

NCSS2007

ve vícerozměrné statistické analýze

Metoda hlavních komponent PCA

Analýza hlavních komponent (PCA)

Zaměření metody PCA: PCA jedna z nejstarších a nejvíce používaných metod vícerozměrné analýzy.

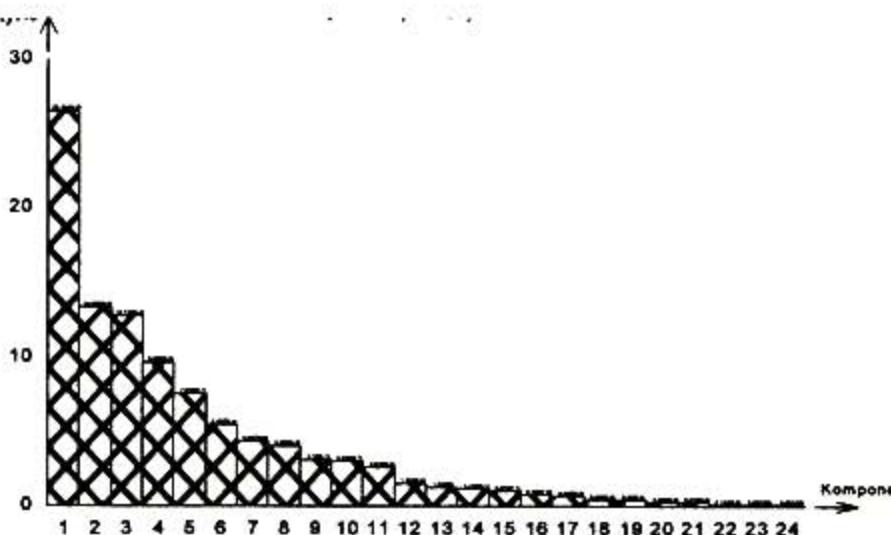
Zavedena: Pearsonem již v roce 1901 a nezávisle Hotellingem v roce 1933.

Cílem PCA: zjednodušení popisu lineárně závislých tj. korelovaných znaků, a to rozkladem matice dat do **matice strukturní** (= využité hlavní komponenty) a do **matice šumové** (= nevyužité hlavní komponenty).

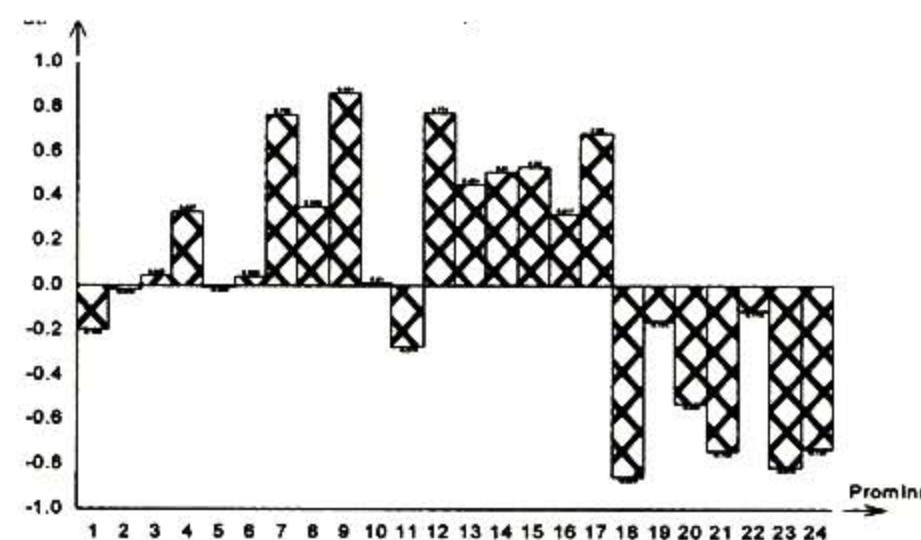
Popis metody PCA:

- 1) Lineární transformace původních znaků na nové, nekorelované proměnné, zvané *hlavní komponenty*.
- 2) Základní charakteristikou každé hlavní komponenty je její míra variability čili *rozptyl*.
- 3) Hlavní komponenty jsou *seřazeny dle důležitosti*, tj. dle klesajícího rozptylu, od největšího k nejmenšímu.
- 4) *Většina informace* o variabilitě dat je přitom soustředěna do první komponenty a nejméně informace je obsaženo v poslední komponentě.

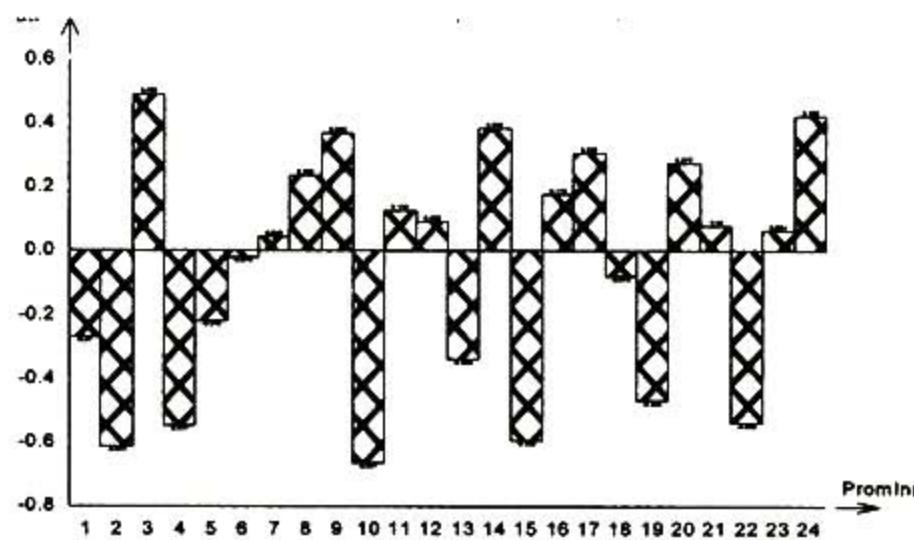
- 5) ***Pravidlo:*** má-li nějaký původní znak malý či dokonce žádný rozptyl, není schopen přispívat k rozlišení mezi objekty.
- 6) Využitím PCA je *snižení dimenze úlohy* čili redukce počtu znaků bez velké ztráty informace, užitím pouze prvních několika hlavních komponent.
- 7) *Nevyužité hlavní komponenty* obsahují malé množství informace, protože jejich rozptyl je příliš malý.
- 8) Hlavní komponenty jsou nekorelované.
- 9) První hlavní komponenta je například vhodným ukazatelem jakosti.
- 10) První dvě resp. první tři hlavní komponenty se využívají především jako techniky zobrazení vícerozměrných dat v projekci do roviny (nebo do prostoru).



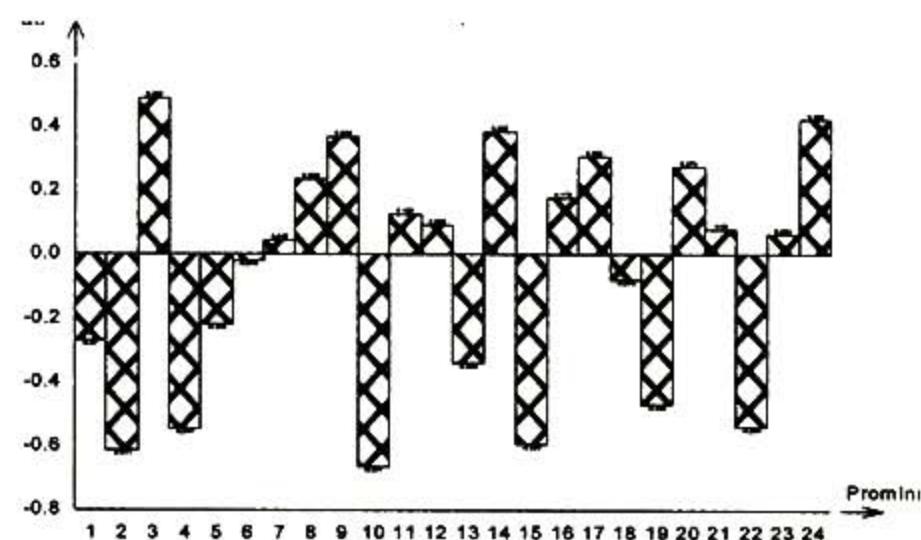
Obr. 1. Sloupcový diagram indexového grafu úpatí pro 38 objektů a 24 původních proměnných zdrojové matice Wine.



Obr. 2. Složení 1. hlavní komponenty z původních proměnných pro 38 objektů a 24 původních proměnných zdrojové matice Wine.



Obr. 3. Složení 2. hlavní komponenty z původních proměnných pro 38 objektů a 24 původních proměnných zdrojové matice Wine.



Obr. 4. Složení 3. hlavní komponenty z původních proměnných pro 38 objektů a 24 původních proměnných zdrojové matice Wine.

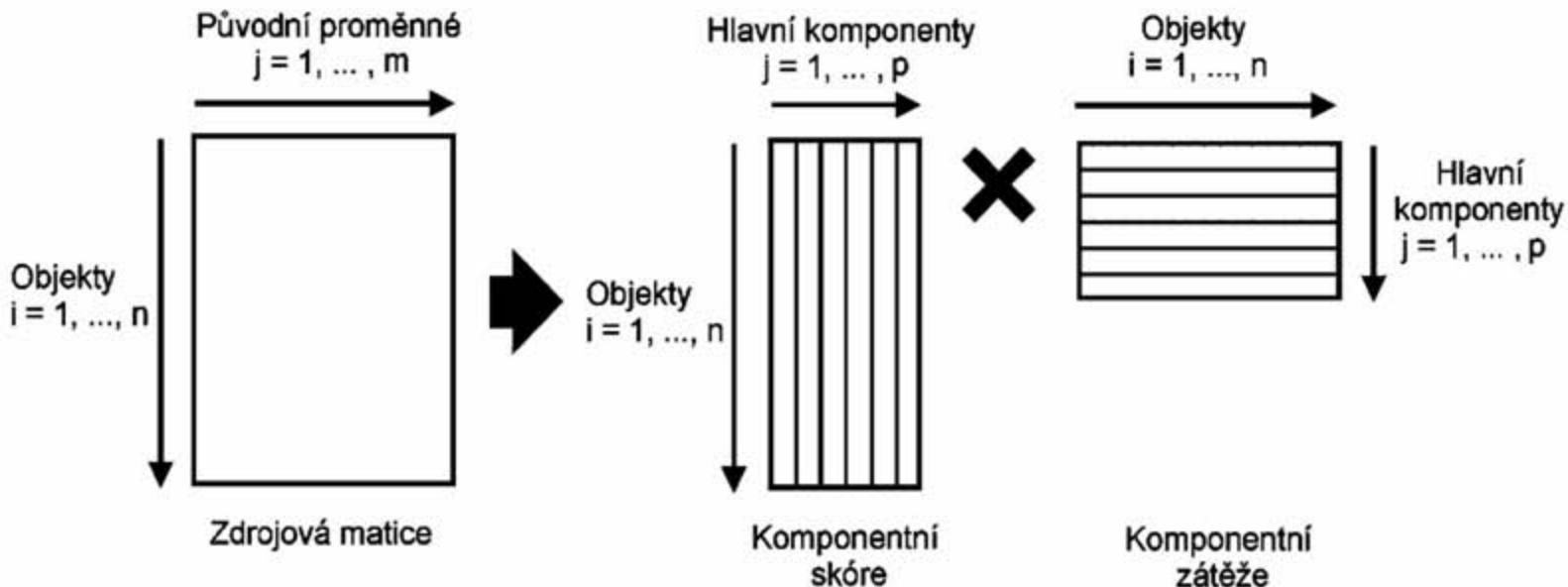
Analýza hlavních komponent (PCA)

Cílem PCA: transformace dat z původních proměnných x_j , $j=1, \dots, m$, do menšího počtu latentních proměnných y_j .

PCA vystihují téměř *proměnlivost* původních proměnných a jsou vzájemně nekorelované.

První hlavní komponenta y_1 popisuje největší část proměnlivosti čili rozptylu původních dat,

Druhá hlavní komponenta y_2 popisuje zase největší část rozptylu neobsaženého v y_1 atd.



Obr. 4.1 Schéma maticových výpočtů v PCA.

- **Maximální počet hlavních komponent:** počet efektivních hlavních komponent se rovná *hodnosti zdrojové matice X*.
- **X = Struktura + šum:** všechny hlavní komponenty jsou vzájemně ortogonální a souvisí postupně se snižující hodnotou rozptylu objektů.

