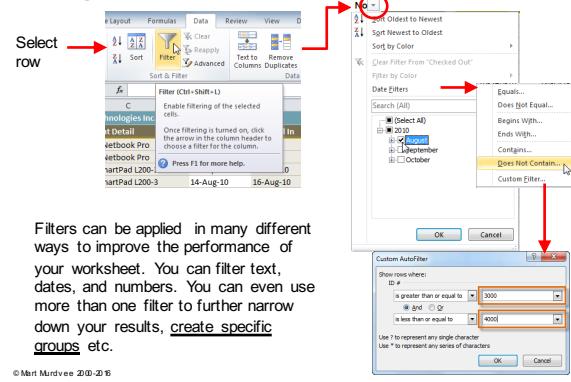




# Using filter

© Mart Murdvee 2000-2013

### Using filter



Filters can be applied in many different ways to improve the performance of your worksheet. You can filter text, dates, and numbers. You can even use more than one filter to further narrow down your results, create specific groups etc.

© Mart Murdvee 2000-2013

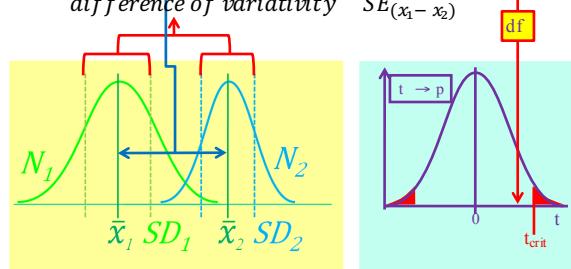
Are means of two sets of data really different?  
**Student's t-test**



The t statistic was introduced in 1908 by **William Sealy Gosset** (1876-1937), a statistician working for the Guinness brewery in Dublin, Ireland ("Student" was his pen name). Gosset had been hired due to Claude Guinness's innovative policy of recruiting the best graduates from Oxford and Cambridge to apply biochemistry and statistics to Guinness' industrial processes. Gosset devised the t-test as a way to cheaply monitor the quality of beer.

© Mart Murdvee 2000-2013

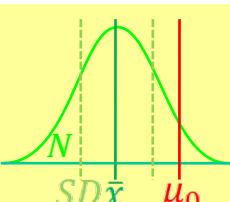
### t-statistik kahe keskmise võrdlemiseks



$$t = \frac{\text{difference of means}}{\text{difference of variativity}} = \frac{x_1 - x_2}{SE(x_1 - x_2)}$$

© Mart Murdvee 2000-2013

**One-sample t-test:**  
Student's t-test for comparison of mean width standard



Whether the sample mean differs from a standard  $\mu_0$ ?

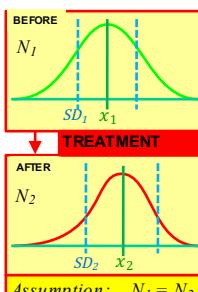
$$t = \frac{\bar{x} - \mu_0}{SD / \sqrt{N}}$$

$$df = N - 1$$

$t = ABS(X-M) / (SD / SQRT(N))$   
 $p = TDIST(t; (N-1)/2)$   
 $p = TDIST(ABS(X-M_0) / (SD / SQRT(N)); (N-1); 2)$

© Mart Murdvee 2000-2013

**Paired t-test:**  
**Compares (before and after) means of two paired groups**



Used to compare means on the same or related subject over time or in differing circumstances; subjects are often tested in a before-after situation.

Given two paired sets  $X_1$  and  $X_2$  of  $N$  measured values, the paired t-test determines if they differ from each other in a significant way.

$$t = \frac{\bar{x}_d}{SD_d / \sqrt{N}}$$

Where:

$$\bar{x}_d = x_1 - x_2$$

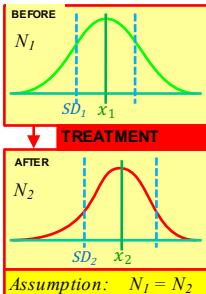
$$x_d = \frac{\sum x_1 - x_2}{N}$$

$$SD_d = \sqrt{\frac{(x_d - \bar{x}_d)^2}{N-1}}$$

© Mart Murdvee 2000-2013

Mart Murdvee:  
Research Methods and Data Analysis - exercise

**Paired t-test:**  
**Compares (before and after) means of two paired groups**



Function:  
 $=T.TEST(arrayX_1, arrayX_2, tails, type)$   
 $=TTEST(arrayX_1, arrayX_2, tails, type)$   
Arguments:  

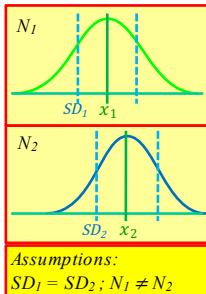
- Array1 - The first data set; X1.
- Array2 - The second data set; X2.
- Tails - Specifies the number of distribution tails
  - If = 1, one-tailed distribution is used.
  - If = 2, two-tailed distribution is used.
- **Type = 1 (Paired)**

Returns the probability (p) associated with a Student's t-Test.

