

Osnova přednášky Korelační analýza

1. Jednoduchá korelace

- 1.1. Pearsonův koeficient korelace a jeho vlastnosti
- 1.2. Výběrový koeficient korelace
- 1.3. Test hypotézy o nezávislosti
- 1.4. Příklad
- 1.5. Interval spolehlivosti pro koeficient korelace
- 1.6. Příklad

2. Vícenásobná korelace

- 2.1. Varianční, korelační a kovarianční matice
- 2.2. Odhady pro jeden náhodný vektor
- 2.3. Příklad
- 2.4. Odhady pro dva náhodné vektory
- 2.5. Příklad
- 2.6. Koeficient vícenásobné korelace a jeho vlastnosti
- 2.7. Výběrový koeficient vícenásobné korelace
- 2.8. Test hypotézy o nevýznamnosti koeficientu vícenásobné korelace
- 2.9. Příklad

3. Parciální korelace

- 3.1. Koeficient parciální korelace
- 3.2. Výběrový koeficient parciální korelace
- 3.3. Test hypotézy o nevýznamnosti koeficientu parciální korelace
- 3.4. Příklad

1. Jednoduchá korelace

1.1. Pearsonův koeficient korelace

Definice: Necht' 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**.

Vlastnosti Pearsonova koeficientu korelace

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

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

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.

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 .)

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!

1.2. 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:

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

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

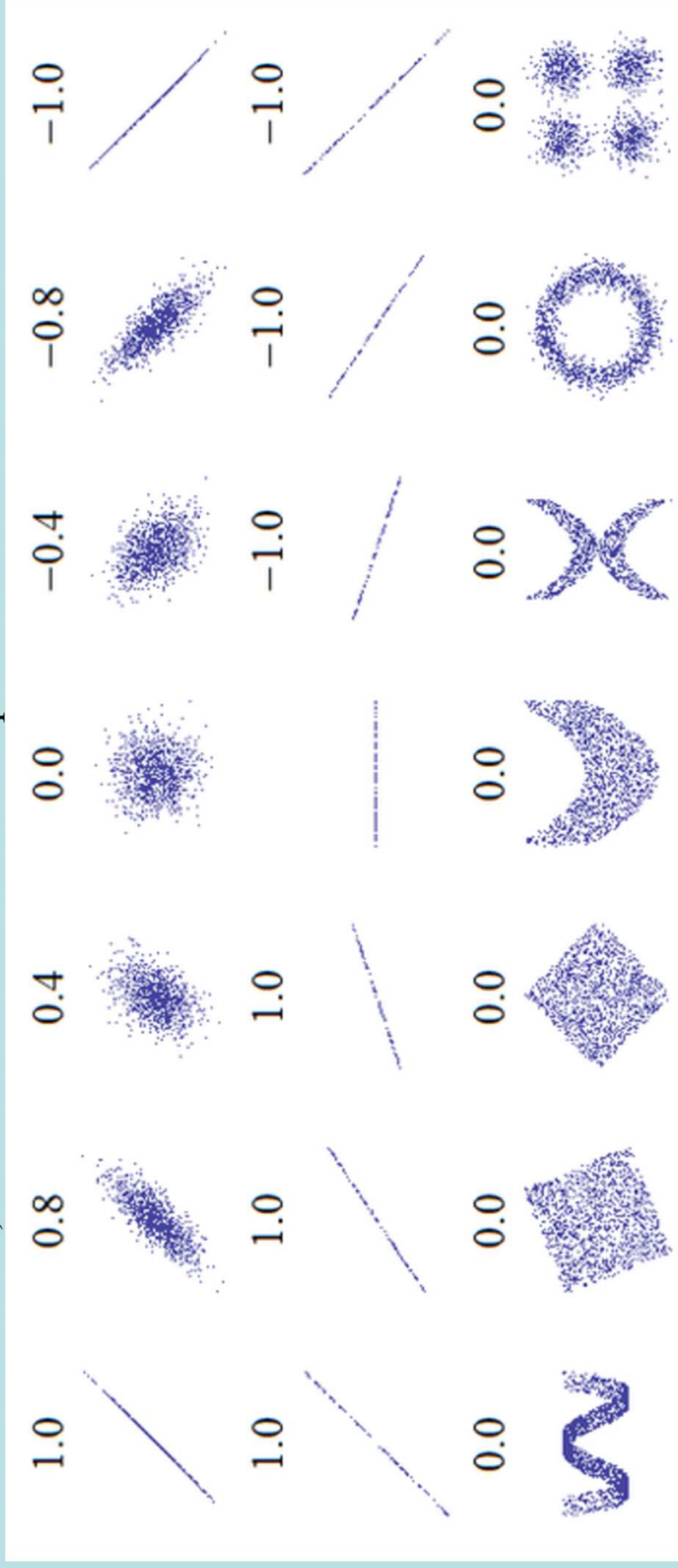
výběrovou kovarianci $S_{12} = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_1)(Y_i - M_2)$ a s jejich pomocí zavedeme

