

Klasifikace objektů

hledání struktury a vzájemných vazeb v objektech

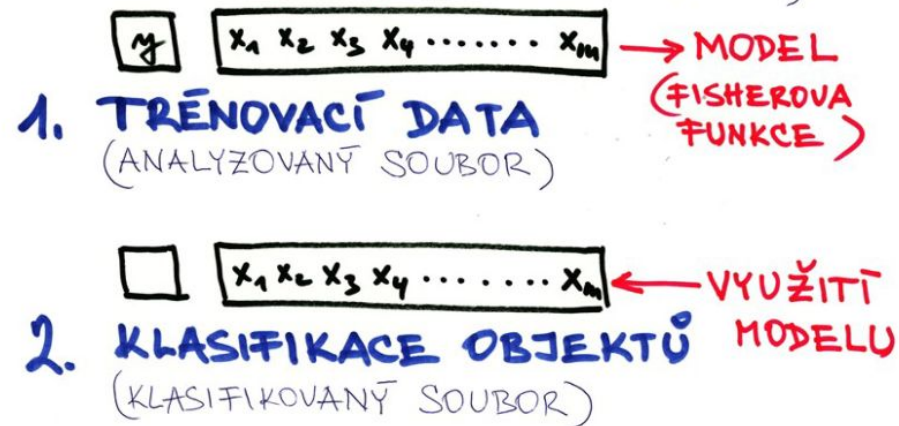
- 1) **Diskriminační analýza DA:** nový objekt se zařadí do již existující třídy.
- 2) **Analýza shluků CLU:** neuspořádanou skupinu objektů lze uspořádat do několika vnitřně sourodých tříd či shluků
- 3) **Vicerozměrné škálování MDS:** hledá strukturu a vazby mezi objekty na základě jejich podobnosti.

Analyzovaný výběr jsou trénovací data, která obsahují pro každý objekt jak výstup y , tak i hodnoty všech znaků x , $\{y; x_1, x_2, \dots, x_m\}$

Klasifikovaný výběr: na základě analyzovaného výběru sestavit predikční model, který umožní zařazení nových objektů do tříd.

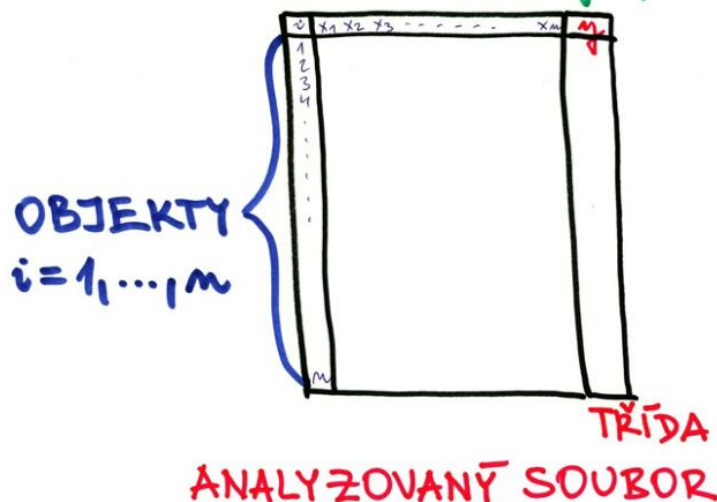
DISKRIMINAČNÍ ANALÝZA

TRÉNOVÁNÍ S UČITELEM (SUPERVISED LEARNING)



DISKRIMINAČNÍ ANALÝZA

DISKRIMINÁTORY, $j=1, \dots, m$



Příklad: použití LDA pro dvě třídy $m=2$. Vychází se ze známých matic:

X_1 rozměru $n_1 \times m$ pro třídu 1, X_2 rozměru $n_2 \times m$ pro třídu 2.

Jednotlivé objekty v matici X všech dat se zařadí do tříd podle výstupu G .

Postup:

1) Vychází se výběrové průměry \bar{x}_1 a \bar{x}_2 a společná kovarianční matice

$$S = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2}$$

2a) Nejjednodušší je předpoklad $\pi_1 = \pi_2 = 0.5$.

2b) Pokud je výběr informativní a byl pořízen jako celek a pak rozdělen do skupin, je možné použít relativních četností

$$\hat{\pi}_1 = n_1 / (n_1 + n_2) \quad \text{a} \quad \hat{\pi}_2 = n_2 / (n_1 + n_2)$$

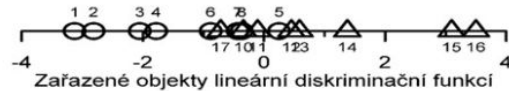
3) Za normality se určí koeficienty Fisherovy lineární diskriminační funkce z odhadů

$$\hat{a} = (\bar{x}_1 - \bar{x}_2) S \quad \text{a} \quad \hat{b} = -0.5 a^T (\bar{x}_1 - \bar{x}_2) - \ln(\hat{\pi}_2 / \hat{\pi}_1)$$

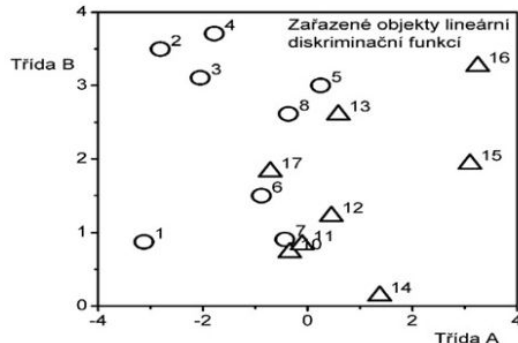
4) Při zařazování nových objektů s hodnotami znaků x_0 se použije pravidlo, že do první skupiny je objekt klasifikován pokud platí $\hat{a}^T x_0 + \hat{b} \geq 0$.

5) V opačném případě se klasifikuje do druhé skupiny.

Ukázka 1: zařazování Fisherovu lineární diskriminační funkcí:



Ukázka 2: zařazování Fisherovu lineární diskriminační funkcí:



1) **Analyzovaný výběr:** na trénovacích datech vytvoříme třídy
První třída (13 lebek z hrobů v Sikkimu) vede ke středním hodnotám
 [174.82, 139.35, 132.00, 69.82, 130.35]

a kovarianční matici $S_1 = \begin{bmatrix} 45.53 & & & & \\ 25.22 & 57.81 & & & \\ 12.39 & 11.88 & 36.09 & & \\ 22.15 & 7.52 & -0.31 & 20.94 & \\ 27.97 & 48.06 & 1.41 & 16.77 & 66.21 \end{bmatrix}$.

Druhá třída (15 lebek z bojišť v Lhase) vede ke středním hodnotám
 [185.73, 138.73, 134.77, 76.47, 137.50]

a kovarianční matici $S_2 = \begin{bmatrix} 74.42 & & & & \\ -9.52 & 37.35 & & & \\ 22.74 & -11.26 & 36.32 & & \\ 17.79 & 0.70 & -10.72 & 15.30 & \\ 11.13 & 9.46 & 7.20 & 8.66 & 17.96 \end{bmatrix}$.

Příklad 4.22 Třídění lebek Tibeťanů lineární diskriminační funkcí

Databáze lebek na pohřebištích v Tibetu svědčí o dvou skupinách lidí:

prvních 13 bylo nalezeno v hrobech v Sikkimu a okolí,

druhých 15 lebek na bojištích okolo Lhasy.

Předpokládejme, že máme data o 2 třídách tibetských lebek.

Data: i index lebky,

x_1 největší délka lebky [mm],

x_2 největší horizontální šířka lebky [mm],

x_3 výška lebky [mm],

x_4 výška horní části obličeje [mm],

x_5 šířka obličeje mezi body lícních kostí [mm].

x_1	x_2	x_3	x_4	x_5
190.5	152.5	145	73.5	136.5
...
...
172.5	132	125.5	63	121
167	130	125.5	69.5	119.5
1825	131	135	68.5	136

Řešení:

Koeficienty diskriminační funkce a_1, \dots, a_5 , jsou vyčísleny podle vztahu

$$\mathbf{a} = S^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) = [-0.09, 0.16, 0.01, -0.18, -0.18]$$

a vedou k průměrům obou tříd: $\bar{Z}_1 = -28.71$ a $Z_2 = -32.21$.

Optimální prahový bod C ,

dle kterého se budou nezařazené objekty třídit do první nebo druhé třídy,
 se vyčíslí jako polosuma obou průměrů dle vztahu

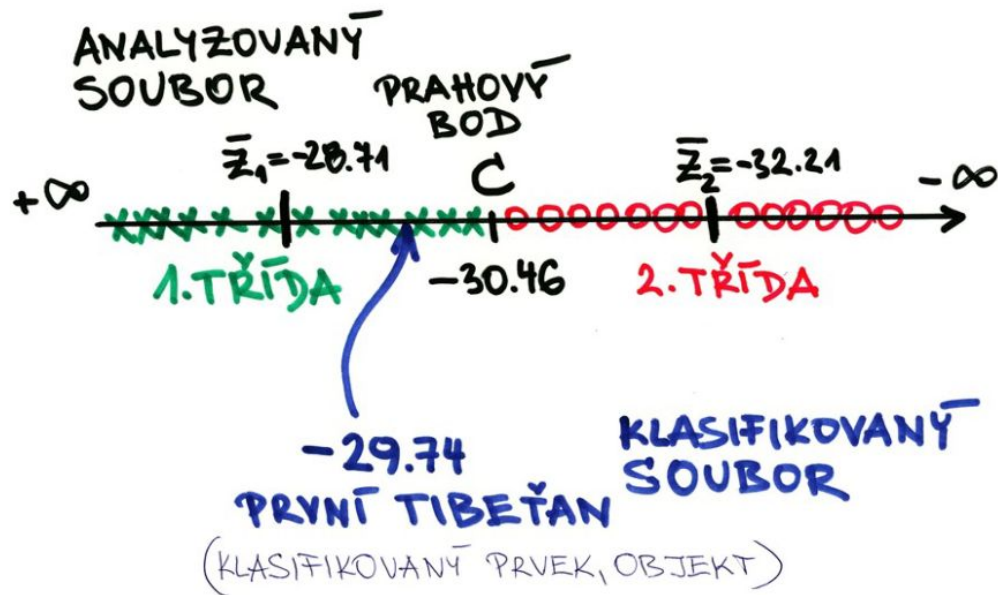
$$C = (\bar{Z}_1 + \bar{Z}_2)/2 = (-28.71 + (-32.21))/2 = -30.46.$$

Diskriminace:

2) **Klasifikovaný výběr:** na nových datech provedeme zařazení dosud nezařazených objektů
 Vezmeme data pro lebku prvního Tibeťana a pokusíme se ji zařadit do 1. nebo do 2. třídy.
 Vyčísleme proto pro ni hodnotu lineární diskriminační funkce

$$Z_1 = -0.09 \times 190.5 + 0.16 \times 152.5 + 0.01 \times 145.0 - 0.18 \times 73.5 - 0.18 \times 136.5 = -29.74,$$

Závěr: Protože lineární diskriminační funkce $Z_1 = -29.74$ je menší než optimální prahový bod $C = -30.46$, patří lebka prvního Tibeťana do první třídy.



Úprava prahového bodu

Volba prahového bodu C poskytuje požadovaný poměr apriorních pravděpodobností π_1 a π_2 .

Optimální volba prahového bodu C je daná vzorcem $C = \frac{\bar{Z}_1 + \bar{Z}_2}{2} + \ln \frac{\pi_1}{\pi_2}$

a když bude $\pi_1 = \pi_2 = 0.5$, bude prahový bod C roven $C = \frac{\bar{Z}_1 + \bar{Z}_2}{2}$

Standardizované koeficienty: hodnoty koeficientů a_1, a_2, \dots, a_p **nejsou** přímo porovnatelné.

Relativní vliv na každou proměnnou v diskriminační funkci získáme ze *standardizovaných diskriminačních koeficientů*.

Standardizované koeficienty se vypočtou vynásobením koeficientů a_i odpovídající směrodatnou odchylkou s_i .

Volba znaků, diskriminátorů

Znaky musí zajistit přesné zařazení objektů do tříd čili diskriminaci.

Principem selekce znaků je zajištění dostatečné separability tříd a maximalizace některé zvolené míry:

- Začneme se všemi znaky. Pak se a vypouštějí se takové znaky, které provedou nedostatečnou separaci.
- V mnoha situacích může být DA užita jako exploratorní pomůcka.
- Nejprve jsou do dat zahrnuty všechny využitelné znaky. Na začátku není známo, které znaky jsou k zařazení objektů do tříd účinné.
- V diskriminační analýze místo testování změny hodnoty čtverce korelačního koeficientu R^2 přidáním nebo odebráním proměnné testujeme změnu Mahalanobisovy vzdálenosti D_M^2 .
- Užívají se stejná testovací kritéria jako při výběru nezávisle proměnných v lineární regresní analýze.
- Krokový výběr znaků* kombinuje jak jejich přidávání, tak i jejich odstraňování.
- První znak, zahrnutý do modelu ve výběrovém kritériu, má největší přijatelnou hodnotu.
- Vybírání znaků* se ukončí, když žádné další znaky nespĺňují zaváděcí nebo odstraňovací kritérium.

Kritéria k vybírání znaků

Wilkovo kritérium λ : když znak poskytuje v diskriminační funkci nejmenší hodnotu Wilkova kritéria λ , je zahrnut do modelu:

- k zavedení nebo odstranění znaku je dovozen jeden krok.
- maximální počet kroků je roven dvojnásobku počtu znaků.

- **tolerance** je mírou stupně lineární asociace mezi znaky. Pro i -tý znak platí vztah $1 - R_i^2$, kde R_i^2 je čtverec vícenásobného korelačního koeficientu, když je uvažován i -tý znak za závisle proměnnou a když je uvažována regresní rovnice mezi tímto i -tým znakem a ostatními znaky.

- **využití tolerance:** malé hodnoty tolerance indikují, že i -tý znak je lineární kombinací ostatních znaků. Znaky s tolerancí menší než 0.001 nejsou do modelu zařazeny.

Testování: významnost změny Wilkova kritéria λ po zavedení znaku do modelu nebo odstranění z modelu je založena na testacím kritériu F .

Na začátku procesu vybírání znaků: tolerance a minimum tolerance jsou položeny rovny 1, protože v modelu zatím nejsou znaky.

