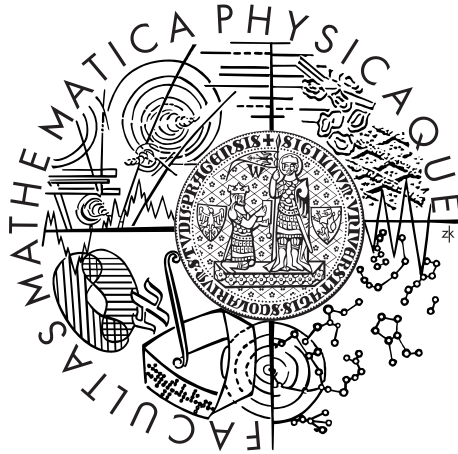


Univerzita Karlova v Praze
Matematicko-fyzikální fakulta

BAKALÁŘSKÁ PRÁCE



Tomáš Rusý

Modely konečných směsí

Katedra pravděpodobnosti a matematické statistiky

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

Studijní program: Matematika

Studijní obor: Obecná matematika

Praha 2015

Poděkování

Rád bych na tomto místě poděkoval mému panu vedoucímu doc. RNDr. Arnoštu Komárkovi, PhD. za cenné rady, věcné připomínky a drobné korektury, které významně pozvedly úroveň této bakalářské práce. Děkuji také panu RNDr. Jaromíru Běláčkovi, CSc. za poskytnutí dat k ilustraci modelu směsi v mé bakalářské práci a za uvedení do problematiky, jichž se data týkají.

Prohlašuji, že jsem tuto bakalářskou práci vypracoval 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 dne

Podpis autora

Název práce: Modely konečných směsí

Autor: Tomáš Rusý

Katedra pravděpodobnosti a matematické statistiky

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

Abstrakt: Tato práce se věnuje konečným směsím a klade si za cíl představit čtenáři použití metody maximální věrohodnosti pro odhad modelu směsi. K tomu využije EM algoritmus, který je v práci detailně zaveden a odvozen. Pro obecný model směsi popíše oba dva kroky algoritmu, E-kroku a M-kroku, a také výpočet nových odhadů části parametrů reprezentující model. Pro specifickou rodinu normálních směsí jsou pak odvozeny explicitní vzorce pro iteraci EM algoritmu. Popsaná teorie je pak aplikována na modelu produkce cytokinu interleukin 10 při napadení člověka paradontózou, díky kterému je názorně ukázána využitelnost modelu v praxi. Na závěr je pak na odhad modelu směsi navázáno teorií shlukování, zabývající se rozřazování prvků do skupin. Podobně jako předchozí teorie je také ilustrována na výše zmíněném modelu produkce cytokinu IL10.

Klíčová slova: konečná směs, normální směs, EM algoritmus, shlukování

Title: Finite Mixture Models

Author: Tomáš Rusý

Department of Probability and Mathematical Statistics

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

Abstract: This work focuses on finite mixture models and aims to introduce the maximum likelihood method as an approach of fitting finite mixtures. For that purpose the EM algorithm is adopted and derived in detail. Both the E and M steps of the EM algorithm are presented and performed for general finite mixture model. We derive new estimates for some of parameters defining the model. All updated estimates of the iteration of the EM algorithm are derived explicitly for the specific family of normal mixtures. The described theory is then applied to a model of the production of cytokine interleukin 10 in human periodontitis attack, which clearly demonstrates an application of the model in practice. Finally, we discuss the theory of clustering, which is based on our previous results. Like the previous theory, this one is also illustrated in the aforementioned model of cytokine IL10 production.

Keywords: finite mixture, normal mixture, EM algorithm, clustering

Obsah

Použité značení	2
Úvod	3
Motivace	4
1 Model konečné směsi	5
1.1 Konečná směs	5
1.2 Interpretace směsi	6
1.3 Parametrizace směsi	6
1.4 Příklad onemocnění paradontózou	9
2 Odhad metodou maximální věrohodnosti	11
2.1 Věrohodnostní funkce	11
2.2 Věrohodnostní rovnice	14
2.3 Identifikovatelnost rozdělení směsi	15
2.4 EM algoritmus	16
2.4.1 E-krok	19
2.4.2 M-krok	19
3 Normální směsi	22
3.1 EM algoritmus pro normální směsi	22
3.2 Odhad modelu produkce cytokinu IL10	25
4 Shlukování	28
4.1 O nezdravých hodnotách zdravých lidí	29
Závěr	32
Literatura	33
Seznam obrázků	34
Seznam tabulek	35

Použité značení

\mathbb{N}	množina přirozených čísel
\mathbb{R}	množina reálných čísel
\mathbb{R}^p	vektorový prostor dimenze p nad reálnými čísly
\mathbf{u}	sloupcový vektor \mathbf{u}
\mathbf{A}^\top	transponovaná matice \mathbf{A}
$\mathbf{A} \otimes \mathbf{B}$	Kroneckerův součin matic \mathbf{A} a \mathbf{B}
$\text{int}(\Psi)$	vnitřek množiny Ψ
$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	normální rozdělení se střední hodnotou $\boldsymbol{\mu}$ a rozptylovou maticí $\boldsymbol{\Sigma}$
$\varphi(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	hustota normálního rozdělení s vektorem středních hodnot $\boldsymbol{\mu}$ a rozptylovou maticí $\boldsymbol{\Sigma}$.
$P(\cdot \mathbf{y})$	$P(\cdot \mathbf{Y} = \mathbf{y})$
$E[\cdot \mathbf{y}]$	$E[\cdot \mathbf{Y} = \mathbf{y}]$
$E[g(\mathbf{Y}) \boldsymbol{\psi}']$	střední hodnota z funkce $g(\mathbf{Y})$, když hustota náhodného vektoru \mathbf{Y} je určena hodnotou parametrického vektoru $\boldsymbol{\psi}'$
MMV	metoda maximální věrohodnosti
IL10	cytokin interleukin 10
IL10_Pi	logaritmizované hodnoty produkce cytokinu IL10, při stimulování kulturou bakterie <i>prevotella intermedia</i>
IL10_nestim	logaritmizované hodnoty produkce cytokinu IL10, nestimulované

Úvod

Konečné směsi představují stále se vyvíjející a často používanou metodu statistické analýzy, jež našla širšího uplatnění zejména díky vývoji výpočetní techniky v minulém století. Úlohy, jejichž řešení bylo dříve nemyslitelné, se se vznikem počítačů staly překonatelnou překážkou, nad kterou bylo možno vystavit teorii zabývající se touto metodou. Za cíl této bakalářské práce si klademe model konečné směsi názorně představit a ilustrovat jeho použití na vhodném praktickém případě.

V úvodní části je představen cíl a struktura práce a zmíníme se i o prvním publikovaném použití modelu směsi ve statistické analýze. Na to navážeme v první kapitole, která je věnována představení modelu směsi. Její součástí je jak rigorózní definice objektu, kterým model směsi nazýváme, tak i abstraktní a praktická ilustrace, jež usnadňují pochopení modelu. Mimo jiné se v této části zmíníme i o reálné interpretaci tohoto teoretického pojmu a situacím, kdy na něj můžeme v obyčejném světě narazit.

Protože ve statistické analýze je zásadním problémem uvažované modely odhadnout a zkoumat jejich vlastnosti, i v našem textu se podobné technice budeme věnovat. Konkrétně se ve druhé kapitole budeme zabývat využitím metody maximální věrohodnosti jako cesty, kterou lze obecný model směsi odhadnout. Zjistíme, že samotnému nalezení odhadu předchází řada komplikací, které se budeme snažit vyřešit. K nalezení odhadů pak budeme používat EM algoritmus, který pro směsi podrobně zavedeme a detailně si vysvětlíme princip, na kterém je založen. Na závěr této kapitoly pak odvodíme v maximální obecnosti vzorce pro část parametrů modelu směsi jedné iterace EM algoritmu.

Odvodit i zbylé vzorce pro ostatní parametry modelu se ukazuje jako nemožné, a proto se ve třetí kapitole omezíme pouze na normální směsi, tj. směsi, jejichž složky mají normální rozdělení. Ukážeme si, že pro tuto specifickou rodinu lze vzorce pro nové odhady parametrů EM algoritmu vyjádřit explicitně, což značně zjednodušuje jeho použití. Tím zakončíme teoretickou část práce, věnující se odhadování pomocí metody maximální věrohodnosti a ukážeme si použití námi odvozené techniky v praxi. K tomu nám poslouží model, jež je představen v první kapitole, zabývající se produkcí obrané látky interleukin 10 při napadení člověka onemocněním zubů - paradontózou.

Nakonec se budeme ve čtvrté kapitole věnovat odpovědi na jednu přirozenou otázku, která v kontextu naší interpretace směsi vyvstává. Tou je, že v případě, kdy se nějaká vlastnost chová jinak ve více skupinách lidí, zda-li nemůžeme

na základě pozorování té vlastnosti u nějakého člověka určit, ze které skupiny pochází. Odborně této technice říkáme shlukování a v rámci kapitoly si představíme jedno pravidlo, kterým se v tomto případě můžeme řídit. Opět pro lepší pochopení ukážeme použití představené látky na modelové situaci uvedené v první kapitole.

Věřím, že celkově bakalářská práce představuje podrobný popis problematiky konečných směsí a komplexní odvození velké většiny vyvozených závěrů. Díky tomu by měla být pochopitelná i širšímu matematickému publiku, jelikož staví pouze na základech teorie pravděpodobnosti.

Motivace

První zmínka o problematice konečných směsí se objevila na konci 19-tého století v práci anglického matematika a biostatistika Karla Pearsona. Pearson (1894) se v ní zabýval daty populace krabů v zálivu u města Neapol. Ty mu byly dodány od biologa Weldona, který analyzoval poměr velikosti jejich čela ku délce těla. Na základě toho pak vyslovil hypotézu, že naměřená data mohou být signálem vzniku nového druhu krabů v zálivu. Tato domněnka byla podpořena faktem, že rozdělení testovaných dat neodpovídalo normálnímu rozdělení.

Pearson se snažil odhadnout toto rozdělení směsí dvou normálních rozdělení momentovou metodou, což při pěti proměnných vedlo na polynomiální rovnici devátého stupně. A ačkoliv se mu tuto úlohu nakonec podařilo vyřešit, následovníků, kteří se vydali v jeho stopách, mnoho nebylo. To proto, že řešení podobných či složitějších rovnic bylo v té době velmi obtížné, pracné, případně nemožné.

Díky pokroku techniky v posledních desetiletích se ale řešení takovýchto úloh stalo snadným, což zapříčinilo značný rozvoj teorie odhadování modelů konečných směsí. Ta posléze našla uplatnění v mnoha odvětvích. Problematika shlukování a vlastnost normálních konečných směsí, že velmi dobře dokáží aproximovat širokou třídu rozdělení, patří mezi ty nejužívanější. Zejména ta druhá je ale už natolik pokročilá, že se jí v této práci věnovat nebudeme.

Kapitola 1

Model konečné směsi

V této kapitole si představíme, co pod konečnou směsí rozumíme a to včetně parametrického vyjádření, které budeme používat v dalších částech práce. Krom toho se navíc zmíníme o interpretaci, kterou tento abstraktní pojem má, a jak je možno na něj nahlížet. Pro lepší představu pak uvedeme jeden praktický příklad, kde náhodný vektor vystihující diskutovaný problém bude mít rozdělení, které je možno modelem směsi popsat.

1.1 Konečná směs

Definice 1 (Konečná směs). *Nechť \mathbf{Y} je p -rozměrný náhodný vektor, $g \in \mathbb{N}$, $\pi_j \in (0,1)$, $j \in \{1, \dots, g\}$, tak, že*

$$\sum_{j=1}^g \pi_j = 1$$

a necht' $f_j(\mathbf{y})$, $j \in \{1, \dots, g\}$, jsou hustoty¹ na \mathbb{R}^p . Necht' \mathbf{Y} má hustotu

$$f(\mathbf{y}) = \sum_{j=1}^g \pi_j f_j(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^p. \quad (1.1)$$

Potom řekneme, že rozdělení náhodného vektoru \mathbf{Y} je g -složková konečná směs. Hustotě (1.1) budeme říkat hustota směsi s g složkami.

Přirozené číslo g představuje počet složek ve směsi, a ačkoliv v definici připouštíme i $g = 1$, ve skutečnosti tento případ není v definici možný, z důvodu podmínek, které klademe na π_j . Nás beztak ale zajímá pouze netriviální případ kdy $g > 1$. Hodnoty π_j vyjadřují poměr zastoupení j -té složky ve směsi a uvažujeme je v otevřeném intervalu $(0,1)$. To proto, že pokud by pro nějaké j platilo $\pi_j = 1$, pak by šlo o 1-složkovou konečnou směs, což je případ, který nás nezajímá. Naopak, v případě kdy by $\pi_j = 0$, pak by se jednalo o $g - 1$ -složkovou směs, kterou popíšeme $g - 1$ složkami.

¹Hustoty v této práci uvažujeme vzhledem k p -rozměrné Lebesgueově míře.

1.2 Interpretace směsi

Nyní si uvedeme jednu z možných interpretací g -složkové konečné směsi. Uvažme náhodnou veličinu U nabývající hodnot $1, \dots, g$ s pravděpodobnostmi

$$P(U = j) = \pi_j, \quad j = 1, \dots, g.$$

Dále předpokládejme, že podmíněná hustota náhodného vektoru \mathbf{Y} při jevu $U = j$ je $f_j(\mathbf{y})$. Pak nepodmíněná hustota náhodného vektoru \mathbf{Y} je

$$f(\mathbf{y}) = \sum_{j=1}^g f(\mathbf{y}|U = j) P(U = j) = \sum_{j=1}^g \pi_j f_j(\mathbf{y}),$$

což je přesně hustota (1.1). Tedy rozdělení konečné směsi \mathbf{Y} , jenž je hustotou (1.1) reprezentováno, lze interpretovat jako chování sledované vlastnosti v g skupinách, kde v každé skupině se vlastnost chová jinak, neboli má jiné rozdělení dané hustotami f_1, \dots, f_g .