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

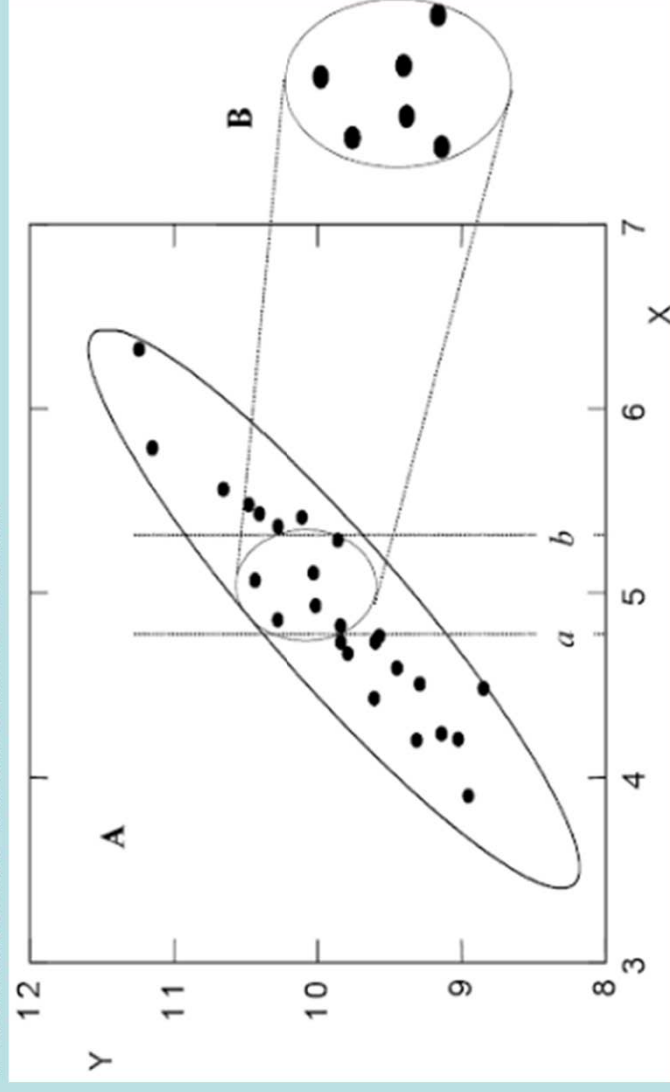
Vlastnosti Pearsonova koeficientu korelace se přenášejí i na výběrový koeficient korelace. (Výběrový koeficient korelace není nestranným odhadem skutečného koeficientu korelace, je odhadem vychýleným. Vychýlení je zanedbatelně malé pro rozsahy výběrů nad 30.)

Upozornění: nadále budeme předpokládat, že $(X_1, Y_1), \dots, (X_n, Y_n)$ je náhodný výběr rozsahu n z dvourozměrného normálního rozložení $N_2\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}\right)$.

