# Quantitative methods

Week #8-9

### Gergely Daróczi

Corvinus University of Budapest, Hungary

### 5 April 2013



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ ○○○

# Outline



### Midterm exam

- Final examination questions
- 3 Descriptive statistics
  - Relation between variables
    - Visual examples
    - Exact types
    - Further examples

## 5 Correlation

- Theoretical background
- Exercises
- · Limitations of the correlation coefficient
- Exercises

Results

### Midterm exam results



Gergely Daróczi (BCE)

< D > < B

▶ < 글 > < 글 >

। 5/4/2013 3/51 Results

### Midterm exam results



Gergely Daróczi (BCE)

< D > < B

▶ < 글 > < 글 >

। ► ≣ ∽ি৭ে 5/4/2013 3/51 Singleton, R. A. Jr. and Bruce C. Straits (1999): Approaches to Social Research. Third Edition. Oxford University Press: New York/Oxford.

Questions:

- What is reliability? How do the main rules concerning the order of survey questions improve the reliability and validity of survey data? (pp. 113-117, 292-296)
- What is meant by probability sampling? How do stratification and multistage cluster sampling affect sampling errors? Why? (pp. 141-142, 145-156)
- What are the main types of non-probability sampling? Explain why these types do not meet the criteria of probability samples. (pp. 157-169)
- What factors affect the desired sample size? (pp. 163-169)

Qualitative and quantitative variables in depth

### Qualitative variables:

- Nominal: exhaustive labels with no intersect (mutual exlcusivity) not in a specific order
- Ordinal: an (possible) ordered variable with exhaustive labels not intersecting

	Nominal	Ordinal	Interval	Ratio
Classification	Х	Х	Х	х
Rank order		Х	Х	Х
Equal intervals			Х	Х
Nonarbitrary zero				Х

### **Quantitative variables:**

- Interval: equal distances between the ordered labels (numbers)
- Ratio: a scale with a zero point

Computation

### For Simple Random Sampling:

• mean: 
$$\overline{x} = \frac{\sum_{i=1}^{n} x_i}{n}$$

- median
- mode
- min, max, range

• standard deviation: 
$$\sigma = \sqrt{\sum_{i=1}^{n} \frac{(x_i - \overline{x})^2}{n}}$$

- variance
- IQR

★ 문 ▶ ★ 문

We have measured 10-10 students in two classrooms.



Which class has higher students based on this small sample? Think about averages as good estimates of population parameters!

# **Descriptive statistics**

Averages





Gergely Daróczi (BCE)

< D > < B

▶ < 글 > < 글 >

# Descriptive statistics

Median and IQR





Gergely Daróczi (BCE)

Quantitative methods, 8-9/13

5/4/2013 9 / 51

#### A visual example



Gergely Daróczi (BCE)

5/4/2013 10 / 51

# The structure of the demo dataset

#### ggplot2/diamonds

Prices of 50,000 round cut diamonds

Description:

A dataset containing the prices and other attributes of almost 54,000 diamonds. The variables are as follows:

- price. price in US dollars (\\$326--\\$18,823)
- carat. weight of the diamond (0.2--5.01)
- cut. quality of the cut (Fair, Good, Very Good, Premium, Ideal)
- colour. diamond colour, from J (worst) to D (best)
- clarity. a measurement of how clear the diamond is (I1 (worst), SI1, SI2, VS1, VS2, VVS1, VVS2, IF (best))
- x. length in mm (0--10.74)
- y. width in mm (0--58.9)
- z. depth in mm (0--31.8)
- depth. total depth percentage = z / mean(x, y) = 2 \* z / (x + y) (43--79)
- table. width of top of diamond relative to widest point (43--95)

イロト イポト イヨト イヨト 一日

A visual example



ggplot(diamonds, aes(carat, price)) + geom\_point() + geom\_smooth() + ylab('') + scale\_y\_continuous(formatter="dollar") + theme\_bw() + opts(title="53.940 diamonds")

Gergely Daróczi (BCE)

#### A visual example



ggplot(diamonds, aes(clarity, fill=cut)) + geom\_bar() + ylab("N") + theme\_bw() + opts(title="53.940 diamonds")

5/4/2013 13 / 5

イロト イポト イヨト イヨト 一日

#### A visual example



ggplot(diamonds, aes(clarity)) + geom\_bar() + ylab("N") + facet\_wrap(~ cut) + theme\_bw() + opts(title="53.940 diamonds")

Gergely Daróczi (BCE)