Podobně lze místo náhodné veličiny U uvážit g -rozměrný náhodný vektor $\mathbf{Z} = (Z_1, \dots, Z_g)^\top$, jehož složky nabývají pouze nul a jedné jedničky, přitom jedničky právě v té složce, jejíž index odpovídá skupině, ze které sledovaný prvek pochází. Tj.

$$\mathbf{Z} \sim \text{Mult}(1, \pi_1, \dots, \pi_g).$$

Předpokládejme, že pro podmíněnou hustotu \mathbf{Y} při $\mathbf{Z} = \mathbf{z}$ platí

$$f(\mathbf{y}|\mathbf{z}) = \sum_{j=1}^g z_j f_j(\mathbf{y}) = \prod_{j=1}^g (f_j(\mathbf{y}))^{z_j}.$$

Pak nepodmíněná hustota \mathbf{Y} , analogicky s předchozím příkladem, je dána hustotou (1.1).

Typicky ale náhodný vektor \mathbf{Z} , respektive hodnotu náhodné veličiny U , které popisují, z jaké skupiny sledované populace sledovaný objekt s pozorovanou hodnotou $\mathbf{Y} = \mathbf{y}$ pocházel, nepozorujeme. S využitím Bayesovy věty však zpětně můžeme počítat pravděpodobnost toho, že sledovaný objekt náležel do j -té skupiny, poté co byla pozorována hodnota $\mathbf{Y} = \mathbf{y}$.

$$\begin{aligned} p_j(\mathbf{y}) &= P(\text{objekt náležel do } j\text{-té složky}|\mathbf{y}) \\ &= P(Z_j = 1|\mathbf{y}) \\ &= \pi_j f_j(\mathbf{y})/f(\mathbf{y}), \quad j = 1, \dots, g. \end{aligned} \tag{1.2}$$

1.3 Parametrizace směsi

Dále se budeme zabývat odhadováním parametrů modelu konečné směsi, a k tomu budeme potřebovat hustotu směsi v řeči nějakých parametrů vyjádřit. Ačkoliv je možné, aby se sledovaná vlastnost chovala v různých skupinách úplně odlišně, my se v této práci omezíme jen na konečné směsi, jejíž složky pocházejí ze stejné parametrické rodiny.

Definice 2 (Parametrická směs). *Nechť \mathbf{Y} je p -rozměrný náhodný vektor, $g \in \mathbb{N}$, $\pi_j \in (0,1)$, $j \in \{1, \dots, g\}$, tak, že*

$$\sum_{j=1}^g \pi_j = 1$$

a nechť $f(\mathbf{y}; \boldsymbol{\theta}_j)$, $j \in \{1, \dots, g\}$, jsou hustoty na \mathbb{R}^p s neznámými parametry obsaženými ve vektoru $\boldsymbol{\theta}_j \in \Theta_j \subset \mathbb{R}^{p_j}$. Nechť \mathbf{Y} má hustotu

$$f(\mathbf{y}; \boldsymbol{\psi}) = \sum_{j=1}^g \pi_j f(\mathbf{y}; \boldsymbol{\theta}_j). \quad (1.3)$$

Potom řekneme, že rozdělení náhodného vektoru \mathbf{Y} je g -složková parametrická směs. Hustotě (1.3) budeme říkat parametrická hustota směsi s g složkami.

Vektor $\boldsymbol{\psi}$ nechť označuje vektor obsahující všechny neznámé parametry modelu. Symbolem Ψ budeme značit parametrický prostor pro $\boldsymbol{\psi}$. Označíme $\boldsymbol{\xi}$ vektor obsahující všechny parametry v $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_g$, o kterých připouštíme, že mohou být různé. Pak lze psát

$$\boldsymbol{\psi} = (\pi_1, \dots, \pi_{g-1}, \boldsymbol{\xi}^\top)^\top.$$

V zápise je záměrně opomenut poměr zastoupení π_g směsi g . To proto, že součet π_j je jedna, a tudíž je jeden z nich nadbytečný.

Značení. Hustotu p -dimenzionálního normálního náhodného vektoru \mathbf{Y} se střední hodnotou $\boldsymbol{\mu}$ a kovarianční maticí $\boldsymbol{\Sigma}$ budeme značit $\varphi(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$. Analogicky symbolem $\varphi(y; \mu, \sigma^2)$ budeme značit hustotu jednorozměrného normálního rozdělení se střední hodnotou μ a rozptylem σ^2 .

Definice 3 (Normální směs). *Nechť \mathbf{Y} je p -rozměrný náhodný vektor, $g \in \mathbb{N}$, $\pi_j \in (0,1)$, $j \in \{1, \dots, g\}$, tak, že*

$$\sum_{j=1}^g \pi_j = 1$$

a nechť $\varphi(\mathbf{y}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, $j \in \{1, \dots, g\}$, jsou hustoty p -rozměrného normálního rozdělení s parametry $\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j$. Nechť \mathbf{Y} má hustotu

$$f(\mathbf{y}; \boldsymbol{\psi}) = \sum_{j=1}^g \pi_j \varphi(\mathbf{y}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j). \quad (1.4)$$

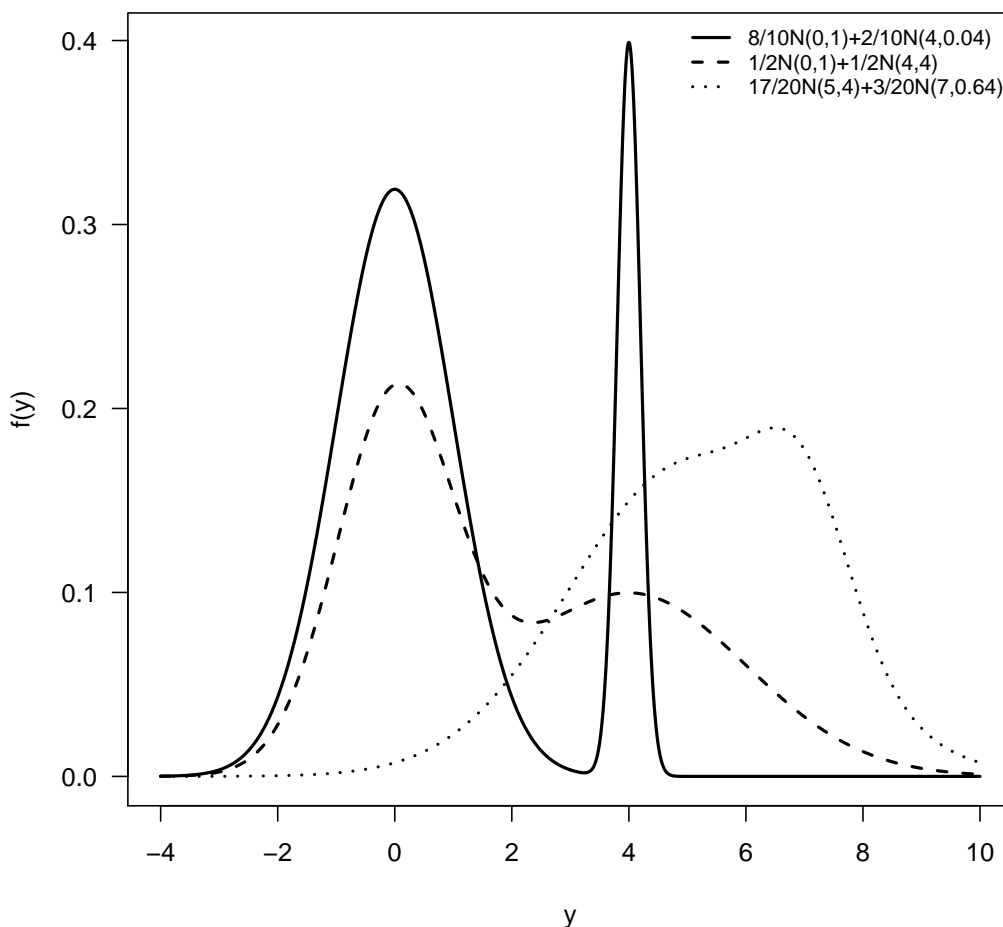
Potom řekneme, že rozdělení náhodného vektoru \mathbf{Y} je g -složková normální směs. Hustotě (1.4) budeme říkat parametrická hustota normální směsi s g složkami. V případě, kdy $p = 1$ budeme mluvit o g -složkové jednorozměrné normální směsi.

UVědomíme si, že pokud použijeme pro normální směs zápis 1.3, pak prvky parametrického vektoru $\boldsymbol{\theta}_j$ jsou všechny parametry normálního rozdělení, u kterých připouštíme, že jsou různé. Typicky se bude jednat zejména o následující dva typy.

Příklad. Heteroskedastickou normální směsí myslíme směs, ve které připouštíme, že rozptylové matice složek mohou být různé. Pak v θ_j jsou obsaženy parametry vektoru středních hodnot μ_j a prvky matice Σ_j obsažené v jejím horním trojúhelníku. Parametrický vektor ψ pak čítá všechny tyto parametry dohromady pro $j = 1, \dots, g$ společně s poměry složek π_j .

Příklad. Homoskedastická normální směs je taková, ve které naopak požadujeme, aby rozptylové matice složek byly shodné. Potom v θ_j jsou společně s prvky parametry vektoru středních hodnot μ_j i prvky určující společnou rozptylovou matici Σ . Vektor ψ pak obsahuje krom prvků μ_j a poměrů složek $\pi_j, j = 1, \dots, g$ i prvky matice Σ , které jsou v jejím horním trojúhelníku, a to právě jednou.

Pouze poznamenejme, že pomocí normálních směsí lze kromě jiného aproximovat širokou škálu rozdělení, jak je empiricky ukázáno na obrázku 1.1, kde jsou navíc uvedeny pouze dvousložkové jednorozměrné normální směsi. Všechny se ale od sebe zdatelně liší. V legendě obrázku symbolem $\pi N(\mu_1, \sigma_1^2) + (1 - \pi)N(\mu_2, \sigma_2^2)$ myslíme směs, kde jsou normální rozdělení definovaná podle příslušných parametrů μ_1, σ_1^2 resp. μ_2, σ_2^2 s poměry zastoupení složek π a $1 - \pi$.



Obrázek 1.1: Příklady několika normálních směsí.

1.4 Příklad onemocnění paradontózou

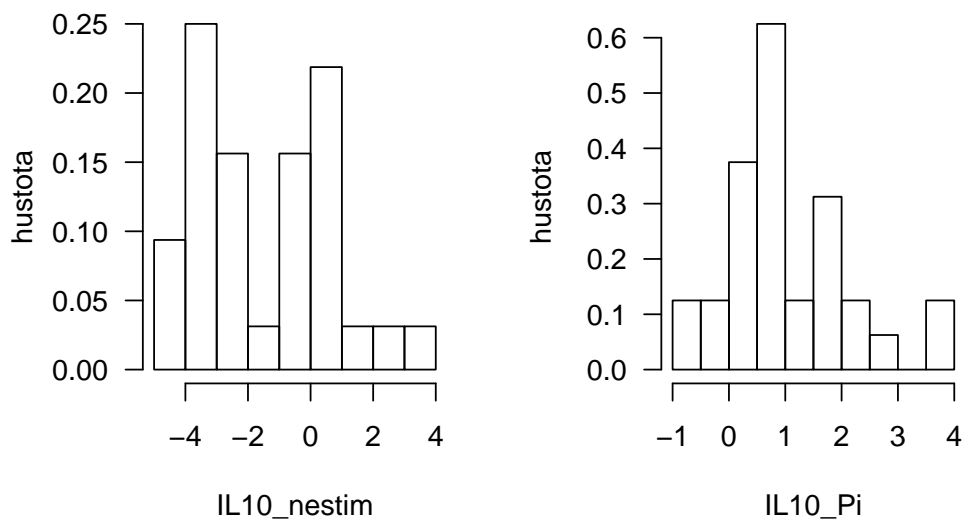
Pro lepší pochopení modelu směsi si zde uvedeme následující příklad, na který budeme v dalších částech práce navazovat. Uvažme populaci lidí, ve které je rozšířena paradontóza, odborně paradontitida. Toto mezi lidmi relativně běžné závažné onemocnění může skončit až ztrátou zubů. Jeho včasná detekce pak často zachrání mnoha lidem jak chrup, tak i velké množství peněz, které musí být na případnou náhradu vynaloženy.

Data diskutovaná v této bakalářské práci jsou výsledkem měření produkce cytokinu *interleukin 10* (dále IL10), což je protein, který se významně účastní v imunitní odpovědi těla na případný útok nějaké nemoci. Jeho zvýšená produkce tedy může být signálem nastupujícího onemocnění. Při měření pak kultura obsahující tento cytokin v prvním případě nebyla nijak stimulována a byla měřena jeho samotná ničím nezrychlená produkce a v druhém případě byla stimulována přítomností bakterie *prevotella intermedia*. To je jedna z těch, které se běžně vyskytují v dutině ústní, a má významný podíl na vzniku paradontitidy.

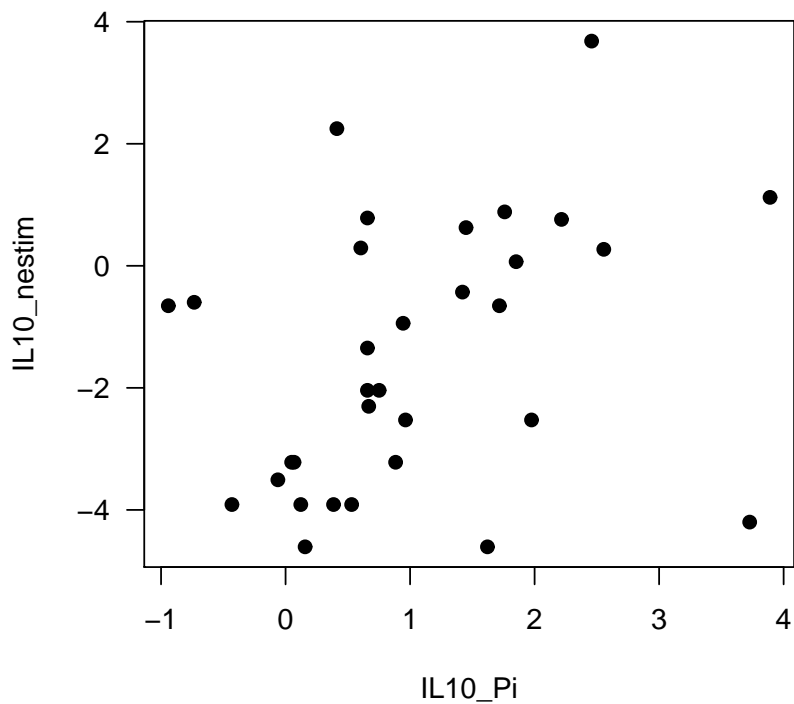
V rámci statistické analýzy nás bude zajímat zejména odpověď na otázku, zda-li je produkce cytokinu IL10 u pacientů s paradontózou vyšší než u lidí, kteří touto nemocí netrpí. Budeme předpokládat, že produkce se v každé z těchto dvou skupin řídí jiným rozdělením, které se pokusíme odhadnout, a na závěr pak budeme zkoumat, zda-li data odpovídají patřičné směsi.

