

## Jednoduchá, mnohonásobná a parciální korelace

### Pearsonův koeficient korelace

Nechť  $X, Y$  jsou náhodné veličiny se středními hodnotami  $E(X), E(Y)$  a rozptyly  $D(X), D(Y)$ .

Číslo

$$R(X, Y) = \begin{cases} E\left(\frac{X - E(X)}{\sqrt{D(X)}} \cdot \frac{Y - E(Y)}{\sqrt{D(Y)}}\right) = \frac{C(X, Y)}{\sqrt{D(X)}\sqrt{D(Y)}} \text{ pro } \sqrt{D(X)}\sqrt{D(Y)} > 0 \\ 0 \text{ jinak} \end{cases}$$

se nazývá **Pearsonův koeficient korelace**.

(Pro výpočet Pearsonova koeficientu korelace musíme znát simultánní distribuční funkci  $\Phi(x, y)$  v obecném případě resp. simultánní hustotu pravděpodobnosti  $\varphi(x, y)$  ve spojitém případě resp. simultánní pravděpodobnostní funkci  $\pi(x, y)$  v diskrétním případě.)

### Vlastnosti Pearsonova koeficientu korelace

a)  $R(a_1, Y) = R(X, a_2) = R(a_1, a_2) = 0$

b)  $R(a_1 + b_1X, a_2 + b_2Y) = \text{sgn}(b_1b_2) R(X, Y) = \begin{cases} R(X, Y) \text{ pro } b_1b_2 > 0 \\ -R(X, Y) \text{ pro } b_1b_2 < 0 \end{cases}$

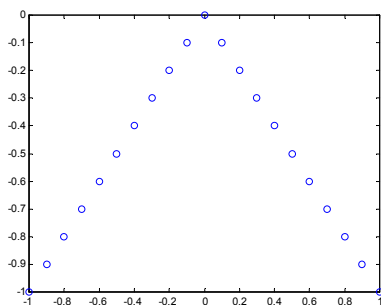
c)  $R(X, X) = 1$  pro  $D(X) \neq 0$ ,  $R(X, X) = 0$  jinak

d)  $R(X, Y) = R(Y, X)$

e)  $|R(X, Y)| \leq 1$  a rovnost nastane tehdy a jen tehdy, když mezi veličinami  $X, Y$  existuje s pravděpodobností 1 úplná lineární závislost, tj. existují konstanty  $a, b$  tak, že pravděpodobnost  $P(Y = a + bX) = 1$ . Přitom  $R(X, Y) = 1$ , když  $b > 0$  a  $R(X, Y) = -1$ , když  $b < 0$ . (Uvedená nerovnost se nazývá Cauchyova – Schwarzova – Buňakovského nerovnost.)

Z vlastností Pearsonova koeficientu korelace vyplývá, že se hodí pouze k měření těsnosti lineárního vztahu veličin  $X$  a  $Y$ . Při složitějších závislostech může dojít k paradoxní situaci, že Pearsonův koeficient korelace je nulový.

Ilustrace:



## Definice nekorelovanosti

Je-li  $R(X, Y) = 0$ , pak řekneme, že náhodné veličiny jsou **nekorelované**. (Znamená to, že mezi  $X$  a  $Y$  neexistuje žádná lineární závislost. Jsou-li náhodné veličiny  $X, Y$  stochasticky nezávislé, pak jsou samozřejmě i nekorelované.)

Je-li  $R(X, Y) > 0$ , pak řekneme, že náhodné veličiny jsou **kladně korelované**. (Znamená to, že s růstem hodnot veličiny  $X$  rostou hodnoty veličiny  $Y$  a s poklesem hodnot veličiny  $X$  klesají hodnoty veličiny  $Y$ .)

Je-li  $R(X, Y) < 0$ , pak řekneme, že náhodné veličiny jsou **záporně korelované**. (Znamená to, že s růstem hodnot veličiny  $X$  klesají hodnoty veličiny  $Y$  a s poklesem hodnot veličiny  $X$  rostou hodnoty veličiny  $Y$ .)

## Výběrový koeficient korelace

Nechť  $(X_1, Y_1), \dots, (X_n, Y_n)$  náhodný výběr rozsahu  $n$  z dvourozměrného rozložení daného distribuční funkcí  $\Phi(x, y)$ . Z tohoto dvourozměrného náhodného výběru můžeme stanovit:

$$\text{výběrové průměry } M_1 = \frac{1}{n} \sum_{i=1}^n X_i, \quad M_2 = \frac{1}{n} \sum_{i=1}^n Y_i,$$

$$\text{výběrové rozptyly } S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_1)^2, \quad S_2^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - M_2)^2,$$

výběrovou kovarianci  $S_{12} = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_1)(Y_i - M_2)$  a s jejich pomocí zavedeme

$$\text{výběrový koeficient korelace } R_{12} = \begin{cases} \frac{1}{n-1} \sum_{i=1}^n \frac{X - M_1}{S_1} \cdot \frac{Y - M_2}{S_2} = \frac{S_{12}}{S_1 S_2} & \text{pro } S_1 S_2 > 0 \\ 0 & \text{jinak} \end{cases}.$$

Vlastnosti Pearsonova koeficientu korelace se přenášejí i na výběrový koeficient korelace. (Výběrový koeficient korelace není nestranným odhadem skutečného koeficientu korelace, je odhadem vychýleným. Vychýlení je zanedbatečně malé pro rozsahy výběrů nad 30.)

## Pearsonův koeficient korelace dvourozměrného normálního rozložení

Nechť náhodný vektor  $(X, Y)$  má dvourozměrné normální rozložení s hustotou

$$\varphi(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[ \left( \frac{x-\mu_1}{\sigma_1} \right)^2 - 2\rho \frac{x-\mu_1}{\sigma_1} \frac{y-\mu_2}{\sigma_2} + \left( \frac{y-\mu_2}{\sigma_2} \right)^2 \right]},$$

přičemž  $\mu_1 = E(X)$ ,  $\mu_2 = E(Y)$ ,  $\sigma_1^2 = D(X)$ ,  $\sigma_2^2 = D(Y)$ ,  $\rho = R(X, Y)$ .

