

Numbers are like people - torture them enough and they'll tell you anything.

Descriptive statistics

Descriptive statistics provides simple summaries about the sample and about the observations that have been made.

© Mart Murdvee 2000-2013

Sample size

Valimi maht

B2 =COUNT(Initial data!B:B)

A	B	C	D
No	D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)
		64	63
		63	57

Sample size (N) – how many elements (numbers) are in sample.
For example: how many answers are there.
Function:
 $=COUNT(range)$
 returns numbers in range.

© Mart Murdvee 2000-2013

Minimum

Miinimum

B2 =MIN(Initial data!C:C)

C	D
D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)
39	2

Function:
 $=MIN(range)$
 returns the minimum value of a data set

© Mart Murdvee 2000-2013

Maximum

Maksimum

B2 =MAX(Initial data!C:C)

C	D	
D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)
2	94	4

Function:
 $=MAX(range)$
 returns the maximum value of a data set

Recommendation:
 The minimum and maximum values should be calculated routinely – values of minimum and maximum show major typing errors on data entry!

© Mart Murdvee 2000-2013

Range

Haare

B2 =MAX(Initial data!C:C)-MIN(Initial data!C:C)

B	C	D
D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)
1	55	2

Function:
 $=MAX(range)-MIN(range)$
 returns range of values

© Mart Murdvee 2000-2013

Sum

Summa

B2 =SUM(Initial data!C:C)

C	D	
D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)
113	4692	150

Function:
 $=SUM(range)$
 adds a range of cells

$Sum = \sum x = x_1 + x_2 + \dots + x_n$

© Mart Murdvee 2000-2013

Descriptive statistics

Central tendency

© Mart Murdvee 2000-2013

Mean (arithmetic average) Aritmeetiline keskmene

=AVERAGE(Initial data!C:C)
D1. Gender (1 - male; 2 - female)
D2. Age (years)
D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)

1.77 24.48 2.77

$$\bar{x} = \frac{\sum x}{n}$$

© Mart Murdvee 2000-2013

- Sign: M, \bar{x}
Function: =AVERAGE(range)
Calculates the mean (arithmetic average) of a range of cells.
The arithmetic mean:
- Depends on the value of all observations.
 - Is simple to interpret.
 - Is the most familiar and the most used measure.
 - Is frequently used as an estimator of the mean of the population.
 - Has a value that can be falsified by the outliers.
- 8

Hunting statistician

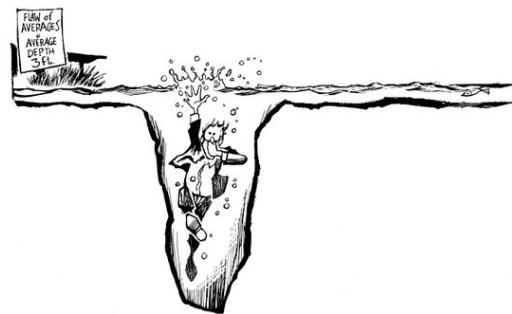
$\bar{x} = \frac{\sum x}{n}$

HIT!

! 9

© Mart Murdvee 2000-2013

Statistician drowned in a lake with an average depth of a half meter.



Mean temperature of patients in hospital

Patients	t^o
Patient 1	41.0
Patient 2	41.0
Patient 3	41.0
Patient 4	41.0
Patient 5	41.0
Patient 6	41.0
Patient 7	41.0
Patient 8	41.0
Patient 9	41.0
Patient 10	40.0
Patient 11	40.0
Patient 12	40.0
Patient 13	40.0
Patient 14	10.0
Patient 15	10.0
MEAN	36.6

11

© Mart Murdvee 2000-2013

Geometric mean Geomeetiline keskmene

$$\bar{x}_{geom} = \sqrt[n]{\prod_{i=1}^n x_i} = \left(\prod_{i=1}^n x_i\right)^{\frac{1}{n}}$$

=GEO_MEAN(range)

© Mart Murdvee 2000-2013

- n member values are multiplied and the nth-degree root is taken.
 - Used (in business and banking) for calculation of the average trends not random processes.
 - More conservative than the arithmetic average which is used for calculation of the probability of processes.
- 12

Harmonic mean

Hamooniline keskmine

$$\bar{x}_{har} = \frac{n}{\sum \frac{1}{n}}$$

=HARMEAN(range)

- Reciprocal value is used.
- Used for calculation of average speed, scalars, vectors in the parallel processes
- Tends strongly toward the least elements of the list, it may (compared to the arithmetic mean) mitigate the influence of large outliers and increase the influence of small values.