Označíme \mathbf{X} dvourozměrný náhodný vektor, jehož první složka X_1 reprezentuje nestimulovanou produkci cytokinu IL10 a druhá složka X_2 produkci IL10 za přítomnosti bakterie *prevotella intermedia*. Pokud budeme předpokládat, že růstové kultury se v každé ze dvou uvažovaných skupin řídí log-normálním rozdělením, pak zřejmě náhodný vektor $\mathbf{Y} = \log \mathbf{X} = (\log X_1, \log X_2)^\top$, bude mít při platnosti naší hypotézy rozdělení dvousložkové dvourozměrné normální směsi.

Pro přehlednost a ilustraci si zde uvedeme ještě histogramy logaritmizovaných hodnot jak nestimulované, tak stimulované produkce cytokinu IL10, (viz obrázek 1.2). Popiskem *IL10_nestim* je označen histogram hodnot nestimulované produkce cytokinu IL10, podobně *IL10_Pi* jsou označeny hodnoty produkce IL10 při přítomnosti bakterie *prevotella intermedia*. Pro celkový náhled, jak vypadají naše data dvojrozměrně, ještě uvádíme bodový graf naměřených hodnot (viz obrázek 1.3). Dva shluky okolo bodů o souřadnicích $[0, -4]$ a $[2, 0]$ v tomto grafu naznačují, že rozdělení by mohlo odpovídat dvousložkové normální směsi. Otázkou zůstává, jaké parametry tato směs má. Jeden z možných způsobů jejich odhadnutí ukážeme v následující kapitole.



Obrázek 1.2: Histogramy hodnot logaritmizované produkce stimulované i nestimulované IL10.



Obrázek 1.3: Bodový záznam dat produkce cytokinu IL10.

Kapitola 2

Odhad metodou maximální věrohodnosti

Existuje několik možností, kterých lze použít pro odhadování parametrů modelu s rozdělením určeným hustotou (1.3). Zmínit můžeme například momentovou metodu, kterou použil Karl Pearson, případně třeba metodu nejmenší vzdálenosti. Těmi se ale v tomto textu zabývat nebudeme a místo toho se budeme soustředit na odhadování pomocí metody maximální věrohodnosti, která je obecně jednou z nejvyužívanějších metod odhadu. Ta spočívá v sestavení věrohodnostní funkce a její následné maximalizaci. V této kapitole mimo jiné rozebereme i úskalí, která na nás při použití této metody v kontextu směsí čekají.

Pro potřeby této práce budeme značit $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ náhodný výběr o rozsahu n ze spojitého rozdělení s pravděpodobnostní hustotou $f(\mathbf{y})$, $\mathbf{y} \in \mathbb{R}^p$. Pak tedy je \mathbf{Y}_i , $i \in \{1, \dots, n\}$, p -rozměrný náhodný vektor. Navíc označíme $\mathbf{Y}_a = (\mathbf{Y}_1^\top, \dots, \mathbf{Y}_n^\top)^\top$ náhodný vektor reprezentující celý výběr. \mathbf{Y}_a je tedy $(n \cdot p)$ -rozměrný náhodný vektor. Realizace náhodných vektorů budeme značit odpovídajícím malým písmenem. Tedy symbolem \mathbf{y}_1 budeme značit realizaci náhodného vektoru \mathbf{Y}_1 , podobně symbolem \mathbf{y}_a budeme značit realizaci celého náhodného výběru.

V kontextu konečných směsí a při zachování značení z kapitoly 1 označíme pro náhodný vektor \mathbf{Y}_i jeho příslušný náhodný vektor \mathbf{Z}_i určující, ze které skupiny měřený prvek pochází. Navíc ještě zavedeme značení pro následující náhodné vektory $\mathbf{Z}_a = (\mathbf{Z}_1^\top, \dots, \mathbf{Z}_n^\top)^\top$ a $\mathbf{Y}_c = (\mathbf{Y}_a^\top, \mathbf{Z}_a^\top)^\top$.

2.1 Věrohodnostní funkce

Pro náhodný výběr $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ z rozdělení s hustotou (1.3) sestavíme věrohodnostní funkci jakožto sdruženou hustotu tohoto náhodného výběru. Potom pro $\boldsymbol{\psi} \in \Psi$ má věrohodnostní funkce tvar

$$L(\boldsymbol{\psi}) = \prod_{i=1}^n f(\mathbf{Y}_i; \boldsymbol{\psi}) = \prod_{i=1}^n \left(\sum_{j=1}^g \pi_j f(\mathbf{Y}_i; \boldsymbol{\theta}_j) \right).$$

Analogicky můžeme sestavit i log-věrohodnostní funkci

$$\ell(\boldsymbol{\psi}) = \log L(\boldsymbol{\psi}) = \sum_{i=1}^n \log f(\mathbf{Y}_i; \boldsymbol{\psi}) = \sum_{i=1}^n \log \left(\sum_{j=1}^g \pi_j f(\mathbf{Y}_i; \boldsymbol{\theta}_j) \right). \quad (2.1)$$

Odhad metodou maximální věrohodnosti (dále budeme značit MMV) je potom prvek $\hat{\boldsymbol{\psi}}_{\text{MLE}} \in \Psi$ splňující

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

Je zřejmé, že aby tato rovnice měla řešení, musí být věrohodnostní funkce omezená. Tento požadavek ale, jak ukazuje následující příklad, není v případě směsi vždy splněn.

Příklad. Uvažme jednorozměrnou normální směs s dvěma komponentami. Nechť tedy

$$f(y; \boldsymbol{\psi}) = \pi_1 \varphi(y; \mu_1, \sigma_1^2) + (1 - \pi_1) \varphi(y; \mu_2, \sigma_2^2).$$

Pak $\boldsymbol{\psi} = (\pi_1, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)^\top$. Položme $\Psi = (0, 1) \times \mathbb{R} \times (0, \infty) \times \mathbb{R} \times (0, \infty)$. Nechť Y_1, \dots, Y_n je náhodný výběr z rozdělení s touto hustotou. Pak log-věrohodnostní funkce pro odhad parametrického vektoru $\boldsymbol{\psi}$ při realizovaném náhodném výběru $Y_1 = y_1, \dots, Y_n = y_n$ je tvaru

$$\ell(\boldsymbol{\psi}) = \sum_{i=1}^n \log \left(\pi_1 \varphi(y_i; \mu_1, \sigma_1^2) + (1 - \pi_1) \varphi(y_i; \mu_2, \sigma_2^2) \right). \quad (2.2)$$

Položme $\boldsymbol{\psi}_\delta = (\frac{1}{2}, y_1, \delta^2, 0, 1)^\top$, $\delta > 0$. Pak $\forall \delta > 0 : \boldsymbol{\psi}_\delta \in \Psi$. Pro toto $\boldsymbol{\psi}_\delta$ nyní vyjádříme log-věrohodnostní funkci (2.2)

$$\ell(\boldsymbol{\psi}_\delta) = \sum_{i=1}^n \log \left(\frac{1}{2} \frac{1}{\sqrt{2\pi}\delta} \exp \left(-\frac{(y_i - y_1)^2}{2\delta^2} \right) + \frac{1}{2} \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{y_i^2}{2} \right) \right). \quad (2.3)$$

Nyní hodnotu $\ell(\boldsymbol{\psi}_\delta)$ odhadneme zespona tak, že pro $i = 1$ vynecháme druhý člen součtu v argumentu logaritmu a pro ostatní i vynecháme první člen. Tím hodnotu na pravé straně (2.3) snížíme, protože oba sčítance jsou vždy kladné. Tedy

$$\begin{aligned} \ell(\boldsymbol{\psi}_\delta) &> \log \left(\frac{1}{2} \frac{1}{\sqrt{2\pi}\delta} \exp \left(-\frac{(y_1 - y_1)^2}{2\delta^2} \right) \right) + \sum_{i=2}^n \log \left(\frac{1}{2} \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{y_i^2}{2} \right) \right) \\ &> -\log \left(2\sqrt{2\pi}\delta \right) - (n-1) \log \left(2\sqrt{2\pi} \right) + \sum_{i=2}^n \frac{-y_i^2}{2} \\ &> -\log(\delta) - n \log \left(2\sqrt{2\pi} \right) - \sum_{i=2}^n \frac{y_i^2}{2}. \end{aligned} \quad (2.4)$$

Nyní je z (2.4) patrné, jak se $\ell(\boldsymbol{\psi}_\delta)$ chová když $\delta \rightarrow 0$. Jinými slovy

$$\lim_{\delta \rightarrow 0} \ell(\boldsymbol{\psi}_\delta) > \lim_{\delta \rightarrow 0} \left(-\log(\delta) - n \log \left(2\sqrt{2\pi} \right) - \sum_{i=2}^n \frac{y_i^2}{2} \right) = \infty.$$

Z toho vyplývá, že věrohodnostní funkce této směsi není omezená, a tedy neexistuje odhad metodou maximální věrohodnosti. Navíc tento postup lze zobecnit pro vícerozměrné, vícesložkové i mnohé další nenormální směsi prostým umístěním střední hodnoty jedné složky do jednoho pozorování a uvážením rozptylu blížícího se k nule, případně rozptylové matice k singulární. Maximálně věrohodné odhady tedy v podstatě odpovídají směsi s degenerovanými komponentami. □

Hustota složená z nekonečných výběžků v místech pozorovaných dat by ale nepředstavovala jejich žádně zjednodušení, které od modelu čekáme. Aitkin (2001, str. 289) navrhuje způsob, kterým lze problém neomezené věrohodnosti řešit. Stačilo by omezit parametrický prostor Ψ tak, aby byl uzavřený a omezený. Toho se dá docílit například určením dolní hranice δ pro směrodatné odchylky složek směsi. Nicméně takovéto opatření způsobí že výsledný odhad bude na δ závislý, což může vadit zejména v případě mnoha odlehlých pozorování.

Na druhou stranu v případě modelu jednorozměrné homoskedastické normální směsi, tj. směsi, jejichž všechny složky mají stejný rozptyl, lze ukázat, že věrohodnostní funkce, v případě dostatečně velkého náhodného výběru, omezená je. To bude předmětem následující věty.

Věta 1. *Nechť Y_1, \dots, Y_n je náhodný výběr obsahující alespoň $g+1$ různých pozorování z jednorozměrné normální směsi s hustotou*

$$f(y; \boldsymbol{\psi}) = \sum_{j=1}^g \pi_j \varphi(y; \mu_j, \sigma^2). \quad (2.5)$$

Pak je věrohodnostní funkce $L(\boldsymbol{\psi})$ shora omezená.

Důkaz. Uvažme maximální možný parametrický prostor který pro normální směr (2.5) připadá v úvahu, tj $\mu_j \in \mathbb{R}$, $\sigma^2 \in (0, \infty)$, $\pi_j \in (0, 1)$, $j = 1, \dots, g$. Pro pozorovaný náhodný výběr $Y_1 = y_1, \dots, Y_n = y_n$ vyjádříme a shora odhadneme věrohodnostní funkci náhodného výběru:

$$\begin{aligned} L(\boldsymbol{\psi}) &= \prod_{i=1}^n \left(\sum_{j=1}^g \pi_j \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2\sigma^2}(y_i - \mu_j)^2\right) \right) \\ &\leq \prod_{i=1}^n \left(\sum_{j=1}^g \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2}(y_i - \mu_j)^2\right) \right). \end{aligned} \quad (2.6)$$

Dále využijeme faktu, že máme více různých pozorování než složek směsi, tudíž

$$\exists \varepsilon > 0, \exists i_0 \in \{1, \dots, n\} : \forall j \in \{1, \dots, g\} |y_{i_0} - \mu_j| > \varepsilon.$$

Bez újmy na obecnosti položme $i_0 = 1$. Dále roztrhneme součin v (2.6) na dvě části, nejprve odhadneme příspěvek do věrohodnostní funkce pro $i = 1$, a poté příspěvek součinu zbytku,

$$\begin{aligned} \sum_{j=1}^g \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2}(y_i - \mu_j)^2\right) &< \sum_{j=1}^g \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2}\varepsilon^2\right) = \frac{g}{\sigma} \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right), \\ \prod_{i=2}^n \left(\sum_{j=1}^g \frac{1}{\sigma} \exp\left(-\frac{(y_i - \mu_j)^2}{2\sigma^2}\right)\right) &< \prod_{i=2}^n \left(\sum_{j=1}^g \frac{1}{\sigma} \exp\left(-\frac{0^2}{2\sigma^2}\right)\right) = \left(\frac{g}{\sigma}\right)^{n-1}. \end{aligned} \quad (2.7)$$

Nyní dosadíme rovnice (2.7) do (2.6) a vyjádříme:

$$L(\boldsymbol{\psi}) < \frac{g}{\sigma} \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right) \cdot \left(\frac{g}{\sigma}\right)^{n-1} = \left(\frac{g}{\sigma}\right)^n \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right). \quad (2.8)$$

Položme

$$f(\sigma) = \left(\frac{g}{\sigma}\right)^n \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right), \quad \sigma \in (0, \infty).$$

Nyní si stačí uvědomit, že funkce $f(\sigma)$ je spojitá na $(0, \infty)$, a navíc

$$\lim_{\sigma \rightarrow \infty} f(\sigma) = 0 \quad \& \quad \lim_{\sigma \rightarrow 0} f(\sigma) = 0,$$

z čehož je patrné, že $f(\sigma)$ je shora omezená, a tudíž z (2.8) je i věrohodnostní funkce $L(\boldsymbol{\psi})$ shora omezená. □

2.2 Věrohodnostní rovnice

Je patrné, že odhad metodou maximální věrohodnosti často nemusí existovat a i v případě, kdy řešení existuje, nemusí být vyhovující. Místo maximalizace věrohodnostní funkce se však můžeme zaměřit na hledání řešení rovnice

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

případně ekvivalentně na hledání řešení rovnice

$$\partial \ell(\boldsymbol{\psi}) / \partial \boldsymbol{\psi} = \mathbf{0}. \quad (2.9)$$

Rovnici (2.9) říkáme věrohodnostní rovnice. Je zřejmé, že pokud existuje řešení $\hat{\boldsymbol{\psi}}_{\text{MLE}} \in \text{int}(\Psi)$, pak je řešením (2.9). Naopak pokud najdeme prvek $\hat{\boldsymbol{\psi}} \in \Psi$ splňující (2.9), které je lokálním maximem věrohodnostní funkce (2.1), můžeme ho považovat za kandidáta na odhad metodou maximální věrohodnosti. Dokonce ho má smysl za jistých okolností jako odhad MMV přijmout, i přesto, že není globálním maximem. Lze totiž ukázat, že má podobné vlastnosti jako odhad MMV. My se touto teorií v práci zabývat nebudeme, více na toto téma je řečeno například v McLachlan a Peel (2000, kapitola 2.2).