Výsledky diskriminační funkční analýzy (B23)						
Počet prom. v modelu: 5; grupovací: B23x6-GROUP (3 skup)						
Wilk. lambda: ,10863 přibliž F (10,272)=55,328 p<0,0000						
N=143	Wilk. Lambda	Parc. Lambda	F na vyj (2,136)	p-hodn.	Toler.	1-toler. R^2
B23x1-RELWT	0,117426	0,925075	5,50754	0,005012	0,745200	0,254800
B23x2-GLUFAST	0,149708	0,725594	25,71637	0,000000	0,149227	0,850773
B23x3-GLUTEST	0,199867	0,543500	57,11504	0,000000	0,143100	0,856900
B23x4-INSTEST	0,122473	0,886949	8,66731	0,000286	0,813954	0,186046
B23x5-SSPG	0,112112	0,968922	2,18109	0,116853	0,562221	0,437779

↓
Test statistické významnosti
dotvčného diskriminátoru

Výsledky diskriminační funkční analýzy (73Iris)						
Počet prom. v modelu: 4; grupovací: Trida (3 skup)						
Wilk. lambda: ,02344 přibliž F (8,288)=199,15 p<0,0000						
N=150	Wilk. Lambda	Parc. Lambda	F na vyj (2,144)	p-hodn.	Toler.	1-toler. R^2
SEPALLENGTH	0,024976	0,938464	4,72115	0,010329	0,347993	0,652007
SEPALWIDTH	0,030580	0,766480	21,93593	0,000000	0,608859	0,391141
PETALLENGTH	0,035025	0,669206	35,59018	0,000000	0,365126	0,634874
PETALWIDTH	0,031546	0,743001	24,90433	0,000000	0,649314	0,350686

Raovo V kritérium, známé jako **Lawley-Hotellingova stopa**, je

$$V = (n - g) \sum_{i=1}^m \sum_{j=1}^m w_{ij} \sum_{k=1}^g (\bar{x}_{ik} - \bar{x}_i)(\bar{x}_{jk} - \bar{x}_j)$$

kde p udává počet znaků v modelu,

g značí počet tříd,

n_k je velikost k -té třídy,

\bar{x}_{ik} je střední hodnota i -tého znaku v k -té třídě,

\bar{x} je průměr i -tého znaku pro všechny třídy kombinované dohromady a

w_{ij} představuje prvek inverzní kovarianční matice mezi třídami B^{-1} .

Testování: čím větší jsou rozdíly mezi středními hodnotami (průměry) tříd, tím větší je hodnota Raova V .

Může se však stát, že znak po zařazení do modelu sníží hodnotu Raova V .

F-test významnosti každého znaku: hodnota F pro změnu Wilkova kritéria

λ při přidání znaku do modelu se vyčíslí dle

$$F_{změny} = \frac{n - g - p}{g - 1} \left(\frac{1 - \lambda_{p+1}}{\lambda_p} \right)$$

kde n je celkový počet objektů,

g udává počet tříd,

p je počet znaků,

λ_p značí Wilkovo lambda před přidáním a

λ_{p+1} je Wilkovo lambda po přidání znaku do modelu.

Testování: Do modelu je zařazen ten znak, který způsobuje nejmenší hodnotu Wilkova kritéria λ

Mahalanobisova vzdálenost $D_{1,2}^2$ je zobecněná míra vzdálenosti mezi dvěma třídami

1 a 2 definovaná vztahem

$$D_{1,2}^2 = (n - g) \sum_{i=1}^m \sum_{j=1}^m w_{ij} (\bar{x}_{i1} - \bar{x}_{i2})(\bar{x}_{j1} - \bar{x}_{j2})$$

kde m udává počet znaků v modelu,

\bar{x}_{ii} je průměr i -tého znaku ve třídě 1,

w_{ij} je prvek inverzní kovarianční matice B^{-1} .

Testování: kritérium všech párů tříd vyčísleno je jako první. Znak, který měl pro dvě od začátku největší třídy, nejmenší hodnotu $D_{1,2}^2$, je zařazen do modelu.

F-testační kritérium mezi třídami testuje nulovou hypotézu H_0 : dva vektory středních hodnot tříd objektů jsou stejné.

Kritérium je postaveno na Mahalanobisově vzdálenosti $D_{1,2}^2$.

Testační kritérium F je definováno vztahem

$$F = \frac{(n - 1 - p) n_1 n_2}{p(n - 2) (n_1 + n_2)} D_{1,2}^2$$

a může být použito k výběru znaků.

Testování: V každém kroku je zařazen do modelu ten znak, který vykazuje největší hodnotu kritéria F .

Po zavedení dalších znaků do modelu jsou sledovány změny hodnoty testačního kritéria $F_{změny}$.

Proces přidávání znaků do modelu buď pokračuje, nebo je zastaven terminačním kritériem.

Vlastnosti testu:

Test je citlivý na velikost výběru, protože velké výběry jsou náchylné snadněji vykazovat statistickou významnost než výběry malé.

I když se totiž velikosti výběru zvyšují, menší klasifikační poměr bude stále významný.

Například, pro $n = 50$, $n_s = 42$ a $k = 2$ bude $q = (50 - 42 \times 2)^2 / (50(2 - 1)) = 23.12$ a kritická hodnota pro $\alpha = 0.01$ je $\chi^2_{1-\alpha}(1) = 6.63$.

Závěr: Predikce jsou statisticky významnější než apriorní pravděpodobnost π_i , která uvádí správnou klasifikaci pro 50 %.

Například zvětšíme velikost výběru na 100 a klasifikační poměr zůstane 84 %, Pressovo q se zvýší na 46.24.

Statistické míry klasifikace spolehlivosti diskriminace:

1) **Pressovo q -kritérium** vyjadřuje míru porovnání počtu správných klasifikací vůči celkové velikosti výběru a počtu tříd. Vyčíslí se vztahem

$$q = \frac{(n - n_s k)^2}{n(k - 1)}$$

kde n je velikost výběru,
 n_s udává počet objektů správně klasifikovaných a
 k značí počet tříd,

Testování: vyšetřuje diskriminační sílu klasifikační matice v porovnání s modelem pravděpodobnosti: vypočtené q se porovnává s $\chi^2_{1-\alpha}(1)$ při dané α , a to když q překročí $\chi^2_{1-\alpha}(1)$, klasifikační matice se jeví statisticky lepší než klasifikační pravděpodobnost.

5. krok: Interpretace výsledků

Tři metody určují relativní důležitost znaku:

1) **Standardizované diskriminační koeficienty:** interpretace diskriminačních funkcí vyšetřuje znaménko a velikost *standardizovaných diskriminačních koeficientů* $\mathbf{a}^T = [a_0, a_1, \dots, a_p]$, které představují relativní příspěvek svého znaku do Fisherovy lineární diskriminační funkce:

Diskuze:

1) Znaky s *relativně velkými koeficienty* přispívají více do diskriminační síly diskriminační funkce než znaky s menšími koeficienty.

2) Znaménko ukazuje, že znak dává buď kladný, nebo záporný příspěvek.

3) Malý koeficient indikuje buď, že odpovídající znak je nevýznamný k určování vztahu nebo je neúplným vztahem, protože je zde vysoký stupeň multikolinearity.

4) Problémem je také značná nestabilita diskriminačních koeficientů.

Proměnná	Standardiz. koeficienty (73Iris) pro kanonické proměnné	
	Kořen1	Kořen2
SEPALLENGTH	0,42695	0,012408
SEPALWIDTH	0,52124	0,735261
PETALLENGTH	-0,94726	-0,401038
PETALWIDTH	-0,57516	0,581040
Vlastní	32,19193	0,285391
KumPodíl	0,99121	1,000000

6. krok: Ověření výsledků

Konečné stadium DA se týká potvrzení diskriminačních výsledků.

DA má tendenci přeceňovat hit poměr.

Metoda dělení do skupin (cross-validation) je velmi užitečná.

Rozdělení výběru:

- 1) Soubor je náhodně rozdělen na *analyzovaný výběr* a na *klasifikovaný výběr*.
- 2) Místo dělení na *analyzovaný* a *klasifikovaný* můžeme zcela náhodně rozdělit soubor několikrát.

Testujeme potvrzení diskriminační funkce pomocí *klasifikační matice* a *hit poměru*.

Když najdeme znaky, které mají největší vliv na diskriminaci mezi třídami, je dalším krokem profilování charakteristiky tříd, založené na třídnicích průměrech.

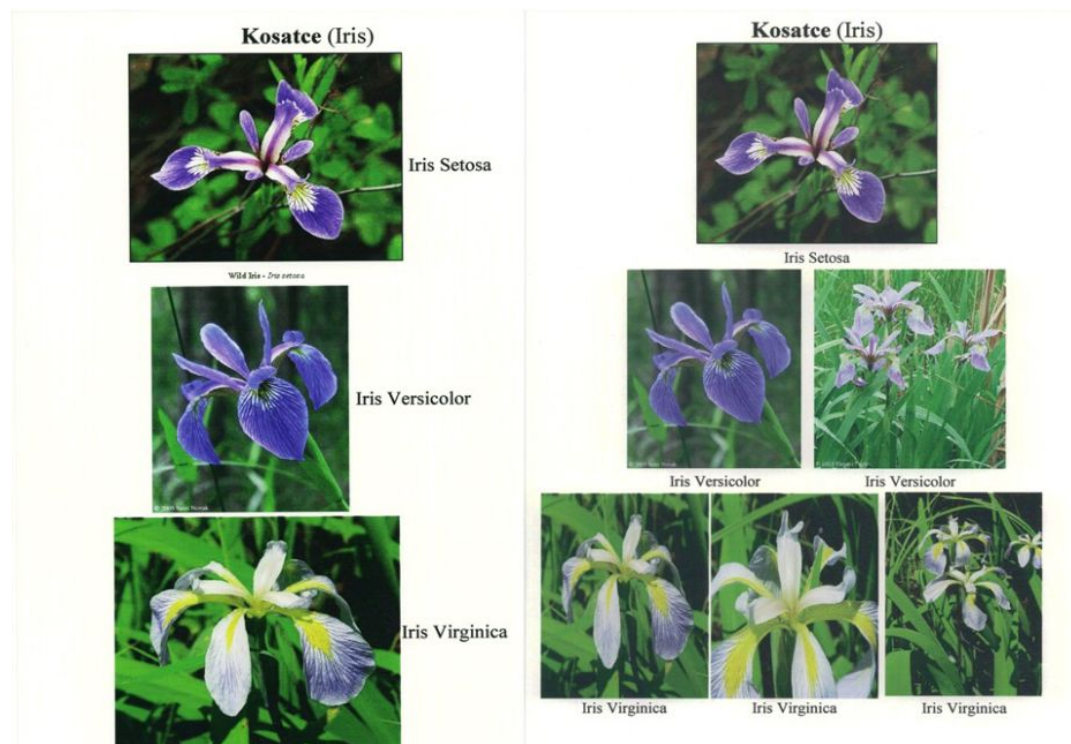
Příklad 4.24 Botanické třídění kosatců diskriminační analýzou

Ve Fisherově úloze o rozměrech okvětních lístků u 150 kosatců se analyzují květy tří základních tříd: (1) *Iris setosa*, (2) *Iris versicolor*, (3) *Iris virginica*. Květy kosatců jsou popsány čtyřmi znaky: délkou kališních lístků v mm anglicky *lsepal* a jejich šířkou *wsepal*, dále délkou korunních plátků v mm *lpetal* a jejich šířkou *wpetal*.

Každý objekt je popsán $p = 4$ znaky, a to *SepalLength*, *SepalWidth*, *PetalLength*, *PetalWidth*.

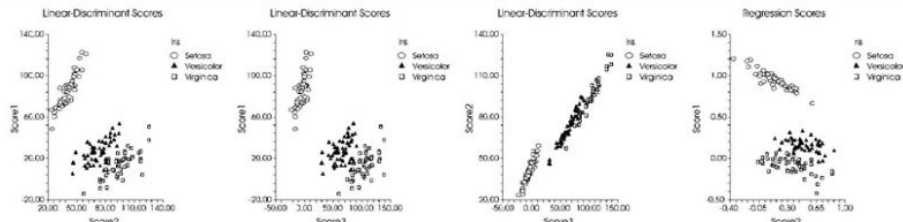
Data: Data jsou v pořadí proměnných *E418a*, *E418b*, *E418c*, *E418d*:

5.1	3.5	1.4	0.2	5.4	3.9	1.3	0.4	5.2	4.1	1.5	0.1	5.3	3.7	1.5	0.2
4.9	3.0	1.4	0.2	5.1	3.5	1.4	0.3	5.5	4.2	1.4	0.2	5.0	3.3	1.4	0.2
4.7	3.2	1.3	0.2	5.7	3.8	1.7	0.3	4.9	3.1	1.5	0.2	7.0	3.2	4.7	1.4
4.6	3.1	1.5	0.2	5.1	3.8	1.5	0.3	5.0	3.2	1.2	0.2	6.4	3.2	4.5	1.5
5.0	3.6	1.4	0.2	5.4	3.4	1.7	0.2	5.5	3.5	1.3	0.2	6.9	3.1	4.9	1.5
5.4	3.9	1.7	0.4	5.1	3.7	1.5	0.4	4.9	3.6	1.4	0.1	5.5	2.3	4.0	1.3
4.6	3.4	1.4	0.3	4.6	3.6	1.0	0.2	4.4	3.0	1.3	0.2	6.5	2.8	4.6	1.5
5.0	3.4	1.5	0.2	5.1	3.3	1.7	0.5	5.1	3.4	1.5	0.2	5.7	2.8	4.5	1.3
4.4	2.9	1.4	0.2	4.8	3.4	1.9	0.2	5.0	3.5	1.3	0.3	6.3	3.3	4.7	1.6
4.9	3.1	1.5	0.1	5.0	3.0	1.6	0.2	4.5	2.3	1.3	0.3	4.9	2.4	3.3	1.0
5.4	3.7	1.5	0.2	5.0	3.4	1.6	0.4	4.4	3.2	1.3	0.2	6.6	2.9	4.6	1.3
4.8	3.4	1.6	0.2	5.2	3.5	1.5	0.2	5.0	3.5	1.6	0.6	5.2	2.7	3.9	1.4
4.8	3.0	1.4	0.1	5.2	3.4	1.4	0.2	5.1	3.8	1.9	0.4	5.0	2.0	3.5	1.0
4.3	3.0	1.1	0.1	4.7	3.2	1.6	0.2	4.8	3.0	1.4	0.3	5.9	3.0	4.2	1.5
5.8	4.0	1.2	0.2	4.8	3.1	1.6	0.2	5.1	3.8	1.6	0.2	6.0	2.2	4.0	1.0
5.7	4.4	1.5	0.4	5.4	3.4	1.5	0.4	4.6	3.2	1.4	0.2	6.1	2.9	4.7	1.4



6. Klasifikační grafy.

Užijeme grafy: (a) lineárních diskriminačních skóre, (b) regresních skóre nebo (c) kanonických skóre.