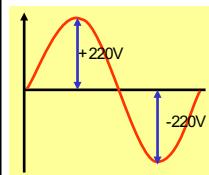
© Mart Murdvee 2000-2013

13

Quadratic mean

Ruutkeskmene

$$\bar{x}_{quad} = \sqrt{\frac{\sum x^2}{n}}$$



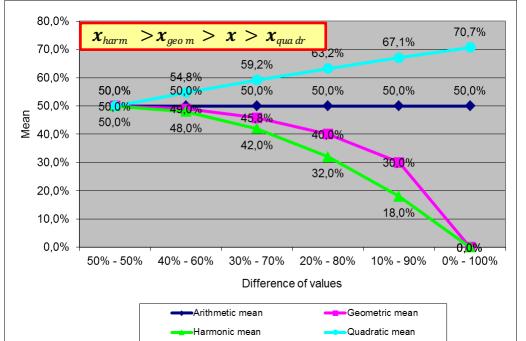
© Mart Murdvee 2000-2013

- Squares of the member's values are used in calculation.
- Used for the data sets where range of data includes data which is less than zero.

- Arithmetic mean = 0 V
- Quadratic mean = 220 V

14

Behavior of the means when there is difference in members values

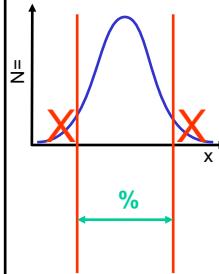


© Mart Murdvee 2000-2013

15

Truncated mean

Trimmitud keskmene



© Mart Murdvee 2000-2013

- Truncated mean is calculated from „survived“ values after extreme parameter values are removed.
- Truncated mean 50% means that mean is calculated from 50% of values, the 25% of high values and 25% low values removed.
- Is used if the parameter extremes can significantly affect the mean.
- In sports: the extreme ratings of judges are removed.

=TRIMMEAN(range;percent)

16

Weighted mean

Kaalutud keskmene

- Used in cases when members of data set have different characteristics that affect the average or occur at different frequencies
- May be calculated using arithmetic or geometric mean.

	Weight	Value
Sample 1	N	x
Sample 1	200	20
Sample 2	20	10
TOTAL / Weighted mean	220	19,09

$$\bar{x}_w = \frac{\sum w x}{\sum w}$$

© Mart Murdvee 2000-2013

17

Weighted mean

Kaalutud keskmene

Excel:

	Weight	Value
Sample 1	N	x
Sample 1	200	20
Sample 2	20	10
TOTAL / Weighted mean	220	19,09

$$\bar{x}_w = \frac{\sum w x}{\sum w}$$

=SUMPRODUCT(range x;range N)

=SUM(range N)

=SUMPRODUCT(range X;range N)/SUM(range N)

© Mart Murdvee 2000-2013

18

Mart Murdvee: Research Methods and Data Analysis - exercise

Mode

Mood

	B	C	D
D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	
2	78	3	

© Mart Murdvee 2000-2016

Function:
=MODE(range)
 returns the most frequently occurring, or repetitive, value in an range of data.

19

Median

Mediaan

	B	C	D
D1. Gender (1 - male; 2 - female)	D2. Age (years)	D3. Education (1 - elementary; 2 - primary; 3 - secondary; 4 - tertiary)	
2	78	3	

© Mart Murdvee 2000-2016

Function:
=MEDIAN(range)
 Returns the median of the given range.

The median is the number in the middle of a set of numbers; half the values in the data set are greater than the median and half are less than the median

20

Quantiles kvantiilid

No	Ordered x	Median	Quartiles	Deciles	Percentiles
		Min = 0.32			
1	0.32			1/10	
2	0.38			Dec1 = 0.83	
3	0.88		1/4	Dec2 = 2.07	Percentile = 1.53
4	1.64			Dec3 = 2.35	
5	2.18			1/10	
6	2.20			Dec4 = 2.49	
7	2.42		1/4	Dec5 = 2.57	
8	2.47			Dec6 = 3.09	
9	2.50			1/10	
10	2.89	Mdn = 3.09	Q1 = 2.20	Dec7 = 3.09	
11	3.29			1/10	
12	4.85			Dec8 = 5.57	
13	6.66		1/4	Dec9 = 6.98	
14	6.80			Dec10 = 7.82	
15	7.42		1/2	1/10	
16	7.56			Dec11 = 7.87	
17	8.87			1/10	
18	9.57		1/4	Dec12 = 9.57	
19	9.59			1/10	
20	9.82	Max = 9.82	Q4 = 9.82	Dec13 = 9.82	

© Mart Murdvee 2000-2016

are a set of 'cut points' that divide a sample of data into groups containing (as far as possible) equal numbers of observations.
 The most frequently used quantiles are:

- **Median**, which divides a distribution of observations into two parts.
- **Quartiles**, which separate a collection of observations into four parts.
- **Deciles**, which separate a collection of observations into ten parts,
- **Centiles**, which separate a collection of observations into a hundred parts.

