

Model Gaussovih mješavina

Vrhovec, Veronika

Master's thesis / Diplomski rad

2020

Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj: **University of Zagreb, Faculty of Science / Sveučilište u Zagrebu, Prirodoslovno-matematički fakultet**

Permanent link / Trajna poveznica: <https://um.nsk.hr/um:nbn:hr:217:024573>

Rights / Prava: [In copyright](#)/[Zaštićeno autorskim pravom.](#)

Download date / Datum preuzimanja: **2024-07-02**



Repository / Repozitorij:

[Repository of the Faculty of Science - University of Zagreb](#)



Model Gaussovih mješavina

Vrhovec, Veronika

Master's thesis / Diplomski rad

2020

Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj: **University of Zagreb, Faculty of Science / Sveučilište u Zagrebu, Prirodoslovno-matematički fakultet**

Permanent link / Trajna poveznica: <https://um.nsk.hr/um:nbn:hr:217:024573>

Rights / Prava: [In copyright](#)/[Zaštićeno autorskim pravom.](#)

Download date / Datum preuzimanja: **2024-06-20**



Repository / Repozitorij:

[Repository of the Faculty of Science - University of Zagreb](#)



SVEUČILIŠTE U ZAGREBU
PRIRODOSLOVNO–MATEMATIČKI FAKULTET
MATEMATIČKI ODSJEK

Veronika Vrhovec

MODEL GAUSSOVIH MJEŠAVINA

Diplomski rad

Voditelj rada:
doc. dr. sc. Azra Tafro

Zagreb, 2020

Ovaj diplomski rad obranjen je dana _____ pred ispitnim povjerenstvom u sastavu:

1. _____, predsjednik
2. _____, član
3. _____, član

Povjerenstvo je rad ocijenilo ocjenom _____.

Potpisi članova povjerenstva:

1. _____
2. _____
3. _____

Sadržaj

Sadržaj	iii
Uvod	2
1 Model konačnih mješavina	3
1.1 Vjerojatnosni prostor i slučajne veličine	3
1.2 Konačne mješavine	10
1.3 Statističko zaključivanje u modelu Gaussovih mješavina	14
2 Algoritam maksimizacije očekivanja	22
2.1 Algoritam i procjena parametara	24
2.2 Praktična primjena algoritma	31
A Računalni kôdovi u Pythonu	37
A.1 Algoritam maksimizacije očekivanja	37
A.2 Algoritam k-srednjih vrijednosti ++	45
Bibliografija	47

Uvod

Ako želimo precizno opisati stvarni svijet koji nas okružuje, moramo ga sagledati iz više kuteva. Tako i kod modeliranja populacija koje opažamo često nije dovoljno koristiti jednodimenzionalnu analizu. A što ako nam ni višedimenzionalni pristup nije dovoljan da tu populaciju točno okarakteriziramo jer previdi postojanje podpopulacija? Intuitivno je tada modelirati svaku podpopulaciju zasebno, ali pritom voditi računa o njihovom međusobnom utjecaju. Teško je egzaktno tvrditi kada smo se prvi puta susreli s ovom vrstom problema u statističkoj analizi, ali već krajem 19. stoljeća nastaju prvi formalizirani modeli koji su primijenjeni na miješane populacije. Takve modele nazivamo modelima konačnih mješavina. Priroda nam je neiscrpan izvor ideja za znanstvena istraživanja pa ne čudi činjenica da je pozamašan broj modela konačnih mješavina primijenjen upravo na Gaussove mješavine, a one su predmet proučavanja ovog diplomskog rada. Mnoge prirodne i društvene značnosti koriste navedene modele mješavina u analizi populacija kako bi mogli usmjeriti istraživanje na komponente, umjesto da zanemare utjecaj podgrupa na distribucijske parametre. Gaussove mješavine naziv su dobile po Gaussovoj distribuciji budući da su njihove komponente normalno distribuirane i *najčešće* ih u modelima ima konačno mnogo. U ovom ćemo radu dati pregled statističkih metoda koje se koriste u njihovom modeliranju. Ovisno o populaciji koju proučavamo, broj komponenata može biti unaprijed poznat ili ga je potrebno procijeniti, a tada se koriste složenije numeričke metode. Dvije su glavne sastavnice modela konačnih mješavina. Prva je procjena parametara vjerojatnosne distribucije svih komponenata mješavine i određivanje pripadnih težina. Težine komponenata još nazivamo proporcijama miješanja jer nam daju uvid u udjele koji svaka komponenta ima u mješavini. Drugi ključan stavak modeliranja je utvrđivanje kojoj komponenti mješavine pripada svaka od jedinki promatrane populacije. Već iz navedenog vidimo da su ove analize međusobno zavisne i da ih moramo proučavati istovremeno. Problematici pristupamo analitički, a princip modeliranja prilagođavamo u ovisnosti od promatranog slučaja. U poglavlju 1 konstruiramo model konačnih mješavina pomoću teorije frekvencionističkog i Bayesovskog zaključivanja, pritom koristeći standardne metode procjene parametara. Teorija matematičke statistike korištena u konstrukciji modela uvedena je potpoglavljem 1.1. Konačne mješavine formalno su definirane u potpoglavlju 1.2 dok je u potpoglavlju 1.3 razrađen teorijski pristup statističkom zaključivanju u modelima Gausso-

vih mješavina. Vidjet ćemo da matematički alat izveden iz statističkog zaključivanja nije dovoljan za rješavanje najzanimljivijeg slučaja modeliranja, a to je slučaj kada zaključke o pripadnosti komponentama donosimo na temelju neopaženih podataka. Rješenje ovog problema nije se, kao u slučajevima s poznatim pridruživanjima, nametnulo samo od sebe. Bio je potreban dugi niz godina, ali i razvoj računarstva kao znanosti da bismo danas taj problem mogli jednostavno riješiti pomoću računalnog algoritma. Poglavlje 2 ovog rada rezervirano je za proučavanje računalnog pristupa modeliranju Gaussovih mješavina. Algoritam koji se koristi za pridruživanje jedinki komponentama je algoritam maksimizacije očekivanja te ga definiramo u potpoglavlju 2.1. Kako princip modeliranja mješavina ne bi ostao na teorijskoj razini, u potpoglavlju 2.2 ilustriramo navedene metode analizirajući stvarnu populaciju.

Poglavlje 1

Model konačnih mješavina

Model konačnih mješavina¹ je statistički model koji pretpostavlja postojanje neopaženih grupa unutar promatrane populacije. Problem modeliranja konačnih mješavina se svodi na analizu grupa i pridruživanje podataka grupama. Za svaku grupu treba odrediti svojstvene značajke te ispitati postoji li teorijska vjerojatnosna distribucija koja *dovoljno dobro* opisuje podatke iz te grupe. Prvi znanstveni radovi na ovu temu nastali su krajem 19. stoljeća te je fokus istraživanja bio modeliranje grupa, dok se pridruživanju podataka grupama nije davao veći značaj. Međutim, činjenica je da je grupiranje neodvojivo od potpune analize mješavine jer je utjecaj pridruživanja podataka na statističko zaključivanje² o tim podacima evidentan. U novije vrijeme, modeli konačnih mješavina prepoznati su u mnogim znanostima upravo zbog informativnosti o pridruživanju podataka grupama. Mogućnosti primjene modela su razne, ali o tome je opsežnije pisano u poglavlju 2. Ovo je poglavlje predviđeno za teorijsko objašnjenje matematičkog aparata koji koristimo pri modeliranju. Uvest ćemo teoriju iz vjerojatnosti i statistike koju ćemo nadograditi kako bismo mogli provesti metode modeliranja različitih populacija konačnih mješavina.

1.1 Vjerojatnosni prostor i slučajne veličine

Definirajmo objekte potrebne za uvođenje vjerojatnosnog prostora na \mathbb{R} i navedimo teoriju matematičke statistike³ koju ćemo koristiti u modeliranju u Poglavlju 2.

Definicija 1.1.1.

I. Neka je $\Omega \subseteq \mathbb{R}$ neprazan skup događaja. Elementi skupa Ω su elementarni događaji.

¹eng. *finite mixture model*

²eng. *statistical inference*

³Definicije i teorijski rezultati preuzeti su iz Teorije vjerojatnosti [13], iz 2. i 9. poglavlja.

II. Neka je $\mathcal{F} \subseteq \Omega$ familija skupova na Ω . \mathcal{F} je σ -algebra ako vrijedi:

$$(S1) \emptyset \in \Omega$$

$$(S2) A \in \Omega \implies A^c \in \Omega$$

$$(S3) (A_n)_{n \in \mathbb{N}} \subseteq \Omega \implies \bigcup_{n=1}^{\infty} A_n \in \Omega$$

Uređeni par (Ω, \mathcal{F}) nazivamo **izmjeriv prostor**.

III. Neka je (Ω, \mathcal{F}) izmjerivi prostor i $\mathbb{P}: \mathcal{F} \rightarrow \mathbb{R}$ preslikavanje sa sljedećim svojstvima zovemo **vjerojatnost**:

$$(P1) \mathbb{P}(A) \geq 0, \text{ za svaki } A \in \Omega \quad \dots \text{ nenegativnost vjerojatnosti}$$

$$(P2) \mathbb{P}(\Omega) = 1 \quad \dots \text{ normiranost vjerojatnosti}$$

$$(P3) \mathbb{P}(\bigcup_{n=1}^{\infty} A_n) = \sum_{n=1}^{\infty} \mathbb{P}(A_n) \quad \dots \sigma\text{-aditivnost vjerojatnosti}$$

Uređenu trojku $(\Omega, \mathcal{F}, \mathbb{P})$ nazivamo **vjerojatnosni prostor**.

Definicija 1.1.2. Borelova σ -algebra na \mathbb{R} je najmanja σ -algebra na \mathbb{R} generirana familijom otvorenih skupova \mathcal{U} na \mathbb{R} .

$$\mathcal{B}_{\mathbb{R}} = \sigma(\mathcal{U}).$$

Definicija 1.1.3. Neka je $(\Omega, \mathcal{F}, \mathbb{P})$ vjerojatnosni prostor. Preslikavanje $X: \Omega \rightarrow \mathbb{R}$ je **slučajna varijabla** ako vrijedi $X^{-1}(\mathcal{B}_{\mathbb{R}}) \subseteq \mathcal{F}$.

Definicija 1.1.4. Neka je $(\Omega, \mathcal{F}, \mathbb{P})$ vjerojatnosni prostor i X slučajna varijabla na Ω . **Funkcija distribucije** slučajne varijable X je preslikavanje $F_X: \mathbb{R} \rightarrow [0, 1]$ definirano s:

$$F_X(x) = \mathbb{P}_X((-\infty, x]) = \mathbb{P}(X^{-1}((-\infty, x])) = \mathbb{P}(\omega \in \Omega : X(\omega) \leq x) = \mathbb{P}(X \leq x), \quad \forall x \in \mathbb{R}$$

Preslikavanje \mathbb{P}_X nazivamo **zakon razdiobe** slučajne varijable X .

Definicija 1.1.5. Slučajna varijabla X na $(\Omega, \mathcal{F}, \mathbb{P})$ je **diskretna** ako postoji prebrojivi skup $D \subseteq \mathbb{R}$ takav da vrijedi $\mathbb{P}(X \in D) = 1$.

Definicija 1.1.6. Slučajna varijabla X na $(\Omega, \mathcal{F}, \mathbb{P})$ je **apsolutno neprekidna** ako za njezinu funkciju distribucije F_X postoji nenegativna Borelova funkcija $f_X: \mathbb{R} \rightarrow \mathbb{R}$ takva da vrijedi

$$F_X(x) = \int_{-\infty}^x f_X(t) d\lambda(t), \quad x \in \mathbb{R}. \quad (1.1)$$

Funkciju f_X tada nazivamo **funkcija gustoće** slučajne varijable X .

