

Microarray Center

APPLIED STATISTICS

Lecture 1

Data Presentation and Numerical Measures

Petr Nazarov petr.nazarov@crp-sante.lu

10-11-2009





The course:

- Enhance understanding of statistical basics
- Gives the methodological tools for the research
- Provides practical skill for fast data analysis





Organization

- 4 12 lectures organized in 6 sections (18 hours)
- 5 tests for self-control
- Lectures are integrated with practical work

Software

- Microsoft[®] Excel
 - with Data Analysis Add-In installed

Web page http://edu.sablab.net/stat



COURSE OVERVIEW

Recommended Literature

Recommended Literature



presentation methodology





examples, advanced topics







Lecture 1

Data and statistic

- elements, variables and observation
- + types of data (qualitative and quantitative) and scales (nominal, ordinal, interval, ratio)

Descriptive statistics: tabular and graphical presentation

- frequency distribution
- pie, bar chart and histogram representation
- cumulative distributions
- exploratory analysis: dot plot, stem-end-leaf
- crosstabulation and scatter diagram

Descriptive statistics: numerical measures

- measures of location: mean, mode, median, quantiles/quartiles/percentiles
- + measure of variability: variance, standard deviation, MAD, coefficient of variation
- other measures: skewness of distribution
- ✤ z-score. Chebyshev's theorem. Detection of outliers.

Exploratory analysis.

- ✤ 5 number summary
- box plot

Measure of association between two variables

- covariance and correlation coefficient
- interpretation of correlation coefficient

DATA AND STATISTICS



Data: Elements, Variables, and Observations



Can we consider the "Place" as element?

$$IFS = 3(\log_{10} N - 4.5)$$



DATA AND STATISTICS

Data Scales and Types



- Ex.1: Male, Female
- Ex.2: Rooms #: 101, 102, 103, ...
- Ex.1: Winners: The 1st, 2nd, 3rd places
- Ex.2: Marks: A, B, C, ...

- Ex.1: Examination score 0 -100
- Ex.2: Internet fame score ©
- Ex.1: Weight

Ex.2: Price



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 | |



TABULAR AND GRAPHICAL PRESENTATION

Frequency Distribution



SUM (data) to get the sum of the values in the "data" area



TABULAR AND GRAPHICAL PRESENTATION

Bar and Pie Charts



In MS Excel use the following steps:

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



Histogram

The following are weights in pounds of 57 children at a day-care center:



Since the smallest number is 12, we may begin our first interval at 10.

| | Weight Interval (lb) | Tally | Frequency | Relative Frequency (%) |
|-------|-------------------------|-------------|-----------|---------------------------|
| | 10–19 | +## | 5 | 8.8 |
| | 20-29 | -++++ -++++ | 19 | 33.3 |
| | 30-39 | | 10 | 17.5 |
| hine | 40-49 | | 4 | 7.0 |
| DITIS | 60–69 | | 4 | 7.0 |
| | 70–79 | | 2 | 3.5 |
| | Total | | 57 | 100.0 |

Chap T. Le, Introductory Biostatistics

TABULAR AND GRAPHICAL PRESENTATION



Histogram



- 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

TABULAR AND GRAPHICAL PRESENTATION



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.

| _ | | | | | | | | | | |
|---|----|----|----|----|----|----|----|----|----|----|
| | 68 | 63 | 42 | 27 | 30 | 36 | 28 | 32 | 79 | 27 |
| | 22 | 23 | 24 | 25 | 44 | 65 | 43 | 25 | 74 | 51 |
| | 36 | 42 | 28 | 31 | 28 | 25 | 45 | 12 | 57 | 51 |
| | 12 | 32 | 49 | 38 | 42 | 27 | 31 | 50 | 38 | 21 |
| | 16 | 24 | 69 | 47 | 23 | 22 | 43 | 27 | 49 | 28 |
| | 23 | 19 | 46 | 30 | 43 | 49 | 12 | | | |
| | | | | | | | | | | |