Marginální hustoty jsou:

$$\varphi_1(x) = \int_{-\infty}^{\infty} \varphi(x, y) dy = \dots = \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}},$$

$$\varphi_2(y) = \int_{-\infty}^{\infty} \varphi(x, y) dx = \dots = \frac{1}{\sigma_2 \sqrt{2\pi}} e^{-\frac{(y-\mu_2)^2}{2\sigma_2^2}}.$$

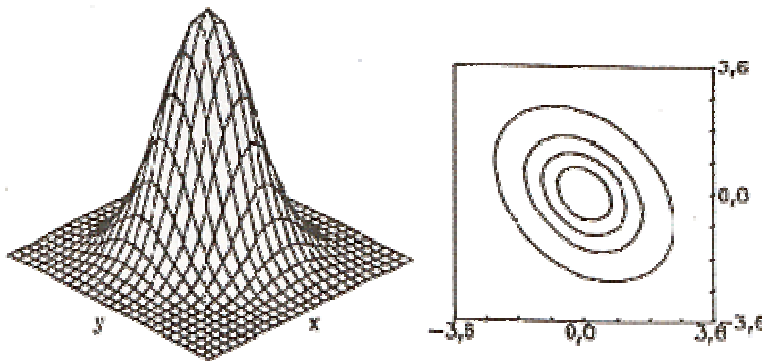
Je-li  $\rho = 0$ , pak pro  $\forall(x, y) \in \mathbb{R}^2 : \varphi(x, y) = \varphi_1(x)\varphi_2(y)$ , tedy náhodné veličiny  $X, Y$  jsou stochasticky nezávislé. Jinými slovy: **stochastická nezávislost složek  $X, Y$  normálně rozloženého náhodného vektoru je ekvivalentní jejich nekorelovanosti**. Pro jiná dvourozměrná rozložení to neplatí!

**Upozornění:** nadále budeme předpokládat, že  $(X_1, Y_1), \dots, (X_n, Y_n)$  je náhodný výběr rozsahu  $n$  z dvourozměrného normálního rozložení

$$N_2\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}\right).$$

Předpoklad dvourozměrné normality lze orientačně ověřit pomocí dvourozměrného tečkového diagramu: tečky by měly zhruba rovnoměrně vyplnit vnitřek elipsoidního obrazce. Vrstevnice hustoty dvourozměrného normálního rozložení jsou totiž elipsy:

Graf hustoty a vrstevnice dvourozměrného normálního rozložení s parametry  $\mu_1 = 0, \mu_2 = 0, \sigma_1^2 = 1, \sigma_2^2 = 1, \rho = -0,75$ :



Do dvourozměrného tečkového diagramu můžeme ještě zakreslit  $100(1-\alpha)\%$  elipsu konstantní hustoty pravděpodobnosti. Bude-li více než  $100\alpha\%$  teček ležet vně této elipsy, svědčí to o porušení dvourozměrné normality. Bude-li mít hlavní osa elipsy kladnou resp. zápornou směrnici, znamená to, že mezi veličinami  $X$  a  $Y$  existuje určitý stupeň přímé resp. nepřímé lineární závislosti.

### Testování hypotézy o nezávislosti

Na hladině významnosti  $\alpha$  testujeme  $H_0$ :  $X, Y$  jsou stochasticky nezávislé náhodné veličiny (tj.  $\rho = 0$ ) proti

- oboustranné alternativě  $H_1$ :  $X, Y$  nejsou stochasticky nezávislé náhodné veličiny (tj.  $\rho \neq 0$ )
- levostranné alternativě  $H_1$ :  $X, Y$  jsou záporně korelované náhodné veličiny (tj.  $\rho < 0$ )
- pravostranné alternativě  $H_1$ :  $X, Y$  jsou kladně korelované náhodné veličiny (tj.  $\rho > 0$ ).

Testová statistika má tvar:  $T_0 = \frac{R_{12} \sqrt{n-2}}{\sqrt{1-R_{12}^2}}$ .

Platí-li nulová hypotéza, pak  $T_0 \sim t(n-2)$ .

Kritický obor pro test  $H_0$  proti

- oboustranné alternativě:  $W = (-\infty, -t_{1-\alpha/2}(n-2)) \cup (t_{1-\alpha/2}(n-2), \infty)$ ,

- levostranné alternativě:  $W = (-\infty, -t_{1-\alpha}(n-2))$ ,

- pravostranné alternativě:  $W = (t_{1-\alpha}(n-2), \infty)$ .

$H_0$  zamítáme na hladině významnosti  $\alpha$ , když  $t_0 \in W$ .

### Příklad: Testování hypotézy o nezávislosti proti oboustranné alternativě

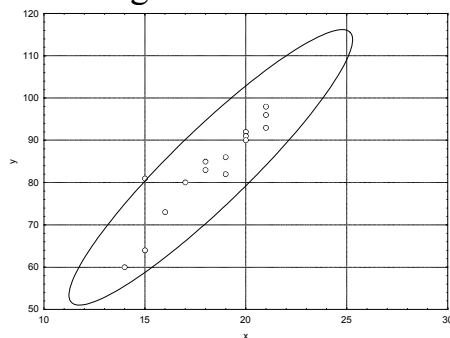
V dílně pracuje 15 dělníků. Byl u nich zjištěn počet směn odpracovaných za měsíc (náhodná veličina X) a počet zhotovených výrobků (náhodná veličina Y):

X 20 21 18 17 20 18 19 21 20 14 16 19 21 15 15

Y 92 93 83 80 91 85 82 98 90 60 73 86 96 64 81.

Orientačně ověřte dvourozměrnou normalitu dat, vypočtěte výběrový koeficient korelace mezi X a Y a na hladině 0,01 testujte hypotézu o nezávislosti X a Y proti oboustranné alternativě.

**Řešení:** Dvourozměrnou normalitu dat ověříme pomocí dvourozměrného tečkového diagramu.



Vidíme, že předpoklad dvourozměrné normality je oprávněný.

Vypočteme realizace

výběrových průměrů:  $m_1 = \frac{1}{n} \sum_{i=1}^n x_i = 18,267$ ,  $m_2 = \frac{1}{n} \sum_{i=1}^n y_i = 83,6$ ,

výběrových rozptylů:  $s_1^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - m_1)^2 = 5,6381$ ,  $s_2^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - m_2)^2 = 121,4$ ,

výběrové kovariance:  $s_{12} = \frac{1}{n-1} \sum_{i=1}^n (x_i - m_1)(y_i - m_2) = 24,2571$ ,

výběrového koeficientu korelace:  $r_{12} = \frac{s_{12}}{s_1 s_2} = 0,927$ .