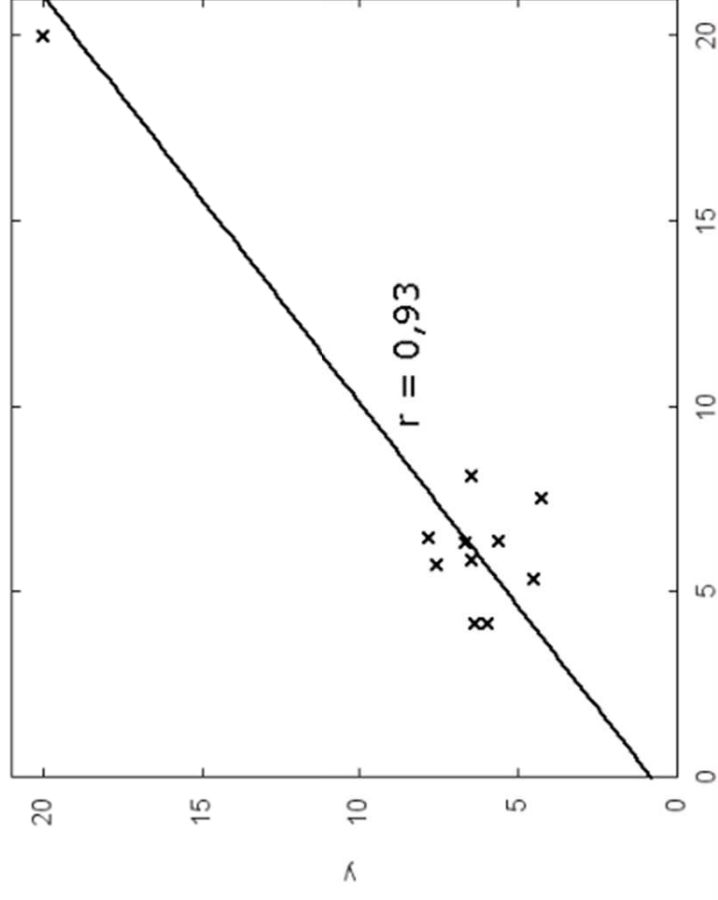
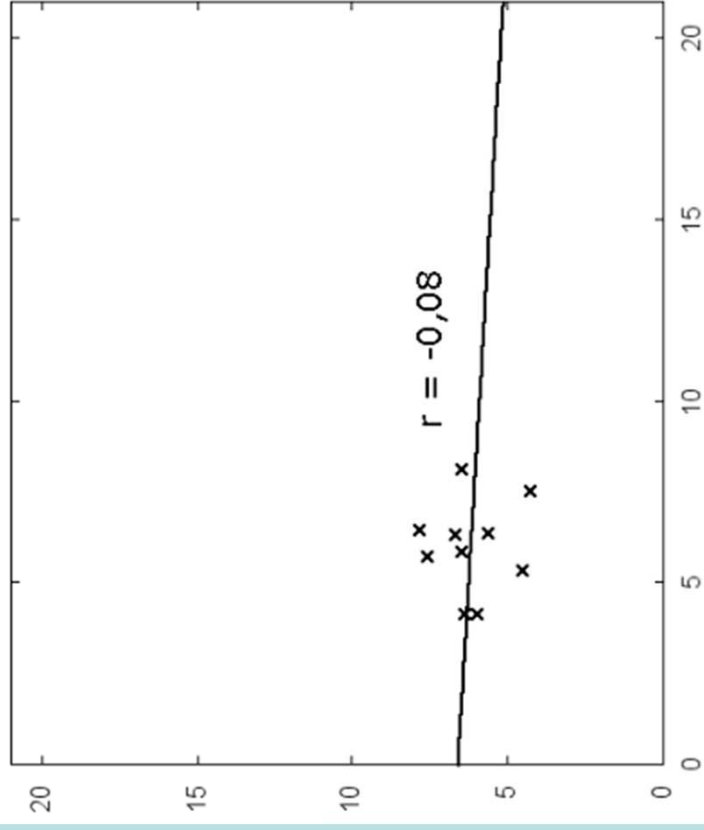
Ilustrace skutečnosti, že koeficient korelace se hodí pouze k měření lineární závislosti:



Ilustrace skutečnosti, že koeficient korelace se korektně používá tehdy, jsou-li data naměřena v dostatečně širokém rozsahu:



Ilustrace skutečnosti, že koeficient korelace je citlivý na vybočující pozorování:



1.3. Testování hypotézy o nezávislosti

Na hladině významnosti α 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 α , když $t_0 \in W$.

1.4. Příklad

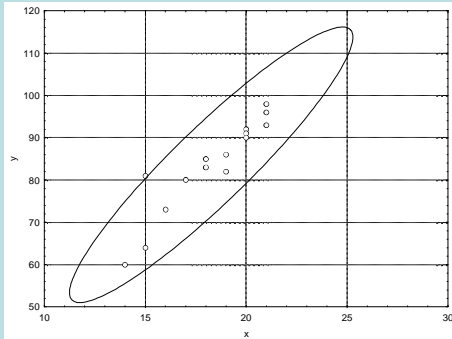
V dílně pracuje 15 dělníků. Byl u nich zjištěn počet směn odpracovaných za měsíc (náhodná veličina X) a počet zhotovených výrobků (náhodná veličina Y):

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, vypočítejte výběrový koeficient korelace mezi X a Y a na hladině 0,01 testujte hypotézu o nezávislosti X a Y.

Řešení: Dvourozměrnou normalitu dat ověříme pomocí dvourozměrného tečkového diagramu.



Vidíme, že předpoklad dvourozměrné normality je oprávněný.

Vypočteme realizace

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

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

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

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

Realizace testové statistiky: $t_0 = \frac{r_{12}\sqrt{n-2}}{\sqrt{1-r_{12}^2}} = \frac{0,927\sqrt{15-2}}{\sqrt{1-0,927^2}} = 8,912$,

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.

