

# Inovace bakalářského studijního oboru Aplikovaná chemie

<http://aplchem.upol.cz>

CZ.1.07/2.2.00/15.0247

Tento projekt je spolufinancován  
Evropským sociálním fondem a státním  
rozpočtem České republiky.



INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

# Korelace



evropský  
sociální  
fond v ČR



EVROPSKÁ UNIE



MINISTERSTVO ŠKOLSTVÍ,  
MLÁDEŽE A TĚLOVÝCHOVY



OP Vzdělávání  
pro konkurenceschopnost



OKRESNÍ HOSPODÁŘSKÁ  
KOMORA OLOMOUC

INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

**Inovace bakalářského studijního  
oboru Aplikovaná chemie**

# Korelace

- závislost dvou náhodných veličin
- těsnost **lineární** závislosti
- kovarianční koeficient
- korelační koeficient
  - Pearsonův
  - Spearmanův

# Výběrová kovariance

$$C_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

$$C_{X,Y} = C_{Y,X} \text{ a } C_{(X,X)} = s^2(X)$$

více – do kovarianční matice

$$C = \begin{pmatrix} s_1^2 & C_{12} & \cdots & C_{1n} \\ C_{21} & s_n^2 & & C_{2n} \\ \vdots & & \ddots & \vdots \\ C_{n1} & C_{n2} & \cdots & s_n^2 \end{pmatrix}$$

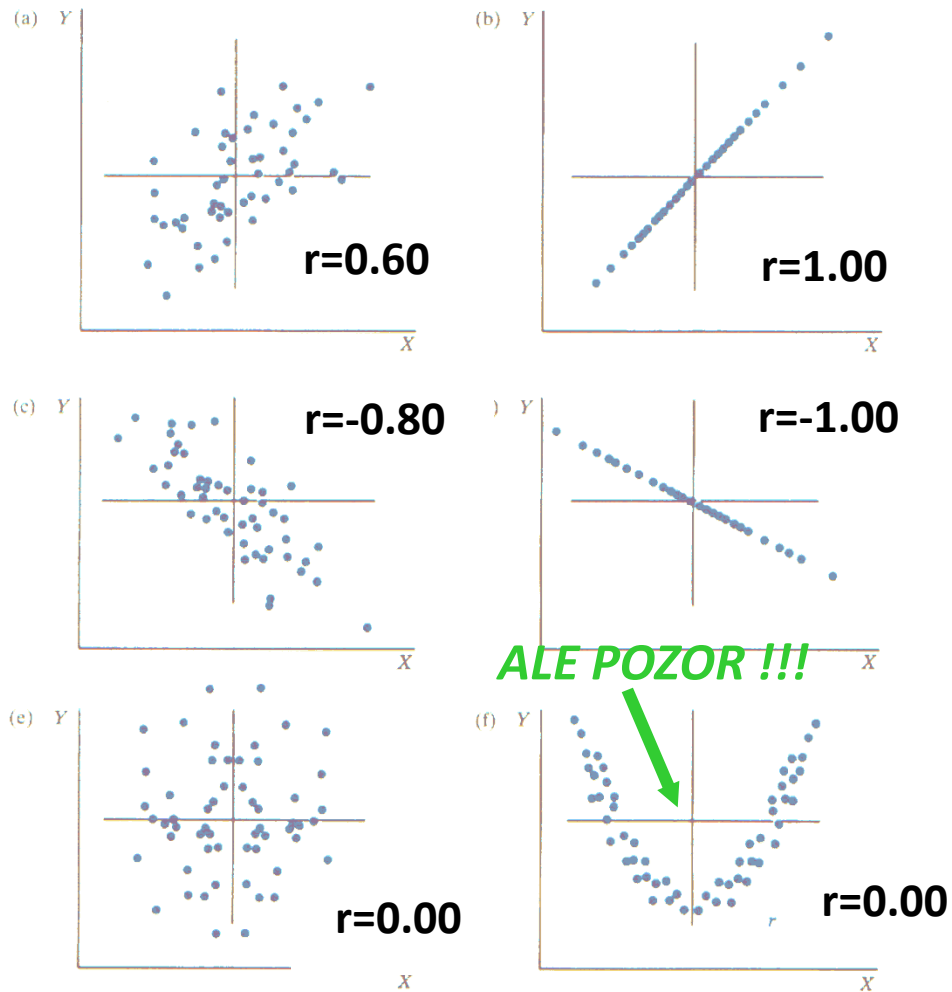
dále se používá u vícerozměrných metod  
např. metody hlavních komponent (PCA)