| Bins | Frequency | Cumulative FD | Relative CFD | Percent CFD |
|-------|-----------|---------------|---------------------|-------------|
| 10-20 | 5 | 5 | 0.09 | 9% |
| 20-30 | 21 | 26 | 0.46 | 46% |
| 30-40 | 8 | 34 | 0.60 | 60% |
| 40-50 | 14 | 48 | 0.84 | 84% |
| 50-60 | 3 | 51 | 0.89 | 89% |
| 60-70 | 4 | 55 | 0.96 | 96% |
| 70-80 | 2 | 57 | 1 | 100% |
| Total | 57 | | | |



Exploratory Data Analysis

Exploratory data analysis

Methods that use simple arithmetic and easy-to-draw graphs to summarize data quickly.





TABULAR AND GRAPHICAL PRESENTATION

Crosstabulation

| pancreatitis.xls | | | | | | | | | |
|------------------|--------|--|----|-----|--|--|--|--|--|
| Smo | king | <i>Disease</i> other pancreatitis Total | | | | | | | |
| Ex-s | mokers | 80 | 13 | 93 | | | | | |
| Neve | er | 56 | 2 | 58 | | | | | |
| Smo | kers | 81 | 38 | 119 | | | | | |
| Tota | | 217 | 53 | 270 | | | | | |

In Excel use the following steps:

- \clubsuit Data \rightarrow Pivot Table and PivotChart \rightarrow MS Office list + Pivot Table
- ✤ Set the range, including the headers of the data
- Select output and set layout by drag-and-dropping the names into the table



TABULAR AND GRAPHICAL PRESENTATION

Scatter Plot





Population and Sample



Lecture 1. Data presentation and numerical measures

910

CENTRE DE RECHERCHE PUBLI



Measures of Location

| Mean A measure of central location computed by summing the data values and dividing by the number of observations. | Median A meas location value in the dat ascend | n sure of centra n provided by n the middle v a are arrange ding order. | l the vhen ed in | Mode A measure of location, defined as the value that occurs with greatest frequency. |
|---|---|--|---------------------------|---|
| $\frac{1}{x} = \frac{\sum x_i}{n}$ | | Weight 12 16 19 | | |
| $\mu = \frac{\sum x_i}{N}$ | | 22 23 23 24 32 | | Mode = 23 Median = 23.5 |
| $p = \frac{\sum(x_i = true)}{n}$ | | 36 42 63 68 | | Mean = 31.7 |



Measures of Location



Lecture 1. Data presentation and numerical measures

Cr

CENTRE DE RECHERCHE PUBLIC



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. | | Variand A meas based o deviatio values a | e ure of y n the s ns of th about th | variabil squarec ne data he mea | ity d | | Si A cc pc va | tandar measu ompute ositive ariance | d devia are of v d by ta square | ation ariabil aking t root c | lity he of the | |
|---|-----------|--|--|--|---------------------|----------|---------------------------|---|--|---------------------------------------|-------------------------|-------------------------------------|
| $IQR = Q_3 - Q_1$ | poj sa | oulation o mple s | $\frac{1}{2} = \frac{\sum_{n=1}^{2}}{n}$ | $\frac{\sum (x_i - \mu)}{N}$ $\frac{(x_i - x)}{n-1}$ | $(\underline{u})^2$ | S Pop | 'ample vulation | r stand n stand | ard de lard de | eviatio eviatio | n = s = $n = \sigma$ | $=\sqrt{s^2}$ $=\sqrt{\sigma^2}$ |
| Weight 12 | 16 1 | 9 22 | 23 | 23 | 24 | 32 | 36 | 42 | 63 | 68 | | |
| <i>IQR</i> = 18 | <u> </u> | | | St. d | ev. = | 17.9 | _ | | | | | |

In Excel use the following functions:

=VAR(data), =STDEV(data)



Measures of Variability

Coefficient of variation

CI

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

$$\left(\frac{Standard\ deviation}{Mean} \times 100\right)\%$$
 $\checkmark CV = 57\%$

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 1 | Cat 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 | MAD | 5 02 | 5.02 |
| 19 | 19 | IVIAD | 5.85 | 5.95 |
| 14 | 14 | | | |
| 13 | 13 | | | |
| 17 | 17 | | | |



Skewness

01

0 05

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

0.15

01

0 05

z-score

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 standard deviations x_i is from the mean.