Model hlavních komponent má tvar

$$X = T P^T + E = \text{struktura dat} + \text{šum}.$$

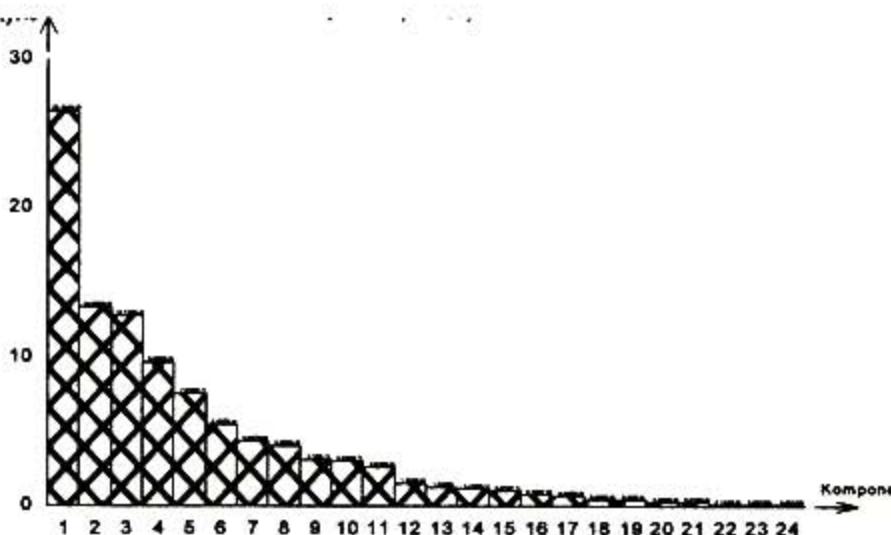
- **Střed modelu:** hlavní komponenty mají společný počátek, který odpovídá *průměrnému objektu* čili těžišti celého shluku objektů (centrování).

● Komponentní váhy, zátěže - vztah mezi X a PC :

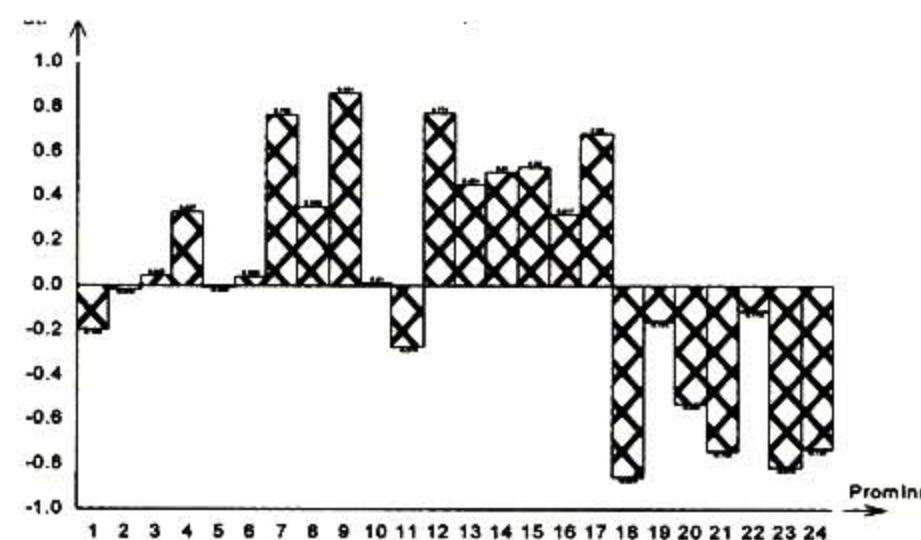
Hlavní komponenta představuje *lineární kombinaci všech m vektorů* v prostoru znaků v m rozměrném prostoru a jejich koeficienty se nazývají *komponentní váhy*.

Komponentní váhy informují o vztahu mezi původními m znaky a hlavními komponentami.

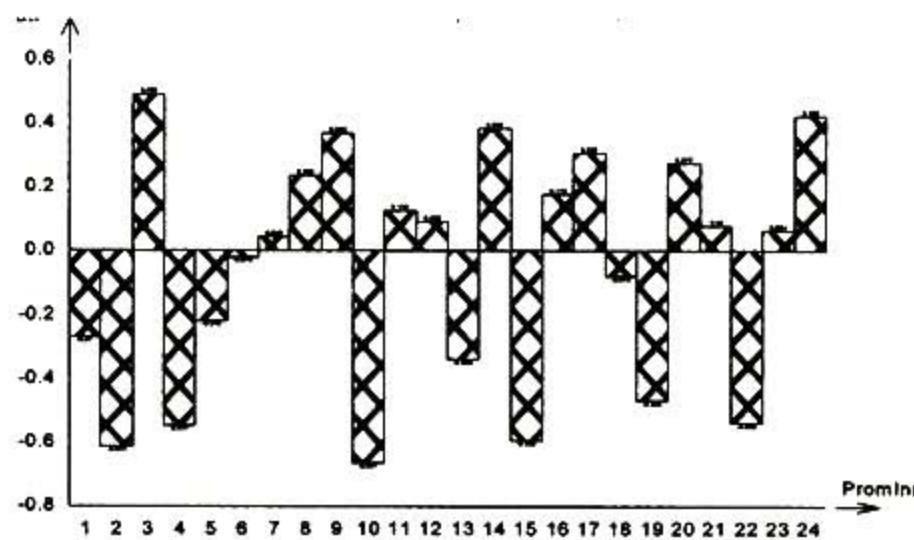
Na grafu komponentních vah p pro $PC1$ a $PC2$ jsou místo objektů jejich znaky a lze tak vyšetřovat závislosti a podobnosti mezi znaky.



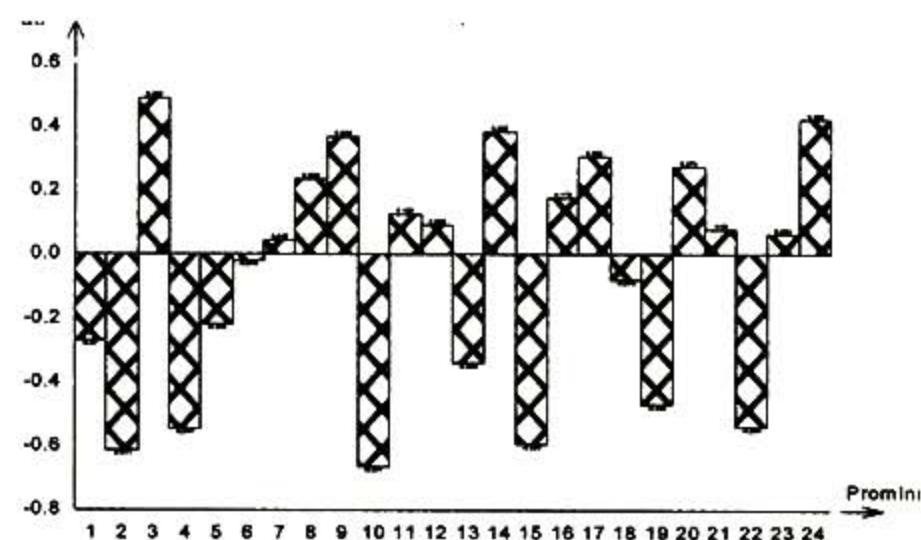
Obr. 1. Sloupcový diagram indexového grafu úpatí pro 38 objektů a 24 původních proměnných zdrojové matice Wine.



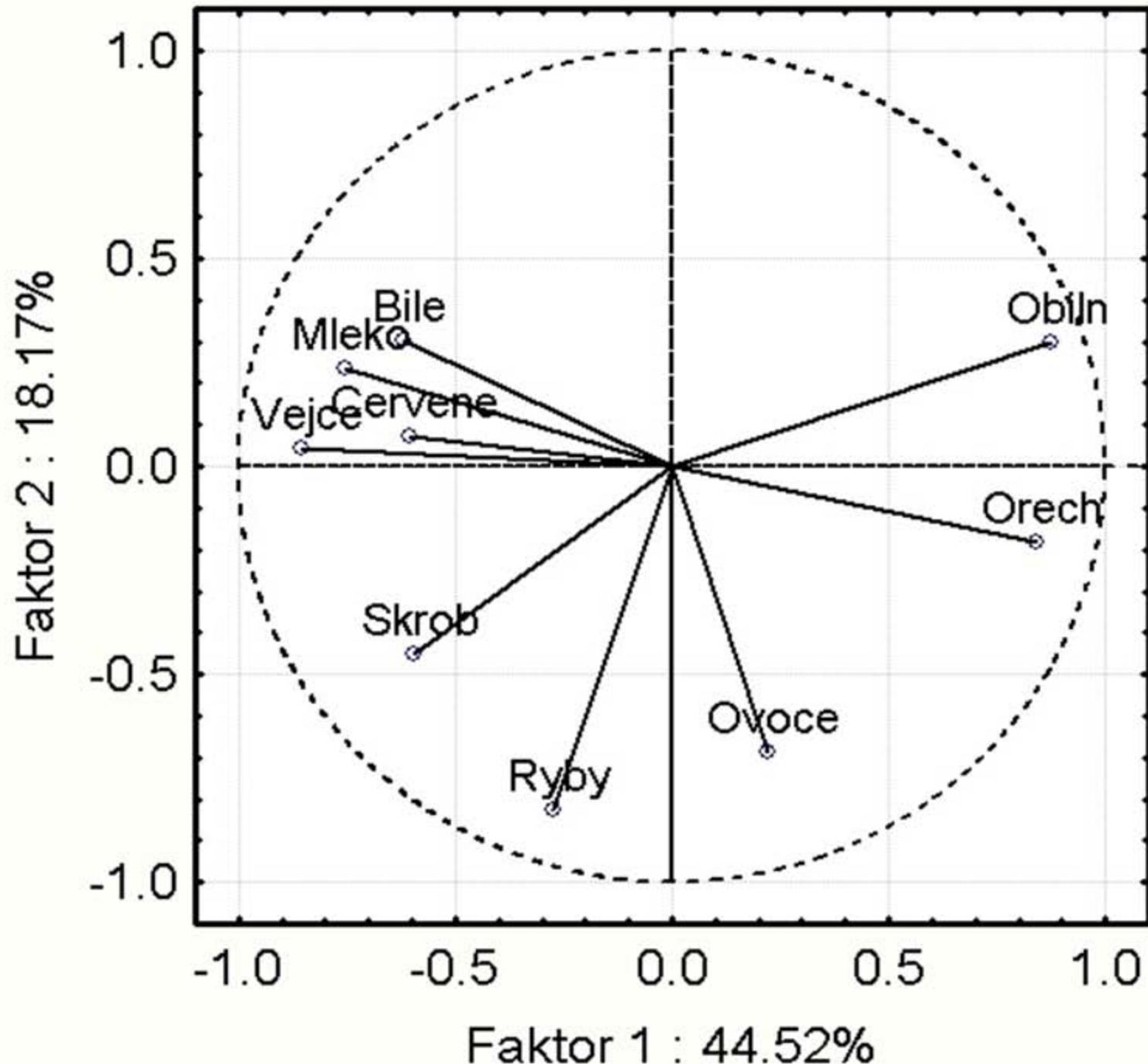
Obr. 2. Složení 1. hlavní komponenty z původních proměnných pro 38 objektů a 24 původních proměnných zdrojové matice Wine.



Obr. 3. Složení 2. hlavní komponenty z původních proměnných pro 38 objektů a 24 původních proměnných zdrojové matice Wine.



Obr. 4. Složení 3. hlavní komponenty z původních proměnných pro 38 objektů a 24 původních proměnných zdrojové matice Wine.



Výklad:

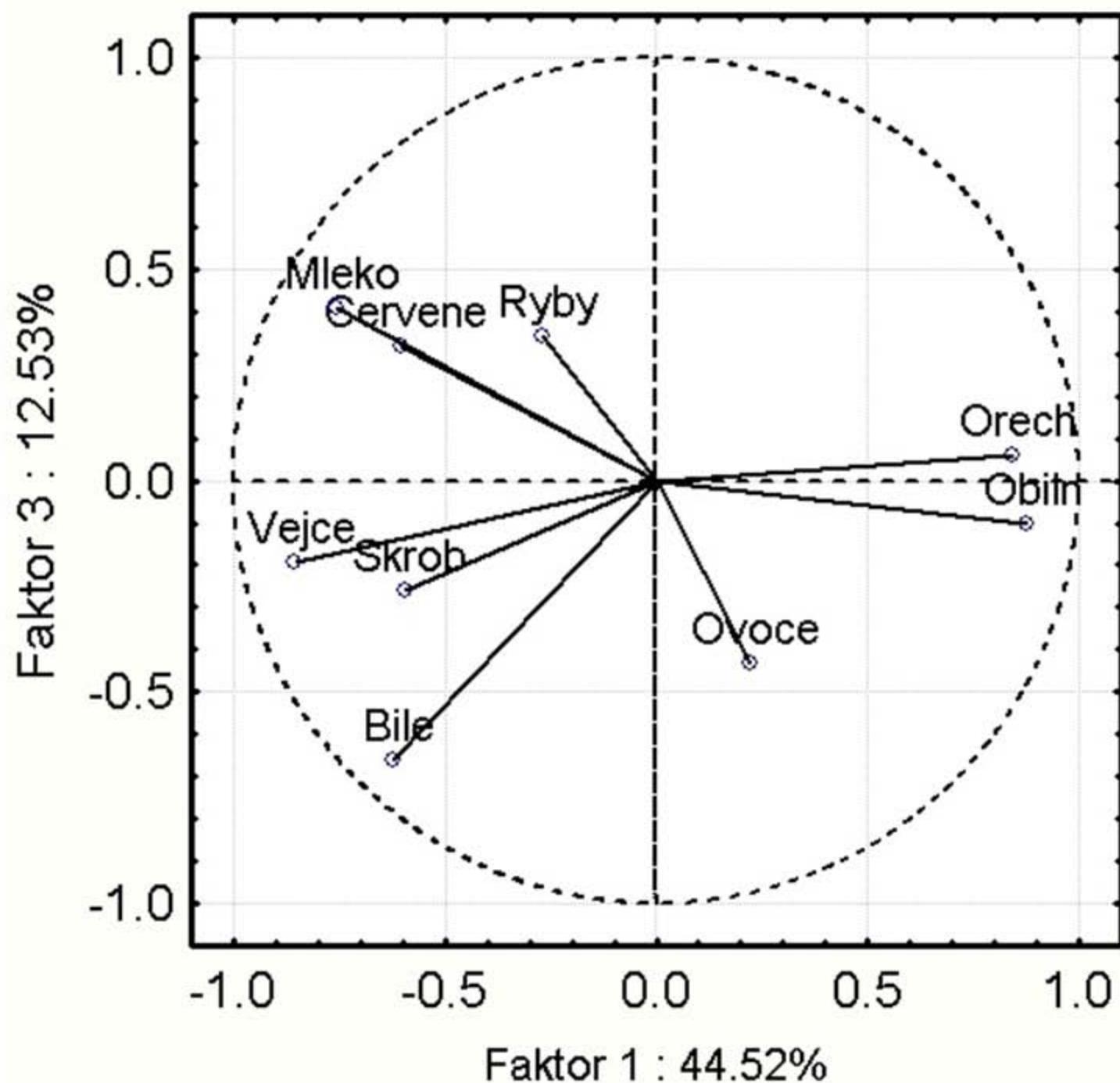
Graf komponentních vah, zátěží (Plot Components Weights)