Definicija 1.1.7. Neka je X slučajna varijabla na vjerojatnosnom prostoru $(\Omega, \mathcal{F}, \mathbb{P})$. Ako red $\sum_{\omega \in \Omega} X(\omega)\mathbb{P}(\omega)$ apsolutno konvergira, onda njegovu sumu zovemo **matematičko očekivanje slučajne varijable X** i označujemo s

$$\mathbb{E}[X] = \sum_{\omega \in \Omega} X(\omega)\mathbb{P}(\omega). \quad (1.2)$$

Ako je X apsolutno neprekidna slučajna varijabla, onda je njeno očekivanje dano s

$$\mathbb{E}[X] = \int_{\Omega} X(\omega)f_X(\omega) d\mathbb{P}(\omega). \quad (1.3)$$

Napomena 1.1.8. Integral u definicijama 1.1.6 i 1.1.7 je Lebesgueov integral funkcije f_X u odnosu na Lebesgueovu mjeru λ .

Definicija 1.1.9. Slučajna varijabla X ima **normalnu distribuciju** s parametrima $\mu \in \mathbb{R}$ i $\sigma \in \mathbb{R}^+$, $X \sim \mathcal{N}(\mu, \sigma^2)$, ako joj je funkcija gustoće dana s

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad x \in \mathbb{R}. \quad (1.4)$$

Definicija 1.1.10 (Višedimenzionalna normalna distribucija). Neka je $\mathbf{X} = (X_1, \dots, X_n)$ slučajni vektor, za $n \in \mathbb{N}$. Označimo s $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n) \in \mathbb{R}^n$ očekivanje vektora \mathbf{X} te sa $\Sigma = \text{Cov}(X_i, X_j) \in M_n(\mathbb{R})$ kovarijacijsku matricu od \mathbf{X} . Stavimo da je $\Sigma^{-1} = [\sigma_{ij}]$ pa uvodimo kvadratnu formu $Q_n: \mathbb{R}^n \rightarrow \mathbb{R}^+$ sljedećim pravilom preslikavanja

$$Q_n(\mathbf{x}) = Q_n(x_1, \dots, x_n) = \sum_{i=1}^n \sum_{j=1}^n \sigma_{ij}(x_i - \mu_i)(x_j - \mu_j) = (\mathbf{x} - \boldsymbol{\mu})\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})^\tau. \quad (1.5)$$

Kažemo da slučajni vektor \mathbf{X} ima **n -dimenzionalnu normalnu distribuciju**, $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ ako je funkcija gustoće od \mathbf{X} dana s

$$f_{\mathbf{X}}(\mathbf{x}) = (2\pi)^{-\frac{n}{2}} (\det \Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}Q_n(\mathbf{x})\right), \quad \mathbf{x} \in \mathbb{R}^n. \quad (1.6)$$

Primijetimo da je ovako definirana kvadratna forma zapravo kvadratna Mahalanobisova udaljenost opažanja \mathbf{x} od promatrane višedimenzionalne normalne distribucije, tj.

$$D_M(\mathbf{x}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})^\tau} = \sqrt{Q_n(\mathbf{x})}. \quad (1.7)$$

Kovarijacijska matrica Σ je po svojoj definiciji uvijek simetrična i pozitivno semidefinitna. U slučaju višedimenzionalne normalne distribucije slučajnog vektora \mathbf{X} , pripadna kovarijacijska matrica Σ je pozitivno definitna te je njezin inverz Σ^{-1} također simetrična i pozitivno definitna matrica.

Metode procjene distribucijskih parametara

Parametri promatrane vjerojatnosne distribucije su obično nepoznati te ih je potrebno odrediti statističkim metodama. Najčešće korištene metode procjene parametara su metoda momenata, metoda maksimalne vjerodostojnosti⁴ i Bayesova procjena. Ovdje uvodimo princip korištenja tih metoda u općenitom slučaju, a u potpoglavlju 1.3 dajemo njihove izvode za procjenu parametara konačne mješavine. Neka je Y slučajna varijabla čija distribucija F_Y dolazi iz familije distribucija $\mathcal{T}(\theta)$ na skupu parametara Θ . Želimo procijeniti nepoznate parametre distribucije F_Y .

Metoda momenata

Neka je $\mathbf{y} = (y_1, \dots, y_N)$ neko opažanje slučajne varijable Y , za $N \in \mathbb{N}$. Parametre distribucije F_Y procjenjujemo metodom momenata tako da izjednačimo uzoračke momente s teorijskim momentima razdiobe.

Definicija 1.1.11. *Teorijski m -ti moment je svojstvo funkcije distribucije slučajne varijable Y koje opisuje graf funkcije distribucije i definiramo ga s $\mathbb{E}[Y^m]$, gdje je $m \in \mathbb{N}$. Ako očekivanje $\mathbb{E}[Y^m]$ postoji, računamo ga formulom*

$$\mathbb{E}[Y^m] = \int_{\Omega} y^m dF_Y(y) = \int_{\Omega} y^m f_Y(y) d\mathbb{P}. \quad (1.8)$$

Označimo s M_m **uzorački m -ti moment** koji odgovara m -tom teorijskom momentu $\mathbb{E}[Y^m]$. Pretpostavimo da je potrebno odrediti $k \in \mathbb{N}$ parametara promatrane distribucije. Tada rješavamo sustav od k jednadžbi oblika

$$M_m = \mathbb{E}[Y^m], \quad m = 1, \dots, k. \quad (1.9)$$

Metoda maksimalne vjerodostojnosti

Uvedimo prvo definiciju funkcije vjerodostojnosti pa ćemo objasniti kako ovom metodom procijeniti nepoznate parametre distribucije.

Definicija 1.1.12. *Neka je $\mathbf{Y} = (Y_1, \dots, Y_N)$, za $N \in \mathbb{N}$ slučajni uzorak iz modela $\mathcal{P} = \{f(\cdot; \theta) : \theta \in \Theta\}$ i $\mathbf{y} = (y_1, \dots, y_N)$ neka njegova realizacija. Tada je **vjerodostojnost** definirana s*

$$L : \Theta \rightarrow \mathbb{R}, \quad L(\theta | \mathbf{y}) = L(\theta) = f_Y(\mathbf{y} | \theta) = \prod_{k=1}^N f(y_k | \theta) \quad (1.10)$$

⁴eng. *maximum likelihood estimation* — kraće **ML metoda**

Budući da je f_Y funkcija gustoće, ona je nenegativna funkcija pa možemo definirati **log-vjerodostojnost** kao logaritam funkcije vjerodostojnosti, tj.

$$l : \Theta \rightarrow \mathbb{R}, \quad l(\theta | \mathbf{y}) = l(\theta) = \log f_Y(\mathbf{y} | \theta) = \sum_{k=1}^N \log f(y_k | \theta) \quad (1.11)$$

Cilj procjene ML metodom je odrediti procjenitelj $\hat{\theta}$ parametara θ tako da određuje niz asimptotski efikasnih i konzistentnih korijena diferencijalne jednačbe

$$\partial l(\theta) / \partial \theta = 0. \quad (1.12)$$

Limes vjerojatnosti korijena jednačbe (1.12) teži u 1 pa ti korijeni odgovaraju lokalnom maksimumu u interioru prostora parametara.⁵ Definirajmo formalno procjenitelje spomenute u prethodnoj diskusiji. Iako procjenjujemo vektor parametara θ , zbog konteksta sljedećih definicija smatrat ćemo da procjenjujemo vrijednost nekog preslikavanja κ definirano na prostoru parametara.

Definicija 1.1.13. Procjenitelj $\hat{K} = K(Y)$ je **nepristrani procjenitelj** vrijednosti preslikavanja $\kappa(\theta)$ ako vrijedi

$$\mathbb{E}[\hat{K}] = \kappa(\theta), \quad \forall \theta \in \Theta. \quad (1.13)$$

Za definiciju efikasnog procjenitelja, potrebno je uvesti Fisherovu informaciju.

Definicija 1.1.14. Fisherova informacija je preslikavanje definirano s

$$I(\theta) = \mathbb{E} \left[\left(\frac{\partial l(Y, \theta)}{\partial \theta} \right)^2 \right]. \quad (1.14)$$

Definicija 1.1.15. Procjenitelj $\hat{K} = \kappa(\hat{\theta})$ vrijednosti funkcije parametara $\kappa(\theta)$ je **efikasan** ako je nepristran za $\kappa(\theta)$, ima konačan drugi moment, te vrijedi

$$\text{Var}(\hat{K}) = c \frac{(\kappa'(\theta))^2}{I(\theta)}, \quad (1.15)$$

za neki $c > 0$. Procjenitelje koji uvjet efikasnosti postižu asimptotski nazivamo **asimptotski efikasni procjenitelji**.

Definicija 1.1.16. Kažemo da je niz procjenitelja $(\hat{K}_n)_{n \in \mathbb{N}}$ **konzistentan** za vrijednosti funkcije parametara $\kappa(\theta)$, ako konvergira⁶ prema $\kappa(\theta)$.

Može se pokazati da takav niz procjenitelja postoji kada su zadovoljeni Cramerovi uvjeti regularnosti koji su, dakle, nužan uvjet da smijemo koristiti ML metodu procjene.

⁵Detaljnije obrazloženje možemo pronaći u članku [The EM Algorithm](#) [9].

⁶Konvergencija može biti po vjerojatnosti, u srednjem reda $p \in \mathbb{N}$ ili gotovo svuda. (kraće pišemo g.s.)

Uvjeti regularnosti:

- Parametar koji procjenjujemo je iz unutrašnjosti parametarskog prostora, $\theta \in \text{Int } \Theta$
- Log–vjerodostojnost je tri puta diferencijabilna na okolini stvarnog parametra θ u parametarskom prostoru Θ , za g.s. $y \in \Omega$
- Treća derivacija log–vjerodostojnosti je ograničena, $\left| \frac{\partial^3 l(\theta|y)}{\partial \theta^3} \right| < M$, za neki $M < \infty$, tj. Fisherova informacija je dobro definirana.

Definicija 1.1.17 (MLE procjenitelj⁷). *Statistika $\hat{\theta} = \hat{\theta}(\mathbf{Y})$ je procjenitelj maksimalne vjerodostojnosti za parametar $\theta^* \in \Theta$ ako vrijedi*

$$L(\hat{\theta} | \mathbf{Y}) = \max_{\theta \in \Theta} L(\theta | \mathbf{Y}) = \max_{\theta \in \Theta} \prod_{k=1}^N f(y_k | \theta). \quad (1.16)$$

Jedna od tehnika određivanja parametra θ koji maksimizira funkciju vjerodostojnosti L je traženje njenih stacionarnih točaka i klasifikacija na sedlastu točku, globalni ili lokalni ekstrem. Funkcija vjerodostojnosti je produkt funkcija gustoća modela \mathcal{P} pa, zbog diferenciranja, račun određivanja točaka ekstrema funkcije postaje zahtjevan. Ovisno o promatranoj funkciji gustoće, moguć je slučaj da se problem ne može riješiti analitički. Funkcija logaritmiranja je neopadajuća na svojoj domeni, stoga funkcije (1.10) i (1.11) postižu maksimum za isti θ , tj.

$$l(\hat{\theta} | \mathbf{Y}) = \max_{\theta \in \Theta} l(\theta | \mathbf{Y}) = \max_{\theta \in \Theta} \sum_{k=1}^N \log(f(y_k | \theta)). \quad (1.17)$$

Zbog jednostavnijeg računa, u nastavku rada ML procjenu parametara distribucije provodimo tako da odredimo MLE procjenitelj funkcije log–vjerodostojnosti. Račun je osobito jednostavan u slučaju da promatrana distribucija pripada eksponencijalnoj familiji distribucija. U procjeni parametara Bayesovom metodom, potrebno nam je statističko zaključivanje bazirano na opaženim podacima iz modela, a izvedeno je iz sljedećeg teorema.

Teorem 1.1.18. (*Bayesov teorem*⁸) *Neka je θ apriorno nepoznati skup parametara promatrane distribucije i Y skup mogućih realizacija opažanja. Funkciju gustoće distribucije uvjetno na opažanje računamo formulom*

$$f(\theta | Y) = \frac{f(Y, \theta)}{f(Y)} = \frac{f(Y | \theta)f(\theta)}{f(Y)}. \quad (1.18)$$

⁷eng. *maximum likelihood estimator*

⁸Formulacija Bayesovog teorema preuzeta iz *Bayesian Statistics and Marketing* [12], iz 2. poglavlja.