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:

- \Rightarrow At least 75 % of the values are within z = 2 standard deviations from the mean
- ✤ At least 89 % of the values are within z = 3 standard deviations from the mean
- ✤ At least 94 % of the values are within z = 4 standard deviations from the mean
- \Rightarrow At least 96% of the values are within z = 5 standard deviations from the mean









z-score

-0.27

NUMERICAL MEASURES

An unusually small or unusually

Detection of Outliers

For bell-shaped distributions:

Cr

CENTRE DE RECHERCHE

- ◆ 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

large data value.





data points with |z|>3 can be considered as outliers.

| Weight | z-score |
|--------|---------|
| 23 | 0.04 |
| 12 | -0.53 |
| 22 | -0.01 |
| 12 | -0.53 |
| 21 | -0.06 |
| 81 | 3.10 |
| 22 | -0.01 |
| 20 | -0.11 |
| 12 | -0.53 |
| 19 | -0.17 |
| 14 | -0.43 |
| 13 | -0.48 |

17





Exploration Data Analysis

Five-number summary

CENTRE DE RECHERCHE PUB

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



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

| open high | Q3 Q3+min(1.5*(Q3-Q1),Max) |
|--------------|-------------------------------|
| low | Q1-max(1.5*(Q3-Q1),Min) |
| 030 | |

In Excel use:

- \clubsuit Chart Wizard \rightarrow Stock \rightarrow Open-high-low-close
- Put "series-in-rows"
- Adjust colors, etc









Covariance

Measure of Association between 2 Variables





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.



Correlation Coefficient





Wikipedia

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

Lecture 1. Data presentation and numerical measures

NUMERICAL MEASURES

Weighted Mean

Weighted mean The mean obtained by assigning each observation a weight that reflects its importance

As an example of the need of weighted mean, consider the following sample of five purchases of a raw material over several months

Cost per Pound (\$)

3.00

3.40

2.80

2.90

3.25

| Note that | t the cost | per pol | und varie | es from | \$2. | 80 to \$3.40, | and quan | tity pure | chase | d ha | as varie | d fro | m 5 | 00, |
|-----------|------------|---------|-----------|---------|------|---------------|-----------|-----------|-------|------|----------|-------|-----|-----|
| to 2750. | Suppose | that m | nanager | asked | for | information | about the | mean | cost | per | pound | of th | e r | aw |
| material. | | | | | | | | | | | | | | |

$$\bar{x} = \frac{1200(3.00) + 500(3.40) + 2750(2.80) + 1000(2.90) + 800(3.25)}{1200 + 500 + 2750 + 1000 + 800} = \frac{18,500}{6250} = 2.96$$

If we would use a simple mean of the cost p.p.:

Purchase

2 3

4

5

(3.00 + 3.40 + 2.80 + 2.90 + 3.25)/5 = 15.35/5 = \$3.07

we overestimate the average cost!

Anderson et al Statistics for Business and Economics

31

Number of Pounds

1200

500

2750

1000

800

 $\overline{x} = \frac{\sum w_i x_i}{\sum w_i}$





Grouped data

CENTRE DE RECHERCHE PUBI

C

Data available in class intervals as summarized by a frequency distribution. Individual values of the original data are not available.



Mean for grouped data



 M_i = the midpoint for class *i* f_i = the frequency for class *i* n = the sample size

Variance for grouped data

$$s^{2} = \frac{\sum_{i}^{k} f_{i} (M_{i} - \overline{x})^{2}}{n-1}$$





Thank you for your attention



to be continued...