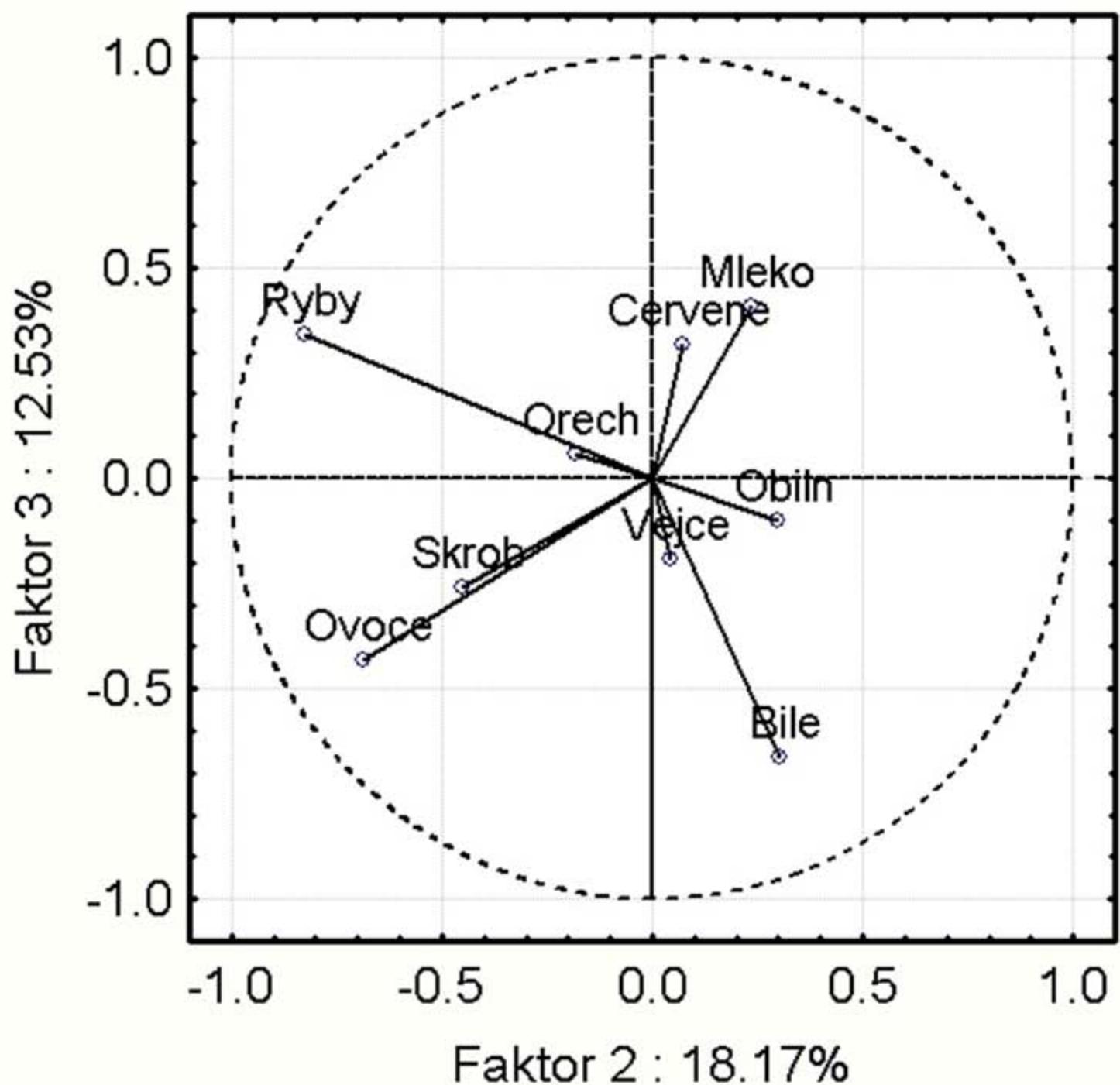
Zobrazí: komponentní váhy

Porovnávají se: vzdálenosti mezi proměnnými.
Krátká znamená silnou korelací.

Nalezneme: shluk podobných proměnných, jež spolu korelují.

Představuje: most mezi původními proměnnými a hlavními komponentami.





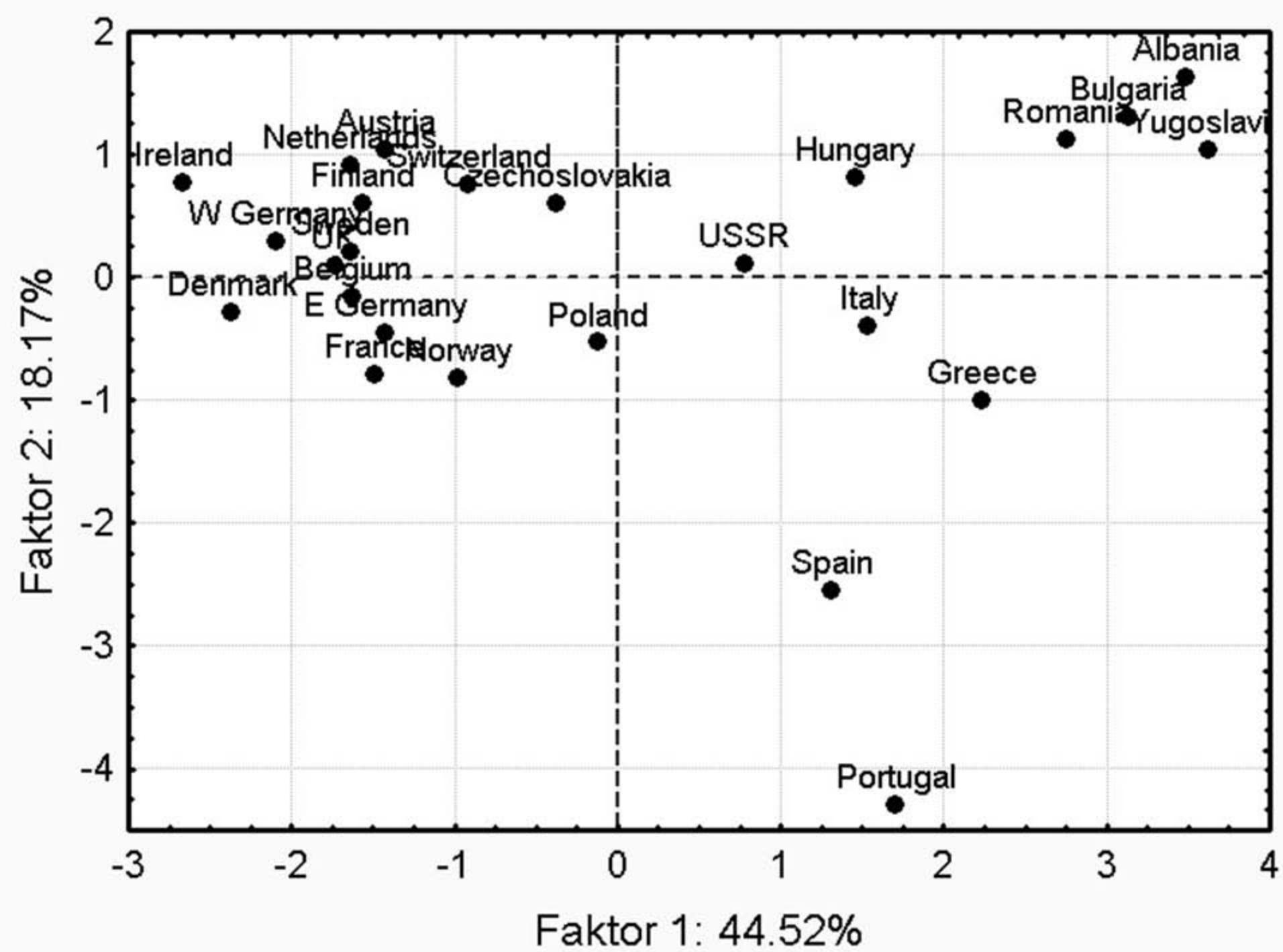
● Komponentní skóre souřadnice objektů v prostoru hlavních komponent:

Souřadnice každého objektu na osách hlavních komponent nazýváme *skóre*.

Graf komponentního skóre je zobrazení dvou skórových vektorů vnesených v systému kartézských os jeden proti druhému.

Pravidlo k volbě grafů:

1. Na x -ové ose užijeme vždy stejnou hlavní komponentu (obvykle první) u všech grafů komponentního skóre: t_1 proti t_2 , t_1 proti t_3 , t_1 proti t_4 , t_1 proti t_5 , ... atd., takže budeme vyšetřovat ostatní hlavní komponenty proti stále stejné, první.
2. Užijeme první hlavní komponentu, protože vykazuje největší hodnotu rozptýlení pro vyšetřovanou úlohu a vyneseme ji na x -ovou osu.

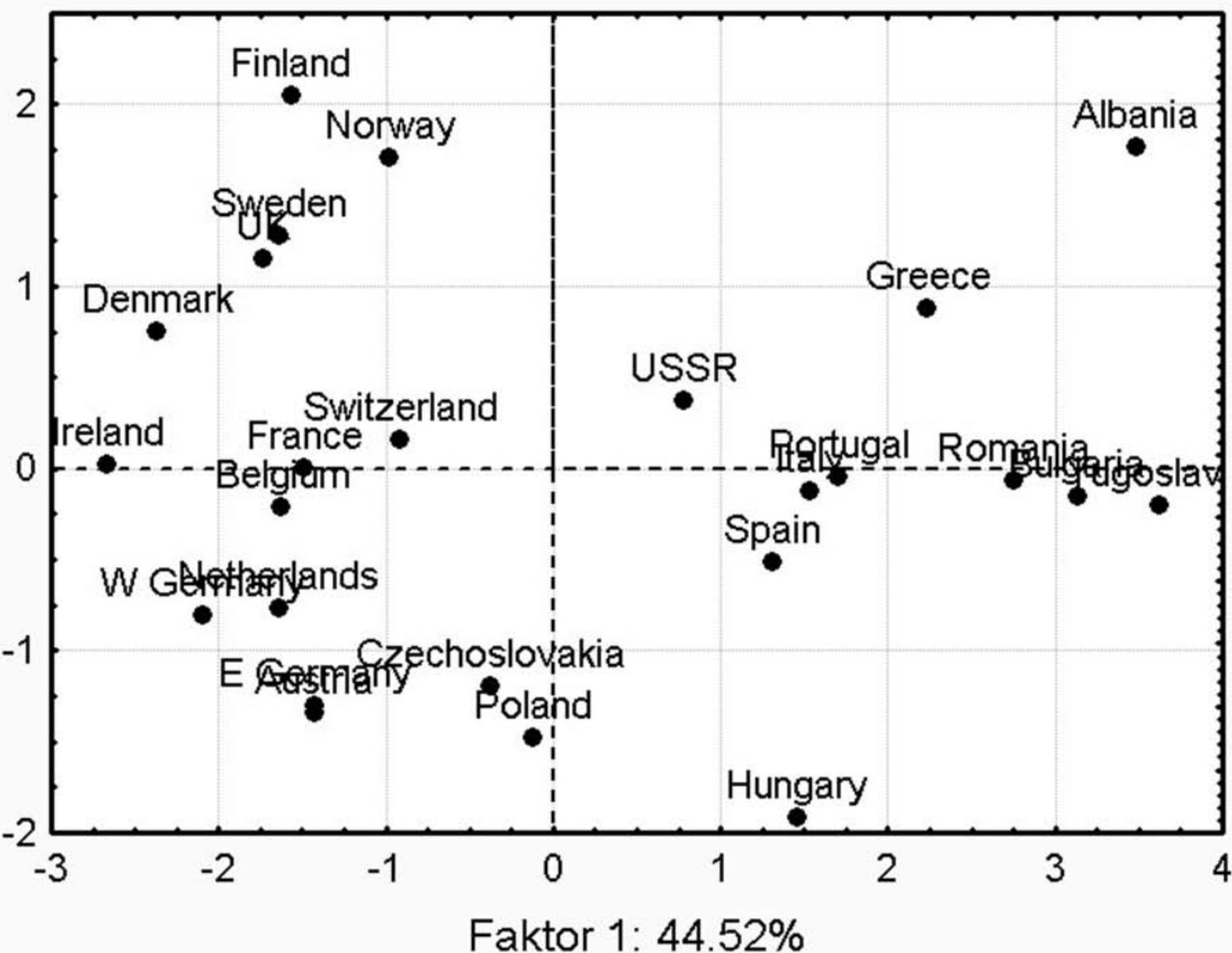


Rozptylový diagram komponentního skóre (Scatterplot)

- 1. Umístění objektů:** daleko od počátku jsou extrémy. Objekty nejblíže počátku jsou nejtypičtější.
- 2. Podobnost objektů:** objekty blízko sebe si jsou podobné, daleko od sebe jsou si nepodobné.
- 3. Objekty v shluku:** umístěné zřetelně v jednom shluku jsou si podobné a nepodobné objektům v ostatních shlucích. Jsou-li shluky blízko sebe, znamená to značnou podobnost objektů.