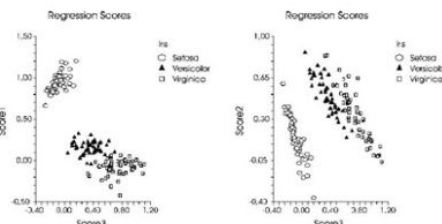


Obr. 7.9a Graf lineárního diskriminačního skóre 1 a 2 u 4 znaků 150 kosateč zdrojové matice dat Iris (STATISTICA).

Obr. 7.9b Graf lineárního diskriminačního skóre 1 a 3 u 4 znaků 150 kosateč zdrojové matice dat Iris (STATISTICA).

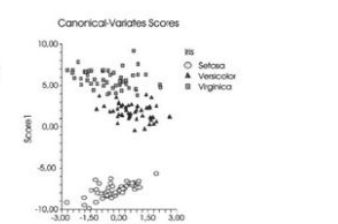
Obr. 7.9c Graf lineárního diskriminačního skóre 2 a 3 u 4 znaků 150 kosateč zdrojové matice dat Iris (STATISTICA).

Obr. 7.9d Graf regresního skóre 1 a 2 u 4 znaků 150 kosateč zdrojové matice dat Iris (STATISTICA).



Obr. 7.9e Graf regresního skóre 1 a 3 u 4 znaků 150 kosateč zdrojové matice dat Iris (STATISTICA).

Obr. 7.9f Graf regresního skóre 2 a 3 u 4 znaků 150 kosateč zdrojové matice dat Iris (STATISTICA).



Obr. 7.9g Graf kanonických proměnných 1 a 2 u 4 znaků 150 kosateč zdrojové matice dat Iris (STATISTICA).

Cvičení v programu STATISTICA

- Načtení dat
- Zadání metody DA.
- Zadání proměnných: grupovací a nezávislé.
- Zadání kódů pro grupovací proměnné.
- Zadání krokové analýzy
- Definice modelu: záložky Základ, Detaily, Popisné statistiky
- Záložka Detaily: okénko Metoda – zvol standardní (ze tří nabízených). Tolerance: 0.01
- Záložka Popisné statistiky: volba Zobrazit popisné statistiky
- Ze tří záložek zvol Detaily: 7 volieb postupně otevřít a prohlédnout.
- STORNO a zpět do okna Popisné statistiky a zobrazit Popisné statistiky zvol záložku Všechny.
- Zadej Celková kovariance a korelace. Korelační matice 4 znaků: Celková korelace.
- Zvol graf korelační matice kliknutím na Graf oskové korelace.
- Zvol Křabivo/ graf průměru všech 4 znaků.
- Upravit křabivový graf.
- 1 krát STORNO až na Definice modelu. Metoda – kroková dopředná. Záložka Detaily a pak Výsledky: pouze shrnutí.
- Zvol Výpočet v záložce Základní nastavení. Seřadí znaky dle statistické významnosti.
- Zpět na záložku Detaily a v ní Kanonická analýza a záložka Základní nastavení.
- Zvol Výpočet: Test Chi kvadrát postupných koef. Hledání ortogonálních diskriminačních funkcí čili kanonických koef. Chi²-testem.
- Komentovat Chi²-test.
- Kanonická analýza a v ní záložka Detaily a v ní záložka pro kanonické proměnné.
- Kanonická analýza a v ní záložka Detaily a v ní Faktorová analýza a pojmenování koef. Koef. 1 je o délkách květů a Koef. 2 je o šířkách květů.
- Zpět a zvol Průměry kanonických proměnných a souřadnice jejich polohy.
- Zpět a záložka Kanonické skóre.
- Záložka Bodový graf kanonických skóre. Upravit graf.
- Zpět a Uložte kanonické skóre za IRISTYPE a udeleť Graf, Bodové grafy. Překopírovat Popisy boů.
- STORNO a pak záložka Klasifikace: zvol a prohlédni si záložky Klasifikační matice, Klasifikace objektů, Mahalanobisovy vzdálenosti, Aposteriorní pravděpodobnosti.

	1	2	3	4	5
	SEPALLENGTH	SEPALWIDTH	PETALLENGTH	PETALWIDTH	Trída
1	5	3,3	1,4	0,2	SETOSA
2	6,4	2,8	5,6	2,2	VIRGINIC
3	6,5	2,8	4,6	1,5	VERSICOL
4	6,7	3,1	5,6	2,4	VIRGINIC
5	6,3	2,8	5,1	1,5	VIRGINIC
6	4,8	3,4	1,4	0,3	SETOSA
7	6,9	3,1	5,1	2,3	VIRGINIC
8	6,2	2,2	4,5	1,5	VERSICOL
9	5,9	3,2	4,8	1,8	VERSICOL
10	4,6	3,6	1	0,2	SETOSA
11	6,1	3	4,6	1,4	VERSICOL
12	6	2,7	5,1	1,6	VERSICOL
13	6,5	3	5,2	2	VIRGINIC
14	5,6	2,5	3,9	1,1	VERSICOL
15	6,5	3	5,5	1,8	VIRGINIC
16	5,8	2,7	5,1	1,9	VIRGINIC
17	6,8	3,2	5,9	2,3	VIRGINIC
18	5,1	3,3	1,7	0,5	SETOSA
19	5,7	2,8	4,5	1,3	VERSICOL
20	6,2	3,4	5,4	2,3	VIRGINIC
21	7,7	3,8	6,7	2,2	VIRGINIC
22	6,3	3,3	4,7	1,6	VERSICOL
23	6,7	3,3	5,7	2,5	VIRGINIC
24	7,6	3	6,6	2,1	VIRGINIC
25	4,9	2,5	4,5	1,7	VIRGINIC
26	5,5	3,5	1,3	0,2	SETOSA
27	6,7	3	5,2	2,3	VIRGINIC
28	7	3,2	4,7	1,4	VERSICOL
29	6,4	3,2	4,5	1,5	VERSICOL
30	6,1	2,8	4	1,3	VERSICOL
31	4,8	3,1	1,6	0,2	SETOSA
32	5,9	3	5,1	1,8	VIRGINIC
33	6,4	2,4	3,8	1,1	VERSICOL
34	5,5	3	6	1,8	VIRGINIC
35	6,1	2,8	4	1,3	VERSICOL
36	6,2	3,4	5,4	2,3	VIRGINIC
37	4,9	2,5	4,5	1,7	VIRGINIC
38	5,4	3	4,5	1,5	VERSICOL
39	7,9	3,8	6,4	2	VIRGINIC
40	4,4	3,2	1,3	0,2	SETOSA
41	6,7	3,3	5,7	2,1	VIRGINIC
42	5	3,5	1,6	0,6	SETOSA
43	5,8	2,6	4	1,2	VERSICOL
44	4,4	3,2	1,3	0,2	SETOSA

1. Fisherova úloha kosatece (iris, 1936) třídí kosatece do tří tříd podle svých 4 znaků, a to děly a šířky kališních lístků a délky a šířky korunních plátků. Poslední sloupec dat je kódovaný znak.

Čtenář u své úlohy zamění okvětní lístky za znaky své úlohy. Poslední sloupec bude vždy kódovací znak.

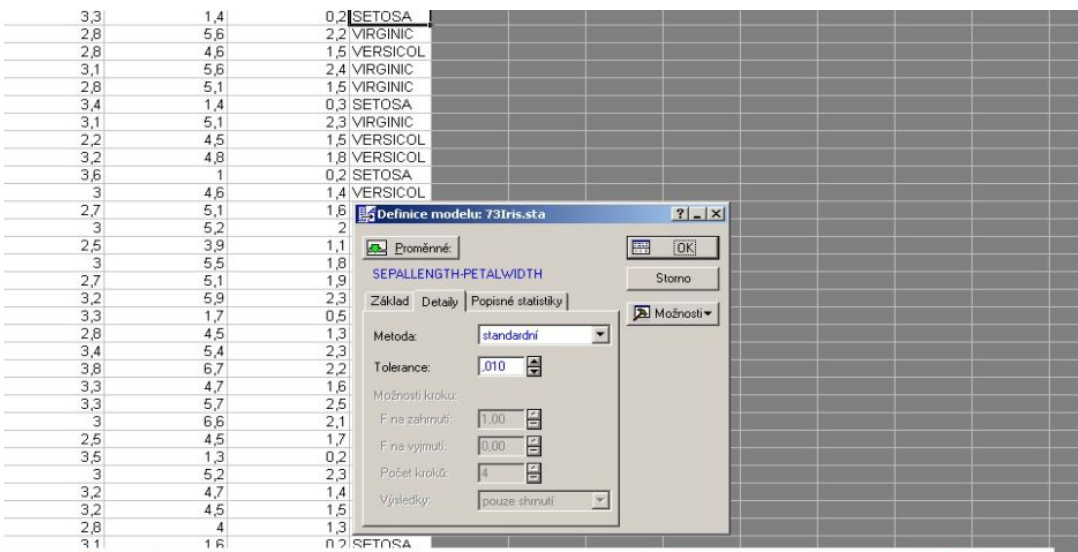
2. Zvolí se *Diskriminační analýza* v bloku *Víceproměnné průzkumné metody*.

3. Zadájí se všechny nezávislé Proměnné a jedna Grupovací proměnná.

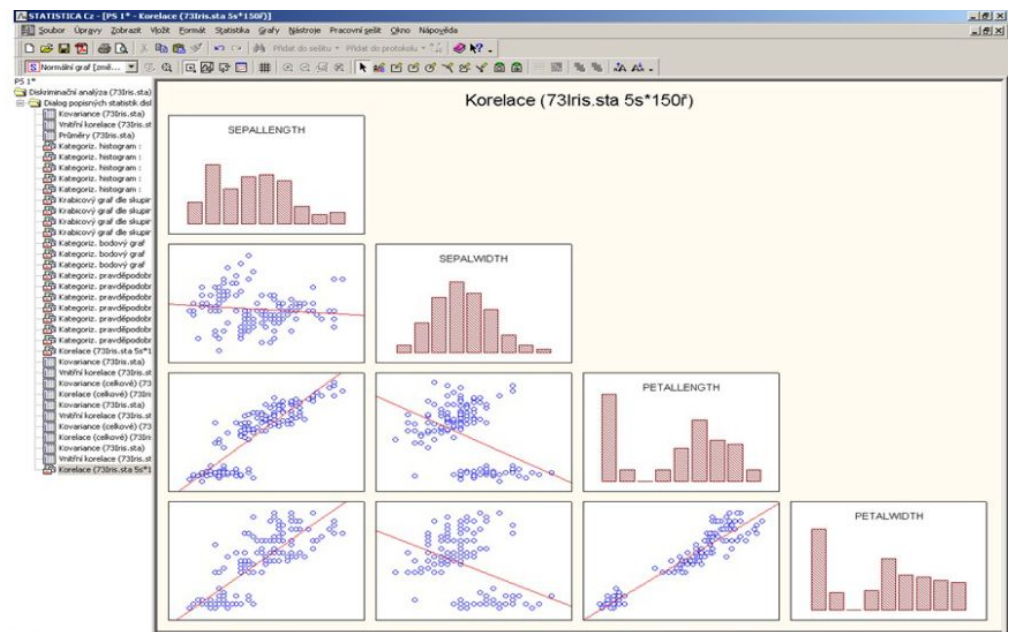
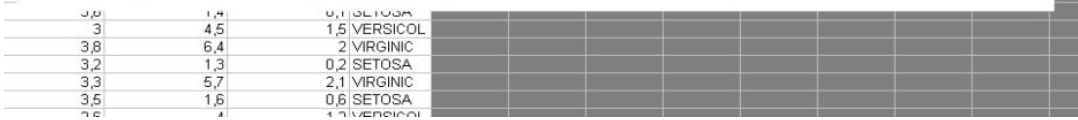
4. Grupovací proměnná má krátký název (zde Trída) a Dlouhý název (zde Three types of iris), které třídí 150 květů kosatců do tříd Setosa, Versicolor a Virginica.

5. V okénku Další možnosti se zaškrtně Kroková analýza. Díry ve zdrojové matici vedou buď k vyškrtnutí řádku, objektu (značeno Celé případy) nebo k dodatečnému vyplnění díry průměrem (značeno Substituce průměrem).

6. Definice modelu má tři záložky to Základ, Detaily, Popisné statistiky, které si nyní postupně prohlédneme.



7. Záložka *Detaily*: zvolí se zde *Metoda* diskriminace a zadají se hodnoty tolerance, pokud je jiná než defaultně zadaná.



44. Červené přímky v oknech diagramu *Celkové korelace* ukazují na lineární závislost mezi dvěma znaky. Osamocené shluky bodů v každém grafu ukazují, že lze diskriminovat kosatce do několika tříd.

Výsledky diskriminační funkční analýzy (73Iris.sta)
 krok 4, poč. prom. v modelu 4, grupovací třída (3 skup.)
 Wilks. lambda: ,02344 přibliž. F (2,289)=199,15 p<0,0000

	Wilks.	Parc.	F na vyl.	Uroveň p	Toler.	1-toler.	R ²
N=150	Lambda	Lambda	(2,144)				
PETALLENGTH	0,035025	0,669206	35,59018	0,000000	0,365126	0,634874	
SEPALWIDTH	0,030580	0,766480	21,93593	0,000000	0,608859	0,391141	
PETALWIDTH	0,031546	0,743001	24,90433	0,000000	0,649314	0,350686	
SEPALLENGTH	0,024976	0,938464	4,72115	0,010329	0,347993	0,652007	

49. Program seřadí znaky a vybere přitom nejprve takový znak, který nejlépe a nejvíce přispívá k diskriminaci objektů do tříd. Jde o znak s největší hodnotou F , větší hodnotou F než bylo F zadáno na začátku, při hledání nejlepších znaků se za počáteční hodnotu totiž volí obvykle nula. Když bychom ale naopak chtěli všechny znaky z modelu odebrat, zadáme na počátku pro F vysokou hodnotu, například 9999.

Tolerance $1-R^2$ je mírou redundance (= nadbytečnosti) znaku. Vykazuje-li totiž znak toleranci třeba 0.01, pak je z 99% nadbytečný vůči dříve zadaným ostatním znakům. Když jeden nebo více znaků jsou příliš redundantní, diskriminaci nelze provést. Položíme proto mezní hodnotu tolerance rovnu 0.01.

Wilksovo λ je standardní statistika testu významnosti diskriminační síly sledovaného znaku v modelu. Nabývá hodnot od 1, značící pak žádnou diskriminační sílu znaku, do 0.0, značící perfektní diskriminační sílu znaku.