Realizace testové statistiky:  $t_0 = \frac{r_{12} \sqrt{n-2}}{\sqrt{1-r_{12}^2}} = 8,912$ ,

kritický obor  $W = (-\infty, -t_{0,995}(13)) \cup (t_{0,995}(13), \infty) = (-\infty, -3,012) \cup (3,012, \infty)$ .

Protože  $t_0 \in W$ , hypotézu o nezávislosti veličin X a Y zamítáme na hladině významnosti 0,01. S rizikem omylu nejvýše 1% jsme tedy prokázali, že mezi počtem směn odpracovaných za měsíc a počtem zhotovených výrobků existuje závislost.

### Výpočet pomocí systému STATISTICA

Vytvoříme datový soubor o dvou proměnných X, Y a 15 případech. Dvourozměrnou normalitu dat ověříme pomocí dvourozměrného tečkového diagramu – viz výše.

Statistiky – Základní statistiky/tabulky – Korelační matice – OK – 1 seznam proměn. – X, Y – OK – na záložce Možnosti vybereme Zobrazit detailní tabulku výsledků – Výpočet.

Korelace (smeny a výrobky.sta) Označ. korelace jsou významné na hlad. p < ,05000 (Celé případy vynechány u ChD)											
Prom. X & prom. Y	Průměr	Sm.Odch.	r(X,Y)	r2	t	p	N	Konst. záv.: Y	Směr. záv.: Y	Konst. záv.: X	Směrnic záv.: X
X	18,26667	2,37447									
X	18,26667	2,37447	1,000000	1,000000			15	0,000000	1,000000	0,000000	1,000000
X	18,26667	2,37447									
Y	83,60000	11,01817	0,927180	0,859663	8,923795	0,000001	15	5,010135	4,302365	1,562407	0,199812
Y	83,60000	11,01817									
X	18,26667	2,37447	0,927180	0,859663	8,923795	0,000001	15	1,562407	0,199812	5,010135	4,302365
Y	83,60000	11,01817									
Y	83,60000	11,01817	1,000000	1,000000			15	0,000000	1,000000	0,000000	1,000000

Výběrový koeficient korelace se realizoval hodnotou 0,92718, testová statistika nabyla hodnoty 8,924, odpovídající p-hodnota je 0,000001, tedy na hladině významnosti 0,01 zamítáme hypotézu o nezávislosti veličin X, Y.

### Výpočet pomocí systému SPSS

Analyze – Correlate – Bivariate – Variables X, Y – OK

#### Correlations

		X	Y
X	Pearson Correlation	1,000	,927**
	Sig. (2-tailed)		,000
	N	15	15
Y	Pearson Correlation	,927**	1,000
	Sig. (2-tailed)	,000	
	N	15	15

\*\* . Correlation is significant at the 0.01 level (2-tailed).

Vidíme, že na rozdíl od systému STATISTICA systém SPSS neposkytuje hodnotu testové statistiky, ve výstupní tabulce lze najít pouze realizaci výběrového koeficientu korelace, odpovídající p-hodnotu a rozsah náhodného výběru.

## Interval spolehlivosti pro korelační koeficient

Náhodná veličina  $Z = \frac{1}{2} \ln \frac{1+R_{12}}{1-R_{12}}$  má přibližně normální rozložení se střední

hodnotou  $E(Z) = \frac{1}{2} \ln \frac{1+\rho}{1-\rho} + \frac{\rho}{2(n-1)}$  (2. sčítanec lze při větším  $n$  zanedbat) a roz-

ptylem  $D(Z) = \frac{1}{n-3}$ . Standardizací veličiny  $Z$  dostaneme veličinu  $U = \frac{Z - E(Z)}{\sqrt{D(Z)}}$ ,

která má asymptoticky rozložení  $N(0,1)$ . Tudíž  $100(1-\alpha)\%$  asymptotický interval

spolehlivosti pro  $\frac{1}{2} \ln \frac{1+\rho}{1-\rho}$  bude mít meze  $Z \pm \frac{u_{1-\alpha/2}}{\sqrt{n-3}}$ . Interval spolehlivosti pro  $\rho$

pak dostaneme zpětnou transformací.

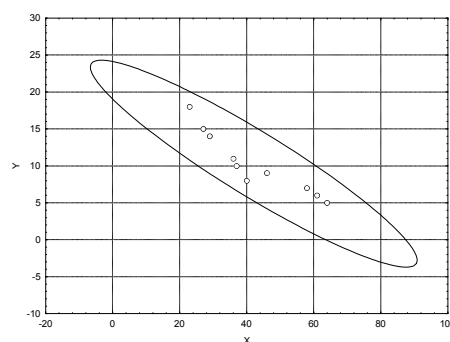
**Poznámka:** Jelikož  $Z = \operatorname{arctgh} R_{12}$ , dostáváme  $R_{12} = \operatorname{tgh} Z$  a meze intervalu spolehlivosti pro  $\rho$  můžeme psát ve tvaru  $\operatorname{tgh}\left(Z \pm \frac{u_{1-\alpha/2}}{\sqrt{n-3}}\right)$ , přičemž  $\operatorname{tgh} x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ .

**Příklad:** Pracovník personálního oddělení určité firmy zkoumá, zda existuje vztah mezi počtem dní absence za rok (veličina  $Y$ ) a věkem pracovníka (veličina  $X$ ). Proto náhodně vybral údaje o 10 pracovnících.

Č.prac.	1	2	3	4	5	6	7	8	9	10
X	27	61	37	23	46	58	29	36	64	40
Y	15	6	10	18	9	7	14	11	5	8

Za předpokladu, že uvedené údaje tvoří číselné realizace náhodného výběru rozsahu 10 z dvourozměrného normálního rozložení, vypočtete výběrový korelační koeficient a na hladině významnosti 0,05 testujte hypotézu, že  $X$  a  $Y$  jsou nezávislé náhodné veličiny. Sestrojte 95% asymptotický interval spolehlivosti pro skutečný korelační koeficient  $\rho$ .