# Výběrová korelace - Pearson

$$\begin{aligned} \langle -1, 1 \rangle \quad r_{XY} &= \frac{\sum^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum^n (X_i - \bar{X})^2 \sum^n (Y_i - \bar{Y})^2}} \\ &= \frac{C_{XY}}{\sqrt{s_X^2 s_Y^2}} \\ &= \frac{\sum^n X_i Y_i - n \bar{X} \bar{Y}}{\sqrt{\sum^n (X_i - \bar{X})^2 \sum^n (Y_i - \bar{Y})^2}}, \end{aligned}$$

$r^2$  koeficient determinace

# Korelační koeficient



# Testování kor. koef.

$$t_e = \frac{r_{XY}}{\sqrt{1 - r_{XY}^2}} \sqrt{n - 2}. \quad (6.3)$$

Nulová hypotéza  $H_0 : \rho_{XY} = 0$  se na hladině  $\alpha$  zamítá ve prospěch oboustranné hypotézy  $H_1 : \rho_{XY} \neq 0$ , právě když je  $|t_e| \geq t_{n-2}(\alpha)$ , pokud se nepoužije absolutní hodnota lze zamítnout ve prospěch levo- či pravostranné hypotézy.

# Řešení v MS Excelu

**CORREL(výběr1;výběr2)**

**COVAR(výběr1;výběr2)**

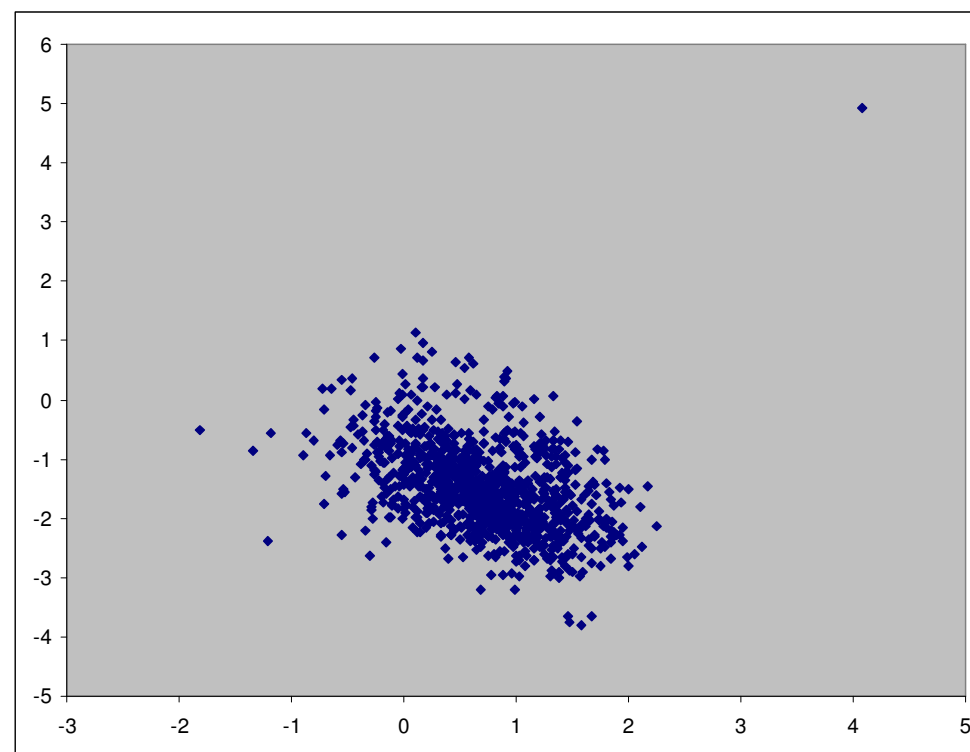
**CORREL(výběr1;výběr2)**

**-0.45807**

**COVAR(výběr1;výběr2)**

**-0.19072**

graf – XY bodový





# Řešení v MS Excelu

**Analýza dat: korelace, kovariance**  
korelační matice, kovarianční matice

proč tu nic není?

