



**MATEMATICKO-FYZIKÁLNÍ
FAKULTA**
Univerzita Karlova

BAKALÁŘSKÁ PRÁCE

Lucie Janečková

**Klasifikace založená na směsových
modelech**

Katedra pravděpodobnosti a matematické statistiky

Vedoucí bakalářské práce: doc. RNDr. Arnošt Komárek, Ph.D.

Studijní program: Obecná matematika

Praha 2023

Prohlašuji, že jsem tuto bakalářskou práci vypracoval(a) samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů. Tato práce nebyla využita k získání jiného nebo stejného titulu.

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 má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle §60 odst. 1 autorského zákona.

V dne

Podpis autora

Zde bych ráda poděkovala mému vedoucímu práce, doc. RNDr. Arnoštu Komárkovi, Ph.D., za cenné rady a čas, který mi při psaní práce věnoval.

Název práce: Klasifikace založená na směsových modelech

Autor: Lucie Janečková

Katedra: Katedra pravděpodobnosti a matematické statistiky

Vedoucí bakalářské práce: doc. RNDr. Arnošt Komárek, Ph.D., Katedra pravděpodobnosti a matematické statistiky

Abstrakt: Tato práce se zabývá klasifikací založenou na směsových modelech, a to převážně na modelech konečných normálních. Nejprve jsou zavedeny základní definice a vlastnosti konečných směsí. Následně je zde popsána metoda maximální věrohodnosti a její úskalí v kontextu konečných směsí, kterou používáme pro odhadování neznámých parametrů. Poté je popsán EM algoritmus, který je používán pro získání maximálně věrohodných odhadů a explicitně spočteny vzorce pro jednu iteraci EM algoritmu. V poslední části je ukázáno, jak lze konečné normální směsi využít ke klasifikaci.

Klíčová slova: konečná směs, normální směs, klasifikace, EM algoritmus, věrohodnost

Title: Classification based on mixture models

Author: Lucie Janečková

Department: Department of Probability and Mathematical Statistics

Supervisor: doc. RNDr. Arnošt Komárek, Ph.D., Department of Probability and Mathematical Statistics

Abstract: This thesis deals with classification based on mixture models, mainly on models finite normal. At first, there are introduced basic definitions and characteristics of finite mixtures. Afterwards there is described the maximum likelihood method and her obstacles in context of finite mixtures, which we are using for unknown parameters estimation. Then there is described EM algorithm, that is used to obtain the maximum likelihood estimator and there are calculated the formulae for one iteration of EM algorithm. In the last part there is shown, how can finite normal mixtures be used for classification.

Keywords: finite mixture, normal mixture, classification, EM algorithm, likelihood

Obsah

Úvod	2
1 Konečná směs	3
1.1 Základní definice	3
1.2 Interpretace konečné směsi a klasifikace	5
2 Metoda maximální věrohodnosti	6
2.1 Omezenost věrohodnostní funkce	6
2.2 Další vlastnosti	9
3 EM algoritmus	10
3.1 Nekompletní data	10
3.2 Věrohodnostní funkce kompletních dat a její vlastnosti	11
3.3 EM algoritmus	13
4 Praktické použití	15
4.1 Popis dat	15
4.2 Klasifikace	18
Závěr	20
Seznam použité literatury	21

Úvod

Uvažujme situaci, kdy máme populaci rozdělenou do předem známého počtu skupin, ale nevíme, který jedinec patří do které skupiny. V každé skupině pozorujeme jiné vzorce chování. Naším cílem tedy bude každého jedince z populace zařadit do příslušné skupiny na základě jeho pozorovaného vzorce chování.

Tuto situaci modelujeme právě pomocí konečné směsi, kterou si v první kapitole zavedeme. Dále se také zmíníme o její možné interpretaci a matematicky popíšeme výše zmíněný cíl práce.

Pro správné zařazení jedinců do skupin budeme potřebovat odhadnout neznámé parametry. K tomu nám poslouží metoda maximální věrohodnosti, kterou popíšeme v druhé kapitole, včetně problémů, které nám použití této metody v kontextu konečných směsí může činit.

Ke konečnému nalezení odhadů použijeme EM algoritmus, který ve třetí kapitole popíšeme a vypočteme konkrétní vzorce pro odhady neznámých parametrů v normální směsi. Teoretická část práce vychází z knihy McLachlan a Peel (2000).

V poslední kapitole pak využijeme nabytých znalostí pro analýzu archeologických dat. Konkrétně z největších délek mozkovny a týlních úhlů budeme chtít roztřídit lebky podle místa nálezů. K tomu nám poslouží právě EM algoritmus.

1. Konečná směs

Mějme populaci sestávající z $g \in \mathbb{N}$ skupin, tu se budeme snažit popsat pomocí konečné, resp. normální směsi, což jsou pojmy, které si v této kapitole zavedeme. Dále si uvedeme jejich základní vlastnosti a popíšeme, jak lze směsi interpretovat a jak se dají použít ke klasifikaci.

1.1 Základní definice

Definice 1. *Nechť $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ je náhodný výběr p -rozměrných náhodných vektorů o rozsahu $n \in \mathbb{N}$, kde $p \in \mathbb{N}$. Nechť pro $k \in \{1, \dots, g\}$ jsou f_k libovolné hustoty vzhledem k p -rozměrné Lebesgueově míře na \mathbb{R}^p , které odpovídají rozdělení se střední hodnotou $\boldsymbol{\mu}_k \in \mathbb{R}^p$ a varianční maticí $\boldsymbol{\Sigma}_k \in \mathbb{R}^{p \times p}$. Pokud má náhodný vektor \mathbf{Y}_j rozdělení s hustotou danou*

$$f(\mathbf{y}) = \sum_{k=1}^g w_k f_k(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^p, \quad (1.1)$$

vzhledem k Lebesgueově míře na \mathbb{R}^p , kde

$$w_k \in (0, 1), \quad \forall k = 1, \dots, g$$

a

$$\sum_{k=1}^g w_k = 1,$$

potom říkáme, že \mathbf{Y}_j má rozdělení, které nazýváme **konečná směs o g -složkách**. Hustotu danou vztahem (1.1) nazýváme **hustota směsi o g -složkách**. Čísla w_k nazýváme **váhy** a hustoty f_k **komponenty směsi**.

Nabízí se otázka, proč uvažujeme váhy v otevřeném intervalu $(0, 1)$ a ne v uzavřeném. Krajní hodnoty jsou totiž pro nás nezajímavé. Kdybychom uvažovali hodnotu 0 u některé z komponent, pak bychom pracovali s $g - 1$ složkovou směsí. Kdybychom uvažovali hodnotu 1, dostali bychom triviální jednosložkovou směs.

Tuto definici můžeme vyložit tak, že celková hustota populace \mathbf{Y}_j je $f(\mathbf{y})$, hustoty jednotlivých skupin jsou $f_k(\mathbf{y})$ a poměrné zastoupení v jednotlivých skupinách jsou právě váhy w_k .

V následující větě určíme střední hodnotu a rozptyl náhodné veličiny s hustotou konečné směsi.