Počítačový výstup

Prom. X & prom. Y	Korelace (smeny a výrobky.sta) Označ. korelace jsou významné na hlad. p < ,05000 (Celé případy vynechány u ChD)										
	Průměr	Sm.Odch.	r(X,Y)	r2	t	p	N	Konst. záv.: Y	Směr. záv: Y	Konst. záv.: X	Směrníc záv.: X
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

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 v systému R:

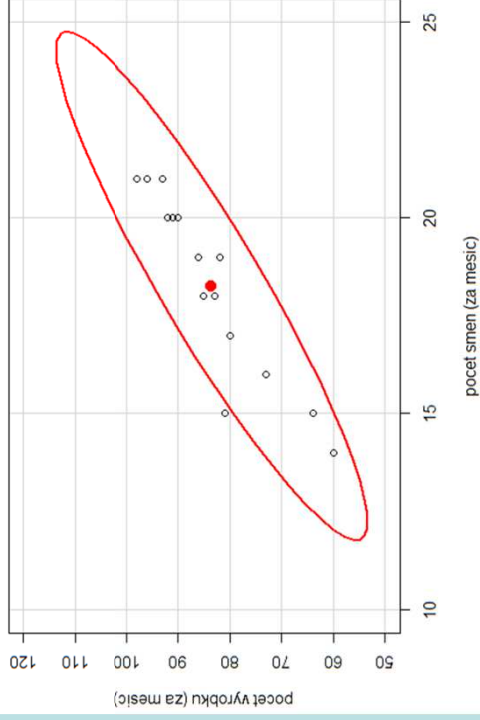
Načteme data a pojmenujeme proměnné:

```
data<-read.delim('smeny a vyrobky.txt', sep=',', header=T)
X<-data$X
Y<-data$Y
```

Načteme knihovnu car:

```
library(car)
```

Vykreslíme dvourozměrný tečkový diagram s 95% elipsou konstantní hustoty pravděpodobnosti:
dataEllipse(X, Y, level=0.95, xlab='pocet smen (za mesic)', ylab='pocet vyrobku (za mesic)', xlim=c(10,25), ylim=c(50,120))



Provedeme test vícerozměrné normality. Nejprve načteme knihovnu `mvnTest`:
`library(mvnTest)`

Provedeme Cramérův – von Misesův test a nakreslíme diagram na vizuální posouzení vícerozměrné normality:

```
CM.test(data, qqplot=T)
```

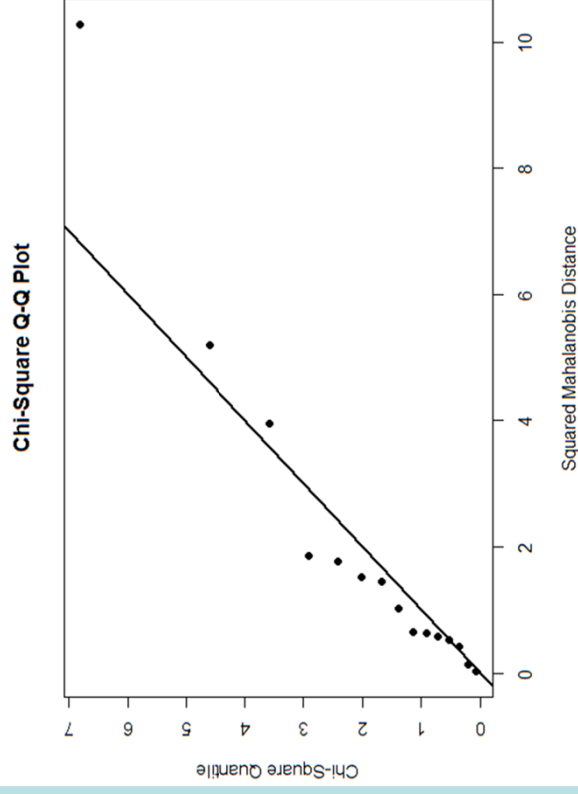
```
Cramer-von Mises test for Multivariate Normality
```

```
data : data
```

```
CM      : 0.1087666
```

```
p-value : 0.3034697
```

```
Result : Data are multivariate normal (sig.level = 0.05)
```



Vypočítáme výběrový koeficient korelace mezi počtem směn a počtem výrobků a provedeme test nezávislosti:

```
cor.test(X,Y)
```

```
Pearson's product-moment correlation
```

```
data: X and Y
```

```
t = 8.9238, df = 13, p-value = 6.629e-07
```

```
alternative hypothesis: true correlation is not equal to 0
```

```
95 percent confidence interval:
```

```
0.7902582 0.9759200
```

```
sample estimates:
```

```
cor
```

```
0.9271801
```

1.5. Interval spolehlivosti pro koeficient korelace

Náhodná veličina $Z = \frac{1}{2} \ln \frac{1+R_{12}}{1-R_{12}}$ (tzv. Fisherova Z-transformace koeficientu korelace) má přibližně nor-

mální 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

rozptylem $D(Z) = \frac{1}{n-3}$.