Výsledky diskriminační funkční analýzy (73Iris.sta)
 krok 4, poč. prom. v modelu 4, grupovací třída (3 skup.)
 Wilks. lambda: ,02344 přibliž. F (2,289)=199,15 p<0,0000

	Wilks.	Parc.	F na vyl.	Uroveň p	Toler.	1-toler.	R ²
N=150	Lambda	Lambda	(2,144)				
PETALLENGTH	0,035025	0,669206	35,59018	0,000000	0,365126	0,634874	
SEPALWIDTH	0,030580	0,766480	21,93593	0,000000	0,608859	0,391141	
PETALWIDTH	0,031546	0,743001	24,90433	0,000000	0,649314	0,350686	
SEPALLENGTH	0,024976	0,938464	4,72115	0,010329	0,347993	0,652007	

Parciální Wilksovo λ je jednotkový příspěvek dotyčného znaku k diskriminaci objektů do tříd. Čím je jeho hodnota bližší k 0, tím větší diskriminační síla tohoto znaku. Výpočet začíná s počáteční hodnotou rovnou 0.0. Wilksovo λ konvertuje do F kritéria a odpovídající vypočtené hladiny významnosti p .

Ukazuje se zde, že *PETALLENGTH* přispívá nejvíce, *PETALWIDTH* jako druhý, pak *SEPALWIDTH* a nejméně přispívá k diskriminaci kosatců do tříd *SEPALLENGTH*. Čím menší je totiž λ , tím více přispívá znak k diskriminaci. Z toho plyne, že *PETAL* je hlavním znakem, který vlastně vůbec umožňuje rozředit kosatce do tříd.

STATISTICA Cz - [PS 1* - Test chí-kvadrát po odstranění post. kořenů (73Iris.sta)]

Soubor Úpravy Zobrazit Vložit Formát Statistika Grafy Nástroje Data Pracovní sešit Okno Nápověda

PS 1*

Diskriminační analýza (73Iris.sta)

Dialog popisných statistik disl

Kořeny odstraněny		Test chí-kvadrát po odstranění post. kořenů (73Iris.sta)				
	Vlastní číslo	Kan. R	Wilks. Lambda	Chi-kv.	sv	Úroveň p
0	32,19193	0,984821	0,023439	546,1153	8	0,000000
1	0,28539	0,471197	0,777973	36,5297	3	0,000000

52. Nejprve určíme, zda obě diskriminační funkce (či kanonické kořeny) jsou statisticky významné. První řádek tabulky *Testu Chi-kvadrát* ukazuje statistickou významnost všech (zde 2) kanonických kořenů. Druhý řádek se týká všech kořenů, ale bez prvního, atd. Tabulka proto odpoví, kolik a které kořeny jsou statisticky ještě významné. Zde jsou obě diskriminační funkce statisticky významné. Když proto budeme u nových kosatců dvě míry okvětních lístků, PETAL a SEPAL, zaručíme, že je pak možné rozřadit kosatce do tříd.

STATISTICA Cz - [PS 1* - Standardiz. koeficienty (73Iris.sta)]

Soubor Úpravy Zobrazit Vložit Formát Statistika Grafy Nástroje Data Pracovní sešit Okno Nápověda

PS 1*

Diskriminační analýza (73Iris.sta)

Dialog popisných statistik disl

Standardiz. koeficienty (73Iris.sta) pro kanonické proměnné		
Proměnná	Kořen1	Kořen2
PETALLENGTH	-0,94726	-0,401038
SEPALWIDTH	0,52124	0,735261
PETALWIDTH	-0,57516	0,581040
SEPALLENGTH	0,42695	0,012408
Vlastní číslo	32,19193	0,285391
Kumulovaný podíl	0,99121	1,000000

53. Tabulka standardizovaných koeficientů: Původní spojení s původními daty poskytují skóre diskriminační funkce. Standardizované koeficienty patří ke standardizovaným proměnným v normovaných vzájemně porovnatelných škálách. První funkce je silně vážená délkami a šířkami PETALS. Zbývající kořeny ale také trochu přispívají. Druhá diskriminační funkce je vážená více znakem SEPALWIDTH a daleko však méně PETALLENGTH.

STATISTICA Cz - [PS 1* - Faktorová strukturní matice (73Iris.sta)]

Soubor Úpravy Zobrazit Vložit Formát Statistika Grafy Nástroje Data Pracovní sešit Okno Nápověda

PS 1*

Diskriminační analýza (73Iris.sta)

Dialog popisných statistik disl

Faktorová strukturní matice (73Iris.sta) Korelační proměnné - Kanonické kořeny (vnitřní korelace)		
Proměnná	Kořen1	Kořen2
PETALLENGTH	-0,706065	0,167701
SEPALWIDTH	0,119012	0,863681
PETALWIDTH	-0,633178	0,737242
SEPALLENGTH	-0,222596	0,310812

54. Faktorová strukturní matice: v tabulce uvedené představují korelace mezi znaky a diskriminačními funkcemi. Proto se k objasnění "významu" diskriminačních funkcí výhodně pojmenovat či názvem přidělit význam těmto kořenům podobně jako se to dělá v FA.

STATISTICA Cz - [PS 1* - Průměry kan. proměnných (73Iris.sta)]

Soubor Úpravy Zobrazit Vložit Formát Statistika Grafy Nástroje Data Pracovní sešit Okno Nápověda

PS 1*

Diskriminační analýza (73Iris.sta)

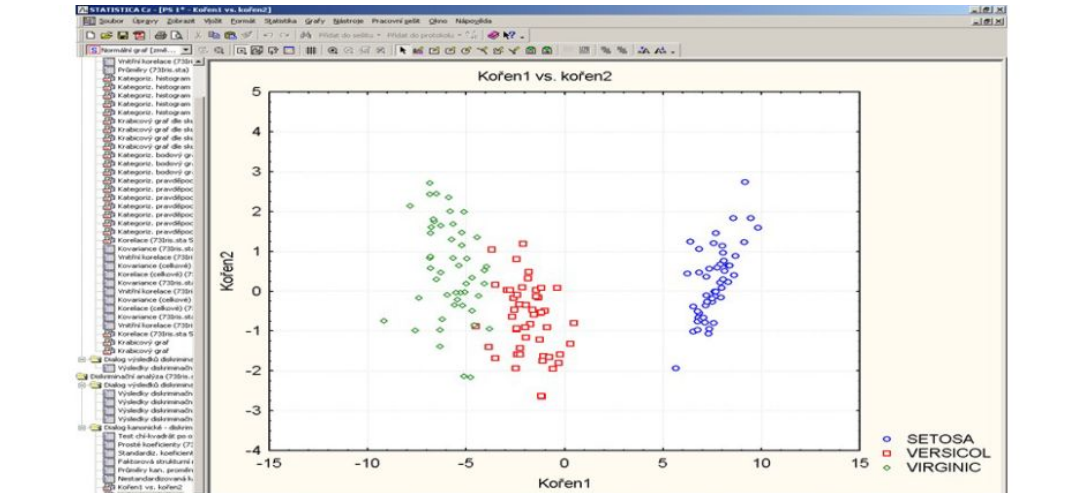
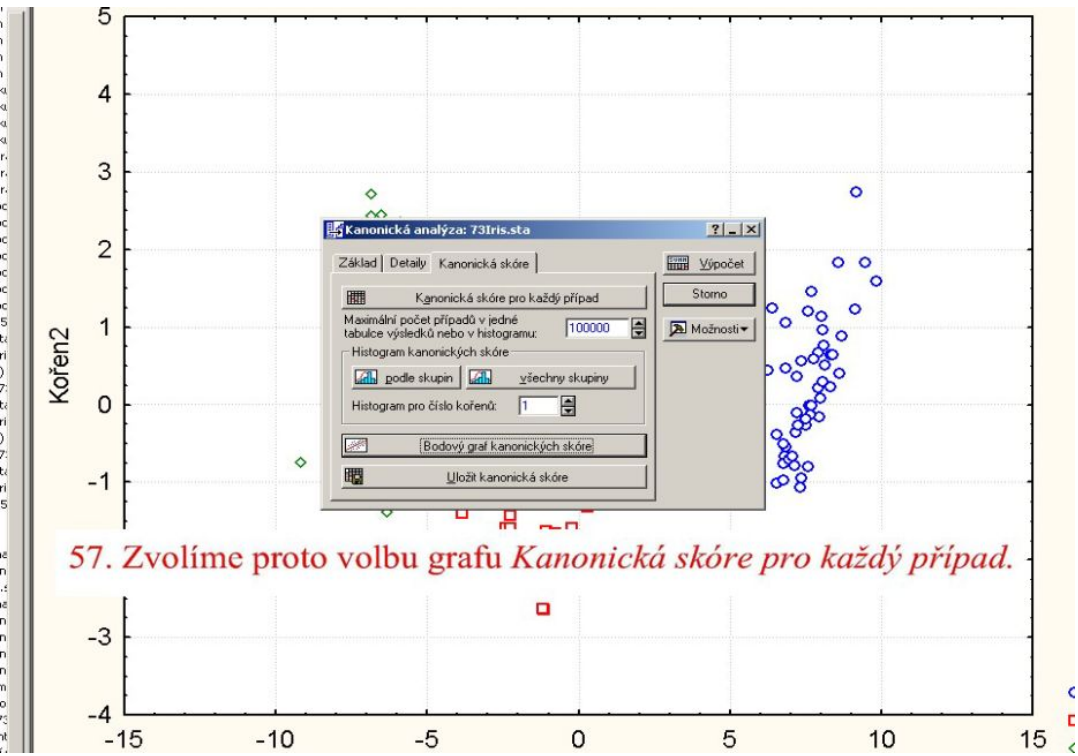
Dialog popisných statistik disl

Průměry kan. proměnných (73Iris.sta)		
Skup.	Kořen1	Kořen2
SETOSA	7,60760	0,215133
VERSICOL	-1,82505	-0,727900
VIRGINIC	-5,78255	0,512767

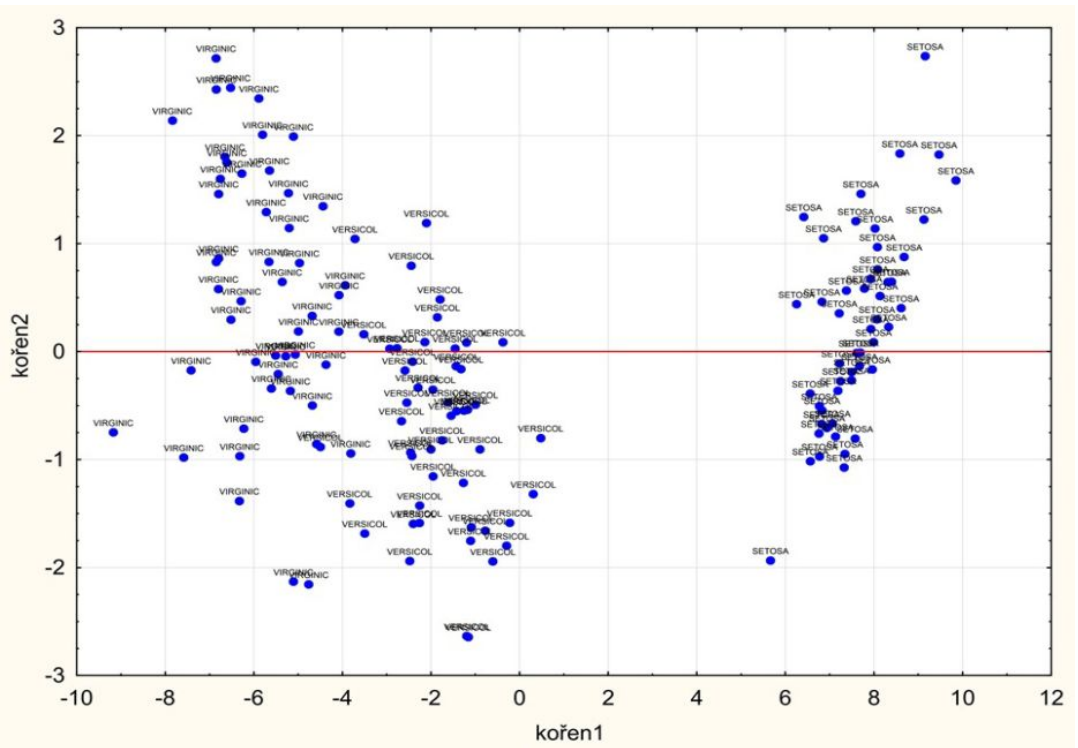
55. Průměry kanonických proměnných: je třeba se podívat na průměry kanonických proměnných v každém kanonickém kořeni. Zde první funkce diskriminuje většinou mezi třídami Setosa a ostatními kosatci. Kanonický průměr Setosy se zcela liší od průměru ostatních kosatců. Druhá diskriminační funkce rozlišuje pak mezi Setosou a ostatními kosatci, avšak tato diskriminace je mnohem méně výrazná než u první funkce.