Věta 1. *Nechť \mathbf{Y} je náhodný vektor s hustotou danou vzorcem (1.1). Označme $\boldsymbol{\Gamma}_k$ matici druhých momentů k -té složky. Pak platí*

$$E(\mathbf{Y}) = \sum_{k=1}^g w_k \boldsymbol{\mu}_k,$$
$$\text{var}(\mathbf{Y}) = \sum_{k=1}^g w_k \boldsymbol{\Gamma}_k - \left(\sum_{k=1}^g w_k \boldsymbol{\mu}_k \right) \left(\sum_{k=1}^g w_k \boldsymbol{\mu}_k \right)^\top.$$

Důkaz. Pro $\mathbf{y} \in \mathbb{R}^p$ počítejme

$$\mathbb{E}(\mathbf{Y}) = \int_{\mathbb{R}^p} \mathbf{y} \sum_{k=1}^g w_k f_k(\mathbf{y}) \, d\mathbf{y} = \sum_{k=1}^g w_k \int_{\mathbb{R}^p} \mathbf{y} f_k(\mathbf{y}) \, d\mathbf{y} = \sum_{k=1}^g w_k \boldsymbol{\mu}_k,$$

$$\mathbb{E}(\mathbf{Y}\mathbf{Y}^\top) = \int_{\mathbb{R}^p} \mathbf{y}\mathbf{y}^\top \sum_{k=1}^g w_k f_k(\mathbf{y}) \, d\mathbf{y} = \sum_{k=1}^g w_k \int_{\mathbb{R}^p} \mathbf{y}\mathbf{y}^\top f_k(\mathbf{y}) \, d\mathbf{y} = \sum_{k=1}^g w_k \boldsymbol{\Gamma}_k,$$

kde používáme linearitu integrálu.

Tedy celkem pro varianční matici platí

$$\text{var}(\mathbf{Y}) = \mathbb{E}(\mathbf{Y}\mathbf{Y}^\top) - (\mathbb{E}(\mathbf{Y}))(\mathbb{E}(\mathbf{Y}))^\top = \sum_{k=1}^g w_k \boldsymbol{\Gamma}_k - \left(\sum_{k=1}^g w_k \boldsymbol{\mu}_k\right) \left(\sum_{k=1}^g w_k \boldsymbol{\mu}_k\right)^\top.$$

□

Pokud v definici 1 vezmeme za f_k hustoty normálního rozdělení, můžeme definovat normální směs.

Definice 2. Necht $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ je náhodný výběr p -rozměrných náhodných vektorů s normálním rozdělením o rozsahu $n \in \mathbb{N}$, kde $p \in \mathbb{N}$. Necht pro $k \in \{1, \dots, g\}$ jsou $\phi_k(\mathbf{y}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ hustoty normálního rozdělení se střední hodnotou $\boldsymbol{\mu}_k \in \mathbb{R}^p$ a varianční maticí $\boldsymbol{\Sigma}_k \in \mathbb{R}^{p \times p}$. Označme $\boldsymbol{\psi} \equiv (w_1, \dots, w_g, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_g)$ sadu neznámých parametrů. Pokud má náhodný vektor \mathbf{Y}_j rozdělení s hustotou danou

$$\phi(\mathbf{y}; \boldsymbol{\psi}) = \sum_{k=1}^g w_k \phi_k(\mathbf{y}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad \mathbf{y} \in \mathbb{R}^p, \quad (1.2)$$

vzledem k Lebesgueově míře na \mathbb{R}^p , kde

$$w_k \in (0,1), \quad \forall k = 1, \dots, g$$

a

$$\sum_{k=1}^g w_k = 1,$$

potom říkáme, že \mathbf{Y}_j má rozdělení, které nazýváme **konečná normální směs o g -složkách**. Hustotu danou (1.2) nazýváme **hustota normální směsi o g -složkách**. Čísla w_k nazýváme **váhy** a hustoty $\phi_k(\mathbf{y}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ **komponenty normální směsi**.

Dále se v práci budeme zabývat hlavně modelem homoskedastické normální směsi, tedy směsi, kde jsou všechny varianční matice $\boldsymbol{\Sigma}_k$ stejné. Směsi heteroskedastické (směsi s různými variančními maticemi $\boldsymbol{\Sigma}_k$) nám totiž mohou dělat problémy při odhadování parametrů pomocí metody maximální věrohodnosti, jak si ukážeme v kapitole 2.

1.2 Interpretace konečné směsi a klasifikace

Zde si uvedeme dvě možnosti, jak lze konečné směsi interpretovat. První z nich je pomocí náhodného vektoru, který nám identifikuje, z které komponenty směsi náš náhodný vektor pochází. Druhá nám popisuje, že chování našeho náhodného vektoru v celkové populaci lze popsat pomocí jeho chování v různých skupinách.

Jednou možností, jak konečnou směs interpretovat, je pomocí náhodného vektoru $\mathbf{Z}_j = (Z_{1j}, \dots, Z_{gj})^\top$ s multinomickým rozdělením, kde

$$Z_{kj} = \begin{cases} 1, & \text{pokud náhodný vektor } \mathbf{Y}_j \text{ náleží do } k\text{-té komponenty směsi,} \\ 0, & \text{jinak.} \end{cases}$$

Jinými slovy

$$\mathbf{Z}_j \sim \text{Mult}_g(\mathbf{1}, \mathbf{w}),$$

kde $\mathbf{w} = (w_1, \dots, w_g)^\top$ je vektor vah. Náhodný vektor \mathbf{Z}_j tedy určuje, z jaké komponenty směsi náhodný vektor \mathbf{Y}_j pochází.

Místo toho lze k interpretaci použít náhodné veličiny $U_j, j = 1, \dots, n$, které nabývají hodnot $1, \dots, g$ s pravděpodobnostmi

$$P(U_j = k) = w_k, \quad k = 1, \dots, g.$$

Pokud podmíněnou hustotu náhodného vektoru \mathbf{Y}_j při jevu $U_j = k$ označíme f_k , pak pro nepodmíněnou hustotu \mathbf{Y}_j z věty o úplné pravděpodobnosti platí

$$f(\mathbf{y}_j) = \sum_{k=1}^g f(\mathbf{y}_j | U_j = k) P(U_j = k) = \sum_{k=1}^g f_k(\mathbf{y}_j) w_k,$$

což odpovídá vztahu (1.1).

Nechť \mathbf{y}_j je napozorovaná hodnota náhodného vektoru \mathbf{Y}_j . Označme $\tau_k(\mathbf{y}_j)$ pravděpodobnost, že \mathbf{Y}_j náleží do k -té komponenty směsi. Pro $\tau_k(\mathbf{y}_j)$ pak platí

$$\tau_k(\mathbf{y}_j) := \frac{w_k f_k(\mathbf{y}_j)}{\sum_{h=1}^g w_h f_h(\mathbf{y}_j)} = P(Z_{kj} = 1 | \mathbf{Y}_j = \mathbf{y}_j) = \frac{P(U_j = k) f_k(\mathbf{y}_j)}{\sum_{h=1}^g P(U_j = h) f_h(\mathbf{y}_j)}, \quad (1.3)$$

což vychází z Bayesovy věty.

V samostatné klasifikaci právě tyto hodnoty $\tau_k(\mathbf{y}_j)$ hrají klíčovou roli. Ta probíhá následovně:

- Vezmeme pozorovanou hodnotu \mathbf{y}_j a spočteme pro ni všechny pravděpodobnosti $\tau_k(\mathbf{y}_j)$, $k = 1, \dots, g$.
- Spočteme $\arg \max_{k=1, \dots, g} \tau_k(\mathbf{y}_j)$.
- Zařadíme pozorování \mathbf{y}_j do této skupiny.