© Mart Murdvee 2000-2013

7

**Two-sample unpaired t-test assuming equal variance (homoscedastic t-test)**



$$t = \frac{|x_1 - x_2|}{\sqrt{\frac{N_1 + N_2}{N_1 \times N_2} \times \frac{(N_1 - 1) \times SD_1^2 + (N_2 - 1) \times SD_2^2}{N_1 + N_2 - 2}}}$$

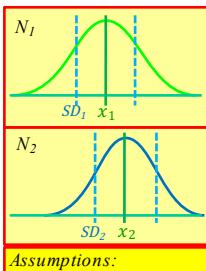
$$df = N_1 + N_2 - 2$$

Assumptions:  
 $SD_1 = SD_2 ; N_1 \neq N_2$

© Mart Murdvee 2000-2013

8

**Two-sample unpaired t-test assuming equal variance (homoscedastic t-test)**



Function:  
 $=T.TEST(arrayX_1, arrayX_2, tails, type)$   
 $=TTEST(arrayX_1, arrayX_2, tails, type)$   
Arguments:  

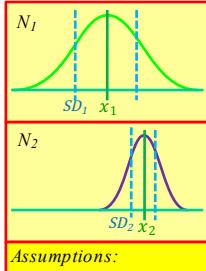
- Array1 - The first data set; X1.
- Array2 - The second data set; X2.
- Tails - Specifies the number of distribution tails
  - If = 1, one-tailed distribution is used.
  - If = 2, two-tailed distribution is used.
- **Type = 2 (Two-sample equal variance (homoscedastic))**

Returns the probability (p) associated with a Student's t-Test.

© Mart Murdvee 2000-2013

9

**Two-sample unpaired t-test assuming unequal variance (heteroscedastic)**



if the variance in the two groups are extremely different, e.g. the two samples are of very different sizes

$$t = \frac{|x_1 - x_2|}{\sqrt{\frac{SD_1^2}{N_1} + \frac{SD_2^2}{N_2}}}$$

$$df = \frac{\left( \frac{SD_1^2}{N_1} + \frac{SD_2^2}{N_2} \right)^2}{\frac{\left( SD_1^2 \right)^2}{N_1} + \frac{\left( SD_2^2 \right)^2}{N_2}} - 2$$

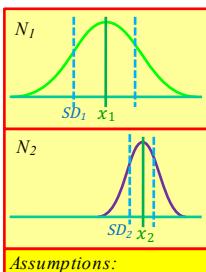
round up to integer

Assumptions:  
 $SD_1 \neq SD_2 ; N_1 \neq N_2$

© Mart Murdvee 2000-2013

10

**Two-sample unpaired t-test assuming unequal variance (heteroscedastic)**



Function:  
 $=T.TEST(arrayX_1, arrayX_2, tails, type)$   
 $=TTEST(arrayX_1, arrayX_2, tails, type)$   
Arguments:  

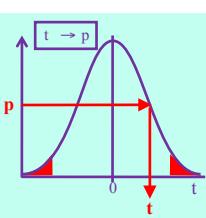
- Array1 - The first data set; X1.
- Array2 - The second data set; X2.
- Tails - Specifies the number of distribution tails
  - If = 1, one-tailed distribution is used.
  - If = 2, two-tailed distribution is used.
- **Type = 3 (Two-sample unequal variance (heteroscedastic))**

Returns the probability (p) associated with a Student's t-Test.

© Mart Murdvee 2000-2013

11

**Obtaining t-value from p-value**



Function:  
 $=T.INV(probability;deg\_freedom)$

returns the t-value of the Student's t-distribution as a function of the probability and the degrees of freedom

© Mart Murdvee 2000-2013

12

Mart Murdvee:  
Research Methods and Data Analysis - exercise

### t-test using Data Analysis Package

The screenshot shows the 'Data Analysis' dialog box with 't-Test: Two-Sample Assuming Unequal Variances' selected. Below it is the 't-Test: Two-Sample Assuming Unequal Variances' output table.

	x	y
Mean	14.77	12.94
Variance	6.67	6.30
Observations	63	63
Hypothesized Mean Difference	0.00	
t Stat	4.03	
P(T>=t) one-tail	0.00	
t Critical one-tail	1.66	
P(T<=t) two-tail	0.00	
t Critical two-tail	1.98	

© Mart Murdvee 2000-2013

13

### Effect size for difference of mean

The standardized effect size - Cohen's d:

Effect size is a standardized measure of the difference between two (or more) group means

$$d = \frac{x_1 - x_2}{SD_{1,2}}$$

$$SD_{1,2} = \sqrt{\frac{(n_1 - 1)SD_1^2 + (n_2 - 1)SD_2^2}{n_1 + n_2 - 2}}$$

Verbal Description	Effect Size (d)	Populations overlap
Small	0.20	85%
Medium	0.50	67%
Large	0.80	53%

© Mart Murdvee 2000-2013

14

### ANOVA - Analysis Of Variance

Dispersioonanalüüs

The screenshot shows the 'Data Analysis' dialog box with 'Anova: Single Factor' selected. Below it is the 'Anova: Single Factor' output table.

Groups	Count	Sum	Average	Variance
Grupp 1	10	20	2.0	0.67
Grupp 2	10	25	2.5	0.94
Grupp 3	10	39	3.9	1.43

Source of Variation	SS	df	MS	F	P-value	crit
Between Groups	19.4	2	9.70	9.56	0.00	<.35
Within Groups	27.4	27	1.01			
Total	46.8	29				

© Mart Murdvee 2000-2013

15

### Scatterplot

Hajuvusdiagramm

Relation of parameters

© Mart Murdvee 2000-2013

16

### Scatterplot

Hajuvusdiagramm

The scatterplot shows three data points plotted against x and y axes. The points are labeled (x1, y1), (x2, y2), and (x3, y3).

