

# STATISTICAL DATA ANALYSIS IN EXCEL

# Part 1

# Introduction to Statistics Descriptive Statistics

dr. Petr Nazarov

31-10-2011

petr.nazarov@crp-sante.lu

Statistical data analysis in Excel.

# **COURSE OVERVIEW**



**Objectives** 

# The course

- Reminds statistical basics
- Gives the methodological tools for the research
- Provides practical skill for fast data analysis

# Organization

 $\Rightarrow$  5 topics, 8-9 hours in total = 1 days

PLEASE: ask questions. Understanding is extremely important for later parts

# http://edu.sablab.net/sdae2011

Look for the data: http://edu.sablab.net/data/xls

# **COURSE OUTLINE**



# 1. Introduction

- Descriptive statistics
- Exploratory analysis
- Discrete probability distribution
- Continues probability distribution

# 2. Interval Estimations

- Sampling distribution
- Interval estimation for mean
- Interval estimation for proportion
- Sample size selection

# 3. Testing Hypotheses about Means

- Hypotheses
- Comparing of a mean and a constant
- Unpaired t-test
- Paired t-test

# 4. ANOVA

- 1-way ANOVA
- 2-way ANOVA

# 5. Linear Regression

- Simple linear regression
- Multiple linear regression

Look for the data: http://edu.sablab.net/data/xls

Statistical data analysis in Excel.



# Lecture 1. Reminding of the Basics ③

- descriptive statistics
- numerical measures



С

# **Frequency Distribution**

### **Frequency distribution**

A tabular summary of data showing the number (frequency) of items in each of several nonoverlapping classes.



### In MS Excel use the following functions:

- =COUNTIF(data,element) to get number of "elements" found in the "data" area
- =SUM(data) to get the sum of the values in the "data" area



# **Example: Pancreatitis Study**

The role of smoking in the etiology of pancreatitis has been recognized for many years. To provide estimates of the quantitative significance of these factors, a hospital-based study was carried out in eastern Massachusetts and Rhode Island between 1975 and 1979. **53 patients** who had a hospital discharge diagnosis of **pancreatitis** were included in this unmatched case-control study. The **control group** consisted of 217 patients admitted for **diseases other** than those of the pancreas and biliary tract. Risk factor information was obtained from a standardized interview with each subject, conducted by a trained interviewer.

adapted from Chap T. Le, Introductory Biostatistics

pancreatitis.xls

### Pancreatitis patients:

Smokers	Ex-smokers	Ex-smokers	Smokers	Smokers	Smokers		
Ex-smokers	Smokers	Smokers	Smokers	Smokers	Smokers		
Ex-smokers	Smokers	Smokers	Ex-smokers	Smokers	Smokers		
Ex-smokers	Ex-smokers	Smokers	Ex-smokers	Smokers			
Smokers	Never	Smokers	Ex-smokers	Ex-smokers			
Smokers	Ex-smokers	Smokers	Smokers	Ex-smokers			
Smokers	Smokers	Smokers	Smokers	Smokers			
Ex-smokers	Smokers	Smokers	Smokers	Smokers			
Smokers	Smokers	Smokers	Smokers	Smokers			
Smokers	Never	Smokers	Smokers	Smokers			

### Statistical data analysis in Excel.





# **Frequency Distribution**

### **Frequency distribution**

A tabular summary of data showing the number (frequency) of items in each of several nonoverlapping classes.