Jinými slovy určíme predikce hodnot náhodných veličin Z_{kj} , $k = 1, \dots, g$, $j = 1, \dots, n$, neboli $\hat{Z}_{kj} = \begin{cases} 1, & \text{pokud } \arg \max_{k=1, \dots, g} \tau_k(\mathbf{y}_j) = k, \\ 0, & \text{jinak.} \end{cases}$

Mohlo by se stát, že by se největší pravděpodobnosti $\tau_k(\mathbf{y}_j)$ mohli rovnat. Potom bychom vybrali libovolnou z nich.

2. Metoda maximální věrohodnosti

Pro klasifikaci napozorovaných subjektů z populace popsanou v první kapitole potřebujeme odhadnout neznámé parametry. Jeden z nejčastěji používaných způsobů hledání odhadů parametrů je metoda maximální věrohodnosti, jejíž použití pro model konečné normální směsi si v této kapitole popíšeme. Pak si ukážeme, kdy je věrohodnostní, resp. log-věrohodnostní funkce omezená a kdy nikoli.

Nechť $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ je náhodný výběr z rozdělení s hustotou (1.2). Označme $\boldsymbol{\psi} \equiv (w_1, \dots, w_g, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_g)$ sadu neznámých parametrů, kde w_g lze vynechat, neboť jej můžeme dopočítat z ostatních hodnot jako $1 - \sum_{k=1}^{g-1} w_k$. Sestavíme věrohodnostní funkci pro $\boldsymbol{\psi} \in \boldsymbol{\Psi}$:

$$L(\boldsymbol{\psi}) = \prod_{j=1}^n \phi(\mathbf{Y}_j; \boldsymbol{\psi}) = \prod_{j=1}^n \left(\sum_{k=1}^g w_k \phi_k(\mathbf{Y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right).$$

Obdobně sestavíme i log-věrohodnostní funkci pro $\boldsymbol{\psi}$:

$$\begin{aligned} \ell(\boldsymbol{\psi}) &= \log(L(\boldsymbol{\psi})) = \log \left(\prod_{j=1}^n \left(\sum_{k=1}^g w_k \phi_k(\mathbf{Y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \right) \\ &= \sum_{j=1}^n \left(\log \left(\sum_{k=1}^g w_k \phi_k(\mathbf{Y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \right). \end{aligned}$$

Odhad $\boldsymbol{\psi}$ metodou maximální věrohodnosti označíme jako $\hat{\boldsymbol{\psi}}$, pokud existuje. Nalezneme ho jako

$$\hat{\boldsymbol{\psi}} = \arg \max_{\boldsymbol{\psi} \in \boldsymbol{\Psi}} L(\boldsymbol{\psi}) = \arg \max_{\boldsymbol{\psi} \in \boldsymbol{\Psi}} \ell(\boldsymbol{\psi}).$$

2.1 Omezenost věrohodnostní funkce

Pro existenci $\hat{\boldsymbol{\psi}}$ potřebujeme, aby věrohodnostní, resp. log-věrohodnostní funkce byla omezená. To ale, jak ukazuje následující tvrzení, není vždy splněno.

Tvrzení 2. *Nechť Y_1, \dots, Y_n je náhodný výběr z rozdělení s hustotou jednorozměrné normální směsi o dvou složkách, tedy*

$$\phi(y; \boldsymbol{\psi}) = w \phi_1(y; \mu_1, \sigma_1^2) + (1 - w) \phi_2(y; \mu_2, \sigma_2^2),$$

kde $\boldsymbol{\psi} = (w, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2)^\top$, $y \in \mathbb{R}$, $w \in (0, 1)$, $\mu_1, \mu_2 \in \mathbb{R}$, $\sigma_1^2, \sigma_2^2 \in (0, \infty)$. Označme y_1, \dots, y_n realizace Y_1, \dots, Y_n . Pokud $\mu_1 = y_1$ a $\sigma_1^2 \rightarrow 0$, tak je log-věrohodnostní funkce neomezená.

Důkaz. Sestavíme log-věrohodnostní funkci pro $\boldsymbol{\psi}$:

$$\ell(\boldsymbol{\psi}) = \sum_{j=1}^n \log \left(w \phi_1(y_j; \mu_1, \sigma_1^2) + (1 - w) \phi_2(y_j; \mu_2, \sigma_2^2) \right), \quad (2.1)$$

kde

$$\phi_k(y_j; \mu_k, \sigma_k^2) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(y_j - \mu_k)^2}{2\sigma_k^2}\right), \quad k = 1, 2.$$

Dosaďme do (2.1) $\mu_1 = y_1$ a dostáváme

$$\ell(\boldsymbol{\psi}) = \sum_{j=1}^n \log\left(w \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(y_j - y_1)^2}{2\sigma_1^2}\right) + (1-w) \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(y_j - \mu_2)^2}{2\sigma_2^2}\right)\right).$$

Chceme $\ell(\boldsymbol{\psi})$ zespoda omezit tak, aby nám vyšel výraz tvaru

$$-\log(\sigma_1^2) + Z,$$

kde Z je funkce nezávislejší na σ_1^2 . Pak totiž

$$\lim_{\sigma_1^2 \rightarrow 0} (-\log(\sigma_1^2) + Z) = \infty,$$

což chceme dokázat.

V argumentu logaritmu vezmeme pro $j = 1$ pouze první člen a pro $j = 2, \dots, n$ druhý člen, který na σ_1^2 nezávisí:

$$\begin{aligned} \ell(\boldsymbol{\psi}) &> \log\left(w \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(y_1 - y_1)^2}{2\sigma_1^2}\right)\right) \\ &\quad + \sum_{i=2}^n \log\left((1-w) \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(y_i - \mu_2)^2}{2\sigma_2^2}\right)\right). \end{aligned}$$

Označíme

$$Z := \sum_{i=2}^n \log\left((1-w) \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(y_i - \mu_2)^2}{2\sigma_2^2}\right)\right),$$

pak máme

$$\begin{aligned} \ell(\boldsymbol{\psi}) &> \log\left(w \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp(0)\right) + Z \\ &= \left(\log\left(\frac{\sqrt{2\pi\sigma_1^2}}{w}\right)\right)^{-1} + Z \\ &= -\log\left(\frac{\sqrt{2\pi\sigma_1^2}}{w}\right) + Z, \end{aligned}$$

což je požadovaný tvar. □

Naopak další věta ukazuje, že věrohodnostní funkce homoskedatické normální směsi, tedy směsi, kde všechny její složky mají stejný rozptyl, omezená je.

Věta 3. Necht Y_1, \dots, Y_n je náhodný výběr z rozdělení s hustotou g -složkové homoskedatické normální směsi, kde $g < n$, tedy s hustotou

$$\phi(y; \boldsymbol{\psi}) = \sum_{k=1}^g w_k \phi_k(y; \mu_k, \sigma^2),$$

kde $\boldsymbol{\psi} = (w_1, \dots, w_g, \mu_1, \dots, \mu_g, \sigma^2)^\top$, $y \in \mathbb{R}$, $w_k \in (0, 1)$, $\sum_{k=1}^g w_k = 1$, $\mu_k \in \mathbb{R}$, $k = 1, \dots, g$, $\sigma^2 \in (0, \infty)$. Označme y_1, \dots, y_n realizace Y_1, \dots, Y_n . Pak je věrohodnostní funkce shora omezená.