U formuli (1.18) uveli smo *apriornu gustoću distribucije* $f(Y)$ i *apriornu vjerodostojnost* $f(Y | \theta)$. Funkciju $f(\theta | Y)$ zovemo *aposteriorna gustoća distribucije*. Bayesovsko zaključivanje koristi se u slučajevima kada želimo izvesti zaključak o mogućim vrijednostima parametra θ na temelju podataka opaženih u uzorku. Primijetimo da je nazivnik u jednadžbi (1.18) konstantan jer ne ovisi o θ . Aposteriornu gustoću stoga analiziramo zanemarujući nazivnik u jednadžbi (1.18), tj. prema relaciji proporcionalnosti

$$f(\theta | Y) \propto f(Y | \theta)f(\theta). \quad (1.19)$$

Konstanta proporcionalnosti je marginalna gustoća uzorka koju računamo formulom

$$f(Y) = \int_{\Theta} f(Y, \hat{\theta})d\hat{\theta} = \int_{\Theta} f(Y | \hat{\theta})f(\hat{\theta})d\hat{\theta}. \quad (1.20)$$

Bayesova metoda

Neka je $f(\theta)$ apriorna distribucija nepoznatog parametra θ . Označimo s $\hat{\theta}$ procjenitelj od θ te uvedimo *funkciju gubitka*⁹ $\lambda = \lambda(\theta, \hat{\theta})$. Funkcija gubitka je svojevrsna mjera udaljenosti procijenjenog parametra od njegove stvarne vrijednosti. Često se koriste mjere apsolutnog i srednjekvadratnog odstupanja koje računamo formulama:

$$\lambda(\theta, \hat{\theta}) = \begin{cases} |\theta - \hat{\theta}|, & \text{apsolutno odstupanje} \\ (\theta - \hat{\theta})^2, & \text{srednjekvadratno odstupanje} \end{cases} \quad (1.21)$$

Bayesova funkcija rizika definirana je kao očekivanje funkcije gubitka, $R(\theta) = \mathbb{E}[\lambda(\theta, \hat{\theta})]$.

Definicija 1.1.19. Procjenitelj $\hat{\theta}$ je *Bayesov procjenitelj* ako minimizira Bayesovu funkciju rizika na prostoru parametara Θ .

Prije no što nekom od metoda odredimo parametre distribucije, trebalo bi razmisliti je li uzorak na osnovu čijih opažanja donosimo zaključke dovoljno dobar. Je li moguće iz prirodnih pretpostavki promatranog modela izvesti dobre i *točne* zaključke? Odgovore na ta pitanja daje nam *raspoznatljivost modela*. Raspoznatljivost modela je svojstvo da iz beskonačno mnogo provedenih opažanja možemo jednoznačno odrediti stvarne teorijske vrijednosti inherentnih parametara modela. Svojstvo jednoznačnog određivanja parametara ekvivalentno je svojstvu da različite vrijednosti parametara modela generiraju različite vjerojatnosne distribucije opaženih podataka. To znači da je na raspoznatljivim mješavinama moguće provesti točnu procjenu parametara i shodno tome statističko zaključivanje na temelju opaženih podataka.

⁹eng. *loss function*; u literaturi se najčešće označava s L , ali u ovom je radu ta oznaka već rezervirana za funkciju vjerodostojnosti

Definicija 1.1.20. Neka je \mathcal{P} parametarska familija distribucija indeksirana parametrom $\vartheta \in \Theta$ koji je definiran na prostoru \mathcal{Y} . Familija \mathcal{P} je **raspoznatljiva**¹⁰ ako vrijedi

$(\forall \vartheta, \vartheta^* \in \Theta)$ parametri ϑ i ϑ^* definiraju istu vjerojatnost na \mathcal{Y} ako i samo ako je $\vartheta = \vartheta^*$.

Raspoznatljivost familije distribucija možemo primijeniti i na funkcije gustoće pa vrijedi

$$f(\mathbf{y} | \vartheta) = f(\mathbf{y} | \vartheta^*), \text{ za g.s. } \mathbf{y} \in \mathcal{Y} \implies \vartheta = \vartheta^*. \quad (1.22)$$

Skup $\mathcal{U}(\vartheta) = \{\vartheta^* \in \Theta: f(\mathbf{y} | \vartheta^*) = f(\mathbf{y} | \vartheta), \text{ za g.s. } \mathbf{y} \in \mathcal{Y}\} \subseteq \Theta$ je **neraspoznatljiv skup** ako nije jednočlan, a pripadna familija distribucija $\mathcal{T}_{\mathcal{U}}(\vartheta)$ je tada **neraspoznatljiva**.

1.2 Konačne mješavine

Kako bismo konstruirali model koji najbolje opisuje populaciju koju proučavamo, trebaju nam početne pretpostavke o prirodi tih podataka, o mogućoj pripadnosti nekoj vjerojatnosnoj distribuciji ili o postojanju smislene podjele na podpopulacije koje još nazivamo grupe ili komponente. Komponente bi trebale dijeliti populaciju na način da podaci budu homogeni unutar i heterogeni izvan njih. Za proizvoljni slučajni uzorak opažamo neku realizaciju i želimo znati kojoj komponenti populacije ona pripada, ali tu informaciju ne dobivamo izravno. Problem je riješen uvođenjem modela mješavina koje još nazivamo modelima miješanih gustoća¹¹. Prvi rad u kojem je spomenuto pridruživanje opažanja nekoj od grupa populacije napisao je W. Feller 1943. godine u članku *On a General Class of "Contagious" Distributions* [3] u kontekstu modeliranja zaraze populacije ličinki. U ovom radu modeliramo konačne Gaussove mješavine¹², tj. pretpostavljamo da slučajni uzorak pripada populaciji s konačno mnogo normalno distribuiranih podpopulacija. Postoje modeli miješanih gustoća koji imaju beskonačno mnogo komponenata, a komponente općenito mogu pripadati bilo kojoj diskretnoj ili neprekidnoj razdiobi. Svi izvodi i zaključivanja se lako prilagođavaju diskretnom slučaju korištenjem brojeće mjere. Glavni cilj analiziranja konačnih mješavina je određivanje parametara distribucije svake komponente kako bi se daljnje statističko zaključivanje moglo prilagoditi predmetnoj populaciji.

Definicija 1.2.1. Neka je $\mathbf{Y} = (Y_1, \dots, Y_r)$ r -dimenzionalni slučajni vektor definiran u vjerojatnosnom prostoru $(\mathcal{Y}, \mathcal{F}, \mathbb{P})$, gdje je $r \in \mathbb{N}$. Promatramo $N \in \mathbb{N}$ nezavisnih opažanja slučajnog vektora \mathbf{Y} koja se sastoje od r značajki. Označimo ih s $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$, gdje je $\mathbf{y}_i = (y_{i1}, \dots, y_{ir})$. Kažemo da \mathbf{Y} dolazi iz **distribucije konačne mješavine** ako njena vjerojatnosna funkcija gustoće $f_{\mathbf{Y}}$ ima formu gustoće mješavine

¹⁰eng. *identifiable*; definicija preuzeta iz *Finite Mixture and Markov Switching Models* [4], str. 14.–15.

¹¹eng. *mixture models*

¹²Definicije i izvodi modela konačnih mješavina preuzeti su iz *Finite Mixture and Markov Switching Models* [4], iz 2. poglavlja.

$$f_Y(\mathbf{y}) = \eta_1 f_1(\mathbf{y}) + \cdots + \eta_K f_K(\mathbf{y}), \quad \text{za svaki } \mathbf{y} \in \mathbf{Y}, \quad (1.23)$$

gdje je $K \in \mathbb{N}$ broj komponenata, a f_k su komponentne gustoće. Vektor $\boldsymbol{\eta} = (\eta_1, \dots, \eta_K)$ je distribucija težina, gdje je težina komponente \mathcal{K}_k označena parametrom η_k . Komponente vektora $\boldsymbol{\eta}$ su nenegativne i normirane, tj. takve da vrijedi

$$\eta_k \geq 0, \quad \forall k \in \{1, \dots, K\} \quad (1.24)$$

$$\eta_1 + \cdots + \eta_K = 1. \quad (1.25)$$

Parametar težine k -te komponente možemo, stoga, interpretirati kao vjerojatnost da slučajno odabrana jedinka Y iz te populacije pripada komponenti \mathcal{K}_k . Stavimo $\mathbb{P}(Y \in \mathcal{K}_k) = \eta_k$. Zbog svojstva (1.25) normiranosti vektora težina, ne procjenjujemo svih K težina komponenata već je dovoljno procijeniti njih $K - 1$.

Neka je slučajni vektor $\mathbf{Y} = (Y_1, \dots, Y_r)$ distribuiran iz mješavine K višedimenzionalnih normalnih distribucija, tj. $Y_k \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$, za $k \in \{1, \dots, K\}$. Tada je $\boldsymbol{\mu}_k \in \mathbb{R}^r$, $\Sigma_k \in M_r(\mathbb{R})$ te je komponentna funkcija gustoće

$$f_k(\mathbf{y}) = (2\pi)^{-\frac{r}{2}} (\det \Sigma_k)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_k) \Sigma_k^{-1} (\mathbf{y} - \boldsymbol{\mu}_k)^\tau\right), \quad \mathbf{y} \in \mathbb{R}^r. \quad (1.26)$$

Parametri Gaussove mješavine koje procjenjujemo su $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K, \eta_1, \dots, \eta_K)$, gdje je $\boldsymbol{\theta}_k = (\boldsymbol{\mu}_k, \Sigma_k)$. Vektor očekivanja od \mathbf{Y} je $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)$, gdje je $\boldsymbol{\mu}_k = (\mu_{k1}, \dots, \mu_{kr})$ očekivanje komponente \mathcal{K}_k , za $k \in \{1, \dots, K\}$. U definiciji 1.1.10 smo uveli oznaku Σ_k za kovarijacijsku matricu k -te komponente mješavine pa vektor kovarijacijskih matrica komponenata označimo sa $\boldsymbol{\Sigma} = (\Sigma_1, \dots, \Sigma_K)$.

Pri modeliranju konačnih mješavina, osim procjene parametara modela, želimo odrediti pripadnost opažanja pojedinim komponentama. Pritom se susrećemo s opaženim i skrivenim varijablama.¹³

Definicija 1.2.2. *Manifestna ili opažena varijabla je slučajna varijabla čiju realizaciju opažamo iz uzorka i koja upućuje na prisustvo neke latentne¹⁴ varijable.*

Definicija 1.2.3. *Latentna ili skrivena varijabla je slučajna varijabla čiju realizaciju ne opažamo izravno. Njeno postojanje se pretpostavlja s ciljem objašnjavanja manifestnih varijabli. Model koji objašnjava opažene varijable u terminima latentnih, naziva se **model latentnih varijabli**.*

¹³Definicija preuzeta iz materijala kolegija Strojno učenje [18]

¹⁴lat. *lateo* – biti skriven

Komponente modela konačnih mješavina $\mathcal{K}_1, \dots, \mathcal{K}_K$ indeksiramo konačnim diskretnim skupom indikatora $\mathcal{I} = \{1, \dots, K\}$, gdje je $K \in \mathbb{N}$. Pretpostavimo da za slučajni vektor \mathbf{Y}_i iz promatrane populacije apriorno znamo da pripada komponenti \mathcal{K}_k . Njegova je funkcija gustoće $f(\mathbf{y}_i | \boldsymbol{\theta}_k)$ tada dana uvjetno na grupu k , gdje je $\boldsymbol{\theta}_k$ vektor parametara distribucije komponente \mathcal{K}_k . Zajedničku gustoću sada možemo izraziti formulom

$$f(\mathbf{y}_i, s_i = k) = f(\mathbf{y}_i | s_i = k)\mathbb{P}(s_i = k) = f(\mathbf{y}_i | s_i = k)\eta_k. \quad (1.27)$$