### In MS Excel use the following functions:

- =COUNTIF(data,element) to get number of "elements" found in the "data" area
- =SUM(data) to get the sum of the values in the "data" area

### pancreatitis.xls

### **Frequency distribution:**

	Smoking	Cases	Controls
	Never	2	56
≯	Ex-smokers	13	80
, 	Smokers	38	81
	Total	53	217

#### **Relative frequency distribution:**

Smoking	Cases	Controls
Never	0.038	0.258
Ex-smokers	0.245	0.369
Smokers	0.717	0.373
Total	1	1



# **TABULAR AND GRAPHICAL PRESENTATION**

# **Bar and Pie Charts**



### In MS Excel use the following steps:

- Chart Wizard  $\rightarrow$  Columns  $\rightarrow$  Set data range (both columns of Percent freq. distribution)
- Chart Wizard  $\rightarrow$  Pie  $\rightarrow$  Set data range (one columns of Percent freq. distribution)



# **TABULAR AND GRAPHICAL PRESENTATION**

### **Mice Data Series**



### Statistical data analysis in Excel.



# Histogram

The following are weights in grams for 970 mice:



Sorted weights show that the values are in the 10 – 49.6 grams. Let us divide the weight into the "bins"

	Weight,g	Frequency
	>=10	1
	→10-20	237
hins	20-30	417
<b>NITO</b>	30-40	124
	40-50	11
<u>_N</u>	/lore	0
Statistical data analys	sis in Exce	el. 1. Introduc

# **TABULAR AND GRAPHICAL PRESENTATION**



# Histogram



# In Excel use the following steps:

- Specify the column of bins (interval) upper-limits
- ◆ Tools → Data Analysis → Histrogram → select the input data, bins, and output (Analysis ToolPak should be installed)
- $\clubsuit$  use Chart Wizard  $\rightarrow$  Columns to visualize the results



# **Cumulative Frequency Distribution**

### **Cumulative frequency distribution**

A tabular summary of quantitative data showing the number of items with values less than or equal to the upper class limit of each class.





# **TABULAR AND GRAPHICAL PRESENTATION**

### **Scatter Plot**



Let us look on mutual dependency of the Starting and Ending weights.



### In Excel use the following steps:

Select the data region

• Use Chart Wizard  $\rightarrow$  XY (Scatter)



# **TABULAR AND GRAPHICAL PRESENTATION**

# **Crosstabulation**



### In Excel use the following steps:

- Set the range, including the headers of the data
- Select output and set layout by drag-and-dropping the names into the table



### **Population and Sample**



C C S ANT É CENTRE DE RECHERCHE PUBLIC

Statistical data analysis in Excel.



### **Measures of Location**

Mean A measure of central location computed by summing the data values and dividing by the number of observations.	Median A measure of central location provided by the value in the middle whe the data are arranged ascending order.	ModeA measure of location, defined as the value that occurs with greatestI in
$\frac{1}{x} = m = \frac{\sum x_i}{n}$ $\mu = \frac{\sum x_i}{N}$ $\mu = \frac{\sum (x_i = true)}{n}$	Weight 12 16 19 22 23 23 23 24 32 36 42 63	Mode = 23 Median = 23.5 Mean = 31.7



# **Measures of Location**





### Statistical data analysis in Excel.



### **Quantiles, Quartiles and Percentiles**





### **Measures of Variability**

Interquartile range (IQR) A measure of variability, defined to be the difference between the third and first quartiles.	Variance A measure of variability based on the squared deviations of the data values about the mean.					Standard deviation A measure of variability computed by taking the positive square root of the variance.						
$IQR = Q_3 - Q_1$ population $\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$ Sample samplesample $s^2 = \frac{\sum (x_i - \bar{x})^2}{N}$ Population standard deviation = $\sigma = \sqrt{\sigma}$						$=\sqrt{s^2}$ $=\sqrt{\sigma^2}$						
Weight 12	16 1	9 22	23	23	24	32	36	42	63	68		
<i>IQR</i> = 18	Variance = 320.2							St. d	ev. =	17.9		

### In Excel use the following functions:

=VAR(data), =STDEV(data)



### **Measures of Variability**

### **Coefficient of variation**

C

CENTRE DE RECHERCHE PUBLIC

A measure of relative variability computed by dividing the standard deviation by the mean.