5/4/2013 14 / 51

▶ < 글 > < 글 >

< 口 > < 🗇

#### A visual example



ggplot(diamonds, aes(carat, price, color=clarity)) + geom\_point() + ylab('') + scale\_y\_continuous(formatter="dollar") + theme\_bw() + opts(title="53.940 diamonds")

Gergely Daróczi (BCE)

5/4/2013 15 / 51

イロト イポト イヨト イヨト

#### A visual example



ggplot(diamonds, aes(carat, price, color=cut)) + geom\_point() + ylab('') + facet\_wrap(~ clarity,nr scale\_y\_continuous(formatter="dollar") + theme\_bw() + opts(title="53.940 diamonds")

5/4/2013 16 / 5

イロト イポト イヨト イヨト

# Test your knowledge!

Reliability and validity



A survey was taken place about diamonds available for sale on the Internet. What do you think of the reliability and validity of this research?

Gergely Daróczi (BCE)

Quantitative methods, 8-9/13

5/4/2013 17 / 51

# Types of variables and attributes in practice

Relationship between variables

### Possible relationship between variables:

- association,
- correlation,
- spurious relationship,
- influence,
- direction of influence,
- causality.

### STATISTICALLY SIGNIFICANT

#### Direction of influence



Gergely Daróczi (BCE)

5/4/2013 19 / 51

Diamonds' colors

GIA	Status: current	AGS	Status: coment	AGS	Status: historical: pre 1995	CIBJO Status: current	IDC Status: current		Scan. D.N.	Status: current	Old World Terms	Status: historical				
græd	e and deacription <sup>[8]</sup>	gra	ade and electronic colorimeter zcale <sup>[7]</sup>	9	rade and electronic colorimeter scale <sup>[7]</sup>	grade <sup>[9]</sup>	grade and description <sup>[8]</sup>		grade for .50ct and over	grade for under .50ct	series 1 scale[1]	zeries 2 scale <sup>[8]</sup>				
D		0	00.49	0	0-0.75	Exceptional white +	Exceptional white +	Coloring	River		Fin and Minite	Jager				
E	Goloriess	0.5	0.5-0.99	1	0.76-1.35	Exceptional white	Exceptional white	contineas	NO.		rines, white					
F		1.0	1.0-1.49		4 00 0 00	Rare white +	Rare white +			White	C	nore!				
Ģ		1.5	1.5-1.99	<u>*</u>	1.36-2.00	Rare white	Rare white	Colorfess when viewed through the crown	Top Wesselton		Fine white	Top Wesselton				
н		2.0	2.0-2.49	3	2.01-2.50	White	White		Wesselton		White	Wesselton				
1	Near Coloness	2.5	2.5-2.99	4	2.51-3.0				Top Crystal		Commercial White	Top Crystal				
J		3.0	3.0-3.49	5	3.01-3.75	Slightly tinted white	Slightly binted white		Crystal	Slightly binted white	Top silver cape	Crystal				
к		3.5	3.5-3.99									Slightly colored				
L	Faint Yellow	4.0	4.0-4.49	l°	3.70-4.5	3./0-4.5 Tinted write		I Inted white		Tintes white	Silver Cape	Top cape				
м		4.5	4.5-4.99			Tinted color 1			Cape		Light cape	Cape				
Ν		6.0	5.0-5.49	14	4.01-0.00		Tinted color 2					Low Cape				
0		6.5	5.5-5.99			1			Light yellow		Cape	Very light yellow				
Р	Very Light Yellow	6.0	6.0-6.49	8	5.51-7.0	Tinted color 2										
Q		6.5	6.5-6.99	1												
R		7.0	7.0-7.49			1										
s		7.5	7.5-7.99	9	7.01-8.5											
т		8.0	8.0-8.49	1			Tinted color	Slightly colored to colored		Tinted color						
U		8.5	8.5-8.99				1						Light yellow			
v		9.0	9.0-9.49	10	8.51-10				Yellow		Dark cape					
w	Light Yellow	9.5	9.5-9.99	1		Tinted color 3										
X Y Z		10	10+		10+											

Source: http://en.wikipedia.org/wiki/Diamond\_color

≣▶ ≣ •⁄) ৭.৫ 5/4/2013 20/51

イロト イポト イヨト イヨト

High correlation



A high correlation can be pointed out. Please explain!