Uveli smo pojam modela latentnih varijabli, a model konačnih mješavina ćemo promatrati kao model hijerarhijske latentne varijable. Distribucija opažanja $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$ ovisi o realizaciji $s = (s_1, \dots, s_N)$ latentnog diskretnog vektora $\mathbf{S} = (S_1, \dots, S_N)$. Slučajni vektor \mathbf{S} nazivamo *pridruživanje* jer sadrži informaciju o pridruživanju svakog od N opažanja odgovarajućoj komponenti. Skup $\{\mathbf{y}, s\}$ zovemo *potpun skup podataka*, dok je $\{\mathbf{y}\}$ *nepotpun skup podataka*. Kod modeliranja konačnih mješavina, autori članka *Finite Mixture Models* [7] zaključuju da vektor pridruživanja \mathbf{S} možemo gledati kao N latentnih jednako distribuiranih slučajnih varijabli. Potrebna je pretpostavka da su pridruživanja svakog od opažanja međusobno nezavisna te se može pokazati da je njihova razdioba multinomijalna,

$$S_k \sim \text{Mult}_K(1, \boldsymbol{\eta}), \quad \forall k \in \{1, \dots, N\}.$$

Marginalna gustoća slučajne veličine \mathbf{Y} je, u slučaju da apriorno nemamo pridruživanje \mathbf{S} , dana sljedećom linearnom kombinacijom funkcija gustoća:

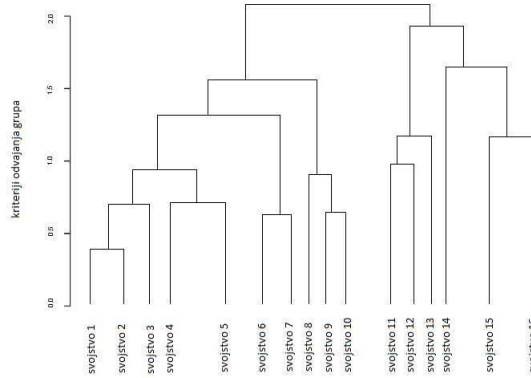
$$f(\mathbf{y}) = \sum_{k=1}^K f(\mathbf{y}, s = k) = \eta_1 f(\mathbf{y} | \boldsymbol{\theta}_1) + \dots + \eta_K f(\mathbf{y} | \boldsymbol{\theta}_K), \quad (1.28)$$

gdje je s oznaka komponente kojoj opažanje \mathbf{y} pripada. U slučaju Gaussove mješavine, komponentna funkcija gustoće $f(\mathbf{y} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ je gustoća višedimenzionalne normalne distribucije $\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ uvedene jednažbom (1.26).

Hijerarhijski modeli

U ovom radu nepoznate parametre miješane gustoće procjenjujemo pomoću Bayesovskog zaključivanja i maksimizacijom log-vjerodostojnosti. Model Gaussove mješavine pripada klasi hijerarhijskih modela. To su modeli grupiranja u kojima sortiramo prema hijerarhijskoj relaciji po slojevima, tako da je za svaki par slojeva poznat uređaj među njima. Grafički možemo odnose između grupa u modelu prikazati dendrogramom. Hijerarhijske modele prema načinu grupiranja dijelimo na aglomerativne i divizivne. Aglomerativni modeli u prvom koraku grupiraju svako opažanje u zasebnu komponentu. Komponente se, kroz iteracije algoritma grupiranja, postepeno stapaju u zajedničke grupe te je na kraju njihov broj manji nego na početku. Divizivni modeli imaju obrnutu logiku i kreću od jedne zajedničke grupe koju dijele na komponente. Broj komponenata divizivnog modela je po

završetku grupiranja veći nego što je bio na početku. U poglavlju 2 objasnit ćemo detaljno algoritam maksimizacije očekivanja¹⁵ koji je primjer aglomerativnog hijerarhijskog modela i koji koristimo za grupiranje konačnih mješavina.



Slika 1.1: dendrogram

Agglomerativni model prikazan ovim dendrogramom možemo interpretirati tako da postoje tri osnovne podgrupe populacije te da je svojstvo 14 netipična vrijednost.

Definicija 1.2.4. *Bayesov hijerarhijski model* je statistički model procjene parametara aposteriorne vjerojatnosne distribucije korištenjem Bayesovskog zaključivanja na apriornim vjerojatnosnim distribucijama prvog i drugog sloja modela.

U prvom sloju dana je zajednička apriorna distribucija opažanja uvjetno na apriorno poznatu realizaciju pridruživanja s :

$$f(\mathbf{y} | \mathbf{s}, \boldsymbol{\vartheta}) = \prod_{i=1}^N f(\mathbf{y}_i | s_i, \boldsymbol{\vartheta}) = \prod_{i=1}^N f(\mathbf{y}_i | \boldsymbol{\theta}_{S_i}). \quad (1.29)$$

Drugi sloj modela daje funkciju gustoće zajedničke distribucije $f(\mathbf{s} | \boldsymbol{\vartheta})$ svih, u tom trenu nepoznatih realizacija s i to je diskretna distribucija duž rešetke \mathcal{S}_K .

$$\mathcal{S}_K = \{(S_1, \dots, S_N) : S_i \in \{1, \dots, K\}, i = 1, \dots, N\} \quad (1.30)$$

Uz pretpostavku da su indikatori s_1, \dots, s_N međusobno nezavisni, imamo da je vjerojatnost $\mathbb{P}(s_i = k | \boldsymbol{\eta}) = \eta_k$ pa je zajednička funkcija distribucije uvjetno na poznate težine

¹⁵eng. *expectation maximization algorithm*, kraće **EM algoritam**

komponentata dana izrazom

$$f(\mathbf{s} | \boldsymbol{\eta}) = \prod_{i=1}^N f(s_i | \boldsymbol{\eta}). \quad (1.31)$$

Konačno, koristimo Bayesovo pravilo (1.18) kao rezultat Bayesovog teorema 1.1.18 da izvedemo aposteriornu gustoću zajedničke distribucije $f(\mathbf{s} | \boldsymbol{\vartheta}, \mathbf{y})$ pomoću vjerodostojnosti $f(\mathbf{y} | \mathbf{s}, \boldsymbol{\vartheta})$ i apriorne gustoće zajedničke distribucije $f(\mathbf{s} | \boldsymbol{\vartheta})$:

$$f(\mathbf{s} | \boldsymbol{\vartheta}, \mathbf{y}) \propto f(\mathbf{y} | \mathbf{s}, \boldsymbol{\vartheta})f(\mathbf{s} | \boldsymbol{\vartheta}). \quad (1.32)$$

1.3 Statističko zaključivanje u modelu Gaussovih mješavina

U ovom potpoglavlju izvodimo statističko zaključivanje za model konačnih mješavina, na način da će, uz račun za općeniti model, biti izveden i rezultat za konačnu Gaussovu mješavinu. Pretpostavka za oba modela je da se mješavine sastoje od $K \in \mathbb{N}$ distribucija iz predmetne familije distribucije. Neka je $\mathcal{T}(\boldsymbol{\theta})$ familija distribucija na skupu parametara Θ i neka je $\mathcal{G}(\boldsymbol{\theta}) \subseteq \mathcal{T}(\boldsymbol{\theta})$ familija normalnih distribucija. Označimo s $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$, $N \in \mathbb{N}$ promatrane realizacije slučajnog uzorka iz konačne mješavine. Apriorne funkcije gustoće $f(\mathbf{y} | \boldsymbol{\theta})$ indeksirane su parametrom $\boldsymbol{\theta} \in \Theta$ te ih računamo po formuli

$$f(\mathbf{y}_i | \boldsymbol{\vartheta}) = \sum_{k=1}^K \eta_k f_k(\mathbf{y}_i | \boldsymbol{\theta}_k). \quad (1.33)$$

U slučaju modeliranja višedimenzionalne Gaussove mješavine, komponentna funkcija gustoće ima formulu (1.26), a njezin parametar koji procjenjujemo je $\boldsymbol{\theta}_k = (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. Problemu procjene parametara distribucije za svaku od komponenata pristupamo u tri koraka, kao što je navedeno u potpoglavlju 1.2:

1. Određivanje broja komponenata mješavine, $K \in \mathbb{N}$.
2. Procjena parametara $\boldsymbol{\theta}$ komponenata mješavine i vektora težine $\boldsymbol{\eta}$.
3. Pridruživanje svakog opažanja iz uzorka nekoj komponenti te zaključivanje na temelju latentnog indikatora \mathbf{S} .

Implementirajući model konačnih mješavina na stvarni poslovni ili znanstveni problem, najčešće na početku studije nemamo informaciju o broju komponenata populacije

koju modeliramo te ga moramo procijeniti. Neke od često korištenih metoda su testiranje prikladnosti modela¹⁶ dobivenog metodom momenata te Bayesovsko zaključivanje uz pretpostavku neodređenosti modela. Nakon što je broj komponenta populacije određen, krećemo s procjenom parametara mješavine i pridruživanjem opažanja. Po završetku modeliranja, potrebno je izračunati utjecaj pretpostavljenog broja komponenta K na pogrešku u zaključivanju. Pritom se za ocjenu pogreške najčešće koriste Akaikeov, odnosno Bayesov informacijski kriterij (procjenitelj pogreške označavamo AIC, odn. BIC). U nastavku rada pretpostavljamo da je broj komponenta promatrane populacije unaprijed poznat i označimo ga s $K \in \mathbb{N}$. Za procjenu distribucijskih parametara komponenta mješavine koristimo EM algoritam. Opažanja pridružujemo odgovarajućim komponentama mješavine, uz dani procijenjeni vektor $\boldsymbol{\vartheta}$, tako da izračunamo aposteriornu vjerojatnost zajedničke funkcije distribucije $f(\mathbf{s} | \boldsymbol{\vartheta})$. Svakom od $N \in \mathbb{N}$ opažanja dodjeljujemo oznaku one komponente čija je vjerojatnost $f(s_i | \boldsymbol{\vartheta}_j)$ najveća, za sve indekse $i \in \{1, \dots, N\}$, $j \in \{1, \dots, K\}$. Statističko zaključivanje moramo prilagoditi dostupnim informacijama o populaciji koju modeliramo, stoga je ovo potpoglavlje podijeljeno u tri dijela.¹⁷ Prvo analiziramo najjednostavniji slučaj, kada su poznati parametri $\boldsymbol{\vartheta}$ te je potrebno pridružiti opažanja komponentama. Zatim obrađujemo procjenu parametara u slučaju kada je poznato pridruživanje vektora opažanja. Zadnji je generalizirani slučaj kada ne postoje apriorne pretpostavke pa je tako i pridruživanje opažanja apriorno nepoznato.

Pridruživanje opažanja u slučaju poznatih komponenta

Pretpostavimo da je poznat $\boldsymbol{\vartheta}$, tj. da su poznati vektor težina $\boldsymbol{\eta}$ i parametri distribucije $\boldsymbol{\theta}$ konačne mješavine koja se sastoji od K komponenta. Tada je apriorna funkcija gustoće $f(\mathbf{y} | \boldsymbol{\theta})$ egzaktno određena formulom (1.33). Neka je $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$ vektor opažanja slučajnog uzorka iz promatrane mješavine. Budući da su poznati parametri komponenta modela, preostalo je još samo odrediti pridruživanje vektora \mathbf{y} komponentama. Problem pridruživanja opažanja riješen je korištenjem Bayesovog pravila (1.18). Uveli smo vektor pridruživanja $\mathbf{s} = (s_1, \dots, s_N)$ kao realizaciju diskretnog slučajnog vektora indikatora s vrijednostima u skupu $\{1, \dots, K\}$. Realizacija slučajne varijable S_i označava kojoj komponenti pripada opažanje \mathbf{y}_i , tj.

$$s_i = k \cdot \mathbb{1}_{\{\mathbf{y}_i \in \mathcal{K}_k\}}. \quad (1.34)$$

Pridruživanje vektora opažanja \mathbf{y} ekvivalentno je pridruživanju svakog pojedinog opažanja u slučaju kada su poznati svi parametri mješavine. Problem je složeniji u slučaju da su parametri neke komponente nepoznati, ali je rješiv korištenjem Bayesovog klasifikatora. U

¹⁶eng. *goodness of fit test*

¹⁷Svi izvodi potpoglavlja 1.3 prate konstrukciju zaključivanja iz knjige *Finite Mixture and Markov Switching Models* [4] iz 2. poglavlja.