**Řešení:** Předpoklad o dvourozměrné normalitě dat ověříme orientačně pomocí dvourozměrného tečkového diagramu.



Vzhled diagramu svědčí o tom, že předpoklad je oprávněný.

Testujeme  $H_0: \rho = 0$  proti  $H_1: \rho \neq 0$ . Vypočítáme  $R_{12} = -0,9325$ , tedy mezi věkem pracovníka a počtem dnů pracovní neschopnosti existuje silná nepřímá lineární závislost. Testová statistika:  $T = -7,3053$ , kvantil  $t_{0,975}(8) = 2,306$ , kritický obor  $W = (-\infty, -2,306) \cup (2,306, \infty)$ . Jelikož  $T \in W$ , zamítáme na hladině významnosti 0,05 hypotézu o nezávislosti veličin X a Y.

Vypočítáme  $Z = \frac{1}{2} \ln \frac{1+R_{12}}{1-R_{12}} = \frac{1}{2} \ln \frac{1-0,9325}{1+0,9325} = -1,6772$ . Meze 95% asymptotického

intervalu spolehlivosti pro  $\rho$  jsou  $\text{tgh}\left(-1,6772 \pm \frac{1,96}{\sqrt{7}}\right)$ , tedy  $-0,9842 < \rho < -0,7336$  s pravděpodobností přibližně 0,95.

### Výpočet pomocí systému STATISTICA:

Ve STATISTICE vypočteme meze  $100(1-\alpha)\%$  asymptotického intervalu spolehlivosti pro koeficient korelace  $\rho$  tak, že otevřeme nový datový soubor se dvěma proměnnými (pojmenujeme je DM a HM) a jedním případem.

Do Dlouhého jména proměnné DM zapíšeme příkaz

=  $\text{TanH}(0,5 * \log((1-0,9325)/(1+0,9325)) - \text{VNormal}(0,975;0;1)/\text{sqrt}(7))$

a do Dlouhého jména proměnné HM zapíšeme příkaz

=  $\text{TanH}(0,5 * \log((1-0,9325)/(1+0,9325)) + \text{VNormal}(0,975;0;1)/\text{sqrt}(7))$

	1 DM	2 HM
1	-0,98425	-0,73358

95% asymptotický interval spolehlivosti pro koeficient korelace  $\rho$  má tedy meze  $-0,98425$  a  $-0,73358$ . (Protože nepokrývá hodnotu 0, zamítáme hypotézu o nezávislosti veličin X, Y na asymptotické hladině významnosti 0,05.)

### Výpočet pomocí systému SPSS

Vzhledem k tomu, že v SPSS není dostupná funkce TanH, využijeme toho, že

$\text{tgh}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ . Dále pro výpočet kvantilů rozložení  $N(0,1)$  použijeme funkci

IDF.NORMAL(x, mean, std).

K datovému souboru přidáme 4 proměnné pojmenované pom1, pom2, dm, hm.

Transform – Compute Variable – Target Variable pom1 =

$0,5 * \ln((1-0,9325)/(1+0,9325)) - \text{IDF.Normal}(0,975,0,1)/\text{sqrt}(7)$

OK – OK

Transform – Compute Variable – Target Variable pom2 =

$0,5 * \ln((1-0,9325)/(1+0,9325)) + \text{IDF.Normal}(0,975,0,1)/\text{sqrt}(7)$

OK - OK

Transform – Compute Variable – Target Variable dm =

$(\exp(\text{pom1}) - \exp(-\text{pom1})) / (\exp(\text{pom1}) + \exp(-\text{pom1}))$

OK - OK

Transform – Compute Variable – Target Variable hm =  
(exp(pom2)-exp(-pom2))/(exp(pom2)+exp(-pom2))

OK – OK

Dostaneme výsledek dm = -0,9842 a hm = -0,7336.

### Varianční, kovarianční a korelační matice

Nechť  $\mathbf{X} = (X_1, \dots, X_p)'$  je náhodný vektor. Označme

$\mu_i = E(X_i)$  střední hodnotu náhodné veličiny  $X_i$ ,

$\sigma_i^2 = D(X_i)$  rozptyl náhodné veličiny  $X_i$ ,

$\sigma_{ij} = C(X_i, X_j)$  kovarianci náhodných veličin  $X_i, X_j$

$\rho_{ij} = R(X_i, X_j)$  koeficient korelace náhodných veličin  $X_i, X_j$

Vektor  $E(\mathbf{X}) = (\mu_1, \dots, \mu_p)'$  se nazývá **vektor středních hodnot** náhodného vektoru  $\mathbf{X}$ .

Čtvercová matice řádu  $p$   $\text{var}(\mathbf{X}) = (\sigma_{ij})_{i,j=1, \dots, p}$  se nazývá **varianční matice** náhodného vektoru  $\mathbf{X}$ .

Čtvercová matice řádu  $p$   $\text{cor}(\mathbf{X}) = (\rho_{ij})_{i,j=1, \dots, p}$  se nazývá **korelační matice** náhodného vektoru  $\mathbf{X}$ .

Je zřejmé, že varianční matice a korelační matice jsou symetrické.

Nechť  $\mathbf{X} = (X_1, \dots, X_p)'$  a  $\mathbf{Y} = (Y_1, \dots, Y_q)'$  jsou náhodné vektory.

Matice typu  $p \times q$   $\text{cov}(\mathbf{X}, \mathbf{Y}) = (C(X_i, Y_j))$  se nazývá **kovarianční matice** vektorů  $\mathbf{X}, \mathbf{Y}$ .

Matice typu  $p \times q$   $\text{cor}(\mathbf{X}, \mathbf{Y}) = (\rho(X_i, Y_j))$  se nazývá **korelační matice** vektorů  $\mathbf{X}, \mathbf{Y}$ .

### Odhady vektoru středních hodnot, varianční a korelační matice jednoho náhodného vektoru $\mathbf{X}$

Nechť  $\mathbf{X}$  je náhodný vektor, který má  $p$ -rozměrné rozložení s vektorem středních hodnot  $\boldsymbol{\mu}$ , varianční maticí  $\text{var}(\mathbf{X})$  a korelační maticí  $\text{cor}(\mathbf{X})$ . Nechť je dán náhodný výběr  $\mathbf{X}_1 = (X_{11}, \dots, X_{1p})', \dots, \mathbf{X}_n = (X_{n1}, \dots, X_{np})'$  rozsahu  $n$  z tohoto rozložení.