• Scatterplot graphs are used to show trends in data.  
• The data is displayed as a collection of points, each having the value of one variable determining the position on the horizontal axis (x) and the value of the other variable determining the position on the vertical axis (y).

© Mart Murdvee 2000-2013

17

### Scatterplot

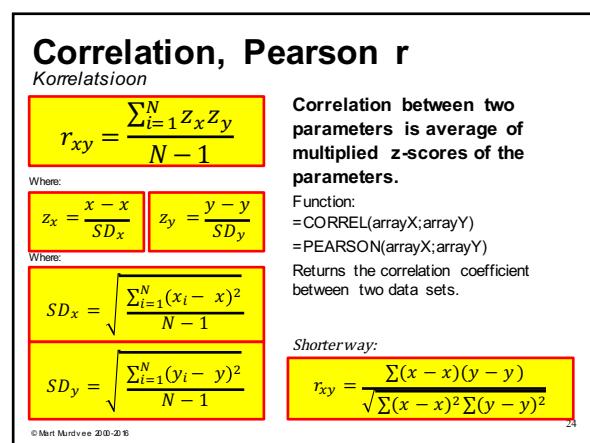
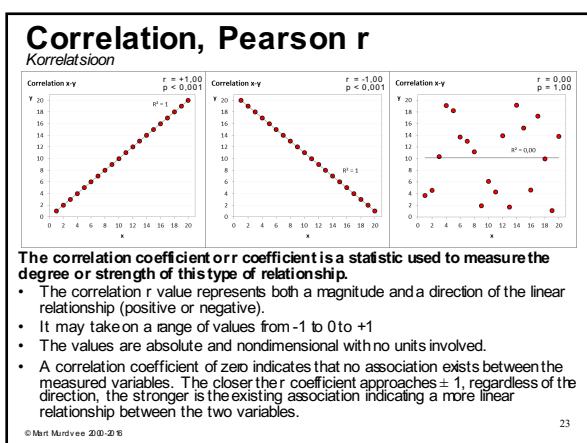
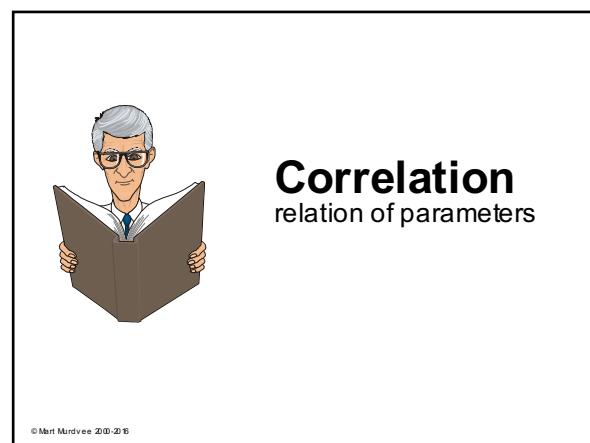
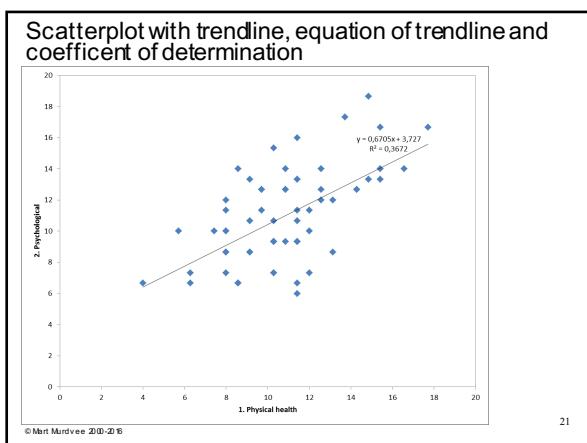
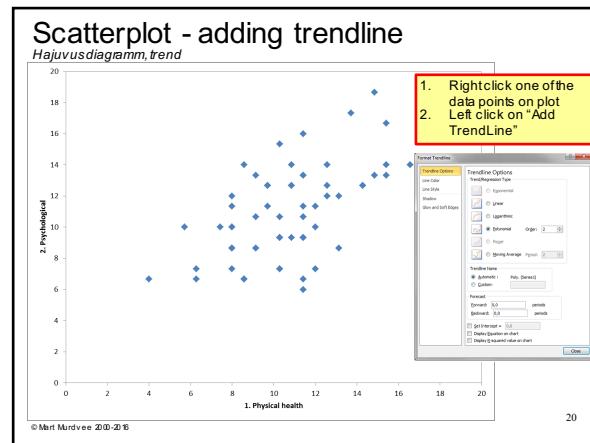
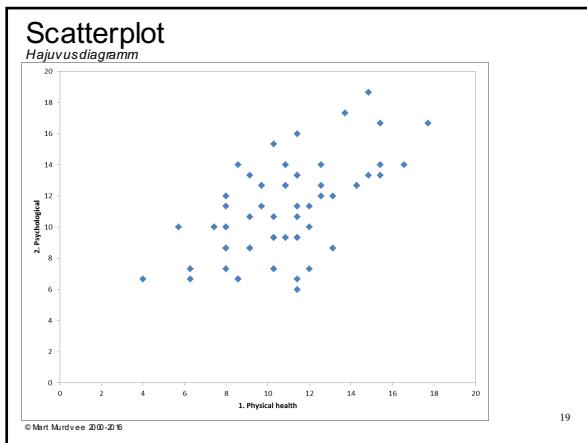
Hajuvusdiagramm

Select some data → Scatter → Scatter dialog box (with 'Series 1' and 'Series 2' fields) → 'Sarja redigeerimine' dialog box (with 'X-sarja väljatsead:' and 'Y-sarja väljatsead:' fields).