Gergely Daróczi (BCE)

Quantitative methods, 8-9/13

5/4/2013 21 / 51

High correlation



### A high correlation can be pointed out. So what?

Gergely Daróczi (BCE)

5/4/2013 22 / 51

#### No correlation. No relationship?



Source: http://xkcd.com/323/

Gergely Daróczi (BCE)

Quantitative methods, 8-9/13

イロト 人間 とくほ とくほ とう

#### Correlation coefficient



Positive (direct: R = 1), negative (inverse: R = -1), linear, curvilinear and uncorrelated (R = 0) relationships R: correlation coefficient

Big shoes and smart kids (example)

We made a small research on the shoe size of some students in an elementary shool, where we also conducted a math exam. See detailed results below:

	Shoe size	Math result
1	29.75	26.67
2	29.75	33.33
3	29.75	41.67
4	31.50	35.00
5	31.50	46.67
6	31.50	63.33
7	31.50	70.00
8	33.25	30.00
9	33.25	38.33
10	33.25	56.67
11	35.00	26.67
12	35.00	40.00
13	35.00	43.33
14	35.00	46.67
15	35.00	53.33
16	38.50	55.00
17	40.25	45.00
18	42.00	58.33
19	42.00	76.67
20	42.00	77.50
21	42.00	100.00
22	43.75	70.83

Big shoes and smart kids (example)



Gergely Daróczi (BCE)

5/4/2013 26 / 51

Big shoes and smart kids (example)

We made a small research on the **age** and shoe size of some students in an elementary shool, where we also conducted a math exam. See detailed results below:

	Shoe size	Math result	Age
1	29.75	26.67	3
2	29.75	33.33	7
3	29.75	41.67	5
4	31.50	35.00	8
5	31.50	46.67	10
6	31.50	63.33	11
7	31.50	70.00	12
8	33.25	30.00	7
9	33.25	38.33	7
10	33.25	56.67	12
11	35.00	26.67	6
12	35.00	40.00	8
13	35.00	43.33	6
14	35.00	46.67	10
15	35.00	53.33	11
16	38.50	55.00	9
17	40.25	45.00	9
18	42.00	58.33	9
19	42.00	76.67	16
20	42.00	77.50	18
21	42.00	100.00	19
22	43.75	70.83	14

Big shoes and smart kids (example)



Gergely Daróczi (BCE)

5/4/2013 28 / 51

Big shoes and smart kids (example)



Big shoes and smart kids (example)



Big shoes and smart kids (example)



Gergely Daróczi (BCE)

5/4/2013 28 / 51

Big shoes and smart kids (example)



Big shoes and smart kids (example)



5/4/2013 28 / 51

Big shoes and smart kids (example)

# Partial correlation:

 $r_{math,size \cdot age} = 0.11$ 

 $r_{math,age\cdot size} = 0.87$ 

 $r_{size,age\cdot math} = 0.22$ 

Gergely Daróczi (BCE)

Quantitative methods, 8-9/13

5/4/2013 28 / 51

Covariation

For x and y variables the joint variability could be computed by :

$$COV(xy) = \sum_{i=1}^{n} \frac{(x_i - \overline{x})(y_i - \overline{y})}{n-1}$$
remember :  $\sigma = \sqrt{\sum_{i=1}^{n} \frac{(x_i - \overline{x})^2}{n}}$ 





Covariation



Ezekiel, M. (1930):

#### Henderson & Velleman (1981): Building multiple regression models interactively

< D > < B

30 / 51 5/4/2013

E ▶ < E >

Correlation

$$r_{xy} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2 \sum_{i=1}^{n} (y_i - \bar{y})^2}}$$

$$\frac{1.0 \quad 0.8 \quad 0.4 \quad 0.0 \quad -0.4 \quad -0.8 \quad -1.0}{(1.0 \quad 1.0 \quad 1.0 \quad -1.0 \quad -1.0 \quad -1.0)}$$

$$\frac{0.0 \quad 0.0 \quad 0.0}{(0.0 \quad 0.0 \quad 0$$

< D > < B

→ < Ξ → < Ξ</p>

5/4/2013 31 / 51

Partial correlation

$$\hat{r}_{XY\cdot \mathbf{Z}} = \frac{N\sum_{i=1}^{N} r_{X,i} r_{Y,i} - \sum_{i=1}^{N} r_{X,i} \sum_{i=1}^{N} r_{Y,i}}{\sqrt{N\sum_{i=1}^{N} r_{X,i}^{2} - (\sum_{i=1}^{N} r_{X,i})^{2}} \sqrt{N\sum_{i=1}^{N} r_{Y,i}^{2} - (\sum_{i=1}^{N} r_{Y,i})^{2}}}$$

so for three variables:

$$\hat{r}_{XY \cdot \mathbf{Z}} = \frac{r_{XY} - r_{X\mathbf{Z}}r_{Y\mathbf{Z}}}{\sqrt{(1 - r_{X\mathbf{Z}}^2)(1 - r_{Y\mathbf{Z}}^2)}}$$

Gergely Daróczi (BCE)

Quantitative methods, 8-9/13

5/4/2013 32 / 51

글 🕨 🖌 글

## Exercises

- What is correlation and partial correlation?
- Building upon your findings, compute the possible pairs of correlation coefficients on the below dataset!
- Also look for partial correlation and comment on your results!

