

Univerzita Karlova v Praze  
Matematicko-fyzikální fakulta

## BAKALÁŘSKÁ PRÁCE



Lenka Goduřová

### **Vybrané problémy a metody při zpracování mnohorozměrných finančních dat**

Katedra pravěpodobnosti a matematické statistiky

Vedoucí bakalářské práce: RNDr. Jitka Zichová, Dr.

Studijní program: Matematika

Studijní obor: Finanční matematika

Praha 2014

Na tomto mieste by som sa rada poďakovala vedúcej mojej bakalárskej práce RNDr. Jitke Zichovej, Dr. za cenné rady a ochotu, s ktorou mi venovala svoj čas a za poskytnuté materiály, mojim rodičom a súrodencom za všestrannú podporu počas celého môjho doterajšieho štúdia a Petronelle Antoniewiczovej za trpezlivosť a podporu.

Prohlašuji, že jsem tuto bakalářskou práci vypracovala samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Univerzita Karlova v Praze má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle § 60 odst. 1 autorského zákona.

V Praze dne 22.5.2014

Lenka Godul'ová

**Název práce:** Vybrané problémy a metody při zpracování mnohorozměrných finančních dat

**Autor:** Lenka Goduřová

**Katedra:** Katedra pravděpodobnosti a matematické statistiky

**Vedoucí bakalářské práce:** RNDr. Jitka Zichová, Dr.

**Abstrakt:** Předložená bakalářská práce se zabývá zpracováním vícerozměrných dat. Úkolem bylo aplikovat vybrané metody na finanční data. Skládá se z teoretické části a zpracování konkrétní databáze. V prvních čtyřech kapitolách jsou shrnuté základní vztahy a pojmy týkající se náhodného vektora, náhodné veličiny, vícerozměrných dat a testu nezávislosti v kontingenční tabulce. Následující část je věnována popisu vybraných metod, kterými jsou shluková analýza a diskriminační analýza. V praktické části jsou tyto metody aplikované na databázi klientů německé banky.

**Klíčová slova:** náhodný vektor, mnohorozměrné rozdělení, mnohorozměrný náhodný výběr, kontingenční tabulka, shluková analýza, diskriminační analýza

**Title:** Selected problems and methods in multivariate data analysis

**Author:** Lenka Goduřová

**Department:** Department of Probability and Mathematical Statistics

**Supervisor:** RNDr. Jitka Zichová, Dr.

**Abstract:** The bachelor thesis deals with processing multidimensional data. The task was to apply selected methods on financial data. The thesis is composed of the theoretical section and the analysis of a particular database. The first four chapters deal with basic relations and definitions concerning random vector and variable, multidimensional data and the independence test in a contingency table. The following section is devoted to defining the particular methods selected: cluster analysis and discriminant analysis. In the practical section these methods are applied to a database of clients of a German bank.

**Keywords:** random vector, multivariate distribution, multivariate random variable, contingency table, cluster analysis, discriminant analysis.

**Názov práce:** Vybrané problémy a metódy pri spracovaní mnohorozmerných finančných dát

**Autor:** Lenka Goduľová

**Katedra:** Katedra pravděpodobnosti a matematické statistiky

**Vedúcí bakalárskej práce:** RNDr. Jitka Zichová, Dr.

**Abstrakt:** Predložená bakalárska práca sa zaoberá spracovaním mnohorozmerných dát. Úlohou bolo aplikovať vybrané metódy na finančné dáta. Pozostáva z teoretického základu a spracovania konkrétnej databáze. V prvých štyroch kapitolách sú zhrnuté základné vzťahy a pojmy týkajúce sa náhodného vektora, náhodnej veličiny, mnohorozmerných dát a testu nezávislosti v kontingenčnej tabuľke. Následná časť je venovaná popisu vybraných metód, ktorými sú zhluková analýza a diskriminačná analýza. V praktickej časti sú tieto metódy aplikované na databázu klientov nemeckej banky.

**Kľúčové slová:** náhodný vektor, mnohorozmerné rozdelenie, mnohorozmerný náhodný výber, kontingenčná tabuľka, zhluková analýza, diskriminačná analýza

# Obsah

Úvod	1
<b>1 Vektory a matice</b>	<b>2</b>
1.1 Vektor a matice . . . . .	2
1.1.1 Vektor . . . . .	2
1.1.2 Matice . . . . .	4
<b>2 Náhodná veličina a náhodný vektor</b>	<b>5</b>
2.1 Náhodná veličina . . . . .	5
2.2 Náhodný vektor . . . . .	7
2.3 Distribučná funkcia . . . . .	8
<b>3 Mnohorozmerné dáta</b>	<b>10</b>
3.1 Mnohorozmerné normálne rozdelenie . . . . .	10
3.2 Mnohorozmerný náhodný výber . . . . .	11
3.3 Dátová matica . . . . .	11
3.4 Objekty a typy premenných . . . . .	12
<b>4 Mnohorozmerné štatistické metódy</b>	<b>15</b>
4.1 Zhluková analýza . . . . .	15
4.1.1 Hierarchická aglomeratívna metóda . . . . .	17
4.1.2 Nehierarchická metóda . . . . .	18
4.2 Diskriminačná analýza . . . . .	18
<b>5 Spracované dáta</b>	<b>22</b>
5.1 Kontingenčná tabuľka . . . . .	22
5.2 Zhluková analýza . . . . .	26

5.2.1	Hierarchická aglomeratívna metóda . . . . .	26
5.2.2	Nehierarchická K-means metóda . . . . .	30
5.3	Diskriminačná analýza . . . . .	35
	<b>Bibliografia</b>	<b>39</b>
	<b>Zoznam obrázkov</b>	<b>40</b>
	<b>Zoznam tabuliek</b>	<b>41</b>
	<b>Zoznam príloh</b>	<b>42</b>

# Úvod

V dnešnej dobe existuje mnoho metód na analýzu dát. V tejto práci si ukážeme dve metódy spracovania mnohorozmerných dát. Dáta tohoto typu používajú napríklad poisťovne, či banky používajúce rozsiahle databáze klientov. Popíšeme si zhlukovú a diskriminačnú analýzu.

Pre porozumenie zhrnieme v troch prvých kapitolách teóriu, prostredníctvom ktorej budeme nazerať na skúmaný problém. Postupne sa budeme venovať vektorom a maticiam, náhodným veličinám či vektorom a v poslednej rade aj mnohorozmerným dátam.

V štvrtej kapitole jednotlivé metódy popíšeme. Rozoberieme spôsob ich použitia a uvedieme niektoré ich vlastnosti. Posledná kapitola je venovaná aplikácii spomínaných analýz na dátach nemeckej banky pomocou softvérových produktov Mathematica, Microsoft Office Excel a NCSS. Výpočty a výstupné protokoly z jednotlivých softvérov nájdeme na priloženom CD.



# Kapitola 1

## Vektory a matice

V prvej kapitole definujeme základné pojmy, ktoré budeme následne využívať. Oboznámime sa s pojmami vektor a matica. Vymenujeme si ich niektoré vlastnosti. Teoretické poznatky pochádzajú z [6], [7] a [8].

### 1.1 Vektor a matice

#### 1.1.1 Vektor

Telesom  $T$  nazveme aspoň dvojprvkovú množinu spĺňajúcu axiómy pre operácie násobenia a sčítania. Týmito axiómami sú myslené komutatívita sčítania, asociatívne zákony, existencia nulového, jednotkového a opačného prvku, inverzný prvok a distributívne zákony. Ak navyše platí komutatívita pre násobenie, potom hovoríme o komutatívnom telese.

Vektorovým priestorom nad telesom  $T$  je neprázdna množina  $V$ , ktorej všetky prvky spĺňajú nasledujúce axiómy:

$$\begin{aligned}
\forall \mathbf{a}, \mathbf{b} \in V : \quad & \mathbf{a} + \mathbf{b} = \mathbf{b} + \mathbf{a} \\
\forall \mathbf{a}, \mathbf{b}, \mathbf{c} \in V : \quad & (\mathbf{a} + \mathbf{b}) + \mathbf{c} = \mathbf{a} + (\mathbf{b} + \mathbf{c}) \\
\forall \mathbf{a}, \mathbf{b} \in V, \forall \alpha \in T : \quad & \alpha(\mathbf{a} + \mathbf{b}) = \alpha\mathbf{a} + \alpha\mathbf{b} \\
\forall \mathbf{a} \in V, \forall \alpha, \beta \in T : \quad & (\alpha + \beta)\mathbf{a} = \alpha\mathbf{a} + \beta\mathbf{a} \\
\forall \mathbf{a} \in V, \forall \alpha, \beta \in T : \quad & \alpha(\beta\mathbf{a}) = (\alpha\beta)\mathbf{a} \\
\forall \mathbf{a} \in V, \exists 1 \in T : \quad & 1\mathbf{a} = \mathbf{a} \\
\forall \mathbf{a} \in V, \exists 0 \in T, \exists \mathbf{0} \in V : \quad & 0\mathbf{a} = \mathbf{0}.
\end{aligned} \tag{1.1}$$

Prvky telesa  $T$  nazývame skaláry a prvkom vektorového priestoru hovoríme vektor. Príkladom vektorového priestoru je vektorový priestor nad telesom reálnych čísel  $\mathbb{R}$ . V texte práce ním bude  $k$ -rozmerný reálny priestor, ktorý označíme  $\mathbb{R}^k$ . Vektor dĺžky  $k$  si môžeme predstaviť ako usporiadaný rad čísel  $x_1, x_2, \dots, x_k$ . Stĺpcový vektor  $\mathbf{0}$ , ktorý je uvedený v (1.1), nazveme nulovým vektorom. Má všetky zložky rovné nule. Stĺpcový vektor budeme značiť  $\mathbf{x} = (x_1, \dots, x_k)^T$ .

Vektory sa vyznačujú rôznymi vlastnosťami a charakteristikami. Nás bude zaujímať lineárna závislosť. Hovoríme, že  $k$ -rozmerné vektory  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_s$  sú lineárne závislé, ak existujú reálne čísla  $c_1, c_2, \dots, c_s$  (aspoň jedno nenulové) také, že

$$c_1\mathbf{x}_1 + c_2\mathbf{x}_2 + \dots + c_s\mathbf{x}_s = \mathbf{0}. \tag{1.2}$$

V určitých prípadoch je možné lineárnu závislosť určiť na základe toho, ako jednotlivé vektory vyzerajú. V prípade, že sa medzi skúmanými vektormi nachádza nulový vektor alebo nejaký vektor je násobkom iného, sú tieto vektory lineárne závislé.

## 1.1.2 Matice

Maticou  $\mathbf{A}$  typu  $m \times n$  budeme rozumieť súbor prvkov:

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \dots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}, \quad (1.3)$$

$\mathbf{A} = (a_{i,j})$ , kde  $i = 1, \dots, m$  a  $j = 1, \dots, n$ ,  $a_{i,j} \in R$ .

Vektor  $(a_{i1}, a_{i2}, \dots, a_{in})$  je  $i$ -tým riadkom a  $(a_{1j}, a_{2j}, \dots, a_{mj})^T$  je  $j$ -tým stĺpcom matice  $A$ . Transponovaná matica  $A^T$  typu  $n \times m$  je matica, vytvorená vzájomnou výmenou riadkov a stĺpcov matice  $A$ . Za platnosti  $A = A^T$ , hovoríme o matici symetrickej. Štvorcová matica  $A$  je špeciálnym typom matice  $A$ , kde počet riadkov je rovný počtu stĺpcov.

Matice môžeme sčítať alebo násobiť ľubovoľnou konštantou. Rovnosť matíc  $A$  a  $B$  nastáva práve vtedy keď  $a_{ij} = b_{ij}$ ,  $i = 1, 2, \dots, m$  a  $j = 1, 2, \dots, n$ . Matica s prvkami  $a_{ij} = 0$ ,  $i = 1, 2, \dots, m$  a  $j = 1, 2, \dots, n$  sa nazýva nulová matica. Prvky  $a_{ii}$  tvoria (hlavnú) diagonálu štvorcovej matice  $A$ . Štvorcovú maticu  $A$  typu  $n \times n$ , ktorá má mimo diagonálu samé nuly, nazveme diagonálnou maticou. Ak sú ešte navyše na diagonále samé jednotky, tak hovoríme o jednotkovej matici. Štvorcovú maticu  $A$  nazveme pozitívne semidefinitnou (definitnou), ak platí  $\mathbf{x}^T A \mathbf{x} \geq 0$  ( $\mathbf{x}^T A \mathbf{x} > 0$ ) pre  $\mathbf{x} \neq 0$ . Determinant matice  $A$  budeme značiť symbolom  $|A|$ . Viac o determinantoch a vlastnostiach matíc môžeme nájsť v [7].