- 4. Osamělé objekty:** izolované objekty mohou být odlehlé.
- 5. Odlehlé objekty:** ideálně bývají objekty rozptýlené po celé ploše diagramu. V opačném případě je špatný model.
- 6. Pojmenování objektů:** výstižná jména objektů slouží k hledání hlubších souvislostí mezi objekty a vystihneme tak jejich fyzikální či biologický vztah.
- 7. Vysvětlení místa objektu:** umístění objektu na ploše v diagramu může být porovnáváno s komponentními vahami původních proměnných ve dvojném grafu.

Faktor 3: 12.53%



Rozptylový diagram komponentního skóre (Scatterplot)

Zobrazuje: *komponentní skóre* čili hodnoty obyčejně prvních dvou hlavních komponent u všech objektů.

Dokonalé rozptýlení objektů v rovině: vede k rozlišení objektů při jejich popisu pomocí y_1 a y_2 .

Lze nalézt: shluk vzájemně podobných objektů a dále objekty odlehlé a silně odlišné od ostatních.

Využití: k identifikaci odlehlých objektů, identifikaci trendů, tříd, shluků objektů, k objasnění podobnosti objektů atd.

Nemožné analyzovat všechny diagramy: jich je mnoho:
např. $m < n$

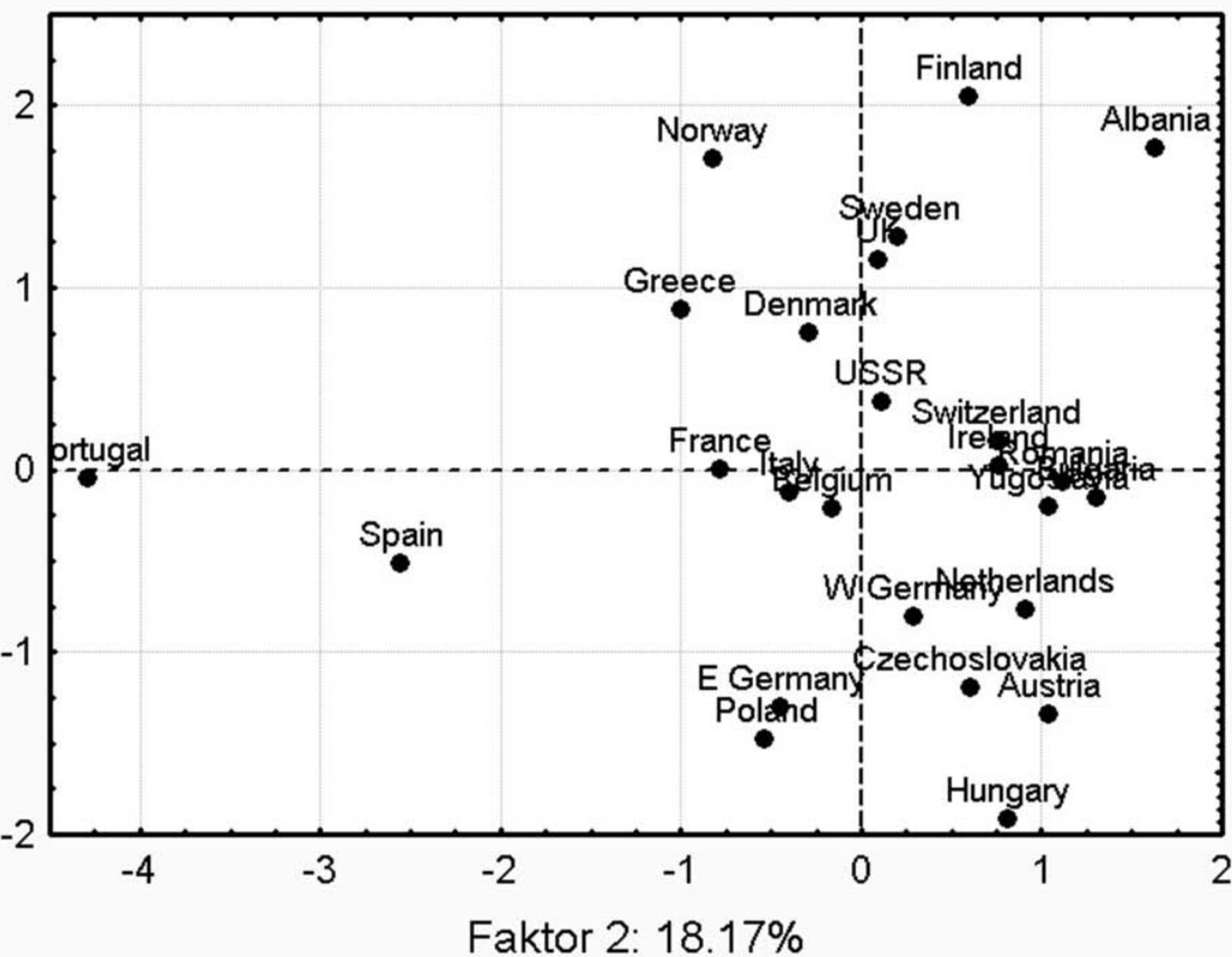
pro $m = 10$ existuje $m(m-1)/2 = 45$ diagramů,

pro $m = 11$ pak 55 diagramů,

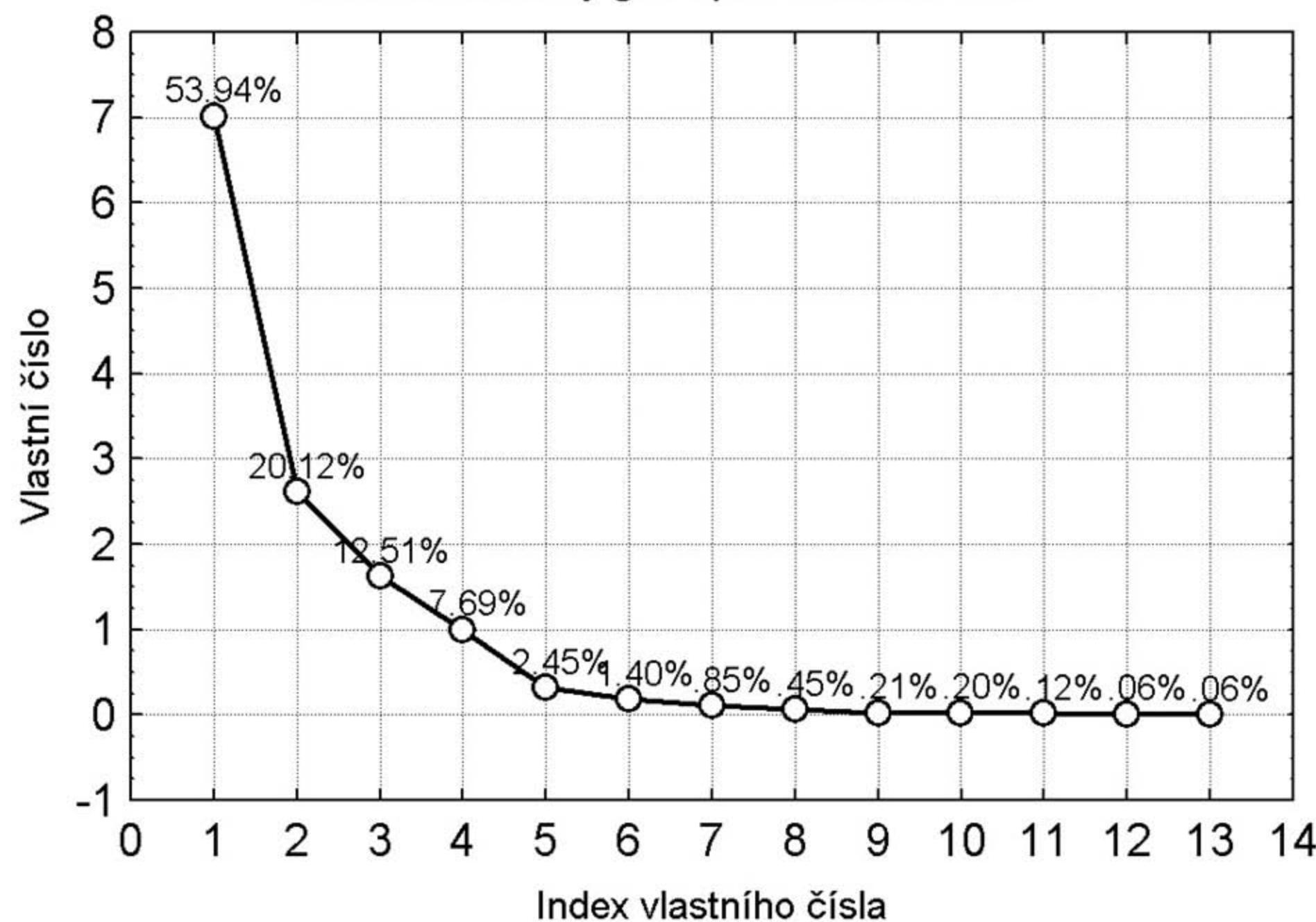
pro $m = 12$ pak 66 diagramů, atd.

Vybíráme: diagramy y_1 vs. y_2 , y_1 vs. y_3 , y_1 vs. y_4 atd. - držíme se první hlavní komponenty y_1 , protože v ní bývá největší míra proměnlivosti v datech.

Faktor 3: 12.53%



Cattelův indexový graf úpatí vlastních čísel



Výklad:

Indexový graf úpatí vlastních čísel (Scree Plot)

Je to sloupcový diagram vlastních čísel proti indexu A .