Nestranný odhad vektoru  $\boldsymbol{\mu}$  je **vektor výběrových průměrů**  $\mathbf{M} = (M_1, \dots, M_p)'$ ,

kde  $M_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$  je výběrový průměr  $j$ -tého výběru,  $j = 1, \dots, p$ .

Nestranný odhad matice  $\text{var}(\mathbf{X})$  je **výběrová varianční matice**  $\mathbf{S} = (S_{ij}) =$

$\frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \mathbf{M})(\mathbf{X}_i - \mathbf{M})'$  řádu  $p$ .

Vychýlený odhad matice  $\text{cor}(\mathbf{X})$  je **výběrová korelační matice**  $\mathbf{R} = (R_{ij})$ , kde  $R_{ij}$  je výběrový korelační koeficient  $i$ -té a  $j$ -té složky vektoru  $\mathbf{X}$ , tedy

$R_{ij} = \frac{S_{ij}}{\sqrt{S_{ii}} \sqrt{S_{jj}}}$ ,  $i, j = 1, \dots, p$ . (Je zřejmé, že diagonální prvky matice  $\mathbf{R}$  jsou jed-

ničky a matice  $\mathbf{R}$  je symetrická.)



## Odhady kovarianční a korelační matice dvou náhodných vektorů $\mathbf{X}$ , $\mathbf{Y}$

Nechť náhodný vektor  $\mathbf{X}$  má  $p$ -rozměrné rozložení a nechť  $\mathbf{X}_1, \dots, \mathbf{X}_n$  je náhodný výběr z tohoto rozložení. Nechť náhodný vektor  $\mathbf{Y}$  má  $q$ -rozměrné rozložení a nechť  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  je náhodný výběr z tohoto rozložení. Předpokládejme, že obě rozložení mají konečné druhé momenty. Nechť  $\text{cov}(\mathbf{X}, \mathbf{Y})$  je kovarianční matice těchto vektorů a  $\text{cor}(\mathbf{X}, \mathbf{Y})$  je korelační matice těchto vektorů. Označme

$$M_{X_j} = \frac{1}{n} \sum_{i=1}^n X_{ij}, j = 1, \dots, p, M_{Y_j} = \frac{1}{n} \sum_{i=1}^n Y_{ij}, j = 1, \dots, q,$$

$$\mathbf{M}_X = (M_{X_1}, \dots, M_{X_p})', \mathbf{M}_Y = (M_{Y_1}, \dots, M_{Y_q})'.$$

Nestranným odhadem kovarianční matice  $\text{cov}(\mathbf{X}, \mathbf{Y})$  vektorů  $\mathbf{X}$ ,  $\mathbf{Y}$  je **výběrová kovarianční matice** vektorů  $\mathbf{X}$ ,  $\mathbf{Y}$  definovaná vzorcem  $\mathbf{S}_{XY} = (S_{ij}) =$

$$\frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \mathbf{M}_X)(\mathbf{Y}_i - \mathbf{M}_Y)', i = 1, \dots, p, j = 1, \dots, q.$$

Vychýleným odhadem korelační matice  $\text{cor}(\mathbf{X}, \mathbf{Y})$  vektorů  $\mathbf{X}$ ,  $\mathbf{Y}$  je **výběrová korelační matice** vektorů  $\mathbf{X}$ ,  $\mathbf{Y}$  definovaná vzorcem  $\mathbf{R}_{XY} = (R_{ij})$ , kde  $R_{ij}$  je výběrový korelační koeficient  $i$ -té a  $j$ -té složky vektorů  $\mathbf{X}$ ,  $\mathbf{Y}$ ,  $i = 1, \dots, p, j = 1, \dots, q$ .

## Koeficient mnohonásobné korelace a výběrový koeficient mnohonásobné korelace

Intenzitu lineární závislosti mezi náhodnou veličinou  $Y$  a náhodným vektorem  $\mathbf{X} = (X_1, \dots, X_p)'$  měříme pomocí **koeficientu mnohonásobné korelace**  $\rho_{Y, \mathbf{X}}$ . Jeho druhá mocnina je dána vzorcem

$$\rho_{Y, \mathbf{X}}^2 = \text{cor}(Y, \mathbf{X}) \text{cor}(\mathbf{X})^{-1} \text{cor}(\mathbf{X}, Y).$$

Má tyto vlastnosti:

- $\rho_{Y, \mathbf{X}} \geq 0$
- $\rho_{Y, \mathbf{X}} \geq |\rho(Y, X_i)|$  pro  $\forall i = 1, \dots, p$
- $\rho_{Y, X_1, \dots, X_p} \geq \dots \geq \rho_{Y, X_1, X_2} \geq \rho(Y, X_1)$
- $\rho_{Y, \mathbf{X}} = 1 \Leftrightarrow$  existují konstanty  $\beta_0, \beta_1, \dots, \beta_p$  tak, že  $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ .

Nechť náhodný vektor  $(Y, X_1, \dots, X_p)'$  má  $(p+1)$ -rozměrné rozložení s koeficientem mnohonásobné korelace  $\rho_{Y, \mathbf{X}}$ .

Nechť je dán náhodný výběr  $(Y_1, X_{11}, \dots, X_{1p})', \dots, (Y_n, X_{n1}, \dots, X_{np})'$  rozsahu  $n$  z tohoto rozložení. Pak jako odhad  $\rho_{Y, \mathbf{X}}$  slouží **výběrový koeficient mnohonásobné korelace**  $r_{Y, \mathbf{X}}$ , jehož druhá mocnina je dána vzorcem

$$r_{Y, \mathbf{X}}^2 = \mathbf{R}_{YX} \mathbf{R}^{-1} \mathbf{R}_{XY},$$

kde  $\mathbf{R}_{YX}$  je výběrová korelační matice veličiny  $Y$  a vektoru  $\mathbf{X}$  (v tomto případě se redukuje na vektor  $(r_{YX_1}, \dots, r_{YX_p})$ ) a  $\mathbf{R}$  je výběrová korelační matice vektoru  $\mathbf{X}$ .

Vlastnosti (a), (b), (c), (d) koeficientu mnohonásobné korelace se přenášejí i na výběrový koeficient mnohonásobné korelace.