	<i>Sloupec 1</i>	<i>Sloupec 2</i>	<i>Sloupec 3</i>	<i>Sloupec 4</i>
Sloupec 1	1			
Sloupec 2	-0.50655	1		
Sloupec 3	-0.06564	0.055605	1	
Sloupec 4	-0.01541	-0.01958	-0.02002	1

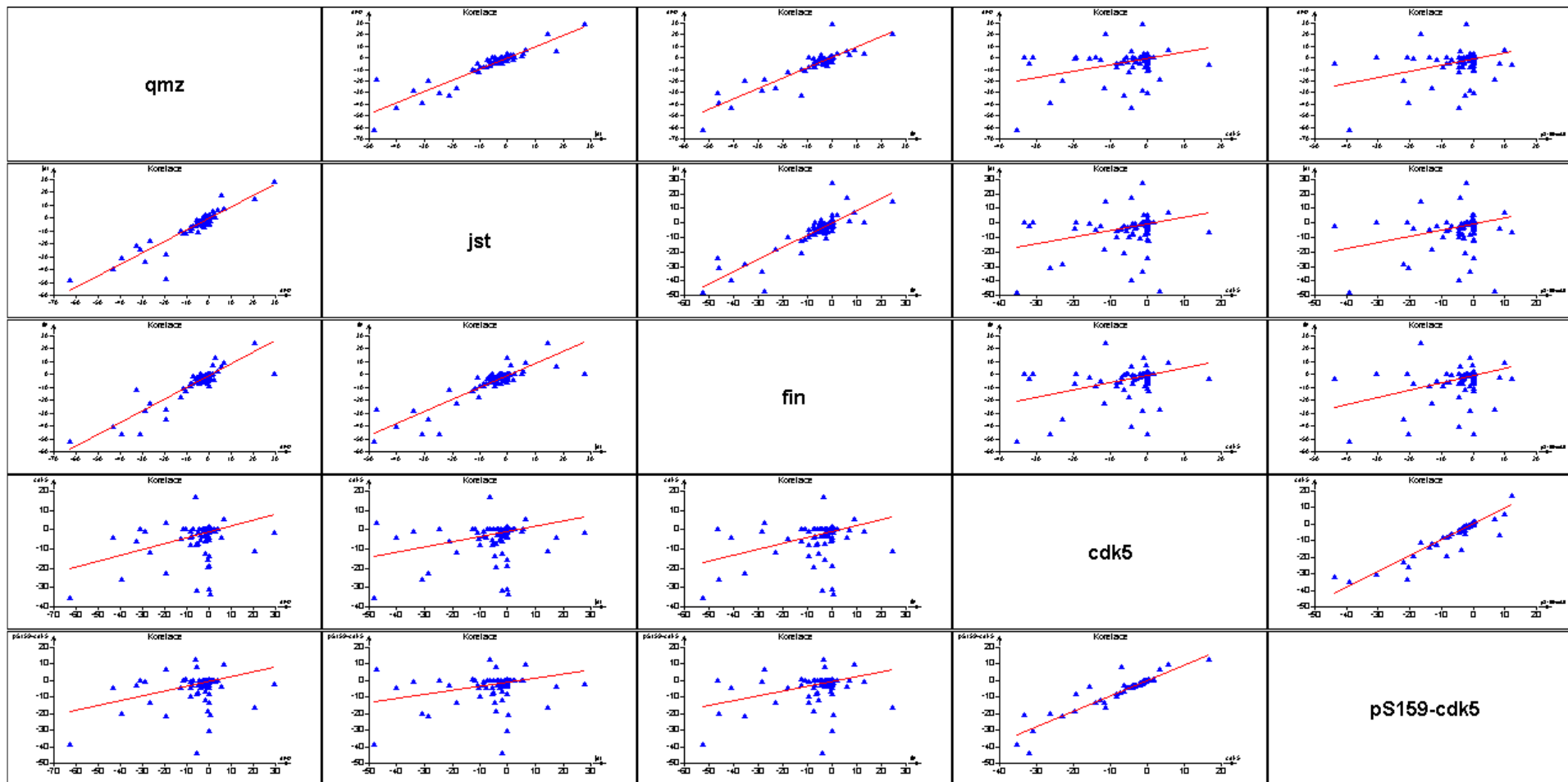
$$r_{XY} = r_{YX}$$

	<i>Sloupec 1</i>	<i>Sloupec 2</i>	<i>Sloupec 3</i>	<i>Sloupec 4</i>
Sloupec 1	1.000	-0.507	-0.066	-0.015
Sloupec 2	-0.507	1.000	0.056	-0.020
Sloupec 3	-0.066	0.056	1.000	-0.020
Sloupec 4	-0.015	-0.020	-0.020	1.000

# QCExpert

vícerozměrné metody - korelace