slučaju pridruživanja jednog opažanja y_i , odmah krećemo s računanjem aposteriorne vjerojatnosti događaja $\{s_i = k\}$, $k \in \mathcal{I}$ uvjetno na opažanje $\{Y = y_i\}$ i na poznavanje distribucije mješavine. Vjerojatnost računamo formulom

$$\mathbb{P}(s_i = k | y_i, \boldsymbol{\vartheta}) = \frac{\mathbb{P}(Y = y_i | s_i = k, \boldsymbol{\vartheta})\mathbb{P}(s_i = k | \boldsymbol{\vartheta})}{\sum_{j=1}^K \mathbb{P}(Y = y_i | s_i = j, \boldsymbol{\vartheta})\mathbb{P}(s_i = j | \boldsymbol{\vartheta})}. \quad (1.35)$$

Pritom je poznata apriorna vjerojatnost da opažanje y_i pripada komponenti \mathcal{K}_k i ona je prema definiciji 1.2.1 jednaka $\mathbb{P}(s_i = k | \boldsymbol{\vartheta}) = \eta_k$. Uvrstimo η_k u jednadžbu (1.35):

$$\mathbb{P}(s_i = k | y_i, \boldsymbol{\vartheta}) = \frac{\mathbb{P}(Y = y_i | s_i = k, \boldsymbol{\vartheta})\eta_k}{\sum_{j=1}^K \mathbb{P}(Y = y_i | s_i = j, \boldsymbol{\vartheta})\eta_j} = h_{ik}. \quad (1.36)$$

Primijetimo da izraz u nazivniku ne ovisi o konkretnom k . Za izvedenu aposteriornu vjerojatnost pridruživanja i -tog opažanja k -toj komponenti je uvedena oznaka h_{ik} . Ova vjerojatnost naziva se još i *odgovornost*¹⁸. Budući da opažanje y_i alociramo onoj komponenti \mathcal{K}_k čija je vjerojatnost (1.36) najveća, dovoljno je, kao u izvodu Bayesovog teorema 1.1.18, promatrati relaciju

$$\mathbb{P}(s_i = k | y_i, \boldsymbol{\vartheta}) \propto \mathbb{P}(Y = y_i | s_i = k, \boldsymbol{\vartheta})\eta_k. \quad (1.37)$$

Klasifikacijski model koji ovdje koristimo je *Bayesov naivni klasifikator*¹⁹. Radi se o parametarskom modelu pridruživanja opažanja komponentama koji koristi Bayesovsko zaključivanje, a termin *naivni* dolazi od pretpostavke da su sva opažanja međusobno nezavisna uvjetno na promatranu komponentu \mathcal{K} , tj. vrijedi

$$\mathbb{P}(y_i | y_j, \mathcal{K}) = \mathbb{P}(y_i | \mathcal{K}), \quad \forall i, j \in \{1, \dots, N\}, \quad i \neq j. \quad (1.38)$$

Izvedimo sad postupak pridruživanja čitavog vektora opažanja $\mathbf{y} = (y_1, \dots, y_N)$ koji ćemo zvati *zajedničko pridruživanje*. Zanima nas vjerojatnost događaja $\{S_1 = s_1, \dots, S_N = s_N\}$, za sva moguća pridruživanja (s_1, \dots, s_N) svakog od N opažanja nekoj od K komponenata. Relacija vjerojatnosti pridruživanja jednog opažanja (1.37) poopćena na čitav uzorak je

$$f(\mathbf{s} | \mathbf{y}, \boldsymbol{\vartheta}) \propto f(\mathbf{y} | \mathbf{s}, \boldsymbol{\vartheta})f(\mathbf{s} | \boldsymbol{\vartheta}), \quad (1.39)$$

gdje je $f(\mathbf{s} | \boldsymbol{\vartheta})$ apriorna gustoća zajedničke distribucije vektora indikatora $\mathbf{S} = (S_1, \dots, S_N)$ trenutno nepoznatih pridruživanja i ona je poznata prije određivanja vektora opažanja \mathbf{y} . Vjerojatnost $f(\mathbf{y} | \mathbf{s}, \boldsymbol{\vartheta})$ je vjerodostojnost uvedena u relaciji (1.32) i ona je gustoća uzoračke distribucije čitavog vektora opažanja \mathbf{y} pod uvjetom poznatog pridruživanja \mathbf{s} i poznatih komponenata mješavine $\boldsymbol{\vartheta}$. Koristimo pretpostavku nezavisnosti slučajne varijable \mathbf{Y} koja pripada promatranoj mješavini pa onda vjerodostojnost računamo formulom

¹⁸eng. *responsibility*

¹⁹eng. *naïve Bayes' classifier*; Teorija je izvedena iz materijala kolegija Strojno učenje [18].

$$f(\mathbf{y} | \mathbf{s}, \boldsymbol{\vartheta}) = f(\mathbf{y} | \mathbf{s}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K) = \prod_{i=1}^N f(y_i | \boldsymbol{\theta}_{s_i}). \quad (1.40)$$

Uz pretpostavku da su pridruživanja opažanja također međusobno nezavisna, apriorna gustoća zajedničke distribucije je onda

$$f(\mathbf{s} | \boldsymbol{\vartheta}) = \prod_{i=1}^N f(s_i | \boldsymbol{\vartheta}). \quad (1.41)$$

Konačno imamo sve potrebno za primjenu Bayesovog pravila (1.18) i izračun aposteriorne gustoće zajedničke distribucije $f(\mathbf{s} | \boldsymbol{\vartheta}, \mathbf{y})$

$$f(\mathbf{s} | \boldsymbol{\vartheta}, \mathbf{y}) = \prod_{i=1}^N f(s_i | \boldsymbol{\vartheta}, y_i). \quad (1.42)$$

Prema diskusiji s početka ove cjeline, aposteriornu gustoću zajedničke distribucije vektora indikatora \mathbf{S} uvjetno na poznate parametre mješavine $\boldsymbol{\vartheta}$ i vektor opažanja \mathbf{y} računamo kao produkt nezavisnih pojedinačnih vjerojatnosti.

Procjena parametara u slučaju poznatog pridruživanja opažanja

Pretpostavimo sad da su parametri komponenta mješavine kao i vektor težina $\boldsymbol{\eta}$ nepoznati. U ovoj cjelini procjenjujemo parametre konačne mješavine (1.33) za slučaj kada uz proizvoljno opažanje $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$ slučajnog vektora \mathbf{Y} , opažamo i odgovarajuće pridruživanje $\mathbf{s} = (s_1, \dots, s_N)$. Procjenu parametara provodimo na potpunom skupu podataka $\{\mathbf{y}, \mathbf{s}\}$ i to dvjema standardnim metodama: metodom maksimizacije vjerodostojnosti i Bayesovskim zaključivanjem. Obje metode su temeljene na računu s **funkcijom potpune vjerodostojnosti**²⁰ $L^c(\mathbf{y}, \mathbf{s} | \boldsymbol{\vartheta})$ koja odgovara uzoračkoj distribuciji potpunog skupa podataka uvjetno na nepoznati parametar $\boldsymbol{\vartheta}$. Vjerojatnost $L^c(\mathbf{y}, \mathbf{s} | \boldsymbol{\vartheta})$ izvodimo koristeći hijerarhijski model latentne varijable uveden u potpoglavlju 1.2 i to iz jednadžbi (1.29) i (1.31). Tražena vjerojatnost je onda

$$L^c(\mathbf{y}, \mathbf{s} | \boldsymbol{\vartheta}) = f(\mathbf{y} | \mathbf{s}, \boldsymbol{\vartheta})f(\mathbf{s} | \boldsymbol{\vartheta}) = \prod_{i=1}^N f(y_i | s_i, \boldsymbol{\vartheta})f(s_i | \boldsymbol{\vartheta}). \quad (1.43)$$

Primijetimo da je $f(\mathbf{y}_i | s_i = k, \boldsymbol{\vartheta}) = f(\mathbf{y}_i | \boldsymbol{\theta}_k)$ i $f(s_i = k | \boldsymbol{\vartheta}) = \eta_k$. Stavimo još da je $N_k(\mathbf{s}) = \#\{s_i = k\}$ broj opažanja pridruženih komponenti k . Uvrstimo spomenute vjerojatnosti i oznake u izvod potpune vjerodostojnosti pa dobijemo

²⁰eng. *complete data likelihood*

$$L^c(\mathbf{y}, \mathbf{s} \mid \boldsymbol{\vartheta}) = \prod_{i=1}^N \prod_{k=1}^K (f(\mathbf{y}_i \mid \boldsymbol{\theta}_k) \eta_k)^{\mathbb{1}_{\{s_i=k\}}} = \prod_{k=1}^K \left(\prod_{i: S_i=k} f(\mathbf{y}_i \mid \boldsymbol{\theta}_k) \right) \left(\prod_{k=1}^K \eta_k^{N_k(\mathbf{s})} \right). \quad (1.44)$$

Provedimo sad procjenu parametara funkcije (1.44) metodom maksimalne vjerodostojnosti i Bayesovom metodom.

Procjena parametara ML metodom na potpunom skupu podataka

Funkcija potpune vjerodostojnosti izvedena u (1.44) sastoji se od $K + 1$ uzastopnih produkata faktoriziranih po parametrima. Prvih K faktora ovise samo o parametru $\boldsymbol{\theta}_k$ komponente \mathcal{K}_k , a posljednji ovisi samo o vektoru težina $\boldsymbol{\eta}$. Takva faktorizacija olakšava određivanje parametara modela jer omogućava da procjeni svakog parametra pristupamo izravno i neovisno od ostalih parametara. Kako bismo dodatno olakšali račun procjene, logaritmirat ćemo funkciju (1.44) pa koristiti metodu maksimalne log-vjerodostojnosti potpunih podataka. Dobivenu funkciju zatim maksimiziramo po vektoru parametara $\boldsymbol{\vartheta}$. Uvedimo oznaku indikatora pridruživanja $S_{ik} = \mathbb{1}_{\{s_i=k\}}$ i izvedimo potpunu log-vjerodostojnost prema navedenoj tehnici.

$$\begin{aligned} \log L^c(\mathbf{y}, \mathbf{S} \mid \boldsymbol{\vartheta}) &= \sum_{i=1}^N \sum_{k=1}^K S_{ik} \log(f(\mathbf{y}_i \mid \boldsymbol{\theta}_k) \eta_k^{N_k(\mathbf{s})}) = \sum_{i=1}^N \sum_{k=1}^K S_{ik} (\log(f(\mathbf{y}_i \mid \boldsymbol{\theta}_k)) + \log(\eta_k^{N_k(\mathbf{s})})) \\ &= \sum_{k=1}^K \sum_{i: s_i=k} \log(f(\mathbf{y}_i \mid \boldsymbol{\theta}_k)) + \sum_{k=1}^K N_k(\mathbf{s}) \log \eta_k. \end{aligned} \quad (1.45)$$

Izvest ćemo procjenu parametara $\boldsymbol{\theta}_k$ i η_k tako da izraz (1.45) deriviramo po traženom parametru. Dobivene su jednačbe

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}_k} \log L^c(\mathbf{y}, \mathbf{s} \mid \boldsymbol{\vartheta}) &= \frac{\partial}{\partial \boldsymbol{\theta}_k} \sum_{j=1}^K \sum_{i: s_i=j} \log(f(\mathbf{y}_i \mid \boldsymbol{\theta}_j)) = \frac{\partial}{\partial \boldsymbol{\theta}_k} \sum_{i: s_i=k} \log(f(\mathbf{y}_i \mid \boldsymbol{\theta}_k)) \\ &= \sum_{i: s_i=k} \frac{\frac{\partial}{\partial \boldsymbol{\theta}_k} f(\mathbf{y}_i \mid \boldsymbol{\theta}_k)}{f(\mathbf{y}_i \mid \boldsymbol{\theta}_k)} \end{aligned} \quad (1.46)$$

te

$$\frac{\partial}{\partial \eta_k} \log L^c(\mathbf{y}, \mathbf{s} \mid \boldsymbol{\vartheta}) = \frac{\partial}{\partial \eta_k} \sum_{j=1}^K N_j(\mathbf{s}) \log(\eta_j) = \frac{\partial}{\partial \eta_k} N_k(\mathbf{s}) \log(\eta_k) = \frac{N_k(\mathbf{s})}{\eta_k}. \quad (1.47)$$

Primijetimo da je $N = N_1(s) + \dots + N_K(s)$ i da procjena komponenta vektora težina $\boldsymbol{\eta}$ ovisi samo o broju opažanja pridruženih svakoj komponenti. Sada je procjenitelj komponente težine maksimalne log–vjerodostojnosti potpunog skupa podataka

$$\hat{\eta}_k = \frac{N_k(s)}{N}, \quad \forall k \in \mathcal{I}. \quad (1.48)$$