# Kapitola 2

## Náhodná veličina a náhodný vektor

Náhoda je niečo, čo vopred nevieme ovplyvniť. V elementárnej teórii pravdepodobnosti pracujeme s pojmom náhodný pokus, ktorého výsledkom sú náhodné javy. My budeme pracovať s elementárnymi javmi  $\omega$ . Sú to všetky možné, ďalej nezjednodušiteľné výsledky náhodného pokusu. Náhodná veličina potom priradí náhodným javom číselné hodnoty. V tejto kapitole si ukážeme spôsob, prostredníctvom ktorého sa dostaneme k náhodným veličinám pomocou teórie pravdepodobnosti a zavedieme pojmy stredná hodnota a rozptyl náhodnej veličiny. Pokračovať budeme ďalším dôležitým pojmom náhodný vektor, rozoberieme si podrobnejšie rozdelenie pravdepodobnosti a distribučnú funkciu. Teoretické poznatky pochádzajú z [2], [3], [4], [8] a [11].

### 2.1 Náhodná veličina

Trojicu  $(\Omega, \mathcal{A}, P)$  nazveme pravdepodobnostný priestor. Ľubovoľná množina  $\Omega$  je priestor elementárnych javov  $\omega$ , na ktorom je definovaná  $\sigma$  - algebra  $\mathcal{A}$  ako neprázdny systém podmnožín priestoru  $\Omega$  uzavretý na spočítateľné zjednotenie a doplnok. Symbol  $P$  predstavuje pravdepodobnostnú mieru. Borelovskú  $\sigma$  - algebru na  $\mathbb{R}$  označíme  $\mathcal{B}_1$ . Je to najmenšia  $\sigma$  - algebra obsahujúca otvorené intervaly. Merateľné zobrazenie z  $(\Omega, \mathcal{A}, P)$  do  $(\psi, \mathcal{B})$  nazveme náhodnou veličinou  $X$ , kde  $\psi$  je výberový priestor a  $\mathcal{B}$  je  $\sigma$  - algebra na  $\psi$ . Pre  $\psi = \mathbb{R}$  a  $\mathcal{B} = \mathcal{B}_1$

reprezentuje každá hodnota  $X(\omega)$  reálne číslo a platí  $\{\omega \in \Omega : X(\omega) \in B\} \in \mathcal{A}$ ,  $\forall B \in \mathcal{B}_1$ . Pravdepodobnostné rozdelenie náhodnej veličiny  $X$  je určené systémom  $\{P[X \in B], B \in \mathcal{B}_1\}$ . Špeciálne máme pre  $B = (-\infty, x]$  funkciu

$$F_x(x) = P[X \leq x] = P\{\omega \in \Omega : X(\omega) \leq x\} \quad x \in R. \quad (2.1)$$

Funkcia  $F_x(x)$  sa nazýva distribučná funkcia náhodnej veličiny  $X$ . Pozrieme sa na niektoré charakteristiky náhodných veličín. Stredná hodnota náhodnej veličiny  $X$  je definovaná predpisom:

$$EX = \int_{\Omega} X(\omega) dP(\omega). \quad (2.2)$$

Z vlastností integrálu plynie  $E(a + bX) = a + bEX$  pre ľubovoľné  $a, b \in R$ . Obecný moment  $k$ -tého rádu je

$$\mu'_k = EX^k, \quad k = 0, 1, \dots \quad (2.3)$$

ak existuje integrál na pravej strane. Pre  $k = 1$  máme strednú hodnotu.

Číslo

$$\mu_k = E(X - EX)^k, \quad k = 0, 1, \dots \quad (2.4)$$

nazveme centrálnym momentom rádu  $k$  v prípade, že integrál na pravej strane existuje. Najvýznamnejším centrálnym momentom je  $\mu_2$ . Nazývame ho rozptyl. Označíme ho  $\sigma^2$  alebo  $\text{var}X$ . Platí  $\text{var}(a + bX) = b^2\text{var}X$  pre ľubovoľné  $a, b \in R$ . Odmocninou z rozptylu je smerodajná odchýlka.

Vzťah medzi dvoma náhodnými veličinami nám priblíži kovariancia a korelácia. Predpokladajme konečné a nenulové rozptyly náhodných veličín  $X$  a  $Y$ . Potom kovarianciou týchto náhodných veličín budeme rozumieť

$$\text{cov}(X, Y) = E(X - EX)(Y - EY). \quad (2.5)$$

Z kovariancie dvoch normovaných náhodných veličín dostaneme korelačný koeficient

$$\rho(X, Y) = \text{cov}\left(\frac{X - EX}{\sqrt{\text{var}X}}, \frac{Y - EY}{\sqrt{\text{var}Y}}\right) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}X \text{var}Y}}. \quad (2.6)$$

Je možné ukázať, že  $\rho(X, Y) = \pm 1$  práve vtedy, keď  $Y = a \pm bX$ ,  $b > 0$ . Čím je  $\rho(X, Y)$  bližšie k 1 respektíve -1, tým je spomenutá lineárna závislosť náhodných veličín  $X$  a  $Y$  silnejšia.

## 2.2 Náhodný vektor

Nech  $R^k$  je  $k$ -rozmerný reálny priestor a  $\mathcal{B}_k$  je systém jeho borelovských podmnožín. Náhodným vektorom  $\mathbf{X}$  nazveme merateľné zobrazenie z  $(\Omega, \mathcal{A}, P)$  do  $(\psi, \mathcal{B})$ , kde  $\psi = R^k$  a  $\mathcal{B} = \mathcal{B}_k$ . Je to stĺpcový vektor o dĺžke  $k$ , ktorého zložky sú náhodné veličiny. Značíme  $\mathbf{X} = (X_1, X_2, \dots, X_k)^T$ .

Strednou hodnotou náhodného vektora  $\mathbf{X}$  budeme rozumieť stĺpcový vektor stredných hodnôt jednotlivých náhodných veličín,  $E\mathbf{X} = (EX_1, EX_2, \dots, EX_k)^T$ . Platí  $E(\mathbf{a} + \mathbf{B}\mathbf{X}) = \mathbf{a} + \mathbf{B}E\mathbf{X}$  pre ľubovoľný  $m$ -rozmerný reálny vektor  $\mathbf{a}$  a ľubovoľnú maticu  $\mathbf{B}$  typu  $m \times k$ . Variančnou maticou náhodného vektora  $\mathbf{X}$  je

$$\Sigma = \text{var}\mathbf{X} = \begin{pmatrix} \text{var}X_1 & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_k) \\ \text{cov}(X_2, X_1) & \text{var}X_2 & \dots & \text{cov}(X_2, X_k) \\ \dots & \dots & \dots & \dots \\ \text{cov}(X_k, X_1) & \text{cov}(X_k, X_2) & \dots & \text{var}X_k \end{pmatrix}. \quad (2.7)$$

Matica  $\Sigma$  je symetrická a pozitívne semidefinitvna. Pozitívna semidefinitnosť vyplýva z vlastnosti  $\text{var}(\mathbf{a} + \mathbf{B}\mathbf{X}) = \mathbf{B}\text{var}(\mathbf{X})\mathbf{B}^T$  a ďalej z predpokladu, že  $\text{var}Z$  náhodnej veličiny  $Z$  je nezáporné číslo (ak existuje). Teda pre ľubovoľné  $\mathbf{c} \in R^k$  a pro  $Z = \mathbf{c}^T\mathbf{X}$ , platí  $\text{var}Z = \text{var}(\mathbf{c}^T\mathbf{X}) = \mathbf{c}^T\text{var}(\mathbf{X})\mathbf{c} \geq 0$ .

Ďalej je možné zaviesť korelačnú maticu náhodného vektora  $\mathbf{X}$  predpisom

$$\Gamma = \text{cor}\mathbf{X} = \begin{pmatrix} \rho(X_1, X_1) & \rho(X_1, X_2) & \dots & \rho(X_1, X_k) \\ \rho(X_2, X_1) & \rho(X_2, X_2) & \dots & \rho(X_2, X_k) \\ \dots & \dots & \dots & \dots \\ \rho(X_k, X_1) & \rho(X_k, X_2) & \dots & \rho(X_k, X_k) \end{pmatrix}. \quad (2.8)$$

Je možné ukázať, že  $\rho(X_i, X_i) = 1$ ,  $i = 1, \dots, k$  a  $\text{cor}\mathbf{X} = \mathbf{A}\text{var}\mathbf{X}\mathbf{A}$ , kde  $\mathbf{A} = \text{diag}\left\{\frac{1}{\sqrt{\text{var}X_1}}, \dots, \frac{1}{\sqrt{\text{var}X_k}}\right\}$ .

## 2.3 Distribučná funkcia

Na konci tejto kapitoly sa ešte raz vrátíme k distribučnej funkcii, ktorú sme definovali v (2.1). V matematickej štatistike sa používajú dva základné typy distribučných funkcií:

- Predpokladajme, že  $F$  je funkcia skokov, ktorá má konečný alebo spočítateľný počet bodov nespojitosti a  $p_j$  je veľkosť skoku funkcie  $F$  v bode  $x_j$ , kde  $\sum p_j = 1$ . Potom postupnosť pravdepodobností

$$p_j = P[X = x_j], \quad j = 1, 2, \dots \quad (2.9)$$

definuje diskkrétne rozdelenie pravdepodobnosti náhodnej veličiny  $X$ .

- Nech  $F$  je absolutne spojitá a existuje funkcia  $f(x)$  taká, že

$$F(x) = \int_{-\infty}^x f(t)dt, \quad (2.10)$$

potom funkcia  $f(x)$  sa nazýva hustota a definuje spojité rozdelenie náhodnej veličiny  $X$ .

Ak namiesto náhodnej veličiny  $X$  uvažujeme náhodný vektor  $\mathbf{X} = (X_1, X_2, \dots, X_k)^T$ , definujeme jeho distribučnú funkciu predpisom

$$F_{\mathbf{X}}(x_1, \dots, x_k) = P[X_1 \leq x_1, \dots, X_k \leq x_k]. \quad (2.11)$$

- Náhodný vektor s diskkrétne rozdelenými zložkami má združené rozdelenie určené pravdepodobnosťami

$$P[X_1 = x_1, X_2 = x_2, \dots, X_k = x_k], \quad (2.12)$$

pre všetky možné kombinácie hodnôt  $x_1, \dots, x_k$  jeho zložiek.

- Náhodný vektor so spojito rozdelenými zložkami má hustotu  $f_{\mathbf{X}}(x)$ , pre ktorú platí

$$F_{\mathbf{X}}(x_1, \dots, x_k) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_k} f_{\mathbf{X}}(t_1, \dots, t_k) dt_1, \dots, dt_k. \quad (2.13)$$

Marginálne rozdelenie jednotlivých zložiek  $X_1, X_2, \dots, X_k$  a rozdelenie ľubovoľného podvektora náhodného vektora  $\mathbf{X}$  je možné odvodiť zo združeného rozdelenia. Pre jednozložkový podvektor  $X_i$  platí

$$F_{X_i}(t) = \lim_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k \rightarrow \infty} F_{\mathbf{X}}(x_1, \dots, x_{i-1}, t, x_{i+1}, \dots, x_k), \quad (2.14)$$

$$P[X_i = t] = \Sigma \dots \Sigma P[X_1 = x_1, \dots, X_{i-1} = x_{i-1}, X_i = t, X_{i+1} = x_{i+1}, \dots, X_k = x_k], \quad (2.15)$$

$$f_{X_i}(t) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_{\mathbf{X}}(x_1, \dots, x_{i-1}, t, x_{i+1}, \dots, x_k) dx_1, \dots, dx_{i-1} dx_{i+1}, \dots, dx_k. \quad (2.16)$$

Vo vzťahoch (2.15) a (2.16) sčítame, respektíve integrujeme cez všetky možné hodnoty zložiek  $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k$ .



# Kapitola 3

## Mnohorozmerné dáta

Predmetom tejto kapitoly je pojem dátová matica. Vysvetlíme si jej rôzne úpravy a pozrieme sa na jej členenie. V praktickej časti budeme pracovať s kontingenčnou tabuľkou, ktorú taktiež popíšeme. Najprv si však definujeme mnohorozmerné normálne rozdelenie. Teoretické poznatky pochádzajú z [4], [8] a [11].

### 3.1 Mnohorozmerné normálne rozdelenie

Jednou z rozoberaných štatistických metód bude diskriminačná analýza. Bližšie bude popísaná v kapitole 4.2. Pre použitie tejto analýzy budeme potrebovať mnohorozmerné normálne rozdelenie, ktoré si teraz definujeme.

Máme náhodný vektor  $\mathbf{X} = (X_1, X_2, \dots, X_k)^T$ , vektor  $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_k)^T$  a symetrickú pozitívne semidefinitnú maticu  $\boldsymbol{\Sigma}$  typu  $k \times k$ . Potom náhodný vektor  $\mathbf{X}$  má  $k$ -rozmerné normálne rozdelenie s parametrami  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , ak pre ľubovoľný vektor  $\mathbf{s} \in \mathbb{R}$  platí  $\mathbf{s}^T \mathbf{X} \sim N(\mathbf{s}^T \boldsymbol{\mu}, \mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s})$ . V prípade, že matica  $\boldsymbol{\Sigma}$  je regulárna, existuje hustota, ktorá má tvar

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{1}{2}k} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}. \quad (3.1)$$

## 3.2 Mnohorozmerný náhodný výber

Uvažujme nezávislé a rovnako rozdelené náhodné vektory  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , kde  $\mathbf{X}_i = (X_{i1}, \dots, X_{ik})^T$  pre  $i = 1, \dots, n$ . Nech  $E\mathbf{X}_i = \boldsymbol{\mu}$ ,  $\text{var}\mathbf{X}_i = \boldsymbol{\Sigma}$  a  $\text{cor}\mathbf{X}_i = \boldsymbol{\Gamma}$  pre  $i = 1, \dots, n$ .

