esmaspäev
9	9	Q9	teisipäev
10	10	Q10	kolmapäev
11	11	Q11	neljapäev
12	12	Q12	reede
13	13	Q13	laupäev
14	14	Q14	pühapäev
15	15	Q15	esmaspäev
16	16	Q16	teisipäev
17	17	Q17	kolmapäev
18	18	Q18	neljapäev
19	19	Q19	reede
20	20	Q20	laupäev
21	21	Q21	pühapäev

Excel is good at expanding series automatically. Apart from simple numerical series (e.g. 1, 2, 3 ...), Excel knows the days of the week etc.

- Type the first two members of the series in a column or row, select the cells and drag the small black square at the bottom left corner of the selection to the right or down respectively.

© Mart Murdvee 2000-2013

21

Series

	A	B	C
1	1	Q1	esmaspäev
2	1	Q1	esmaspäev
3	1	Q1	esmaspäev
4	1	Q1	esmaspäev
5	1	Q1	esmaspäev
6	1	Q1	esmaspäev
7	1	Q1	esmaspäev
8	1	Q1	esmaspäev
9	1	Q1	esmaspäev
10	1	Q1	esmaspäev
11	1	Q1	esmaspäev
12	1	Q1	esmaspäev
13	1	Q1	esmaspäev
14	1	Q1	esmaspäev
15	1	Q1	esmaspäev
16	1	Q1	esmaspäev
17	1	Q1	esmaspäev
18	1	Q1	esmaspäev
19	1	Q1	esmaspäev
20	1	Q1	esmaspäev
21	1	Q1	esmaspäev

If You does not like series:

- Copy the cell(s) (CTR+C)
- Select range
- Paste cells (CTR+V)

© Mart Murdvee 2000-2013

22

Keep Your research (literature, data, analysis) in order and meaningful!

- Folder system and folder names
- File names
- Sheet names
- Parameter / variable names
- Apply a system in names

Record important manipulations and changes!

© Mart Murdvee 2000-2013

23

	A	B	C	D	E	F
1	No	D1: Gender (1 - female)	D2: Age (years)	D3: Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	D4: Marital status (1 - single; 2 - married; 3 - cohabiting; 4 - separated; 5 - divorced; 6 - widow)	D5: How would you rate your quality of life? (1 - very poor; 2 - poor; 3 - neither poor nor good; 4 - good; 5 - very good)
2	1	2	31	3	6	2
3	2	2	75	2	6	3
4	3	2	87	2	6	4
5	4	2	88	4	6	4
6	5	2	71	2	5	2
7	6	2	74	3	6	2
8	7	2	65	3	6	2
9	8	1	39	3	1	3
10	9	2	79	4	1	3
11	10	2	77	4	1	3
12	11	2	59	3	1	3
13	12	2	86	3	6	5
14	13	2	81	2	6	1
15	14	2	77	2	1	3
16	15	2	73	3	1	2
17	16	1	84	2	6	1
18	17	2	79	3	1	3
19	18	2	65	4	1	3
20	19	2	70	2	1	2
21	20	2	85	3	6	3
22	21	2	85	2	6	3
23	22	2	82	2	6	3
24	23	2	78	3	6	3
25	24	2	92	6	6	1
26	25	2	77	4	6	3
27	26	2	77	4	6	3
28	27	2	69	4	6	3
29	28	2	69	4	6	3
30	29	2	81	4	6	3
31	30	2	67	4	6	3
32	31	2	85	3	6	3
33	32	1	68	3	5	3

Meaningful names of sheets / data sets

Practical recommendation for analysis:

- „Clean“ data
- Use different sheets.
- Use meaningful labels for variables, sheets, tables, and graphs.