Svaki parametar $\boldsymbol{\theta}_k$ k -te komponente procjenjujemo iz opažanja pridruženih komponenti \mathcal{K}_k , za sve $k \in \mathcal{I}$. U slučaju kada komponente pripadaju nekoj Gaussovoj višedimenzionalnoj distribuciji, procjenitelji parametara maksimalne log–vjerodostojnosti potpunog skupa podataka su uzoračka aritmetička sredina i uzoračka kovarijacijska matrica

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{N_k(s)} \sum_{i: s_i=k} \mathbf{y}_i, \quad (1.49)$$

$$\hat{\boldsymbol{\Sigma}}_k = \frac{1}{N_k(s)} \sum_{i: s_i=k} (\mathbf{y}_i - \boldsymbol{\mu}_k)^\tau (\mathbf{y}_i - \boldsymbol{\mu}_k). \quad (1.50)$$

Ako je broj opažanja pridruženih nekoj komponenti jako mali, može se dogoditi da MLE procjenitelj ne postoji (ili je degeneriran) ili da leži u rubu parametarskog prostora Θ . U slučaju da je komponenta \mathcal{K}_k prazna, a budući da smo pretpostavili postojanje K podpopulacija, procjenitelj komponente težine η_k leži na rubu parametarskog prostora. Tada nemamo zadovoljene uvjete regularnosti pa nismo niti smjeli provoditi procjenu metodom maksimalne vjerodostojnosti. Kod modeliranja stvarnih podataka je zato važno što bolje ocijeniti mogući broj komponentata te, prije provođenja procjene ML metodom, provjeriti je li neka komponenta zapravo prazna. Najbolje procjene dobiju se ako sve komponente imaju *velik* broj pridruženih opažanja.

Bayesova procjena parametara na potpunom skupu podataka

Kako bismo odredili aposteriornu gustoću $f(\boldsymbol{\vartheta} \mid \mathbf{y}, s)$, primijenit ćemo Bayesov teorem 1.1.18 na apriornu gustoću $f(\boldsymbol{\vartheta})$ i funkciju potpune vjerodostojnosti $L^c(\mathbf{y}, s \mid \boldsymbol{\vartheta})$. Relacija proporcionalnosti koju koristimo je

$$f(\boldsymbol{\vartheta} \mid \mathbf{y}, s) \propto L^c(\mathbf{y}, s \mid \boldsymbol{\vartheta})f(\boldsymbol{\vartheta}). \quad (1.51)$$

Faktorizacijom potpune vjerodostojnosti (1.44), dobili smo $K + 1$ umnožaka. Na sličan način raspišemo apriornu gustoću pa dobijemo

$$f(\boldsymbol{\vartheta}) = f(\boldsymbol{\eta}) \prod_{k=1}^K f(\boldsymbol{\theta}_k). \quad (1.52)$$

Faktoriziranu apriornu gustoću uvrstit ćemo u relaciju Bayesove proporcionalnosti (1.51). Znamo da faktorizacijom po parametrima aposteriornu gustoću potpunih podataka možemo izraziti s

$$f(\boldsymbol{\vartheta} | \mathbf{s}, \mathbf{y}) = \prod_{k=1}^K f(\boldsymbol{\eta} | \mathbf{s}) f(\boldsymbol{\theta}_k | \mathbf{y}, \mathbf{s}). \quad (1.53)$$

Faktore iz jednadžbe (1.53) možemo neovisno analizirati po parametrima $\boldsymbol{\theta}$, odnosno $\boldsymbol{\eta}$ iz sljedećih relacija:

$$f(\boldsymbol{\eta} | \mathbf{s}) \propto \prod_{k=1}^K f(\boldsymbol{\eta}) \eta_k^{N_k(\mathbf{s})}, \quad (1.54)$$

odnosno

$$f(\boldsymbol{\theta}_k | \mathbf{y}, \mathbf{s}) \propto \prod_{i: S_i=k} f(\mathbf{y}_i | \boldsymbol{\theta}_k) f(\boldsymbol{\theta}_k). \quad (1.55)$$

Ovakva faktorizacija pogodna je za procjenu parametara mješavine jer, kao i u slučaju procjene korištenjem ML metode, procjeni svakog parametra pristupamo izravno. Dobivene izraze koristimo u metodi Bayesove procjene parametara opisane u potpoglavlju 1.1.

Procjena parametara u slučaju nepoznatog pridruživanja opažanja

U posljednjem dijelu potpoglavlja 1.3 razmatramo slučaj procjene parametara mješavine kada je poznat broj komponenata K , ali kada apriorno nemamo informaciju o realizaciji pridruživanja. Budući da parametre procjenjujemo na nepotpunom skupu podataka, a nije zanemariv utjecaj koji pridruživanje opažanja komponentama ima na zaključivanje, potrebne su složenije tehnike procjene. Ovaj problem je stoga daleko kompliciraniji od prethodnih. Razlog tomu je što, čak u slučaju da komponente mješavine pripadaju distribucijama koje se općenito lako modeliraju, ne postoje izravne statističke metode, već su potrebne metode koje koriste računalne algoritme za procjenu parametara. Jednom kada su parametri mješavine procijenjeni, koristimo zaključivanje s početka ovog potpoglavlja kako bismo pridružili opažanja odgovarajućim komponentama.

Neka je $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$ vektor opažanja slučajne varijable \mathbf{Y} iz konačne distribucije (1.33) te neka je $\boldsymbol{\vartheta} = (\boldsymbol{\theta}, \boldsymbol{\eta})$ vektor parametara koji procjenjujemo. Metoda momenata je početkom 20. stoljeća bila najčešće korištena metoda procjene parametara. Čuveni engleski matematičar Karl Pearson koristi ju još 1894. u svome radu *Contributions to the Mathematical Theory of Evolution* [11]. Pearson je modelirao binormalnu populaciju rakova kojoj je procijenio parametre distribucije. Jedan je od prvih autora znanstvenog članka na temu miješanih gustoća. Navodimo princip korištenja metode momenata u procjeni parametara.

Metoda momenata na nepotpunom skupu podataka

Neka je $H_j(Y)$ neka funkcija²¹ slučajne varijable Y koju koristimo u procjeni parametra $\boldsymbol{\vartheta}_j$, za $j \in \{1, \dots, d\}$, gdje je $d = \dim \boldsymbol{\vartheta}$. Uvjetno očekivanje računamo formulom

$$\mathbb{E} [H_j(\mathbf{Y} | \boldsymbol{\vartheta})] = \int_{\mathcal{Y}} H_j(y) f(y | \boldsymbol{\vartheta}) dy. \quad (1.56)$$

Želimo odrediti parametre $\boldsymbol{\vartheta}$ tako da teoretski momenti $\mathbb{E} [H_j(\mathbf{Y} | \boldsymbol{\vartheta})]$ odgovaraju uzoračkim $\mathbb{E} [H_j^*(\mathbf{Y} | \boldsymbol{\vartheta})]$. Tražena jednakost je funkcija nepoznatog parametra

$$H_j^* = \sum_{k=1}^K \mathbb{E} [H_j(\mathbf{Y} | \boldsymbol{\theta}_k) \eta_k]. \quad (1.57)$$

Vektor parametara $\boldsymbol{\vartheta}$ onda odredimo kao rješenje sustava jednažbi pa je za postojanje jedinstvenog rješenja nužno da imamo d linearno nezavisnih jednažbi.

Procjena parametara ML metodom na nepotpunom skupu podataka

Razvojem računala i numeričkih metoda u drugoj polovici 20. stoljeća, iterativne metode potiskuju teorijske. Tako je fokus s metode momenata prebačen na već spomenutu metodu maksimalne vjerodostojnosti. Izvest ćemo *vjerodostojnost mješavine* $L(\mathbf{y} | \boldsymbol{\vartheta})$ kao zajedničku distribuciju opažanja $\mathbf{y}_1, \dots, \mathbf{y}_N$ uvjetno na vektor parametara $\boldsymbol{\vartheta}$. Uvrstimo u definiciju vjerodostojnosti (1.10) raspis po komponentama mješavine:

$$L(\mathbf{y} | \boldsymbol{\vartheta}) = \prod_{i=1}^N f(\mathbf{y}_i | \boldsymbol{\vartheta}) = \prod_{i=1}^N \left(\sum_{k=1}^K f(\mathbf{y}_i | \boldsymbol{\theta}_k) \eta_k \right). \quad (1.58)$$

Primijetimo da se funkcija vjerodostojnosti mješavine sastoji od N faktora, od kojih je svaki sadrži sumu K pribrojnika. Usporedbe radi, funkciju potpune vjerodostojnosti mogli smo raščlaniti na $K + 1$ faktor te smo procjeni svakog parametra komponente pristupali izravno, a to u ovom slučaju nećemo moći. Formulu (1.58) možemo rastaviti u zbroj od K^N pribrojnika, ali ju ne možemo značajnije pojednostaviti. Slijede standardni postupci u procjeni parametara ML metodom: logaritmiranje i deriviranje izraza (1.58). Zbog velikog broja operacija koje smo do sad napravili, vidimo da je dobar smjer u daljnjoj analizi mješavine korištenje računalnih algoritama. Algoritam koji pritom koristimo je već ranije spomenuti EM algoritam. Ovdje započetu procjenu parametara mješavine u slučaju nepoznatog vektora pridruživanja dovršit ćemo u poglavlju 2 nakon što objasnimo EM algoritam.

²¹Primjerice $H_j(Y) = Y$ ili u slučaju računa m -tog momenta: $H_j(Y) = (Y - \mu)^m$.

Poglavlje 2

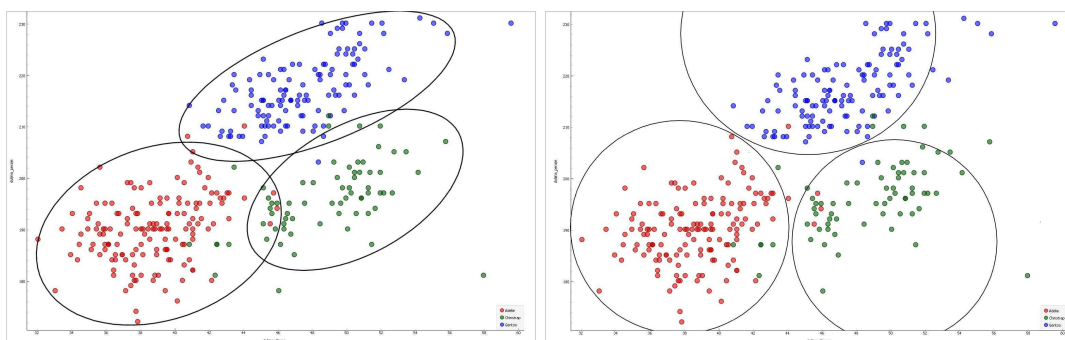
Algoritam maksimizacije očekivanja

Problem modeliranja konačnih mješavina dugo se niz godina rješavao metodama klasičnog statističkog zaključivanja koje smo objasnili u poglavlju 1. Preokret u pristupu procjeni parametara dogodio se uvođenjem računala u proces analize podataka. Razvojem numeričkih metoda u računarstvu, nametnula se ideja da metodu maksimalne vjerodostojnosti prilagodimo implementaciji na računalu. Formaliziranjem postupka procjene parametara u konačnici nastaje algoritam maksimizacije očekivanja. EM Algoritam su 1977. godine predstavili Dempster, Laird i Rubin¹ u radu *Maximum Likelihood from Incomplete Data via the EM Algorithm* [1], gdje su predložili postupak koji se sastoji od dva naizmjenično ponavljajuća koraka. Prema operacijama koje se izvode u svakom koraku, Dempster et al. su i iskovali naziv algoritma. Prvi je korak očekivanja (**E-korak**) u kojemu se, za trenutni procjenitelj parametara $\hat{\theta}$, računa očekivanje potpune log-vjerodostojnosti izvedene u izrazu (1.45). Korak maksimizacije (**M-korak**) određuje novi procjenitelj parametara koji očekivanje iz E-koraka maksimizira. U svakom se koraku algoritma približavamo lokalnom maksimumu funkcije log-vjerodostojnosti jer algoritam nužno konvergira u smislu da niz procjenitelja parametara dobivenih u iteracijama algoritma konvergira prema stvarnoj vrijednosti parametara mješavine. Uz blage zahtjeve regularnosti na funkciju L^c , osiguravamo nužnu konvergenciju niza procjenitelja u njezin lokalni maksimum. Međutim, postizanje konvergencije može biti vrlo sporo, osobito ako početna vrijednost vektora parametara nije izabrana pažljivo. Algoritme grupiranja dijelimo prema načinu pridruživanja podataka komponentama na klase tvrdog i mekog grupiranja.² Algoritmi tvrdog grupiranja ne dozvoljavaju preklapanje grupa i dodjeljivanje oznaka više grupa jednoj jedinki. S druge strane, algoritmi mekog grupiranja ne određuju striktno granice između komponentata te računaju vjerojatnost da jedinka pripada svakoj od komponentata te pridružuju na temelju najvjerojatnijeg događaja. EM algoritam je stoga tipičan primjer algoritma mekog

¹Dempster et al.