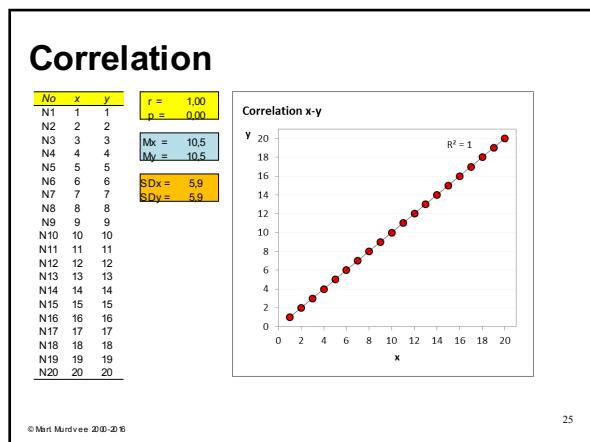
© Mart Murdvee 2000-2013

18

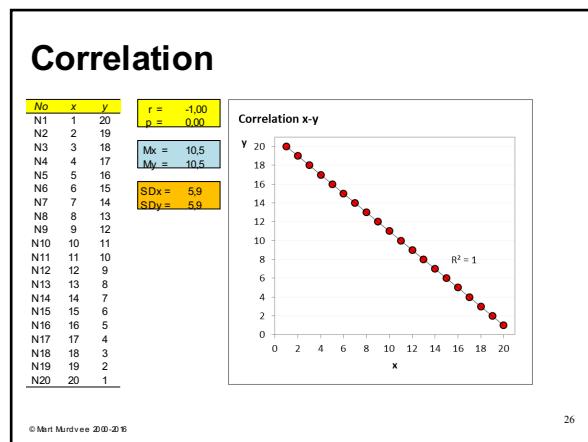
Mart Murdvee:  
**Research Methods and Data Analysis - exercise**



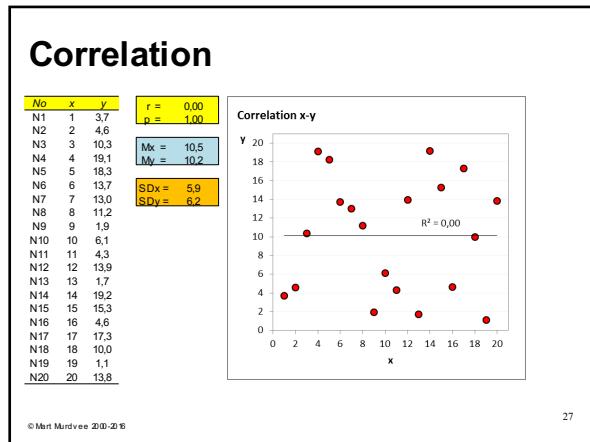
Mart Murdvee:  
Research Methods and Data Analysis - exercise



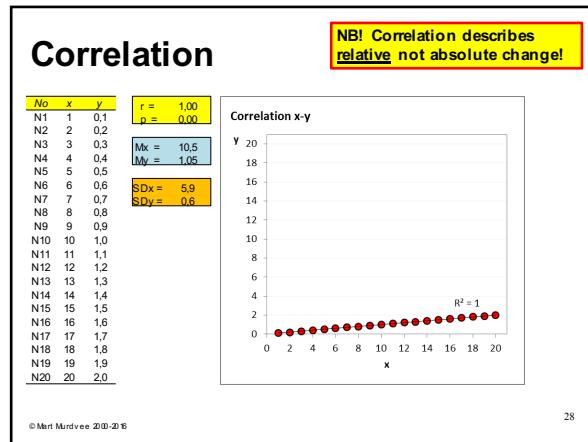
25



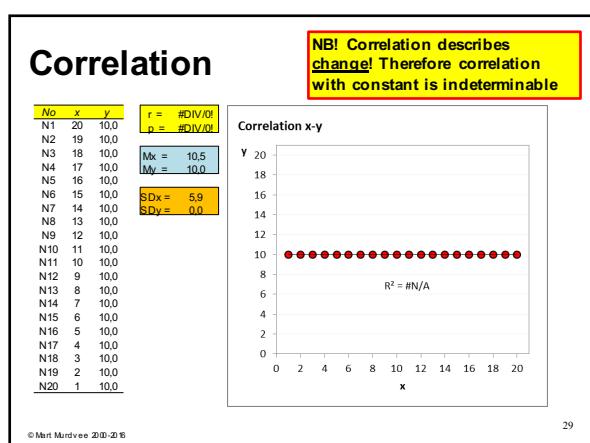
26



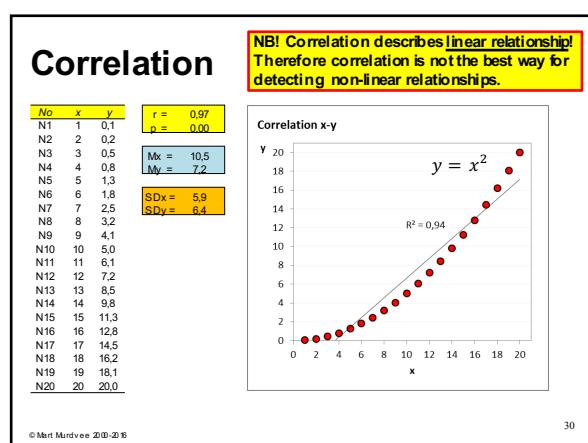
27



28



29



30

## Correlation

NB! Correlation describes linear relationship!  
Therefore correlation is not the best way for detecting non-linear relationships.