2.3 Identifikovatelnost rozdělení směsi

Nyní se budeme věnovat jinému problému, který vyvstává při odhadování modelů konečných směsí, a to je problém identifikovatelnosti. Jinak řečeno, chceme odhadnout rozdělení náhodné veličiny, které je reprezentováno její hustotou závisící na parametrickém vektoru. Má tedy smysl si klást podmínku, aby různé hodnoty parametrických vektorů v našem parametrickém prostoru určovaly jiná rozdělení. Přesněji řečeno, aby platilo $\forall \boldsymbol{\psi}_1, \boldsymbol{\psi}_2 \in \Psi$:

$$f(\mathbf{y}; \boldsymbol{\psi}_1) = f(\mathbf{y}; \boldsymbol{\psi}_2) \text{ s.v.} \Leftrightarrow \boldsymbol{\psi}_1 = \boldsymbol{\psi}_2.$$

Jak se ale ukazuje, v kontextu směsí toto nemusí být pravda.

Příklad. Uvažme homoskedastickou normální směs o dvou složkách tak, že

$$f(y; \boldsymbol{\psi}) = \pi_1 \varphi(y; \mu_1, \sigma^2) + (1 - \pi_1) \varphi(y; \mu_2, \sigma^2).$$

Pak $\boldsymbol{\psi} = (\pi_1, \mu_1, \mu_2, \sigma^2)^\top$. Při přirozené volbě $\Psi = (0, 1) \times \mathbb{R}^2 \times (0, \infty)$, položme $\boldsymbol{\psi}_1 = (0.3, 1, 5, 1)^\top$ a $\boldsymbol{\psi}_2 = (0.7, 5, 1, 1)^\top$. Zřejmě $\boldsymbol{\psi}_1, \boldsymbol{\psi}_2 \in \Psi, \boldsymbol{\psi}_1 \neq \boldsymbol{\psi}_2$ a přitom pro s.v. $y \in \mathbb{R}$

$$f(y; \boldsymbol{\psi}_1) = 0.3 \cdot \varphi(y; 1, 1) + 0.7 \cdot \varphi(y; 5, 1) = 0.7 \cdot \varphi(y; 5, 1) + 0.3 \cdot \varphi(y; 1, 1) = f(y; \boldsymbol{\psi}_2).$$

Při takovéto volbě parametrů došlo pouze k prohození první a druhé složky, což výslednou hustotu neovlivní. □

Zřejmě v případě g složek existuje $g!$ permutací indexů složek a tudíž i $g!$ parametrických vektorů $\boldsymbol{\psi}$, které reprezentují totožné rozdělení (za předpokladu že všechny jsou prvkem parametrického prostoru).

Definice 4 (Identifikovatelnost rozdělení v modelu směsi). *Řekneme, že model směsi $\{f(\mathbf{y}; \boldsymbol{\psi}) : \boldsymbol{\psi} \in \Psi\}$ je identifikovatelný, pokud $\forall \boldsymbol{\psi}, \boldsymbol{\psi}' \in \Psi$ takové, že*

$$f(\mathbf{y}; \boldsymbol{\psi}) = \sum_{j=1}^g \pi_j f(\mathbf{y}; \boldsymbol{\theta}_j), \quad f(\mathbf{y}; \boldsymbol{\psi}') = \sum_{j=1}^{g'} \pi'_j f(\mathbf{y}; \boldsymbol{\theta}'_j)$$

a

$$f(\mathbf{y}; \boldsymbol{\psi}) = f(\mathbf{y}; \boldsymbol{\psi}') \text{ s.v.}$$

platí, že $g = g'$ a existuje permutace indexů τ složek směsi tak, že

$$\pi_j = \pi'_{\tau(j)}, \quad \boldsymbol{\theta}_j = \boldsymbol{\theta}'_{\tau(j)}, \quad j = 1, \dots, g.$$

V případě, kdy budeme diskutovat výsledky nějakých parametrických odhadů, budeme diskutovat případ pouze pro jedno uspořádání směsi. Kde to bude možné, budeme směsi řadit buď dle velikosti jejich poměru zastoupení případně podle velikosti první složky střední hodnoty. Na závěr bych chtěl ještě dodat, že tento fenomén nám nebude činit žádné obtíže při počítání odhadu MMV EM algoritmem, který se pro výpočet modelu směsi běžně používá, a bude představen v následující části.

2.4 EM algoritmus

Jedna z možností, jak hledat odhady metodou maximální věrohodnosti, je pomocí EM algoritmu. V této sekci se nejdříve budeme věnovat zavedení EM algoritmu pro směsi, vyjasníme si, proč a jak funguje, a dokážeme si základní vlastnost, která zaručuje, že po jedné iteraci EM algoritmu dostaneme „lepší“ odhad parametrů, než jsme doposud měli. V další části pak pro obecný model směsi krok algoritmu pro část parametrů vyjádříme.

Nechť $g \in \mathbb{N}$, $\pi_j \in (0,1)$, $j \in \{1, \dots, g\}$, tak, že

$$\sum_{j=1}^g \pi_j = 1$$

a necht' $f(\mathbf{y}; \boldsymbol{\theta}_j)$, $j \in \{1, \dots, g\}$, jsou hustoty na \mathbb{R}^p s neznámými parametry obsaženými ve vektoru $\boldsymbol{\theta}_j$. Mějme náhodný výběr $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ z rozdělení s hustotou

$$f(\mathbf{y}; \boldsymbol{\psi}) = \sum_{j=1}^g \pi_j f(\mathbf{y}; \boldsymbol{\theta}_j).$$

Nechť

$$\mathbf{Z}_i \sim \text{Mult}(1, \pi_1, \dots, \pi_g), \quad i = 1, \dots, n,$$

jsou náhodné vektory určující, z jaké skupiny populace objekt pochází. Připomeňme si, že značíme \mathbf{Y}_a celý náhodný výběr, podobně také $\mathbf{Z}_a = (\mathbf{Z}_1^\top, \dots, \mathbf{Z}_n^\top)^\top$ a $\mathbf{Y}_c = (\mathbf{Y}_a^\top, \mathbf{Z}_a^\top)^\top$.

Myšlenka, proč se EM algoritmus využívá pro počítání odhadů MMV v případech směsí, je jednoduchá. V zásadě je totiž velmi obtížné hledat maxima log-věrohodnostní funkce (2.1), kterou navíc ani při speciálních případech konkrétních hustot nelze nějak jednoduše maximalizovat. EM algoritmus tento problém tzv. nekompletních dat \mathbf{Y}_a , kdy neznáme příslušné vektory určující, ze kterých skupin naměřené objekty pocházely, převádí na hledání maxima věrohodnostní funkce odpovídající tzv. kompletním datům \mathbf{Y}_c .

Uvažme nyní hustoty těchto náhodných vektorů:

$$\begin{aligned} f(\mathbf{Y}_c; \boldsymbol{\psi}) &= f(\mathbf{Y}_a, \mathbf{Z}_a; \boldsymbol{\psi}) = f(\mathbf{Y}_c | \mathbf{Y}_a; \boldsymbol{\psi}) f(\mathbf{Y}_a; \boldsymbol{\psi}), \\ \log f(\mathbf{Y}_a; \boldsymbol{\psi}) &= \log f(\mathbf{Y}_c; \boldsymbol{\psi}) - \log f(\mathbf{Y}_c | \mathbf{Y}_a; \boldsymbol{\psi}). \end{aligned} \quad (2.10)$$

Pokud vyjádříme

$$f(\mathbf{Y}_c; \boldsymbol{\psi}) = \prod_{i=1}^n \prod_{j=1}^g \left(\pi_j f(\mathbf{Y}_i; \boldsymbol{\theta}_j) \right)^{Z_{ij}}, \quad f(\mathbf{Y}_a; \boldsymbol{\psi}) = \prod_{i=1}^n \sum_{j=1}^g \pi_j f(\mathbf{Y}_i; \boldsymbol{\theta}_j)$$

a dosadíme do (2.10) získáme rovnici ve tvaru

$$\ell(\boldsymbol{\psi}) = \ell_c(\boldsymbol{\psi}) - \log f(\mathbf{Y}_c | \mathbf{Y}_a; \boldsymbol{\psi}), \quad (2.11)$$

kde $\ell_c(\boldsymbol{\psi})$ je věrohodnostní funkce odpovídající kompletním datům \mathbf{Y}_c , tj.

$$\ell_c(\boldsymbol{\psi}) = \log \left(\prod_{i=1}^n \prod_{j=1}^g (\pi_j f(\mathbf{Y}_i; \boldsymbol{\theta}_j))^{Z_{ij}} \right) = \sum_{i=1}^n \sum_{j=1}^g Z_{ij} (\log \pi_j + \log f(\mathbf{Y}_i, \boldsymbol{\theta}_j)). \quad (2.12)$$

Nyní uvažme pozorovaný náhodný výběr $\mathbf{Y}_a = \mathbf{y}_a$ a spočtěme střední hodnoty v rovnosti (2.11) za podmínky, že známe \mathbf{y}_a a $\boldsymbol{\psi}'$. Pak s využitím, že

$$\mathbb{E} [\ell(\boldsymbol{\psi}) | \mathbf{y}_a, \boldsymbol{\psi}'] = \mathbb{E} \left[\sum_{i=1}^n \log f(\mathbf{y}_i; \boldsymbol{\psi}) \middle| \mathbf{y}_a, \boldsymbol{\psi}' \right] = \sum_{i=1}^n \log f(\mathbf{y}_i; \boldsymbol{\psi}) \mathbb{E} [1 | \mathbf{y}_a, \boldsymbol{\psi}'] = \ell(\boldsymbol{\psi}),$$

získáme

$$\ell(\boldsymbol{\psi}) = Q(\boldsymbol{\psi} | \boldsymbol{\psi}') - H(\boldsymbol{\psi} | \boldsymbol{\psi}'), \quad (2.13)$$

kde

$$Q(\boldsymbol{\psi} | \boldsymbol{\psi}') = \mathbb{E} [\ell_c(\boldsymbol{\psi}) | \mathbf{y}_a, \boldsymbol{\psi}'], \quad H(\boldsymbol{\psi} | \boldsymbol{\psi}') = \mathbb{E} [\log f(\mathbf{y}_a, \mathbf{Z}_a; \boldsymbol{\psi}) | \mathbf{y}_a, \boldsymbol{\psi}'].$$

Lemma 2 (Dempster a kol., 1977, str. 6). *Pro všechny dvojice $(\boldsymbol{\psi}, \boldsymbol{\psi}') \in \Psi \times \Psi$ platí $H(\boldsymbol{\psi} | \boldsymbol{\psi}') \leq H(\boldsymbol{\psi}' | \boldsymbol{\psi}')$.*

Důkaz. Zvol pevné $\boldsymbol{\psi}' \in \Psi$, pak pro libovolné $\boldsymbol{\psi} \in \Psi$ platí

$$\begin{aligned} H(\boldsymbol{\psi} | \boldsymbol{\psi}') - H(\boldsymbol{\psi}' | \boldsymbol{\psi}') &= \\ &= \mathbb{E} [\log f(\mathbf{y}_a, \mathbf{Z}_a; \boldsymbol{\psi}) | \mathbf{y}_a, \boldsymbol{\psi}'] - \mathbb{E} [\log f(\mathbf{y}_a, \mathbf{Z}_a; \boldsymbol{\psi}') | \mathbf{y}_a, \boldsymbol{\psi}'] \\ &= \mathbb{E} \left[\log \frac{f(\mathbf{y}_a, \mathbf{Z}_a; \boldsymbol{\psi})}{f(\mathbf{y}_a, \mathbf{Z}_a; \boldsymbol{\psi}')} \middle| \mathbf{y}_a, \boldsymbol{\psi}' \right] \\ &\leq \mathbb{E} \left[\frac{f(\mathbf{y}_a, \mathbf{Z}_a; \boldsymbol{\psi})}{f(\mathbf{y}_a, \mathbf{Z}_a; \boldsymbol{\psi}')} - 1 \middle| \mathbf{y}_a, \boldsymbol{\psi}' \right] \\ &= \mathbb{E} [1 | \mathbf{y}_a, \boldsymbol{\psi}] - \mathbb{E} [1 | \mathbf{y}_a, \boldsymbol{\psi}'] = 1 - 1 = 0. \end{aligned}$$

Ve třetím řádku v nerovnosti bylo využito, že $\forall x \in (0, \infty), \log x \leq x - 1$. Stojí za zmínku, že rovnost nastává právě tehdy, když $f(\mathbf{y}_a, \mathbf{z}_a | \mathbf{y}_a; \boldsymbol{\psi}) = f(\mathbf{y}_a, \mathbf{z}_a | \mathbf{y}_a; \boldsymbol{\psi}')$ skoro všude. □

Z (2.13) je patrné, že pokud pro pevné $\boldsymbol{\psi}' \in \Psi$ najdeme prvek $\boldsymbol{\psi} \in \Psi$ splňující nerovnost $Q(\boldsymbol{\psi} | \boldsymbol{\psi}') \geq Q(\boldsymbol{\psi}' | \boldsymbol{\psi}')$, pak

$$\ell(\boldsymbol{\psi}) \geq \ell(\boldsymbol{\psi}'). \quad (2.14)$$

To proto, že

$$\ell(\boldsymbol{\psi}) = Q(\boldsymbol{\psi} | \boldsymbol{\psi}') - H(\boldsymbol{\psi} | \boldsymbol{\psi}') \geq Q(\boldsymbol{\psi}' | \boldsymbol{\psi}') - H(\boldsymbol{\psi}' | \boldsymbol{\psi}') = \ell(\boldsymbol{\psi}').$$

Nyní už je zřejmé, jak bude algoritmus postupovat. Pro libovolnou počáteční hodnotu $\boldsymbol{\psi}^{(0)} \in \Psi$ budeme hledat $\boldsymbol{\psi}^{(1)} \in \Psi$ takové, že $Q(\boldsymbol{\psi}^{(1)} | \boldsymbol{\psi}^{(0)}) \geq Q(\boldsymbol{\psi}^{(0)} | \boldsymbol{\psi}^{(0)})$, a pak $\boldsymbol{\psi}^{(2)}$ a tak dále. To nám zaručí, že $\ell(\boldsymbol{\psi}^{(2)}) \geq \ell(\boldsymbol{\psi}^{(1)}) \geq \ell(\boldsymbol{\psi}^{(0)})$ a budeme tak získávat věrohodnější odhady parametrů. Takovému postupu říkáme zobecněný EM algoritmus. Samotný EM algoritmus spočívá v hledání maxima funkce $Q(\boldsymbol{\psi} | \boldsymbol{\psi}')$, $\boldsymbol{\psi} \in \Psi$. Přesněji vše popíšeme v následující definici.