²eng. *hard and soft clustering*

grupiranja. Međutim, za određivanje inicijalnih vrijednosti parametara, dobar je izbor neki od algoritama tvrdog grupiranja, osobito neka inačica algoritma k -srednjih vrijednosti³. Komponente dobivene algoritmom k -srednjih vrijednosti su, u slučaju da mješavina pripada eksponencijalnoj familiji distribucija, disjunktne kugle u prostoru podataka. Grupe određene EM algoritmom često se preklapaju, te čine oblik elipsoida kao što je prikazano na slici 2.1.



Slika 2.1: meko i tvrdo grupiranje

Bojom su označene tri komponente kojima pripadaju podaci. Lijevi dijagram rasipanja prikazuje raspored grupa dobivenih mekim grupiranjem, dok desni prikazuje razmještaj dobiven tvrdim grupiranjem. Prikazana je prostorna raspršenost populacije koju modeliramo u potpoglavlju 2.2.

Osim procijenjenih parametara konačne mješavine, ovaj algoritam daje informaciju kojoj je komponenti pridruženo koje opažanje. Na konačne mješavine zato često indirektno nailazimo u brojnim znanstvenim disciplinama. Tako u biomedicini rješavamo problem klasifikacije genetskih promjena u organizmu ili određujemo značajke neke populacije živih bića. U poslovnoj ekonomiji postoji potreba za segmentacijom tržišta i kupaca te se prirodno nameće model miješanih gustoća. S gledišta strojnog učenja, EM algoritam pripada generativnim algoritmima te se proučava iz konteksta optimizacije učenja sistema. Osim navedenih područja, mnoge druge prirodne i društvene znanosti bave se modeliranjem konačnih mješavina u nekom njima svojstvenom kontekstu. Puno više o praktičnoj primjeni algoritma možemo pročitati u [4]. Frühwirth-Schnatter daje u 7. poglavlju navedene knjige pregled različitih problema analize podataka koji koriste modele konačnih mješavina te nas upućuje na znanstvene radove u kojima su pojedine teme detaljnije obrađene. U ovom nam poglavlju predstoji dovršiti statističko zaključivanje započeto u potpoglavlju 1.3. Potpoglavlje 2.1 bavi se formaliziranjem računarskog pristupa modeliranju konačnih mješavina, dok je potpoglavlje 2.2 rezervirano za analizu i modeliranje stvarnih podataka.

³eng. *k-means algorithm*

2.1 Algoritam i procjena parametara

Promatramo populaciju neke Gaussove konačne mješavine koja se sastoji od $K \in \mathbb{N}$ različitih normalno distribuiranih komponenata. Kao u definiciji (1.26) mješavina ima $r \in \mathbb{N}$ značajki koje promatramo unutar slučajnog vektora $\mathbf{Y} = (Y_1, \dots, Y_r)$. Iz mješavine je izabran proizvoljan uzorak $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)$ koji se sastoji od $N \in \mathbb{N}$ nezavisnih opažanja. Na temelju tih opažanja želimo izvesti zaključke o čitavoj populaciji. Problem koji rješavamo je određivanje parametara distribucije za svaku od K komponenata u slučaju kada vektor pridruživanja promatranog opažanja nije unaprijed poznat. Procjenu parametara provodimo ML metodom, a pritom rješavamo i problem pridruživanja opažanja komponentama. Prisjetimo se funkcije gustoće promatrane Gaussove mješavine:

$$f_{\mathbf{Y}}(\mathbf{y}) = \sum_{k=1}^K \eta_k \left((2\pi)^{-\frac{r}{2}} (\det \Sigma_k)^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu}_k) \Sigma_k^{-1} (\mathbf{y} - \boldsymbol{\mu}_k)^{\tau} \right) \right), \quad \forall \mathbf{y} \in \mathbf{Y}, \quad (2.1)$$

ML metoda procjene parametara određuje procjenitelj parametara koji maksimizira vjerodostojnost. Funkciju nepotpune vjerodostojnosti (1.58) logaritmiramo pa je nepotpuna log-vjerodostojnost

$$\log L(\mathbf{y} | \boldsymbol{\vartheta}) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \eta_k f(\mathbf{y}_i | \boldsymbol{\theta}_k) \right). \quad (2.2)$$

Funkcija gustoće komponente \mathcal{K}_k uvjetno na odgovarajući parametar $\boldsymbol{\theta}_k$ je Gaussova komponentna gustoća koju vidimo u gustoći mješavine (2.1).

EM algoritam

Algoritam maksimizacije očekivanja je ponavljajući algoritam koji parametre mješavine procjenjuje maksimizirajući očekivanje potpune log-vjerodostojnosti. Budući da od ranije nemamo informacije o latentnim varijablama, potpunu log-vjerodostojnost ćemo maksimizirati u terminima njezinog uvjetnog očekivanja uz dani poznati vektor opažanja \mathbf{y} i trenutnu vrijednost parametara $\boldsymbol{\vartheta}^{(w)}$. Indeks $w \in \mathbb{N}_0$ označava da provodimo $(w+1)$ -u iteraciju algoritma. U potpoglavlju 2.8. knjige *Finite Mixture Models* [8] detaljno je objašnjena primjena algoritma na modele konačnih mješavina te je opravdano korištenje uvjetnog očekivanja potpune log-vjerodostojnosti. To ćemo zaključivanje izvesti pri definiranju E-koraka. Označimo s $\boldsymbol{\vartheta}^*$ stvarnu vrijednost vektora parametara koji maksimizira potpunu log-vjerodostojnost. Algoritam započinje inicijalizacijom vektora parametara $\boldsymbol{\vartheta}^{(0)}$. Izbor prvih vrijednosti parametara može biti slučajaj, ali je bolje napraviti analizu funkcije $\log L^c(\mathbf{y}, \mathbf{s} | \boldsymbol{\vartheta})$ i odabrati vrijednost koja je relativno blizu stvarne vrijednosti $\boldsymbol{\vartheta}^*$. Također,

u slučaju da log–vjerodostojnost ima višestruke nultočke, algoritam bi trebalo provesti na širokom izboru početnih vrijednosti.

E–korak

Uvedimo oznaku za uvjetno očekivanje potpune log–vjerodostojnosti uz dani poznati vektor opažanja \mathbf{y} i trenutnu vrijednost parametara $\boldsymbol{\vartheta}^{(w)}$ u $(w + 1)$ –oj iteraciji:

$$Q(\boldsymbol{\vartheta} \mid \boldsymbol{\vartheta}^{(w)}) = \mathbb{E}_{\boldsymbol{g}^{(w)}} [\log L^c(\boldsymbol{\vartheta}, \mathbf{s}) \mid \mathbf{y}]. \quad (2.3)$$

U jednadžbi (2.3) oznaka $\boldsymbol{\vartheta}^{(w)}$ označava da u $(w + 1)$ –om koraku na očekivanje uvjetuje i trenutna vrijednost procijenjenog vektora parametara. Kao što je najavljeno, računamo uvjetno očekivanje potpune log–vjerodostojnosti iako smo u slučaju kada je pridruživanje \mathbf{S} apriorno nepoznato. Primijetimo da je funkcija potpune log–vjerodostojnosti $\log L^c$ linearna s obzirom na nepoznatu realizaciju pridruživanja. Računamo uvjetno očekivanje pridruživanja uz dano opažanje i trenutnu vrijednost parametara:

$$\mathbb{E}_{\boldsymbol{g}^{(w)}} [S_i = s_k \mid \mathbf{y}] = \mathbb{P}_{\boldsymbol{g}^{(w)}}(s_{ik} = 1 \mid \mathbf{y}), \quad \forall i \in \{1, \dots, N\}, \forall k \in \mathcal{I}, \quad (2.4)$$

gdje je \mathcal{I} skup indikatora komponenata. Primijetimo da smo dobili aposteriornu vjerojatnost pridruživanja i –tog opažanja komponenti \mathcal{K}_k . Navedena vjerojatnost je odgovornost h_{ik} uvedena izrazom (1.36). Uvrstimo u izraz uvjetnog očekivanja potpune log–vjerodostojnosti (2.3) odgovornost iz (2.4) pa dobijemo

$$Q(\boldsymbol{\vartheta} \mid \boldsymbol{\vartheta}^{(w)}) = \sum_{i=1}^N \sum_{k=1}^K h_{ik} (\log \eta_k + \log f_k(\mathbf{y}_i \mid \boldsymbol{\theta}_k)). \quad (2.5)$$

Dobivenu funkciju maksimiziramo u sljedećem koraku. Budući da su parametri težina i komponenata odvojeni operacijom zbrajanja, njihovoj procjeni pristupamo neovisno jer deriviranjem jednadžbe (2.5) po nekom određenom parametru, ostaje samo izraz s tim parametrom.

M–korak

Algoritam je još uvijek u $(w + 1)$ –oj iteraciji i potrebno je odrediti novu vrijednost vektora parametara za koju funkcija (2.5) postiže maksimum. Izvedimo prvo procjenitelj vektora težina. Na potpunom skupu podataka izveli smo procjenu parametara težina u jednadžbi (1.48) kao količnik ukupnog broja opažanja pridruženih promatranj komponenti $N_k(s)$ i ukupnog broja opažanja N . Sada međutim, radimo s nepotpunim skupom podataka i ne znamo koja je realizacija pridruživanja. Zbog konstrukcije E–koraka, ideja je zamijeniti broj opažanja $N_k(s)$ s uvjetnim očekivanjem h_{ik} uz trenutne vrijednosti parametara. Nova vrijednost procjenitelja parametara težine je onda

$$\eta_k^{(w+1)} = \frac{1}{N} \sum_{i=1}^n h_{ik}^{(w)}, \quad \forall k \in \mathcal{I}. \quad (2.6)$$

Utjecaj svakog od opažanja \mathbf{y}_i na procjenu težine komponente η_k nije zanemariv i on odgovara trenutnoj vrijednosti aposteriorne vjerojatnosti $h_{ik}^{(w)}$. Zbog normiranosti vektora težina (1.25), dovoljno je procijeniti $K - 1$ parametar težine u trenutnoj iteraciji, a jednu težinu odrediti pomoću ostalih. Izbor težine komponente koju ne procjenjujemo već računamo koristeći normiranost je proizvoljan. Ovdje je izbor pao na procjenitelj težine $\eta_K^{(w+1)}$:

$$\eta_K^{(w+1)} = 1 - \sum_{k=1}^{K-1} \eta_k^{(w+1)}. \quad (2.7)$$

Preostalo nam je još odrediti procjenitelje distribucijskih parametara komponenata u $(w + 1)$ -oj iteraciji koristeći ML metodu procjene na uvjetnom očekivanju (2.5). Za svaku ćemo komponentu odrediti parametre koji maksimiziraju uvjetno očekivanje tako da odredimo rješenje diferencijalne jednadžbe

$$\frac{\partial}{\partial \boldsymbol{\theta}_k} Q(\boldsymbol{\vartheta} | \boldsymbol{\vartheta}^{(w)}) = \sum_{i=1}^N \sum_{j=1}^K h_{ij} \frac{\partial (\log f_j(\mathbf{y}_i | \boldsymbol{\theta}_j))}{\partial \boldsymbol{\theta}_k} = 0. \quad (2.8)$$