Důkaz. Sestavíme věrohodnostní funkci pro $\boldsymbol{\psi}$:

$$L(\boldsymbol{\psi}) = \prod_{j=1}^n \left(\sum_{k=1}^g w_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_j - \mu_k)^2}{2\sigma^2}\right) \right).$$

Platí:

$$L(\boldsymbol{\psi}) < \prod_{j=1}^n \left(\sum_{k=1}^g \frac{1}{\sigma} \exp\left(-\frac{(y_j - \mu_k)^2}{2\sigma^2}\right) \right).$$

Z předpokladů věty máme více pozorování než je složek směsi, a tedy:

$$\exists \delta > 0 \quad \exists n_0 \in 1, \dots, n : \forall k \in 1, \dots, g : |y_{n_0} - \mu_k| > \delta.$$

Upravíme zvlášť první člen součinu a zbylých $n - 1$ členů.

Bez újmy na obecnosti necht $n_0 = 1$ (jinak bychom upravovali jiný člen zvlášť).

$$\sum_{k=1}^g \frac{1}{\sigma} \exp\left(-\frac{(y_1 - \mu_k)^2}{2\sigma^2}\right) < \sum_{k=1}^g \frac{1}{\sigma} \exp\left(-\frac{\delta^2}{2\sigma^2}\right) = \frac{g}{\sigma} \exp\left(-\frac{\delta^2}{2\sigma^2}\right),$$

$$\prod_{j=2}^n \left(\sum_{k=1}^g \frac{1}{\sigma} \exp\left(-\frac{(y_1 - \mu_k)^2}{2\sigma^2}\right) \right) < \prod_{j=2}^n \left(\sum_{k=1}^g \frac{1}{\sigma} \exp\left(-\frac{0^2}{2\sigma^2}\right) \right) = \left(\frac{g}{\sigma}\right)^{n-1}.$$

Celkem tedy

$$L(\boldsymbol{\psi}) < \left(\frac{g}{\sigma}\right)^n \exp\left(-\frac{\delta^2}{2\sigma^2}\right).$$

A platí

$$\lim_{\sigma^2 \rightarrow \infty} \left(\frac{g}{\sigma}\right)^n \exp\left(-\frac{\delta^2}{2\sigma^2}\right) = 0,$$

$$\lim_{\sigma^2 \rightarrow 0} \left(\frac{g}{\sigma}\right)^n \exp\left(-\frac{\delta^2}{2\sigma^2}\right) = 0.$$

Z toho plyne, že věrohodnostní funkce je shora omezená. □

2.2 Další vlastnosti

Odhad metodou maximální věrohodnosti můžeme také dostat jako řešení věrohodnostních rovnic, tedy

$$\frac{\partial L(\boldsymbol{\psi})}{\partial \boldsymbol{\psi}} = \mathbf{0},$$

či ekvivalentně

$$\frac{\partial \ell(\boldsymbol{\psi})}{\partial \boldsymbol{\psi}} = \mathbf{0},$$

a tak také budeme v kapitole 3 postupovat.

Zbývá nám odpovědět na otázku, zda existuje maximum pouze jedno. Jednou z vlastností konečných směsí je to, že je invariantní na „přečíslování“ složek směsi. Tedy pokud máme g -složkovou konečnou směs, existuje $g!$ permutací indexů složek, neboli $g!$ sad neznámých parametrů $\boldsymbol{\psi}$, které reprezentují stejné rozdělení, a tedy máme nejméně $g!$ (stejně hodnotných) maxim (to budeme nazývat sada maxim).

I přesto ale maxim, respektive sad maxim, je většinou více. Samotný EM algoritmus, který budeme pro odhad maxim v práci používat, zavedený v kapitole 3, se používá k nalezení ne nutně globálního maxima. Přesto se s tímto výsledkem spokojíme, neboť lze ukázat, že vlastnosti těchto lokálních maxim jsou velmi podobné.

3. EM algoritmus

EM algoritmus je jednou z možností hledání maximálně věrohodného odhadu. Pochází z anglického Expectation-Maximization. Je to iterativní algoritmus, který postupně optimalizuje parametry modelu na základě pozorovaných dat. Dělí se na dva kroky: E-krok (výpočet střední hodnoty) a M-krok (maximalizace). V E-kroku se vypočítávají pravděpodobnosti příslušnosti každého pozorování k jednotlivým komponentám. V M-kroku se pak aktualizují parametry komponent na základě těchto pravděpodobností. Tento algoritmus byl poprvé popsán v práci Dempstera, Lairdové a Rubina (Dempster, Laird a Rubin, 1977).

V této kapitole si nejprve představíme problém nekompletních dat, poté si řekneme pár potřebných vlastností věrohodnostní funkce kompletních dat a spočteme odhady neznámých parametrů homoskedatické normální směsi. Nakonec si shrneme, jak EM algoritmus funguje.

3.1 Nekompletní data

Uvažujme náhodný výběr $\mathbf{Y}_c = \left(\left(\begin{array}{c} \mathbf{Y}_1 \\ \mathbf{Z}_1 \end{array} \right), \dots, \left(\begin{array}{c} \mathbf{Y}_n \\ \mathbf{Z}_n \end{array} \right) \right)^\top$, kde $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_n)^\top$ je náhodný výběr z rozdělení s hustotou homoskedastické normální směsi a $\mathbf{Z} = (\mathbf{Z}_1^\top, \dots, \mathbf{Z}_n^\top)^\top$ je tedy náhodný vektor určující, z jaké komponenty směsi náhodné vektory $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ pochází popsány v kapitole 1.2. Tento náhodný výběr nazýváme kompletní data. Náhodný výběr $(\mathbf{Y}_1, \dots, \mathbf{Y}_n)^\top$ nazveme nekompletní data.

Mějme $\mathbf{y}_c = \left(\left(\begin{array}{c} \mathbf{y}_1 \\ \mathbf{z}_1 \end{array} \right), \dots, \left(\begin{array}{c} \mathbf{y}_n \\ \mathbf{z}_n \end{array} \right) \right)^\top$ potenciálně napozorovaný kompletní náhodný výběr, kde $\mathbf{z}_j = (z_{1j}, \dots, z_{gj})^\top$. Dále označme $\boldsymbol{\xi} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma})$ sadu vybraných neznámých parametrů a $\boldsymbol{\psi} = (w_1, \dots, w_g, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma})$ sadu všech neznámých parametrů.

Sestavíme log-věrohodnostní funkci pro $\boldsymbol{\psi} = (w_1, \dots, w_g, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma})$ a náhodný výběr Y_1, \dots, Y_n

$$\ell(\boldsymbol{\psi}) = \sum_{j=1}^n \log \left(\sum_{k=1}^g w_k \phi_k(\mathbf{Y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) \right).$$

Maximalizace $\ell(\boldsymbol{\psi})$ je obtížná (neboť s logaritmem součtu se velmi špatně pracuje). Proto v této situaci používáme právě EM algoritmus. Dále také proto, že je schopen z maximalizace nekompletních dat přejít na maximalizaci dat kompletních.