vizuální analýza: korelační grafy



Inovace bakalářského studijního oboru Aplikovaná chemie

INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

# Autokorelace

## 6.1.1 Autokorelace

Autokorelační koeficient 1. řádu  $r_1$  lze chápat jako „korelaci“ 1. hodnoty se 2., 2. se 3., 3. se 4. atd. Obecněji je autokorelační koeficient  $k$ -tého řádu  $r_k$  je definován jako

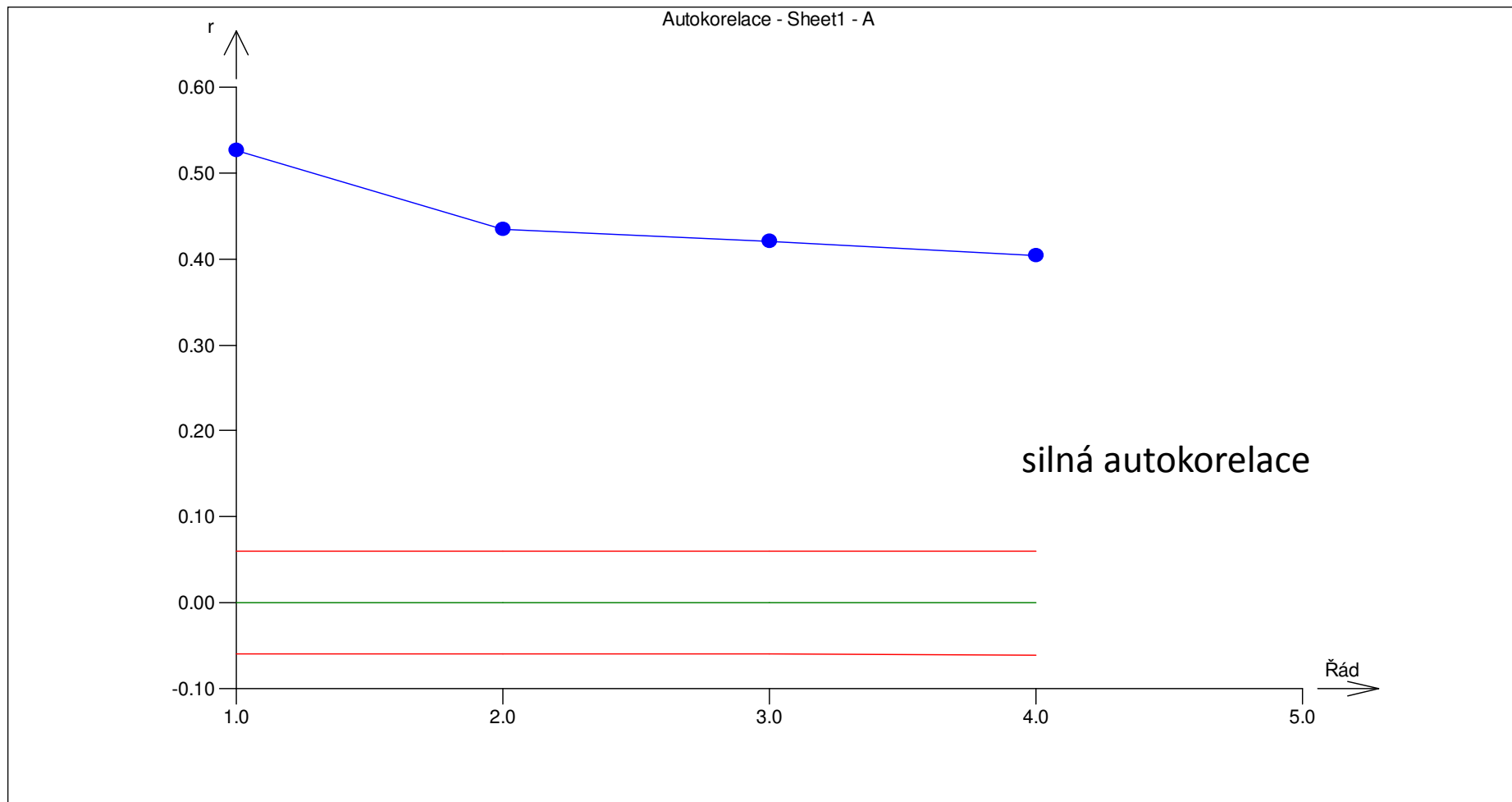
$$r_k = \frac{\sum_{i=1}^{N-k} (y_i - \bar{y})(y_{i+k} - \bar{y})}{\sum_{i=1}^N (y_i - \bar{y})^2}. \quad (6.4)$$