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 ρ pak dostaneme zpětnou transformací.

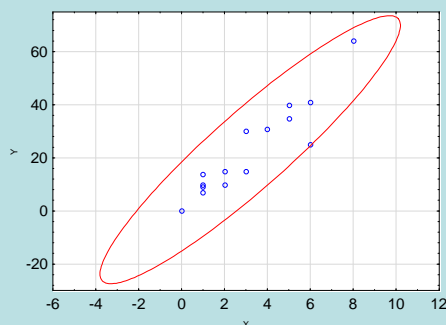
Poznámka: Jelikož $Z = \operatorname{arctgh} R_{12}$, dostáváme $R_{12} = \operatorname{tgh} Z$ a meze intervalu spolehlivosti pro ρ můžeme psát ve tvaru $\operatorname{tgh} \left(Z \pm \frac{u_{1-\alpha/2}}{\sqrt{n-3}} \right)$, přičemž $\operatorname{tgh} x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$.

1.6. Příklad: Učitel tělocviku zjišťoval, zda existuje vztah mezi počtem shybů (veličina X) a počtem kliků (veličina Y) u 15 náhodně vybraných chlapců:

Číslo chlapce	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Počet shybů	1	3	2	0	5	6	14	3	5	6	2	1	1	8	
Počet kliků	10	15	15	0	4	0	25	7	31	30	35	41	10	14	9

Za předpokladu, že uvedené údaje tvoří číselné realizace náhodného výběru rozsahu 15 z dvourozměrného normálního rozložení, vypočítejte 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 ρ .

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 $R_{12} = 0,9276$, tedy mezi počtem shybů a počtem kliků existuje silná přímá lineární závislost. Testová statistika: $T = 8,9511$, kvantil $t_{0,975}(13) = 2,1604$, kritický obor $W = (-\infty, -2,1604) \cup (2,1604, \infty)$. Jelikož $T \in W$, zamítáme na hladině významnosti 0,05 hypotézu o nezávislosti veličin X a Y.

Vypočítáme $Z = \frac{1}{2} \ln \frac{1+R_{12}}{1-R_{12}} = \frac{1}{2} \ln \frac{1+0,9276}{1-0,9276} = 1,6409$. Meze 95% asymptotického intervalu spolehlivosti pro ρ jsou

$\operatorname{tgh}\left(1,6409 \pm \frac{1,96}{\sqrt{12}}\right)$, tedy $0,7914 < \rho < 0,9761$ s pravděpodobností přibližně 0,95.

Počítačový výstup

	Odhad intervalu Jedna korelace, t-test
	Hodnota
Pozorovaný korel. koef. R	0,9276
Korelace dle nulové hypotézy (R ₀)	0,0000
Oboustranná p-hodnota	0,0000
Velikost vz. ve skup. (N)	15,0000
Interval spolehlivosti	0,9500
Meze spolehlivosti (Fisher. Z původní):	
R ₀ :	
Dolní mez	0,7914
Horní mez	0,9761

95% asymptotický interval spolehlivosti pro koeficient korelace ρ má tedy meze 0,7914 a 0,9761. (Protože nepokrývá hodnotu 0, zamítáme hypotézu o nezávislosti veličin X, Y na asymptotické hladině významnosti 0,05.)

2. Vícenásobná korelace

2.1. Varianční, korelační a kovarianční 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 (přitom $\sigma_{ii} = \sigma_i^2$)

$\rho_{ij} = R(X_i, X_j)$ 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} .

2.2. Odhady pro jeden náhodný vektor

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 $\mathbf{X}_1 = (X_{11}, \dots, X_{1p})'$, ..., $\mathbf{X}_n = (X_{n1}, \dots, X_{np})'$ 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** $\mathbf{S} = (S_{ij})_{i,j=1}^p$, což je symetrická matice řádu p .

Vychýlený odhad matice $\text{cor}(\mathbf{X})$ je **výběrová korelační matice** $\mathbf{R} = (R_{ij})_{i,j=1}^p$ 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 jedničky a matice \mathbf{R} je symetrická.)

2.3. Příklad: U 28 náhodně vybraných osob byly zjišťovány tyto údaje:

Sex ... 1 – muž, 2 – žena (mužů i žen bylo po 14)

výška (v cm), hmotnost (v kg), boty (číslo bot).