Naším cílem je dostat $\hat{\boldsymbol{\psi}} = (\hat{w}_1, \dots, \hat{w}_g, \hat{\boldsymbol{\mu}}_1, \dots, \hat{\boldsymbol{\mu}}_g, \hat{\boldsymbol{\Sigma}})$ (maximálně věrohodný odhad parametru $\boldsymbol{\psi}$) jako řešení soustavy věrohodnostních rovnic

$$\frac{\partial \ell(\boldsymbol{\psi})}{\partial \boldsymbol{\psi}} = \mathbf{0}.$$

3.2 Věrohodnostní funkce kompletních dat a její vlastnosti

Zásadní úlohu v EM algoritmu hraje log-věrohodnostní funkce pro $\boldsymbol{\psi}$ z kompletních dat. Pro tu si nejprve vyjádříme hustotu náhodného vektoru $(\mathbf{Y}_j^\top, \mathbf{Z}_j^\top)^\top$, pro $j = 1, \dots, n$:

$$\begin{aligned} f_c(\mathbf{Y}_j, \mathbf{Z}_j; \boldsymbol{\psi}) &= f(\mathbf{Y}_j | \mathbf{Z}_j; \boldsymbol{\psi}) f(\mathbf{Z}_j; \boldsymbol{\psi}) \\ &= \sum_{k=1}^g Z_{kj} \phi_k(\mathbf{Y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) \sum_{k=1}^g Z_{kj} w_k = \prod_{k=1}^g (\phi_k(\mathbf{Y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) w_k)^{Z_{kj}}. \end{aligned}$$

Z toho pro log-věrohodnostní funkci pro $\boldsymbol{\psi}$ z kompletních dat platí

$$\ell_c(\boldsymbol{\psi}) = \sum_{k=1}^g \sum_{j=1}^n Z_{kj} (\log w_k + \log (\phi_k(\mathbf{Y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}))).$$

V následujícím lemmatu si shrneme základní vlastnosti $\ell_c(\boldsymbol{\psi})$, které budeme dále potřebovat.

Lemma 4. *Nechť $\tau_k(\mathbf{y}_j; \boldsymbol{\psi}) = \frac{w_k \phi_k(\mathbf{y}_j; \boldsymbol{\psi})}{\sum_{h=1}^g w_h \phi_h(\mathbf{y}_j; \boldsymbol{\psi})}$, pak platí*

1.

$$E_{\boldsymbol{\psi}}[Z_{kj} | \mathbf{Y}_j = \mathbf{y}_j] = \tau_k(\mathbf{y}_j; \boldsymbol{\psi}),$$

2.

$$\begin{aligned} E_{\boldsymbol{\psi}}[\ell_c(\boldsymbol{\psi}) | \mathbf{Y} = \mathbf{y}] &= \\ \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) &\left(\log w_k - \frac{1}{2} \log (2\pi)^p - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{y}_j - \boldsymbol{\mu}_k)^\top (\boldsymbol{\Sigma})^{-1} (\mathbf{y}_j - \boldsymbol{\mu}_k) \right), \end{aligned}$$

kde $|\boldsymbol{\Sigma}|$ značí determinant matice $\boldsymbol{\Sigma}$,

3.

$$\begin{aligned} \arg \max_{\boldsymbol{\psi}} E_{\boldsymbol{\psi}}[\ell_c(\boldsymbol{\psi}) | \mathbf{Y} = \mathbf{y}] &= \\ \left(\frac{1}{n} \sum_{j=1}^n \tau_1(\mathbf{y}_j; \boldsymbol{\psi}), \frac{1}{n} \sum_{j=1}^n \tau_2(\mathbf{y}_j; \boldsymbol{\psi}), \dots, \frac{1}{n} \sum_{j=1}^n \tau_g(\mathbf{y}_j; \boldsymbol{\psi}), \right. \\ \frac{\sum_{j=1}^n \tau_1(\mathbf{y}_j; \boldsymbol{\psi}) \mathbf{y}_j}{\sum_{j=1}^n \tau_1(\mathbf{y}_j; \boldsymbol{\psi})}, \frac{\sum_{j=1}^n \tau_2(\mathbf{y}_j; \boldsymbol{\psi}) \mathbf{y}_j}{\sum_{j=1}^n \tau_2(\mathbf{y}_j; \boldsymbol{\psi})}, \dots, \frac{\sum_{j=1}^n \tau_g(\mathbf{y}_j; \boldsymbol{\psi}) \mathbf{y}_j}{\sum_{j=1}^n \tau_g(\mathbf{y}_j; \boldsymbol{\psi})}, \\ \left. \frac{1}{n} \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) (\mathbf{y}_j - \boldsymbol{\mu}_k) (\mathbf{y}_j - \boldsymbol{\mu}_k)^\top \right)^\top. \end{aligned}$$

Důkaz.

1.

$$E_{\boldsymbol{\psi}}[Z_{kj} | \mathbf{Y}_j = \mathbf{y}_j] = P_{\boldsymbol{\psi}}(Z_{kj} = 1 | \mathbf{Y}_j = \mathbf{y}_j) = \tau_k(\mathbf{y}_j; \boldsymbol{\psi}),$$

kde druhá rovnost plyne ze vztahu (1.3).

2. Vypočteme hodnotu této funkce a upravíme do tvaru vhodného k derivování.

$$\begin{aligned}
\mathbf{E}_\psi[\ell_c(\boldsymbol{\psi})|\mathbf{Y} = \mathbf{y}] &= \sum_{k=1}^g \sum_{j=1}^n \mathbf{E}_\psi[Z_{kj}|\mathbf{Y}_j = \mathbf{y}_j] (\log w_k + \log(\phi_k(\mathbf{y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}))) \\
&= \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) (\log w_k + \log(\phi_k(\mathbf{y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}))) \\
&= \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \left(\log w_k + \log \left(\frac{1}{\sqrt{(2\pi)^p |\boldsymbol{\Sigma}|}} \exp \left(\frac{1}{2} (\mathbf{y}_j - \boldsymbol{\mu}_k)^\top (\boldsymbol{\Sigma})^{-1} (\mathbf{y}_j - \boldsymbol{\mu}_k) \right) \right) \right) \\
&= \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \left(\log w_k - \frac{1}{2} \log (2\pi)^p - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{y}_j - \boldsymbol{\mu}_k)^\top (\boldsymbol{\Sigma})^{-1} (\mathbf{y}_j - \boldsymbol{\mu}_k) \right).
\end{aligned}$$

3. Chceme najít řešení rovnice $\hat{\boldsymbol{\psi}} = \arg \max_{\boldsymbol{\psi}} \mathbf{E}_\psi[\ell_c(\boldsymbol{\psi})|\mathbf{Y} = \mathbf{y}]$. Parametr $\boldsymbol{\psi}$ lze rozdělit na $\mathbf{w} = (w_1, \dots, w_g)^\top$ a $\boldsymbol{\xi} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma})$. Stejně tak lze $\mathbf{E}_\psi[\ell_c(\boldsymbol{\psi})|\mathbf{Y} = \mathbf{y}]$ rozdělit na část závisící na \mathbf{w} a část závisící na $\boldsymbol{\xi}$ jako

$$\arg \max_{\boldsymbol{\psi}} \mathbf{E}_\psi[\ell_c(\boldsymbol{\psi})|\mathbf{Y} = \mathbf{y}] = S(\mathbf{w}, \boldsymbol{\psi}) + T(\boldsymbol{\xi}, \boldsymbol{\psi}),$$

kde

$$S(\mathbf{w}, \boldsymbol{\psi}) = \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \log w_k$$

a

$$T(\boldsymbol{\xi}, \boldsymbol{\psi}) = \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \log(\phi_k(\mathbf{y}_j; \boldsymbol{\mu}_k, \boldsymbol{\Sigma})).$$

Tím můžeme maximalizaci $\mathbf{E}_\psi[\ell_c(\boldsymbol{\psi})|\mathbf{Y} = \mathbf{y}]$ rozdělit na dvě části. Maximalizací $S(\mathbf{w}, \boldsymbol{\psi})$ získáme nový odhad \mathbf{w} a maximalizací $T(\boldsymbol{\xi}, \boldsymbol{\psi})$ získáme nový odhad $\boldsymbol{\xi}$.