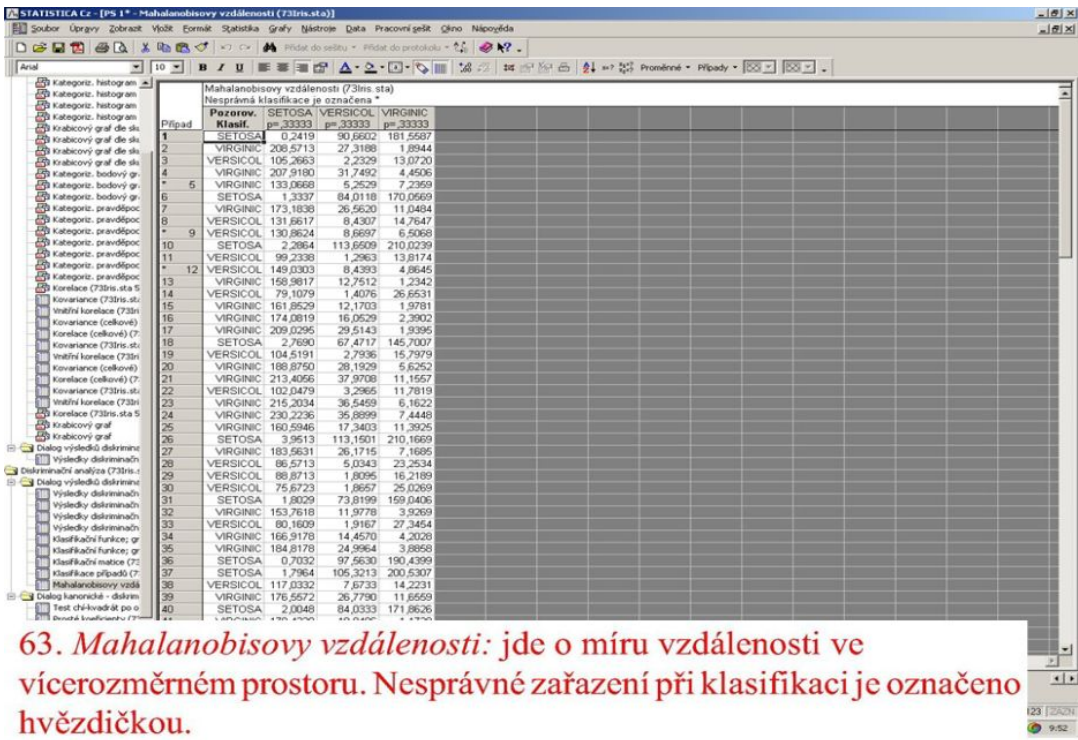
Nestandardizované kanonické skóre (73iris.sta)				Nestandardizované kanonické skóre (73iris.sta)				Nestandardizované kanonické skóre (73iris.sta)			
Případ	Skup.	Kořen1	Kořen2	Případ	Skup.	Kořen1	Kořen2	Případ	Skup.	Kořen1	Kořen2
1	SETOSA	7.67197	-0.13489	53	SETOSA	6.55993	-0.38922	106	SETOSA	7.58233	-0.00922
2	VIRGINIC	-6.80215	0.58950	54	SETOSA	6.77136	-0.07063	107	SETOSA	7.91794	0.67512
3	SETOSA	7.21252	0.35584	55	SETOSA	6.84994	-1.05954	108	VIRGINIC	-5.30071	0.64672
4	VIRGINIC	-6.65309	1.80532	56	VIRGINIC	-5.10748	-2.15099	109	VERSICOL	-2.14258	0.68619
5	VIRGINIC	-3.81516	-0.94299	57	VIRGINIC	-5.51910	0.29688	110	VIRGINIC	-6.32774	-0.84730
6	VERSICOL	-2.54868	0.47220	58	SETOSA	7.57247	-0.82546	111	VERSICOL	-2.14258	0.68619
7	VIRGINIC	-5.10559	1.99218	59	SETOSA	9.84994	-1.05954	112	VERSICOL	-2.14258	0.68619
8	VERSICOL	-3.49695	-1.68486	60	SETOSA	6.85444	1.05165	113	SETOSA	7.34000	-0.84730
9	VERSICOL	-3.71950	-1.04651	61	VERSICOL	-1.20717	0.08444	114	VIRGINIC	-4.08387	0.18954
10	SETOSA	6.68104	0.67160	62	SETOSA	6.84994	-1.05954	115	SETOSA	6.39733	0.64738
11	VERSICOL	-2.29249	-0.33286	63	SETOSA	7.12969	-0.76996	116	VERSICOL	-1.32553	-0.16287
12	VERSICOL	-4.49847	-0.88275	64	SETOSA	8.09180	0.30042	117	VERSICOL	-3.83628	-1.40589
13	VERSICOL	-6.68104	0.67160	65	SETOSA	-4.08774	-0.32144	118	VERSICOL	-2.25741	-1.42619
14	VERSICOL	-1.09043	-1.62925	66	VERSICOL	-5.77895	-1.05919	119	VERSICOL	-2.25741	-1.42619
15	VIRGINIC	-5.06601	-0.02627	67	SETOSA	6.82308	0.48301	120	VERSICOL	-4.07597	-0.79991
16	SETOSA	8.68104	0.87769	68	SETOSA	7.18677	-0.39099	121	VERSICOL	-1.54848	0.50300
17	VERSICOL	-2.29249	-0.33286	69	VERSICOL	-2.59693	0.17481	122	VIRGINIC	-6.85053	0.62883
18	VERSICOL	-4.49847	-0.88275	70	VERSICOL	0.30746	-1.31887	123	VIRGINIC	-5.50123	1.14471
19	VIRGINIC	-4.96774	0.82114	71	SETOSA	9.16624	2.73780	124	VIRGINIC	-6.85053	0.62883
20	VIRGINIC	-5.88575	2.34050	72	SETOSA	6.13231	0.91440	125	VIRGINIC	-5.50123	1.14471
21	VIRGINIC	-4.96774	0.82114	73	VIRGINIC	-6.79972	0.86320	126	SETOSA	8.55243	1.93448
22	VERSICOL	-1.09043	-1.62925	74	VIRGINIC	-6.52450	2.44854	127	SETOSA	7.78075	0.58434
23	VIRGINIC	-6.84736	2.42895	75	VIRGINIC	-6.50748	-0.03581	128	SETOSA	-6.85281	2.71759
24	VIRGINIC	-5.06601	-0.02627	76	VERSICOL	-1.61616	-0.47070	129	SETOSA	7.44407	1.34724
25	VIRGINIC	-5.20748	-0.03681	77	VERSICOL	-1.61616	-0.47070	130	VERSICOL	-2.59693	0.17481
26	SETOSA	6.25140	0.43970	78	VIRGINIC	-4.58372	-0.85862	131	VERSICOL	-2.59693	0.17481
27	VIRGINIC	-6.79602	1.48069	79	SETOSA	6.37486	0.95954	132	VERSICOL	-2.59693	0.17481
28	SETOSA	6.25140	0.43970	80	SETOSA	9.12935	1.72463	133	VERSICOL	-2.59693	0.17481
29	VERSICOL	-2.42997	-0.96613	81	VIRGINIC	-6.26201	0.46718	134	VERSICOL	-0.31637	0.98954
30	VIRGINIC	-6.88882	1.75164	82	VIRGINIC	-4.99550	0.18777	135	VERSICOL	-2.41696	-0.09278
31	VERSICOL	-2.44848	0.79596	83	VIRGINIC	-3.83628	0.11423	136	VERSICOL	-2.41696	-0.09278
32	VIRGINIC	-6.84736	2.42895	84	VERSICOL	-0.22267	-1.58467	137	VERSICOL	-2.41696	-0.09278
33	VIRGINIC	-7.41817	-0.17312	85	VERSICOL	-2.02559	-0.90542	138	SETOSA	-5.68034	0.53711
34	VIRGINIC	-4.67800	-0.49910	86	VERSICOL	-1.18196	-0.93797	139	SETOSA	7.21930	-0.10960
35	VIRGINIC	-5.11259	-0.36348	87	SETOSA	9.46768	1.82523	140	SETOSA	7.21930	-0.10960
36	SETOSA	8.61367	0.40325	88	SETOSA	7.06201	-0.69340	141	VERSICOL	-1.43768	-0.13442
37	VIRGINIC	-5.64500	1.67772	89	VIRGINIC	-9.17147	-1.74600	142	VERSICOL	-1.43768	-0.13442
38	VERSICOL	-1.45928	0.02854	90	VIRGINIC	-4.76454	-2.15574	143	SETOSA	8.01638	0.96658
39	VERSICOL	-1.79771	0.48439	91	SETOSA	7.70195	-1.48172	144	SETOSA	8.02097	1.14620
40	VERSICOL	-0.99761	-0.48053	92	VERSICOL	-1.75203	-0.82113	145	VERSICOL	-2.26247	-1.58725
41	SETOSA	6.78995	-0.79900	93	VERSICOL	-1.95842	-0.35196	146	VERSICOL	-2.26247	-1.58725
42	VIRGINIC	-4.68315	-0.33203	94	VERSICOL	-2.10361	-1.19157	147	VERSICOL	-2.26247	-1.58725
43	VERSICOL	-1.10669	-1.75225	95	SETOSA	7.60529	-0.01863	148	VERSICOL	-2.26247	-1.58725
44	VIRGINIC	-5.17956	0.36348	96	SETOSA	6.50095	-1.01916	149	SETOSA	7.58682	-0.18838
45	VIRGINIC	-7.58120	-0.98072	97	VERSICOL	-1.19378	-2.63446	144	SETOSA	7.58682	-0.18838
46	VIRGINIC	-4.37150	-0.12130	98	VERSICOL	-0.80552	-1.84288	145	SETOSA	7.58682	-0.18838
47	SETOSA	6.33042	0.23813	99	VERSICOL	-0.80552	-1.84288	146	SETOSA	7.58682	-0.18838
48	VERSICOL	-2.40197	-1.59423	100	SETOSA	7.48883	-0.26538	147	VERSICOL	-3.51848	0.16028
49	VIRGINIC	-5.27916	-0.04246	101	VERSICOL	6.81320	-0.67063	148	VIRGINIC	-7.83947	2.13071
50	SETOSA	6.76470	-0.50575	102	VIRGINIC	-6.27264	1.04948	149	SETOSA	8.31445	0.64489
51	SETOSA	6.03192	0.76330	103	VERSICOL	-1.42159	-0.95174	150	VERSICOL	-0.28918	-1.78672
52	VIRGINIC	-4.07704	0.52034	104	VIRGINIC	-6.22824	-0.71272				

56. Nestandardizované kanonické skóre: Tabulka nestandardizovaných kanonických kořenů ukazuje na diskriminaci všech 150 kosatců. Samozřejmě by bylo zde názornější provést grafické zobrazení.

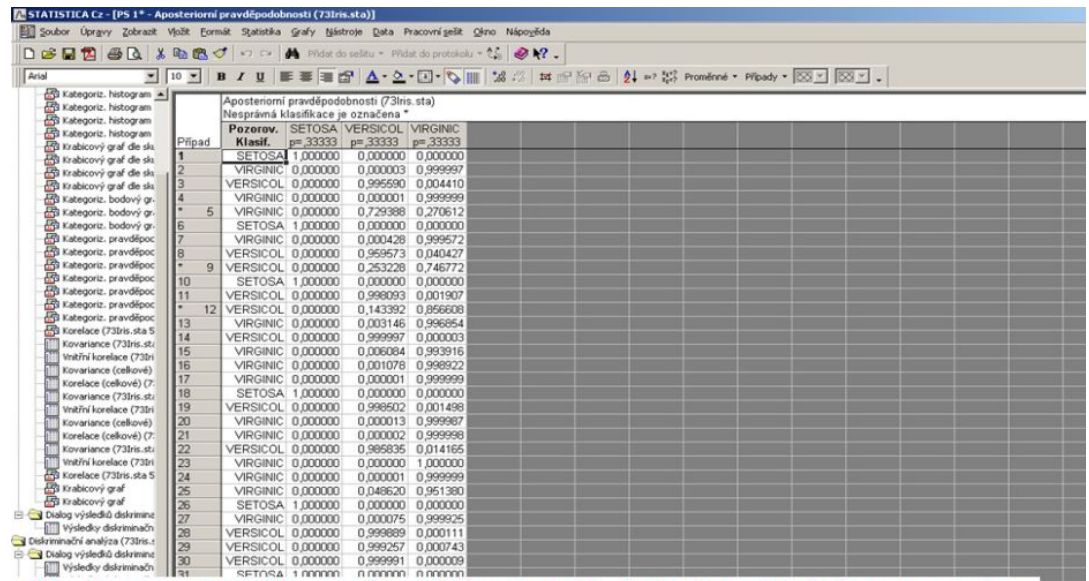


58. Graf kanonických skóre: v tomto grafu jsou kosatce Setosa zobrazeny modře daleko vpravo od ostatních. Kořen 1 (čili první diskriminační funkce) silně diskriminuje Setosu od ostatních kosatců. Kořen 2 (druhá diskriminační funkce) poskytuje diskriminaci mezi Versicolor se zápornými hodnotami Kořene 2 a ostatními kosatci, převážně kladného Kořene 2. Diskriminace však není tak zřetelná jako je tomu u prvního Kořene 1.

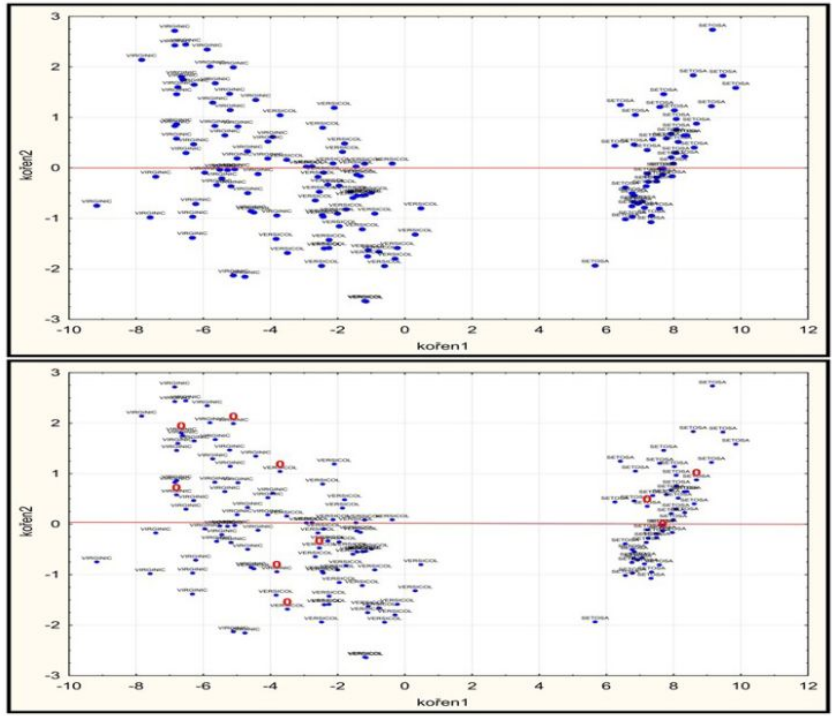




63. Mahalanobisovy vzdálenosti: jde o míru vzdálenosti ve vícerozměrném prostoru. Nesprávné zařazení při klasifikaci je označeno hvězdičkou.



64. A posteriorní pravděpodobnosti: Každá třída je dokumentována svou vypočtenou číli a posteriorní pravděpodobností. Největší hodnota pravděpodobnosti ve sloupečku ukazuje, že kosatec náleží do třídy dotyčného sloupce.