- Odhadom vektora stredných hodnôt  $\boldsymbol{\mu}$  je výberový priemer

$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \quad (3.2)$$

so složkami

$$\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}, \quad j = 1, \dots, k. \quad (3.3)$$

- Odhadom variančnej matice  $\boldsymbol{\Sigma}$  je výberová variančná matica

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^T \quad (3.4)$$

s prvkami

$$s_{rj} = \frac{1}{n-1} \sum_{i=1}^n (X_{ir} - \bar{X}_r)(X_{ij} - \bar{X}_j), \quad j, r = 1, \dots, k. \quad (3.5)$$

- Odhadom korelačnej matice  $\boldsymbol{\Gamma}$  je výberová korelačná matica

$$\mathbf{C} = \hat{\mathbf{A}}\mathbf{S}\hat{\mathbf{A}}, \quad (3.6)$$

kde  $\hat{\mathbf{A}} = \text{diag}\{\frac{1}{\sqrt{s_{11}}}, \dots, \frac{1}{\sqrt{s_{kk}}}\}$  s prvkami

$$r_{ij} = \frac{s_{ij}}{\sqrt{s_{ii}s_{jj}}}, \quad i, j = 1, \dots, k. \quad (3.7)$$

## 3.3 Dátová matica

Predpokladajme maticu typu  $n \times k$ , kde riadky určujú nejaké študované objekty a stĺpce zisťované znaky na týchto objektoch. Označíme ju

$$\mathbf{D} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \dots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix}, \quad (3.8)$$

kde prvok  $x_{ij}$  predstavuje hodnotu  $j$ -tého znaku na  $i$ -tom objekte,  $i = 1, 2, \dots, n$  a  $j = 1, 2, \dots, k$ . Hodnoty  $x_{ij}$  môžeme taktiež nazvať pozorovaniami či meraniami. Riadky matice  $\mathbf{D}$  sú  $k$ -rozmerné vektory  $\mathbf{x}_i$ . Nazývame ich viacrozmernými pozorovaniami. Pozorovania sú konkrétne realizácie náhodných veličín  $X_{ij}$ , respektíve náhodných vektorov  $\mathbf{X}_i$ . Riadky matice  $\mathbf{D}$  sú teda realizáciou  $k$ -rozmerného náhodného výberu  $\mathbf{X}_1, \dots, \mathbf{X}_n$ .

Dátová matica sa rozpadá horizontálne pri niektorých štatistických metódach, kde pozorujeme situácie rozčlenenia objektov do niekoľkých podmnožín. Jednotlivé objekty v danej množine majú podobné vlastnosti. Takéto rozčlenenie je niekedy dané a cieľom je analýza jednotlivých skupín (analýza rozptylu, môžeme nájsť v [4]) alebo je cieľom samotné rozčlenenie (zhluková analýza). Zhluková analýza bude popísaná v kapitole 4.1. Ak sú sledované znaky rôznej povahy alebo merané v rôznych jednotkách, doporučuje sa aplikovať zhlukovú analýzu na štandardizované dáta

$$x_{ij}^* = \frac{x_{ij} - \bar{x}_j}{\sqrt{s_{jj}}} \quad i = 1, 2, \dots, n \quad , j = 1, 2, \dots, k \quad (3.9)$$

kde  $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$  je priemerná hodnota  $j$ -tého znaku cez všetky namerané realizácie a  $\sqrt{s_{jj}}$  je výberová smerodatná odchýlka, teda prvok diagonály výberovej variančnej matice (3.4). Štandardizované dáta sú realizáciou náhodného výberu s nulovou strednou hodnotou a jednotkovým rozptylom.

### 3.4 Objekty a typy premenných

Záznamy o objektoch tvoria riadky dátovej matice. Pri použití rôznych štatistických metód predpokladáme primerane rozsiahly súbor objektov. V odbornej literatúre sa objekt nazýva aj ako prvok, jednotka respektíve štatistická jednotka a znak ako premenná. Na každom z  $n$  objektov pozorujeme realizáciu  $k$ -rozmerného náhodného vektora. Poznávanie závislosti medzi jednotlivými objektami prostredníctvom meraných znakov je typickým cieľom štatistickej analýzy. Klasifikácia premenných zo štatistického hľadiska je podrobnejšie popísaná napríklad v knihe [5].

Znaky je možné klasifikovať z rôznych hľadísk. Základné delenie je:

- *Kvalitatívne*, kde nie je možné jednoznačné vyjadrenie číslom
- *Kvantitatívne*, ktoré sa dajú vyjadriť jednoznačne číselne

Ďalšou charakteristikou premenných je počet hodnôt, ktoré môžu nadobúdať.

Rozlišujeme premenné :

- *Spojité*
- *Nespojité*

Spojité nadobúdajú nespočetne mnoho hodnôt. Najčastejšie vznikajú vážením a meraním objektov. Majú teda číselné vyjadrenie. V prípade vzniku nespojitých resp. diskretných premenných dochádza k nadobúdaníu spočítateľne alebo konečne mnoho hodnôt. Tieto premenné delíme na nominálne (ak nemá zmysel porovnávať ich hodnoty, respektíve zoradiť hodnoty podľa veľkosti), ordinálne (ak má zmysel ich hodnoty porovnávať). Poslednou skupinou nespojitých premenných sú alternatívne (inak binárne alebo dichotomické), ktoré nadobúdajú len dve hodnoty.

Varianty znakov sú hodnoty, ktoré môžu jednotlivé diskretné znaky nadobúdať. V prípade jedného študovaného znaku ( $k = 1$ ) s malým počtom hodnôt vytvoríme tabuľku, ktorá bude obsahovať všetky varianty a početnosť ich výskytu v súbore  $n$  objektov. Početnosť  $l$ -tej varianty značíme  $n_l$  pre  $l = 1, 2, \dots, r$  a pre celý rozsah súboru  $n$  je

$$n = \sum_{l=1}^r n_l. \quad (3.10)$$

Hore popísaná tabuľka sa nazýva tabuľka rozdelenia početností a proces jej vytvárania je triedenie prvého stupňa resp. jednorozmerné triedenie. Pre  $k = 2$  vznikne kontingenčná tabuľka s dvomi výstupmi. Kontingenčná tabuľka pre znaky  $Y, Z$  vyzerá nasledovne:

Y	Z	$\Sigma$
	1 ... c	
1	$n_{11} \dots n_{1c}$	$n_{1.}$
...	...	...
r	$n_{r1} \dots n_{rc}$	$n_{r.}$
$\Sigma$	$n_{.1} \dots n_{.c}$	$n$

Tabuľka 3.1: Kontingenčná tabuľka

Udáva prehľad variánt prvého a druhého znaku. Riadky zodpovedajú možným hodnotám prvého znaku a slúpce zobrazujú hodnoty druhého znaku. V bunkách tabuľky nájdeme pozorované početnosti kombinácií hodnôt znakov  $Y$  a  $Z$ . Na krajoch sa nachádzajú pozorované marginálne početnosti definovanej ako riadkový a stĺpcový súčet. Celkový rozsah súboru je

$$n = \sum_{i=1}^r n_{i.} = \sum_{j=1}^c n_{.j} = \sum_{i=1}^r \sum_{j=1}^c n_{ij}. \quad (3.11)$$

Tento proces vytvárania nazývame triedenie druhého druhu alebo dvojrozmerné triedenie. Nezávislosť náhodných veličín  $Y$  a  $Z$  nastáva práve vtedy, keď

$$p_{ij} = p_{i.}p_{.j} \quad pre \quad i = 1, 2, \dots, r, \quad j = 1, 2, \dots, c \quad (3.12)$$

Hypotézu  $H_0$  nezávislosti  $Y$  a  $Z$  je možné testovať pomocou kontingenčnej tabuľky. Za platnosti  $H_0$  má veličina

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(n_{ij} - \frac{n_{i.}n_{.j}}{n})^2}{\frac{n_{i.}n_{.j}}{n}} \quad (3.13)$$

asymptoticky rozdelenie chí-kvadrát s  $(r-1)(c-1)$  stupňami voľnosti. Ak prekročí testová štatistika  $\chi^2$   $(1-\alpha)$ -kvantil rozdelenia chí-kvadrát s  $(r-1)(c-1)$  stupňami voľnosti, zamietame  $H_0$  na hladine  $\alpha$ . Aproximácia pomocou chí-kvadrát rozdelenia v praxi dáva rozumné výsledky, ak pre teoretické početnosti platí

$$\frac{n_{i.}n_{.j}}{n} > 5 \quad pre \quad i = 1, 2, \dots, r, \quad j = 1, 2, \dots, c. \quad (3.14)$$

Odvodenie testu nezávislosti alebo informácie o  $p$ -rozmernej kontingenčnej tabuľke nájdeme v [4], [9].

# Kapitola 4

## Mnohorozmerné štatistické metódy

Táto kapitola sa zaoberá dvoma prístupmi k spracovaniu mnohorozmerných dát. Prvým z nich bude zhluková analýza, ktorá triedi objekty do skupín s podobnými vlastnosťami z hľadiska sledovaných znakov. Problém zaradenia nového objektu do už vybraných skupín rieši diskriminačná analýza, ktorá je druhou metódou rozoberanou v nasledujúcom texte. Teoretické poznatky pochádzajú z [1], [4] a [8].

### 4.1 Zhuková analýza

Pri zhukovej analýze skúmame  $n$  objektov. Sú charakterizované  $k$  znakmi či už diskkrétnej alebo spojitej povahy. Cieľom je nájsť skupiny s podobnými objektami. Budeme predpokladať standardizované dáta (3.9), ktoré nám vytvoria dátovú maticu  $\mathbf{D}$  zavedenú v (3.8).

Zhuky podobných objektov budeme hľadať pomocou vzdialeností. Uvažujeme  $i$ -tý a  $j$ -tý objekt, ktorým prislúchajú jednotlivé riadky  $\mathbf{x}_i$  a  $\mathbf{x}_j$  matice  $\mathbf{D}$ . Dva základné postupy výpočtu vzdialenosti dvoch objektov sú :

- Euklidovská vzdialenosť

$$d_E(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{r=1}^k (x_{ir} - x_{jr})^2}. \quad (4.1)$$

- **Manhattanská vzdialenosť**

$$d_M(\mathbf{x}_i, \mathbf{x}_j) = \sum_{r=1}^k |x_{ir} - x_{jr}|. \quad (4.2)$$

V závislosti na type dát sa rozhodujeme, ktorú vzdialenosť použijeme. Ak sú merané znaky nezávislé a normálne rozdelené s rovnakými rozptylmi, tak použijeme euklidovskú vzdialenosť. Toto je doporučené i pre prípad vecne podobných znakov, ktoré sú rovnako dôležité pri klasifikácii objektov. Pri existencii odľahlých pozorovaní, v prípade porušenia podmienky normality rozdelenia vektora znakov je vhodnejšou voľbou Manhattanská vzdialenosť. Tú je možné zvoliť i v prípade diskretných celočíselných znakov.

Následujúce metódy popisujú vzájomnú vzdialenosť zhlukov  $S_1$  a  $S_2$  :

- **Metóda najbližšieho suseda**

$$D_{min}(S_1, S_2) = \min_{\mathbf{x}_i \in S_1, \mathbf{x}_j \in S_2} d(\mathbf{x}_i, \mathbf{x}_j), \quad (4.3)$$

- **Metóda najvzdialnejšieho suseda**

$$D_{max}(S_1, S_2) = \max_{\mathbf{x}_i \in S_1, \mathbf{x}_j \in S_2} d(\mathbf{x}_i, \mathbf{x}_j), \quad (4.4)$$

- **Metóda priemernej vzdialenosti**

$$D_{priem}(S_1, S_2) = d(\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2), \quad (4.5)$$

kde  $d$  je miera vzdialenosti objektov a

$$S_1 = \{\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_r}\},$$

$$S_2 = \{\mathbf{x}_{i_{r+1}}, \mathbf{x}_{i_{r+2}}, \dots, \mathbf{x}_{i_m}\} \quad \text{pre } 1 \leq r < m \leq n,$$

$$\bar{\mathbf{x}}_1 = \frac{1}{r} \sum_{j=1}^r \mathbf{x}_{i_j},$$

$$\bar{\mathbf{x}}_2 = \frac{1}{m-r} \sum_{j=r+1}^m \mathbf{x}_{i_j}.$$

Zhlukovacie algoritmy delíme na

- **Hierarchické**
- **Nehierarchické**

### 4.1.1 Hierarchická aglomeratívna metóda

Najprv popíšeme hierarchickým aglomeratívnym algoritmom. Pre vytváranie zhlukov je potrebné mať k dispozícii maticu vzdialeností  $\Delta$ , ktorej prvky sú  $d_{ij} = D(S_i, S_j)$ .