© Mart Murdvee 2000-2013

24

A	B	C	D	E	F
No	Gender (1-male, 2-female)	Age (years)	Education (1-primary, 2-secondary, 3-tertiary)	Marital status (1-single, 2-married, 3-cohabiting, 4-separated, 5-divorced, 6-widow)	Verbal ability (101-109: very good, 110-119: good, 120-129: fair, 130-139: poor, 140-149: very poor)
1	2	81	3	6	2
2	1	72	3	3	3
3	1	2	2	6	2
4	2	2	2	6	2
5	2	2	2	6	2
6	2	2	2	6	2
7	2	74	1	1	1
8	1	72	2	2	5
9	1	39	1	1	1
10	1	71	2	2	1
11	2	74	4	1	3
12	1	50	3	1	3
13	2	86	3	6	6
14	3	2	1	2	6
15	2	74	3	6	6
16	1	73	1	1	2
17	1	84	2	6	1
18	1	79	3	1	1
19	1	65	4	1	3
20	1	73	2	1	2
21	2	80	3	6	5
22	2	80	3	6	5
23	2	80	2	6	5
24	2	78	3	6	1
25	2	74	2	6	1
26	2	79	2	2	2
27	2	79	2	2	2
28	2	77	4	6	5
29	2	77	4	6	5
30	2	67	3	6	1
31	3	67	3	6	1
32	1	65	3	6	1
33	1	65	3	5	3

When things go wrong:

Questionnaire

35. Kuidas käitute Teie, kui
Tellige oma lemmiklooma põlvkonda vastavate ja kirjutuse sse kättäntvõttelise sel olemala
nõjale rääkida. Kirjutage number iga variandi ette.

```
1= väga tihti  
2= tihti  
3= mõnikord  
4= harva  
5= mitte kunagi
```

_____ Vähem eaka teema tippu
_____ kiiresti, mis tal väga on
_____ et tee midagi
_____ sõltumata ajast hõlmamisi
_____ täpselt kindla-koostajakannet eka teasta ja võimalikke järeleid
_____ kausta kannatust ja kausta eka vaigistamiseks võimalist jõudu
_____ suure eaka inimene püstitatakse või teadlase rooli ning rõhkotat
_____ tema liikumisel arvestatakse või rahuldaga
_____ kaustana ohukera tun, see ekkal on figus vabardusa ja vihanne olla
_____ et tee midagi kindlat, sünnaga ega oskama
_____ pakun ekkade tegevust, et muudaks kindlamini, juhatamist
_____ kausta kannatust ja kiiret elat
_____ palendusel ekkade ega, veelis temaga, puudutan teda, olen tema
_____ lähedal
_____ , annan rõhutamist reaktiooni
_____ mõineme mõeld mulsi välja (mõeld)? _____

Data input

35. Kuidas käitute Teie, kui
eakas kahtub kõrvaleseerivalt?

1,2,5,3,2,5,3,4,4,1,2,2

Remember:

**ONE CELL –
ONE NUMBER**

© Matt McDevree 2018-2019

26

Excel: Elements of formulas

=FUNCTION_NAME(cell;range;number;"text";function)

- formula begins with equal sign =
- followed by FUNCTION_NAME
- formula operators are in brackets
- formula operators are separated by semicolon sign ;
- formula operators:
 - + ; - ; * ; / ; ^ ; = ; > ; < ; >= ; <= ; < > ;
 - cell „address“
 - range of cells (cell:cell)
 - number
 - „text“ (between quotation marks)
 - function operators
 - nesting functions

© Matt Mordvee 2010

Arithmetic operators

To perform basic mathematical operations such as addition, subtraction, or multiplication; combine numbers; and produce numeric results, use the following arithmetic operators.

Arithmetic operator	Meaning (Example)
+ (plus sign)	Addition (3+3)
- (minus sign)	Subtraction (3-1)
* (asterisk)	Multiplication (3*3)
/ (forward slash)	Division (3/3)
% (percent sign)	Percent (20%)
^ (caret)	Exponentiation (3^2)

© Matt Murdree 2020-2021

28

Comparison operators

You can compare two values with the following operators. When two values are compared by using these operators, the result is a logical value either TRUE or FALSE.

Comparison operator	Meaning (Example)
= (equal sign)	Equal to ($A1=B1$)
> (greater than sign)	Greater than ($A1>B1$)
< (less than sign)	Less than ($A1<B1$)
>= (greater than or equal to sign)	Greater than or equal to ($A1\geq B1$)
<= (less than or equal to sign)	Less than or equal to ($A1\leq B1$)
<> (not equal to sign)	Not equal to ($A1\neq B1$)

© Matt McDevine 2002-2016

29

Text concatenation operator

Use the ampersand (&) to join, or concatenate, one or more text strings to produce a single piece of text.