## Testování hypotézy o nezávislosti veličiny Y a vektoru X

### Popis testu

Nechť náhodný výběr  $(Y_1, X_{11}, \dots, X_{1p})', \dots, (Y_n, X_{n1}, \dots, X_{np})'$  pochází z  $(p+1)$ -rozměrného normálního rozložení, které má koeficient mnohonásobné korelace  $\rho_{Y, X}$ . Musí platit  $n > p+1$ .

Testujeme hypotézu  $H_0: \rho_{Y, X} = 0$  proti  $H_1: \rho_{Y, X} \neq 0$ . Vzhledem k tomu, že se jedná o výběr z  $(p+1)$ -rozměrného normálního rozložení, testujeme, zda existuje závislost mezi veličinou Y a vektorem X. (Je-li  $\rho_{Y, X} = 0$ , pak z vlastnosti (b) plyne, že  $\rho(Y, X_i) = 0$  pro všechna  $i = 1, \dots, p$ , tudíž náhodné veličiny Y a  $X_i$  jsou stochasticky nezávislé pro všechna  $i = 1, \dots, p$ .)

Testová statistika  $F = \frac{n-p-1}{p} \cdot \frac{r_{Y, X}^2}{1-r_{Y, X}^2}$  se řídí rozložením  $F(p, n-p-1)$ , pokud  $H_0$

platí. Kritický obor:  $W = \langle F_{1-\alpha}(p, n-p-1), \infty \rangle$ . Jestliže  $F \in W$ ,  $H_0$  zamítáme na hladině významnosti  $\alpha$ .

### Příklad

Při zkoumání závislosti hodinové výkonnosti dělníka (veličina Y – v kusech) na jeho věku (veličina  $X_1$  – v letech) a době zapracovanosti (veličina  $X_2$  – v letech) byly u 10 náhodně vybraných dělníků zjištěny tyto údaje:

Y	67	65	75	66	77	84	69	60	70	66
$X_1$	43	40	49	46	41	41	48	34	32	42
$X_2$	6	8	14	14	8	12	16	1	5	7

Za předpokladu, že uvedené hodnoty představují číselné realizace náhodného výběru rozsahu 10 ze třírozměrného normálního rozložení, testujte na hladině významnosti 0,05 hypotézu, že výkon dělníka nezávisí na jeho věku a době zapracovanosti.

### Řešení:

$r_{YX_1} = 0,2287, r_{YX_2} = 0,4538, r_{X_1X_2} = 0,847, R_{YX} = (0,2287, 0,4538)'$ ,

$$\mathbf{R} = \begin{pmatrix} 1 & 0,847 \\ 0,847 & 1 \end{pmatrix},$$

$$r_{Y, X}^2 = (0,2287, 0,4538) \begin{pmatrix} 1 & 0,847 \\ 0,847 & 1 \end{pmatrix}^{-1} \begin{pmatrix} 0,2287 \\ 0,4538 \end{pmatrix} = 0,2917$$

Interpretace: Variabilita výkonů dělníků je z 29% vysvětlena jejich věkem a dobou zapracovanosti.

$$r_{Y, X} = \sqrt{0,2917} = 0,5401.$$

Testujeme hypotézu  $H_0: \rho_{Y, X} = 0$  proti  $H_1: \rho_{Y, X} \neq 0$ .

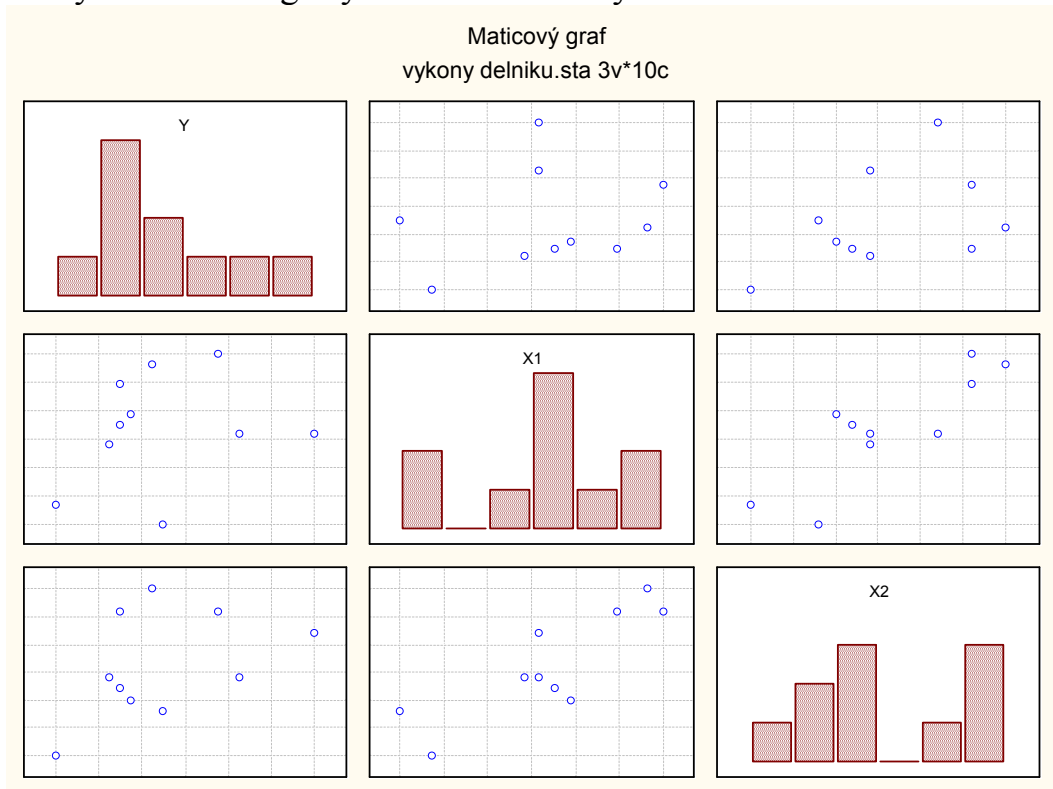
Testová statistika:  $F = \frac{10-2-1}{2} \cdot \frac{0,2917}{1-0,2917} = 1,441$ , kvantil  $F_{0,95}(2,7) = 4,737$ , kri-

tický obor  $W = \langle 4,737, \infty \rangle$ , tedy  $F \notin W$  a  $H_0$  nezamítáme na hladině významnosti 0,05.