No	x	y
N1	1	0
N2	2	4
N3	3	8
N4	4	11
N5	5	13
N6	6	15
N7	7	17
N8	8	18.5
N9	9	19.5
N10	10	20
N11	11	20
N12	12	19.5
N13	13	18.5
N14	14	17
N15	15	15
N16	16	13
N17	17	11
N18	18	8
N19	19	4
N20	20	0

$r = 0.00$   
 $p = 1.00$   
 $Mx = 10.5$   
 $My = 12.8$   
 $SDx = 5.9$   
 $SDy = 6.4$

Correlation x-y

Recommendation: use scatterplot.

© Mart Murdvee 2000-2013

31

## Elements of correlation matrix

Parameter names

	a	b	c	d	e	f
a	$r_{aa}=1$	$r_{ab}$	$r_{ac}$	$r_{ad}$	$r_{ae}$	$r_{af}$
b	$r_{ba}$	$r_{bb}=1$	$r_{bc}$	$r_{bd}$	$r_{be}$	$r_{bf}$
c	$r_{ca}$	$r_{cb}$	$r_{cc}=1$	$r_{cd}$	$r_{ce}$	$r_{cf}$
d	$r_{da}$	$r_{db}$	$r_{dc}$	$r_{dd}=1$	$r_{de}$	$r_{df}$
e	$r_{ea}$	$r_{eb}$	$r_{ec}$	$r_{ed}$	$r_{ee}=1$	$r_{ef}$
f	$r_{fa}$	$r_{fb}$	$r_{fc}$	$r_{fd}$	$r_{fe}$	$r_{ff}=1$

Symmetrical ( $r_{cb} = r_{bc}$ ;  $r_{db} = r_{bd}$ ; etc.)

Correlation coefficient for correlation between parameters f and a

© Mart Murdvee 2000-2013

32

## Correlation

Data Analysis →

Analysis Tools

- Anova: Two-Factor Without Replication
- Anova: Two-Factor With Replication
- Covariance
- Descriptive Statistics
- Exponential Smoothing
- F-Test: Two Sample for Variances
- Fourier Analysis
- Histogram
- Moving Average
- Random Number Generation

Correlation

Input: Input Range: [ ]  
Grouped By: Columns  
Labels in First Row  
Output options: Output Range: [ ]  
New Worksheet By:  
New Workbook

© Mart Murdvee 2000-2013

33

## Correlation matrix

This part of table is intentionally left blank.

Read

© Mart Murdvee 2000-2013

34

## Statistical significance of correlations

N=	$r_{crit}$ ( $p=0.05$ )
5	0.8784
10	0.6319
50	0.2788
100	0.1966
1000	0.0620
2000	0.0439

$r_{crit}$  vs N

© Mart Murdvee 2000-2013

35

## Statistical significance of correlations - calculations

Using t-test:

$$t = \frac{|r|}{\sqrt{\frac{1-r^2}{N-2}}}$$

$$t = ABS(r)/SQRT((1-POWER(r;2))/(N2))$$

$$p = TDIST(t;(N-1)2)$$

$$p = TDIST(ABS(t)/SQRT((1-POWER(r;2))/(N-2));(N-1)2)$$

Using z-test:

$$z = r\sqrt{n-1}$$

$$z = r*SQRT(N-1)$$

$$p = NORM.S.DIST(z;FALSE)$$

$$p = NORM.S.DIST(r*SQRT(N-1);FALSE)$$

© Mart Murdvee 2000-2013

36

## Calculating t-value from r

	A	B	C	D
1	D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	
2	D1. Gender (1 - male; 2 - female)	#DIV/0!		
3	D2. Age (years)	3,71		
4	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	0,39		
5	D4. Marital status (1 - single; 2 - married; 3 - cohabiting; 4 - separated; 5 - divorced; 6 - widow)	2,35		

$$t = \frac{|r|}{\sqrt{\frac{1-r^2}{N-2}}}$$

37

## Calculating p-value from t

	A	B	C
1	D1. Gender (1 - male; 2 - female)	D2. Age (years)	
2	D1. Gender (1 - male; 2 - female)	#DIV/0!	0,00
3	D2. Age (years)	0,00	#DIV/0!
4	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	0,37	0,38

38

## Filtering significant correlations

	A	B	C
1	0,05	D1. Gender (1 - male; 2 - female)	D2. Age (years)
2	D1. Gender (1 - male; 2 - female)	#DIV/0!	0,42
3	D2. Age (years)	0,42	#DIV/0!
4	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	-	-
5	D4. Marital status (1 - single; 2 - married)	-	-

39

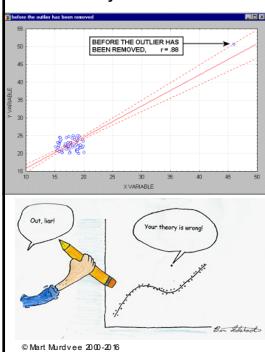
## Generating text

	A	B	C	D
1	D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	D4. Marital status (1 - single; 2 - married)
2	D1. Gender (1 - male; 2 - female)	#DIV/0!	i = 0,42 (p = 0,00)	i = 0
3	D2. Age (years)	i = 0,42 (p = 0,00)	#DIV/0!	i = 0
4	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	-	-	#DIV/0!
5	D4. Marital status (1 - single; 2 - married; 3 - cohabiting; 4 - separated; 5 - divorced; 6 - widow)	i = 0,28 (p = 0,03)	i = 0,32 (p = 0,01)	-

NB! The expression p = 0,00 must be later replaced with p < 0,01

40

## Outlier effect Eemalasuja efekt



An outlier is a value that is very different from the other data in data set.