Pro hledání \hat{w}_k použijeme metodu Lagrangeových multiplikátorů s vazebnou podmínkou

$$\sum_{k=1}^g w_k = 1.$$

Řešíme rovnici

$$\frac{\partial}{\partial w_k} \left(\sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \log w_k + \lambda \left(\sum_{k=1}^g w_k - 1 \right) \right) = 0, \quad (3.1)$$

kde $\lambda \in \mathbb{R}$ je Lagrangeův multiplikátor. Z (3.1) plyne, že

$$\sum_{j=1}^n \frac{\tau_k(\mathbf{y}_j; \boldsymbol{\psi})}{\hat{w}_k} + \lambda = 0,$$

tedy

$$\sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) + \lambda \hat{w}_k = 0$$

a neboť

$$\sum_{k=1}^g \hat{w}_k = 1$$

a

$$\sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) = n,$$

tak

$$n + \lambda * 1 = 0 \implies \lambda = -n.$$

Z čehož dostáváme

$$\hat{w}_k = \frac{1}{n} \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}), \quad k = 1, \dots, g.$$

Pro hledání $\hat{\boldsymbol{\mu}}_k$ zderivujeme $T(\boldsymbol{\xi}, \boldsymbol{\psi})$ podle $\boldsymbol{\mu}_k$, položíme rovno nule a vyřešíme

$$\sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \boldsymbol{\Sigma}^{-1} (\mathbf{y}_j - \hat{\boldsymbol{\mu}}_k) = 0, \quad k = 1, \dots, g,$$

a protože je $\boldsymbol{\Sigma}$ regulární, upravíme rovnici do tvaru

$$\sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \mathbf{y}_j = \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \hat{\boldsymbol{\mu}}_k.$$

Z čehož dostáváme

$$\hat{\boldsymbol{\mu}}_k = \frac{\sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \mathbf{y}_j}{\sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi})}, \quad k = 1, \dots, g.$$

Obdobně pro hledání $\hat{\boldsymbol{\Sigma}}$ zderivujeme $T(\boldsymbol{\xi}, \boldsymbol{\psi})$ podle $\boldsymbol{\Sigma}^{-1}$, položíme rovno nule a vyřešíme:

$$\sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) \left(\hat{\boldsymbol{\Sigma}} - (\mathbf{y}_j - \hat{\boldsymbol{\mu}}_k)(\mathbf{y}_j - \hat{\boldsymbol{\mu}}_k)^\top \right) = 0,$$

z toho

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}) (\mathbf{y}_j - \hat{\boldsymbol{\mu}}_k)(\mathbf{y}_j - \hat{\boldsymbol{\mu}}_k)^\top.$$

□

3.3 EM algoritmus

Nyní si představíme samotný EM algoritmus a jeho fungování na základě informací uvedených výše.

Začneme s počátečním odhadem parametrů $\boldsymbol{\psi}^{(0)}$, který volíme libovolně z parametrického prostoru $\boldsymbol{\Psi}$. Provedeme jeden iterační krok EM algoritmu (E-krok a M-krok) a dostaneme $\boldsymbol{\psi}^{(1)} \in \boldsymbol{\Psi}$ takové, že $L(\boldsymbol{\psi}^{(0)}) \leq L(\boldsymbol{\psi}^{(1)})$. Postupujeme dále až do té doby, než dostaneme $\boldsymbol{\psi}^{(m)}$ takové, že

$$|L(\boldsymbol{\psi}^{(m)}) - L(\boldsymbol{\psi}^{(m-1)})| \leq \alpha \tag{3.2}$$

pro $\alpha \in (0, \infty)$ předem dané dostatečně malé.

Definujeme funkci $Q(\boldsymbol{\psi}, \boldsymbol{\psi}^{(p)})$ jako

$$Q(\boldsymbol{\psi}, \boldsymbol{\psi}^{(p)}) = \mathbb{E}_{\boldsymbol{\psi}^{(p)}}[\ell_c(\boldsymbol{\psi})|\mathbf{y}].$$

V E-kroku hodnotu této funkce spočteme a v M-kroku maximalizujeme $Q(\boldsymbol{\psi}, \boldsymbol{\psi}^{(p)})$ přes $\boldsymbol{\Psi}$.

Konkrétně je hodnota funkce $Q(\boldsymbol{\psi}, \boldsymbol{\psi}^{(p)})$ je dle druhé části lemmatu 4 rovna

$$\sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}^{(p)}) \left(\log w_k - \frac{1}{2} \log (2\pi)^p - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{y}_j - \boldsymbol{\mu}_k)^\top (\boldsymbol{\Sigma})^{-1} (\mathbf{y}_j - \boldsymbol{\mu}_k) \right)$$

V M-kroku chceme najít řešení rovnice $\boldsymbol{\psi}^{(p+1)} = \arg \max_{\boldsymbol{\psi} \in \boldsymbol{\Psi}} Q(\boldsymbol{\psi}, \boldsymbol{\psi}^{(p)})$.

Dle třetí části lemmatu 4 jsou odhady vah rovny

$$w_k^{(p+1)} = \frac{1}{n} \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}^{(p)}), \quad k = 1, \dots, g.$$

To lze interpretovat tak, že každé pozorování přispívá k odhadu w_k takovou vahou, jaká je jeho pravděpodobnost, že do skupiny k patří.

Dále odhady průměrů jsou rovny

$$\boldsymbol{\mu}_k^{(p+1)} = \frac{\sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}^{(p)}) \mathbf{y}_j}{\sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}^{(p)})}, \quad k = 1, \dots, g,$$

což je vážený průměr pozorování příslušících ke skupině k .

A nakonec odhad varianční matice je

$$\boldsymbol{\Sigma}^{(p+1)} = \frac{1}{n} \sum_{k=1}^g \sum_{j=1}^n \tau_k(\mathbf{y}_j; \boldsymbol{\psi}^{(p)}) (\mathbf{y}_j - \boldsymbol{\mu}_k^{(p+1)}) (\mathbf{y}_j - \boldsymbol{\mu}_k^{(p+1)})^\top,$$

což je analogií vážené výběrové varianční matice.

Za zmínku stojí, že pro odhady vah $w_k^{(p+1)}$ bychom se k tomuto výsledku dostali i kdybychom uvažovali obecnější případ, tedy libovolnou konečnou směs. Naopak ale pro odhady středních hodnot $\boldsymbol{\mu}_k^{(p+1)}$ a varianční matice $\boldsymbol{\Sigma}^{(p+1)}$ bychom v obecném případě konečných směsí nijak derivaci funkce $T(\boldsymbol{\xi}, \boldsymbol{\psi}^{(p)})$ upravit nemohli, a tedy ani dostat explicitní odhady. Případ normálních směsí je jedním z mála, kde zvládneme explicitně určit odhady všech neznámých parametrů.

4. Praktické použití

V poslední části práce aplikujeme teorii popsanou v prvních třech kapitolách na klasifikaci reálných dat.