Definice 5. Nechť M je zobrazení, $M : \Psi \rightarrow \Psi$, takové, že $\forall \psi' \in \Psi$

$$Q(M(\psi')|\psi') \geq Q(\psi'|\psi').$$

Pak algoritmus se zobrazením M nazveme zobecněný EM algoritmus. Pokud navíc platí, že $\forall \psi' \in \Psi$

$$M(\psi') = \arg \max_{\psi \in \Psi} Q(\psi|\psi'),$$

řekneme, že algoritmus se zobrazením M je EM algoritmus.

Zobrazení M nám po k iteracích vytvoří pro libovolnou počáteční hodnotu parametru prvek $\psi^{(k)}$. Z (2.14) víme, že $\ell(\psi^{(k+1)}) \geq \ell(\psi^{(k)})$. Tudíž, v případě, že je věrohodnostní funkce omezená, musí posloupnost $\ell(\psi^{(k)})$ konvergovat k nějakému reálnému číslu. Konvergence může ale nastat i v případě, že věrohodnostní funkce omezená není. Dempster a kol., 1977 ukázali, že za jistých podmínek pro výše uvedené funkce, zejména jejich diferencovatelnosti, pak

$$\ell^* = \lim_{k \rightarrow \infty} \ell(\psi^{(k)})$$

bude lokální maximum log-věrohodnostní funkce $\ell(\psi)$, případně nějaký sedlový bod. Samotný algoritmus ukončíme v okamžiku kdy rozdíl $\ell(\psi^{(k+1)}) - \ell(\psi^{(k)})$ bude menší než nějaká předem určená mez.

Pro použití EM algoritmu tedy bude první nutné zvolit několik různých počátečních hodnot a z nich EM algoritmus spustit. Pro všechny tyto pozice pak algoritmus dokonverguje k nějaké hodnotě funkce $\ell(\psi)$, které v drtivé většině případů bude jejím lokálním maximem. Jejich počet pak bude záviset na tvaru log-věrohodnostní funkce. Parametrické vektory, které tyto maxima nabývají, pak pro nás budou kandidáty na hledaný odhad. O jeho rozdělení, případně jiných vlastnostech ale nebudeme schopni nic říci, neboť ten nejspíše nebude klasickým odhadem metodou maximální věrohodnosti a zkoumání jeho vlastností by bylo, jak jsme si uvedli v sekci 2.2, nad rámec této bakalářské práce.

V následující části se budeme snažit pro konečné směsi takový algoritmus najít. Nechť $\psi^{(k)} \in \Psi$, $k \in \mathbb{N}_0$, budeme hledat $\psi^{(k+1)} \in \Psi$ takové, že

$$\psi^{(k+1)} = \arg \max_{\psi \in \Psi} Q(\psi|\psi^{(k)}).$$

Hledání tohoto prvku rozdělíme do dvou kroků. V prvním, v takzvaném E-kroku, spočítáme střední hodnotu $E[\ell_c(\psi)|\mathbf{y}_a, \psi^{(k)}]$ v závislosti na parametru ψ a poté, ve druhém kroku zvaném M-krok budeme tuto funkci maximalizovat pro $\psi \in \Psi$.

2.4.1 E-krok

Nyní tedy spočítáme střední hodnotu $E [\ell_c(\boldsymbol{\psi}) | \mathbf{y}_a, \boldsymbol{\psi}^{(k)}]$. Za $\ell_c(\boldsymbol{\psi})$ dosadíme z (2.12) a vyjádříme:

$$\begin{aligned}
 Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) &= E [\ell_c(\boldsymbol{\psi}) | \mathbf{y}_a, \boldsymbol{\psi}^{(k)}] \\
 &= E \left[\sum_{i=1}^n \sum_{j=1}^g Z_{ij} \left(\log \pi_j + \log f(\mathbf{y}_i, \boldsymbol{\theta}_j) \right) \middle| \mathbf{y}_a, \boldsymbol{\psi}^{(k)} \right] \\
 &= \sum_{i=1}^n \sum_{j=1}^g E [Z_{ij} | \mathbf{y}_a, \boldsymbol{\psi}^{(k)}] \left(\log \pi_j + \log f(\mathbf{y}_i, \boldsymbol{\theta}_j) \right) \\
 &= \sum_{i=1}^n \sum_{j=1}^g p_j(\mathbf{y}_i; \boldsymbol{\psi}^{(k)}) \left(\log \pi_j + \log f(\mathbf{y}_i, \boldsymbol{\theta}_j) \right). \tag{2.15}
 \end{aligned}$$

V poslední rovnosti jsme využili toho, že

$$E [Z_{ij} | \mathbf{y}_a, \boldsymbol{\psi}^{(k)}] = P (Z_{ij} = 1 | \mathbf{y}_a, \boldsymbol{\psi}^{(k)}) = p_j(\mathbf{y}_i; \boldsymbol{\psi}^{(k)}).$$

Zde symbolem $p_j(\mathbf{y}_i; \boldsymbol{\psi}^{(k)})$ značíme zpětnou pravděpodobnost toho, že \mathbf{y}_i náleželo do j -té skupiny, pokud hustota směsi odpovídá parametru $\boldsymbol{\psi}^{(k)}$. Výpočet byl proveden v (1.2). Tedy

$$p_j(\mathbf{y}_i; \boldsymbol{\psi}^{(k)}) = \frac{\pi_j^{(k)} f(\mathbf{y}_i; \boldsymbol{\theta}_j^{(k)})}{f(\mathbf{y}_i; \boldsymbol{\psi}^{(k)})}. \tag{2.16}$$

Dále budeme značit symbolem

$$p_{ij}^{(k)} = p_j(\mathbf{y}_i; \boldsymbol{\psi}^{(k)}).$$

2.4.2 M-krok

Ted' bude naším cílem najít řešení rovnice $\boldsymbol{\psi}^{(k+1)} = \arg \max_{\boldsymbol{\psi} \in \Psi} Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$. Označíme-li

$$Q_1(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) = \sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \log \pi_j, \quad Q_2(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) = \sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \log f(\mathbf{y}_i, \boldsymbol{\theta}_j),$$

z (2.15) lze ověřit, že $Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) = Q_1(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) + Q_2(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$. Navíc v $Q_1(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$ se z parametrů vyskytují pouze poměry složek π_j a v $Q_2(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$ naopak pouze parametry hustot $\boldsymbol{\theta}_j$. To nám umožní maximalizaci rozdělit do dvou částí. V první, pomocí maximalizace $Q_1(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$, najdeme nové odhady $\pi_j^{(k+1)}$ a ve druhé se poté zaměříme na hledání nových odhadů $\boldsymbol{\theta}_j^{(k+1)}$.

Nyní se tedy budeme věnovat hledání $\pi_j^{(k+1)}$. Označme

$$\boldsymbol{\pi} = (\pi_1, \dots, \pi_g)^\top, \quad P = \left(\boldsymbol{\pi} \in (0,1)^g : \sum_{j=1}^g \pi_j = 1 \right), \quad P_{1j}^{(k)} = \sum_{i=1}^n p_{ij}^{(k)}.$$

Pak P je parametrický prostor pro $\boldsymbol{\pi}$. Pomocí Lagrangeovy věty o multiplikátoru nyní vyřešíme rovnici

$$\boldsymbol{\pi}^{(k+1)} = \arg \max_{\boldsymbol{\pi} \in P} \sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \log \pi_j = \arg \max_{\boldsymbol{\pi} \in P} \sum_{j=1}^g P_{1j}^{(k)} \log \pi_j, \quad (2.17)$$

Protože P není kompaktní prostor, musíme ho nějak uzavřít. To ale nebude obtížné, protože uzavřený není jen když $\pi_j \rightarrow 0$ pro nějaké j a v tomto případě se hodnota sumy blíží k $-\infty$. Formálně položíme $c = \min_{j=1, \dots, g} P_{1j}^{(k)}$, $c > 0$ a uvažme $\boldsymbol{\pi} \in P$ takové, že $\exists j_0 : \pi_{j_0} \leq \left(\frac{c}{n}\right)^{gn/c}$. Potom, protože sčítance sumy jsou záporné

$$\sum_{j=1}^g P_{1j}^{(k)} \log \pi_j < P_{1j_0}^{(k)} \log \pi_{j_0} \leq c \log \left(\frac{c}{n}\right)^{gn/c} = gn \log \frac{c}{n}.$$

Nyní je zřejmé, že pokud maximum sumy bude větší než $gn \log \frac{c}{n}$, musí být splněna následující podmínka:

$$\frac{\partial}{\partial \pi_j} \left(\sum_{j=1}^g P_{1j}^{(k)} \log \pi_j + \lambda \left(\sum_{j=1}^g \pi_j - 1 \right) \right) (\pi_j^{(k+1)}) = 0, \quad j = 1, \dots, g,$$

kde λ je Lagrangeův multiplikátor. Z čehož máme:

$$P_{1j}^{(k)} \frac{1}{\pi_j^{(k+1)}} + \lambda = 0 \quad \Leftrightarrow \quad P_{1j}^{(k)} + \lambda \pi_j^{(k+1)} = 0, \quad j = 1, \dots, g. \quad (2.18)$$

Sečtením rovnic (2.18) pro všechna j a s využitím, že

$$\sum_{j=1}^g \pi_j^{(k+1)} = 1 \quad \& \quad \sum_{j=1}^g P_{1j}^{(k)} = n,$$

dostaneme $\lambda = -n$. Dosazením do (2.18) získáme, že

$$\pi_j^{(k+1)} = \frac{P_{1j}^{(k)}}{n}, \quad j = 1, \dots, g.$$

Poslední, co zbývá ukázat je:

$$\sum_{j=1}^g P_{1j}^{(k)} \log \frac{P_{1j}^{(k)}}{n} > \sum_{j=1}^g n \log \frac{c}{n} = gn \log \frac{c}{n}.$$

Tedy můžeme říci, že

$$\pi_j^{(k+1)} = \frac{P_{1j}^{(k)}}{n} = \frac{1}{n} \sum_{i=1}^n p_{ij}^{(k)}, \quad j = 1, \dots, g, \quad (2.19)$$

maximalizuje (2.17). Zajímavá je interpretace těchto vzorců. Každé pozorování přispívá do odhadu π_j velikostí pravděpodobnosti (při aktuální hodnotě $\boldsymbol{\psi}^{(k)}$), se kterou do j -té komponenty patří, což je velmi přímočaré a odpovídá intuici.

Zbývá nám maximalizovat $Q_2(\boldsymbol{\psi}|\boldsymbol{\psi}^{(k)})$. Tedy najít nový odhad $\boldsymbol{\xi}^{(k+1)}$ parametru $\boldsymbol{\xi}$. To se v takovéto obecnosti ukazuje jako nemožné. Jedna z možných cest je hledání řešení rovnice

$$\frac{\partial Q_2(\boldsymbol{\psi}|\boldsymbol{\psi}^{(k)})}{\partial \boldsymbol{\xi}} = \sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \frac{\partial \log f(\mathbf{y}_i, \boldsymbol{\theta}_j)}{\partial \boldsymbol{\xi}} = \mathbf{0}.$$

Tyto rovnice však už v obecném tvaru nejdou dále nějak pěkně upravit. Ukážeme si ale, že v případě normálních směsí lze najít explicitní řešení. To bude předmětem následující kapitoly.

Kapitola 3

Normální směsi

V této kapitole se zaměříme na rodinu směrů, které v jejím samotném kontextu hrají prominentní roli. Těmi jsou normální směsi. V předchozí kapitole jsme odvodili pouze část vzorců pro iteraci EM algoritmu a nyní si ukážeme, že v případě normálních směrů lze explicitně odvodit vzorce i pro ostatní parametry modelu.

Mějme tedy p -rozměrnou normální směs ve tvaru (1.4), tj. směs jejíž hustoty složek jsou tvaru

$$f(\mathbf{y}; \boldsymbol{\theta}_j) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_j|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1}(\mathbf{y} - \boldsymbol{\mu}_j)\right), \quad j = 1, \dots, g,$$

kde vektor $\boldsymbol{\theta}_j$ obsahuje vektor středních hodnot $\boldsymbol{\mu}_j$ a prvky matice $\boldsymbol{\Sigma}_j$, které jsou v horním trojúhelníku. Ty budeme značit $\text{vec}(\boldsymbol{\Sigma}_j)$. Tj $\boldsymbol{\theta}_j = (\boldsymbol{\mu}_j^\top, \text{vec}(\boldsymbol{\Sigma}_j)^\top)^\top$.