Text operator	Meaning (Example)
& (ampersand)	Connects or concatenates two values to produce one continuous text value ("North"&"wind")

Reference operators

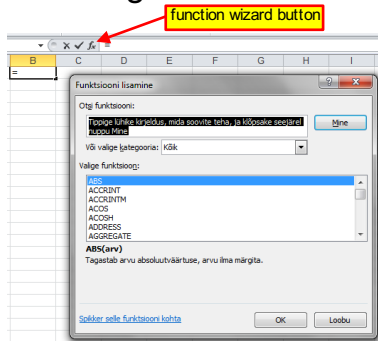
Combine ranges of cells for calculations with the following operators.


Reference operator	Meaning (Example)
: (colon)	Range operator, which produces one reference to all the cells between two references, including the two references (B5:B15)
; (semicolon)	Union operator, which combines multiple references into one reference (=SUM(B5:B15;D5:D15)) Separator between operators in formulas (=COUNTIF(B:B;2))

© Matt Mordvee 2010-2016

30

Entering formulas



1. Click on a cell where you will paste a function
2. Write formula or Click  function wizard button and use function wizard menu.

© Mart Murdvee 2010-2016

31

Text compilation

Operator &

or

Function

=CONCATENATE(cell;"text";...)

© Mart Murdvee 2010-2016

32



Making sense & Cleaning data

© Mart Murdvee 2010-2016

33

Common sources of error in databases are:

- missing data coded as "999"
- 'not applicable' or 'blank' coded as "0"
- typing errors on data entry (use Data Validation)
- column shift - data for one variable column was entered under the adjacent column
- fabricated data - 'made up' or contrived
- coding errors
- measurement and interview errors

Most errors will be detected using three procedures:

- descriptive statistics (MIN; MAX)
- scatterplots
- histograms

© Mart Murdvee 2010-2016

34

Coding missing data

2	2	2
3	3	3
1	1	1
3	3	3
3	3	3
4	4	4
5	5	5
3	3	3
3	3	3
2	2	2
2	2	2
1	1	1
3	3	3
3	3	3
3	3	3
4	4	4

- **Recommendation for use of „999“ and other „meaningless“ number for coding missing data IS FOR FOOLS!**
- Missing values in the database must be marked, and **not left blank**.
- Use "X", ":", etc. - signs, which does not represent anything calculable.
- **NB! Excel plotting and regression functions treat missing values (empty cells) as zeroes.**

© Mart Murdvee 2010-2016

35

Problem: missing data

- Missing data can be meaningful, while it contains information about:
 - questionnaire quality (eg poorly worded question).
 - respondents - some questions can be unpleasant, too sensitive, and therefore, less answered.
- **The person's data set can be considered usable, when more than 80% of the questions are answered.**
- Missing data (in order to obtain a sufficiently large sample for analysis) may be replaced with the parameter average, median, mode, or the parameter value can be derived from other parameters by using regression equation.

© Mart Murdvee 2010-2016

36

How many?

Function	Does
=COUNT(range)	Counts the number of cells containing numbers in a range
=COUNTA(range)	Counts the number of non-blank cells within a range
=COUNTBLANK(range)	Counts the number of blank cells within a range
=COUNTIF(range;value)	Counts the number of cells in range that are the same as value / criteria.
=COUNTIFS(criteria_range1; criteria1;[criteria_range2, criteria2];...)	Counts the number of times all criteria are met.

© Mart Murdvee 2010-2016

37

Logical operation IF:

Function:
=IF(logical_test;value_if_true;value_if_false)

Counting empty cells:

Function:
=COUNTBLANK(range)

Replace empty cells with mode of a range and let full cells with initial value:
=IF(COUNTBLANK(reference cell)=1;MODE(range);reference cell)

Counting "-"-marked cells:

Function:
=COUNTIF(range;"-")

Replace "-"-marked cells with mean of a range and let not "-"-marked cells with initial value:
=IF(COUNTIF(reference cell;"-")=1;AVERAGE(range)reference cell)

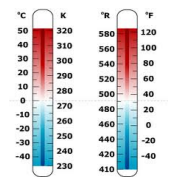
© Mart Murdvee 2010-2016

38

Scale



a gradation pattern of measurement (rule), which determines the numeric value of the result of measurement.



Recommendation:
Use standard and controlled scales!

© Mart Murdvee 2010-2016

39

Direction of scale

Confusing Scale!