Weight
12
16
19
22
23
23
24
32
36
42
63
68

(Standard deviation
$$\times 100$$
 %
 $\sim CV = 57\%$ 

Median absolute deviation (MAD) MAD is a robust measure of the variability of a univariate sample of quantitative data.

$$MAD = median(|x_i - median(x)|)$$

Set 1	Set 2			
23	23			
12	12			
22	22		Cat 4	Cot 0
12	12		Set	Set 2
21	21	Mean	17.3	22.2
18	81	Median	18	19
22	22			
20	20	St.dev.	4.23	18.18
12	12	ΜΛΠ	5 02	5.02
19	19		0.80	0.93
14	14			
13	13			
17	17			



#### **Skewness**

CENTRE DE RECHERCHE PUBLIC

A measure of the shape of a data distribution. Data skewed to the left result in negative skewness; a symmetric data distribution results in zero skewness; and data skewed to the right result in positive skewness.





adapted from Anderson et al Statistics for Business and Economics

# **Chebyshev's theorem** For any data set, at least $(1 - 1/z^2)$ of the data values must be within *z* standard deviations from the mean, where *z* – any value > 1.

### For ANY distribution:

standard deviations  $x_i$  is from the mean.

✤ At least 75 % of the values are within z = 2 standard deviations from the mean

A value computed by dividing the deviation about the mean  $(x_i - x)$  by the standard

deviation s. A **z-score** is referred to as a standardized value and denotes the number of

- $\clubsuit$  At least 89 % of the values are within z = 3 standard deviations from the mean
- $\clubsuit$  At least 94 % of the values are within z = 4 standard deviations from the mean
- At least 96% of the values are within z = 5 standard deviations from the mean





22



z-score

### Statistical data analysis in Excel.

# 1. Introduction

# NUMERICAL MEASURES

### **Detection of Outliers**

### For bell-shaped distributions:

For bell-shaped distributions data points with |z|>3 can be

considered as outliers.

**z-score** 0.04

-0.53

-0.01

-0.53

-0.06

-0.01

-0.11

-0.53

-0.17

-0.43

-0.48

-0.27

Weight

23 12

22

12

21

81

22

20

12 19

14

13

17

- Approximately 68 % of the values are within 1 st.dev. from mean
- Approximately 95 % of the values are within 2 st.dev. from mean
- Almost all data points are inside 3 st.dev. from mean

### Outlier

An unusually small or unusually large data value.

### **Example: Gaussian distribution**







# **Exploration Data Analysis**

### **Five-number summary**

An exploratory data analysis technique that uses five numbers to summarize the data: smallest value, first quartile, median, third quartile, and largest value



CENTRE DE RECHERCHE PUBL



#### In Excel use:

✤ Tool → Data Analysis → Descriptive Statistics



# **Example: Mice Weight**

### Example

Build a box plot for weights of male and female mice

1. Build 5 number summaries for males and females

	Female	Male
Min	10.0	12.0
Q1	17.2	23.8
Q2	20.7	27.1
Q3	23.3	31.2
Max	41.5	49.6

2. Combine the numbers into the following order

1. Introduction

open	Q3
high	Q3+min(1.5*(Q3-Q1),Max)
low	Q1-max(1.5*(Q3-Q1),Min)
close	Q1

### In Excel use:

- ♦ Chart Wizard → Stock → Open-high-low-close
- Put "series-in-rows"

Adjust colors, etc









**Covariance** 

### **Measure of Association between 2 Variables**

#### indicate a positive relationship; negative values indicate a negative relationship. population sample $\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{\sum (x_i - \mu_y)}$ S<sub>xy</sub> n-60 In Excel use function: mice.xls 50 =COVAR(data) 40 Ending weight 30 $s_{xy} = 39.8$ 20 Ending weight VS. 10 Starting weight hard to 0 10 20 30 40 50 0 interpret Starting weight

A measure of linear association between two variables. Positive values

### Statistical data analysis in Excel.

#### 1. Introduction

#### 26



### **Measure of Association between 2 Variables**

### **Correlation (Pearson product moment correlation coefficient)**

A measure of linear association between two variables that takes on values between -1 and +1. Values near +1 indicate a strong positive linear relationship, values near -1 indicate a strong negative linear relationship; and values near zero indicate the lack of a linear relationship.