Grade (mean)	Scholarship (in HUF)	Money spent on books (in HUF)
3.05	22000	3500
3.2	25000	3000
3.35	27000	2800
3.35	24000	3700
3.45	25000	2200
3.55	28000	3200
3.7	28000	3700
45	30000	4100
3.8	27000	4000
3.8	29000	3800

# Exercises

Solution

					$x_i - \overline{x}$			$(x_i - \overline{x})^2$				$(x_i - \overline{x})(y_i - \overline{x})$	$\overline{y}$ )	
	Grade	Scholarship	Books	Grade	Scholarship	Books	Grade	Scholarship	Books	G	irade-sch	Sch-books	Grade-books	
	3.05	22000.00	35000.00	-0.48	-4500.00	1000.00	0.225625	20250000	1000000		2137.5	-4500000	-475	
	3.20	25000.00	30000.00	-0.33	-1500.00	-4000.00	0.105625	2250000	16000000		487.5	6000000	1300	
	3.35	27000.00	28000.00	-0.18	500.00	-6000.00	0.030625	250000	36000000		-87.5	-3000000	1050	
	3.35	24000.00	37000.00	-0.18	-2500.00	3000.00	0.030625	6250000	9000000		437.5	-7500000	-525	
	3.45	25000.00	22000.00	-0.07	-1500.00	-12000.00	0.005625	2250000	144000000		112.5	18000000	900	
	3.55	28000.00	32000.00	0.02	1500.00	-2000.00	0.000625	2250000	4000000		37.5	-3000000	-50	
	3.70	28000.00	37000.00	0.18	1500.00	3000.00	0.030625	2250000	9000000		262.5	4500000	525	
	4.00	30000.00	41000.00	0.48	3500.00	7000.00	0.225625	12250000	49000000		1662.5	24500000	3325	
	3.80	27000.00	40000.00	0.28	500.00	6000.00	0.075625	250000	36000000		137.5	3000000	1650	
	3.80	29000.00	38000.00	0.28	2500.00	4000.00	0.075625	6250000	16000000		687.5	10000000	1100	
Min	3.05	22000.00	22000.00			Sum	0.80625	54500000	320000000		5875	48000000	8800	
Max	4.00	30000.00	41000.00			St. dev.	0.29930475	2460.80384	5962.84794	r	0.89	0.36	0.55	
Range	0.95	8000.00	19000.00							$r^2$	0.79	0.13	0.30	
Mean	3.53	26500.00	34000.00						partial c	orr	0.96	-0.24	0.46	
Median	3.50	27000.00	36000.00									•		







・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

- Correlation and causality
- Lazarsfeld paradigm
- Correlation and linearity

# Limitations of the correlation coefficient

Correlation does not imply causation!



Source: http://xkcd.com/552

< ロ ト < 同

▶ < 프 ▶ < 프 ▶</p>

Correlation does not imply causation! - Theoretical background

Aristotle: logic, syllogism – if  $(A \rightarrow B)\&(B \rightarrow C) \Rightarrow A \rightarrow C$ 

David Hume: scepticism

- "only correlation can actually be perceived [not causality]"
- see: our belief that the sun will rise tomorrow
- see: "If I see a billiard ball moving towards another, on a smooth table, I can easily conceive to stop upon contact."

Popper: falsification

Pearl, J. - *Causality: Models, Reasoning, and Inference*, Cambridge University Press, 2000

Lazarsfeld paradigm

# Stouffer: The American Soldier

Soldiers in branches with higher promotion rates are happier than soldiers in branches with lower rates of promotion. Lazarsfeld paradigm

Stouffer: The American Soldier

 $H_0$ : Soldiers in branches with higher promotion rates are happier than soldiers in branches with lower rates of promotion. **BUT:** 

"Soldiers in branches with higher promotion rates were more pessimistic about their own chances of being promoted than soldiers in branches with lower rates of promotion." Lazarsfeld paradigm

Stouffer: The American Soldier

 $H_0$ : Soldiers in branches with higher promotion rates are happier than soldiers in branches with lower rates of promotion. **BUT:** 

"Soldiers in branches with higher promotion rates were more pessimistic about their own chances of being promoted than soldiers in branches with lower rates of promotion."