Vypočtěte realizaci výběrové varianční matice a výběrové korelační matice. (Soubor udaje_o_lidech_1.sta)

Řešení:

Výběrová varianční matice Výběrová korelační matice

Proměnná	vyska	hmotnost	boty
vyska	112,8611	161,0926	41,45370
hmotnost	161,0926	248,4709	61,99206
boty	41,4537	61,9921	16,40608

Proměnná	vyska	hmotnost	boty
vyska	1,000000	0,961979	0,963360
hmotnost	0,961979	1,000000	0,970948
boty	0,963360	0,970948	1,000000

Z výběrové varianční matice plyne, že největší variabilitu má hmotnost, pak výška a nakonec číslo bot.

Z výběrové korelační matice plyne, že mezi všemi třemi dvojicemi proměnných existuje velmi silná přímá lineární závislost, nejsilnější je mezi hmotností a velikostí bot.

Výpočet v systému R:

Načteme data:

```
data<-read.table('udaje_o_lidech_1.txt',sep=',',header=T)
```

Vypočteme výběrovou varianční matici:

```
cov(data[,2:4], method='pearson')
```

	vyska	hmotnost	boty
vyska	112.8611	161.09259	41.45370
hmotnost	161.0926	248.47090	61.99206
boty	41.4537	61.99206	16.40608

Vypočteme výběrovou korelační matici:

```
cor(data[,2:4], method='pearson')
```

	vyska	hmotnost	boty
vyska	1.0000000	0.9619791	0.9633602
hmotnost	0.9619791	1.0000000	0.9709477
boty	0.9633602	0.9709477	1.0000000

2.4. Odhady pro dva náhodné vektory

Nechť náhodný vektor \mathbf{X} má p -rozměrné rozložení a necht' $\mathbf{X}_1, \dots, \mathbf{X}_n$ je náhodný výběr z tohoto rozložení. Necht' náhodný vektor \mathbf{Y} má q -rozměrné rozložení a necht' $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ je náhodný výběr z tohoto rozložení. Předpokládejme, že obě rozložení mají konečné druhé momenty. Necht' $\text{cov}(\mathbf{X}, \mathbf{Y})$ 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 $\mathbf{S}_{XY} = (S_{ij}),$ kde S_{ij} je výběrová kovariance i -té a j -té složky vektorů $\mathbf{X}, \mathbf{Y}, i = 1, \dots, p, j = 1, \dots, q.$

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 $\mathbf{X}, \mathbf{Y}, i = 1, \dots, p, j = 1, \dots, q.$

2.5.Příklad: Necht' vektor $\mathbf{X} = (X_1, X_2, X_3)'$ obsahuje údaje o výšce, hmotnosti a číslu bot mužů, vektor $\mathbf{Y} = (Y_1, Y_2)'$ obsahuje údaje výšce a hmotnosti žen. Vypočtete realizace výběrové kovarianční a výběrové korelační matice vektorů \mathbf{X} , \mathbf{Y} . (Soubor udaje_o_lidech_2.sta)

Řešení:

Výběrová kovarianční matice Výběrová korelační matice

Efekt	Sloup.4 Vyska_z	Sloup.5 Hmotnost_z
Vyska_m	10,81319	17,39560
Hmotnost_m	15,70879	15,22527
Boty_m	4,43407	5,13736

Efekt	Sloup.4 Vyska_z	Sloup.5 Hmotnost_z
Vyska_m	0,467318	0,767160
Hmotnost_m	0,514047	0,508409
Boty_m	0,560289	0,662427

Výpočet v systému R:

Načteme data:

```
data<-read.table('udaje_o_lidech_2.txt',sep=',',header=T)
```

Vypočteme výběrovou varianční matici všech proměnných a uložíme ji do objektu S:

```
S<-cov(data, method='pearson')
```

```
S
      Vyska_m Hmotnost_m      Boty_m  Vyska_z Hmotnost_z      Boty_z
Vyska_m  28.417582 30.637363  8.4505495 10.813187 17.395604  3.2857143
Hmotnost_m 30.637363 49.565934 10.3131868 15.708791 15.225275  4.4725275
Boty_m      8.450549 10.313187  3.3241758  4.434066  5.137363  0.8021978
Vyska_z    10.813187 15.708791  4.4340659 18.840659 14.258242  3.7912088
Hmotnost_z 17.395604 15.225275  5.1373626 14.258242 18.093407  4.1208791
Boty_z      3.285714  4.472527  0.8021978  3.791209  4.120879  1.6043956
```

Z matice S vybereme příslušné řádky a sloupce:

```
S<-S[1:3,4:5]
```

```
> S
```

```
      Vyska_z Hmotnost_z
Vyska_m 10.813187 17.395604
Hmotnost_m 15.708791 15.225275
Boty_m      4.434066  5.137363
```