Outliers have a profound influence on the slope of the regression line and consequently on the value of the correlation coefficient.

- Identify outliers by examining a scatterplot of each important correlation.
- Useful: exclude values that are outside the range of  $\pm 2$  standard deviations
- Question: who or what is outlier?
  - Error?
  - Real very different person or event?

41

## Strength of correlation

	Perfect	Perfekte	$\pm 1,0$
	Strong	Tugev	$\pm 0,8$
	Moderate	Mõõdukas	$\pm 0,6$
	Weak	Nõrk	$\pm 0,2$
	No relation	Seos puudub	$\pm 0,0$
	Or		
	Strong	+/- 0,7 ... +/- 1,0	
	Moderate	+/- 0,5 ... +/- 0,70	
	Weak	0 ... +/- 0,5	

In social sciences the correlation is considered meaningful if the absolute value of correlation is  $\geq 0,2$  (Disputable!)

42

## Importance of small statistical effects

It is important to avoid the error of assuming that small statistical effects necessarily translate into small practical or public health effects. There are many circumstances in which statistically small effects have large practical consequences, especially when small effects accumulate over time and over large proportions of the relevant population.

Huesmann, Taylor 2005

© Mart Murdvee 2000-2013

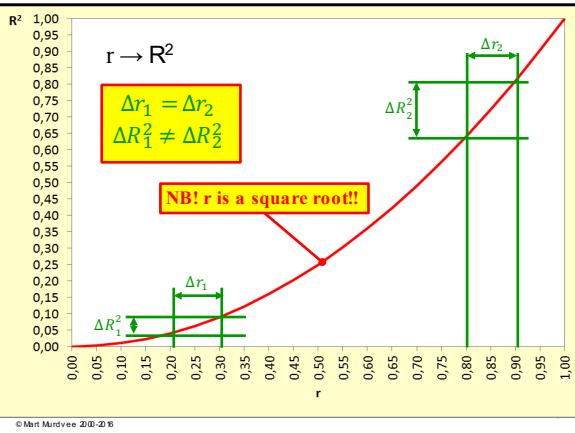
43



HELDUR MEERITS,  
investor ja LHV nõukogu liige

Heldur Meerits: kas raha teeb riigikogu? Postimees, 16.12.2012  
<http://arvamus.postimees.ee/1076188/heldur-meerits-kas-raha-teeb-riigikogu>  
© Mart Murdvee 2000-2013

44



## Effect size

$R^2$	Character of size of effect
$0,01 < R^2 < 0,09$	Small
$0,09 < R^2 < 0,25$	Moderate
$0,25 \leq R^2$	Strong

Cohen, 1988

© Mart Murdvee 2000-2013

47

## Coefficient of Determination

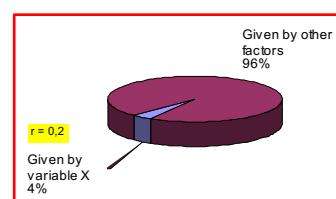
Determinatsioonikordaja

$r$	$R^2$
1	100,0%
0,9	81,0%
0,8	64,0%
0,7	49,0%
0,6	36,0%
0,5	25,0%
0,4	16,0%
0,3	9,0%
0,2	4,0%
0,1	1,0%

=RSQ(arrayX,arrayY)  
returns the square of the Pearson product moment correlation coefficient

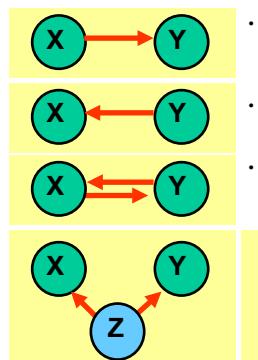
© Mart Murdvee 2000-2013

is the square of the correlation between two variables ( $R^2$ ). It expresses the amount of common variation between the two variables. Special case of effect size.



48

## Correlation does not imply causation!



- Establishing a correlation between two variables is not a sufficient condition to establish a causal relationship (in either direction).
- Correlation analysis measures a relationship or association; it does not define the explanation or its basis.
- Causes and explanations of the relationships must be found using other methods.

## Correlation does ...



Recent surveys shows that 100% of people who drinks water, dies.



That's a fact.



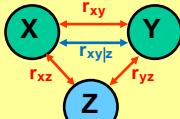
- It is proven that the celebration of birthdays is healthy. Statistics show that those people who celebrate the most birthdays become the oldest.
- Children with bigger feet spell better.
- Height of humans is negatively correlated with hair length.
- IMPORTANT WARNING** for those who have been drawn unsuspectingly into the use of bread:
  - More than 98 percent of convicted felons are bread users.
  - More than 90 percent of violent crimes are committed within 24 hours of eating bread.
  - Fully HALF of all children who grow up in bread-consuming households score below average on standardized tests.