Máme k dispozici archeologické popisy lebek z článku Howells (Howells, 1996). Konkrétně se jedná o 156 ženských lebek, u kterých jsou známy jejich největší délky mozkovny v milimetrech (*gol*) a týlní úhly ve stupních (*oca*). Dále máme informaci o místě jejich nálezu: Austrálie, Berg v Rakousku nebo Burjati na Sibiři.

Pro ilustraci klasifikace zapomeneme, že známe rozdělení lebek dle místa nálezu a budeme se je snažit správně zařadit do skupin pouze na základě dat o největších délkách mozkovny a týlních úhlech. Následně porovnáme naše výsledky se skutečností.

4.1 Popis dat

Nejprve se podíváme na základní popisné statistiky proměnných *gol* a *oca* jak pro skupinu všech lebek, tak pro lebky v jednotlivých skupinách. Tyto statistiky jsou shrnuty v tabulkách 4.2 a 4.1.

Pro lepší přehlednost rozdílů mezi skupinami jsme vytvořili boxploty pro proměnné *gol* a *oca*, které jsou prezentovány v grafech 4.1 a 4.2.

Jak můžeme vidět, skupina Austrálie má výrazně větší hodnoty největší délky mozkovny než skupina Berg a Burjati, které mají hodnoty podobné. U týlního úhlu pozorujeme menší rozdíly mezi skupinami, přičemž skupina Austrálie vykazuje nižší hodnoty než skupiny Berg a Burjati.

Můžeme tedy předpokládat, že se nám nejlépe podaří oddělit skupinu Austrálie od zbytku, ale mezi skupiny Berg a Burjati se nám správně klasifikovat nejspíš nepovede.

V neposlední řadě si vykreslíme graf 4.3 ukazující vztah mezi *gol* a *oca* a barevně odděluje příslušníky jednotlivých skupin.

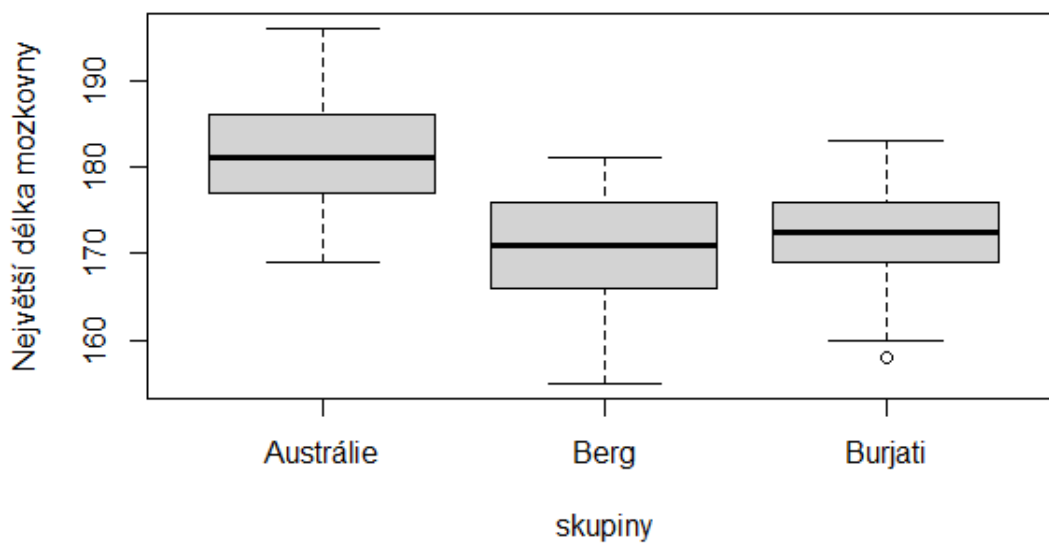
Na základě směrodatných odchylek a kovariancí mezi proměnnými v jednotlivých skupinách můžeme předpokládat homoskedasticitu.

	kovariance
všichni	-19,95
Austrálie	-12,76
Berg	-22,44
Burjati	-10,07

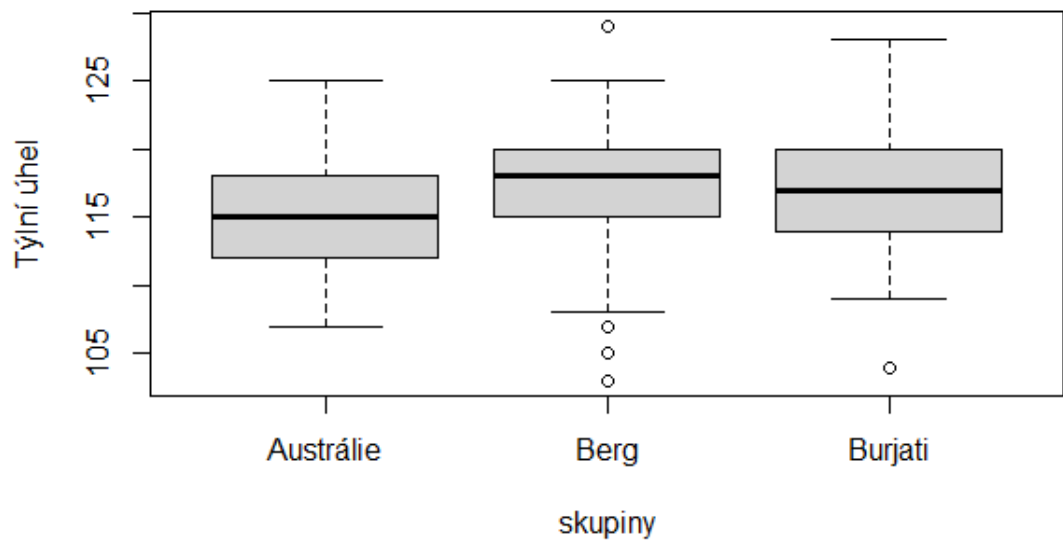
Tabulka 4.1: Kovariance mezi *gol* a *oca*.

	minimum	medián	výb. průměr	maximum	směr. odchylka
Největší délka mozkovny [mm]					
všichni	155	175,0	174,3	196	7,64
Austrálie	169	181,0	181,1	196	6,36
Berg	155	171,0	170,5	181	6,51
Burjati	158	172,5	171,8	183	5,43
Týlní úhel [stupně]					
všichni	103	117,0	116,3	129	4,89
Austrálie	107	115,0	114,7	125	4,06
Berg	103	118,0	117,0	129	5,56
Burjati	104	117,0	117,0	128	4,61

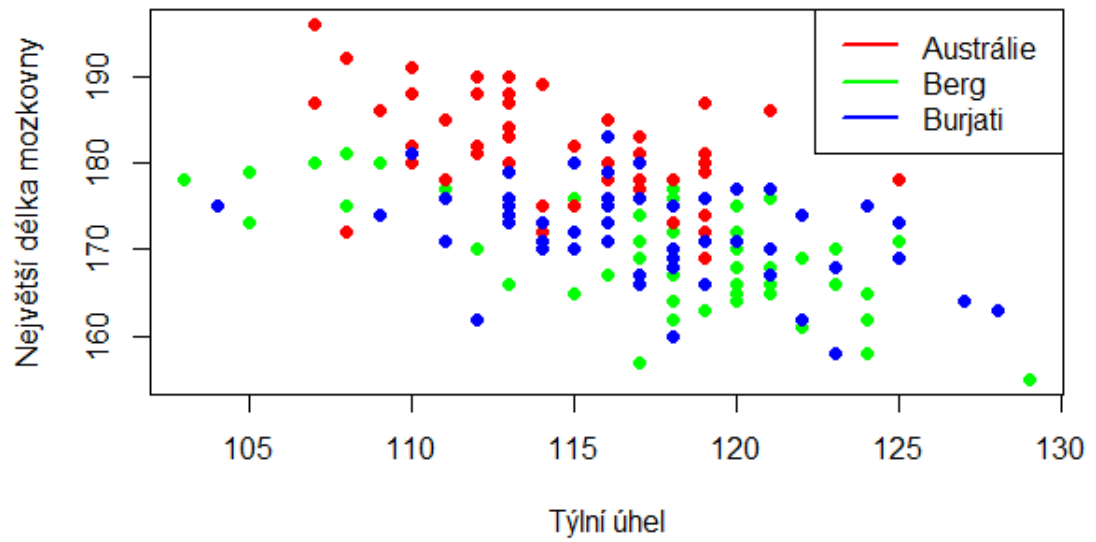
Tabulka 4.2: Základní popisné statistiky *gol* a *oca*.