závislost  $n$ -té a  $n+k$ -té hodnoty, kde  $k$  je řád autokorelace

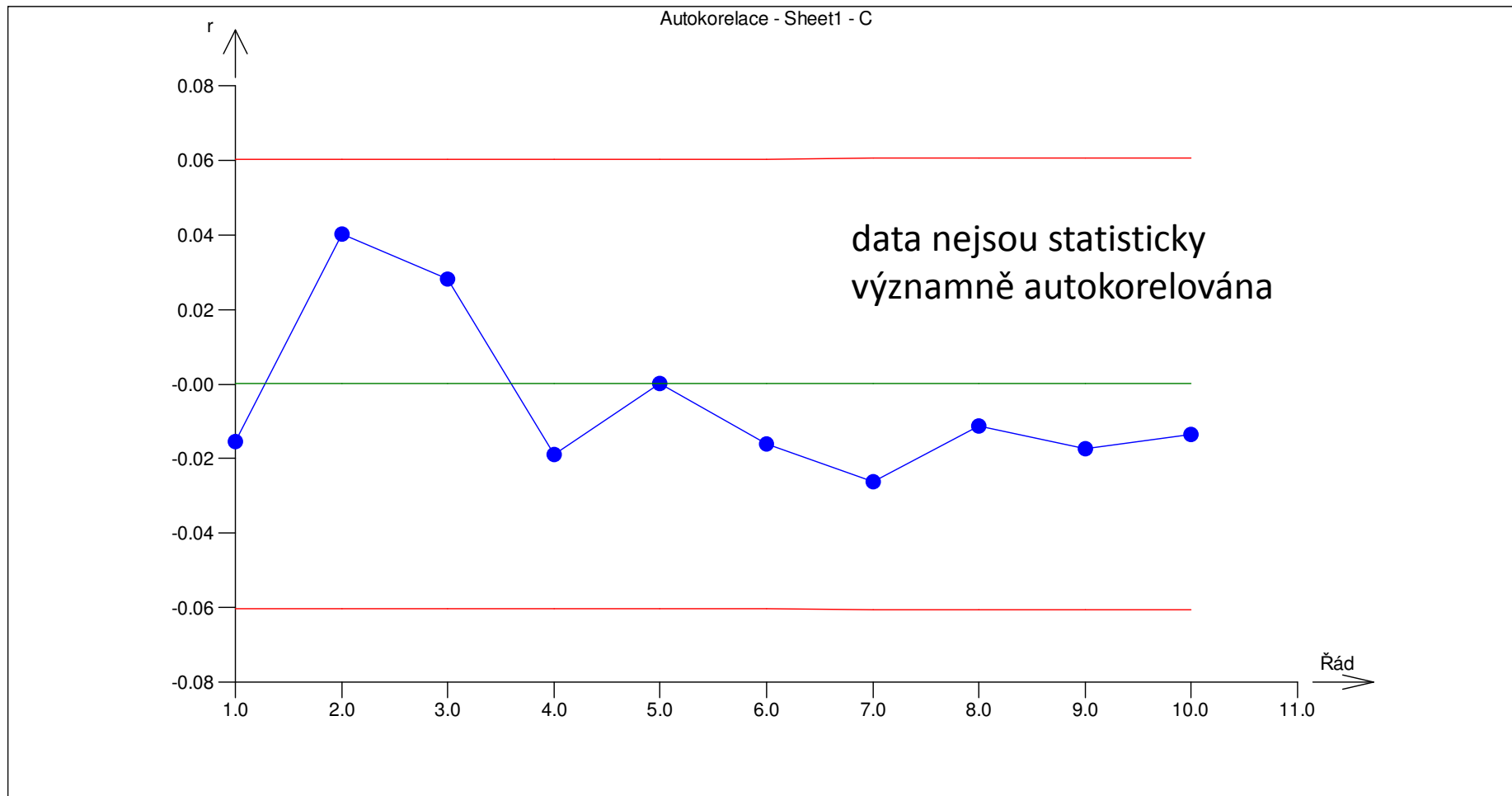
QCExpert – základní statistika

EXCEL – posunuté sloupce, a fce. CORREL

# Autokorelační graf



# Autokorelační graf



# Spearmanův korelační koeficient

pořadová korelace

$$r_S = 1 - \frac{6 \sum^n q_i^2}{n^3 - n},$$

rozdíl v pořadí

**Příklad 6.1** Máme  $n$  různých chemických látek a máme rozhodnout, zdali existuje korelace mezi interakční energií enzym-inhibitor, vypočtenou na základě molekulové mechaniky, a inhibičním účinkem ( $IC_{50}$ ) těchto látek na příslušný enzym.

látka	energie	$IC_{50}$	pořadí podle E	pořadí podle $IC_{50}$	rozdíl pořadí
n	kcal·mol <sup>-1</sup>	μM	k	l	q <sup>2</sup>
A	-25	12	1	1	0
B	-26	11	2	2	0
C	-28	10	3	3	0
D	-30	9	4	4	0
E	-31	7	5	6	1
F	-32	8	6	5	1
G	-33	1	7	8	1
H	-35	5	8	7	1

Odtud plyne, že  $r_S = 1 - \frac{24}{504} = 0,9524$ . Pro zajímavost uvádíme i  $r_{XY} = 0,8449$ . Na základě srovnání s kritickou hodnotou Spearmanova korelačního koeficientu, na hladině významnosti  $\alpha = 0,05$ ,  $r_S = 0,6905$  plyne, že existuje korelace mezi vypočtenou interakční energií a inhibičním účinkem. □

# Závislost nominálních veličin



evropský  
sociální  
fond v ČR



EVROPSKÁ UNIE



MINISTERSTVO ŠKOLSTVÍ,  
MLÁDEŽE A TĚLOVÝCHOVY



OP Vzdělávání  
pro konkurenceschopnost



OKRESNÍ HOSPODÁŘSKÁ  
KOMORA OLOMOUC

INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

**Inovace bakalářského studijního  
oboru Aplikovaná chemie**

# Kontingenční tabulky

$$N_{i.} = \sum_j N_{ij}, \quad N_{.j} = \sum_i N_{ij}, \quad \text{marginální četnosti}$$

$$X^2 = \sum_i \sum_j \frac{(N_{ij} - N_{i.}N_{.j}/n)^2}{N_{i.}N_{.j}/n}. \quad (6.7)$$