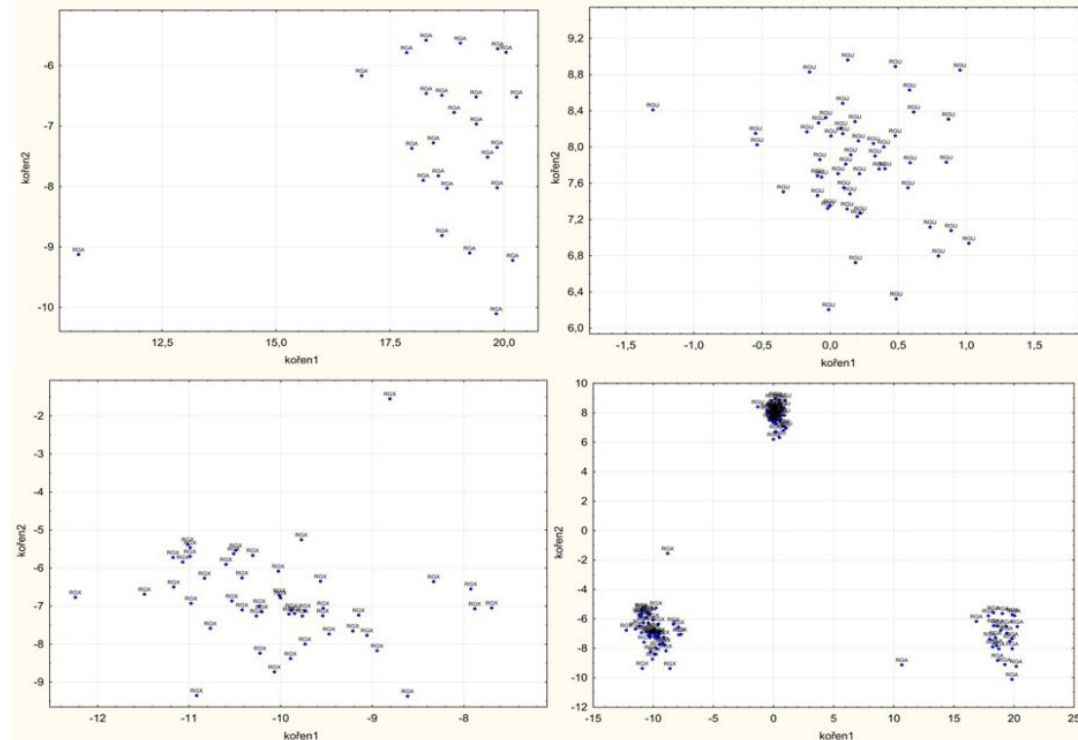
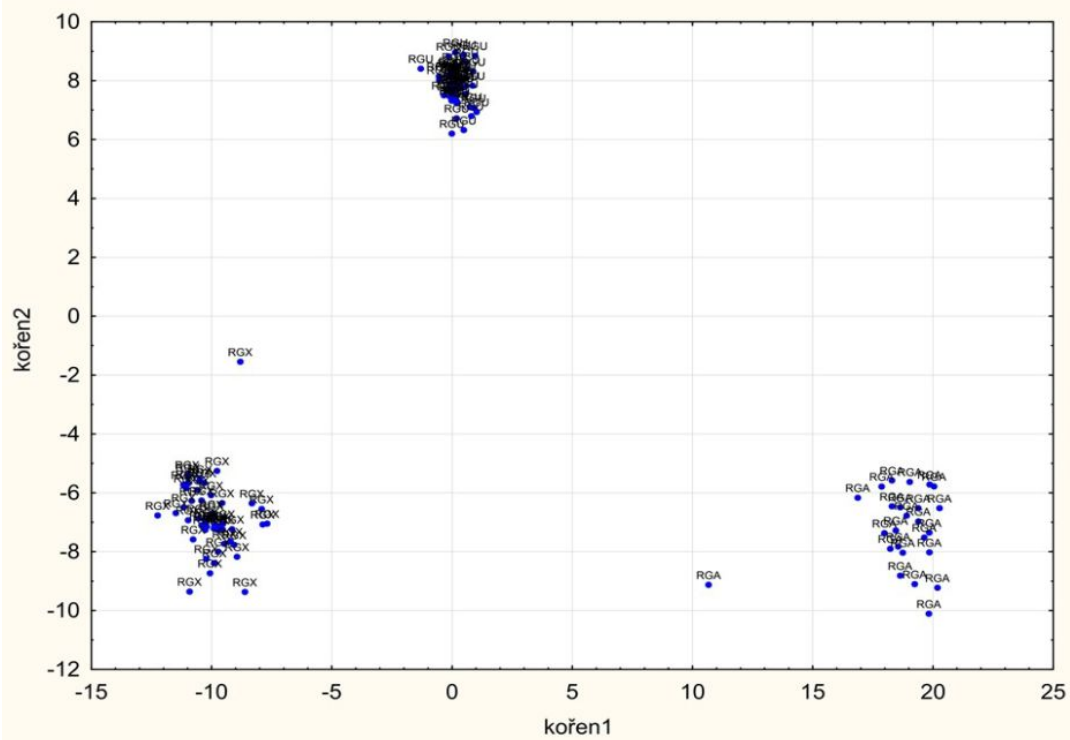


Úloha C32 Klasifikace a struktura znaků oxidu titaničitého dle chemických a fyzikálních vlastností (DA)

Výsledky analýz tří skupin bílého pigmentu oxidu titaničitého. Proved'te klasifikaci diskriminační analýzou (DA).

- **Data:** Popis znaků ve sloupcích: Data **ID**, Zařazení do tříd **Klas**, tj. typy RGA, RGX, RGU tvoří tréninková data. Diskriminátory:
- Ti** obsah Ti (%),
- TL** obsah těkavých látek (%),
- SpO** je spotřeba oleje (g/100 g pigmentu),
- Barv** značí barvivot,
- Podt** značí podtón,
- Si** značí obsah Si (%),
- Al** značí obsah Al (%).

ID	Ti	TL	SpO	Barv	Podt	Si	Al	Klas
1	93.78	0.5	24.2	1870	10	2.174	3.511	RGA
.....
144	98.21	0.2	19	1860	13	0.354	1.014	RGX



Úloha B23 Porovnání analýzy krve diabetiků a zdravých jedinců

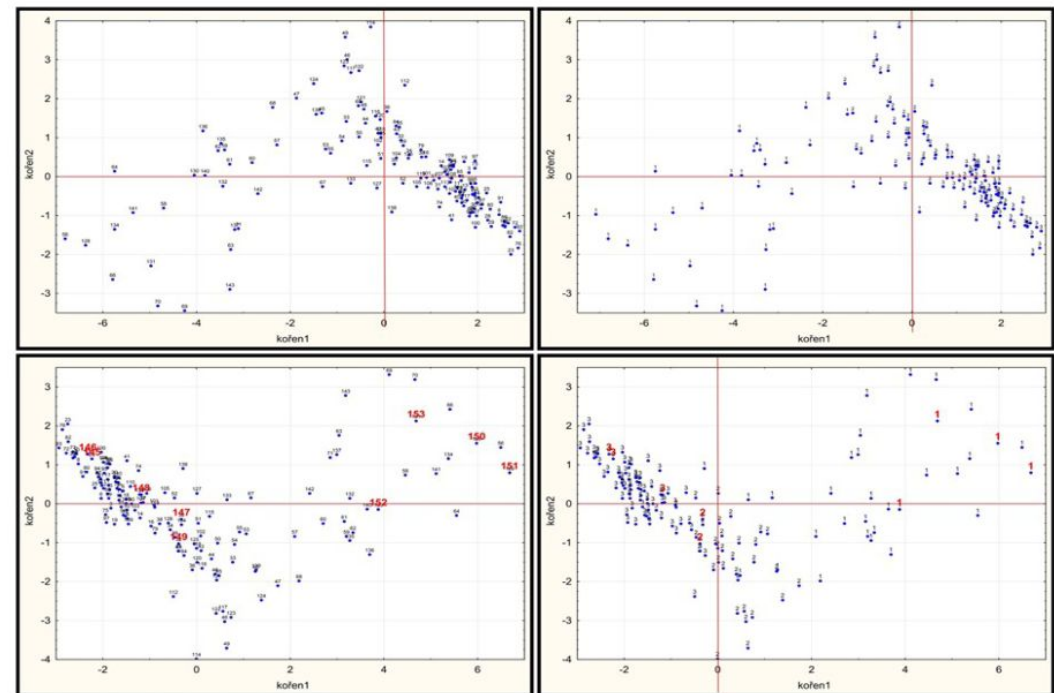
Jsou porovnány výsledky rozboru krve zdravých jedinců a dvou skupin diabetiků. Postavte soubor diskriminačních funkcí, které mohou předpovědět třídu a který rozlišuje mezi třídami podle hodnot ostatních kvantitativních proměnných.

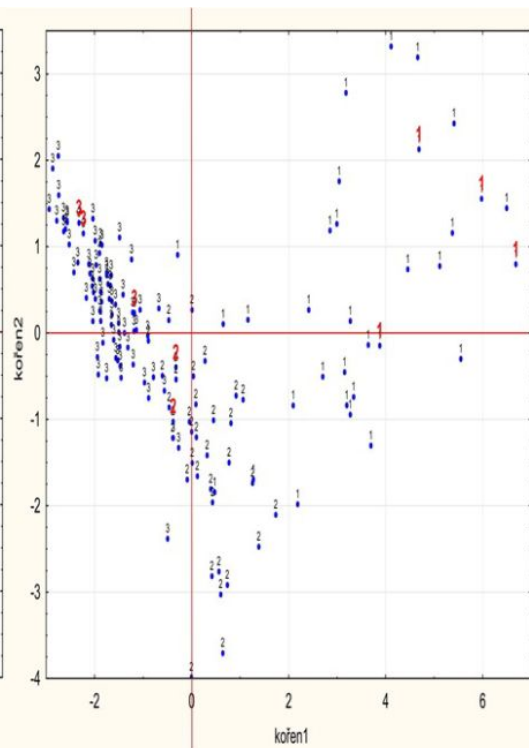
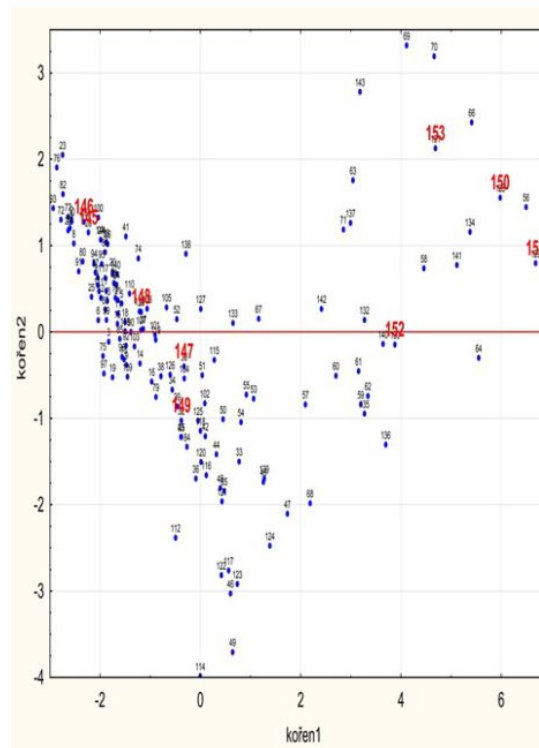
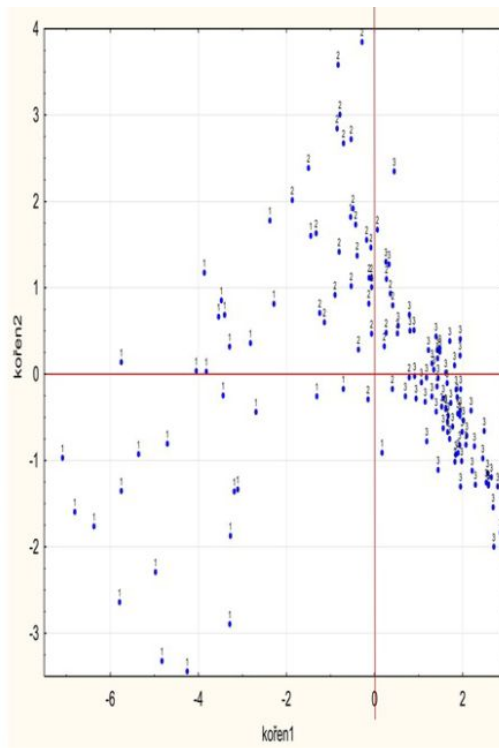
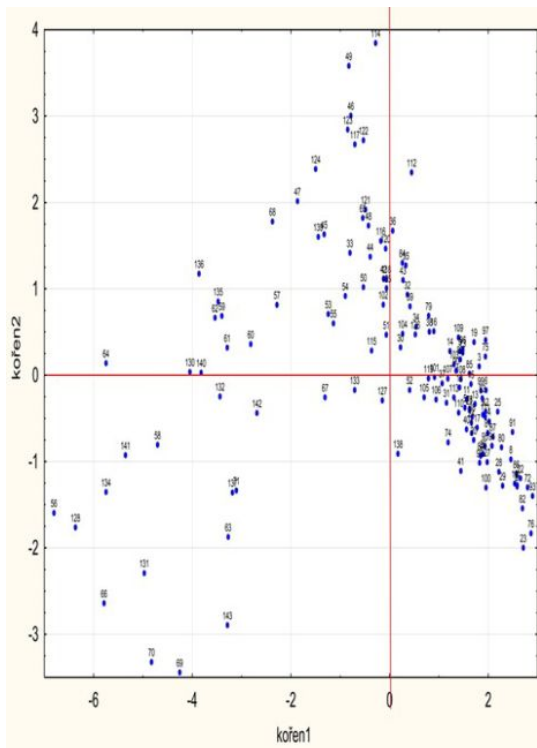
○ **Data:** Výběr DIABETES se týká vyšetření 144 jedinců pro 6 znaků:

- i (ID) identifikační číslo pacienta,
- x_1 (RELWT) značí relativní hmotnost pacienta,
- x_2 (GLUFAST) značí půst plasmové glukózy,
- x_3 (GLUTEST) značí prověrku plasmové glukózy,
- x_4 (INSTEST) značí plasmový inzulin v průběhu prověrky,
- x_5 (SSPG) značí ustálený rovnovážný stav plasmové glukózy,
- x_6 (GROUP) značí klinickou skupinu diabetes:

- 1 je zřetelný, vyložený diabetik,
- 2 je chemický,
- 3 je normální zdravý člověk.

i	x_1	x_2	x_3	x_4	x_5	x_6
1	0.81	80	356	124	55	3
...
144	1.11	328	1246	124	442	1





Úloha B40 Vliv vlákninové diety ve snídaňových lupínkách na nadváhu žen

Ženy konzumují lupínky různých typů, např. otrubové lupínky, gumové, obojí a konečně i s kontrolním druhem lupínek bez vlákniny k snídani tolik, na kolik mají hlad. Kontingenční tabulka ukazuje vztah 4 různých druhů lupínek a 4 hladin nevolnosti z přesyčení. χ^2 -test může otestovat, zda je přesyčení nezávislé na druhu lupínek.