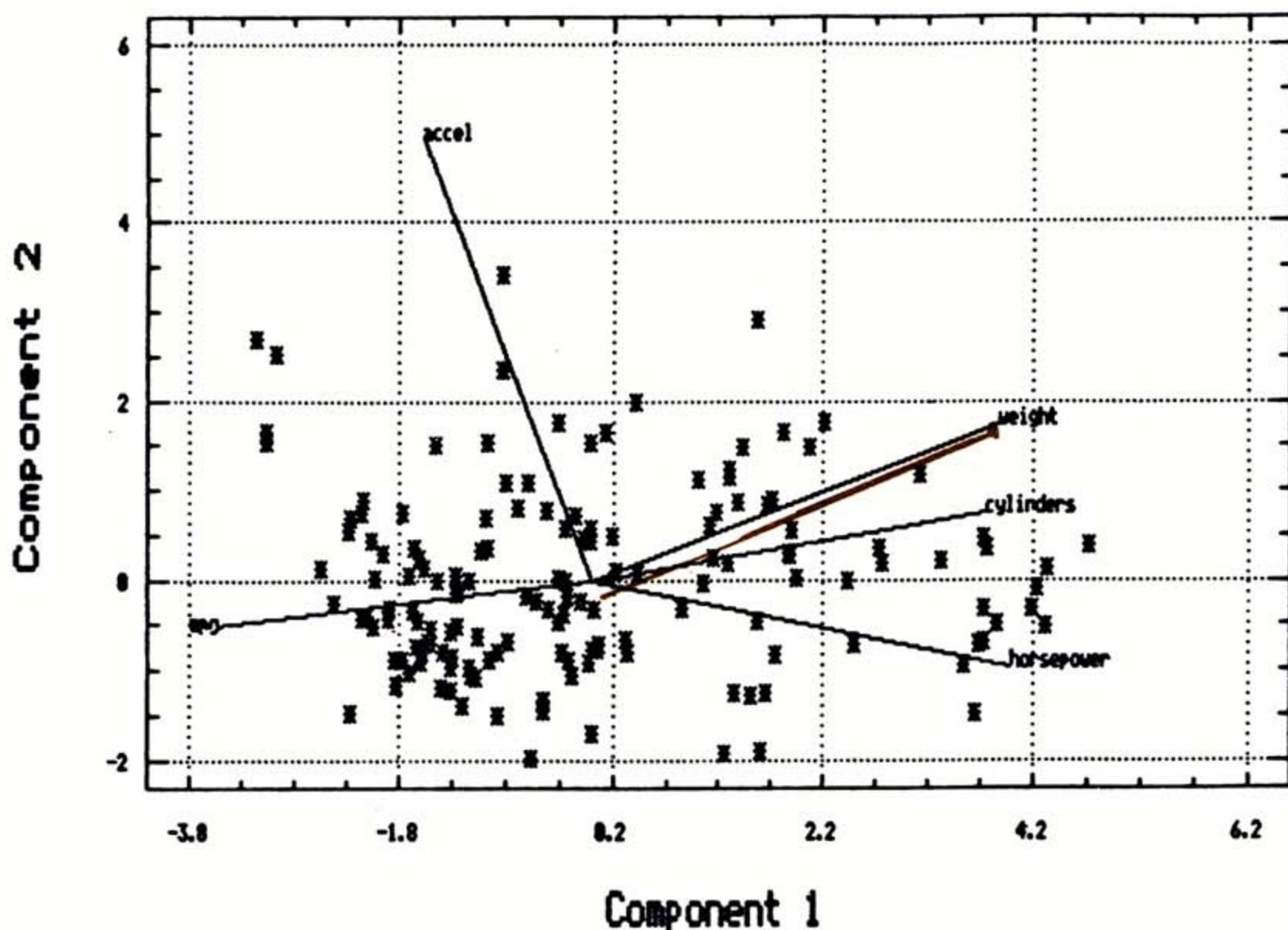
Zobrazuje: relativní velikost jednotlivých vlastních čísel.

Využití: k určení počtu A "užitečných" hlavních komponent.

Graf úpatí se jeví neobjektivnějším kritériem.

Kritérium "1": hrubším kritériem PC, jejichž vlastní číslo je větší než jedna. Graf úpatí se však jeví objektivnějším.

Biplot for First Two Principal Components



Dvojný graf (Biplot) kombinuje oba předchozí grafy.

Úhel mezi průvodiči x_j a x_k : je nepřímo úměrný velikosti korelace mezi proměnnými x_j a x_k : *Cím je menší úhel, tím je větší korelace.*

Délka souřadnice bodu: je úměrná příspěvku x_j do hlavní komponenty, čili je úměrná komponentní váze.

Interakce: objekt je v blízkosti proměnné x_j , tzn. objekt "obsahuje" hodně této proměnné.

Umístění vpravo (nahoře) od nuly: Interakce proměnných a objektů vysvětuje umístění objektů vpravo od nuly na ose y_1 (či vlevo od nuly), resp. umístění nahoře od nuly (či dole od nuly) na ose y_2 .

Diagnostika metody PCA

Maticový graf rozptylových diagramů znaků slouží k získání počáteční informace o datech, zda data potřebují škálování. V PCA postupně provádíme:

1. Vyšetření indexového grafu úpatí vlastních čísel – z hrany úpatí v tomto diagramu se určí vhodný počet hlavních komponent.

2. Výpočet vlastních vektorů – vedle číselných hodnot se užívá i názorný čárový diagram hodnot vlastních vektorů, který přehledně informuje o relativním zastoupení původních znaků x_j , $j = 1, \dots, m$, v hlavních komponentách.

3. Výpočet komponentních vah – matice párových korelačních koeficientů obsahující korelace původních znaků s hlavními komponentami. Uživatel nyní vybere pouze prvních k hlavních komponent a vytvoří tak model PCA.

4. Vyšetření grafu komponentních vah.

5. Vyšetření rozptylového diagramu komponentního skóre.

6. Vyšetření dvojného grafu.

7. Vyšetření reziduí – rezidua objektů a rezidua proměnných by měla prokazovat dostatečnou těsnost proložení.

8. Určení významných původních znaků – je výhodné vyhledávat významné znaky, protože klasická metoda PCA umožňuje sice redukci počtu hlavních komponent, ale každá komponenta zůstává stále kombinací všech původních znaků.

Příklad 4.1/str. 70

This screenshot shows the NCSS software interface for Principal Components Analysis. A data spreadsheet is open on the left, showing 10 rows of data with columns labeled MA, Osoba, Vyska, Hmotnost, Vlasy, Boty, Vek, Prijem, Pivo, Vino, Sex, Plavani, Puved, IQ, C14, C15, C16, C17, C18, and C19.

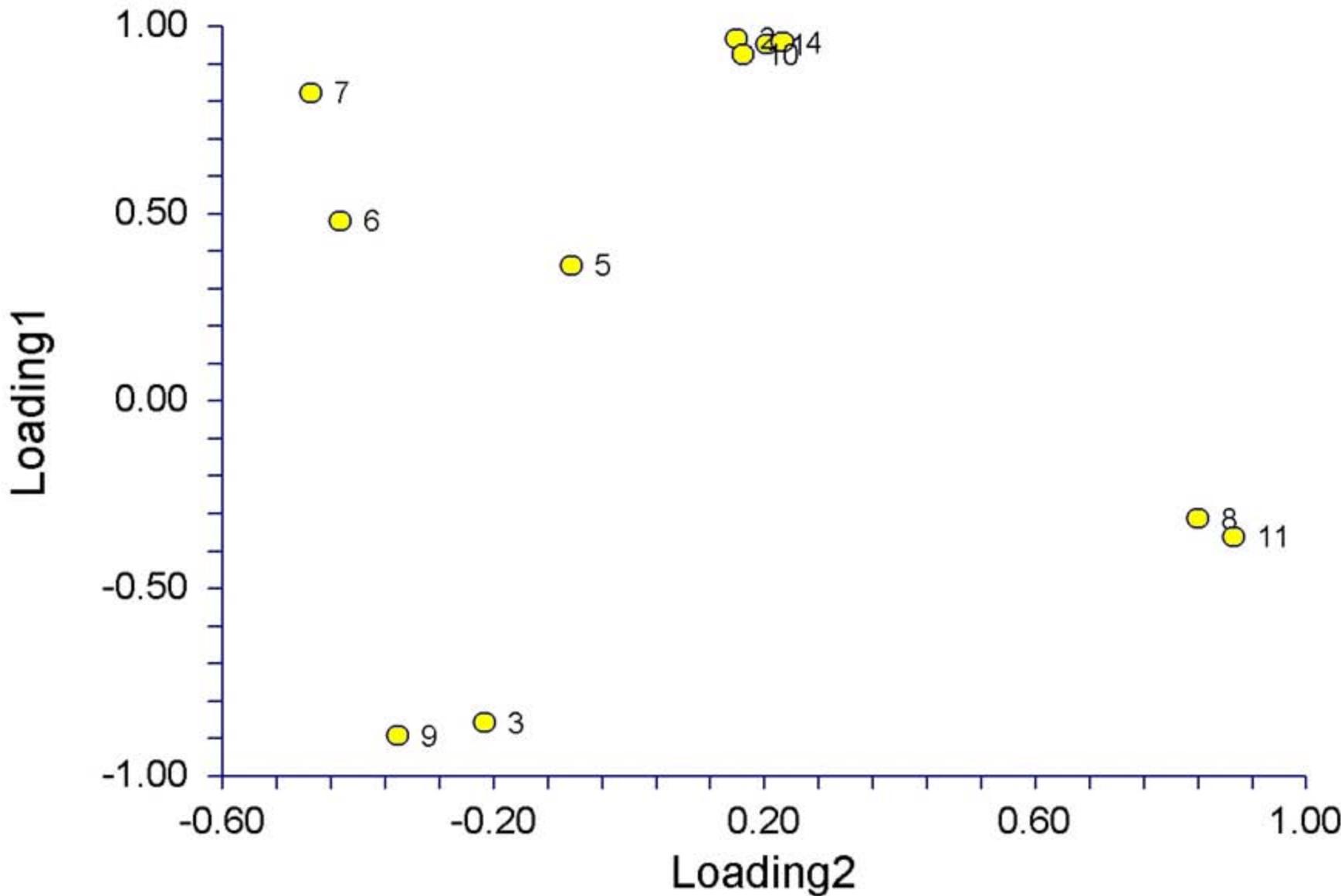
The main window displays the "NCSS: Principal Components Analysis" dialog. Key settings visible include:

- Input Variables:** Variable selected is "VYSLKA-PUVOD". Data input format is "Regular Data".
- Covariance Estimation Options:**
 - Robust Covariance Matrix Estimation: Unchecked.
 - Robust Weight: Value 4.0.
 - Missing Value Estimation: None.
 - Maximum Iterations: 1000.
- Type of Matrix Used in Analysis:**
 - Matrix Type: Correlation.
- Factor (Component) Options:**
 - Factor Rotation: None.
 - Factor Selection - Method: Percent of Eigenvalues.
 - Factor Selection - Value: 100.

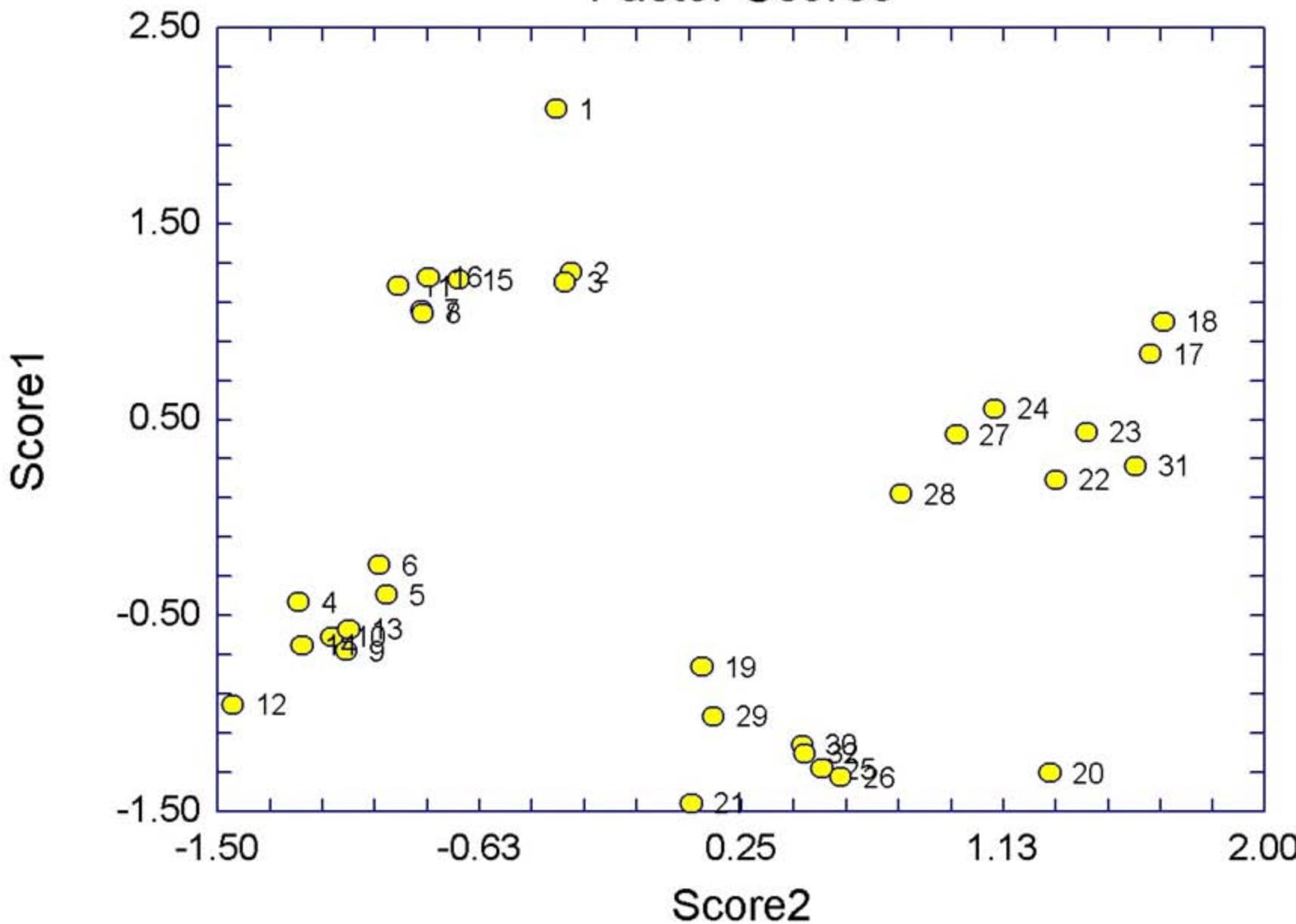
A context menu is open on the right side of the dialog, with the "QUICK ACCESS BUTTON" option highlighted. The menu includes instructions for loading procedures and assigning them to buttons.

The status bar at the bottom indicates: "This is the spreadsheet that lets you enter and edit your data."