Vypočteme výběrovou korelační matici všech proměnných a uložíme ji do objektu R:

```
R<-cor(data, method='pearson')
```

```
R
```

		Vyska_m	Hmotnost_m	Boty_m	Vyska_z	Hmotnost_z	Boty_z
Vyska_m	1.0000000	0.8163308	0.8694596	0.4673176	0.7671597	0.4866094	
Hmotnost_m	0.8163308	1.0000000	0.8034507	0.5140473	0.5084091	0.5015401	
Boty_m	0.8694596	0.8034507	1.0000000	0.5602890	0.6624270	0.3473631	
Vyska_z	0.4673176	0.5140473	0.5602890	1.0000000	0.7722497	0.6895630	
Hmotnost_z	0.7671597	0.5084091	0.6624270	0.7722497	1.0000000	0.7648460	
Boty_z	0.4866094	0.5015401	0.3473631	0.6895630	0.7648460	1.0000000	

Z matice R vybereme příslušné řádky a sloupce:

```
R<-R[1:3,4:5]
```

```
R
```

	Vyska_z	Hmotnost_z
Vyska_m	0.4673176	0.7671597
Hmotnost_m	0.5140473	0.5084091
Boty_m	0.5602890	0.6624270

2.6. Koeficient vícenásobné korelace

Intenzitu lineární závislosti mezi náhodnou veličinou Y a náhodným vektorem $\mathbf{X} = (X_1, \dots, X_p)'$ měříme pomocí **koeficientu vícenásobné 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:

a) $\rho_{Y, \mathbf{X}} \geq 0$

b) $\rho_{Y, \mathbf{X}} \geq |\rho(Y, X_i)|$ pro $\forall i = 1, \dots, p$

c) $\rho_{Y, X_1 \dots X_p} \geq \dots \geq \rho_{Y, X_1 X_2} \geq \rho(Y, X_1)$

d) $\rho_{Y, \mathbf{X}} = 1 \Leftrightarrow$ existují konstanty $\beta_0, \beta_1, \dots, \beta_p$ tak, že $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$.

2.7. Výběrový koeficient vícenásobné korelace

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 vícenásobné korelace** $r_{Y, \mathbf{X}}$, jehož druhá mocnina je dána vzorcem

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

kde $\mathbf{R}_{Y\mathbf{X}}$ 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 koeficientu vícenásobné korelace se přenášejí i na výběrový koeficient vícenásobné korelace.

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

Vypočtete výběrový koeficient vícenásobné korelace $r_{Y,(X_1,X_2)}$ popisující závislost hodinové výkonnosti dělníka na jeho věku a době zapracovanosti.

Řešení:

Koeficient $r_{Y,(X_1,X_2)}$ najdeme v záhlaví výstupní tabulky pod označením $R = 0,54$.

Výsledky regrese se závislou proměnnou : Y (vykony delniku.sta)						
R= ,54005243 R2= ,29165662 Upravené R2= ,08927280						
F(2,7)=1,4411 p<,29913 Směrod. chyba odhadu : 6,6491						
N=10	b*	Sm.chyba z b*	b	Sm.chyba z b	t(7)	p-hodn.
Abs.člen			86,74217	25,32397	3,425299	0,011056
X1	-0,550937	0,598452	-0,70031	0,76071	-0,920604	0,387883
X2	0,920415	0,598452	1,35062	0,87817	1,537994	0,167937

Jeho druhá mocnina (ozn. R^2) nám říká, že variabilita výkonů dělníků je z 29 % vysvětlena jejich věkem a dobou zapracovanosti.

2.8. Test hypotézy o nevýznamnosti koeficientu vícenásobné korelace

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 vícenásobné 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 \mathbf{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 α .

2.9. Příklad

Předpokládáme, že údaje o výkonnosti 10 náhodně vybraných dělníků, jejich věku a době zapracovanosti představují číselné realizace náhodného výběru rozsahu 10 ze třírozměrného normálního rozložení. Na hladině významnosti 0,05 testujte hypotézu, že výkon dělníka nezávisí na jeho věku a době zapracovanosti.

Řešení:

Výsledky regrese se závislou proměnnou : Y (vykony delniku.sta) R= ,54005243 R2= ,29165662 Upravené R2= ,08927280 F(2,7)=1,4411 p<,29913 Směrod. chyba odhadu : 6,6491						
N=10	b*	Sm.chyba z b*	b	Sm.chyba z b	t(7)	p-hodn.
Abs.člen			86,74217	25,32397	3,425299	0,011056
X1	-0,550937	0,598452	-0,70031	0,76071	-0,920604	0,387883
X2	0,920415	0,598452	1,35062	0,87817	1,537994	0,167937

Hodnota testové statistiky pro test nevýznamnosti koeficientu vícenásobné 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 v systému R:

Načteme data a pojmenujeme proměnné:

```
data<-read.delim('vykony delniku.txt',sep=' ', header=T)
Y<-data$Y
X1<-data$X1
X2<-data$X2
```

Vytvoříme model lineární regrese a vypíšeme jeho výstup:

```
vystup<-lm(Y~X1+X2)
summary(vystup)
Call:
lm(formula = Y ~ X1 + X2)

Residuals:
    Min       1Q   Median       3Q      Max
-7.4367 -4.4717 -0.9345  3.3150  9.7630

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  86.7422    25.3240   3.425  0.0111 *
X1          -0.7003     0.7607  -0.921  0.3879
X2           1.3506     0.8782   1.538  0.1679
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.649 on 7 degrees of freedom
Multiple R-squared:  0.2917, Adjusted R-squared:  0.08927
F-statistic: 1.441 on 2 and 7 DF, p-value: 0.2991
```


3. Parciální korelace

3.1. Koeficient parciální korelace

Nechť Y, Z jsou náhodné veličiny a $\mathbf{X} = (X_1, \dots, X_p)'$ je náhodný vektor. Koeficient korelace $\rho(Y, Z)$ 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 regresi funkcí

$$\hat{Y} = b_0 + b_1 X_1 + \dots + b_p X_p.$$

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

$$\hat{Z} = a_0 + a_1 X_1 + \dots + a_p X_p,$$

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

Koeficient korelace 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 $\rho_{Y,Z,X}$.

Tedy $\rho_{Y,Z,X} = \rho(Y - \hat{Y}, Z - \hat{Z})$. Počítá se podle vzorce

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

3.2. Výběrový koeficient parciální korelace

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

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**

$r_{Y,Z.X}$:

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

3.3. Test hypotézy o nevýznamnosti koeficientu parciální korelace

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 \mathbf{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 α .

3.4. Příklad

Pro data z příkladu o výkonnosti dělníků vypočítejte 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:

Nejprve vypočteme párový koeficient korelace mezi výkonem a věkem.

Proměnná	X1
Y	0,2287

Dále vypočteme parciální korelační koeficient mezi výkonem a věkem při vyloučení vlivu doby zapracovanosti a otestujeme jeho významnost.

Proměnná	Y	X1
Y	1,0000	-,3286
	p= ---	p=,388
X1	-,3286	1,0000
	p=,388	p= ---

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.

Odpovídající p-hodnota je 0,388, tedy na hladině významnosti 0,05 nezamítáme hypotézu o nevýznamnosti

ρ_{Y,X_1,X_2} .

Nyní vypočteme koeficient korelace mezi výkonem a dobou zapracovanosti:

Proměnná	X2
Y	0,4538

Dále vypočteme parciální korelační koeficient mezi výkonem a dobou zapracovanosti při vyloučení vlivu věku pracovníka a otestujeme jeho významnost.

Proměnná	Y	X2
Y	1,0000	,5026
	p= ---	p=,168
X2	,5026	1,0000
	p=,168	p= ---

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.

Odpovídající p-hodnota je 0,168, tedy na hladině významnosti 0,05 nezamítáme hypotézu o nevýznamnosti $\rho_{Y, X_2 \cdot X_1}$.

Výpočet v systému R:

Vypočteme nejprve párové korelační koeficienty:

```
cor(data)
```

```
      Y      x1      x2
Y  1.000000 0.228680 0.4537570
x1 0.228680 1.000000 0.8470271
x2 0.453757 0.8470271 1.0000000
```

Načteme knihovnu ppcor:

```
library(ppcor)
```

Vypočteme parciální korelační koeficienty:

```
pcor(data)
```

```
$estimate (výběrové parciální koeficienty korelace)
```

```
      Y      x1      x2
Y  1.0000000 -0.3286296 0.5025635
x1 -0.3286296 1.0000000 0.8567741
x2 0.5025635 0.8567741 1.0000000
```

```
$p.value (p-hodnoty pro test nevýznamnosti parciálních korelačních koeficientů)
```

```
      Y      x1      x2
Y  0.0000000 0.38788283 0.16793749
x1 0.3878828 0.00000000 0.00317348
x2 0.1679375 0.00317348 0.00000000
```

```
$statistic (testové statistiky pro test nevýznamnosti parciálních korelačních koeficientů)
```

```
      Y      x1      x2
Y  0.0000000 -0.9206035 1.537994
x1 -0.9206035 0.0000000 4.395670
x2 1.5379945 4.3956702 0.000000
```