Aglomeratívny algoritmus má nasledujúci priebeh:

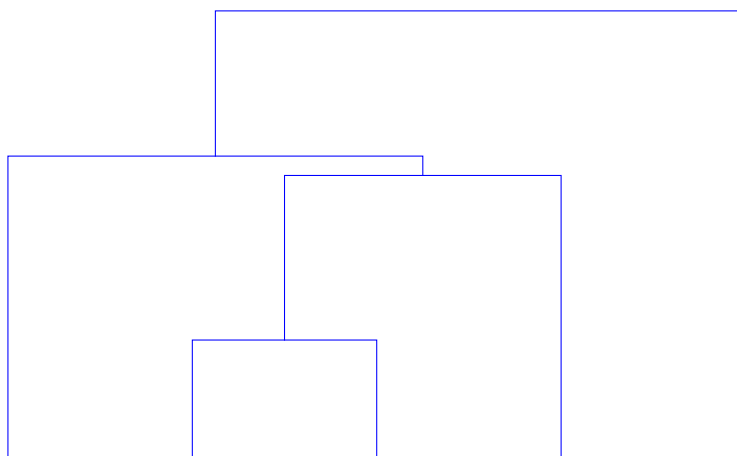
Počiatok: Každý objekt tvorí samostatný zhluk.

Krok 1 : Určíme maticu vzdialeností objektov  $\Delta^{(1)}$ , nájdeme  $d_{rs}^{(1)} = \min(d_{ij}^{(1)}; i, j = 1, 2, \dots, n)$  a spojíme objekty  $r, s$  do skupiny.

Krok  $t$  : Pre  $1 < t < n$  analogicky určíme maticu vzdialeností objektov a zhlukov  $\Delta^{(t)}$ , nájdeme  $d_{rs}^{(t)} = \min(d_{ij}^{(t)}; i, j = 1, 2, \dots, n - t + 1)$  a spojíme zhluky  $r, s$  do novej skupiny.

Koniec nastáva, ak sa všetky objekty nachádzajú v jednom zhluk. Algoritmus je možné ukončiť skôr.

Priebeh algoritmu je zachytený grafom nazvaným dendrogram.



Obr. 4.1: Jednoduchý príklad dendrogramu

Na obrázku 4.1 máme zobrazenú jednoduchú simuláciu hierarchického algoritmu vytvorenú pomocou softvéru Mathematica. Na začiatku pozorujeme 5 objektov (5 bodov na spodnej časti obrázku), ktoré sú postupom času na základe euklidovskej vzdialenosti spájané do zhlukov. Očíslujeme ich od 1 do 5 zľava doprava. Prvým spojením, respektíve najbližšie sú k sebe body 2 a 3. Najvzdialenejším bodom a posledným pridaným je 5. Pri pohľade na dendrogram môžeme rozhodnúť o skoršom ukončení algoritmu. Vzdialenosť určujú vertikálne čiary.



Vidíme, že piaty objekt je ďaleko od ostatných. Do zhľuku ho nezaradíme a algoritmus ukončíme po pridaní objektu 1 do zhľuku. Týmto spôsobom odlúčime netypický objekt 5 od skupiny ostatných podobných objektov. Zhľukovanie je však možné ukončiť aj pred pridaním objektu 4 do zhľuku (2,3), lebo vzdialenosť objektu 4 od zhľuku (2,3) je výrazne väčšia ako vzdialenosť objektov 2 a 3. Okamih ukončenia algoritmu je záležitosťou subjektívnej voľby analytika. Týmto dôjde k rozdeleniu dendrogramu na dve časti a to na realizované spojenia a neuskutočnené spojenia. Rozdelenie je určené priamkou vedenou v úrovni, ktorej súradnicu napríklad v programe NCSS nazývame cluster cutoff.

### 4.1.2 Nehierarchická metóda

Pre nehierarchické algoritmy máme k dispozícii predom určený počet skupín, do ktorých budeme objekty triediť. Každá skupina má svojho zástupcu. Nazýva sa centroid. Môže ním napríklad byť aritmetický priemer (3.2) vektorov reprezentujúcich jednotlivé objekty v skupine.

Jednou z používaných metód v rámci tohto prístupu je metóda K - means clustering, ktorá s vopred určeným počtom zhľukov vytvára optimálny rozklad objektov. Je založená na minimalizácii súčtov štvorcov  $K^{(q)}$  cez všetky možné usporiadania objektov do skupín. Ak predpokladáme usporiadanie do  $h$  zhľukov, potom je

$$K^{(q)} = \sum_{r=1}^h \sum_{i: x_i \in P_r^{(q)}} \sum_{j=1}^k (x_{ij} - c_{rj}^{(q)})^2, \quad (4.6)$$

kde  $q$  je číslo usporiadania objektov do zhľukov,

$P_r^{(q)}$  znamená  $r$  -tý zhľuk pri  $q$ -tom usporiadaní objektov,

$x_{ij}$  predstavuje hodnotu  $j$ -tého znaku v  $i$  -tom objekte,

$c_{rj}^{(q)}$  určuje hodnotu  $j$ -tého znaku v centroide  $r$ -tého zhľuku  $P_r^{(q)}$ .

## 4.2 Diskriminačná analýza

Diskriminačná analýza sa zaoberá zaradením nového objektu do už existujúcich skupín, ktoré boli vytvorené napríklad na základe podobných hodnôt znakov zhľukovou analýzou. Znova predpokladáme dátovú maticu definovanú

predpisom (3.8) a  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ik})^T$ ,  $i = 1, 2, \dots, n$  je  $i$ -tý objekt s nameranými  $k$  znakmi. Objekty sú potom zaradené do tried  $T_1, T_2, \dots, T_p$ ,  $p < n$  a tým je vytvorená testovacia databáza. Pre nový objekt  $\mathbf{x}_i$ ,  $i > n$  vytvoríme rozhodovacie pravidlo, diskriminačný skór, na základe ktorého tento nový objekt zaradíme do niektorej triedy. Pre skupiny  $T_1, T_2, \dots, T_p$  definujeme diskriminačné skóry  $D_1^{(i)}, D_2^{(i)}, \dots, D_p^{(i)}$  ako funkcie znakov nameraných na  $i$ -tom objekte, ktoré minimalizujú strednú hodnotu straty pri zaradení objektu do nesprávnej triedy. Stratu označíme  $z_{rs}$ , ak objekt prislúcha  $r$ -tej skupine a zaradíme ho do  $s$ -tej skupiny. Prepokladajme spojité znaky. Za podmienky, že objekt patrí do  $r$ -tej skupiny, je stredná hodnota straty

$$L_r = z_{r1} \int_{A_1} f_r(\mathbf{x}) d\mathbf{x} + \dots + z_{rp} \int_{A_p} f_r(\mathbf{x}) d\mathbf{x} = \sum_{s=1}^p z_{rs} \int_{A_s} f_r(\mathbf{x}) d\mathbf{x}, \quad (4.7)$$

kde  $f_r$  je hustota vektora znakov  $\mathbf{x}$  v  $r$ -tej triede a  $A_d$  pre  $d = 1, 2, \dots, p$  je rozklad  $k$ -rozmerného reálneho priestoru  $R^k$  do disjunktných borelovských množín. Predpokladajme, že  $\pi_r > 0$  je pravdepodobnosť, že objekt patrí do  $r$ -tej triedy. Potom nepodmienená stredná hodnota straty je

$$L = \pi_1 L_1 + \dots + \pi_p L_p. \quad (4.8)$$

Ak definujeme

$$q_s(\mathbf{x}) = \sum_{r=1}^p \pi_r z_{rs} f_r(\mathbf{x}), \quad s = 1, \dots, p, \quad (4.9)$$

potom (4.8) je možné upraviť pomocou (4.7) a (4.9) na

$$\begin{aligned} L &= \pi_1 L_1 + \dots + \pi_p L_p \\ &= \pi_1 \sum_{s=1}^p z_{1s} \int_{A_s} f_1(\mathbf{x}) d\mathbf{x} + \dots + \pi_p \sum_{s=1}^p z_{ps} \int_{A_s} f_p(\mathbf{x}) d\mathbf{x} \\ &= \sum_{r=1}^p \pi_r \sum_{s=1}^p z_{rs} \int_{A_s} f_r(\mathbf{x}) d\mathbf{x} \\ &= \sum_{s=1}^p \sum_{r=1}^p \int_{A_s} \pi_r z_{rs} f_r(\mathbf{x}) d\mathbf{x} \\ &= \sum_{s=1}^p \int_{A_s} \sum_{r=1}^p \pi_r z_{rs} f_r(\mathbf{x}) d\mathbf{x} \\ &= \sum_{s=1}^p \int_{A_s} q_s(\mathbf{x}) d\mathbf{x}. \end{aligned} \quad (4.10)$$

Hľadáme optimálny rozklad priestoru  $R^k$ , aby bolo  $L$  minimálne. Nech  $A_1^*, \dots, A_p^*$  je rozklad priestoru  $R^k$ , pre ktorý platí

$$\mathbf{x} \in A_t^* \Rightarrow q_t(\mathbf{x}) \leq q_s(\mathbf{x}), s = 1, \dots, p. \quad (4.11)$$

Ak označíme

$$L^* = \sum_{s=1}^p \int_{A_s^*} q_s(\mathbf{x}) d\mathbf{x}, \quad (4.12)$$

potom platí

$$L \geq L^*. \quad (4.13)$$

Teraz sa bližšie pozrieme na túto nerovnosť:

$$\begin{aligned} L &= \sum_{s=1}^p \int_{A_s} q_s(\mathbf{x}) d\mathbf{x} = \sum_{s=1}^p \sum_{t=1}^p \int_{A_s \cap A_t^*} q_s(\mathbf{x}) d\mathbf{x} = \sum_{t=1}^p \sum_{s=1}^p \int_{A_t^* \cap A_s} q_s(\mathbf{x}) d\mathbf{x} \geq \\ &\geq \sum_{t=1}^p \sum_{s=1}^p \int_{A_t^* \cap A_s} q_t(\mathbf{x}) d\mathbf{x} = \sum_{t=1}^p \int_{A_t^*} q_t(\mathbf{x}) d\mathbf{x} = L^*. \end{aligned} \quad (4.14)$$

Rovnosti v (4.14) sú zrejmé a nerovnosť vyplýva z (4.11).

Špeciálne volíme  $z_{ss} = 0$  a  $z_{rs} = 1$  pre  $r \neq s$ . Ak označíme  $c(\mathbf{x}) = \sum_{r=1}^p \pi_r f_r(\mathbf{x})$ , dostávame  $q_s(\mathbf{x}) = c(\mathbf{x}) - \pi_s f_s(\mathbf{x})$ . Minimalizáciou podielu chybné zaradených objektov dostaneme minimálne  $L$ . Pri danom  $\mathbf{x}$  a  $t$  platí  $q_t(\mathbf{x}) \leq q_s(\mathbf{x}) \Leftrightarrow \pi_t f_t(\mathbf{x}) \geq \pi_s f_s(\mathbf{x})$  pre  $s = 1, \dots, p$ . Ak nájdeme taký vektor, pre ktorý platí  $\pi_t f_t(\mathbf{x}) > \pi_s f_s(\mathbf{x})$ ,  $s \neq t$ , potom tento objekt reprezentovaný vektorom  $\mathbf{x}$  zaradíme do skupiny  $t$ . Ak nastáva rovnosť, môžeme vybrať akýkoľvek z maximalizujúcich indexov.

Teraz dosadíme za  $f_s(\mathbf{x})$  hustotu  $k$ -rozmerného normálneho rozdelenia (3.1) so strednou hodnotou  $\boldsymbol{\mu}_s$  a regulárnou variančnou maticou  $\boldsymbol{\Sigma}_s$  pre  $s = 1, \dots, p$ . Nerovnosť  $\pi_t f_t(\mathbf{x}) > \pi_s f_s(\mathbf{x})$ ,  $s \neq t$  je ekvivalentná s nerovnosťou

$$\ln \pi_t + \ln f_t > \ln \pi_s + \ln f_s, \quad s \neq t. \quad (4.15)$$

Ak označíme

$$D_s = -\frac{1}{2} \ln |\boldsymbol{\Sigma}_s| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_s)^T \boldsymbol{\Sigma}_s^{-1} (\mathbf{x} - \boldsymbol{\mu}_s) + \ln \pi_s, \quad (4.16)$$

tak nerovnosť (4.15) platí práve vtedy keď  $D_t > D_s$  pre  $s \neq t$ . Nový objekt  $x_{n+1}$  zaradíme do triedy  $T_t$  ak

$$D_t^{(n+1)} = \max(D_s^{(n+1)}, s = 1, 2, \dots, p) \quad (4.17)$$

Za odhad  $\hat{\pi}_s$  pravdepodobnosti  $\pi_s$  je možné vziať relatívnu početnosť  $\frac{n_s}{n}$ , kde  $n_s$  je počet objektov testovacej databáze v triede  $T_s$  a  $n_1 + n_2 + \dots + n_p = n$ . Ďalšou možnou voľbou je  $\hat{\pi}_s = \frac{1}{p}$  pre  $s = 1, 2, \dots, p$ . Priemer meraných znakov v triede  $T_s$  je

$$\bar{\mathbf{x}}_s = \frac{1}{n_s} \sum_{k: \mathbf{x}_k \in T_s} \mathbf{x}_k. \quad (4.18)$$

Pre variančnú maticu znakov v  $s$ -tej triede je odhadom výberová variančná matica

$$\mathbf{S}_s = \frac{1}{n_s - 1} \sum_{k: \mathbf{x}_k \in T_s} (\mathbf{x}_k - \bar{\mathbf{x}}_s)(\mathbf{x}_k - \bar{\mathbf{x}}_s)^T. \quad (4.19)$$

Predpokladajme normálne rozdelenie vektorov znakov, potom môžeme hovoriť o dvoch typoch diskriminačných skórov:

- **Kvadratický diskriminačný skór** počítaný pre objekt  $\mathbf{x}_i$  a triedu  $T_s$ , ktorý vznikne dosadením odhadov (4.18) a (4.19) do (4.16)

$$D_s^{(i)} = -\frac{1}{2} \ln |\mathbf{S}_s| - \frac{1}{2} (\mathbf{x}_i - \bar{\mathbf{x}}_s)^T \mathbf{S}_s^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_s) + \ln \hat{\pi}_s. \quad (4.20)$$

Kvadratický diskriminačný skór určuje obecný prístup. Ak sa triedy odlišujú len strednou hodnotou a majú rovnakú variančnú maticu, používa sa

- **Lineárny diskriminačný skór**

$$D_s^{(i)} = \bar{\mathbf{x}}_s^T \mathbf{S}^{-1} \mathbf{x}_i - \frac{1}{2} \bar{\mathbf{x}}_s^T \mathbf{S}^{-1} \bar{\mathbf{x}}_s + \ln \hat{\pi}_s, \quad (4.21)$$

pre  $s$ -tu triedu a  $i$ -tý objekt, kde

$$\mathbf{S} = \frac{1}{n-p} \sum_{j=1}^p (n_j - 1) \mathbf{S}_j.$$

V knihe [4] sa lineárna diskriminačná analýza nedoporučuje pre prípad viac ako dvoch tried objektov.