3.1 EM algoritmus pro normální směsi

Nechť $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ je náhodný výběr z p -rozměrné g -složkové normální směsi $f(\mathbf{y}, \boldsymbol{\psi})$, kde vektor $\boldsymbol{\psi}$ obsahuje všechny neznámé parametry modelu. Z minulé kapitoly víme, že naším cílem je maximalizovat funkci

$$Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) = \mathbf{E} [\ell_c(\boldsymbol{\psi}) | \mathbf{y}_a, \boldsymbol{\psi}^{(k)}]$$

pro předem dané $\boldsymbol{\psi}^{(k)}$. Označili jsme si

$$\boldsymbol{\psi}^{(k+1)} = \arg \max_{\boldsymbol{\psi} \in \Psi} Q(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) = \arg \max_{\boldsymbol{\psi} \in \Psi} \left(Q_1(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) + Q_2(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)}) \right)$$

a pomocí maximalizace funkce $Q_1(\boldsymbol{\psi} | \boldsymbol{\psi}^{(k)})$ jsme získali nový odhad (2.19) poměrů složek směsi

$$\pi_j^{(k+1)} = \frac{P_{1j}^{(k)}}{n}, \quad j = 1, \dots, g.$$

Zbývá nám dopočítat parametry vektoru $\boldsymbol{\xi}$, což v případě námi uvažovaných normálních směrů znamená nové odhady středních hodnot a rozptylových matic.

Vyjádříme-li $Q_2(\boldsymbol{\psi}|\boldsymbol{\psi}^{(k)})$ pro normální směs, dostaneme

$$\begin{aligned}\boldsymbol{\xi}^{(k+1)} &= \arg \max \sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \left(-\frac{1}{2} \log((2\pi)^k |\boldsymbol{\Sigma}_j|) - \frac{1}{2} (\mathbf{y}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_j) \right) \\ &= \arg \max \sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \left(\log(|\boldsymbol{\Sigma}_j^{-1}|) - (\mathbf{y}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_j) \right).\end{aligned}\quad (3.1)$$

Uvědomíme si, že nalezení nových středních hodnot je jednoduché. To z toho důvodu, že (3.1) je kvadratickou funkcí $\boldsymbol{\mu}_j$, navíc se záporným koeficientem u druhé mocniny. Tudíž pro $\boldsymbol{\mu}_j^{(k+1)}$ nutně musí platit a je i postačující, že

$$\frac{\partial}{\partial \boldsymbol{\mu}_j} \left(\sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \left(\log(|\boldsymbol{\Sigma}_j^{-1}|) - (\mathbf{y}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_j) \right) \right) (\boldsymbol{\mu}_j^{(k+1)}) = \mathbf{0}.$$

Derivováním dle $\boldsymbol{\mu}_j$, využitím, že $\boldsymbol{\Sigma}_j^{-1}$ je regulární matice a následnou úpravou potom dostaneme, že

$$\begin{aligned}\sum_{i=1}^n p_{ij}^{(k)} 2\boldsymbol{\Sigma}_j^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_j^{(k+1)}) &= \mathbf{0}, \\ \sum_{i=1}^n p_{ij}^{(k)} \mathbf{y}_i &= \sum_{i=1}^n p_{ij}^{(k)} \boldsymbol{\mu}_j^{(k+1)}, \\ \boldsymbol{\mu}_j^{(k+1)} &= \frac{\sum_{i=1}^n p_{ij}^{(k)} \mathbf{y}_i}{\sum_{i=1}^n p_{ij}^{(k)}}.\end{aligned}\quad (3.2)$$

Podobně jako vzorec pro výpočet poměrů složek směsi, vzorec (3.2) má také přirozenou interpretaci. Každé měření přispívá do odhadu dané střední hodnoty proporcí, se kterou si v daném okamžiku myslíme, že do té složky patří. Tento součet pak vydělíme předpokládanou velikostí složky ve směsi. Jde tudíž o analogii váženého výběrového průměru.

Konečně nový odhad rozptylové matice $\boldsymbol{\Sigma}_j$ spočítáme obdobně. Vyjádříme

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_j^{-1}} \left(\sum_{i=1}^n \sum_{j=1}^g p_{ij}^{(k)} \left(\log(|\boldsymbol{\Sigma}_j^{-1}|) - (\mathbf{y}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_j) \right) \right) (\boldsymbol{\Sigma}_j^{(k+1)}) = \mathbf{0}.$$

Odsud dostaneme, že je nutné řešit rovnici

$$\sum_{i=1}^n p_{ij}^{(k)} \left(\boldsymbol{\Sigma}_j^{(k+1)} - (\mathbf{y}_i - \boldsymbol{\mu}_j^{(k+1)}) (\mathbf{y}_i - \boldsymbol{\mu}_j^{(k+1)})^\top \right) = \mathbf{0}.$$

Jednoduchou úpravou poté dostaneme výsledný vzorec

$$\boldsymbol{\Sigma}_j^{(k+1)} = \frac{\sum_{i=1}^n p_{ij}^{(k)} (\mathbf{y}_i - \boldsymbol{\mu}_j^{(k+1)}) (\mathbf{y}_i - \boldsymbol{\mu}_j^{(k+1)})^\top}{\sum_{i=1}^n p_{ij}^{(k)}}, \quad j = 1, \dots, g.\quad (3.3)$$

Zbývá ukázat, že se opravdu jedná o maximum funkce. K tomu nám bude stačit zjištění, že se jedná o konkávní funkci matice Σ_j^{-1} . S využitím pravidel maticového kalkulu o derivaci inverzní matice můžeme počítat

$$\frac{\partial^2}{\partial(\Sigma_j^{-1})^2} Q_2(\boldsymbol{\psi}|\boldsymbol{\psi}^{(k)}) = \sum_{i=1}^n -p_{ij}^{(k)}(\Sigma_j \otimes \Sigma_j) = -P_{1j}^{(k)}(\Sigma_j \otimes \Sigma_j), \quad (3.4)$$

kde symbol \otimes značí Kroneckerův maticový součin. Potřebujeme, aby matice na pravé straně v (3.4) byla negativně definitní pro všechna Σ_j^{-1} , k čemuž stačí ukázat, že matice $\Sigma_j \otimes \Sigma_j$ je pozitivně definitní. Pro naše účely tuto vlastnost ukážeme pouze pro matice Σ_j rozměru 2×2 . Předpokládejme, že prvky matice Σ_j jsou postupně a_{11} , a_{12} , a_{21} a a_{22} . Ze Schurovy komplementární podmínky pro pozitivní definitnost matice (Dattorro, 2015, strana 527) nám stačí ukázat, že matice $a_{11}\Sigma_j$ a $a_{22}\Sigma_j - a_{12}\Sigma_j \cdot \frac{1}{a_{11}}\Sigma_j^{-1} \cdot a_{21}\Sigma_j$ jsou pozitivně definitní. To je ale přímý důsledek toho, že Σ_j je pozitivně definitní díky čemuž ze Sylvestrova kritéria víme, že $a_{11} > 0$ a $a_{11}a_{22} > a_{12}a_{21}$.

Znovu se jedná o analogii výpočtu výběrové rozptylové matice, které jsou navíc váženy odhadovanou pravděpodobností, že daný prvek náleží do j -té skupiny.

Označíme-li navíc

$$\mathbf{P}_{2j}^{(k)} = \sum_{i=1}^n p_{ij}^{(k)} \mathbf{y}_i \quad \& \quad \mathbf{P}_{3j}^{(k)} = \sum_{i=1}^n p_{ij}^{(k)} \mathbf{y}_i \mathbf{y}_i^\top,$$

můžeme všechny parametry pro $(k+1)$ ní iteraci EM algoritmu normální směsi odvozené v (2.19),(3.2) a (3.3) vyjádřit jako

$$\begin{aligned} \pi_j^{(k+1)} &= \frac{P_{1j}^{(k)}}{n}, \quad \boldsymbol{\mu}_j^{(k+1)} = \frac{\mathbf{P}_{2j}^{(k)}}{P_{1j}^{(k)}}, \\ \Sigma_j^{(k+1)} &= \frac{\mathbf{P}_{3j}^{(k)} - \left(P_{1j}^{(k)}\right)^{-1} \mathbf{P}_{2j}^{(k)} \left(\mathbf{P}_{2j}^{(k)}\right)^\top}{P_{1j}^{(k)}}, \quad (j = 1, \dots, g). \end{aligned} \quad (3.5)$$

Zobrazení definované vzorci (3.5) pro $\boldsymbol{\psi} \in \Psi$ je tedy EM algoritmem alespoň pro dvourozměrné normální směsi. Obecně v literatuře, která se odvozením vzorců EM algoritmu pro normální směsi zabývá, bylo zvykem odvodit toto zobrazení pouze z nutných podmínek existence extrému. To, zda-li tyto parametry skutečně maximalizují (2.15), nebylo, pokud nám je známo, nikde detailně ukázáno a tento fakt se tiše předpokládá.

V případě homoskedastické směsi, když předpokládáme, že rozptylové matice jsou stejné, můžeme postupovat obdobně pomocí derivace funkce $Q_2(\boldsymbol{\psi}|\boldsymbol{\psi}^{(k)})$. Po jednoduchých úpravách získáme vážený průměr rozptylových matic (3.3).

$$\begin{aligned} \Sigma^{(k+1)} &= \frac{\sum_{j=1}^g P_{1j}^{(k)} \Sigma_j^{(k+1)}}{n} = \frac{\sum_{j=1}^g \sum_{i=1}^n p_{ij}^{(k)} (\mathbf{y}_i - \boldsymbol{\mu}_j^{(k+1)}) (\mathbf{y}_i - \boldsymbol{\mu}_j^{(k+1)})^\top}{n} \\ &= \frac{1}{n} \sum_{j=1}^g \left(\mathbf{P}_{3j}^{(k)} - \left(P_{1j}^{(k)}\right)^{-1} \mathbf{P}_{2j}^{(k)} \left(\mathbf{P}_{2j}^{(k)}\right)^\top \right). \end{aligned}$$

3.2 Odhad modelu produkce cytokinu IL10

Nyní navážeme na model produkce cytokinu IL10, který jsme uvedli v sekci 1.4. V předchozích kapitolách jsme popsali způsob odhadování modelu pomocí metody maximální věrohodnosti s využitím EM algoritmu, jenž nyní aplikujeme. Jakým způsobem analýzu provedeme, popíšeme na následujících řádcích.

K analýze použijeme software R a volně dostupný balík *EMCluster* (Chen a kol., 2012), v němž je již EM algoritmus naprogramován. Cílem analýzy bude najít vhodný model rozdělení produkce cytokinu, který odpovídá jak naměřeným datům tak i lidské intuici. Uvažujme tudíž dvě skupiny, jednu nemocnou a druhou zdravou, čemuž odpovídá volba $g = 2$. Pro spuštění EM algoritmu bude potřeba vybrat počáteční podmínky pro algoritmus. Ty se pokusíme zvolit tak, abychom objevili co možná všechna lokální maxima log-věrohodnostní funkce (2.1). Zvolme tedy 9 středních hodnot rovnoměrně rozdělených po ploše, na které jsou umístěna data, a dále 9 různých rozptylových matic. Samotný algoritmus spustíme z odlišných kombinací těchto středních hodnot a rozptylových matic při shodném poměru složek směsi 0,5. Tímto postupem vytvoříme celkem 2916 různých počátečních pozic, ze kterých postupně algoritmus odstartujeme. Ukázalo se, že všech těchto počátečních pozic algoritmus dokonvergoval vždy k jednomu ze sedmi parametrických vektorů, určujících rozdělení směsí uvedených v tabulce 3.1.

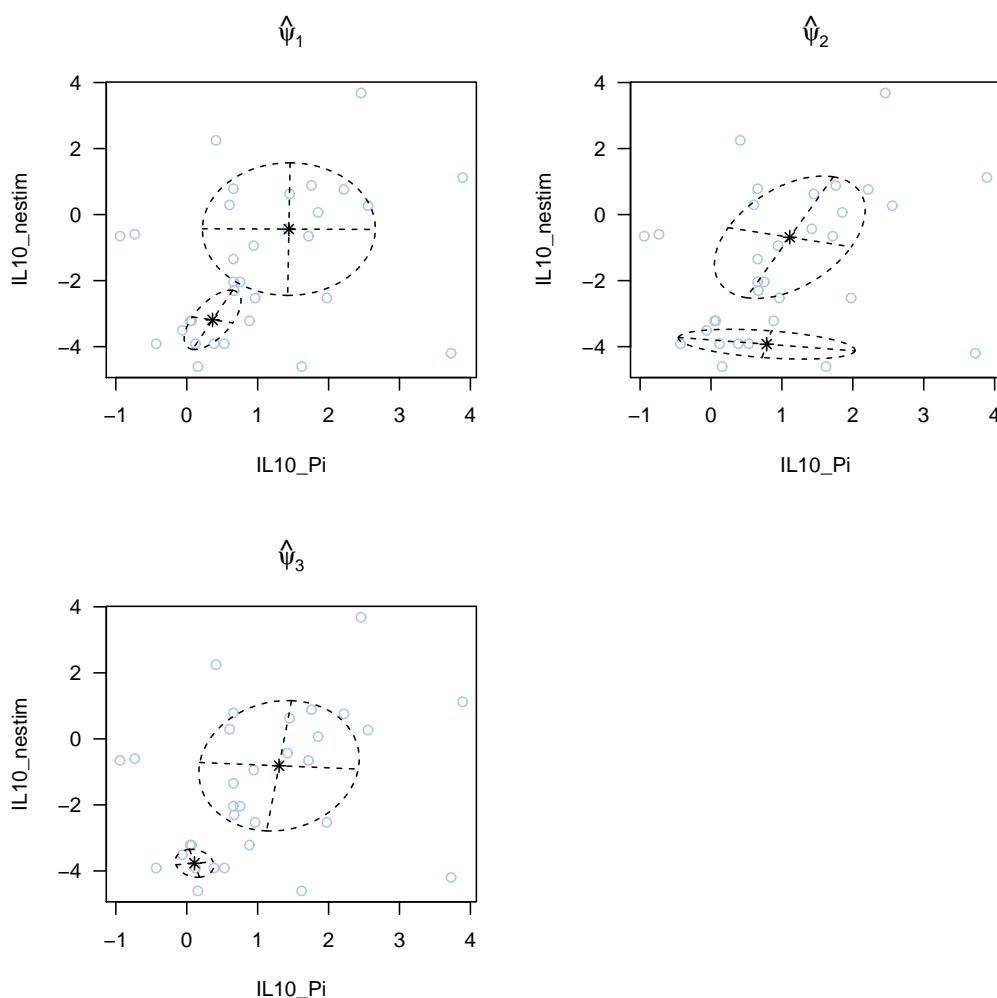
	π	μ_1	μ_2	σ_1	σ_2	ρ	$\ell(\boldsymbol{\psi})$
$\hat{\boldsymbol{\psi}}_1$	0,62	1,44	-0,44	1,22	2,01	0,01	-108,7
	0,38	0,36	-3,19	0,40	0,90	0,63	
$\hat{\boldsymbol{\psi}}_2$	0,75	1,11	-0,69	1,06	1,85	0,44	-111,6
	0,25	0,79	-3,93	1,25	0,45	-0,41	
$\hat{\boldsymbol{\psi}}_3$	0,77	1,30	-0,81	1,13	1,98	0,11	-109,2
	0,23	0,11	-3,77	0,27	0,42	-0,16	
$\hat{\boldsymbol{\psi}}_4$	0,94	0,85	-1,48	0,89	2,10	0,45	41,4
	0,06	3,81	-1,54	0,08	2,66	1,00	
$\hat{\boldsymbol{\psi}}_5$	0,90	1,06	-1,44	0,95	2,19	0,59	-109,9
	0,10	0,76	-1,89	2,05	1,60	-1,00	
$\hat{\boldsymbol{\psi}}_6$	0,84	0,97	-1,83	0,94	2,17	0,30	-113,5
	0,16	1,34	0,26	1,73	0,67	0,93	
$\hat{\boldsymbol{\psi}}_7$	0,94	0,96	-1,52	1,04	2,05	0,55	12,7
	0,06	2,07	-0,98	1,66	3,22	-1,00	

Pro $\hat{\boldsymbol{\psi}}_i$ tabulka ukazuje odhad parametrů dvousložkové normální směsi jemu odpovídající. π je poměr zastoupení složky ve směsi a $\mu_1, \mu_2, \sigma_1, \sigma_2, \rho$ pak její parametry. (V každé směsi jsou dvě složky.)