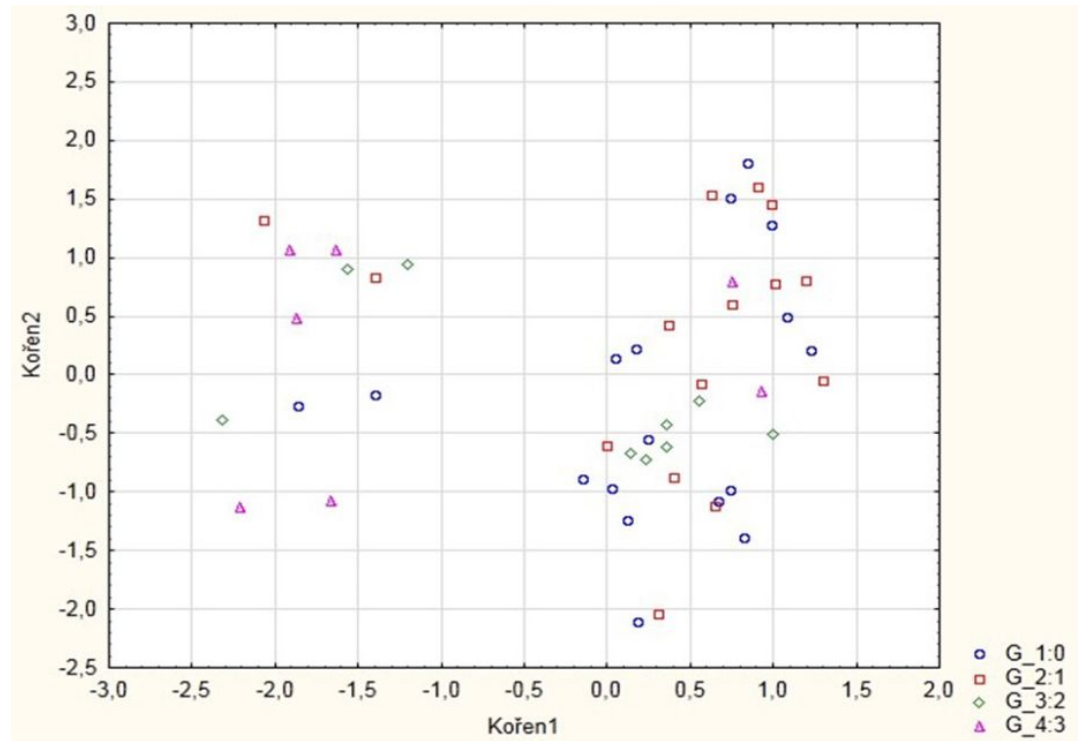
U diskriminační analýzy vezměte za závisle proměnnou údaje o stupni přesyčení x_5 (BLOAT) a ostatní znaky pak za nezávisle diskriminátory.

Data: Celkem 12 žen v řádcích bylo testováno na 3 diety ve sloupcích a na kontrolní dietu. Před každým jídlem otlyky jedly lupínky a) s otrubami, b) s gumovými vláknky, c) s kombinací obou, a konečně d) bez vláken. Celkem bylo provedeno 48 pokusů a užitó 5 sledovaných znaků.

Datový výběr CRACKER obsahuje:

- i (ID) pořadové číslo ženy,
- x_1 (CRACKER) značí typ vlákniny v lupínkách,
- x_2 (DIET) značí jednu ze 4 diet čili typu lupínek,
- x_3 (IDFEM) index sledované ženy,
- x_4 (DIGESTED) značí strávené kalorie,
- x_5 (BLOAT) značí stupeň přesyčení a nadýmání po jídle: 0 žádný, 1 malý, 2 střední a 3 vysoký.

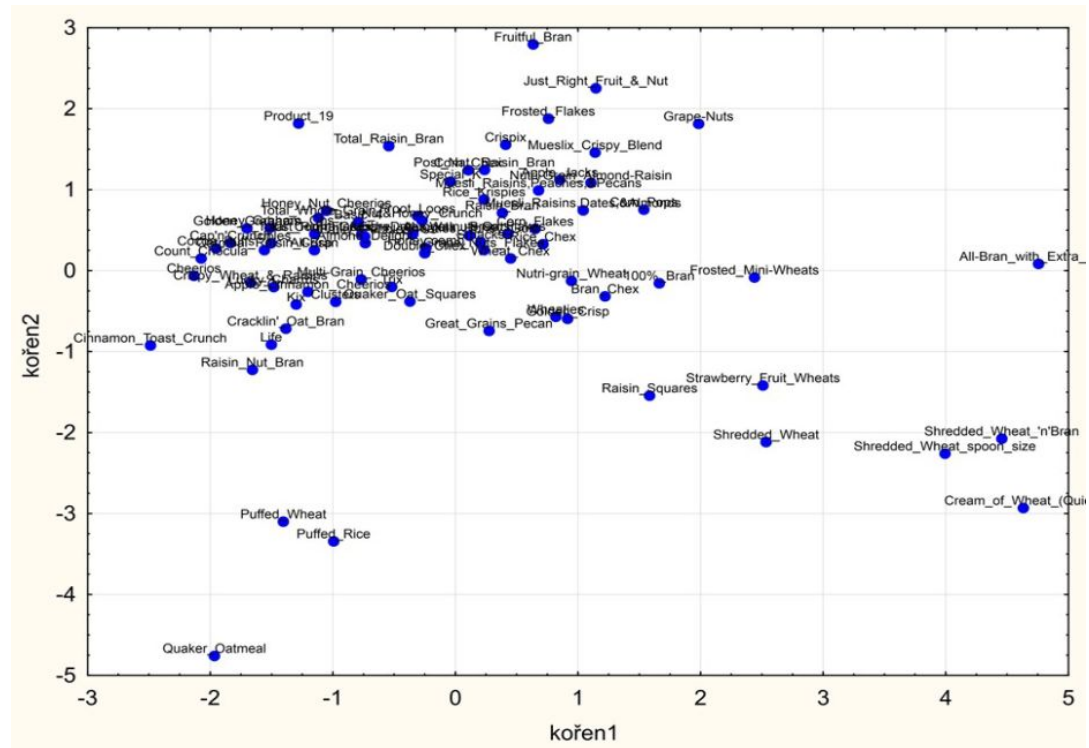
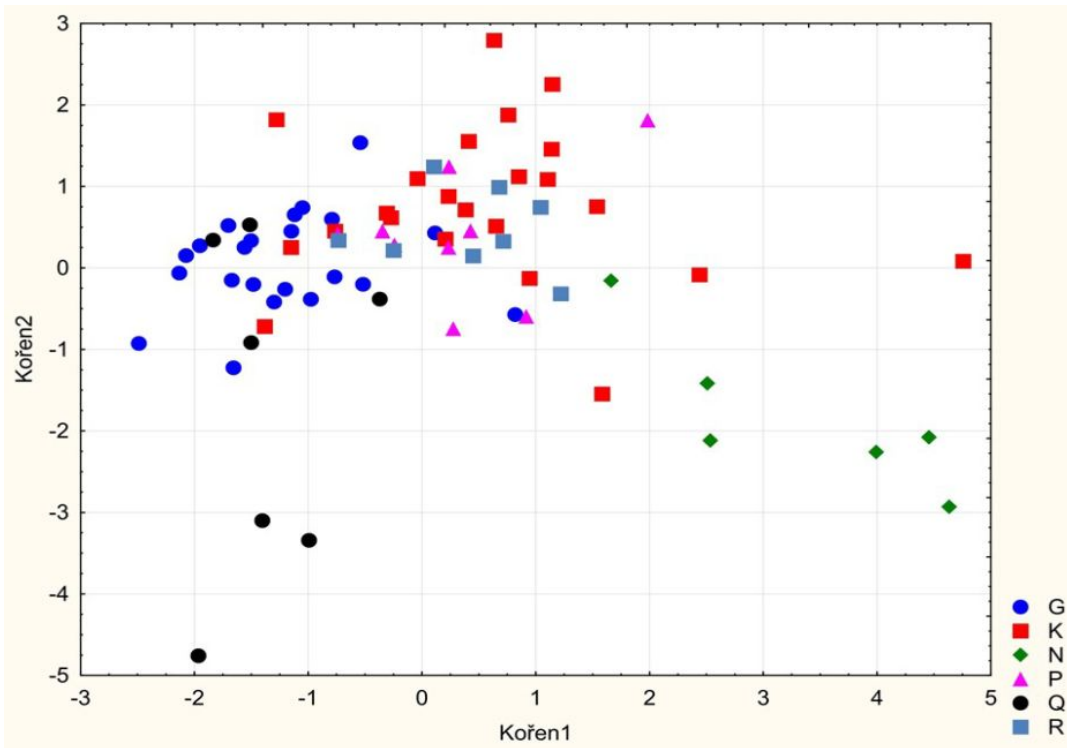
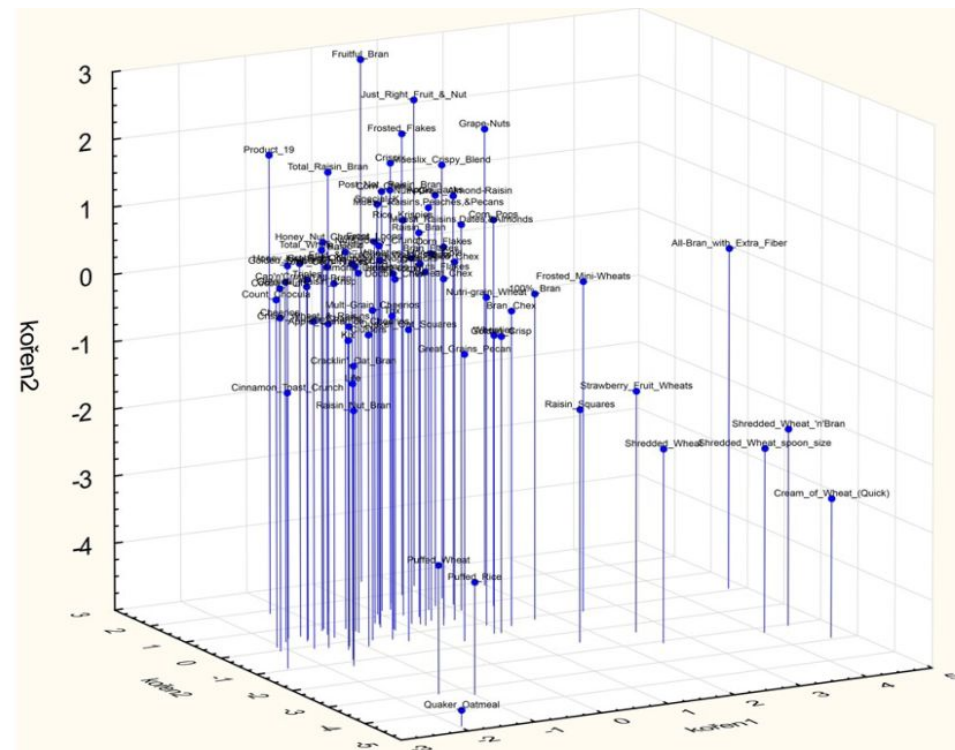
i	x_1	x_2	x_3	x_4	x_5
1	control	1	3	1772.84	0
...
12	bran	4	12	2287.52	0

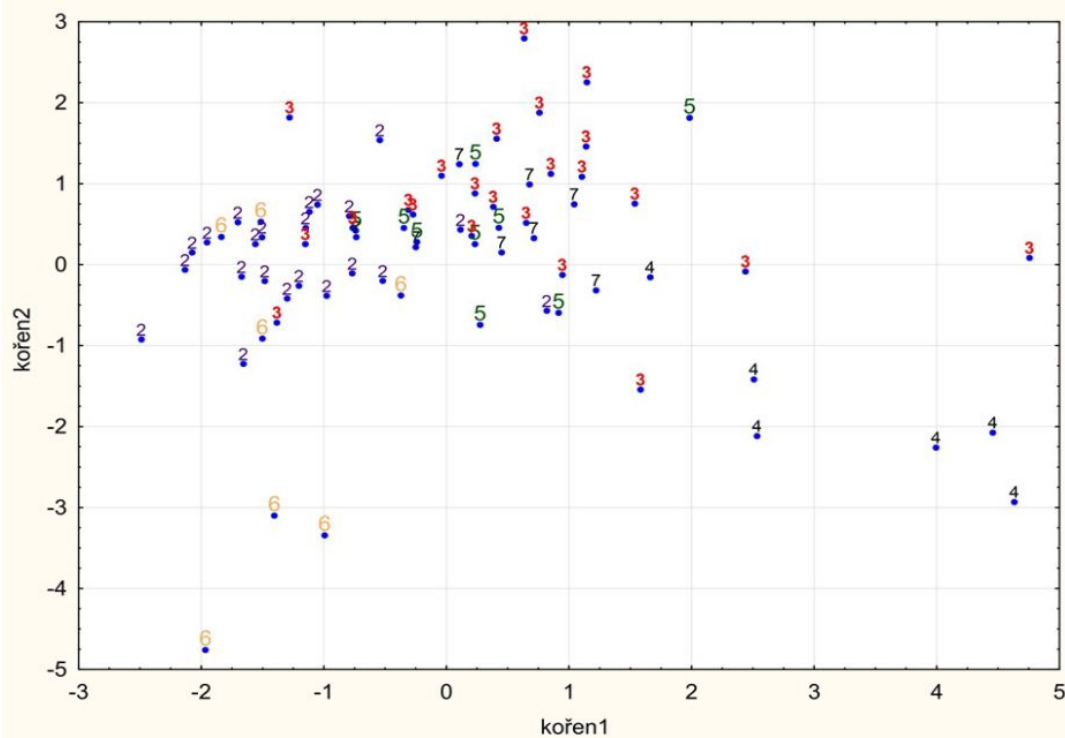


Úloha E22 Zdravotní klasifikace 77 druhů obilninových lupínek k snídani

Zdravotní norma u lupínek požaduje, že dospělí by měli denně zkonsumovat ne více než 30% kalorií ve formě tuku, muži potřebují okolo 63 g bílkovin a ženy 50 g bílkovin a zbytek poskytnou uhlohydráty (1 cal = 4.184 joulu). 1 g tuku obsahuje asi 9 kalorií, bílkoviny a uhlohydráty 4 kalorie na 1 g. Vhodná dieta by měla obsahovat 20 - 35 g dietní vlákniny. Níže analyzovaná data obsahují počet kalorií na 1 porci, gramy bílkovin a tuků, mg Na⁺ a K⁺iontů, gramy vlákniny, uhlohydrátů a cukru, typické procento FDA RDA vitaminů, hmotnost jedné porce a počet šálků na 1 porci podávanou. Data obsahují umístění na regále v prodejně: nahore, ve středu či dole u země.