# Kapitola 5

## Spracované dáta

V praktickej časti budeme demonštrovať použitie popísaných metód na voľne dostupných dátach istej nemeckej banky [10], pričom konkrétne budeme pracovať s databázou žiadateľov o úver v počte 1000. Celkovo táto databáza tvorí dátovú maticu typu (3.8). Objektami sú klienti banky s jednotlivými priradenými znakmi, ktoré sú diskretného alebo spojitého charakteru. V nasledujúcich podkapitolách budú jednotlivé znaky popísané.

### 5.1 Kontingenčná tabuľka

Najprv sa budeme zaoberať testom nezávislosti v kontingenčných tabuľkách (tabuľka 3.1) pre tri diskretné znaky. Vybranými znakmi sú bilancia na účte, informácia o pôvode a vlastníctvo. Pod jednotlivými znakmi sa nachádzajú kategórie, do ktorých sú klienti rozdelení:

- *Bilancia na účte*
  - Žiadna bilancia alebo dlh
  - Kladná bilancia
  - Žiaden účet
- *Cudzinec*
  - Áno
  - Nie

- *Vlastníctvo*
  - Dom alebo pozemok
  - Životná poisťka
  - Auto
  - Iné

Nasleduje prehľad kontingenčných tabuliek pozorovaných početností  $n_{ij}$ , teoretických početností  $\frac{n_{i.}n_{.j}}{n}$  a hodnôt jednotlivých sčítancov veličín  $\chi^2$  zo vzorca (3.13).

		Vlastníctvo				
Cudzinec		Dom alebo pozemok	Životná poisťka	Auto	Nič	Celkom
Áno		19	14	2	2	<b>37</b>
Nie		263	218	330	152	<b>963</b>
Celkom		<b>282</b>	<b>232</b>	<b>332</b>	<b>154</b>	

Tabuľka 5.1: Kontingenčná tabuľka pozorovaných početností znakov Cudzinec a Vlastníctvo

		Vlastníctvo				
Cudzinec		Dom alebo pozemok	Životná poisťka	Auto	Iné	Celkom
Áno		10,43	8,58	12,28	5,70	37
Nie		271,57	223,42	319,72	148,30	963
Celkom		282	232	332	154	

Tabuľka 5.2: Teoretické početnosti znakov Cudzinec a Vlastníctvo

		Vlastníctvo				
Cudzinec		Dom alebo pozemok	Životná poisťka	Auto	Nič	Celkom
Áno		7,03	3,42	8,61	2,40	<b>21,46</b>
Nie		0,27	0,13	0,33	0,09	<b>0,82</b>
Celkom		<b>7,30</b>	<b>3,55</b>	<b>8,94</b>	<b>2,49</b>	<b>22,28</b>

Tabuľka 5.3: Hodnoty sčítancov  $\chi^2$  rozdelenia znakov Cudzinec a Vlastníctvo

<b>Bilancia na účte</b>				
<b>Vlastníctvo</b>	<b>Žiaden dlh alebo zostatok</b>	<b>Kladný zostatok</b>	<b>Žiaden účet</b>	<b>Celkom</b>
<b>Dom alebo pozemok</b>	70	75	137	<b>282</b>
<b>Životná poisťka</b>	80	55	97	<b>232</b>
<b>Auto</b>	74	92	166	<b>332</b>
<b>Nič</b>	50	47	57	<b>154</b>
<b>Celkom</b>	<b>274</b>	<b>269</b>	<b>457</b>	

Tabuľka 5.4: Kontingenčná tabuľka pozorovaných početností znakov Bilancia na účte a Vlastníctvo

<b>Bilancia na účte</b>				
<b>Vlastníctvo</b>	<b>Žiaden dlh alebo zostatok</b>	<b>Kladný zostatok</b>	<b>Žiaden účet</b>	<b>Celkom</b>
<b>Dom alebo pozemok</b>	77,27	75,86	128,87	<b>282</b>
<b>Životná poisťka</b>	63,57	62,41	106,02	<b>232</b>
<b>Auto</b>	90,97	89,31	151,72	<b>332</b>
<b>Nič</b>	42,20	41,43	70,38	<b>154</b>
<b>Celkom</b>	<b>274</b>	<b>269</b>	<b>457</b>	

Tabuľka 5.5: Teoretické početnosti znakov Bilancia na účte a Vlastníctvo

<b>Bilancia na účte</b>				
<b>Vlastníctvo</b>	<b>Žiaden dlh alebo zostatok</b>	<b>Kladný zostatok</b>	<b>Žiaden účet</b>	<b>Celkom</b>
<b>Dom alebo pozemok</b>	0,68	0,01	0,51	<b>1,21</b>
<b>Životná poisťka</b>	4,25	0,88	0,77	<b>5,89</b>
<b>Auto</b>	3,16	0,08	1,34	<b>4,59</b>
<b>Nič</b>	1,44	0,75	2,54	<b>4,74</b>
<b>Celkom</b>	<b>9,54</b>	<b>1,72</b>	<b>5,17</b>	<b>16,43</b>

Tabuľka 5.6: Hodnoty sčítancov  $\chi^2$  rozdelenia znakov Bilancia na účte a Vlastníctvo

<b>Bilancia na účte</b>				
<b>Cudzinec</b>	<b>Žiaden dlh alebo zostatok</b>	<b>Kladný zostatok</b>	<b>Žiaden účet</b>	<b>Celkom</b>
<b>Áno</b>	15	6	16	<b>37</b>
<b>Nie</b>	259	263	441	<b>963</b>
<b>Celkom</b>	<b>274</b>	<b>269</b>	<b>457</b>	

Tabuľka 5.7: Kontingenčná tabuľka pozorovaných početností znakov Bilancia na účte a Cudzinec

<b>Bilancia na účte</b>				
<b>Cudzinec</b>	<b>Žiaden dlh alebo zostatok</b>	<b>Kladný zostatok</b>	<b>Žiaden účet</b>	<b>Celkom</b>
<b>Áno</b>	10,14	9,95	16,91	<b>37</b>
<b>Nie</b>	263,86	259,05	440,09	<b>963</b>
<b>Celkom</b>	<b>274</b>	<b>269</b>	<b>457</b>	

Tabuľka 5.8: Teoretické početnosti znakov Bilancia na účte a Cudzinec

<b>Bilancia na účte</b>				
<b>Cudzinec</b>	<b>Žiaden dlh alebo zostatok</b>	<b>Kladný zostatok</b>	<b>Žiaden účet</b>	<b>Celkom</b>
<b>Áno</b>	2,33	1,57	0,05	<b>3,95</b>
<b>Nie</b>	0,09	0,06	0,00	<b>0,15</b>
<b>Celkom</b>	<b>2,42</b>	<b>1,63</b>	<b>0,05</b>	<b>4,10</b>

Tabuľka 5.9: Hodnoty sčítancov  $\chi^2$  rozdelenia znakov Bilancia na účte a Cudzinec

V tabuľkách 5.3, 5.6 a 5.9 nájdeme sčítance  $\chi^2$  rozdelenia daných znakov a v pravom dolnom rohu hodnotu testovej štatistiky  $\chi^2$  zo vzorca (3.13). Hladina testu je 0.05. Pri dvojici znakov Cudzinec a Vlastníctvo porovnáваме hodnotu  $\chi^2 = 22,28$  s 0,95-kvantilom  $\chi^2$  rozdelenia s tromi stupňami voľnosti, ktorá je  $\chi^2 = 7,82$ . Testová štatistika prekračuje kvantil, teda zamietame nezávislosť znakov. Rovnako je to pri druhej dvojici znakov Vlastníctvo a Bilancia na účte. Testovú štatistiku  $\chi^2 = 16,43$  porovnáваме s 0,95-kvantilom  $\chi^2$  rozdelenia o šesť stupňov voľnosti 12,59. V poslednom prípade k prekročeniu 0,95-kvantilu nedôjde. Testová štatistika  $\chi^2 = 4,1$  je menšia ako 0,95-kvantil  $\chi^2$  rozdelenia s dvoma stupňami voľnosti 5,99. Takže na základe našich dát nemôžeme zamietnuť nezávislosť znakov Cudzinec a Bilancia na účte.

Závislosť spojitých znakov bude vyšetrená na základe korelačnej matice (3.6). Spojité znaky, ktoré budeme posudzovať, sú :

- *Vek*
- *Výška úveru*
- *Doba splatnosti úveru*

Doba splatnosti úveru je v intervale od 4 mesiacov do 72 mesiacov. Výška úveru sa pohybuje v rozmedzí 250 - 18 424 nemeckých mariek a vek klientov je od 19 do 75 rokov.



Korelačná matica spojitých znakov je

$$\begin{pmatrix} 1. & -0.04 & 0.03 \\ -0.04 & 1. & 0.62 \\ 0.03 & 0.62 & 1. \end{pmatrix}, \quad (5.1)$$

kde stĺpce, respektíve riadky reprezentujú spojité znaky v poradí Vek, Výška úveru a Doba splatnosti úveru. Najsilnejšia lineárna závislosť je medzi znakom Výška úveru a Doba splatnosti úveru, čo odráža úverovú politiku bánk. Korelácia medzi znakmi Vek a Doba splatnosti úveru a aj Vek a Výška úveru je veľmi slabá.

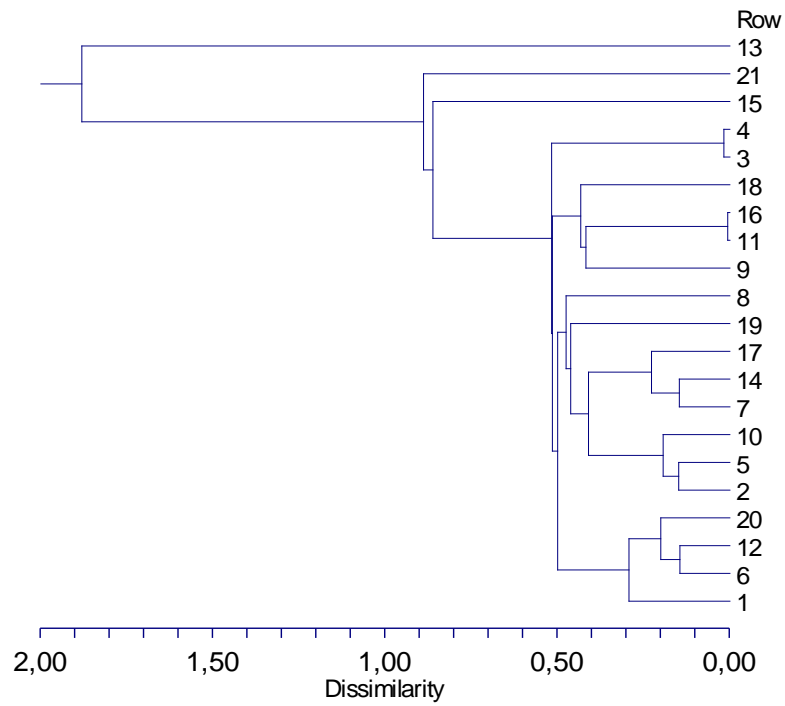
## 5.2 Zhluková analýza

V tejto časti budeme aplikovať zhlukovú analýzu popisovanú v kapitole 4.1 na tri spojité znaky. Budeme používať rovnaké spojité znaky ako v kapitole 5.1. Vzhľadom k lepšej názornosti bude dátová matica obsahovať obmedzený počet klientov. V tejto kapitole aplikujeme na dáta hierarchickú aglomeratívnu metódu a nehierarchickú zhlukovú metódu (K-means clustering). Použitými softvérmi budú Mathematica, NCSS a MS Excel.