### population





### sample

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{\sum (x_i - \overline{x})(y_i - \overline{y})}{s_x s_y (n-1)}$$

### In Excel use function:

=CORREL(data)

$$r_{xv} = 0.94$$

# mice.xls

Statistical data analysis in Excel.

# **Correlation Coefficient**





Wikipedia

If we have only 2 data points in *x* and y datasets, what values would you expect for correlation b/w

x and y?

### Statistical data analysis in Excel. 1. Introduction



# Discrete and continuous probability distributions

- discrete probability distribution
- continuous probability distribution
- normal probability distribution





Random variable A numerical description of the outcome of an experiment.

A random variable is always a numerical measure.

Roll a die

ENTRE DE RECHERCH



Discrete random variable

A random variable that may assume either a finite number of values or an infinite sequence of values.

**Continuous random variable** A random variable that may assume any numerical value in an interval or collection of intervals.

Number of calls to a reception per hour



Time between calls to a reception



Volume of a sample in a tube



Weight, height, blood pressure, etc



Statistical data analysis in Excel.





# **Discrete Probability Distribution**

### **Probability distribution**

A description of how the probabilities are distributed over the values of the random variable.

### **Probability function**

A function, denoted by f(x), that provides the probability that x assumes a particular value for a discrete random variable. Number of cells under microscope Random variable X:

. . .













### **Probability density function**

A function used to compute probabilities for a continuous random variable. The area under the graph of a probability density function over an interval represents probability.



### Statistical data analysis in Excel.



# **Normal Probability Distribution**

### Normal probability distribution

A continuous probability distribution. Its probability density function is bell shaped and determined by its mean  $\mu$  and standard deviation  $\sigma$ .



### In Excel use the function:

NORMDIST(x,m,s,false) for probability density function

• = NORMDIST(x,m,s,true) for cumulative probability function of normal distribution (area from left to x)

### Statistical data analysis in Excel. 1. Introduction



# **CONTINUOUS PROBABILITY DISTRIBUTIONS**

## **Standard Normal Probability Distribution**

**Standard normal probability distribution** A normal distribution with a mean of zero and a standard deviation of one.





### In Excel use the function:

 $\Rightarrow$  = NORMSDIST(z)



# **Dose Selection**

### Example

Assume that you have developed an extremely efficient chemical treatment for glioblastoma. During tests on animal models it was found that the substance X, which you use, is able to kill all tumor cells (theoretically), but being given at high concentration it leads to the death of a patient due to intoxication. As the survived cancer cells fast evolve into resistant form, the efficiency of the treatment is significantly reduced if the second course is given. Therefore the treatment should be performed in one injection.

The experimental data suggest that the average concentration needed for the positive treatment is 1  $\mu$ g/kg. The concentration needed for effective treatment is, of course, a random variable. Being presented in log10 scale and in g/kg, it can be approximated by a normal random variable with mean of –6 and standard deviation of 0.4.

The 50% lethal dose for human is 35  $\mu$ g/kg. And the tests on animals suggest that in log10 scale it has a normal distribution as well with the standard deviation of 0.3.

<b>parameter</b>	ug/kg	log scale	
mean positive treatment	1	-6	
std positive treatment	x	0.4	
mean lethal dose	35	-4.456	
std lethal dose	x	0.3	
atistical data analysis in Excel 1 Introdu	uction		









### In Excel use the function:

= NORMDIST(x,mean,std,FALSE)



# **CONTINUOUS PROBABILITY DISTRIBUTIONS**

# **Dose Selection**

- Probability to die from disease = inverse probability to treat
- Over-dose and disease behaviors are independent =>

$$P_{survive} = (1 - P_{lethal \ disease}) \cdot (1 - P_{lethal \ treatment})$$







# Thank you for your attention



to be continued...