Obrázek 4.1: Největší délka mozkovny v jednotlivých skupinách.



Obrázek 4.2: Týlní úhel v jednotlivých skupinách.



Obrázek 4.3: Vztah mezi *gol* a *oca*.

4.2 Klasifikace

Nyní zapomeneme, že skutečnou příslušnost lebek do skupin známe a pokusíme se ji určit pomocí klasifikační metody popsané v prvních třech kapitolách.

Nejprve potřebujeme zvolit počáteční hodnoty neznámých parametrů. Za počáteční odhady vah vezmeme pro všechny skupiny stejnou hodnotu, tedy

$$w_1^{(0)} = w_2^{(0)} = w_3^{(0)} = \frac{1}{3}.$$

Jak jsme uvedli v kapitole 4.1, u proměnné *gol* se pohybujeme mezi 155mm a 196mm a u proměnné *oca* mezi 103° a 129°. Tyto rozsahy rovnoměrně rozdělíme a zaokrouhlíme na celá čísla, tím získáme hodnoty

$$\boldsymbol{\mu}_1^{(0)} = \begin{pmatrix} 186 \\ 123 \end{pmatrix}, \quad \boldsymbol{\mu}_2^{(0)} = \begin{pmatrix} 176 \\ 116 \end{pmatrix}, \quad \boldsymbol{\mu}_3^{(0)} = \begin{pmatrix} 165 \\ 109 \end{pmatrix},$$

jako počáteční odhady středních hodnot. Konečně za počáteční hodnotu varianční matice vezmeme

$$\boldsymbol{\Sigma}^{(0)} = \begin{pmatrix} 58,37 & -19,95 \\ -19,95 & 23,91 \end{pmatrix},$$

kteřou jsme sestavili pomocí kvadrátů směrodatných odchylek a kovariance pro všechna pozorování.

Pro tyto hodnoty spustíme EM algoritmus s $\alpha = 10^{-6}$ z (3.2) a obdržíme konečné odhady neznámých parametrů.

Tabulka 4.3 porovnává odhady se skutečnými hodnotami neznámých parametrů, kde odhad skutečné varianční matice jsme spočítali následovně:

$$\frac{1}{n-3} \left(\sum_{i=1}^{n_1} (\mathbf{Y}_{1i} - \boldsymbol{\mu}_1)(\mathbf{Y}_{1i} - \boldsymbol{\mu}_1)^\top + \sum_{i=1}^{n_2} (\mathbf{Y}_{2i} - \boldsymbol{\mu}_2)(\mathbf{Y}_{2i} - \boldsymbol{\mu}_2)^\top + \sum_{i=1}^{n_3} (\mathbf{Y}_{3i} - \boldsymbol{\mu}_3)(\mathbf{Y}_{3i} - \boldsymbol{\mu}_3)^\top \right),$$

kde \mathbf{Y}_{ki} je i -té pozorování příslušící k -té skupině a n_k je počet pozorování v k -té skupině, pro $k = 1, 2, 3$.

parametr	odhad	skutečnost
w_1	0,75	0,31
w_2	0,08	0,34
w_3	0,17	0,35
$\boldsymbol{\mu}_1$	$\begin{pmatrix} 174,63 \\ 116,02 \end{pmatrix}$	$\begin{pmatrix} 181,10 \\ 114,70 \end{pmatrix}$
$\boldsymbol{\mu}_2$	$\begin{pmatrix} 173,97 \\ 108,43 \end{pmatrix}$	$\begin{pmatrix} 170,50 \\ 117,00 \end{pmatrix}$
$\boldsymbol{\mu}_3$	$\begin{pmatrix} 173,00 \\ 121,21 \end{pmatrix}$	$\begin{pmatrix} 171,80 \\ 117,00 \end{pmatrix}$
$\boldsymbol{\Sigma}$	$\begin{pmatrix} 57,69 & -18,90 \\ -18,90 & 14,58 \end{pmatrix}$	$\begin{pmatrix} 37,32 & -15,12 \\ -15,12 & 23,05 \end{pmatrix}$

Tabulka 4.3: Porovnání odhadů se skutečností.

Nakonec se podívejme na to, jak se nám podařilo zařadit jednotlivá pozorování do skupin.

Vzhledem k tomu, že metoda je schopna pouze rozdělit populaci na tři skupiny a nikoli je rozřadit do „pravých“ skupin, zvolili jsme pojmenování jednotlivých skupin jako takovou kombinaci odhadu a skutečnosti, která maximalizuje počet správně zařazených.

Tabulce 4.4 zachycuje kolik lebek se nám podařilo zařadit správně. Z celkových 156 lebek byla metoda úspěšná u 58 z nich.

Odhadnutá skupina \ Skutečná skupina	Austrálie	Berg	Burjati
	Austrálie	44	42
Berg	1	6	3
Burjati	4	5	8

Tabulka 4.4: Porovnání zařazení se skutečností.

Závěr

V této práci jsme se věnovali tomu, jak lze konečné směsi využít ke klasifikaci populace do skupin. Důraz jsme kladli na odhadování neznámých parametrů, které jsou ke klasifikaci stěžejní. Proto jsme popsali metodu maximální věrohodnosti, a to i s problémy, které zde mohou nastat. Zjistili jsme také, že hledat maximálně věrohodný odhad pro směšové modely není vůbec snadné, a proto jsme použili EM algoritmus, který je zde popsán. Zde by se mohla práce dále rozvinout o určité vlastnosti tohoto algoritmu, které by nám objasnily, proč je tak výhodné ho v této situaci použít.

S pomocí EM algoritmu jsme se snažili roztrždit lebky na základě jejich největší délky mozkovny a týlních úhlů do správných skupin podle místa nálezu. Výsledky ukázaly, že tato metoda je velmi přínosná pro klasifikaci určitého typu dat.

Věřím, že tato práce čtenáře obohatila a usnadnila mu vhléd do problematiky využití konečných směsí ke klasifikaci.

Seznam použité literatury

DEMPSTER, A. P., LAIRD, N. M. a RUBIN, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, **39**(1), 1–38.

HOWELLS, W. W. (1996). Howells' craniometric data on the internet. *American journal od physical anthropology*, **101**, 441–442.

MC LACHLAN, G. J. a PEEL, D. (2000). *Finite mixture models*. Wiley Series in Probability and Statistics, New York.