### 5.2.1 Hierarchická aglomeratívna metóda

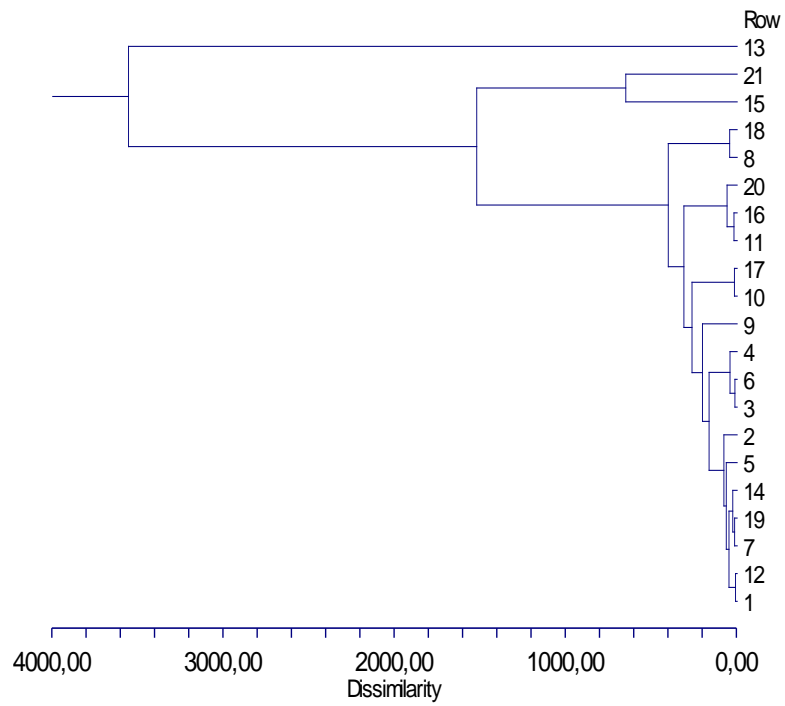
Sledovanými znakmi sú Vek, Doba splatnosti úveru a Výška úveru. Budeme sa zaoberať klientami, ktorí sú viac ako 60-roční a menej ako 64-roční. Počet klientov sme znížili na 21. Celkovo máme dátovú maticu s tromi znakmi a dvadsaťjeden objektami. Zhluková analýza bude aplikovaná na štandardizovanú a neštandardizovanú dátovú maticu s euklidovskou respektíve manhattanskou vzdialenosťou. Ako metódu na určenie vzájomnej vzdialenosti zhlukov sme zvolili metódu najbližšieho suseda. Vytvorením dendrogramov v softvéri NCSS bolo zistené, že voľba vzdialenosti objektov nemá vplyv na tvar dendrogramu. Uvedieme dendrogramy s euklidovskou vzdialenosťou na štandardizované a neštandardizované dáta:

### Dendrogram



Obr. 5.1: Dendrogram štandardizovaných dát

### Dendrogram



Obr. 5.2: Dendrogram neštandardizovaných dát

Čísla v pravej časti dendrogramov predstavujú riadky v dátovej matici.

Z dendrogramu usúdime, že algoritmus je vhodné ukončiť pri štandardizovaných dátach napríklad voľbou cluster cutoff = 0.6 a pri neštandardizovaných dátach je možné použiť cluster cutoff = 1000. Ak vztýčíme v bodoch cluster cutoff na vodorovnej ose zvislú priamku, rozdelíme dendrogram tak, že spojenie naľavo od tejto kolmice nebude realizované. Tento fakt neskôr podložíme číselnými hodnotami. Konkrétne hodnoty datovej matice sú:

Klient	Doba splatnosti úveru	Výška úveru	Vek
1	18	1239	61
2	12	1655	63
3	24	2032	60
4	24	1940	60
5	15	1520	63
6	12	2012	61
7	10	1364	64
8	9	3832	64
9	24	2384	64
10	10	781	63
11	24	2924	63
12	12	1255	61
<b>13</b>	<b>60</b>	<b>13756</b>	<b>63</b>
14	13	1409	64
<b>15</b>	<b>20</b>	<b>6468</b>	<b>60</b>
16	24	2957	63
17	6	753	64
18	24	3757	62
19	6	1338	62
20	12	3059	61
<b>21</b>	<b>30</b>	<b>7596</b>	<b>63</b>

Tabuľka 5.10: Použitá dátová matica

Aplikácia algoritmu na štandardizované dáta je naznačená v tabuľke 5.11.

Krok	Vzdialenosť	Objekty
1	0,0063	<b>11,16</b>
2	0,0176	<b>3,4</b>
3	0,1450	<b>6,12</b>
4	0,1467	<b>7,14</b>
...	...	...
17	0,5174	1,6,12,20,2,5,10,7,14,17,19,8,9,11,16,18, <b>3,4</b>
18	0,8620	1,6,12,20,2,5,10,7,14,17,19,8,9,11,16,18,3,4, <b>15</b>
19	0,8893	1,6,12,20,2,5,10,7,14,17,19,8,9,11,16,18,3,4,15, <b>21</b>
20	1,8810	1,6,12,20,2,5,10,7,14,17,19,8,9,11,16,18,3,4,15, 21, <b>13</b>

Tabuľka 5.11: Vzďialenosti zhluokov pre štandardizované dáta

Najbližšie sú k sebe objekty 11 a 16 so vzdialenosťou 0,0063. Je to prvé spojenie dvoch objektov do skupiny. Pri pohľade na konkrétne znaky vidíme, že títo klienti majú rovnakú dobu splatnosti úveru 24 a vek 63. Malý rozdiel pozorujeme vo výške úveru, kde v prvom prípade je to 2 924 nemeckých mariek a v druhom 2 957 nemeckých mariek. Druhá a tretia v poradí najmenšia vzdialenosť sa týka objektov 3 a 4 respektíve 6 a 12. U týchto klientov je zhoda v znaku Vek a Doba splatnosti úveru, no majú vyšší rozdiel vo výške úveru. Zmena nastáva v štvrtej spojenej dvojici v poradí, kde sa u klientov 7 a 14 okrem výšky úveru líši aj doba splatnosti. Celkovo však považujeme objekty týchto spojení za podobné.

Pri náhľade na tabuľku 5.11 vidíme, že posledné spojenie je so vzdialenosťou 1,8810. Došlo k pripojeniu klienta 13. Pozorujeme značné vychýlenie v hodnotách znakov Doba splatnosti úveru 60 mesiacov a Výška úveru 13 756 nemeckých mariek oproti všetkým ostatným klientom v tabuľke v tabuľke 5.10. Predchádzajúce dve spojenia majú podobnú vzdialenosť 0,8893 a 0,8620. Tu boli pripojení klienti 21 a 15. U klienta 21 stále pozorujeme značný rozdiel v hodnotách znakov Doba splatnosti a Výška úveru, no pri klientovi 15 sa rozdiel v znaku Doba splatnosti vytratí. Avšak rozdiel v znaku Výška úveru je stále značný. Následne by sme už videli zmenšujúcu sa vzdialenosť, ktorá najviac závisí na výške úveru.

Algoritmus ukončíme po sedemnástom kroku. Vzdialenosť spojovaných zhluokov je vtedy 0,5175. Ak by sme pridali ešte osemnásty krok, tak by sa vzdialenosť zmenila na 0,8620 čo predstavuje veľký skok. Prvých sedemnást krokov teda tvorí zhluok a posledné tri objekty ostanú nezaraďené.

V prípade neštandardizovaných dát bude tabuľka vzdialeností vyzerať takto:

Krok	Vzdialenosť	Objekty
1	9,87	<b>1,12</b>
2	13,48	<b>3,6</b>
3	15,23	<b>7,19</b>
4	16,34	<b>10,17</b>
...	...	...
17	403,05	1,12,7,19,14,5,2,3,6,4,9,10,17,11,16,20, <b>8,18</b>
18	651,28	<b>15,21</b>
19	1521,91	1,12,7,19,14,5,2,3,6,4,9,10,17,11,16,20,8,18, <b>15,21</b>
20	3556,52	1,12,7,19,14,5,2,3,6,4,9,10,17,11,16,20,8,18,15,21, <b>13</b>

Tabuľka 5.12: Vzďialenosťi zhľukov pre neštandardizované dáta

Na rozľišenie odľahľých objektov nemala vplyv štandardizácia dát. V oboch prípadoch by sme po sedemnástich krokoch algoritmu nezaradili rovnakých troch klientov a to 13, 21 a 15. Vyznačujú sa netypicky veľkou výškou úveru.

### 5.2.2 Nehierarchická K-means metóda

Túto metódu budeme znova aplikovať na rovnaké tri spojité znaky. Budú nás zaujímať klienti, ktorí majú viac ako 60 rokov. Ďalej je potrebné predom vybrať počet zhľukov. Hierarchická metóda nám rozdelila dáta na dva zhľuky, preto budeme predpokladať dva zhľuky. Ukážeme si, čo sa zmení, ak zvýšime počet zhľukov na štyri a päť. Dátová matica je väčšia než sme použili v hierarchickej analýze. Obsahuje 50 objektov a 3 znaky. V programe NCSS bola využitá funkcia K-means Clustering popísaná v kapitole 4.1.2.

Softvér rozdelil 50 objektov na 47 a 3 (obrázok 5.3), ak predpokladáme dva zhľuky. Rozptyl v zhľuku pozostávajúcich z troch objektov je v znaku Doba splatnosti úveru 0.

Znak	1. zhľuk	2. zhľuk
Doba splatnosti úveru	8,46	0,00
Výška úveru	2 583,41	4 321,99
Vek	4,36	1,73

Tabuľka 5.13: Smerodajné odchýľky v zhľukoch pre dané znaky

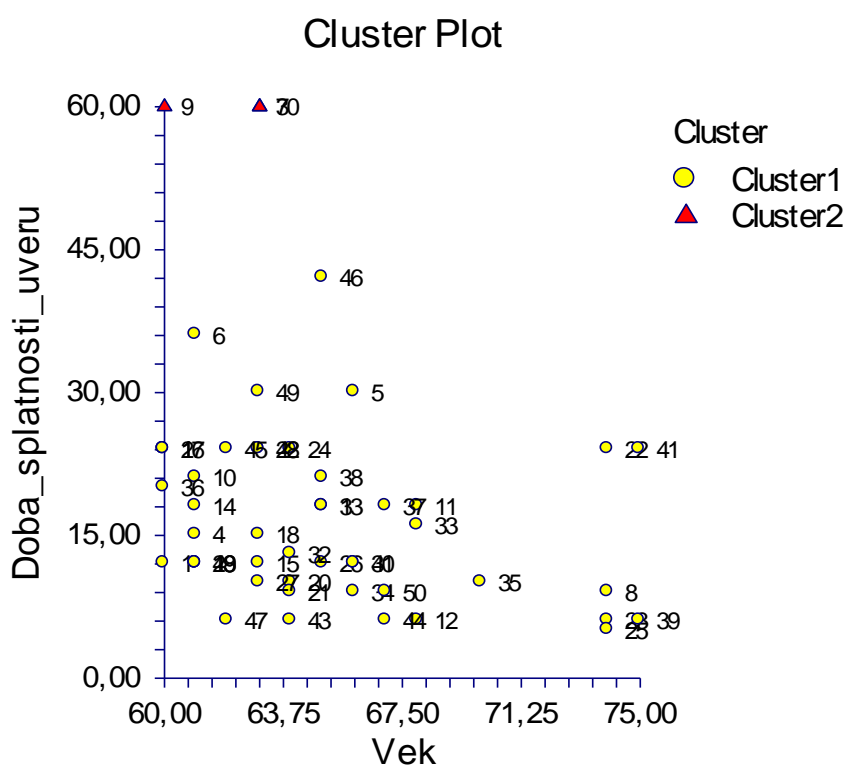
Je zjavné, že tieto tri objekty majú rovnakú Dobu splatnosti úveru. V priemeroch v tabuľke 5.14 je viditeľný veľký rozdiel medzi zhľukmi v Dobe splatnosti

úveru a Výške úveru.