### Výpočet pomocí systému STATISTICA:

Závislost mezi dvojicemi proměnných posoudíme pomocí dvourozměrných tečkových diagramů:

Grafy – Maticové grafy – Proměnné – Vybrat vše – OK



Ve všech třech případech existuje mezi dvojicemi proměnných určitý stupeň přímé lineární závislosti.

Statistiky – Vícenásobná regrese – Proměnné – Závislá proměnná Y, seznam nezáv. proměnných X1, X2 – OK – OK.

Koeficient  $r_{Y,(X_1,X_2)}$  najdeme v záhlaví výstupní tabulky pod označením Vícenás.

$R = 0,54$

Hodnota testové statistiky pro test nevýznamnosti koeficientu mnohonásobné korelace  $\rho_{Y,(X_1,X_2)}$  je 1,4411, počet stupňů volnosti čitatele je 2, jmenovatele 7, odpovídající p-hodnota je 0,2991, tedy na hladině významnosti 0,05 nezamítáme hypotézu, že výkon dělníka není závislý na jeho věku a době zapracovanosti.

### Výpočet pomocí systému SPSS

## Analyze – Regression – Linear – Dependent Y, Independent X1, X2 – OK

**Model Summary**

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,540 <sup>a</sup>	,292	,089	6,649

a. Predictors: (Constant), doba zapracovanosti, vek

**ANOVA<sup>b</sup>**

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	127,425	2	63,712	1,441	,299 <sup>a</sup>
	Residual	309,475	7	44,211		
	Total	436,900	9			

a. Predictors: (Constant), doba zapracovanosti, vek

b. Dependent Variable: vykon delnika

### Koeficient parciální korelace

Nechť  $Y, Z$  jsou náhodné veličiny a  $\mathbf{X} = (X_1, \dots, X_p)'$  je náhodný vektor. Korelační koeficient  $\rho(Y, Z)$  udává míru těsnosti lineárního vztahu mezi veličinami  $Y$  a  $Z$ . Ta však může být ovlivněna i tím, že mezi veličinami  $X_1, \dots, X_p$  existují veličiny, které silně korelují jak s  $Y$ , tak se  $Z$ . Zajímá nás proto, jaká je „čistá“ korelace mezi  $Y$  a  $Z$ , když se eliminuje vliv náhodného vektoru  $\mathbf{X}$ .

Pokud se omezíme na lineární vztahy, můžeme vliv vektoru  $\mathbf{X}$  na veličinu  $Y$  popsat lineární regresní funkcí

$$\hat{Y} = \alpha + \boldsymbol{\beta}'\mathbf{X}, \text{ kde } \boldsymbol{\beta} = \text{var}(\mathbf{X})^{-1} \text{cov}(\mathbf{X}, Y), \alpha = E(Y) - \boldsymbol{\beta}'E(\mathbf{X}).$$

Tu část veličiny  $Y$ , kterou vektor  $\mathbf{X}$  nevysvětlí, si můžeme představit jako reziduum  $Y - \hat{Y}$ . Analogicky pro veličinu  $Z$  dostáváme

$$\hat{Z} = \gamma + \boldsymbol{\delta}'\mathbf{X}, \text{ kde } \boldsymbol{\delta} = \text{var}(\mathbf{X})^{-1} \text{cov}(\mathbf{X}, Z), \gamma = E(Z) - \boldsymbol{\delta}'E(\mathbf{X}),$$

tudíž reziduum  $Z - \hat{Z}$  chápeme jako tu část veličiny  $Z$ , kterou vektor  $\mathbf{X}$  nevysvětlí.

Korelační koeficient mezi rezidui  $Y - \hat{Y}$  a  $Z - \hat{Z}$  se nazývá **parciální korelační koeficient** mezi náhodnými veličinami  $Y$  a  $Z$  při pevně daném vektoru  $\mathbf{X}$  a značí se  $\rho_{Y,Z,\mathbf{X}}$ . Tedy  $\rho_{Y,Z,\mathbf{X}} = \rho(Y - \hat{Y}, Z - \hat{Z})$ . Počítá se podle vzorce

$$\rho_{Y,Z,\mathbf{X}} = \frac{\rho(Y, Z) - \text{cov}(Y, \mathbf{X})\text{cor}(\mathbf{X})^{-1} \text{cov}(\mathbf{X}, Z)}{\sqrt{[1 - \text{cov}(Y, \mathbf{X})\text{cor}(\mathbf{X})^{-1} \text{cov}(\mathbf{X}, Y)][1 - \text{cov}(Z, \mathbf{X})\text{cor}(\mathbf{X})^{-1} \text{cov}(\mathbf{X}, Z)]}}.$$

Nechť náhodný vektor  $(Y, Z, X_1, \dots, X_p)'$  pochází z  $(p+2)$ -rozměrného rozložení, které má parciální korelační koeficient  $\rho_{Y,Z,\mathbf{X}}$ . Nechť je dán náhodný výběr

$(Y_1, Z_1, X_{11}, \dots, X_{1p})', \dots, (Y_n, Z_n, X_{n1}, \dots, X_{np})'$  rozsahu  $n$  z tohoto rozložení. Musí platit  $n > p+2$ . Jako odhad  $\rho_{Y,Z.X}$  slouží **výběrový parciální korelační koeficient**  $r_{Y,Z.X}$ :

$$r_{Y,Z.X} = \frac{r_{YZ} - \mathbf{S}_{YX} \mathbf{R}_{XX}^{-1} \mathbf{S}_{XZ}}{\sqrt{[1 - \mathbf{S}_{YX} \mathbf{R}_{XX}^{-1} \mathbf{S}_{XY}] [1 - \mathbf{S}_{ZX} \mathbf{R}_{XX}^{-1} \mathbf{S}_{XZ}]}}$$

## Testování hypotézy o nezávislosti veličin Y a Z při eliminaci vlivu vektoru X

### Popis testu

Budeme předpokládat, že uvedený náhodný výběr pochází z  $(p+2)$ -rozměrného normálního rozložení. Testujeme hypotézu  $H_0: \rho_{Y,Z.X} = 0$  proti  $H_1: \rho_{Y,Z.X} \neq 0$ . Vzhledem k tomu, že se jedná o výběr z normálního rozložení, testujeme, zda existuje závislost mezi Y a Z při eliminaci vlivu X.