Nulovou hypotézu pak budeme zamítat na hladině  $\alpha$ , pokud bude platit

$$X^2 \geq \chi_{(a-1)(b-1)}^2(\alpha). \quad (6.8)$$



# Kontingenční tabulka

**Příklad 6.32.** Rozhodněte zda je ve studované populaci barva lidských vlasů závislá na pohlaví.

pohlaví	barva vlasů				celkem
	černá	hnědá	světlá	zrzavá	
mužské	32	43	16	9	100
ženské	55	65	64	16	200
celkem	87	108	80	25	300

$$X^2 = \sum_i \sum_j \frac{(N_{ij} - N_{i.}N_{.j}/n)^2}{N_{i.}N_{.j}/n}$$

Nejprve vyčíslíme marginální četnosti  $N_{M.} = 100$ ,  $N_{Z.} = 200$ ,  $N_{c.} = 87$ ,  $N_{.h} = 108$ ,  $N_{.s} = 80$ ,  $N_{.z} = 25$ . Následně vypočteme testovací statistiku  $X^2 = \frac{(32-100 \cdot 87/300)^2}{100 \cdot 87/300} + \dots = 8,987$ . Po porovnání s kritickou hodnotou  $\chi^2_{(3)}(0,05) = 7,815$  na hladině  $\alpha = 0,05$  zamítáme nulovou hypotézu o nezávislosti barvy vlasů na pohlaví v dané populaci.  $\square$

# QCExpert

data z předchozího př.

The screenshot displays the TriloByte QC.Expert 2.7 software interface. The main window shows a spreadsheet with a contingency table. The table has columns labeled 'x', 'c', 'h', 's', 'z', and 'F'. The data is as follows:

	x	c	h	s	z	F
1	m	32	43	16	9	
2	f	55	65	64	16	

A dialog box titled "Kontingenční tabulka" (Contingency Table) is open, showing the following settings:

- Název úlohy: Sheet1
- Hladina významnosti: 0.05
- Vyber sloupce četností: x, c, h, s, z
- Názvy řádků: x

The "Testování" (Testing) menu is open, showing the following options:

- Základní statistika...
- Porovnání dvou výběrů...
- Pravděpodobnostní modely...
- Transformace...
- Pravděpodobnostní kalkulátor...
- Testování (selected)
  - Síla a rozsah výběru
  - Testy
  - Kontingenční tabulka...
- Simulace
- ANOVA
- Responsní povrch...
- Vícerozměrné metody
- Regrese
- Kalibrace...
- Shewhartovy diagramy
- Rozšířené diagramy
- Způsobilost...
- Paretův diagram...
- Přejímka
- Grafy...
- Rychlé zobrazení dat...

# QCExpert

TriloByte QC.Expert 2.7 - [PROTOKOL]

Soubor Úpravy Formát Okno Nápověda

Analýza kontingenční tabulky

Název úlohy : Sheet1

Tabulka počtů

	c	h	s	z	Celkem
m	32	43	16	9	100
Teoretické:	( 29 )	( 36 )	( 27 )	( 8 )	
f	55	65	64	16	200
Teoretické:	( 58 )	( 72 )	( 53 )	( 17 )	
Celkem	87	108	80	25	300

Tabulka poměrů a pravděpodobností

	c	h	s	z	Celkem
m	0.106667	0.143333	0.053333	0.03	0.333333
Teoretické:	( 0.096667 )	( 0.12 )	( 0.088889 )	( 0.027778 )	
f	0.183333	0.216667	0.213333	0.053333	0.666667
Teoretické:	( 0.193333 )	( 0.24 )	( 0.177778 )	( 0.055556 )	
Celkem	0.29	0.36	0.266667	0.083333	1

Závěr

Nezávislost proměnných se zamítá

Hladina významnosti 0.05

Stupně volnosti 3

Chi2 statistika 8.967183908

Kritická hodnota 7.814727901

p-hodnota 0.02946176965

Sheet1 Sheet2

4.4.2006 22:40