Median

Mediaan

No	Ordered x	Median
		Min = 0.32
1	0.32	
2	0.38	
3	0.88	
4	1.64	
5	2.18	Q1 = 2.20
6	2.20	
7	2.42	
8	2.47	
9	2.50	
10	2.89	Mdn = 3.09
11	3.29	
12	4.85	
13	6.66	
14	6.80	
15	7.42	1/2
16	7.56	
17	8.87	
18	9.57	
19	9.59	
20	9.82	Max = 9.82

© Mart Murdvee 2000-2016

is the middle score of a set if the scores are organised from the smallest to the largest.

Function:
=MEDIAN(range)
 returns the median of the given range.

The median:

- Is easy to determine because only one data classification is needed.
- Is easy to understand but less used than the arithmetic mean.
- Is not influenced by outliers, which gives it an advantage over the arithmetic mean, if the series really have outliers.

22

Quartiles kUARTIILID

No	Ordered x	Quartiles
		Q0 = 0.32
1	0.32	
2	0.38	
3	0.88	1/4
4	1.64	
5	2.18	Q1 = 2.20
6	2.20	
7	2.42	
8	2.47	1/4
9	2.50	
10	2.89	Q2 = 3.09
11	3.29	
12	4.85	
13	6.66	1/4
14	6.80	
15	7.42	Q3 = 7.45
16	7.56	
17	8.87	
18	9.57	1/4
19	9.59	
20	9.82	Q4 = 9.82

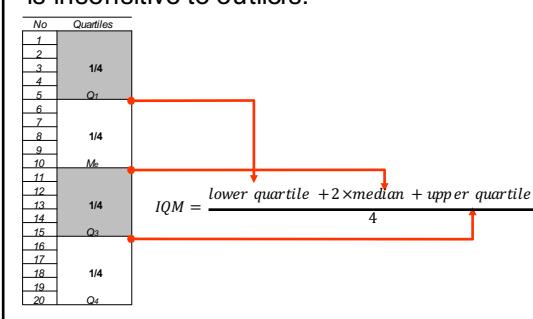
© Mart Murdvee 2000-2016

are values that divide a sample of data into four groups containing equal numbers of observations.
 $=QUARTILE(range;q)$
 $q = 0$ (minimum)
 $q = 1$ (lower quartile)
 $q = 2$ (median)
 $q = 3$ (upper quartile)
 $q = 4$ (maximum)

23

Interquartile mean (IQM), or midmean

is insensitive to outliers.



24

Deciles detsiiliid

No	Ordered x	Deciles
1	0.32	Dec1 = 0.32
2	0.38	1/10 Dec1 = 0.83
3	0.88	1/10
4	1.64	Dec2 = 2.07
5	2.18	1/10
6	2.20	Dec3 = 2.35
7	2.42	1/10
8	2.47	Dec4 = 2.49
9	2.53	1/10
10	2.89	Dec5 = 3.09
11	3.29	1/10
12	4.53	Dec6 = 5.57
13	4.68	1/10
14	6.80	Dec7 = 6.98
15	7.42	1/10
16	7.56	Dec8 = 7.82
17	8.87	1/10
18	9.57	Dec9 = 9.57
19	9.59	1/10
20	9.82	Dec10 = 9.82

are values that divide a sample of data into ten equal groups containing equal numbers of observations.

- The x^{th} decile indicates the value where $10x\%$ of the observations occur below this value and $(100 - 10x)\%$ of the observations occur above this value.
- For example, the eighth decile is the value where 80% of the observations fall below this and 20% occur above it.

25

Percentiles protsentiliid

No	Ordered x	Percentiles
1	0.32	
2	0.38	
3	0.88	Perce15% = 1.53
4	1.64	
5	2.18	
6	2.20	
7	2.42	
8	2.47	
9	2.50	
10	2.89	
11	3.29	
12	4.85	
13	6.66	
14	6.80	
15	7.42	
16	7.56	
17	8.87	
18	9.57	
19	9.59	
20	9.82	

are the values that divide a distribution into 100 equal parts (each part contains the same number of observations).

Function
`=PERCENTILE(range;k)`
 $0 < k < 1$
returns the k-th percentile of values in a range.

For example `=PERCENTILE(range;0,15)` gives value below which are 15% and over which are 85% of sample elements.

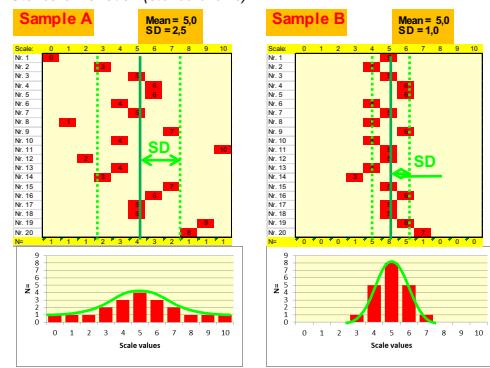
26



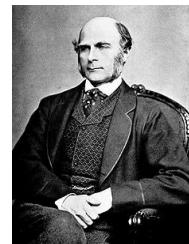
Descriptive statistics Variance, dispersion dispersioon, hajuvus

© Mart Murdvee 2000-2013

Variability of Values Standard Deviation (standardhälve)



29



Francis Galton

1822-1911

- Use of standard questionnaire
- Idea of Central Limit Theorem
- Concepts of standard deviation, regression and correlation



Karl Pearson

1857-1936

- Histogram
- Standard deviation
- Correlation coefficient (Pearsonr)
- χ^2 chi-square test
- Statistics as a general method that is applicable to all sciences

28

Variance and standard deviation Dispersioon ja standardhälve

`f1 = VAR(Korrasstatud andmed!B:B)`

D1: Gender (= male: 2 - female: 1)

`f2 = STDEV(Korrasstatud andmed!B:B)`

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

$$SD = \sqrt{s^2}$$

© Mart Murdvee 2000-2013

- Variance (s^2)** is a measure of the spread or dispersion within a set of data.

Function:
`=VAR(range)`

- Standard deviation (SD)** is a measure of how widely values are dispersed from the average value (the mean).

Function:
`=STDEV(range)`

30

Standard score, Z-score

standardiseeritud väärus, z-skoor

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

indicates by how many standard deviations a value of data set is above or below the mean.

- It is a dimensionless quantity derived by subtracting the population mean from an individual raw score and then dividing the difference by the population standard deviation.
- This conversion process is called standardizing or normalizing.
- Standard scores are also called z-values, z-scores, normal scores, and standardized variables.

Function:
`=STANDARDIZE(x;x;SD)`

© Mart Murdvee 2000-2013

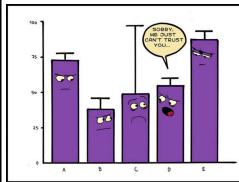
31

Coefficient of variation (CV)

(variatsioonikordaja)

or relative standard error (RSE)

$$CV \text{ or } RSE = \frac{SD}{\bar{x}}$$



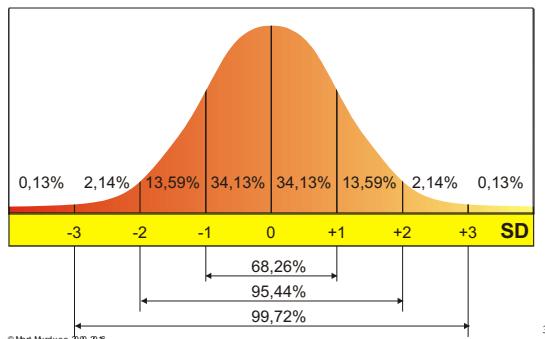
© Mart Murdvee 2000-2013

32

is a normalized measure of dispersion of a distribution and measures the relative dispersion of the studied variable.

- Distributions with $CV < 1$ are considered low-variance,
- Distributions with $CV > 1$ are considered high-variance.
- Relative standard error (RSE) is expressed as a percentage

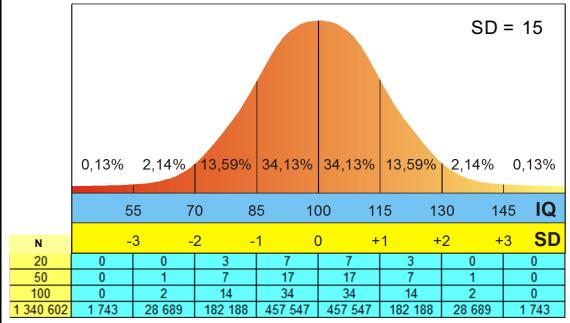
Standard deviation by normal distribution



© Mart Murdvee 2000-2013

33

Example: distribution of IQ



© Mart Murdvee 2000-2013

34

Standard error of the mean

Standardviga

$$SE = \frac{SD}{\sqrt{N}}$$

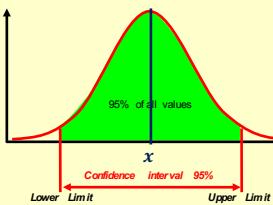
(SE, SEM) is the standard deviation of the values of a given function of the data (parameter) over all possible samples of the same size.

© Mart Murdvee 2000-2013

35

Confidence interval for mean

Usaldusvahemik



Confidence intervals for mean are intervals that will contain the population mean a specified proportion of the time, typically either 95% or 99% of the time. These intervals are referred to as 95% (CI95%) and 99% (CI99%) confidence intervals respectively.

x Function:

`=CONFIDENCE(alpha;SD;N)`

`=CONFIDENCE.NORM(alpha;SD;N)` returns the confidence interval value for a population mean using a normal distribution.

`=CONFIDENCE.T(alpha;SD;N)` Returns the confidence interval value for a population mean using a Student's t-distribution.

© Mart Murdvee 2000-2013

F-test

N₁

N₂

F-test is used to test if the variances of two populations are equal.

To remember:

$$F = \frac{S_1^2}{S_2^2}$$

$$SD = \sqrt{S^2}$$

Function:
=F.TEST(array1;array2)
 returns the result of an F-test, the two-tailed probability (p) that the variances in array1 and array2 are not significantly different.

Function:
=F.INV.RT(p;N1-1;N2-1)
 returns the inverse of the (right-tailed) F probability distribution

© Mart Murdvee 2000-2013

Probability and Statistical Significance

© Mart Murdvee 2000-2013

Sir Ronald Aylmer Fisher (1890-1962)

• ANalysis Of VAriance
 • Fisher's exact test
 • Fisher's z-distribution (F distribution)
 • Fisher transformation of r
 • Fisher's combined probability test – a meta-analysis technique
 • term "variance" in statistics
 • variance ratio (F-test)
 • phrase „test of significance“
 • principles of the design of experiments
 • Books:
 – Statistical Methods for Research Workers (1925)
 – Design of Experiments (1935)

© Mart Murdvee 2000-2013

k=2 = 0
k=1 = 1

X=1	SUM
0	0
1	1
TOTAL	1

Value	Count	Probability %
0	1	50%
1	2	50%
TOTAL	3	100%

k=1

© Mart Murdvee 2000-2013

k=2 = 0
k=1 = 1

X=2	SUM
0	0
1	1
0	1
1	2
TOTAL	3

Value	Count	Probability %
0	1	33%
1	2	67%
2	1	33%
TOTAL	3	100%

k=2

© Mart Murdvee 2000-2013

k=2 = 0
k=1 = 1

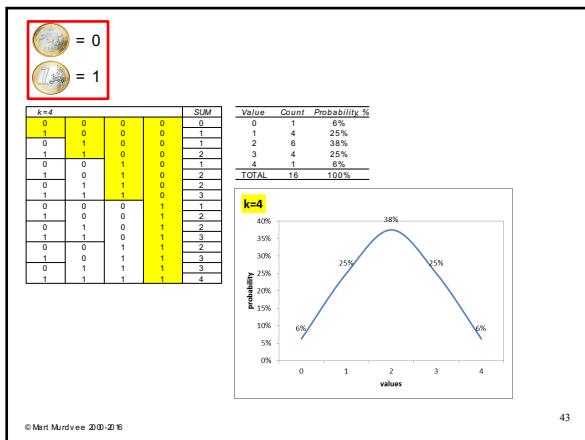
X=3	SUM
0	0
1	0
0	1
1	1
0	0
1	0
0	1
1	1
TOTAL	8