Tabulka 3.1: Parametry rozdělení směsi určené lokálními maximizátory věrohodnostní funkce.

Je zajímavé, že hodnoty log-věrohodnostní funkce v získaných parametrických vektorech jsou od sebe dosti odlišné. Zejména hodnoty $\ell(\hat{\psi}_4)$ a $\ell(\hat{\psi}_7)$ jsou překvapivě vysoké a zřetelně se liší od ostatních, které jsou všechny blízko hodnotě -110 , jak se můžeme přesvědčit v tabulce 3.1. Pokud se ale podrobněji podíváme na korelační koeficient $\hat{\psi}_4$ a $\hat{\psi}_7$ druhé složky v této tabulce, zjistíme, že je 1, resp. -1 . To znamená, že rozptylová matice této složky je singularní. O tom, že podobná situace může nastat, jsme diskutovali v sekci 2.1. Tato degenerovaná rozdělení určená vektory $\hat{\psi}_4$ a $\hat{\psi}_7$ nám zřejmě neposkytnou vhodný model.

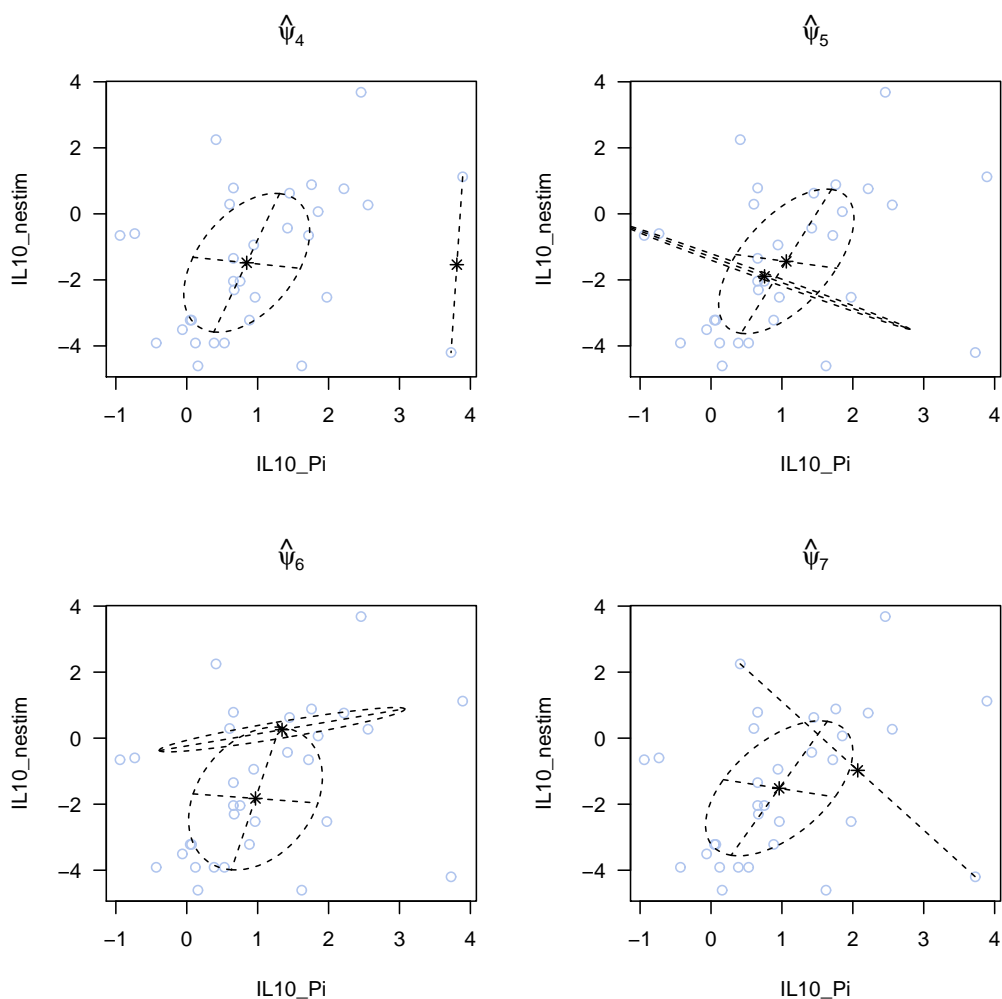
Vidíme, že ani jeden ze zbylých odhadů parametru ψ nebude odhadem metodou maximální věrohodnosti, protože nemaximalizuje věrohodnostní funkci. Avšak jak je z obrázků 3.1 a 3.2, kde jsou rozdělení směsí naznačena, patrné, některá získaná rozdělení by mohla přijatelně vystihovat popisovaný model. Protože jsme v průběhu této bakalářské práce nezmínili žádné možné testy, kterými lze vhodnost modelu testovat, v případě že bychom chtěli vybrat jedno rozdělení, které „sedí nejlépe“, nezbývá nám nic jiného než použít intuici a heuristiku.



Obrázek 3.1: Znázornění prvních tří rozdělení, jejichž odpovídající parametrické vektory lokálně maximalizují věrohodnostní funkci.

Jedním z vodítek, které z rozdělení by mohlo být vhodné, může být informace o poměru podílu složky π_1 . Data, která máme k dispozici, totiž obsahovala i údaj, zda-li měřený pacient trpěl paradontózou či ne. Díky tomu víme, že zkoumáme měření 17 nemocných a 15 zdravých pacientů. Proto lze očekávat, že poměr π_1 by měl být okolo hodnoty 0,55. Pokud se podíváme na parametrické vektory $\hat{\psi}_5$ a $\hat{\psi}_6$, hodnota π_1 je u nich větší než 0,83, což je výrazně více, než je očekáváno. Navíc složka s menším poměrem je u obou rozdělení do značné míry singulární, jak je patrné z obrázku 3.2 i hodnot v tabulce 3.1, a z medicínského hlediska není jediný důvod, proč by měla být produkce cytokinu u jedné takové skupiny takto svázána. Proto se rozděleními reprezentovanými těmito parametrickými vektory nebudeme dále zabývat.

Zbylé 3 rozdělení, definované parametrickými vektory $\hat{\psi}_1$, $\hat{\psi}_2$ a $\hat{\psi}_3$ zatím ponecháme jako realistické a budeme je zkoumat podrobněji v další sekci, kde budeme porovnávat správnost určení jednotlivých pacientů do daných skupin. Zjistíme, jak moc modely odpovídají naměřeným hodnotám, a zkusíme určit, který z nich nejlépe popisuje produkci cytokinu IL10, tak jak byla zjištěna.



Obrázek 3.2: Znázornění zbylých čtyř rozdělení, jejichž odpovídající parametrické vektory lokálně maximalizují věrohodnostní funkci.

Kapitola 4

Shlukování

V některých aplikacích modelů konečných směsí může být cílem zjištění, do kterých skupin pozorované prvky patří. Je přirozené se ptát, pokud předpokládáme, že studovaná populace se dělí řečneme na dvě skupiny, ze které skupiny daný objekt pochází. Dost možná je toto totiž i hlavní cíl, proč model směsi odhadujeme.

Mějme tedy g různých skupin nějaké populace a sledujme na ní rozdělení náhodného vektoru \mathbf{Y} . Nechť tato veličina má v j -té skupině rozdělení definované hustotou $f(\mathbf{y}; \boldsymbol{\theta}_j)$. Pak, jak jsme si ukázali, má náhodný vektor \mathbf{Y} rozdělení dáno hustotou (1.3), tedy hustotou

$$f(\mathbf{y}; \boldsymbol{\psi}) = \sum_{j=1}^g \pi_j f(\mathbf{y}; \boldsymbol{\theta}_j),$$

kde $\boldsymbol{\theta}_j$ a π_j , $j = 1, \dots, g$, $\sum \pi_j = 1$, jsou neznámé parametry obsaženy ve vektoru $\boldsymbol{\psi}$. Naším cílem je na základě pozorování \mathbf{Y} určit, ze které skupiny pozorovaný objekt pochází. Jinými slovy to znamená určit pro \mathbf{y} jeho odpovídající vektor \mathbf{z} , jenž daný jev přesně popisuje.

Abychom tak ale mohli učinit, je první nutné odhadnout samotné rozdělení směsi \mathbf{Y} . To uděláme na základě pozorovaného náhodného výběru $\mathbf{y}_1, \dots, \mathbf{y}_n$, s hustotou (1.3), který je pro odhadnutí modelu nezbytný. Odhad lze pak provést například pomocí dříve uvedené metody maximální věrohodnosti a EM algoritmu. Předpokládejme tedy, že jsme pomocí jisté metody získali odhad $\hat{\boldsymbol{\psi}} \in \Psi$ parametrického vektoru $\boldsymbol{\psi}$. Predikce hodnot vektorů \mathbf{z}_i pak lze provést na základě velikosti zpětné pravděpodobnosti, že pozorovaný vektor \mathbf{y}_i patří do nějaké skupiny. Ty jsou dány pravděpodobnostmi $p_1(\mathbf{y}_i; \hat{\boldsymbol{\psi}}), \dots, p_g(\mathbf{y}_i; \hat{\boldsymbol{\psi}})$, kde používáme značení $p_j(\mathbf{y}_i; \hat{\boldsymbol{\psi}}) = \mathbb{P}(Z_{ij} = 1 | \mathbf{Y}_i = \mathbf{y}_i; \hat{\boldsymbol{\psi}})$, $j = 1, \dots, g$, viz (2.16).

Přesněji, hodnotu vektoru \mathbf{z}_i budeme predikovat vektorem $\hat{\mathbf{z}}_i$, jehož prvek \hat{z}_{ij} definujeme jako

$$\begin{aligned} \hat{z}_{ij} &= 1, & \text{když } j &= \arg \max_{l=1, \dots, g} p_l(\mathbf{y}_i; \hat{\boldsymbol{\psi}}), \\ &= 0, & \text{jinak,} \end{aligned} \tag{4.1}$$

pro $i = 1, \dots, n$ a $j = 1, \dots, g$. Je zřejmé, že v případě rovnosti některých z těchto pravděpodobností by se mohlo stát, že vektor $\hat{\mathbf{z}}_i$ obsahuje více než jednu jedničku.

I v tomto případě bychom pozorování přiřadili pouze jedné skupině, a to libovolné z těch, u kterých je dosažena nejvyšší hodnota $p_j(\mathbf{y}_i; \hat{\boldsymbol{\psi}})$. Označme tento odhad $\hat{z}_i = r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}})$.

V případě, že náš odhad $\hat{\boldsymbol{\psi}}$ přesně odpovídá skutečné hodnotě parametrického vektoru $\boldsymbol{\psi}$, nazývá se toto pravidlo přiřazování (4.1) Bayesovo pravidlo. Protože navíc všechny tyto pravděpodobnosti mají společný jmenovatel $f(\mathbf{y}_i)$, lze se pochopitelně rozhodovat pouze na základě relativních velikostí hustot složek vážených podle poměru zastoupení ve směsi, tj. na základě velikosti $\pi_j f(\mathbf{y}_i; \hat{\boldsymbol{\theta}}_j)$.

Bude užitečné se i zamyslet, jak bude fungovat toto pravidlo v případě, kdy získáme špatný odhad $\hat{\boldsymbol{\psi}}$ parametru $\boldsymbol{\psi}$. Je možné, že i tak budou \hat{z}_i stále dobrými odhady z_i čímž myslíme, že ve většině případů bude platit $r_B(\mathbf{y}_i; \boldsymbol{\psi}) = r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}})$. Ze způsobu určování těchto odhadů je zřejmé, že je podstatné, aby model velmi dobře odpovídal zejména na hranicích určujících různé skupiny, tedy v okolí množiny

$$\{\mathbf{y} : \pi_j f(\mathbf{y}; \boldsymbol{\theta}_j) = \pi_l f(\mathbf{y}; \boldsymbol{\theta}_l), \quad j < l = 2, \dots, g\}. \quad (4.2)$$

Tudíž například pro vzájemně odlehlejší skupiny je důležité jen, aby model přesně odpovídal zejména na krajích rozdělení jednotlivých skupin. Přesnost určení, odkud měření pochází, je tedy dána v podstatě vlastností jak úspěšně dokáže $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}})$ přiřazovat do správných skupin pozorování nejasného původu.