© Mart Murdvee 2000-2016

49

## Partial correlation

Osaomrelatsioon



$$r_{xy|z} = \frac{r_{xy} - (r_{xz} \times r_{yz})}{\sqrt{(1 - r_{xz}^2) \times (1 - r_{yz}^2)}}$$

- Partial correlation is a method used to describe the relationship between two variables (x and y) whilst taking away the effects of another variable (z).
- Possible to find:
  - correlations due to the influence of third variables;
  - hidden correlations, ie, relationships that are masked by the impact of third variables.
- If  $r_{xy} \neq r_{xy|z}$  and the difference of correlation is statistically significant, then the parameter Z is influencing the relationship of XY.

50

## Cronbach's alpha

Cronbach's alpha

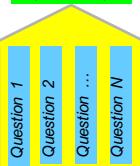


© Mart Murdvee 2000-2016

51

## Cronbach's alpha

CONSTRUCT  
(FACTOR)



is a coefficient of internal consistency of measures. It is commonly used as an estimate of the reliability of a psychometric test.

Cronbach's alpha indicates how well some features (questions) fit into a single construct:

- $\alpha = 1$ , practically same questions
- $\alpha = 0$ , totally incompatible questions

Cronbach's alpha	Internal consistency of construct
$\alpha \geq 0.9$	Excellent
$0.8 \leq \alpha < 0.9$	Good
$0.7 \leq \alpha < 0.8$	Acceptable
$0.6 \leq \alpha < 0.7$	Questionable
$0.5 \leq \alpha < 0.6$	Poor
$\alpha < 0.5$	Unacceptable

52

## Cronbach's alpha

Question 1 -  $s_1^2$   
Question 2 -  $s_2^2$   
Question ... -  $s_n^2$   
Question N -  $s_N^2$

CONSTRUCT  
(FACTOR)  
 $s_{sum}^2$

Cronbach's  $\alpha$

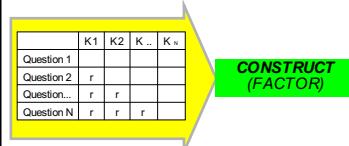
$$\alpha = \frac{N}{N-1} \times 1 - \frac{\sum_{i=0}^N s_i^2}{s_{sum}^2}$$

N - number of questions  
 $s^2$  - variance

© Mart Murdvee 2000-2016

53

## Standardized Cronbach's alpha



Standardized  $\alpha$

$$\alpha = \frac{N \times \bar{r}}{1 + (N-1) \times \bar{r}}$$

N - number of correlations  
 $\bar{r}$  - arithmetical mean of correlations

54



## Linear Regression

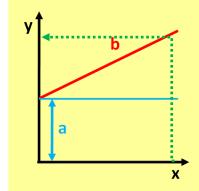
*lineaame regresioon*

relationship of parameters

© Mart Murdvee 2000-2013

### Linear regression analysis

*Lineaame regresioon*



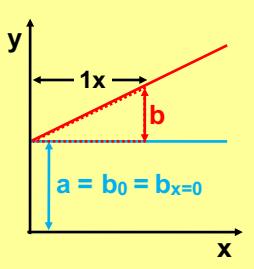
is an approach to modeling the relationship between a scalar dependent variable  $y$  and one or more explanatory variables denoted  $X$ . The case of one explanatory variable is called simple regression. More than one explanatory variable is multiple regression.

Equation:  
 $y = a + bx$   
 $x \rightarrow y$

The relationship between  $x$  and  $y$  is described using linear equation, although there are other options: square, polynomial, exponential, logarithm, etc equations.

© Mart Murdvee 2000-2013

### Linear regression equation

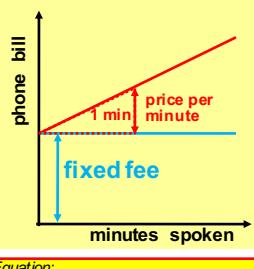


$y = a + bx$

$y$  – dependent variable;  
 $x$  – independent variable;  
 $a$  – constant, value of  $y$  when  $x = 0$ ;  
 $b$  – strength of relationship, when  $x$  changes 1 unit, changes  $y$  on amount of  $b$  (change of  $y$  -  $\Delta y$ ), slope of line.

© Mart Murdvee 2000-2013

### Example: phone bill



$y = a + bx$

where:

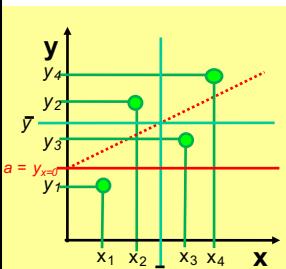
$y$  - the size of a phone bill;  
 $x$  - the number of minutes spoken;  
 $a$  - a fixed fee (minimum bill), phone bill, if the phone is not spoken at all;  
 $b$  - price per minute.

Equation:  
size of a phone bill = fixed fee + price per minute \* number of minutes spoken

© Mart Murdvee 2000-2013

### Least Squares Method

*Vähimruutude meetod*



Johann Carl Friedrich Gauss (1777-1855), A.D. 1794

- Solution:
$$b = \frac{\sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$a = \bar{y} - b\bar{x}$$
- EXCEL functions:
 $b = \text{SLOPE}(Y:Y;X:X)$  returns the slope of the linear regression line  
 $a = \text{INTERCEPT}(Y:Y;X:X)$  returns the intercept of the linear regression line