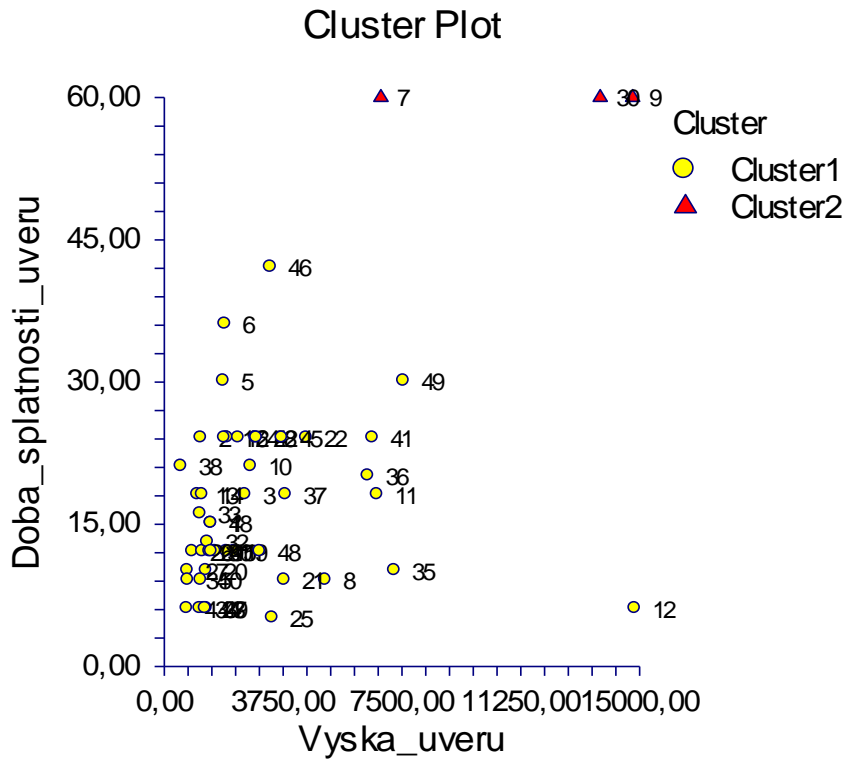
Znak	1. zhluk	2. zhluk
Doba splatnosti úveru	16,34	60,00
Výška úveru	2 830,21	11 791,33
Vek	65,09	62,00

Tabuľka 5.14: Priemery v zhluchoch pre dané znaky

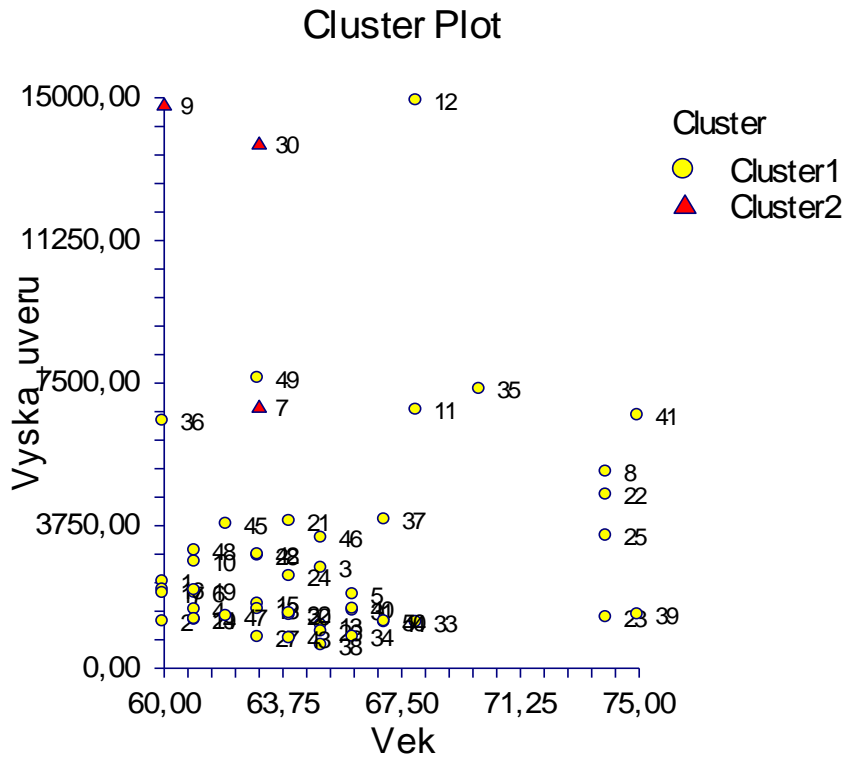
K interpretácii vytvorených zhlučkov môžu pomôcť dvojrozmerné projekcie v nasledujúcich obrázkoch:



Obr. 5.3: Zobrazenie zhlučkov pre znaky Doba splatnosti úveru a Vek



Obr. 5.4: Zobrazenie zhhlukov znaky Doba splatnosti úveru a Výška úveru



Obr. 5.5: Zobrazenie zhhlukov pre znaky Výška úveru a Vek

Na predchádzajúcich grafoch sme mohli vidieť jednotlivé zhluky farebne rozlíšené. Aj nehierarchická metóda nám odčlenila tri objekty, ktoré sa značne odlišujú od ostatných. Sú to klienti 7, 9 a 30. Všetci traja majú netypicky dlhú dobu splatnosti úveru 60 mesiacov, nízky vek a klienti 9 a 30 majú vysokú výšku úveru. Jeden z týchto objektov sa nám už objavil v hierarchickej aglomeratívnej metóde, v tabuľke 5.15 má poradie 30 (v tabuľke 5.11 má číslo 13) a ďalšie dva sú novo pridané do dátovej matice pre nehierarchické zhukovanie. V nižšie uvedenej tabuľke sú zvýraznené. Prehľad vybraných objektov:

Klient	Doba splatnosti úveru	Výška úveru	Vek
4	15	1512	61
5	30	1908	66
6	36	1953	61
<b>7</b>	<b>60</b>	<b>6836</b>	<b>63</b>
<b>9</b>	<b>60</b>	<b>14782</b>	<b>60</b>
12	6	14896	68
24	24	2384	64
<b>30</b>	<b>60</b>	<b>13756</b>	<b>63</b>
33	16	1175	68
35	10	7308	70
38	21	571	65
39	6	1374	75
41	24	6615	75
43	6	753	64
44	6	1169	67
46	42	3394	65
47	6	1338	62
49	30	7596	63

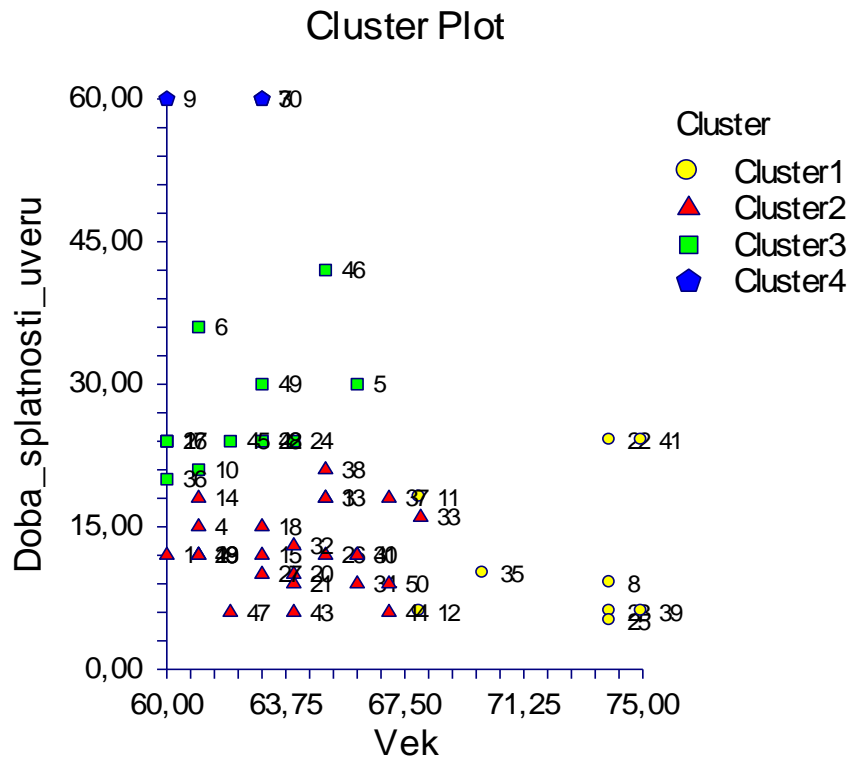
Tabuľka 5.15: Hodnoty znakov vybraných objektov

Tučne zvýraznené sú objekty, ktoré sa nachádzajú v zhluku 2. Na obrázkoch 5.3, 5.4 a 5.5 sú vyznačené červeným trojuholníkom.

Teraz predpokladajme štyri zhluky. Vyberieme jednu dvojicu znakov, na ktorej demoštrujeme, čo sa stane. Vo zvyšných prípadoch je situácia podobná. Voľbou vyššieho počtu zhlukov docielime väčšiu podobnosť znakov v rámci jedného zhuku. Na obrázku 5.6 môžeme vidieť postupné delenie skupiny žltej farby z obrázka 5.3. Pozorujeme rozdelenie do zhlukov, kde je možné charakterizovať každú skupinu. Modrou farbou je, ako sme už spomínali hore, charakteristická dlhá doba splatnosti a nízky vek v zhluke. Naopak pre žltý zhluk pozorujeme vysoký vek, no

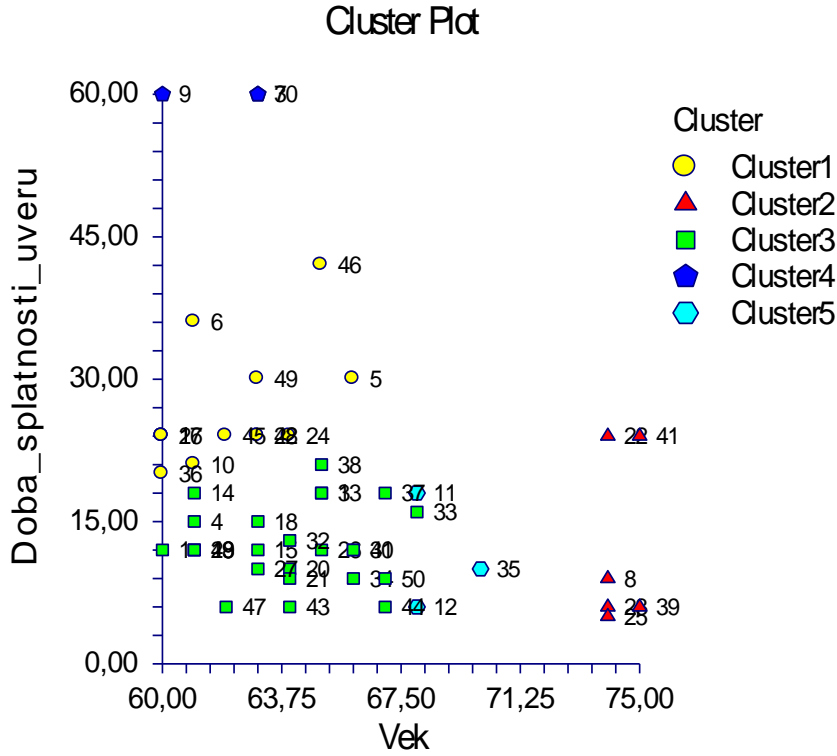


krátku dobu splatnosti. Červenou farbou je zobrazený nízky vek a krátka doba splatnosti úveru a zelenou nízky vek a stredne dlhá doba splatnosti.



Obr. 5.6: Zobrazenie zhhlukov pre znaky Doba splatnosti úroku a Vek

Ak zvýšime počet zhhlukov na päť, potom pozorujeme na obrázku 5.7 rovnaké rozdelenie objektov do zhhlukov až na žltý zhhluk na obrázku 5.6. Tento zhhluk je rozdelený na dva vzhľadom k veku klientov. Sú označený červenou a svetlo modrou farbou.



Obr. 5.7: Zobrazenie zhhlukov pre znaky Doba splatnosti úroku a Vek

### 5.3 Diskriminačná analýza

V poslednej kapitole pridáme podľa teoretického popisu v kapitole 4.2 ďalšieho klienta. Zaradíme ho do jednej zo skupín, ktorá vznikla použitou zhlukovou nehierarchickou metódou. Uvedieme diskriminačnú analýzu pre dve a štyri počiatočné skupiny.