Factor Loadings



Factor Scores



Příklad 4.2/str. 76

NCSS: Principal Components Analysis

Input Variables:

Variables: Aro-Ztr Data Input Format: Regular Data

Covariance Estimation Options:

Robust Covariance Matrix Estimation Robust Weight: 4.0

Missing Value Estimation: None Maximum Iterations: 1000

Type of Matrix Used in Analysis:

Matrix Type: Correlation

Factor (Component) Options:

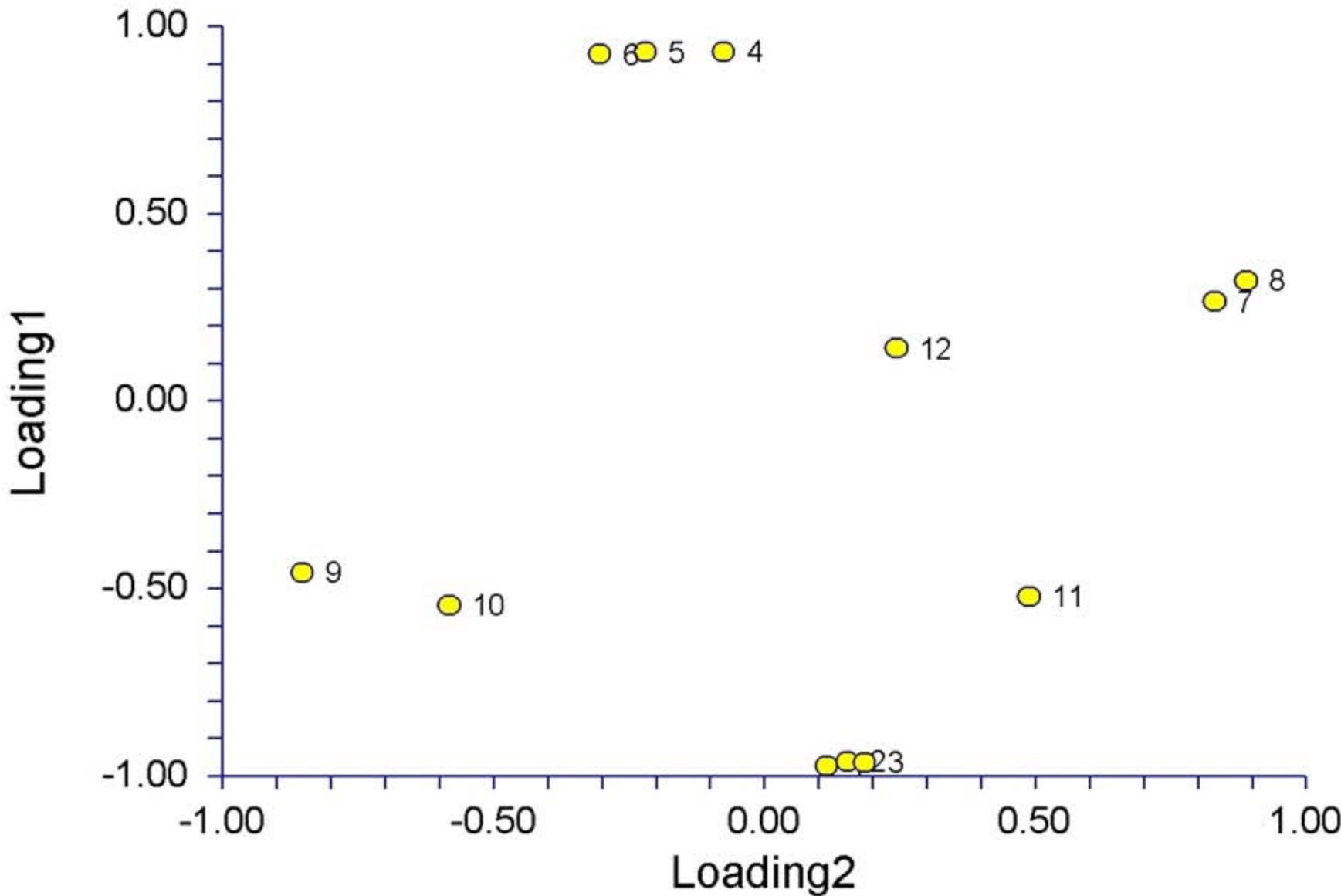
Factor Rotation: None Factor Selection - Method: Percent of Eigenvalues Factor Selection - Value: 100

Opt 2 Template Id: 1

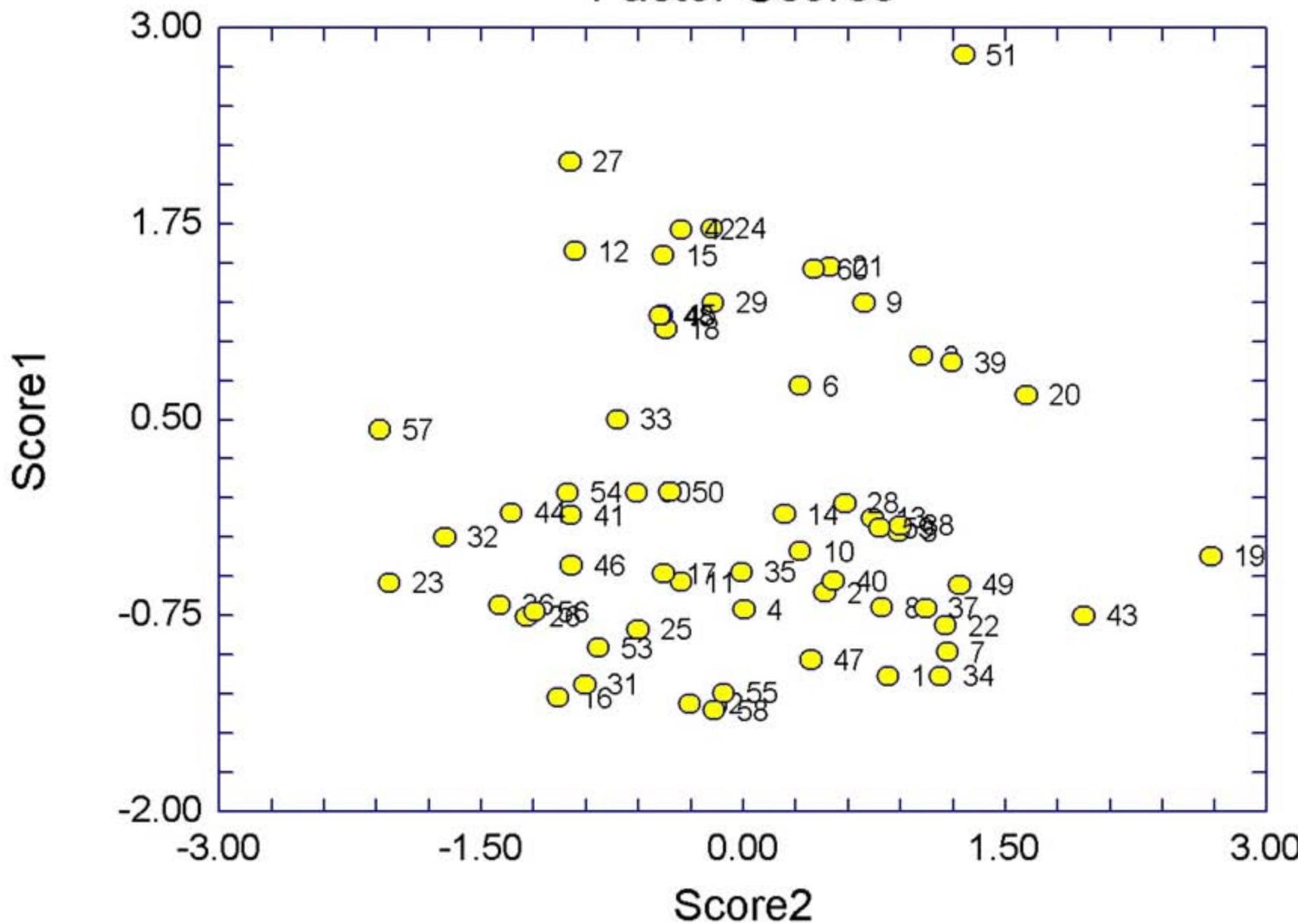
Reset Guide Me

Objekt	Aro	Slad	Med	Bez	Klas	Tvrđ	Bel	Bar1	Bar2	Bar3	Slup	Ztr	C14	C15	C16	C17	C18	C19
B5	6.48	6.66	4.56	2.2	2.91	3.47	4.72	5.585	5.735	5.985	4.26	3.25						
C4	5.75	6.09	3.81	2.32	4.03	3.77	4.17	5.73	5.745	5.325	3.82	3.38						
B2	3.94	4.12	2.44	3.63	6.77	5.39	4.77	6.665	5.105	4.595	3.5	3.03						
D5	6.6	6.17	4.44	1.93	6.71	4.46	4.46	5.74	5.444	5.12	3.12	3.04						

Factor Loadings



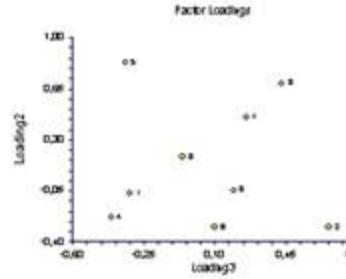
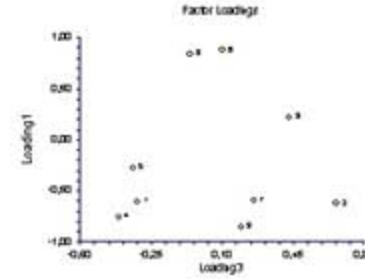
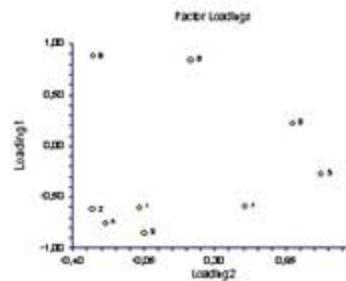
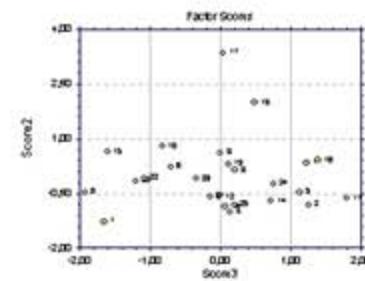
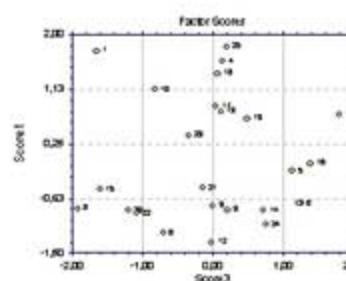
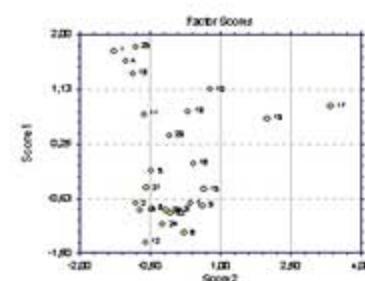
Factor Scores