COUNTRY	FREEDOM ON THE NET STATUS	FREEDOM ON THE NET TOTAL 0-100 Points
Estonia	Free	10
USA	Free	13
Germany	Free	16
Australia	Free	18
UK	Free	25
Italy	Free	26
South Africa	Free	26

Freedom House
FREEDOM ON THE NET 2011
A Global Assessment of Internet and Digital Media

© Mart Murdvee 2010-2016

It is desirable that the numerical values of the scale corresponds to the concept which is studied. A higher rate or more of something should correspond to the bigger numbers. Otherwise, the results are difficult to comprehend.

- For example, the concept of measuring the "exhaustion rate" with question, "I feel used up at the end of the workday," should have numeric scale direction:

1 - never, 2 - sometimes, 3 - always

NOT: 3 - never, 2 - sometimes, 1 - always

40

Inversion of the scale values:

Satisfaction (5 - completely dissatisfied, 4 - dissatisfied, 3 - do not know, 2 - satisfied, 1 - completely satisfied)	Satisfaction (1 - completely dissatisfied, 2 - dissatisfied, 3 - do not know, 4 - satisfied, 5 - completely satisfied)
initial value	inverted value
1	5
2	4
3	3
4	2
5	1

$$X_{\text{inverted}} = (\text{Min}_{\text{scale}} + \text{Max}_{\text{scale}}) - X$$

- when the numeric value of scale does not meet the direction of the concept meaning, or
- the value of the parameter must be included in to the factor with inverted value

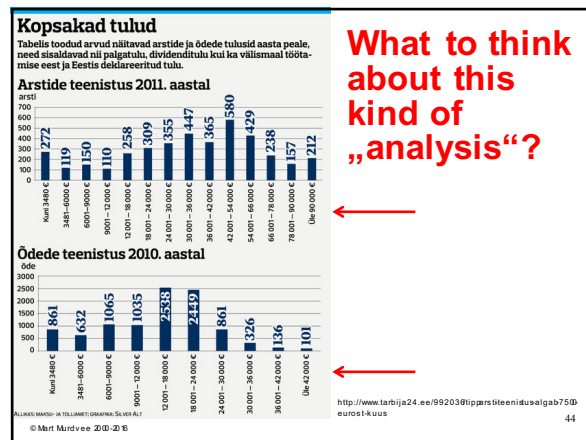
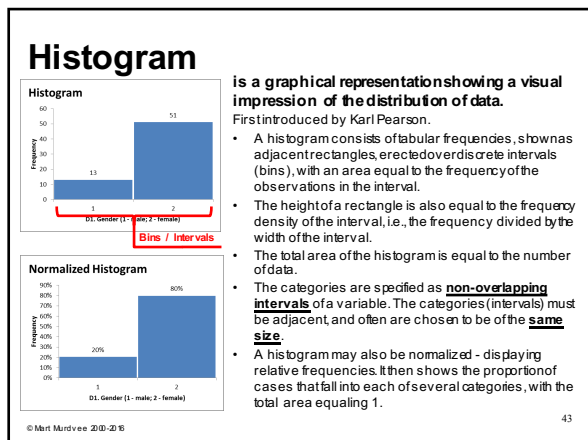
© Mart Murdvee 2010-2016

41

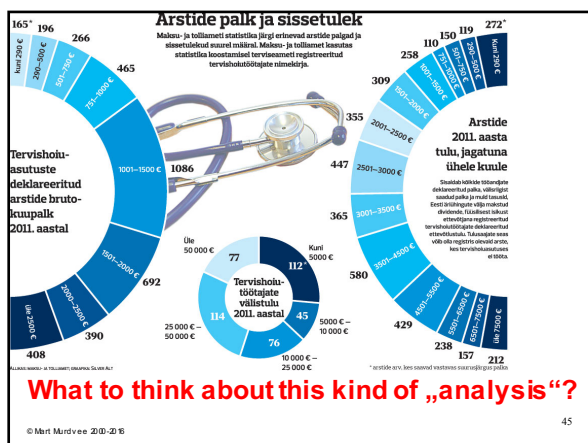


Distribution of values Histogram

© Mart Murdvee 2010-2016



What to think about this kind of „analysis“?