Najprv predpokladajme dve počiatočné skupiny  $T_1$  (47 objektov) a  $T_2$  (3 objekty), ktoré vytvoria testovaciu databázu. Pre správne zaradenie ďalšieho objektu je potrebné napočítať diskriminačné skóre. Keďže predpokladáme dve počiatočné skupiny, budeme počítať lineárne diskriminačné skóre (4.21) :

$$D_s^{(i)} = \bar{\mathbf{x}}_s^T \mathbf{S}^{-1} \mathbf{x}_i - \frac{1}{2} \bar{\mathbf{x}}_s^T \mathbf{S}^{-1} \bar{\mathbf{x}}_s + \ln \hat{\pi}_s \quad s = 1, 2 \quad i = 1, \dots, 51, \quad (5.2)$$

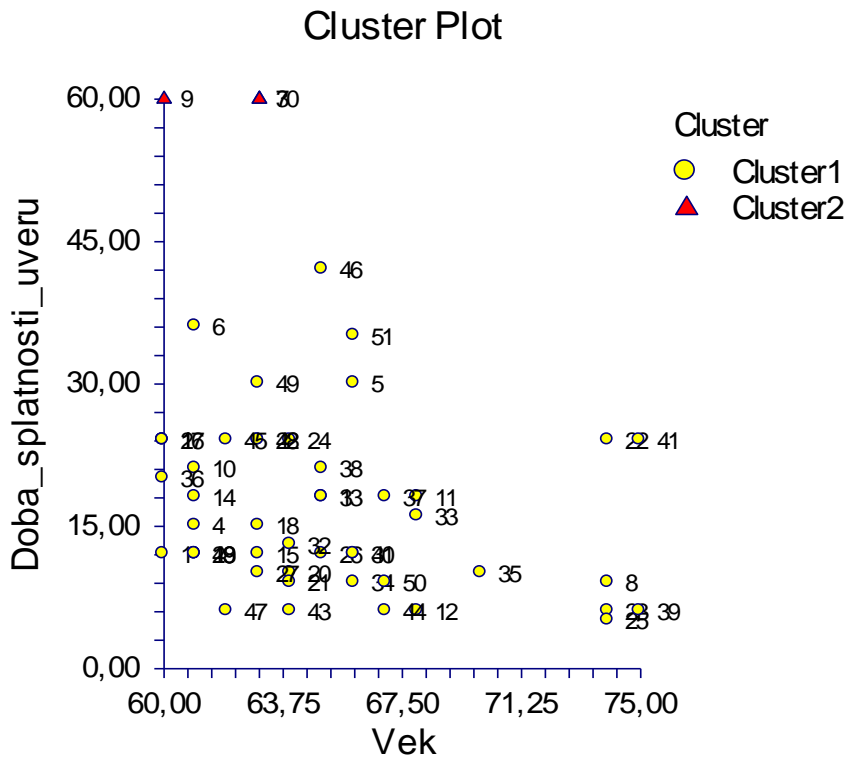
kde priemery  $\bar{\mathbf{x}}_s$  sú uvedené v tabuľke 5.14 a matica  $\mathbf{S}$  vyzerá nasledovne:

$$\begin{pmatrix} 68.51 & 1499.03 & -9.92 \\ 1499.03 & 7.17424 \times 10^6 & 2665.19 \\ -9.92 & 2665.19 & 18.37 \end{pmatrix}. \quad (5.3)$$

Znaky sú v poradí Vek, Výška úveru a Doba splatnosti úveru. Za odhad  $\pi_s$  bolo dané  $\frac{3}{50}$  v prvej skupine a  $\frac{47}{50}$  v druhej. Najprv aplikujeme skóry (5.2) na testovaciu databázu. Napočítaním diskriminačných skórov pomocou softvéru Mathematica a použitím pravidla (4.17) pre 50 objektov bolo zistené, že objekty boli rozdelené do rovnakých skupín ako v nehierarchickej analýze. Grafické zobrazenie skupín je viditeľné na obrázku 5.3. Preukázalo sa správne fungovanie lineárnych diskriminačných skórov.

Pridaným objektom bude klient, ktorý má 66 rokov s výškou úveru 766 mariek na 35 mesiacov. Diskriminačné skóry pre skupiny T1 a T2 sú 165,407 a 153,379. Zaradíme teda nový objekt podľa pravidla (4.17) do prvej skupiny.

Teraz porovnáme zaradenie nového objektu podľa diskriminačnej analýzy s nehierarchickou analýzou. Pre nehierarchickú analýzu budeme predpokladať dáta ako v predchádzajúcej kapitole s vyššie popísaným pridaným objektom 51. Celkovo pracujeme s maticou, ktorá má 51 objektov a 3 znaky. Pri konštrukcii budeme postupovať presne ako v predchádzajúcej kapitole. Pre kontrolu zaradenia už uvediem len graf.



Obr. 5.8: K- means metóda s novým prvkom 51

Predpokladajme štyri skupiny T1, T2, T3 a T4 tak, ako sú vyznačené na obrázku 5.6. Softvér rozdelil 50 objektov na skupiny s počtom 9, 25, 13 a 3 objekty.

Ak máme viac ako dve počiatočné skupiny, budeme počítať kvadratické diskriminačné skóry podľa (4.20):

$$D_s^{(i)} = -\frac{1}{2} \ln |\mathbf{S}_s| - \frac{1}{2} (\mathbf{x}_i - \bar{\mathbf{x}}_s)^T \mathbf{S}_s^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_s) + \ln \hat{\pi}_s, \quad s = 1, 2, 3 \quad i = 1, \dots, 51.$$

Priemery  $\bar{\mathbf{x}}_s$  nájdeme v nasledujúcej tabuľke:

Znak	1. zhluk	2. zhluk	3. zhluk	4. zhluk
Doba splatnosti úveru	12,00	12,52	26,69	60,00
Výška úveru	5 706,22	1 615,40	3 175,31	11 791,33
Vek	72,44	63,96	62,15	62,00

Tabuľka 5.16: Priemery štyroch zhlukov

Matice  $\mathbf{S}_i, i = 1, 2, 3, 4$  vyzerajú nasledovne:

$$\mathbf{S}_1: \begin{pmatrix} 61.75 & 1838.13 & 1.75 \\ 1838.13 & 1.67467 \times 10^7 & -8912.24 \\ 1.75 & -8912.24 & 8.52778 \end{pmatrix} \quad (5.4)$$

$$\mathbf{S}_2: \begin{pmatrix} 16.51 & 453.533 & 0.271667 \\ 453.533 & 776335. & -209.192 \\ 0.271667 & -209.192 & 5.29 \end{pmatrix} \quad (5.5)$$

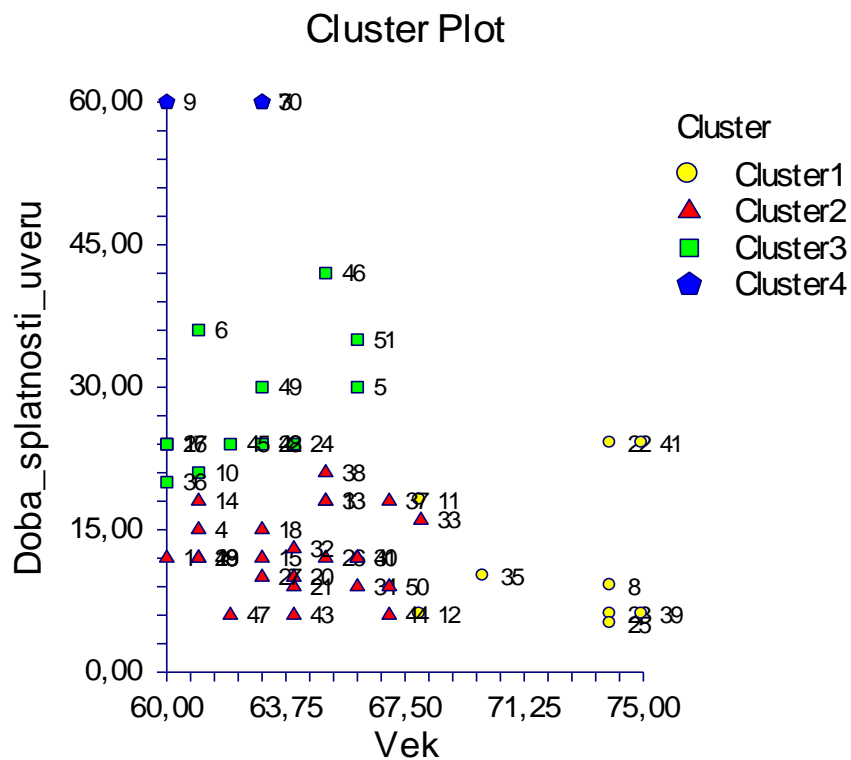
$$\mathbf{S}_3: \begin{pmatrix} 39.2308 & -313.064 & 6.46795 \\ -313.064 & 3.45944 \times 10^6 & 142.532 \\ 6.46795 & 142.532 & 4.14103 \end{pmatrix} \quad (5.6)$$

$$\mathbf{S}_4: \begin{pmatrix} 0. & 0. & 0. \\ 0. & 1.86797 \times 10^7 & -4486. \\ 0. & -4486. & 3. \end{pmatrix} \quad (5.7)$$

Vo štvrtom zhluky sa nachádzajú len 3 pozorovania, ktoré sú netypické vzhľadom ku dobe splatnosti úveru a výške úveru. Znak Vek má v tomto zhluky nulový rozptyl. Preto zúžime testovaciu databázu na 47 objektov a budeme pracovať so skupinami T1, T2 a T3. Za odhady  $\pi_s$  vezmeme  $\frac{p}{47}$ , kde  $p$  je počet objektov

v skupine. Napočítaním kvadratických diskriminačných skórov pre 47 objektov bolo zistené, že objekty boli rozdelené do rovnakých skupín ako v nehierarchickej analýze.

Pridaným objektom bude opäť klient, ktorý má 66 rokov s výškou úveru 766 mariek na 35 mesiacov. Diskriminačné skóry pre skupiny T1, T2 a T3 sú -31.88, -26.5 a -13.99. Nový objekt zaradíme podľa (4.17) do skupiny T3. Aplikáciou nehierarchickej zhlukovej analýzy na dátovú maticu, ktorá má 51 objektov a 3 znaky sa potvrdilo zaradenie nového klienta do zelenej skupiny T3, ako je vidieť na obrázku 5.9



Obr. 5.9: K- means metóda s novým prvkom 51

# Literatúra

- [1] Ajvazjan S., Bežajevová Z., Staroverov O.: *Metody vícerozměrné analýzy*. Praha: SNTL, 1981.
- [2] Anděl J.: *Základy matematické statistiky*. Praha: Matfyzpress, 2007.
- [3] Anděl J.: *Statistické metody*. Praha: Matfyzpress, 2003.
- [4] Anděl J.: *Matematická statistika*. Bratislava: SNTL/Alfa, 1978.
- [5] Anderberg M. R.: *Cluster Analysis for Applications*. New York: Academic Press, 1978.
- [6] Bečvář J.: *Lineární algebra*. Praha: Matfyzpress, 2005.
- [7] Bican L.: *Lineární algebra a geometrie*. Praha: Academia, 2000.
- [8] Hebák P., Hustopecký J.: *Vícerozměrné statistické metody s aplikacemi*. Praha: SNTL/Alfa, 1987.
- [9] Řehák J., Řeháková B.: *Analýza kategorizovaných dat v sociologii*. Praha: Academia, 1985.
- [10] [http : //www.stat.uni - muenchen.de/service/datenarchiv/kredit/kredit\\_e.html](http://www.stat.uni-muenchen.de/service/datenarchiv/kredit/kredit_e.html)
- [11] [http : //www.karlin.mff.cuni.cz/kulich/vyuka/statfpm/doc/statfpmvetnik.pdf](http://www.karlin.mff.cuni.cz/kulich/vyuka/statfpm/doc/statfpmvetnik.pdf)

# Zoznam obrázkov

4.1	Jednoduchý príklad dendrogramu . . . . .	17
5.1	Dendrogram štandardizovaných dát . . . . .	27
5.2	Dendrogram neštandardizovaných dát . . . . .	27
5.3	Zobrazenie zhlukov pre znaky Doba splatnosti úveru a Vek . . . .	31
5.4	Zobrazenie zhlukov znaky Doba splatnosti úveru a Výška úveru .	32
5.5	Zobrazenie zhlukov pre znaky Výška úveru a Vek . . . . .	32
5.6	Zobrazenie zhlukov pre znaky Doba splatnosti úroku a Vek . . . .	34
5.7	Zobrazenie zhlukov pre znaky Doba splatnosti úroku a Vek . . . .	35
5.8	K- means metóda s novým prvkom 51 . . . . .	36
5.9	K- means metóda s novým prvkom 51 . . . . .	38

# Zoznam tabuliek

3.1	Kontingenčná tabuľka . . . . .	14
5.1	Kontingenčná tabuľka pozorovaných početností znakov Cudzinec a Vlastníctvo . . . . .	23
5.2	Teoretické početnosti znakov Cudzinec a Vlastníctvo . . . . .	23
5.3	Hodnoty sčítancov $\chi^2$ rozdelenia znakov Cudzinec a Vlastníctvo .	23
5.4	Kontingenčná tabuľka pozorovaných početností znakov Bilancia na účte a Vlastníctvo . . . . .	24
5.5	Teoretické početnosti znakov Bilancia na účte a Vlastníctvo . . .	24
5.6	Hodnoty sčítancov $\chi^2$ rozdelenia znakov Bilancia na účte a Vlastníctvo	24
5.7	Kontingenčná tabuľka pozorovaných početností znakov Bilancia na účte a Cudzinec . . . . .	24
5.8	Teoretické početnosti znakov Bilancia na účte a Cudzinec . . . . .	25
5.9	Hodnoty sčítancov $\chi^2$ rozdelenia znakov Bilancia na účte a Cudzinec	25
5.10	Použitá dátová matica . . . . .	28
5.11	Vzdialenosti zhlukov pre štandardizované dáta . . . . .	29
5.12	Vzdialenosti zhlukov pre neštandardizované dáta . . . . .	30
5.13	Smerodajné odchýlky v zhlukoch pre dané znaky . . . . .	30
5.14	Priemery v zhlukoch pre dané znaky . . . . .	31
5.15	Hodnoty znakov vybraných objektov . . . . .	33
5.16	Priemery štyroch zhlukov . . . . .	37



# Zoznam príloh

## Príloha č. 1: CD obsahujúce

- DATA.xls - zdrojové dáta
- LenkaGodulova-BP.pdf - elektronická verzia práce
- výstupné protkoly z NCSS
- Súbor Microsoft Office Excel s výpočtami
- Súbor Mathematica s výpočtami