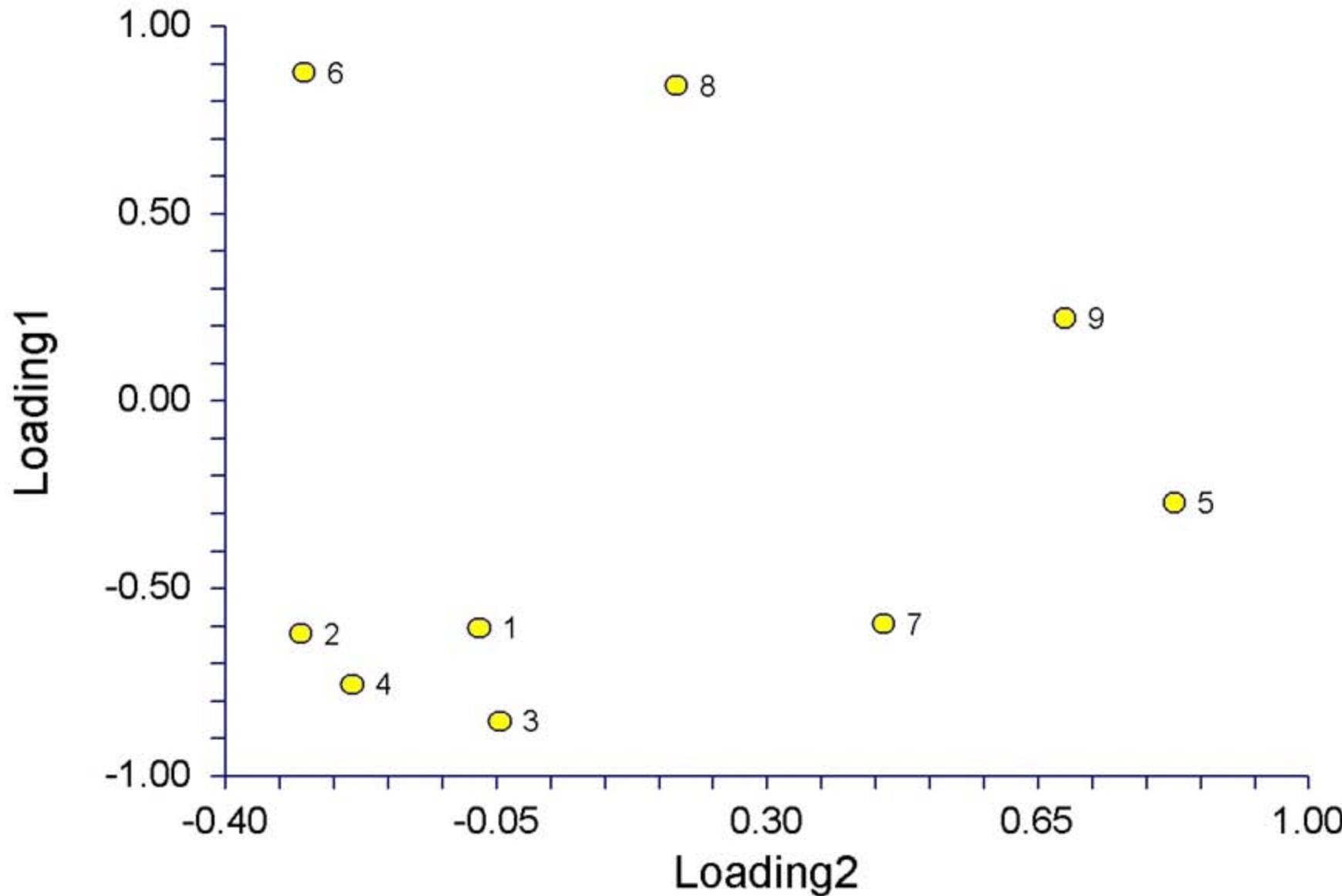


17	0.0892	-0.0012	0.8369
18	-0.2554	0.0789	1.0737
19	-1.1459	-0.8901	-1.5143
20	-1.0384	0.2118	-0.3443
21	0.9830	-1.5521	0.2114
22	-1.2491	0.7015	-0.4191
23	1.1016	0.1526	0.2912
24	-0.8765	-1.0458	-0.0803
25	0.1176	0.5668	0.4740

Plots Section

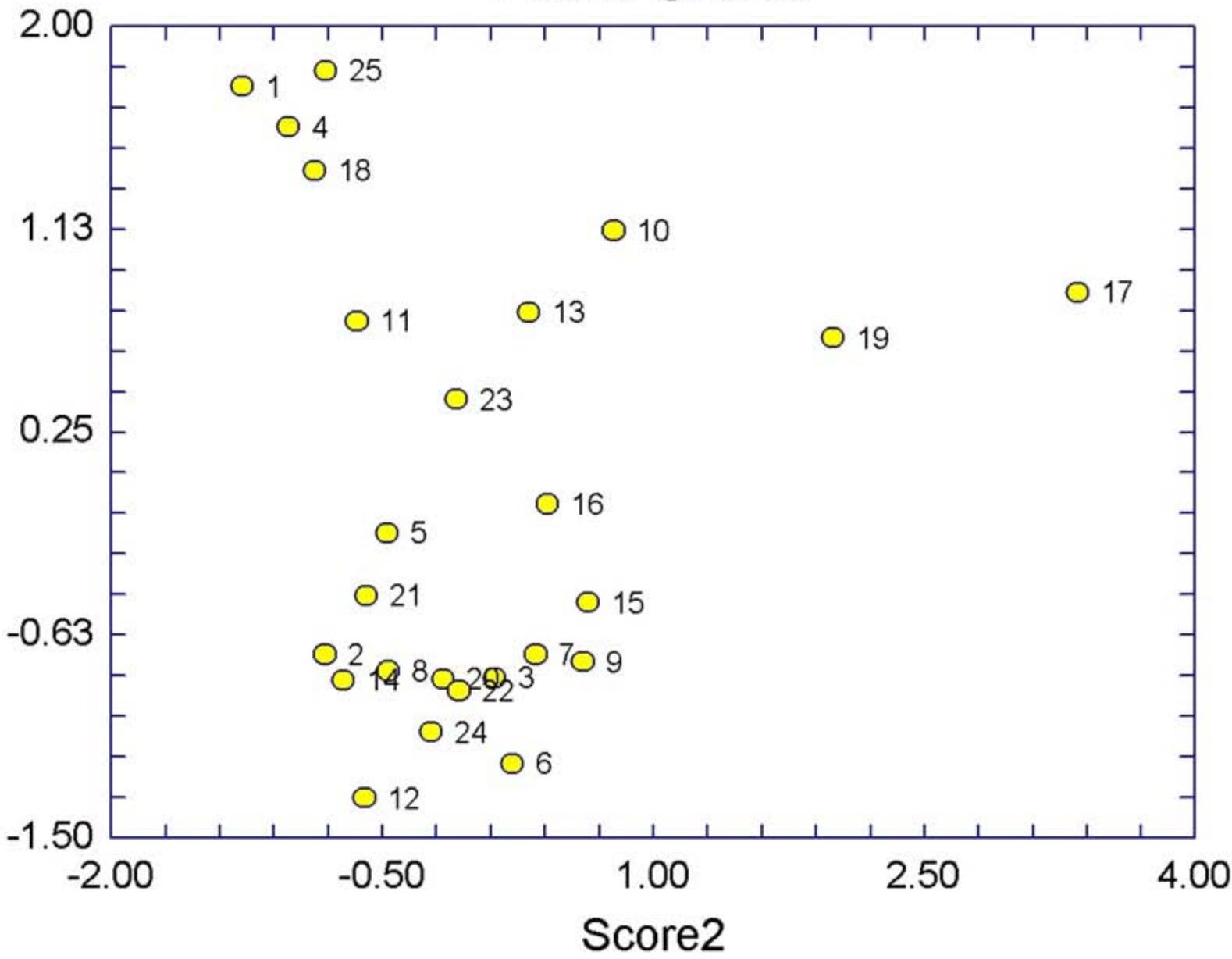


Factor Loadings



Factor Scores

Score1



Score2

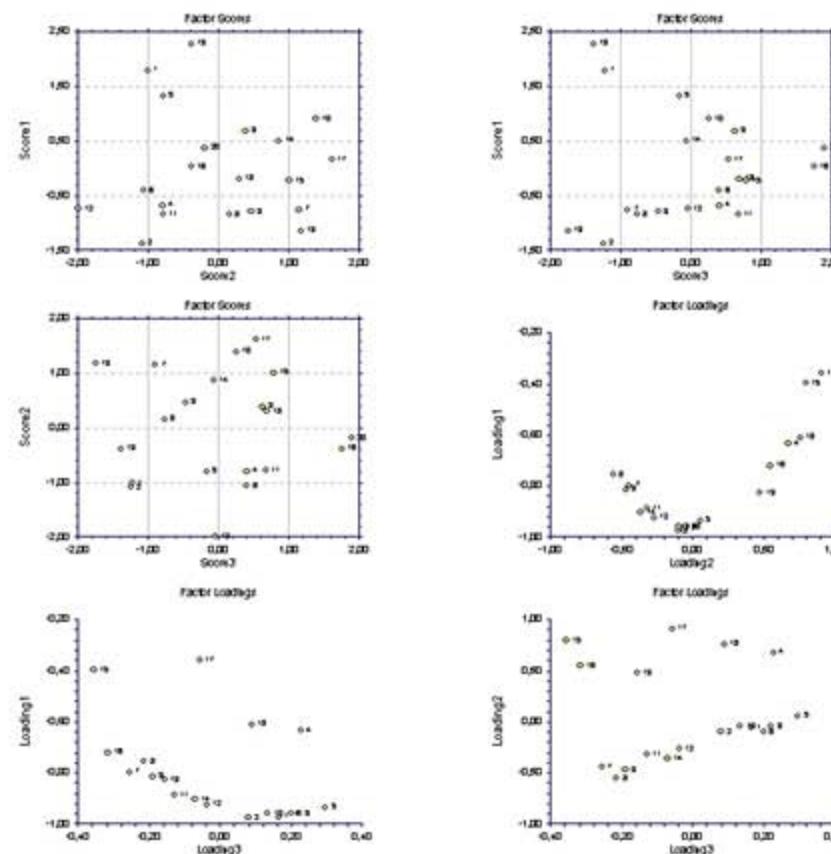


18	0.0562	-0.5232	-0.1865	-1.5234	-0.5220	-0.1620
19	0.0517	-0.1654	0.5658	1.2474	0.8854	0.7163
20	-0.9155	0.3822	-0.2885	0.7986	-0.7211	1.1154

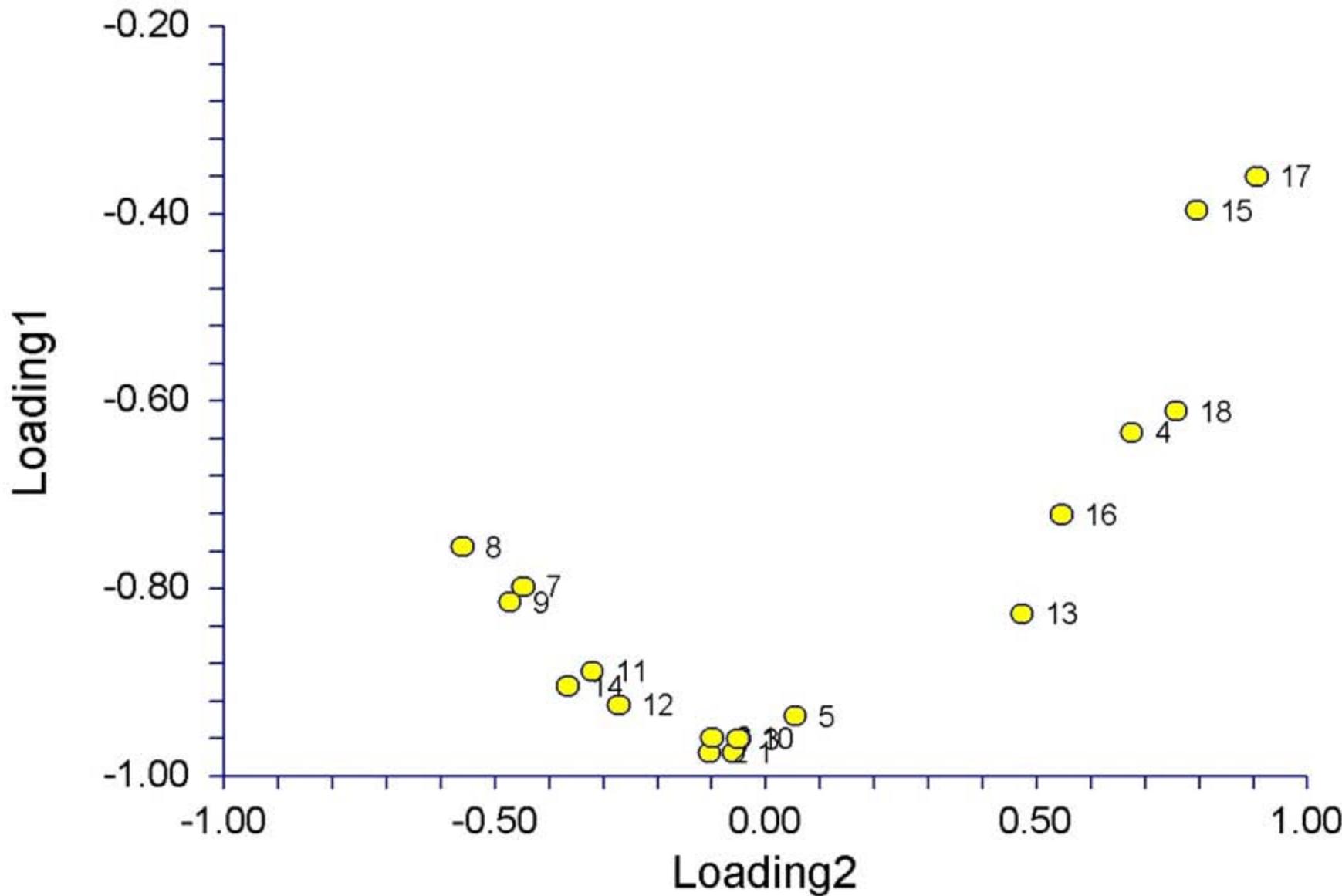
Principal Components Report

Page/Date/Time 24 31.8.2009 13:12:17
 Database E:\Moje\Učebnice 2005 M+M+H\... ATANCSS2002\P405Guiseppe SO