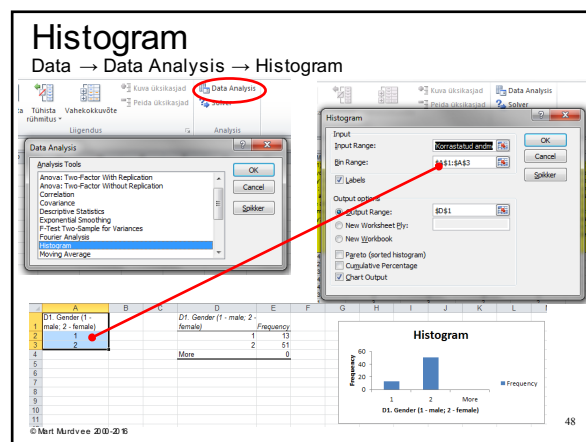
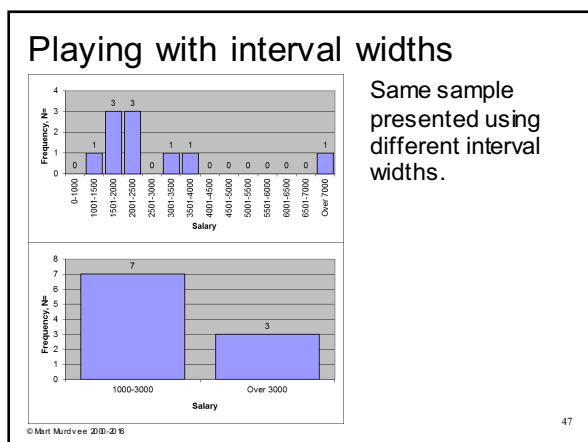


What to think about this kind of „analysis“?

Interval widths (bin widths), some "rules of thumb",

- Sturges's rule:** the number of intervals is close as possible to $1 + \log_2(N)$, where $\log_2(N)$ is the base 2 logarithm of the number of observations. The formula can also be written as $1 + 3.3 \log_{10}(N)$ where $\log_{10}(N)$ is the logarithm base 10 of the number of observations.
 $= 1 + \log_2(N)$ or $= 1 + 3.3 \log_{10}(N)$
- Rice rule:** which is to set the number of intervals to twice the cube root of the number of observations.
 $= 2 * \text{POWER}(N; 1/3)$

Basic rule:
Experiment with different choices of width, and to choose a histogram according to how well it communicates the shape of the distribution.



Histogram

Function:

=FREQUENCY(data_array;bins_array)

- Calculates how often values occur within a range of values, and then returns a vertical array of numbers.
- FREQUENCY is entered as an array formula (Shift+CTR+Enter) after you select a range of adjacent cells into which you want the returned distribution to appear.

	C2						
	A	B	C	D	E	F	G
1	Gender	Bin	N				
2	1 - male	1	13				
3	2 - female	2	51				

© Mart Murdvee 2000-2013

49

Histogram

Function:

=COUNTIF(range;criteria)

Counts the number of cells in range that are the same as value.

	TREND				
	A	B	C	D	E
1		D1. Gender (1 - male; 2 - female)			
2	Bin	Scale	N	%	
3	1 1 - male		=B3/A3	20.3%	
4	2 2 - female		51	79.7%	
5	TOTAL		64	100.0%	

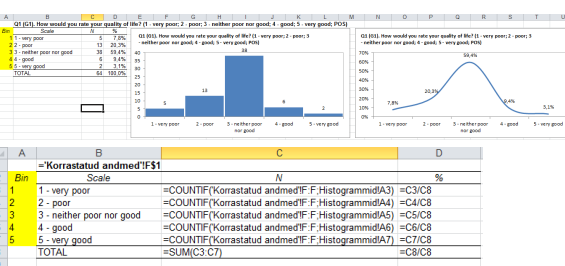
© Mart Murdvee 2000-2013

50

Histogram

Function:

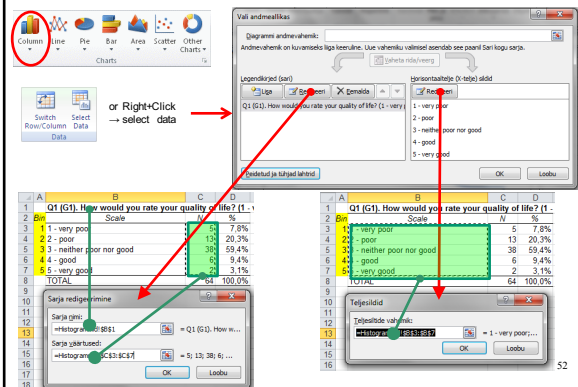
=COUNTIF(range;criteria)