Testová statistika  $T_0 = \frac{r_{Y,Z.X} \sqrt{n-p-2}}{\sqrt{1-r_{Y,Z.X}^2}}$  se řídí rozložením  $t(n-p-2)$ , pokud  $H_0$  platí.

Kritický obor:  $W = (-\infty, t_{1-\alpha/2}(n-p-2)) \cup (t_{1-\alpha/2}(n-p-2), \infty)$ . Jestliže  $T_0 \in W$ ,  $H_0$  zamítáme na hladině významnosti  $\alpha$

### Příklad

Pro data z příkladu o výkonnosti dělníků vypočtete výběrové parciální korelační koeficienty  $r_{Y,X_1,X_2}, r_{Y,X_2,X_1}$ , interpretujte je, porovnejte je s obyčejnými výběrovými korelačními koeficienty  $r_{YX_1}, r_{YX_2}$  a pro  $\alpha = 0,05$  otestujte významnost uvedených parciálních korelačních koeficientů.

### Řešení:

$$r_{YX_1} = 0,2287, r_{YX_2} = 0,4538, r_{X_1X_2} = 0,847$$

$$a) r_{Y,X_1,X_2} = \frac{r_{YX_1} - r_{YX_2} r_{X_1X_2}}{\sqrt{(1-r_{YX_2}^2)(1-r_{X_1X_2}^2)}} = -0,3286$$

**Interpretace:** Korelační koeficient mezi výkonem a věkem vyšel 0,2287, tedy s rostoucím věkem roste výkon. Parciální korelační koeficient mezi výkonem a věkem při vyloučení vlivu doby zapracovanosti vyšel -0,3286, tedy u dělníků se stejnou dobou zapracovanosti klesá s rostoucím věkem výkon.

Testujeme  $H_0: \rho_{Y,X_1,X_2} = 0$  proti  $H_1: \rho_{Y,X_1,X_2} \neq 0$

$$\text{Testová statistika: } T_0 = \frac{r_{Y,X_1,X_2} \sqrt{n-p-2}}{\sqrt{1-r_{Y,X_1,X_2}^2}} = \frac{-0,3286\sqrt{7}}{\sqrt{1-0,3286^2}} = -0,9205$$

Kvantil:  $t_{0,975}(7) = 2,3646$

Kritický obor:  $W = (-\infty, -2,3646) \cup (2,3646, \infty)$

$T_0 \notin W \Rightarrow H_0$  nezamítáme na hladině významnosti 0,05.

$$b) r_{Y, X_2, X_1} = \frac{r_{YX_2} - r_{YX_1} r_{X_1X_2}}{\sqrt{(1 - r_{YX_1}^2)(1 - r_{X_1X_2}^2)}} = 0,5026$$

**Interpretace:** Korelační koeficient mezi výkonem a dobou zapracovanosti vyšel 0,4538, tedy čím delší doba zapracovanosti, tím lepší výkon dělník podává. Parciální korelační koeficient mezi výkonem a dobou zapracovanosti při vyloučení vlivu věku vyšel 0,5026, tedy u stejně starých dělníků je poněkud silnější přímá lineární vazba mezi výkonem a dobou zapracovanosti.

Testujeme  $H_0: \rho_{Y, X_2, X_1} = 0$  proti  $H_1: \rho_{Y, X_2, X_1} \neq 0$

$$\text{Testová statistika: } T_0 = \frac{r_{Y, X_2, X_1} \sqrt{n - p - 2}}{\sqrt{1 - r_{Y, X_2, X_1}^2}} = \frac{0,5026 \sqrt{7}}{\sqrt{1 - 0,5026^2}} = 1,538$$

Kvantil:  $t_{0,975}(7) = 2,3646$

Kritický obor:  $W = (-\infty, -2,3646) \cup (2,3646, \infty)$

$T_0 \notin W \Rightarrow H_0$  nezamítáme na hladině významnosti 0,05.

### Výpočet pomocí systému STATISTICA

Statistiky – Základní statistiky/tabulky – Korelační matice – OK – na záložce Možnosti zaškrtneme Zobrazit r, úroveň p, počty N a zaškrtneme Zobrazit dlouhá jména proměnných, na záložce Detaily zvolíme Parciální korelace – 1. seznam proměnných Y, X1, druhý seznam proměnných X2 – OK

Proměnná	Parciální korelace (vykony delniku) Označ. korelace jsou významné na hlad. p < ,05000 N=10 (Celé případy vynechány u ChD)	
	Y	X1
Y: vykon delnika	1,0000	-,3286
	p= ---	p=,388
X1: vek	-,3286	1,0000
	p=,388	p= ---

Vidíme, že výběrový parciální korelační koeficient  $r_{Y, X_1, X_2}$  je -0,3286, p-hodnota je 0,388, tedy na hladině významnosti 0,05 nezamítáme hypotézu o nevýznamnosti  $\rho_{Y, X_1, X_2}$ .

Analogicky 1. seznam proměnných Y, X2, druhý seznam proměnných X1 – OK

	Parciální korelace (vykony delnika) Označ. korelace jsou významné na hlad. $p < ,05000$ N=10 (Celé případy vynechány u ChD)	
Proměnná	Y	X2
Y: výkon delnika	1,0000	,5026
	p= ---	p=,168
X2: doba zpracovanosti	,5026	1,0000
	p=,168	p= ---

V tomto případě výběrový parciální korelační koeficient  $r_{Y,X_2.X_1}$  je 0,5026, p-hodnota je 0,168, tedy na hladině významnosti 0,05 nezamítáme hypotézu o nevýznamnosti  $\rho_{Y,X_2.X_1}$ .

### Výpočet pomocí systému SPSS

Analyze – Correlate – Partial – Variables Y, X1 – Controlling for X2 – OK

#### Correlations

Control Variables			vykon delnika	vek
doba zpracovanosti	vykon delnika	Correlation	1,000	-,329
		Significance (2-tailed)	.	,388
		df	0	7
vek	vek	Correlation	-,329	1,000
		Significance (2-tailed)	,388	.
		df	7	0