Data: Výběr CEREALS obsahuje 77 druhů obilninových lupínek v rádcích a 15 znaků:
i (ID) je index druhu lupínek, *j* (J) je název druhu obilninových lupínek,
x1 (MFR) značí výrobce lupínek: A je American Home Food Products; G je General Mills; K je Kelloggs; N je Nabisco; P je Post; Q je Quaker Oats; R je Ralston Purina,
x2 (TYPE) značí formu v jaké je lupinka konzumována: C je studená, H je horká,
x3 (CALORIES) značí počet kalorií jedné porce,
x4 (PROTEIN) značí obsah bílkovin v gramech v jedné porci,
x5 (FAT) značí obsah tuku v gramech v jedné porci,
x6 (SODIUM) značí obsah sodných iontů v miligramech v jedné porci,
x7 (FIBER) značí obsah dietní vlákniny v jedné porci,
x8 (CH) značí obsah komplexních uhlohydrátů v gramech,
x9 (SUGARS) značí obsah cukrů v gramech v jedné porci,
x10 (POTASS) značí obsah draselných iontů v miligramech v jedné porci,
x11 (VITAMINS) značí obsah vitaminů a minerálů v jedné porci (0, 25, nebo 100) ukazuje typické % doporučených FDA,
x12 (SHELF) značí umístění na regále, jehož výška je číslována směrem od země 1 (nejníže), 2, nebo 3 (nejvýše),
x13 (WEIGHT) značí hmotnost jedné porce v uncích,
x14 (CUPS) značí počet šálků v jedné porci,
x15 (RATING) značí hodnocení, ocenění druhu lupínek.
 Číslo -1 značí chybějící numerickou informaci.





Kořeny odstraněny	Test chí-kvadrát po odstranění post. kořenů (E22)					
	Vlastní číslo	Kan. R	Wilk. Lambda	Chi-kv.	sv	p-hodn.
0	1,824023	0,803676	0,100675	149,2308	70	0,000000
1	0,971859	0,702043	0,284308	81,7502	52	0,005270
2	0,428362	0,547629	0,560616	37,6167	36	0,395068
3	0,199115	0,407495	0,800763	14,4424	22	0,885038
4	0,041442	0,199482	0,960207	2,6394	10	0,988685

Proměnná	Klasifikační funkce; grupovací : MFR (E22)				
	G_1:2 p=,28947	G_2:3 p=,30263	G_3:4 p=,07895	G_4:5 p=,11842	G_5:6 p=,10526
TYPE	431	429	443	429	431
CALORIES	1444741	1444724	1444645	1444730	1444725
PROTEIN	-21234981	-21234726	-21233572	-21234813	-21234741
FAT	10976930	10976795	10976200	10976840	10976805
SODIUM	353534	353530	353511	353531	353530
FIBER	-22336642	-22336371	-22335156	-22336463	-22336389
CH	-7087424	-7087338	-7086952	-7087368	-7087344
SUGARS	4703311	4703254	4702999	4703273	4703257
POTASS	220361	220359	220347	220360	220359
VITAMINS	332383	332379	332361	332380	332379
SHELF	1966	1966	1963	1966	1967
WEIGHT	41729	41722	41718	41721	41730
CUPS	-5056	-5058	-5059	-5059	-5054
RATING	6487518	6487440	6487088	6487467	6487445
Konstant	-178184354	-178180076	-178160726	-178181531	-178180305

Úloha B77 Rozlišení mono-, di- a trifazické hormonální antikoncepce u žen

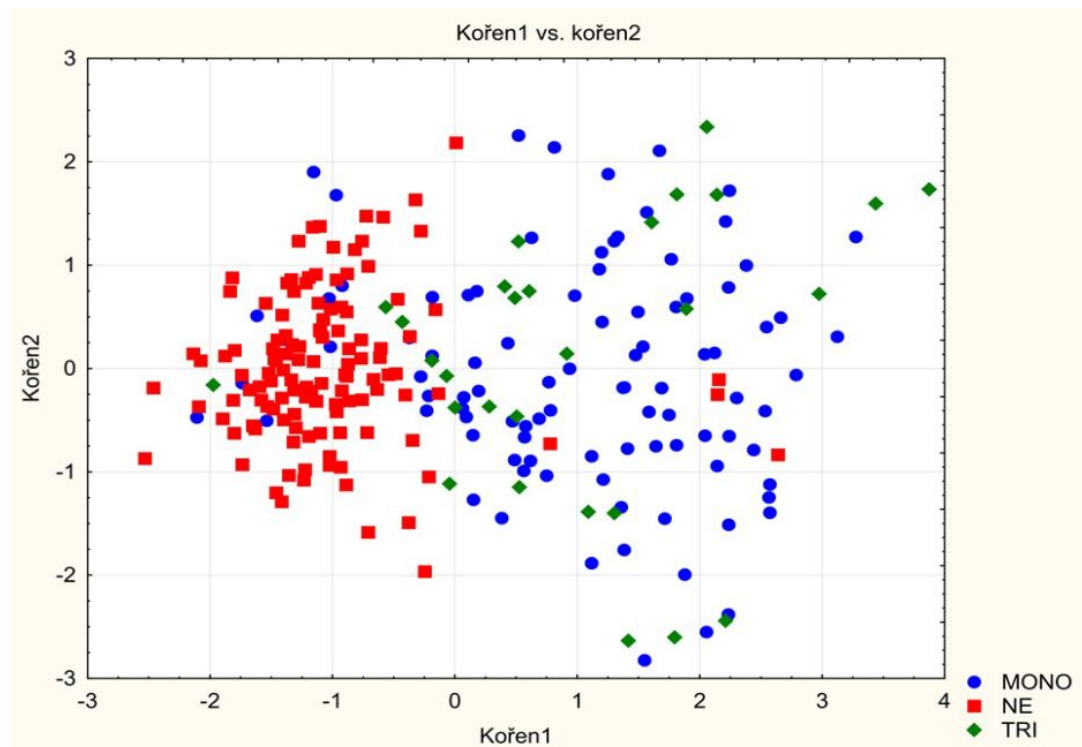
247 žen se dělí do tří tříd podle užívání hormonální antikoncepce.

V 1. třídě je 123 žen, které neužívají hormonální antikoncepci (**HA** = NE) – závisle p., ve 2. třídě je 152 žen užívající monofazickou (**HA** = MONO) antikoncepci, ve 3. třídě je 95 žen, které užívají trifazickou hormonální antikoncepci (**HA** = TRI).

Cílem diskriminační analýzy je zjistit, zda lze na základě diskriminátorů zde steroidních hormonů a proteinem přenášejícím steroidní hormony (**SHBG**) rozpoznat ženy neužívající a užívající hormonální antikoncepci, a dále rozpoznat mono- a trifazickou antikoncepci, lišící se dávkováním účinné látky.

○ Data: Zdrojová data s popisem znaků ve sloupcích 5 diskriminatory (**Kortizol, Testosteron, SHBG, DHEA a DHEAS**).

HA	Kortizol	SHBG	DHEAS	DHEA	Testosteron
i	x1	x2	x3	x4	x5
MONO	1383	85.2	5.68	46.28	1.62
.....
TRI	620	86.5	4.16	12.02	1.68



Kofeny odstraněny	Test chí-kvadrát po odstranění post. kofenů (B68-B77)					
	Vlastní číslo	Kan. R	Wilk. Lambda	Chi-kv.	sv	p-hodn.
0	1,087462	0,721768	0,474520	180,3991	10	0,000000
1	0,009547	0,097245	0,990543	2,2994	4	0,680877

Skup.	Průměry kan. proměnných (B68-B77)	
	Kofen1	Kofen2
MONO	1,02516	-0,076571
NE	-1,04063	0,000854
TRI	1,05541	0,247217

Úloha E19 Klasifikace smrku ztepilého z pokusné plochy v Orlických horách

Při klasifikaci 50 stromů smrku ztepilého (*Picea abies* L.) z pokusné plochy v Orlických horách byl vyšetřován dle 6 charakteristik každý strom zvlášť.

Data: Výběr STROMY:

i (ID) index stromu,

x1 (D) značí výčetní tloušťku stromu v cm, měřenou ve výšce 1.3 m nad zemí,

x2 (RTP) značí roční tloušťkový přírůstek v cm (tamtéž),

x3 (H) značí výšku stromu v m,

***x4* (KRAFT) značí Kraftovu klasifikaci stromů v porostu: 1 předrůstavý, 2 úrovnový, 3 částečně úrovnový,**

x5 (ZDRAV) značí zdravotní stav: 1 olistění 90 – 100 %, 2 olistění 70 – 80 %, 3 olistění pod 60 %, *x6* (RAŠ) značí rašení letorostů: 1 pozdní, 2 střední, 3 časně.

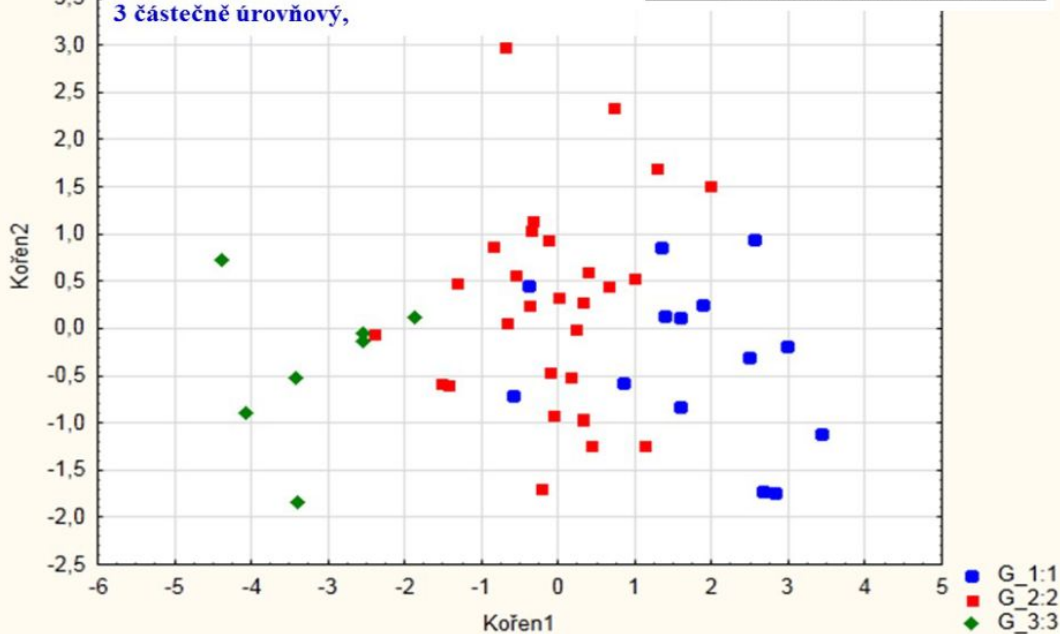
Proměnná	Standardiz. koeficienty (B68-B77) pro kanonické proměnné	
	Kofen1	Kofen2
B77x1	0,545676	0,915162
B77x2	0,723332	-0,865965
B77x3	0,204730	-0,234211
B77x4	-0,277455	-0,206162
B77x5	-0,072968	0,261355
Vlastní	1,087462	0,009547
KumPodíl	0,991297	1,000000

Proměnná	Faktorová strukturální matice (B68-B77) Korelační proměnné - Kanonické kofeny (vnitřní korelace)	
	Kofen1	Kofen2
B77x1	0,705334	0,637291
B77x2	0,840542	-0,458747
B77x3	-0,162699	0,016145
B77x4	-0,157594	0,037997
B77x5	0,045107	0,119113

<i>i</i>	<i>x1</i>	<i>x2</i>	<i>x3</i>	<i>x4</i>	<i>x5</i>	<i>x6</i>
1	8.9	0.8	4.60	2	1	2
...
50	7.0	0.5	5.80	2	3	3

Kraftovu klasifikace stromů v porostu:

- 1 předrůstavý,
- 2 úrovnový,
- 3 částečně úrovnový,



Kofeny odstraněny	Test chí-kvadrát po odstranění post. kofenů (E01-E19)					
	Vlastní číslo	Kan. R	Wilk. Lambda	Chi-kv.	sv	p-hodn.
0	2,425084	0,841449	0,268798	59,12081	10	0,000000
1	0,086183	0,281682	0,920655	3,72014	4	0,445204

Kofeny odstraněny	Test chí-kvadrát po odstranění post. kofenů (E01-E19)					
	Vlastní číslo	Kan. R	Wilk. Lambda	Chi-kv.	sv	p-hodn.
0	2,425084	0,841449	0,268798	59,12081	10	0,000000
1	0,086183	0,281682	0,920655	3,72014	4	0,445204

Skup.	Průměry kan. proměnných (E01-E19)	
	Kofen1	Kofen2
G_1:1	1,75688	-0,314043
G_2:2	-0,08173	0,241716
G_3:3	-3,17518	-0,373307

Proměnná	Faktorová strukturální matice (E01-E19) Korelační proměnné - Kanonické kofeny (vnitřní korelace)	
	Kofen1	Kofen2
E19x1	0,806638	0,521042
E19x2	0,277166	0,507606
E19x3	0,653265	-0,292438
E19x5	0,272048	-0,523211
E19x6	0,033093	0,174083

Proměnná	Klasifikační funkce; grupovací : E19x4 (E01-E19)		
	G_1:1 p=,28000	G_2:2 p=,58000	G_3:3 p=,14000
E19x1	9,1286	7,8943	4,7326
E19x2	-12,7102	-9,8558	-7,5089
E19x3	10,0572	8,4413	7,5413
E19x5	7,4350	5,8572	4,0769
E19x6	2,0369	2,3188	2,1130
Konstant	-84,6280	-62,2406	-35,8244

Skup.	Klasifikační matice (E01-E19) Řádky: pozorované klasifikace Sloupce: předpovězené klasifikace			
	% správných	G_1:1 p=,28000	G_2:2 p=,58000	G_3:3 p=,14000
	G_1:1	71,42857	10	4
G_2:2	89,65517	2	26	1
G_3:3	85,71429	0	1	6
Celkem	84,00000	12	31	7

Proměnná	Standardiz. koeficienty (E01-E19) pro kanonické proměnné	
	Kofen1	Kofen2
E19x1	0,857251	0,679610
E19x2	-0,197606	0,291570
E19x3	0,369694	-0,827552
E19x5	0,449677	-0,368888
E19x6	-0,017054	0,361201
Vlastní	2,425084	0,086183
KumPodíl	0,965681	1,000000