© Mart Murdvee 2000-2013

51

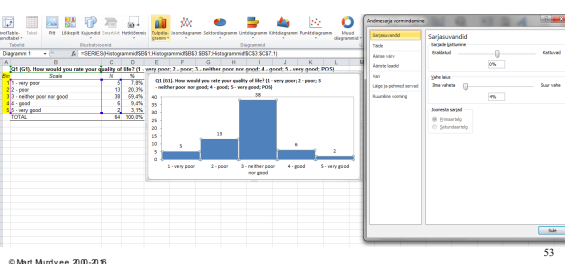
Chart



52

Designing Chart Elements

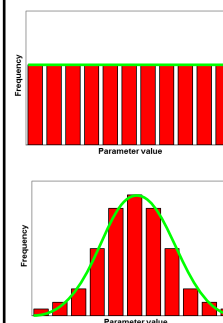
- Select chart element
- RightClick → Format ...
- Use menu



© Mart Murdvee 2000-2013

53

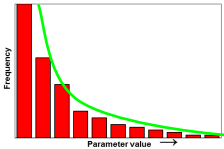
Interpreting histogram



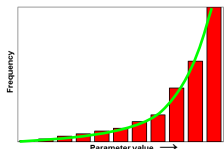
- Parameter values are uniformly distributed. Ideal random process. No trend.
- There is a central tendency. Normal distribution.

54

Interpreting histogram



- a rare phenomenon
- Maybe the range of the scale is too small?
- Logarithmic distribution?



- common phenomenon
- Maybe the range of the scale is too small?
- Square (power) distribution?

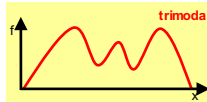
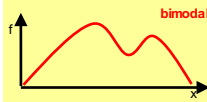
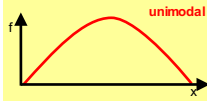
© Mart Murdvee 2000-2013

55

Distribution of parameter

Modality

modaalsus



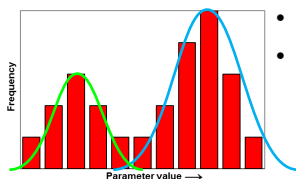
refers to the number of distinct peaks, or areas of cluster, that appear within a distribution, with each such peak being spoken of as a mode.

- A distribution with only one distinct peak is described as unimodal (*unimodaalne*);
- if it has two distinct peaks it is spoken of as bimodal (*bimodaalne*);
- three peaks, trimodal (*trimodaalne*);
- and so on.

© Mart Murdvee 2000-2013

56

Interpreting histogram

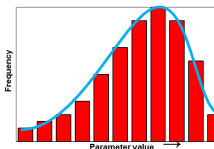


- Bimodal distribution.
- There can be two distinctive groups.

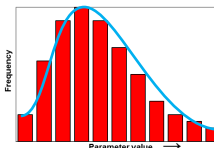
© Mart Murdvee 2000-2013

57

Interpreting histogram



- Distribution skewed toward greater values



- Distribution skewed toward smaller values

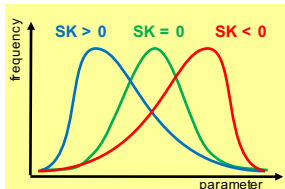
© Mart Murdvee 2000-2013

58

Distribution of parameter

Skewness

Asümeetriakordaja



- a measure of symmetry, or more precisely, the lack of symmetry of distribution:

- **SK > 0** – data that are skewed left, more smaller values;
- **SK = 0** – symmetrical;
- **SK < 0** – data that are skewed right, more greater values.

- =SKEW(range)

Pearson's 1. skewness coefficient:

$$SK = \frac{\bar{x} - M_o}{SD} \quad M_o - \text{mode}$$

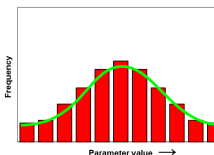
Pearson's 2. skewness coefficient:

$$SK = \frac{3(\bar{x} - M_d)}{SD} \quad M_d - \text{median}$$

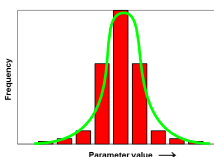
© M

59

Interpreting histogram



- „Flat“ distribution – lot of small and large values (heavy and long tails – (*jaotuse*)*sabad*). Weak central tendency