4.1 O nezdravých hodnotách zdravých lidí

V této poslední části mé bakalářské práce se ještě jednou vrátíme k modelu produkce cytokinu IL10 u lidí trpících paradontózou. Jak jsme si uvedli v sekci 3.2, metodou maximální věrohodnosti s pomocí EM algoritmu jsme získali 7 parametrických vektorů které lokálně maximalizují věrohodnostní funkci (2.1). Po porovnání vlastností, které mají jimi definovaná rozdělení s přirozenou lékařskou představou produkce cytokinu, nám zůstaly tři parametrické vektory $\hat{\boldsymbol{\psi}}_1, \hat{\boldsymbol{\psi}}_2, \hat{\boldsymbol{\psi}}_3$, které by mohly vhodně popisovat náš model.

Každé z těchto rozdělení má dvě dvourozměrné normální složky. Je přirozené předpokládat, že rozdělení, kterým se řídí produkce cytokinu IL10 u zdravých lidí bude to, jehož střední hodnota je v obou složkách menší. V případě rozdělení definované vektory $\hat{\boldsymbol{\psi}}_1, \hat{\boldsymbol{\psi}}_2$ a $\hat{\boldsymbol{\psi}}_3$ toto vždy splňuje složka s menším poměrem ve směsi, jak je patrné například z obrázku 3.1. Nyní navážeme na problematiku diskutovanou na počátku této kapitoly a budeme různým měřením přiřazovat skupinu, do které zkoumaní pacienti patří. To budeme provádět na základě pravidla $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_i)$, $i = 1, 2, 3$.

Protože víme, kteří pacienti trpěli paradontózou, můžeme pro porovnání uvážit i odhad parametrického vektoru $\hat{\boldsymbol{\psi}}_c$, který byl vypočten jako směs normálních rozdělení, které byly odhadnuty metodou maximální věrohodnosti pro každou skupinu zvlášť. Tedy jsme odhadli přímo parametry složek $\boldsymbol{\theta}_j$ a jejich poměry pak byly dopočteny jako poměr zastoupení nemocných/zdravých v datech.

Nejmenší počet šesti chyb zaznamenalo přiřazovací pravidlo $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_1)$, které bylo těsně následováno pravidlem $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_3)$, jak je patrné z tabulky 4.1. Třetí v pořadí skončilo $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_c)$ a nejhorší s více jak třetinou špatně určených pacientů bylo $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_2)$. Je zajímavé, že $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_c)$ skončilo v porovnání s dalšími pravidly o tolik hůře, ačkoliv na výpočet parametrického vektoru určující toto pravidlo bylo použito více informací než pro ostatní odhady. To je ale dáno zejména čtyřmi odlehlými pozorováními „zdravých“ lidí, které ve všech případech byly zařazeny do skupiny nemocných. Pro $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_c)$ měla ale tato pozorování důsledek zvětšení oblasti pro zdravé pacienty, do které se pak vešly i hodnoty pacientů nemocných.

Překvapující je i nula v kolonce špatně určených zdravých lidí u pravidla daného vektorem $\hat{\boldsymbol{\psi}}_3$. To znamená, že každý pacient, u kterého bylo odhadnuto, že je zdravý, byl zdravý. Pokud se podíváme na třetí graf v obrázku 3.1, můžeme si všimnout, že menší složka ze směsi má velmi malý rozptyl a je koncentrovaná u několika málo bodů v levém dolním rohu grafu. V té oblasti se ale nachází pouze 8 měření zdravých pacientů.

Jak bylo zmíněno, cytokin IL10, stejně jako další látky ze skupiny interleukinů jsou produkovány jako imunologická reakce při napadení těla. V našich datech za nemocného pacienta považujeme pacienta trpícího paradontózou. Pochopitelně ale i pro nás zdravý pacient může mít zvýšenou produkci IL10 z důvodu jiného onemocnění, což lze přijmout i za důvod čtyř velmi odlehlých pozorování zdravých pacientů. To by mohlo vést k utvoření hypotézy, že v těchto hodnotách se pohybují čísla úplně zdravých pacientů a že menší složka ze směsi určena odhadem $\hat{\boldsymbol{\psi}}_3$ je odhadem produkce IL10 u této skupiny.

Zmínili jsme, že přiřazovací pravidlo $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_1)$ je ještě o trochu lepší než výše diskutované $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_3)$. Jak je z obrázku 4.1 patrné, výsledné predikované vektory pro jednotlivá pozorování se liší pouze u pěti pacientů. Jejich měření leží relativně někde mezi skupinami s největší koncentrací zdravých a nemocných pacientů, tj. právě v okolí množiny (4.2). Jedná se tedy o měření nejasného původu a přiřazovací pravidla se vlastně liší jen v tom, že jednou tuto malou skupinku přiřadí k nemocným a jednou ke zdravým. Protože v ní jsou ale 3 hodnoty zdravých pacientů, tak pravidlo $r_B(\mathbf{y}_i; \hat{\boldsymbol{\psi}}_1)$ vychází v hodnocení lépe.

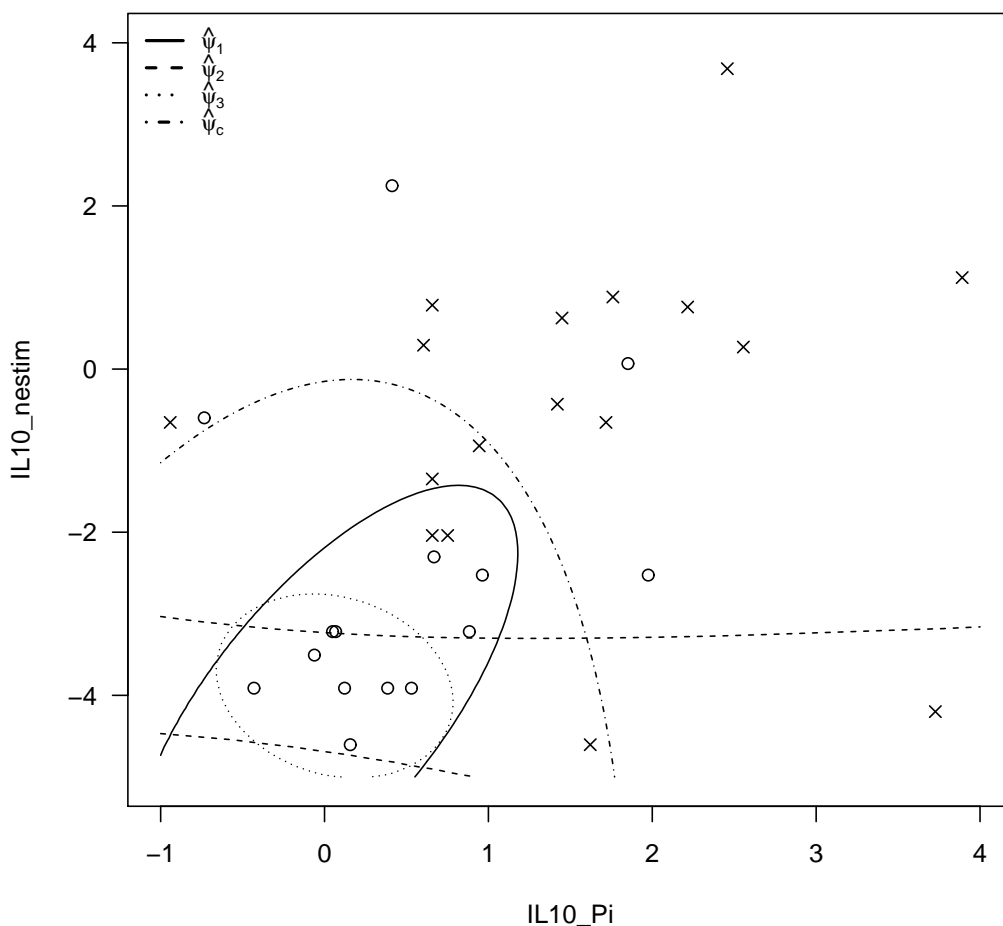
Bylo by předčasné dělat závěry ohledně výsledného tvaru směsi, ty by totiž měly být provedeny až na základě nějakých seriózních testů, které jsme v této práci neodvozovali ani neprováděli. Na druhou stranu si lze na základě naší

Pravidlo	Správně určení		Špatně určení		Počet chyb
	Zdraví	Nemocní	Zdraví	Nemocní	
$\boldsymbol{\psi}_1$	11	15	2	4	6
$\boldsymbol{\psi}_2$	6	15	2	9	11
$\boldsymbol{\psi}_3$	8	17	0	7	7
$\boldsymbol{\psi}_c$	11	12	5	4	9

Tabulka 4.1: Chybovost určování do skupin.

analýzy vytvořit hlubší představu o modelu a povaze pozorování. Pro skutečné použití našeho modelu v praxi by pak navíc byl potřeba odhad na základě náhodného výběru z celé populace, což výběr, se kterým jsme pracovali nespĺňoval. Data totiž jsou primárně určena pro ukázání vztahu zvýšené produkce interleukinů u pacientů trpící paradontózou a v zájmu výzkumných cílů paní doktorky Bártlové, která data naměřila, bylo získat hodnoty produkce u jistého počtu lidí v každé skupině. Proto je podíl nemocných a zdravých lidí v našem datovém souboru podobný, což pochopitelně neodpovídá realitě. Na druhou stranu by tento problém šel vyřešit prostým upravením složek směsi tak, aby u směsi reprezentující skupinu nemocných byl poměr roven velikosti rozšíření nemoci v populaci.

I přesto věřím, že se nám na tomto příkladě podařilo ilustrovat použití modelu směsi jako způsob analýzy výskytu nějaké vlastnosti v populaci.



Legenda: křížek znázorňuje hodnotu pacienta s paradontózou, kolečko zdravého.

Obrázek 4.1: Hranice přiřazovacího pravidla pro rozdělení definované uvažovanými parametrickými vektory ψ_1 , ψ_2 , ψ_3 a ψ_c .

Závěr

Jsme na konci bakalářské práce, která měla za cíl představit model konečné směsi a ukázat cestu, jakým se dá následně v praxi použít. Zejména jsme se věnovali odhadování modelu pomocí metody maximální věrohodnosti, při které ale může nastat množství komplikací. Problematická se nám ukázala hlavně vlastnost, že věrohodnostní funkce náhodného výběru nemusí být omezená, kvůli čemuž jsme se zaměřili na hledání lokálních maxim věrohodnostní funkce pomocí EM algoritmu. Ačkoliv tyto odhady pak nemaximalizují hodnotu věrohodnostní funkce, lze ukázat, že mají podobné vlastnosti jako odhady metodou maximální věrohodnosti. Bohužel už nám nezbyl prostor se touto velmi důležitou teorií zabývat, čímž by se použitelnost konečných směsí ukázala v plné síle. Zde vidím prostor, kam by se tato práce mohla dále rozvíjet.

V průběhu této práce jsme se detailně věnovali EM algoritmu, znalost jehož principu je pro pochopení funkcionality stěžejní. V rámci druhé kapitoly, která se odvození algoritmu věnuje, je zachycena i elegantní myšlenka, kterou algoritmus využívá. V dalších částech textu jsou pak odvozeny vzorce pro iteraci EM algoritmu aplikovaného na nejběžněji používaný model normálních směsí. V rámci toho jsme kladli důraz na splnění postačujících podmínek pro maximum funkce, když veškerá literatura, která se podobnému odvození věnovala a se kterou autor přišel do styku, hledala řešení pouze podmínek nutných. Funkčnost EM algoritmu se nám pak podařilo dokázat pouze pro nejvýše dvourozměrné normální směsi. Ve vyšší dimenzi se o ni můžeme pouze domnívat. To ale bylo postačující s ohledem na naši následnou analýzu v první kapitole uvedeného modelu, kde jsme zkoumali dvourozměrnou normální směs.

Při výběru zadání bakalářské práce autor kladl důraz na možnost prakticky aplikovat vypsání téma, které v ní mělo být náležitě vystiženo. Jsem přesvědčen, že model produkce interleukinu 10 jakožto obranné reakce lidského těla při napadení nějakou nemocí názorně ilustroval jedno z možných použití směsi. Souvislost produkce interleukinů při napadení paradontózou je aktuálně předmětem bádání paní doktorky Bártlové a pana doktora Běláčka z 1. lékařské fakulty Univerzity Karlovy a já věřím, že model směsi jim umožní trochu jiný náhled na jimi zkoumaný problém. Použití modelu směsi k ukázání jistých vztahů mezi produkcí cytokinu a onemocněním paradontózou už by pak byla taková třesnička na dortu této práce. Je ale zřejmé, že pro případné publikování výsledků by bylo nutné v analýze pokračovat, což součástí této práce nebude.

Pevně věřím, že tato práce poskytne všem čtenářům jednoduše pochopitelný přehled o problematice konečných směsí a jejím možném využití a zanechá v nich přesvědčení, že se jí vyplatí, s ohledem na její praktičnost i eleganci, věnovat.

Literatura

- AITKIN, M. (2001). Likelihood and bayesian analysis of mixtures. *Statistical Modelling*, **1**, 287–304.
- CHEN, W.-C., MAITRA, R. a MELNYKOV, V. (2012). EMCluster: EM algorithm for model-based clustering of finite mixture gaussian distribution. R Package, URL <http://cran.r-project.org/package=EMCluster>.
- DATTORRO, J. (2015). *Convex Optimization, Euclidean Distance Geometry*. Meboo Publishing USA. ISBN 978-0-578-16140-2.
- 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.
- MCLACHLAN, G. a PEEL, D. (2000). *Finite Mixture Model*. John Wiley & Sons, Inc., New York. ISBN 0-471-00626-2.
- PEARSON, K. (1894). Contributions to the theory of mathematical evolution. *Philosophical Transactions of the Royal Society of London A*, **185**, 71–110.

Seznam obrázků

1.1	Příklady několika normálních směsí.	8
1.2	Histogramy hodnot logaritmizované produkce stimulované i nestimulované IL10.	10
1.3	Bodový zákres dat produkce cytokinu IL10.	10
3.1	Znázornění prvních tří rozdělení, jejichž odpovídající parametrické vektory lokálně maximalizují věrohodnostní funkci.	26
3.2	Znázornění zbylých čtyř rozdělení, jejichž odpovídající parametrické vektory lokálně maximalizují věrohodnostní funkci.	27
4.1	Hranice přiřazovacího pravidla pro rozdělení definované uvažovanými parametrickými vektory ψ_1, ψ_2, ψ_3 a ψ_c	31

Seznam tabulek

3.1	Parametry rozdělení směsi určené lokálními maximizátory věrohodnostní funkce.	25
4.1	Chybovost určování do skupin.	30