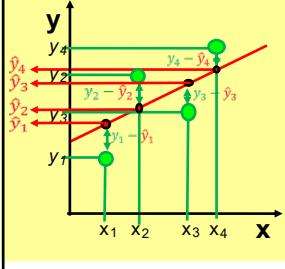
© Mart Murdvee 2000-2013

### How precise is the model? Standard error of regression

- Compute values of  $y$  from real  $x$ -values using regression equation:  
 $\hat{y}_i = a + bx_i$
- Compute differences of real  $y$  and computed  $\hat{y}$ :  
 $\epsilon_i = y_i - \hat{y}_i$
- Compute quadratic mean of differences = standard error of regression:  

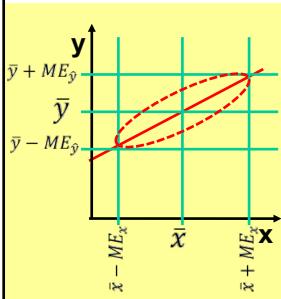
$$SE_{\hat{y}} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - 2}}$$

Excel:  $SE = \text{STEXY}(Y:Y;X:X)$  returns the standard error of the predicted  $y$ -value for each  $x$ -value in a regression



© Mart Murdvee 2000-2013

## Confidence interval of x and y



In what range are 95% or 99% of all values of x and y:

$$ME = Z_{crit} * SE$$

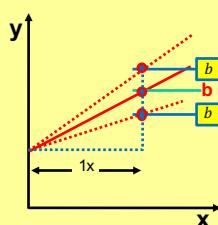
ME - margin of error

$$Z_{crit} 95\% = 1,96$$

$$Z_{crit} 99\% = 2,58$$

© Mart Murdvee 2000-2016

## Confidence interval of b



In which interval is b with probability of 95% or 99%:

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

or

$$SE_b = \frac{\sqrt{\sum_{i=1}^n (y_i - \hat{y}_i)^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} (n-2)}$$

$$ME_b = Z_{crit} * SE_b$$

where:

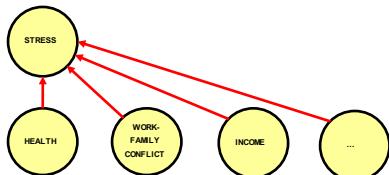
$$Z_{crit} 95\% = 1,96$$

$$Z_{crit} 99\% = 2,58$$

© Mart Murdvee 2000-2016

## Multiple Regression

Mitmene regressioon



- When parameter y is influenced by multiple parameters:  $x_1, x_2, x_3 \dots x_n$
- $y = a + b_1x_1 + b_2x_2 + b_3x_3 + \dots b_nx_n$

© Mart Murdvee 2000-2016

## Multiple Regression

Mitmene regressioon

Example: phone bill

$$y = a + b_1x_1 + b_2x_2$$

where:

- y - the size of a phone bill;  
 a - a fixed fee (minimum bill), phone bill, if the phone is not used at all;  
 $b_1$  - per-minute rate  
 $x_1$  - the number of minutes used;  
 $b_2$  - call setup fee  
 $x_2$  - number of calls

© Mart Murdvee 2000-2016

64

## Multiple Regression

Y	X1	X2	X3
1	2	4	2
3	2	5	4
2	3	3	2
4	2	4	1
1	2	2	1
2	3	2	1
2	3	2	1
2	3	5	1
1	4	5	4

SUMMARY OUTPUT

Regression Statistics

Multiple R

R Square

Adjusted R Square

Standard Error

Observations

Anova

df

SS

MS

F

Significance F

Regression

Residual

Total

Intercept

X1

X2

X3

Coefficients

Standard

Error

t Stat

P-value

Lower 95%

Upper 95%

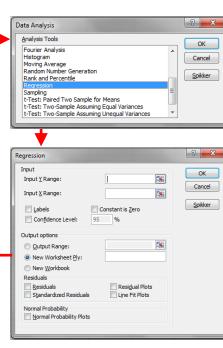
Intercept

X1

X2

X3

© Mart Murdvee 2000-2016



65

## Correlation and regression

The standardized slope of linear regression ( $b$ ) has the same value as the correlation coefficient ( $r$ ).

$$b = r_{xy} \frac{SD_y}{SD_x}$$

$$r_{xy} = \frac{b}{\frac{SD_y}{SD_x}}$$

$$a = \bar{y} - b\bar{x}$$

- If you are interested in simply characterizing the magnitude of the relationship between two variables, use correlation
- If you are interested in predicting or explaining your results in terms of particular values, use regression.

© Mart Murdvee 2000-2016

### Useful:

- Collection of statistical terms and definitions:
  - <http://www.stats.gla.ac.uk/steps/glossary/index.html>
- Internet sources that may be useful for various aspects of statistics:
  - <http://www.amstat.org/> (American Statistical Association),
  - <http://www.stat.ufl.edu> (University of Florida statistics department)
  - <http://www.statssoft.com/textbook/> (covers a wide range of topics, the emphasis is on techniques rather than concepts or mathematics)
  - <http://www.york.ac.uk/depts/maths/histstat/welcome.htm> (some information about the history of statistics)
  - <http://www.isid.ac.in/> (Indian Statistical Institute),
  - <http://www.isi-web.org/30-statsoc/statsoc/282-nsslist> (The International Statistical Institute)
  - <http://www.rss.org.uk> (The Royal Statistical Society),
  - <http://lib.stat.cmu.edu/> (an index of statistical software and routines)

© Mart Murdvee 2000-2013

67

### Recommendations:

When You start research  
– think on statistical methods!

Use different types of data analysis methods, charts, graphs – then the regularities are more easily detectable.

© Mart Murdvee 2000-2013

68