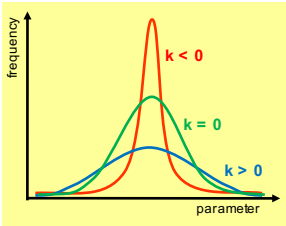


- „Peaked“ distribution – few small and large values (thin and short tails). Strong central tendency

© Mart Murdvee 2000-2013

60

Distribution of parameter
Kurtosis
järskuskordaja, ekstsess



the peakedness of a graph of the distribution compared to normal distribution:

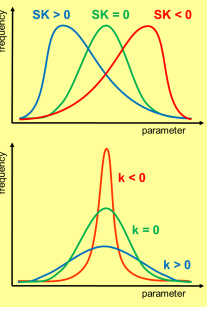
- $k > 0$ – extreme values, "long tails"
- $k = 0$ – normal distribution;
- $k < 0$ – few extreme values, "short tails"

• $=KURT(range)$

$$k = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2)^2} - 3$$

© Mart Murdvee 2000-2013 61

Rule of thumb about normality of variable



PROBLEM! Many of the statistical methods (Z test t-test, ANOVA, regression, covariance, correlation, chi-square test, F test) require the assumption that a variable or variables are normally distributed. Otherwise the statistical methods will produce totally invalid results. (QUESTIONABLE!)

A variable is reasonably close to normal if its skewness and kurtosis have values between -1,0 and +1,0.

- When a variable is not normally distributed, we can create a transformed variable and test it for normality. If the transformed variable is normally distributed, we can substitute it in our analysis.
- Three common transformations are: the logarithmic transformation, the square root transformation, and the inverse transformation.
- OR: use non-parametric tests

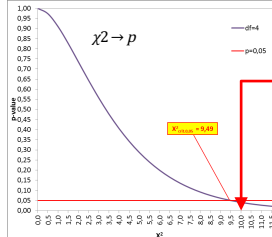
© Mart Murdvee 2000-2013 62

Chi-squared (χ^2) test
Hii-ruut test

tests a 0-hypothesis stating that the frequency distribution of certain events observed in a sample is consistent with a particular theoretical or other distribution.

NB! Chi-square procedures can be applied only if all expected values of are equal to or greater than 5.

Distributions	Observed	Expected
Scale	Group 1	Group 2
1 - not at all	22	6
2 - a little	17	13
3 - a moderate amount	8	25
4 - very much	0	12
5 - an extreme amount		



$\chi^2 = \sum_{i=1}^n \frac{(O_{observed} - E_{expected})^2}{E_{expected}}$

$df = (\text{number of categories}) - 1$

© Mart Murdvee 2000-2013 63

Chi-squared (χ^2) test
Hii-ruut test

Function:
=CHITEST(observed_range;expected_range)

returns the p-value from the chi-squared (χ^2) distribution for the statistic and the appropriate degrees of freedom.

Obs.	j = 1	j = 2	j = ...	j = c
i = 1	1 <i>ij</i>	1 <i>ij</i>	1 <i>ij</i>	1 <i>ij</i>
i = 2	2 <i>ij</i>	2 <i>ij</i>	2 <i>ij</i>	2 <i>ij</i>
i = ...	i <i>ij</i>	i <i>ij</i>	i <i>ij</i>	i <i>ij</i>
i = r	r <i>ij</i>	r <i>ij</i>	r <i>ij</i>	r <i>ij</i>

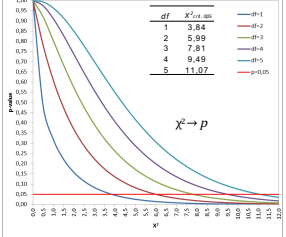
Exp.	j = 1	j = 2	j = ...	j = c
i = 1	1 <i>ij</i>	1 <i>ij</i>	1 <i>ij</i>	1 <i>ij</i>
i = 2	2 <i>ij</i>	2 <i>ij</i>	2 <i>ij</i>	2 <i>ij</i>
i = ...	i <i>ij</i>	i <i>ij</i>	i <i>ij</i>	i <i>ij</i>
i = r	r <i>ij</i>	r <i>ij</i>	r <i>ij</i>	r <i>ij</i>

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where:
O_{ij} = observed frequency in the i-th row, j-th column
E_{ij} = expected frequency in the i-th row, j-th column

$df = (r - 1)(c - 1)$

r = number of rows
c = number of columns



© Mart Murdvee 2000-2013 64