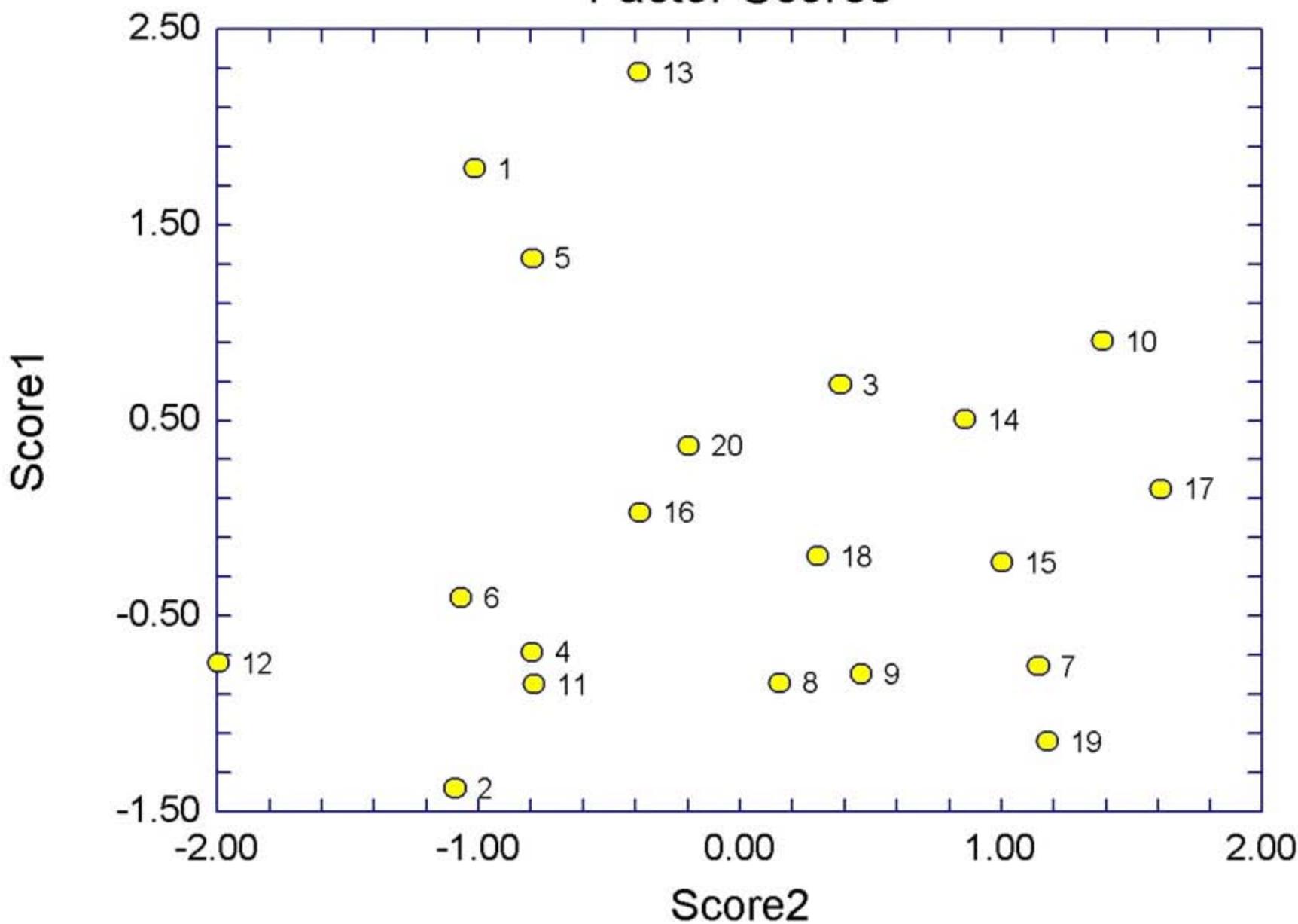
Plots Section



Factor Loadings



Factor Scores





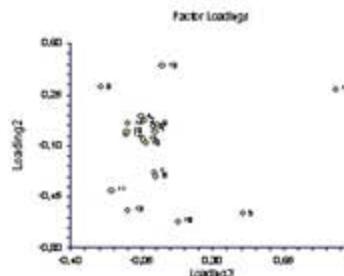
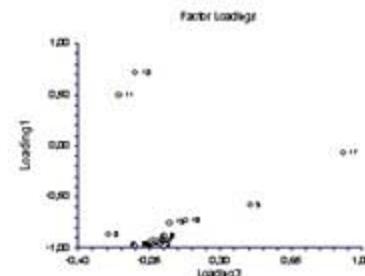
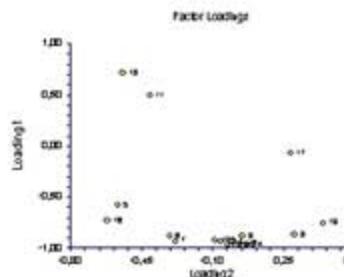
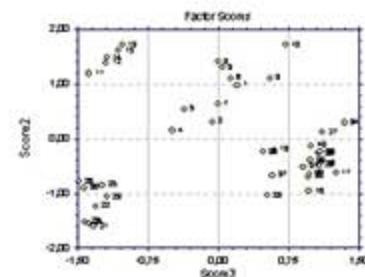
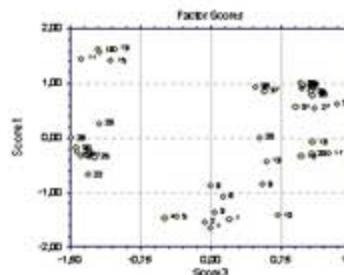
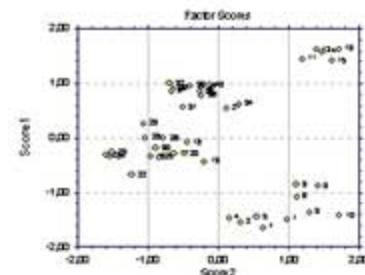
36	-2.1314
37	0.0988
38	-0.5884
39	-1.4358
40	0.6177

Principal Components Report

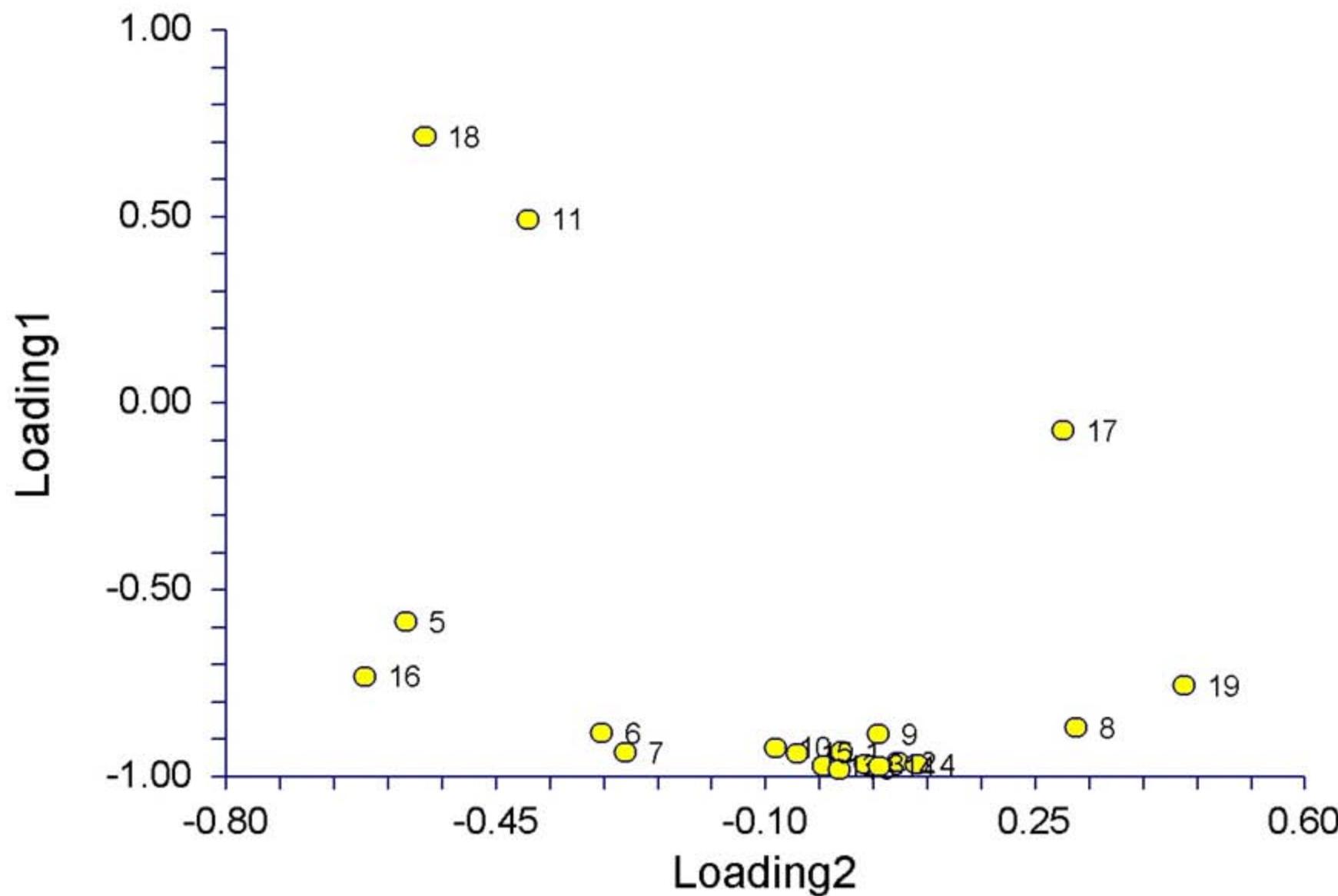
Page/Date/Time 27 31.8.2009 13:13:43

Database E:\Moje\Učebnice 2005 M+M+H\... e\DATA\NCSS2002\P406Msice.S0

Plots Section

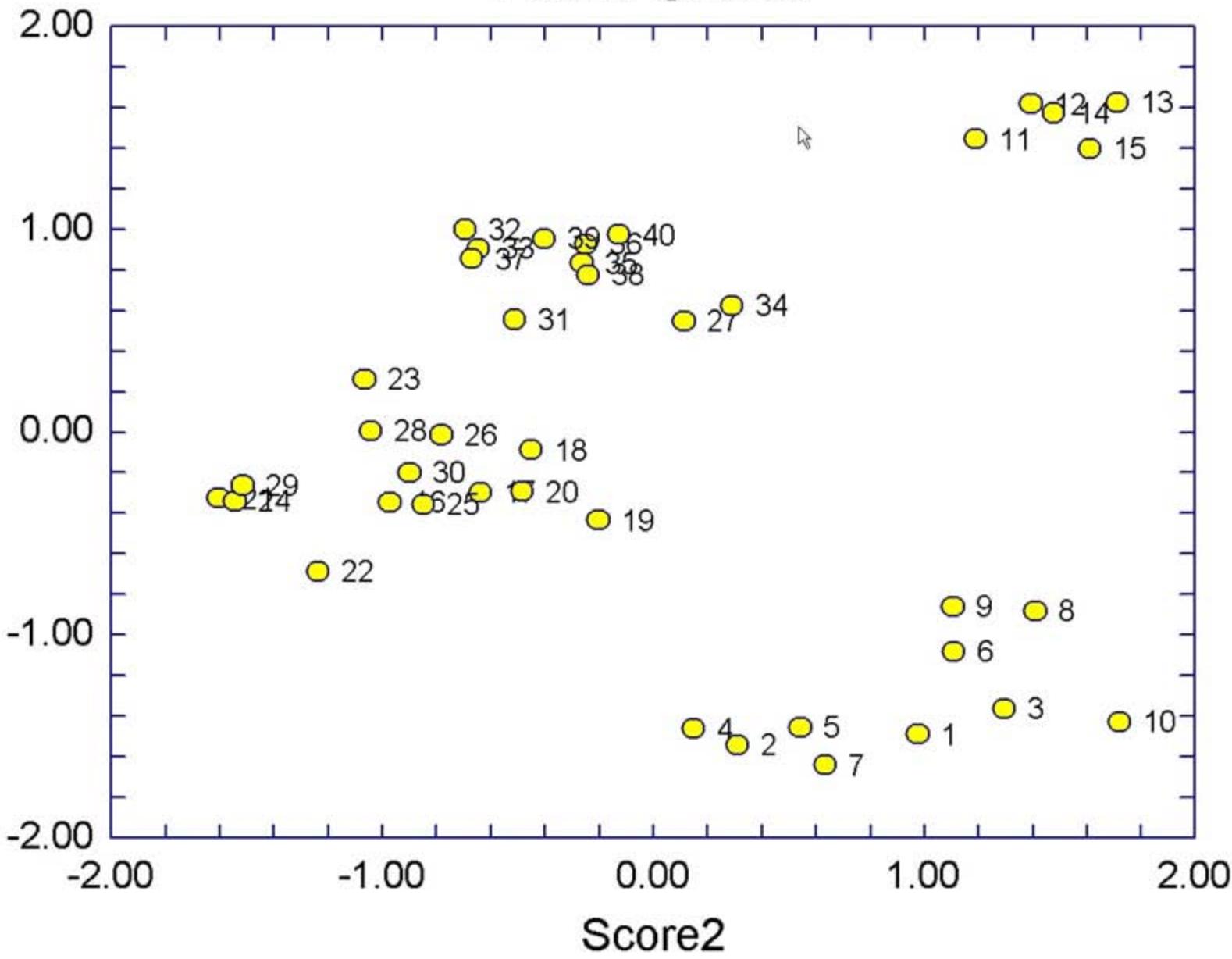


Factor Loadings



Factor Scores

Score1



Score2