Value	Count	Probability %
0	1	13%
1	3	38%
2	3	38%
3	1	13%
TOTAL	8	100%

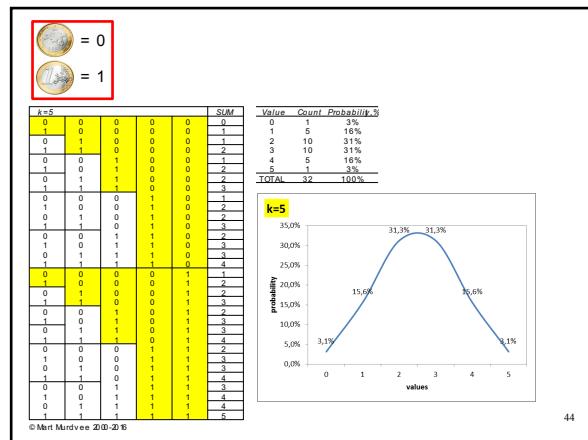
k=3

© Mart Murdvee 2000-2013

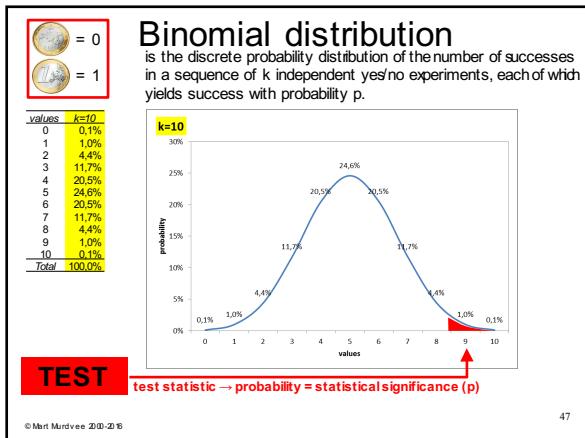
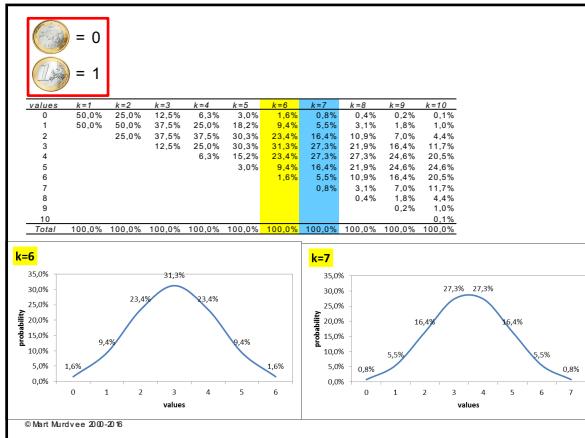
Mart Murdvee:
Research Methods and Data Analysis - exercise



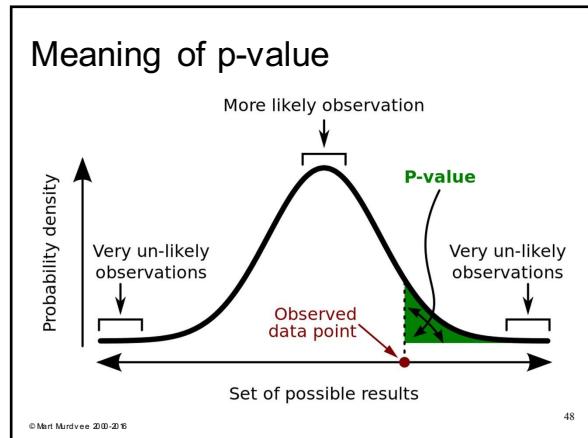
43



44

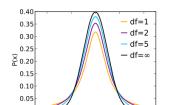


47



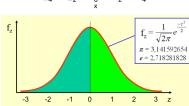
48

Distributions

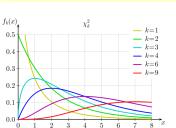


- Student's t-distribution

Normal distribution



- Chi-square Distribution
- etc



49

Confidence level, statistical significance

- Level of confidence (usaaldusnivo)

95% - in 95 times from 100 cases the regularity found in sample is present in population.

99% - in 99 times from 100 cases the regularity found in sample is present in population.

- Significance level - α (olulisuse nivo)

- Statistical significance, size of error, p-value (veasuurus, statistiline olulisus)

$p = 1 - \text{level of confidence}$
($p = 0,05 \dots p = 0,01 \dots p < 0,01$)

NB! The results which size of error (p) is equal to or less than 0,05 are conventionally declared as statistically significant.

50

Levels of statistical significance

p ≤ 0,05	borderline statistically significant	<i>piiripealselt statistiliselt oluline</i>
p ≤ 0,01	statistically significant	<i>statistiliselt oluline</i>
p ≤ 0,005	highly statistically significant	<i>kõrgelt statistiliselt oluline</i>
p ≤ 0,001		

© Mart Murdvee 2000-2013

51