# Keywords: reference group, relative deprivation

Gergely Daróczi (BCE)

# Limitations of the correlation coefficient

Correlation and linearity - Variations of the Same Theme



Source: Anscombe, F. J. (1973) Graphs in statistical analysis. American Statistician, 27, 17-21, C

Gergely Daróczi (BCE)

Correlation



Source: http://www.businessweek.com/magazine/correlation-or-causation-12012011-gfx.html

**Cross-correlation** 



#### 20-Year Lag Time Between Smoking and Lung Cancer

Gergely Daróczi (BCE)

≣▶ ≣ ∽Ω< 5/4/2013 42/51

**Top 25 Countries World Wide** Hungary 130.6 Sevchelles Belgium 126,2 Czech Rep. United States Netherlands UK Poland 115,9 114,4 Canada Russian Fed. Denmark 110.9 Estonia 106,1 Luxembourg Latvia Singapore 100,6 Slovenia Lithuania 98,6 Ireland 94.7 Italy New Zealand 89.1 Uruguay Iceland 85,4 Greece 84.4 Germany Australia

Lung Cancer Death Rate per 100,000 persons

Gergely Daróczi (BCE)

0 20 40 50 80 100 120 140 160

≣▶ ≣ ∽Ω< 5/4/2013 43/51

# Exercise

The data was extracted from the 1974 Motor Trend US magazine, and comprises fuel consumption and 10 aspects of automobile design and performance for 32 automobiles (1973-74 models).

- mpg: Miles/(US) gallon
- cyl: Number of cylinders
- disp: Displacement (cu.in.)
- hp: Gross horsepower
- o drat: Rear axle ratio
- wt: Weight (lb/1000)
- qsec: 1/4 mile time
- vs: V/S
- am: Transmission (0 = automatic, 1 = manual)
- gear: Number of forward gears
- carb: Number of carburetors

Source: Henderson and Velleman (1981), Building multiple regression models interactively. Biometrics, 37, 391-411.

• • = • • = •



Henderson and Velleman (1981), Building multiple regression models interactively. Biometrics, 37, 391-411= 🛌 😑

Gergely Daróczi (BCE)

#### Edgar Anderson's Iris Data



Anderson, Edgar (1935). The irises of the Gaspe Peninsula, Bulletin of the American Iris Society, 59, 2-5.

Gergely Daróczi (BCE)

Edgar Anderson's Iris Data



Anderson, Edgar (1935). The irises of the Gaspe Peninsula, Bulletin of the American Iris Society, 59, 2-5.

Gergely Daróczi (BCE)



Wind (miles per hour)

Gergely Daróczi (BCE)

5/4/2013 48 / 51





Gergely Daróczi (BCE)

5/4/2013 48 / 51

# Exercise #3

Real association?



Gergely Daróczi (BCE)

5/4/2013 48 / 51

# Computation exercises

Required formulas

$$\overline{x} = \frac{\sum_{i=1}^{n} x_i}{n} \quad S_x = \sqrt{\sum_{i=1}^{n} \frac{(x_i - \overline{x})^2}{n}}$$

$$r_{xy} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2 \sum_{i=1}^{n} (y_i - \bar{y})^2}}$$

$$\hat{r}_{XY \cdot \mathbf{Z}} = \frac{r_{XY} - r_{X\mathbf{Z}}r_{Y\mathbf{Z}}}{\sqrt{(1 - r_{X\mathbf{Z}}^2)(1 - r_{Y\mathbf{Z}}^2)}}$$
$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} \quad \hat{\beta} = \frac{COV(x, y)}{VAR(x)} \quad \hat{y} = \beta x + \alpha$$

イロト イポト イヨト イヨト

Compute the correlation and build linear models:

Grade	Monthly scholarship	Spent on books
3	22000 Ft	4000 Ft
4	24000 Ft	3000 Ft
5	27000 Ft	2500 Ft
3.5	24000 Ft	3500 Ft
2	23000 Ft	2000 Ft

▶ < 토 ▶ < 토</p>

# It was a pleasure!

Gergely Daróczi daroczi.gergely@btk.ppke.hu

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで