

## Průzkumová analýza vícerozměrných dat

### Osnova:

- vícerozměrný datový soubor
- vizualizace vícerozměrných dat
- snížení dimenze dat metodou hlavních komponent
- shluková analýza

**Vícerozměrná data:** vyskytují se v situacích, kdy u každého z  $n$  objektů zjišťujeme hodnoty  $p$  znaků  $X_1, \dots, X_p$ .  
**p-rozměrný datový soubor:** matice  $n \times p$ :

$$\begin{pmatrix} X_{11} & \cdots & X_{1p} \\ \cdots & \cdots & \cdots \\ X_{n1} & \cdots & X_{np} \end{pmatrix}.$$

Řádky charakterizují objekty, sloupce znaky.

Např. máme  $n$  sportovců, u každého sledujeme tyto znaky: pohlaví (0 – žena, 1 – muž), tělesná výška (v cm), tělesná hmotnost (v kg), nejlepší výkon ve skoku do dálky (v cm), nejlepší výkon ve skoku do výšky (v cm), nejlepší výkon v běhu na 100 m (v s).

Úkoly průzkumové analýzy vícerozměrných dat:

- odhalit vektory pozorování nebo jejich složky, které se jeví jako vybočující
- postihnout závislosti mezi sloupci datového souboru
- identifikovat shluky v datech, které svědčí o nehomogenitě daného výběru
- posoudit vícerozměrnou normalitu dat.

Omezíme se na dva problémy, a to na vizualizaci dat pomocí hlavních komponent a na shlukovou analýzu dat.

## Vizualizace vícerozměrných dat

Je-li  $p = 2$  nebo  $p = 3$ , můžeme hodnoty znaků chápat jako souřadnice v dvou či třírozměrném prostoru a získáme tak dvourozměrný či třírozměrný tečkový diagram. Ze vzhledu těchto tečkových diagramů lze poznat, zda se v datech vyskytují odlehlá pozorování, zda mezi znaky existuje nějaká závislost nebo zda se objekty sdružují do skupin.

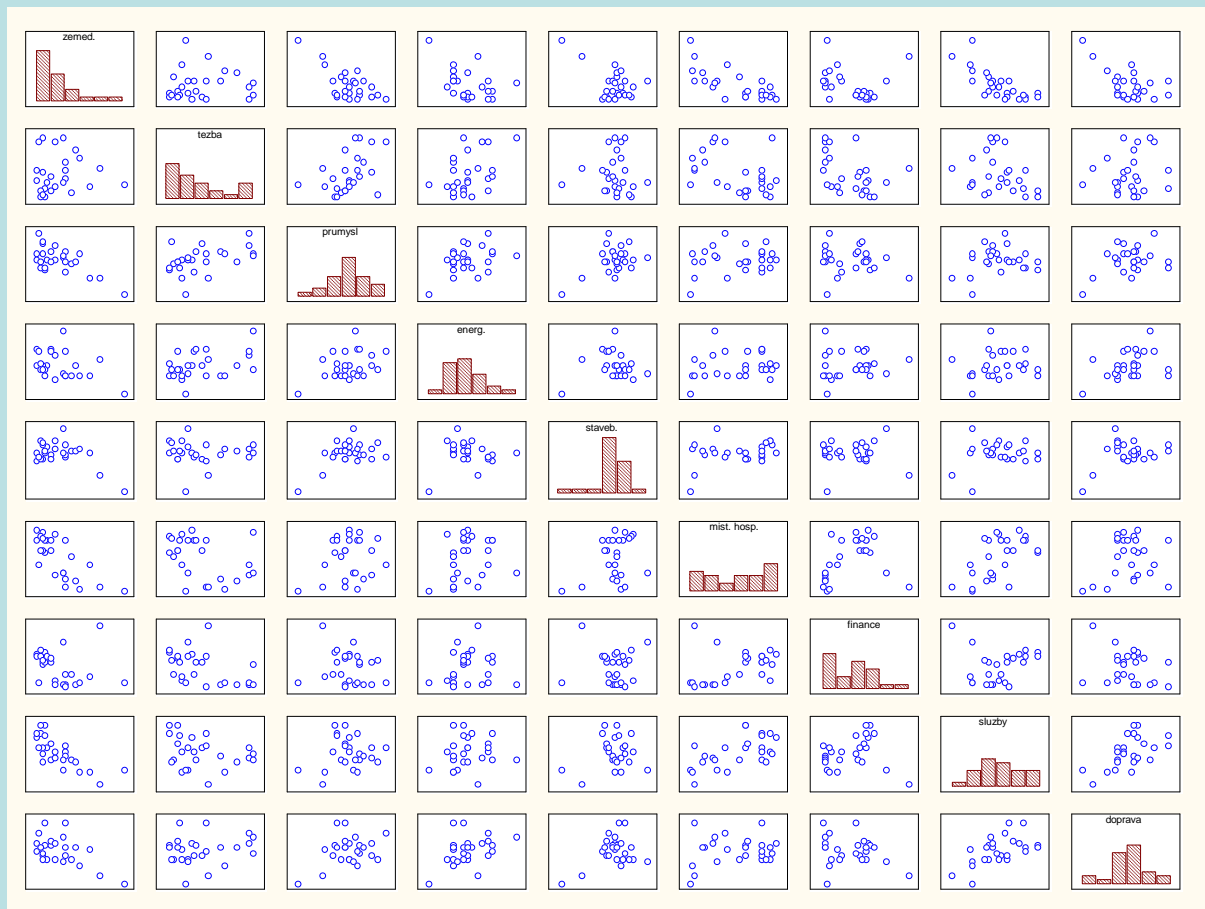
**Příklad:** Máme k dispozici datový soubor z roku 1979 o 26 evropských zemích, který obsahuje údaje o procentuálním zastoupení ekonomicky činného obyvatelstva v různých odvětvích národního hospodářství: zemědělství, těžba, průmyslová výroba, energetika, stavebnictví, místní hospodářství, finanční sektor, služby, doprava a komunikace.

	1	2	3	4	5	6	7	8	9
	zemed.	tezba	prumysl	energ.	staveb.	míst. hosp.	finance	sluzby	doprava
Belgie	3,3	0,9	27,6	0,9	8,2	19,1	6,2	26,6	7,2
Dánsko	9,2	0,1	21,8	0,6	8,3	14,2	6,5	32,2	7,1
Francie	10,8	0,8	27,5	0,9	8,9	16,8	6	22,6	5,7
Záp. Německo	6,7	1,3	35,8	0,9	7,3	14,4	5	22,5	6,1
Irsko	23,2	1	20,7	1,3	7,5	16,8	2,8	20,6	6,1
Itálie	15,9	0,6	27,6	0,5	10	18,1	1,5	20,1	5,7
Lucembursko	7,7	3,1	30,8	0,8	9,2	18,5	4,5	19,2	6,2
Nizozemsko	6,3	0,1	22,5	1	9,9	18	6,9	28,5	6,8
Velká Británie	2,7	1,4	30,2	1,4	6,9	16,9	5,8	28,3	6,4
Rakousko	12,7	1,1	31,4	1,4	8	16,8	4,9	16,7	7
Finsko	13	0,4	25,9	1,3	7,4	14,7	5,5	24,2	7,6
Řecko	41,4	0,6	17,6	0,6	8,1	11,5	2,4	11,1	6,7
Norsko	9	0,5	22,4	0,8	8,6	16,9	4,7	27,7	9,4
Portugalsko	27,8	0,3	24,5	0,6	8,4	13,3	2,7	16,7	5,7
Španělsko	22,9	0,8	28,5	0,7	11,5	9,7	8,5	11,9	5,5
Švédsko	6,1	0,4	25,9	0,8	7,2	14,4	6	32,4	6,8
Švýcarsko	7,7	0,2	37,8	0,8	9,5	17,5	5,3	15,5	5,7
Turecko	66,8	0,7	7,9	0,1	2,8	5,5	1,1	11,9	3,2
Bulharsko	23,6	1,9	32,3	0,6	7,9	8	0,7	18,2	6,8
Československo	16,5	2,9	35,5	1,2	8,7	9,2	0,9	17,9	7,2
Vých. Německo	4,2	2,9	41,2	1,3	7,6	11,2	1,2	22,1	8,3
Maďarsko	21,7	3,1	29,6	1,9	8,2	9,4	0,9	17,2	8
Polsko	31,1	2,5	25,7	0,9	8,4	7,5	0,9	16,1	6,9
Rumunsko	34,7	2,1	30,1	0,6	8,7	5,9	1,3	11,6	5
Sovětský svaz	23,7	1,4	25,8	0,6	9,2	6,1	0,5	23,4	9,3
Jugoslávie	48,7	1,5	16,8	1,1	4,9	6,4	11,3	5,3	4

Vytvořte dvourozměrné tečkové diagramy pro všechny dvojice proměnných.

## Řešení pomocí systému STATISTICA:

Grafy – Maticové grafy – Proměnné – Vybrat vše – OK.



Na hlavní diagonále maticového grafu jsou histogramy jednotlivých proměnných, mimo hlavní diagonálu jsou dvourozměrné tečkové diagramy odpovídajících dvojic proměnných. Vidíme např., že podíl obyvatel zaměstnaných v zemědělství záporně koreluje s podílem obyvatel zaměstnaných v průmyslu, službách či dopravě.

Je-li  $p > 3$ , použijeme k vizualizaci dat **metodu hlavních komponent (principal component analysis)**, která umožňuje vyjádřit informace o variabilitě obsažené v datovém souboru pomocí několika málo nových znaků  $Y_1, \dots, Y_m$  získaných jako lineární kombinace znaků původních  $X_1, \dots, X_p$ ,  $m < p$  :

$$Y_1 = v_{11}X_1 + \dots + v_{1p}X_p,$$

$$Y_2 = v_{21}X_1 + \dots + v_{2p}X_p.$$

·  
·  
·

$$Y_m = v_{m1}X_1 + \dots + v_{mp}X_p.$$

Tyto nové znaky, kterým se říká hlavní komponenty, jsou

- nekorelované,
- uspořádané podle svého klesajícího rozptylu.

Většina informace o variabilitě původních dat je tedy soustředěna v první hlavní komponentě a nejméně informace je obsaženo v poslední hlavní komponentě. Ukazuje se, že pouze několik prvních hlavních komponent má dostatečně velký rozptyl. Ostatní pak můžeme zanedbat, čímž docílíme snížení dimenze dat. V datovém souboru však musí existovat mezi znaky dostatečně silná korelace, aby bylo možno tuto redukci provést.

Analýza hlavních komponent může být chápána jako transformace z původního do nového souřadnicového systému, jehož osy jsou tvořeny hlavními komponentami. Osy procházejí směry maximálního rozptylu, protože podmínka nezávislosti komponent vede ke kolmosti os.

Data pak znázorníme v prostoru prvních dvou či tří hlavních komponent.

Metodu hlavních komponent (Principal Component Analysis – PCA) popsal v r. 1901 Karl Pearson a ve 30. letech 20. století ji dále rozvinul Harold Hotelling.



Harold Hotelling (1895 – 1973), americký matematik a statistik

### Podstata metody hlavních komponent

Uvažme datový soubor, který vznikl tak, že 6 žáků absolvovalo 4 testy, které měří následující veličiny:

$X_1$  – přírodovědné znalosti,

$X_2$  – literární vědomosti,

$X_3$  – schopnost koncentrace,

$X_4$  – logické myšlení.

Testy se hodnotí na škále od 1 do 10 (1 = špatný výsledek, 10 = výborný výsledek)

	1 X1	2 X2	3 X3	4 X4
1	7	9	10	8
2	9	8	8	10
3	4	3	1	2
4	2	3	2	2
5	3	1	2	4
6	1	1	1	4

## Označení

$\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$  – vektor pozorování  $i$ -tého objektu,  $i = 1, 2, \dots, n$

Např. pro  $i = 3$  máme  $\mathbf{x}_3 = (4 \ 3 \ 1 \ 2)^T$

$m_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$  - průměr  $j$ -tého znaku,  $j = 1, 2, \dots, p$ . Např. pro  $j = 1$  máme  $m_1 = \frac{1}{6}(7 + 9 + 4 + 2 + 3 + 1) = 4,3\bar{3}$

$s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - m_j)^2$  - rozptyl  $j$ -tého znaku,  $j = 1, 2, \dots, p$ . Např. pro  $j = 1$  máme  $s_1^2 = \frac{1}{5}[(7 - 4,3\bar{3})^2 + \dots + (1 - 4,3\bar{3})^2] = 9,4\bar{6}$

Datový soubor s průměry, směrodatnými odchylkami a rozptyly:

	1 X1	2 X2	3 X3	4 X4
1	7	9	10	8
2	9	8	8	10
3	4	3	1	2
4	2	3	2	2
5	3	1	2	4
6	1	1	1	4
průměry	4,33	4,17	4,00	5,00
s.o.	3,08	3,49	3,95	3,29
rozptyly	9,47	12,17	15,60	10,80

$z_{ij} = \frac{x_{ij} - m_j}{s_j}$  -  $(i,j)$ -tá standardizovaná hodnota,  $i = 1, 2, \dots, n$ ,  $j = 1, 2, \dots, p$

Např. pro  $i = 1, j = 1$  máme  $z_{11} = \frac{7 - 4,3\bar{3}}{\sqrt{9,4\bar{6}}} = 0,8667$

### Datový soubor standardizovaných hodnot

	1 X1	2 X2	3 X3	4 X4
1	0,866703	1,385674	1,519109	0,912871
2	1,51673	1,098983	1,012739	1,521452
3	-0,10834	-0,33447	-0,75955	-0,91287
4	-0,75836	-0,33447	-0,50637	-0,91287
5	-0,43335	-0,90786	-0,50637	-0,30429
6	-1,08338	-0,90786	-0,75955	-0,30429

$\mathbf{z}_i = (z_{i1}, \dots, z_{ip})^T$  – vektor standardizovaných pozorování i-tého objektu,  $i = 1, 2, \dots, n$

$\mathbf{m} = (m_1, \dots, m_p)^T$  – vektor průměrů

$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{m})(\mathbf{x}_i - \mathbf{m})^T$  - výběrová varianční matice. V našem případě:

Proměnná	Kovariance (pca)			
	X1	X2	X3	X4
X1	9,46667	9,73333	10,60000	8,80000
X2	9,73333	12,16667	13,20000	9,40000
X3	10,60000	13,20000	15,60000	11,60000
X4	8,80000	9,40000	11,60000	10,80000

$\mathbf{R} = \frac{1}{n-1} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T$  - výběrová korelační matice. V našem případě:

Proměnná	Korelace (pca)			
	X1	X2	X3	X4
X1	1,000000	0,906937	0,872258	0,870307
X2	0,906937	1,000000	0,958133	0,820031
X3	0,872258	0,958133	1,000000	0,893684
X4	0,870307	0,820031	0,893684	1,000000

(**S** a **R** jsou čtvercové symetrické matice řádu  $p$ .)



## Základní pojmy

$\mathbf{A}$  - čtvercová matice řádu  $p$ .

**Vlastní číslo matice  $\mathbf{A}$**  – takové číslo  $\lambda$ , které pro libovolný nenulový vektor  $\mathbf{v}$  typu  $p \times 1$  splňuje rovnici  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ .

**Vlastní vektor matice  $\mathbf{A}$**  – vektor  $\mathbf{v}$ .

**Charakteristický polynom matice  $\mathbf{A}$**  - determinant  $|\mathbf{A} - \lambda\mathbf{I}|$ .

**Stopa matice  $\mathbf{A}$**  - součet jejích diagonálních prvků (značí se  $\text{Tr}(\mathbf{A})$ ).

## Výpočet vlastních čísel matice $\mathbf{A}$

Rovnici  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$  upravíme na tvar  $(\mathbf{A} - \lambda\mathbf{I})\mathbf{v} = \mathbf{0}$ . Tato soustava  $p$  rovnic má netriviální řešení, právě když charakteristický polynom matice  $\mathbf{A}$  je roven 0. Dostaneme rovnici  $p$ -tého stupně. Jejím řešením jsou vlastní čísla  $\lambda_1, \dots, \lambda_p$ .

## Vlastnosti vlastních čísel

Jejich součet je roven stopě matice  $\mathbf{A}$ :  $\lambda_1 + \dots + \lambda_p = \text{Tr}(\mathbf{A})$ ,

jejich součin je roven determinantu matice  $\mathbf{A}$ :  $\lambda_1 \dots \lambda_p = \det(\mathbf{A})$ ,

jsou seřazena sestupně:  $\lambda_1 \geq \dots \geq \lambda_p$ .

## Vlastnosti vlastních vektorů

Mají jednotkovou délku:  $\mathbf{v}_i^T \mathbf{v}_i = 1$ ,  $i = 1, \dots, p$ ,

jsou vzájemně ortogonální:  $\mathbf{v}_i^T \mathbf{v}_j = 0$  pro všechna  $i \neq j$

## Získání hlavních komponent

Nechť výběrová varianční matice  $\mathbf{S}$  má vlastní čísla  $l_1, \dots, l_p$  a vlastní vektory  $\mathbf{v}_1, \dots, \mathbf{v}_p$ , přičemž  $\mathbf{v}_j^T \mathbf{v}_j = 1, j = 1, \dots, p$  a  $\mathbf{v}_j^T \mathbf{v}_k = 0$  pro  $j \neq k$ .

Znamená to, že vektory  $\mathbf{v}_1, \dots, \mathbf{v}_p$  jsou ortonormální.

Bez újmy na obecnosti předpokládáme, že  $l_1 > l_2 > \dots > l_p$ .

**1. hlavní komponenta** vznikne jako lineární kombinace znaků  $X_1, \dots, X_p$ , kde koeficienty této lineární kombinace jsou souřadnice vlastního vektoru  $\mathbf{v}_1$ , tedy

$$Y_1 = v_{11}X_1 + \dots + v_{1p}X_p.$$

Její rozptyl je  $l_1$ .

Dosadíme-li za  $X_1, \dots, X_p$  vektory pozorování  $\mathbf{x}_i, i = 1, \dots, n$ , dostaneme **vektor souřadnic**  $\mathbf{y}_1 = (y_{11}, \dots, y_{1n})^T$ , kde  $y_{1i} = \mathbf{v}_1^T \mathbf{x}_i$ .

**2. hlavní komponenta** vznikne jako lineární kombinace znaků  $X_1, \dots, X_p$ , kde koeficienty této lineární kombinace jsou souřadnice vlastního vektoru  $\mathbf{v}_2$ , tedy

$$Y_2 = v_{21}X_1 + \dots + v_{2p}X_p.$$

Její rozptyl je  $l_2$ .

Přitom  $\mathbf{v}_1^T \mathbf{v}_2 = 0$ , tj. 1. a 2. hlavní komponenta jsou lineárně nezávislé.

Dosadíme-li za  $X_1, \dots, X_p$  vektory pozorování  $\mathbf{x}_i, i = 1, \dots, n$ , dostaneme **vektor souřadnic**  $\mathbf{y}_2 = (y_{21}, \dots, y_{2n})^T$ , kde  $y_{2i} = \mathbf{v}_2^T \mathbf{x}_i$ .

.....

**j-tá hlavní komponenta** vznikne jako lineární kombinace znaků  $X_1, \dots, X_p$ , kde koeficienty této lineární kombinace jsou souřadnice vlastního vektoru  $\mathbf{v}_j$ , tedy

$$Y_j = v_{j1}X_1 + \dots + v_{jp}X_p.$$

Její rozptyl je  $l_j$ . Přitom  $\mathbf{v}_j^T \mathbf{v}_k = 0, j = 1, \dots, k-1$ , tj. j-tá hlavní komponenta je lineárně nezávislá se všemi ostatními hlavními komponentami.

Dosadíme-li za  $X_1, \dots, X_p$  vektory pozorování  $\mathbf{x}_i, i = 1, \dots, n$ , dostaneme **vektor souřadnic**  $\mathbf{y}_j = (y_{j1}, \dots, y_{jn})^T$ , kde  $y_{ji} = \mathbf{v}_j^T \mathbf{x}_i$ .

Lze dokázat, že celková variabilita obsažená v datech je rovna stopě matice  $\mathbf{S}$ , tj. součtu vlastních čísel  $l_1 + \dots + l_p$ .

1. hlavní komponenta tedy vyčerpává  $\frac{l_1}{l_1 + \dots + l_p} 100\%$  celkové variability.

Pokud je číslo  $\frac{l_1}{l_1 + \dots + l_p}$  dostatečně blízké 1, znamená to, že 1. hlavní komponenta dobře nahrazuje celý datový soubor. Je-

li toto číslo podstatně menší než 1, musíme vzít tolik hlavních komponent, aby jejich součet dělený stopou matice  $\mathbf{S}$  byl dostatečně blízký 1.

(V mnoha aplikacích se stává, že i při velkém počtu znaků stačí poměrně malý počet hlavních komponent.)

Znázorníme-li rozmístění objektů na ploše prvních dvou hlavních komponent, můžeme poznat, které objekty se řadí do skupin neboli shluků.

(Před provedením metody hlavních komponent je třeba se rozhodnout, zda budeme pracovat s původními hodnotami znaků nebo standardizovanými hodnotami.)

**Důležité upozornění:** Proměnné  $X_1, \dots, X_p$  musí být mezi sebou dostatečně korelované, jinak metoda hlavních komponent nedá dobré výsledky.

**Koeficient korelace**  $i$ -tého znaku  $X_i$  s  $k$ -tou hlavní komponentou  $Y_k$  lze vyjádřit jako  $R(X_i, Y_k) = \frac{v_{ki} \sqrt{l_k}}{s_i}$ .

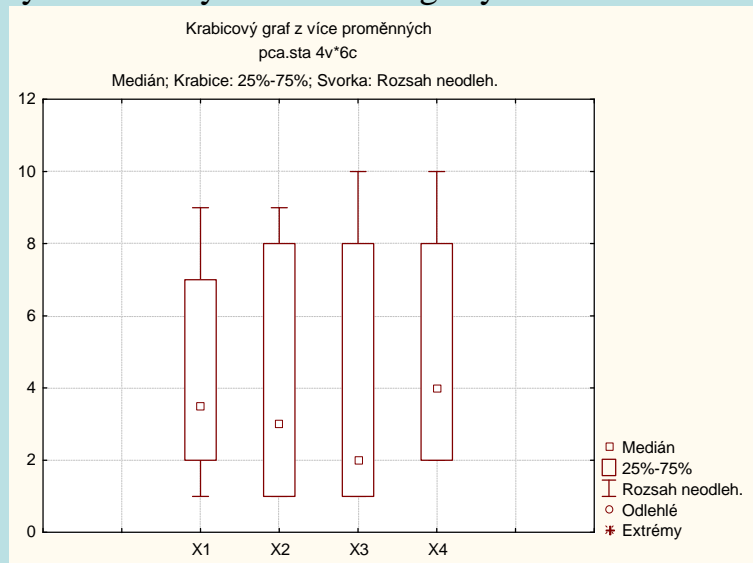
**Reprodukce výchozí kovarianční matice:** platí vzorec  $\mathbf{S} = \sum_{i=1}^p l_i \mathbf{v}_i \mathbf{v}_i^T$  (tzv. spektrální rozklad matice  $\mathbf{S}$ ).

Rozhodneme-li se uvažovat právě  $m$  hlavních komponent ( $m \leq p$ ), pak pomocí tohoto vztahu můžeme posoudit, jak těchto  $m$  hlavních komponent reprodukuje rozptyly a kovariance původních proměnných. Lze posoudit i reziduální matici, tj. matici, kterou získáme jako rozdíl výchozí kovarianční matice a reprodukované kovarianční matice.

## Doporučený postup při analýze hlavních komponent

- a) Provedeme tabulkové a grafické zpracování datového souboru, abychom se blíže seznámili s daty.
- b) Sestavíme korelační matici a prověříme, zda jsou korelace natolik silné, aby mělo smysl provádět analýzu hlavních komponent.
- c) Rozhodneme, kolika hlavními komponentami lze popsat datový soubor bez podstatné ztráty informace. Označme tento vhodný počet jako  $m$ . Při stanovení  $m$  můžeme použít tato pomocná kritéria:
  - **Kaiserovo kritérium** - za  $m$  volíme počet těch vlastních čísel matice  $\mathbf{R}$ , která jsou větší než 1.
  - **Sutinový test** (scree test) – grafická metoda, která spočívá v subjektivním posouzení vzhledu sutinového grafu (scree plot), tj. grafu znázorňujícího velikosti sestupně uspořádaných vlastních čísel matice  $\mathbf{R}$ . Objeví-li se v grafu určité zploštění, pak za  $m$  vezmeme to pořadové číslo, kde se zploštění projevilo.
  - **Kritérium založené na kumulativním procentu vysvětleného rozptylu**. Požadujeme, aby vybrané hlavní komponenty vysvětlily aspoň 70% celkového rozptylu.
  - **Kritérium založené na reziduální korelační či kovarianční matici**. Požadujeme, aby prvky reziduální matice byly co možná nejmenší.
- d) Pokusíme se o interpretaci prvních  $m$  hlavních komponent. Zkoumáme přitom, jak jsou jednotlivé vybrané hlavní komponenty utvořeny z původních znaků a jak s nimi korelují.
- e) Vypočítáme vektory souřadnic a následně sestrojíme dvourozměrné tečkové diagramy.

Pro náš datový soubor obsahující výsledky 6 žáků ve 4 testech nejprve znázorníme data pomocí krabicových diagramů: Grafy – 2D Grafy – Krabicové grafy – zvolíme Vícenásobný – Proměnné - Závisle proměnné X1-X4 – OK – OK



Nyní vypočteme korelační matici: Statistiky – Vícerozměrné průzkumné techniky – Hlavní komponenty & klasifikační analýza – Proměnné X1 až X4, OK – OK – Popisné statistiky – Korelační matice

Proměnná	Korelace (pca.sta)			
	X1	X2	X3	X4
X1	1,000000	0,906937	0,872258	0,870307
X2	0,906937	1,000000	0,958133	0,820031
X3	0,872258	0,958133	1,000000	0,893684
X4	0,870307	0,820031	0,893684	1,000000

Dále vypočteme vlastní čísla a procento vysvětleného rozptylu: na záložce Základní výsledky vybereme Vlastní čísla.

Vlastní čísla korelační matice a související statistiky (pca) Pouze aktiv. proměnné				
Pořadí vl.č.	vl. číslo	% celk. rozptylu	Kumulativ. vl. číslo	Kumulativ. %
1	3,661431	91,53577	3,661431	91,5358
2	0,188636	4,71589	3,850066	96,2517
3	0,134072	3,35181	3,984139	99,6035
4	0,015861	0,39653	4,000000	100,0000

Vidíme, že 1. vlastní číslo  $l_1 = 3,66$ , tedy 1. hlavní komponenta vyčerpává 91,5% variability dat,

2. vlastní číslo  $l_2 = 0,19$ , 2. hlavní komponenta vyčerpává 4,7% variability dat atd.

Podle Kaiserova kritéria by stačilo uvažovat pouze 1. hlavní komponentu, protože pouze první vlastní číslo je větší než 1.

Kvůli znázornění objektů však budeme uvažovat první dvě hlavní komponenty.

Dále vypočítáme vlastní vektory: na záložce Proměnné vybereme Vlastní vektory

Vlastní vektory korelační matice (pca) Pouze aktiv. proměnné				
Proměnná	Faktor 1	Faktor 2	Faktor 3	Faktor 4
X1	-0,498301	-0,000518	0,817131	-0,289816
X2	-0,503657	0,582217	-0,082290	0,632916
X3	-0,508833	0,185043	-0,539021	-0,645217
X4	-0,488994	-0,791696	-0,187036	0,314832

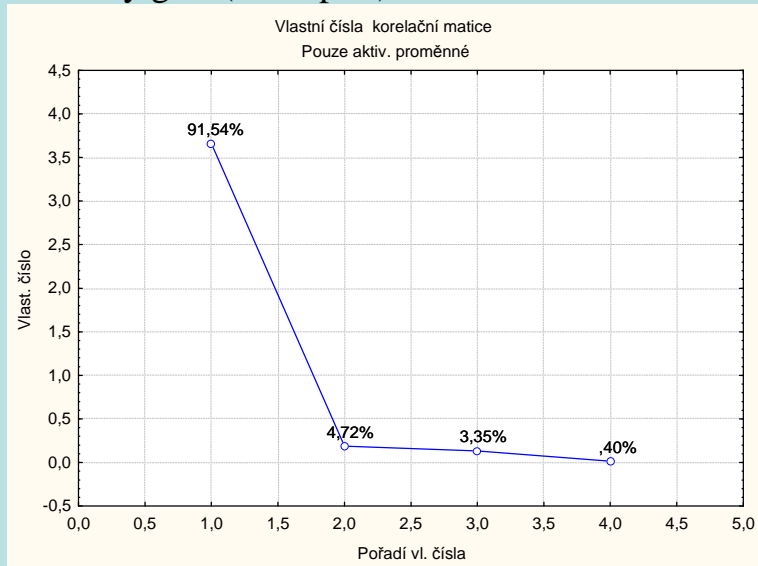
1. hlavní komponenta:

$$Y_1 = -0,49X_1 - 0,5X_2 - 0,51X_3 - 0,49X_4,$$

2. hlavní komponenta:

$$Y_2 = -0,0005X_1 + 0,58X_2 + 0,19X_3 - 0,79X_4 \text{ atd.}$$

## Sutinový graf (scree plot):



V sutinovém grafu nastává výrazné zploštění po 1. vlastním čísle.

Výpočet koeficientů korelace 1. a 2. hlavní komponenty a původních čtyř proměnných: na záložce Proměnné vybereme Korelace faktorů & proměnných

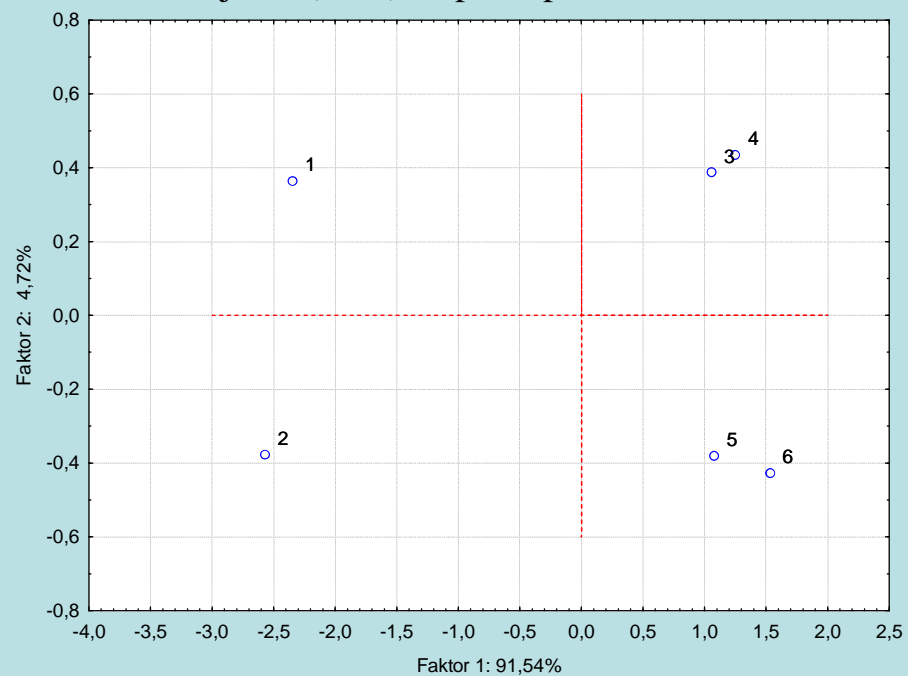
Proměnná	Faktor 1	Faktor 2
X1	-0,953492	-0,000225
X2	-0,963740	0,252869
X3	-0,973645	0,080368
X4	-0,935684	-0,343851

Vidíme, že 1. hlavní komponenta vysoce záporně koreluje se všemi proměnnými. 2. hlavní komponenta slabě kladně koreluje s druhou proměnnou a středně silně záporně koreluje s třetí proměnnou.

Podívejme se rovněž na vektory souřadnic (v systému STATISTICA se jim říká faktorové souřadnice případů): na záložce Případy vybereme Faktorové souřadnice případů.

Případ	Faktor 1	Faktor 2
1	-2,34914	0,364696
2	-2,56859	-0,378068
3	1,05532	0,387487
4	1,25040	0,434674
5	1,07964	-0,381138
6	1,53238	-0,427651

Znázornění objektů (žáků) na ploše prvních dvou hlavních komponent:





## Shluková analýza

### Cíl shlukové analýzy

Cílem shlukové analýzy je rozřídění  $n$  objektů, z nichž každý je popsán  $p$  znaky, do několika pokud možno stejnorodých (homogenních) skupin (shluků, clusterů). Požadujeme, aby objekty uvnitř shluků si byly podobné co nejvíce, zatímco objekty z různých shluků co nejméně. Přesný počet shluků většinou není přesně znám.

Shluková analýza nachází uplatnění v celé řadě oborů, např. v biologii. U  $n$  populací změříme  $p$  biometrických charakteristik a zjišťujeme, zda určité skupiny populací tvoří shluky.

Shluková analýza je ovšem průzkumovou metodou a měla by sloužit jako určité vodítko při dalším zpracování dat.

### Podobnost objektů

Podobnost (či rozdílnost) objektů posuzujeme pomocí různých měr vzdálenosti. Pro znaky intervalového či poměrového typu nejčastěji používáme **euklidovskou vzdálenost**.

Nechť  $k$ -tý objekt je popsán vektorem pozorování  $\mathbf{x}_k = (x_{k1}, \dots, x_{kp})^T$  a  $l$ -tý objekt vektorem  $\mathbf{x}_l = (x_{l1}, \dots, x_{lp})^T$ .

Euklidovská vzdálenost  $k$ -tého a  $l$ -tého objektu:

$$d_{kl} = \sqrt{\sum_{j=1}^p (x_{kj} - x_{lj})^2}.$$

Vzdálenosti vypočtené pro všechny dvojice objektů se uspořádají do **matice vzdáleností**. Je zřejmé, že je to čtvercová symetrická matice, která má na hlavní diagonále nuly.

### Matice euklidovských vzdáleností pro datový soubor s údaji o 6 žácích:

Statistiky – Vícerozměrné průzkumné techniky – Shluková analýza – Spojování (hierarchické shlukování) – OK – Proměnné X1 – X4 – OK – na záložce Details vybereme Shlukovat Případy (řádky) – OK – na záložce Details vybereme Matice vzdáleností.

Případ	Euklid. vzdálenosti (pca)					
	P_1	P_2	P_3	P_4	P_5	P_6
P_1	0,0	3,6	12,7	12,7	12,6	14,0
P_2	3,6	0,0	12,8	13,2	12,5	14,1
P_3	12,7	12,8	0,0	2,2	3,2	4,1
P_4	12,7	13,2	2,2	0,0	3,0	3,2
P_5	12,6	12,5	3,2	3,0	0,0	2,2
P_6	14,0	14,1	4,1	3,2	2,2	0,0

## Hierarchické shlukování

Při aplikacích shlukové analýzy se nejčastěji používá **aglomerativní hierarchická procedura**. Její princip spočívá v postupném slučování objektů, a to nejprve nejbližších a v dalších krocích pak stále vzdálenějších.

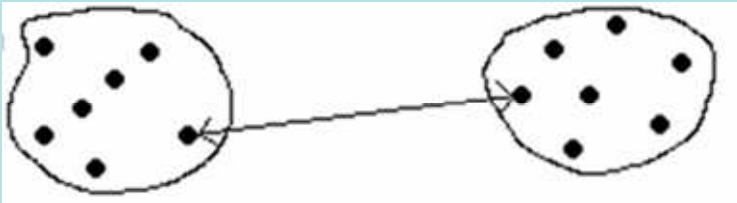
### Algoritmus:

1. krok: Každý objekt považujeme za samostatný shluk.
2. krok: Najdeme dva shluky, jejichž vzdálenost je minimální.
3. krok: Tyto dva shluky spojíme v nový, větší shluk a přepočítáme matici vzdáleností. Její řád se sníží o 1. Vrátime se na 2. krok.

Funkce algoritmu končí, až jsou všechny objekty spojeny do jediného shluku.

Vzdálenost mezi shluky se počítá různými způsoby. Uvedeme tři z nich.

a) **Metoda nejbližšího souseda:** Vzdálenost mezi dvěma shluky je minimem ze všech vzdáleností mezi jejich objekty.



b) **Metoda nejvzdálenějšího souseda:** Vzdálenost mezi dvěma shluky je maximem ze všech vzdáleností mezi jejich objekty.

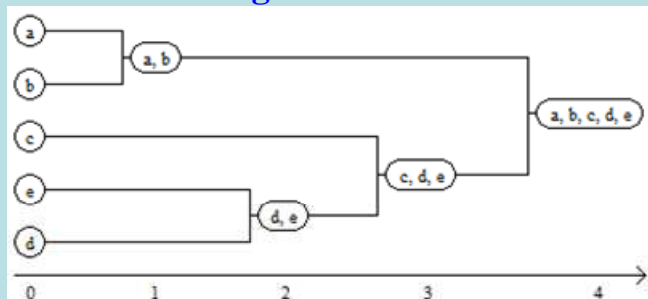


c) **Metoda průměrné vazby:** Vzdálenost mezi dvěma shluky je průměrem ze všech vzdáleností mezi jejich objekty.



Výsledky aglomerativní hierarchické procedury se zpravidla znázorňují pomocí **dendrogramu**. Je to graficky znázorněná posloupnost dvojic  $\{(v_1, S^{(1)}), \dots, (v_n, S^{(n)})\}$ , kde  $\{v_i\}_{i=1}^n$  je neklesající posloupnost úrovní spojování a  $S^{(i)}$  je rozřídění objektů odpovídající úrovni  $v_i$ ,  $i = 1, \dots, n$ .

### Příklad dendrogramu:



V levém sloupci jsou jednotlivé objekty, další sloupce reprezentují shluky, do nichž byly objekty zařazeny a délky čar představují vzdálenosti mezi shluky.

### Poznámka:

Hierarchická shluková analýza může být použita nejen na shlukování objektů, ale též na shlukování znaků.

**Dendrogram podobnosti objektů** je standardní výstup hierarchických shlukovacích metod, z něhož je zjevná struktura objektů ve shlucích.

**Dendrogram podobnosti znaků** odhaluje nejčastěji dvojice či trojice (všeobecně m-tice) znaků, které si jsou velmi podobné a silně spolu korelují. Znaky, které jsou ve společném shluku, si jsou značně podobné a jsou tudíž vzájemně nahraditelné. To má značný význam při plánování experimentu - některé vlastnosti či znaky není zapotřebí vůbec zjišťovat či měřit, protože jsou snadno nahraditelné jinými znaky a nemají velkou vypovídací hodnotu.

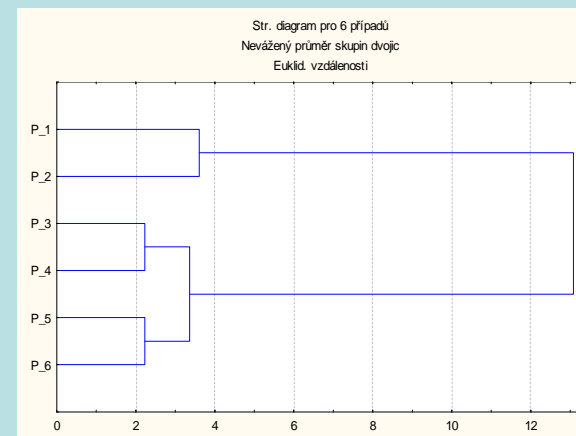
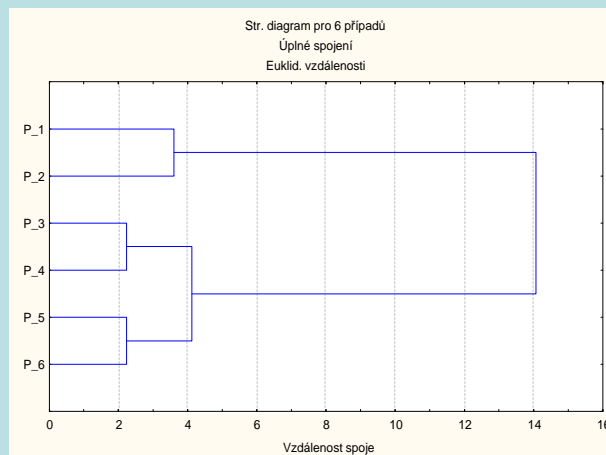
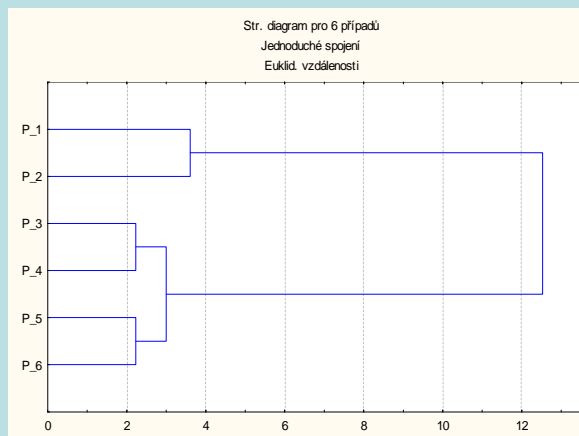
## Vytvoření dendrogramu v systému STATSTICA:

- pro metodu nejbližšího souseda:

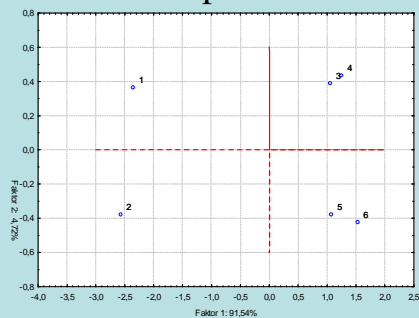
Statistiky – Vícerozměrné průzkumné techniky – Shluková analýza – Spojování (hierarchické shlukování) – OK – Proměnné X1 – X4 – OK – na záložce Details vybereme Shlukovat Případy (řádky), pravidlo slučování ponecháme Jednoduché spojení, míru vzdálenosti ponecháme Euklidovské vzd. – OK – Horizontální graf hierarch. stromu

- pro metodu nejvzdálenějšího souseda: na záložce Details vybereme pravidlo slučování Úplné spojení,

- pro metodu úplné vazby: Na záložce Details vybereme pravidlo slučování Nevážený průměr skupin dvojic.



Vidíme, že výsledky všech tří metod jsou velmi podobné a odpovídají rozmístění objektů (žáků) na ploše prvních dvou hlavních komponent.



**Příklad:** Uvažme datový soubor s údaji o 26 evropských státech. Tento datový soubor budeme analyzovat metodou hlavních komponent a následně provedeme shlukovou analýzu.

### Provedení PCA

Nejprve pomocí korelační matice posoudíme, zda má smysl aplikovat PCA.

Statistiky – Vícerozměrné průzkumné techniky – Hlavní komponenty & klasifikační analýza – Proměnné X1 až X19, OK – OK – Popisné statistiky – Korelační matice.

Proměnná	Korelace (staty1979.sta)								
	X1	X2	X3	X4	X5	X6	X7	X8	X9
X1	1,00	0,04	-0,67	-0,40	-0,53	-0,73	-0,22	-0,75	-0,56
X2	0,04	1,00	0,44	0,41	-0,02	-0,40	-0,44	-0,28	0,16
X3	-0,67	0,44	1,00	0,39	0,48	0,21	-0,15	0,15	0,36
X4	-0,40	0,41	0,39	1,00	0,03	0,20	0,11	0,13	0,37
X5	-0,53	-0,02	0,48	0,03	1,00	0,33	0,01	0,17	0,38
X6	-0,73	-0,40	0,21	0,20	0,33	1,00	0,36	0,57	0,17
X7	-0,22	-0,44	-0,15	0,11	0,01	0,36	1,00	0,11	-0,25
X8	-0,75	-0,28	0,15	0,13	0,17	0,57	0,11	1,00	0,56
X9	-0,56	0,16	0,36	0,37	0,38	0,17	-0,25	0,56	1,00

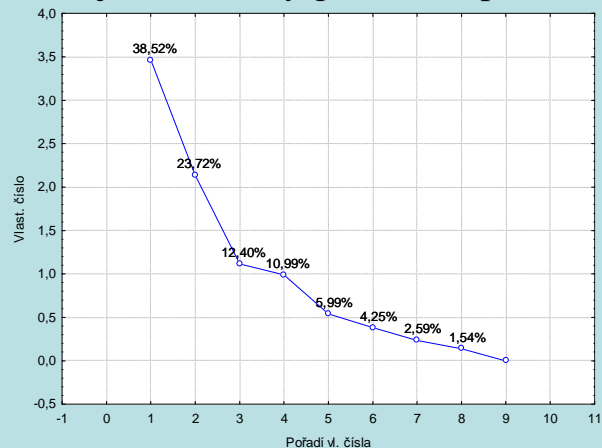
Některé korelační koeficienty jsou v absolutní hodnotě dostatečně velké a zřejmě tedy bude mít smysl provést analýzu hlavních komponent.

Nyní získáme vlastní čísla výběrové korelační matice a procento vysvětleného rozptylu: na záložce Základní výsledky vybereme Vlastní čísla.

Pořadí vl.č.	vl. číslo	% celk. rozptylu	Kumulativ. vl. číslo	Kumulativ. %
1	3,466490	38,51655	3,466490	38,5166
2	2,135004	23,72227	5,601494	62,2388
3	1,115581	12,39534	6,717075	74,6342
4	0,989394	10,99326	7,706468	85,6274
5	0,539211	5,99123	8,245679	91,6187
6	0,382111	4,24568	8,627790	95,8643
7	0,233226	2,59140	8,861015	98,4557
8	0,138985	1,54428	9,000000	100,0000

První hlavní komponenta tedy vysvětluje 38,52% variability obsažené v devíti sledovaných proměnných, druhá 23,72%, třetí 12,40% atd. Celkové procento variability vysvětlené prvními třemi hlavními komponentami je 74,63%.

Sestrojíme sutinový graf (scree plot): na záložce Základní výsledky vybereme Sutinový graf.

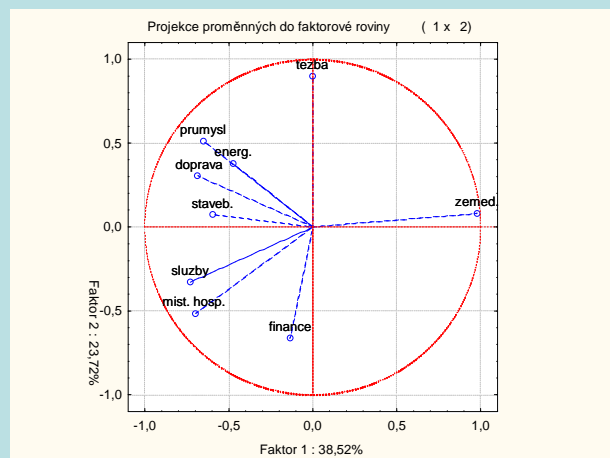


Počet m hlavních komponent zvolíme tři. V nabídce Výsledky hlavních komponent snížíme počet faktorů na 3.

Vypočteme korelační koeficienty prvních tří hlavních komponent a původních devíti proměnných: na záložce Proměnné vybereme Korelace faktorů & proměnných.

Proměnná	Korelace faktorů a proměnných (faktor. zátěže) podle korelací (staty1979.sta)		
	Faktor 1	Faktor 2	Faktor 3
X1	0,978776	0,081725	-0,049455
X2	-0,000898	0,901105	0,216344
X3	-0,652174	0,513343	0,112868
X4	-0,474888	0,378598	0,649962
X5	-0,595263	0,073032	-0,304047
X6	-0,698213	-0,513734	0,119592
X7	-0,136193	-0,663299	0,589451
X8	-0,727506	-0,327637	-0,251642
X9	-0,684094	0,304809	-0,337074

Graficky lze znázornit souvislost mezi novými proměnnými (např. 1. a 2. HK) a původními proměnnými X1, ..., X9 takto: na záložce Proměnné vybereme 2D graf fakt. souřadnic prom. - Osa x: Faktor I, Osa y: Faktor 2 - OK. Na ose x budou souřadnice vstupních proměnných vzhledem k první hlavní komponentě, na ose Y vzhledem ke druhé komponentě.



1. HK vysoce kladně koreluje s proměnnou X1, tj se zemědělstvím a negativně s proměnnou X8 – služby. Jelikož je podíl lidí v zemědělství a ve službách obecně považován za určité měřítko vyspělosti země, můžeme první komponentu interpretovat jako míru zaostalosti/vyspělosti.
2. HK výrazně pozitivně koreluje s těžebním průmyslem, energetikou a zpracovatelským průmyslem. Negativně koreluje se službami a finanční sférou. Budeme ji proto interpretovat jako míru toho, nakolik se země orientuje na průmyslovou výrobu. (Ne vždy mají komponenty takto jasnou interpretaci. Jsou jen jistou matematickou transformací vstupních proměnných, která může a nemusí odrážet nějakou reálnou vlastnost objektů!).

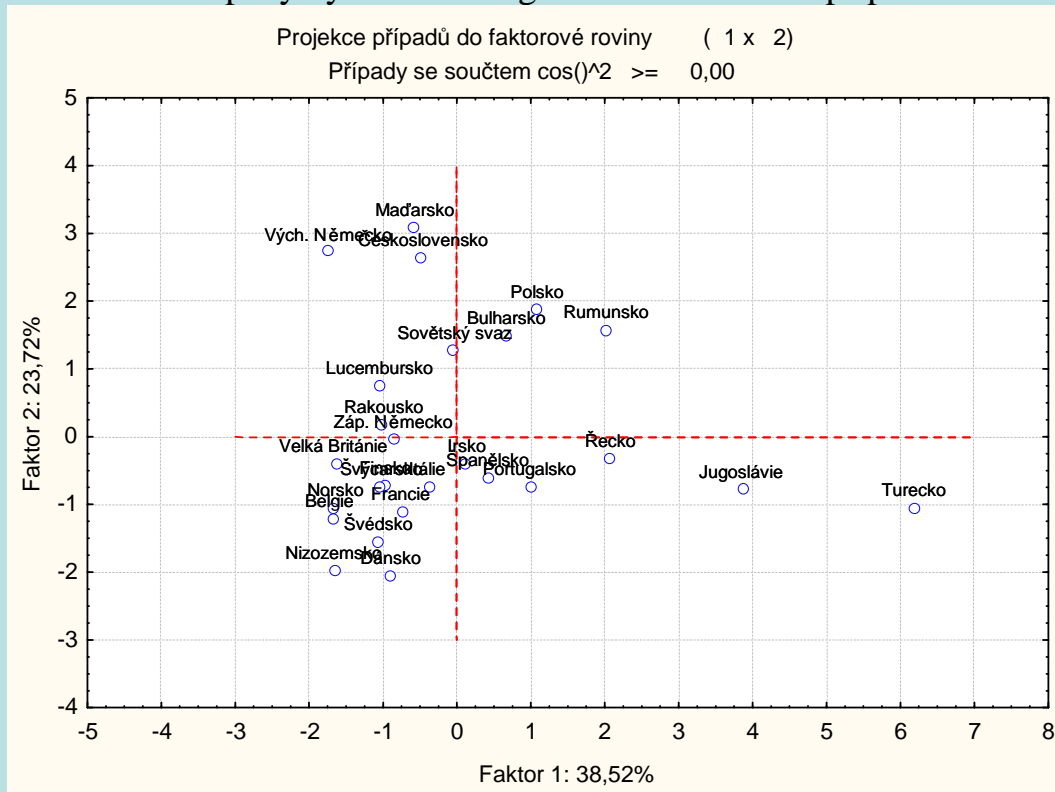


Podívejme se rovněž na vektory souřadnic (v systému STATISTICA se jim říká faktorové souřadnice případů): na záložce Případy vybereme Faktorové souřadnice případů.

Případ	Faktor 1	Faktor 2	Faktor 3
Belgie	-1,68273	-1,20656	0,16668
Dánsko	-0,90831	-2,05598	-0,85147
Francie	-0,74050	-1,11048	0,38553
Záp. Německo	-0,85647	-0,03165	0,56466
Irsko	0,11153	-0,40400	0,53134
Itálie	-0,36366	-0,74902	-1,29050
Lucembursko	-1,04022	0,74294	0,46327
Nizozemsko	-1,65732	-1,98866	-0,08729
Velká Británie	-1,61201	-0,39776	1,35031
Rakousko	-1,01103	0,16508	1,16804
Finsko	-0,97223	-0,73166	0,54475
Řecko	2,07154	-0,33521	-0,92274
Norsko	-1,66538	-1,05092	-1,14341
Portugalsko	0,99709	-0,74259	-0,75474
Španělsko	0,43244	-0,60818	0,31825
Švédsko	-1,07387	-1,55390	-0,22815
Švýcarsko	-1,04031	-0,74707	0,28216
Turecko	6,19519	-1,04930	-0,64265
Bulharsko	0,67558	1,48159	-1,03101
Československo	-0,48005	2,63421	0,07902
Vých. Německo	-1,73669	2,73412	0,26970
Maďarsko	-0,57526	3,07981	1,09460
Polsko	1,08637	1,87264	-0,54684
Rumunsko	2,01536	1,57550	-0,48595
Sovětský svaz	-0,04779	1,26246	-2,30671
Jugoslávie	3,87872	-0,78542	3,07316

1. HK vysoce kladně koreluje s proměnnou  $X_1$  (zemědělství) a záporně se všemi ostatními proměnnými. Tato hlavní komponenta tedy rozlišuje země na zemědělské a průmyslové. Povšimněte si, že souřadnice této hlavní komponenty jsou nejvyšší u Turecka (6,2) a Jugoslávie (3,9).
2. HK vysoce kladně koreluje s proměnnou  $X_2$  (těžba) a podstatně slaběji s proměnnou  $X_3$  (průmyslová výroba). Vysoké hodnoty souřadnic této hlavní komponenty najdeme u Maďarska, Východního Německa a Československa.
3. HK středně silně koreluje s proměnnou  $X_4$  (energetika) a  $X_7$  (finanční sektor). Nejvyšší hodnotu najdeme u Jugoslávie.

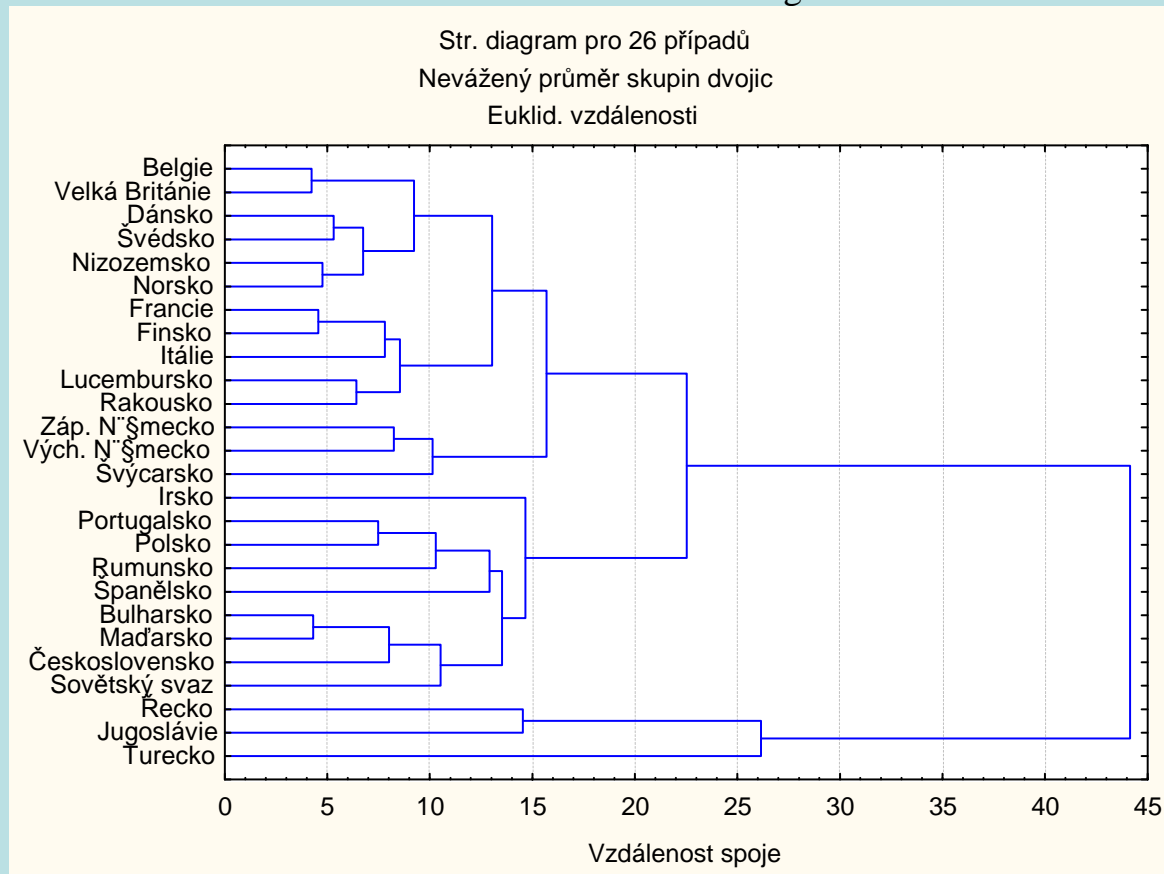
Nyní znázorníme rozmístění zemí na ploše prvních dvou hlavních komponent:  
 Na záložce Případy vybereme 2D graf fakt. Souřadnic příp.



Státy napravo jsou státy s vysokým podílem zemědělství. Vyniká zde zejména Turecko a Jugoslávie. Všechny státy obvykle považované za ekonomicky vyspělé jsou naopak na levé straně. Jsou to státy, kde je nižší podíl osob zaměstnaných v zemědělství, zato vyšší podíl osob pracujících ve službách. Je zde také hezky vidět zaměření zemí tehdejšího socialistického bloku na průmyslovou výrobu - horní část grafu. A naopak severské státy a státy Beneluxu orientované na finanční a další služby v dolní části.

## Provedení shlukové analýzy

Statistiky – Vícerozměrné průzkumné techniky – Shluková analýza - Spojování (hierarchické shlukování) – OK - Proměnné X1 až X4, OK, Detaily - Shlukovat případy (řádky) – Pravidlo slučování: Nevážený průměr skupin dvojic – Míry vzdálenosti: Euklidovské vzdálenosti - OK – Horizontální graf hierarch. stromu.



Ukazuje se, že země se dělí do tří skupin: první skupinu tvoří rozvinuté demokratické země společně s NDR, druhou skupinu socialistické země s Irskem, Portugalskem a Španělskem a třetí Řecko s Jugoslávií. Turecko se chová jako singulární entita.