Izvedimo log-vjerodostojnost k -te komponente Gaussove mješavine kako bismo mogli odrediti procjenitelje parametara komponenata mješavine

$$\begin{aligned} \log f_k(\mathbf{y} | (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) &= \log \left(\frac{\exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_k)\boldsymbol{\Sigma}_k^{-1}(\mathbf{y} - \boldsymbol{\mu}_k)^\tau\right\}}{(2\pi)^{\frac{r}{2}}(\det \boldsymbol{\Sigma}_k)^{\frac{1}{2}}}\right) = \\ &= -\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}_k)\boldsymbol{\Sigma}_k^{-1}(\mathbf{y} - \boldsymbol{\mu}_k)^\tau - \log\left((2\pi)^{\frac{r}{2}}(\det \boldsymbol{\Sigma}_k)^{\frac{1}{2}}\right). \end{aligned} \quad (2.9)$$

Procjenitelje očekivanja i kovarijacijske matrice komponente \mathcal{K}_k određujemo rješavajući sljedeće jednadžbe:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\mu}_k} Q(\boldsymbol{\vartheta} | \boldsymbol{\vartheta}^{(w)}) &= \sum_{i=1}^N h_{ik} \frac{\partial (\log f_k(\mathbf{y}_i | (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)))}{\partial \boldsymbol{\mu}_k} = \\ &= \sum_{i=1}^N h_{ik} \frac{\partial \left(-\frac{1}{2}(\mathbf{y}_i - \boldsymbol{\mu}_k)\boldsymbol{\Sigma}_k^{-1}(\mathbf{y}_i - \boldsymbol{\mu}_k)^\tau - \log\left((2\pi)^{\frac{r}{2}}(\det \boldsymbol{\Sigma}_k)^{\frac{1}{2}}\right)\right)}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N h_{ik} \frac{\partial \left(-\frac{1}{2}(\mathbf{y}_i - \boldsymbol{\mu}_k)\boldsymbol{\Sigma}_k^{-1}(\mathbf{y}_i - \boldsymbol{\mu}_k)^\tau\right)}{\partial \boldsymbol{\mu}_k} = 0, \end{aligned} \quad (2.10)$$

$$\begin{aligned}
 \frac{\partial}{\partial \Sigma_k} Q(\boldsymbol{\vartheta} \mid \boldsymbol{\vartheta}^{(w)}) &= \sum_{i=1}^N h_{ik} \frac{\partial (\log f_k(\mathbf{y}_i \mid (\boldsymbol{\mu}_k, \Sigma_k)))}{\partial \Sigma_k} = \\
 &= \sum_{i=1}^N h_{ik} \frac{\partial \left(-\frac{1}{2} (\mathbf{y}_i - \boldsymbol{\mu}_k) \Sigma_k^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_k)^\tau - \log \left((2\pi)^{\frac{r}{2}} (\det \Sigma_k)^{\frac{1}{2}} \right) \right)}{\partial \Sigma_k} = 0. \quad (2.11)
 \end{aligned}$$

Parametre koje anuliraju jednađbe (2.10) i (2.11) možemo zapisati u zatvorenoj formi:

$$\boldsymbol{\mu}_k^{(w+1)} = \frac{\sum_{i=1}^N h_{ik}^{(w)} \mathbf{y}_i}{\sum_{i=1}^N h_{ik}^{(w)}}, \quad (2.12)$$

$$\Sigma_k^{(w+1)} = \frac{\sum_{i=1}^N h_{ik}^{(w)} (\mathbf{y}_i - \boldsymbol{\mu}_i^{(w+1)}) (\mathbf{y}_i - \boldsymbol{\mu}_i^{(w+1)})^\tau}{\sum_{i=1}^N h_{ik}^{(w)}}. \quad (2.13)$$

Distribucijama koje pripadaju eksponencijalnoj familiji distribucija je zbog forme njihove funkcije gustoće jednostavno provesti procjenu parametara. Vratimo li se na izvod nepotpune vjerodostojnosti mješavine u jednađbi (1.58), vidimo da bi direktnu primjenu ML metode bilo teško egzaktno izvesti. S druge strane, procjenitelje dobivene EM–algoritmom dobili smo nizom jednostavnih operacija. Račun je osobito jednostavan ako ga provodimo na računalu, ali je i tada korisno izbaciti redundantne operacije koje provode iste izračune kroz sve iteracije. Broj računskih operacija ćemo optimizirati uvođenjem dovoljnih statistika za vektor parametara $\boldsymbol{\vartheta}$. Jedan od prijedloga smanjenja broja operacija opisan je u potpoglavlju 3.2 knjige [8] gdje su uvedene statistike

$$\begin{aligned}
 T_{k1}^{(w)} &= \sum_{i=1}^N h_{ik}^{(w)}, \\
 T_{k2}^{(w)} &= \sum_{i=1}^N h_{ik}^{(w)} \mathbf{y}_i \quad \text{i} \\
 T_{k3}^{(w)} &= \sum_{i=1}^N h_{ik}^{(w)} \mathbf{y}_i \mathbf{y}_i^\tau,
 \end{aligned} \quad (2.14)$$

za svaku komponentu \mathcal{K}_k , $k \in \mathcal{I}$. Sada kovarijacijsku matricu k –te komponente u $(w+1)$ –oj iteraciji možemo odrediti po formuli

$$\Sigma_k^{(w+1)} = \frac{T_{k3}^{(w)} - \left(T_{k1}^{(w)} \right)^{-1} T_{k2}^{(w)} \left(T_{k2}^{(w)} \right)^\tau}{T_{k1}^{(w)}}. \quad (2.15)$$

Korištenjem formule (2.15) umjesto formule (2.13) za procjenu kovarijacijske matrice komponente \mathcal{K}_k , za 50% smanjujemo CPU vrijeme potrebno za izvođenje iteracije.⁴ Koraci E i M algoritma naizmjenično se ponavljaju do postizanja konvergencije. Potrebno je stoga zadati uvjet zaustavljanja prije no što pokrenemo algoritam. Tipičan uvjet zaustavljanja je kada je razlika vektora dobivenih uzastopnim iteracijama ili razlika njihovih log–vjerodostojnosti dovoljno mala. Dempster et al. su prvi diskutirali konvergenciju EM algoritma, ali se u formalnom dokazu predstavljenom u njihovom članku [1] potkrala pogreška. Na pogrešku je ukazao Jeff C. F. Wu u svom članku *On the Convergence Properties of EM Algorithm* [17] gdje je i dao dokaz konvergencije algoritma u seriji vezanih teorema.

Konvergencija EM algoritma

Iskazat ćemo glavni teorem konvergencije EM algoritma čiji dokaz možemo pronaći u [17]. Prije toga moramo još eksplicitno postaviti uvjete u kojima algoritam nužno konvergira. Na problem možemo naići u M–koraku ako rješenje diferencijalne jednadžbe (2.5) ne postoji u zatvorenoj formi. U tom su slučaju Dempster et al. predložili uvođenje preslikavanja $M: \Theta \rightarrow \Theta$ za koje vrijedi

$$\boldsymbol{\vartheta}^{(w+1)} = M(\boldsymbol{\vartheta}^{(w)}), \quad \forall w \in \mathbb{N}_0. \quad (2.16)$$

Ovako definirano preslikavanje koristimo u M–koraku za odabir novog procjenitelja vektora parametara, a ono spada u klasu preslikavanja *točke–u–skup*.⁵ Algoritam tada nazivamo **generalizirani EM algoritam**, ili kraće **GEM algoritam**. Rezultat uvođenja preslikavanja (2.16) je da je niz $(L(\boldsymbol{\vartheta}^{(w+1)}))_{w \in \mathbb{N}_0}$ neopadajući. U iskazu teorema koristit ćemo oznaku \mathcal{M} za skup lokalnih maksimuma funkcije vjerodostojnosti, a oznaku \mathcal{S} za skup stacionarnih točaka funkcije vjerodostojnosti na interioru prostora parametara.

Teorem 2.1.1. *Neka je $(\boldsymbol{\vartheta}^{(w+1)})_{w \in \mathbb{N}_0}$ niz parametara određenih pomoću preslikavanja (2.16) u M–koraku GEM algoritma. Pretpostavimo da vrijedi:*

- i. *M je zatvoreno⁶ preslikavanje točke–u–skup na komplementu skupa \mathcal{S} .*
- ii. *$L(\boldsymbol{\vartheta}^{(w+1)}) > L(\boldsymbol{\vartheta}^{(w)})$, za sve $\boldsymbol{\vartheta}^{(w)} \notin \mathcal{S}$.*

Tada su sva gomilišta niza $(\boldsymbol{\vartheta}^{(w)})$ limesi i vrijedi da su stacionarne točke funkcije L . Također, niz $(L(\boldsymbol{\vartheta}^{(w)}))_{w \in \mathbb{N}_0}$ monotono konvergira prema nekoj vrijednosti $\tilde{L} = L(\tilde{\boldsymbol{\vartheta}})$, gdje je $\tilde{\boldsymbol{\vartheta}} \in \mathcal{S}$.

⁴CPU vrijeme je vrijeme potrebno procesoru računala da izvrši zadane naredbe.

⁵eng. *point-to-set map* — u matematičkom programiranju ovakva preslikavanja definiraju algoritme, a iteracije algoritma generiraju niz $\boldsymbol{\vartheta}^{(w)}$ t.d. $\boldsymbol{\vartheta}^{(w+1)} \in M(\boldsymbol{\vartheta}^{(w)})$

⁶eng. *closed map* — zatvorena preslikavanja su ona koja preslikavaju zatvoreni skup u zatvoreni skup

Korolar 2.1.2. *Tvrđnja teorema 2.1.1 vrijedi i u slučaju kada tvrdnje (i.) i (ii.) vrijede na komplementu skupa \mathcal{M} . Tada limes niza vrijednosti funkcija log-vjerodostojnosti konvergira u lokalni maksimum, tj.*

$$\lim_{w \rightarrow \infty} L(\boldsymbol{\vartheta}^{(w)}) = L(\tilde{\boldsymbol{\vartheta}}), \quad \text{gdje je } \tilde{\boldsymbol{\vartheta}} \in \mathcal{M}. \quad (2.17)$$

Sad kad znamo da algoritam konvergira, zanima nas kojom brzinom. Brzinu konvergencije algoritma možemo opisati matricno, a pritom ćemo opet koristiti preslikavanje (2.16) uvedeno ranije. Ako niz $\boldsymbol{\vartheta}^{(w)}$ konvergira prema nekoj fiksnoj točki $\boldsymbol{\vartheta}^\circ$ i ako je preslikavanje M neprekidno, onda je $\boldsymbol{\vartheta}^\circ$ fiksna točka algoritma, tj. vrijedi

$$\boldsymbol{\vartheta}^\circ = M(\boldsymbol{\vartheta}^\circ). \quad (2.18)$$

Primijenimo li Taylorov razvoj preslikavanja M oko točke $\boldsymbol{\vartheta}^{(w)} = \boldsymbol{\vartheta}^\circ$, tada postoji okolina točke $\boldsymbol{\vartheta}^\circ$ na kojoj vrijedi

$$\boldsymbol{\vartheta}^{(w+1)} - \boldsymbol{\vartheta}^\circ \approx J(\boldsymbol{\vartheta}^\circ)(\boldsymbol{\vartheta}^{(w)} - \boldsymbol{\vartheta}^\circ), \quad (2.19)$$

gdje je $J(\boldsymbol{\vartheta})$ d -dimenzionalna Jacobijeva matrica preslikavanja $M(\boldsymbol{\vartheta}) = (M_1(\boldsymbol{\vartheta}), \dots, M_d(\boldsymbol{\vartheta}))$. Stoga u okolini fiksne točke $\boldsymbol{\vartheta}^\circ$ EM algoritam konvergira linearno pa Jacobijevu matricu $J(\boldsymbol{\vartheta}^\circ)$ nazivamo *matricom brzine konvergencije*.⁷ Na prostoru parametara možemo definirati mjeru opažene konvergencije kao globalnu brzinu konvergencije sljedećim izrazom

$$r = \lim_{w \rightarrow \infty} \frac{\|\boldsymbol{\vartheta}^{(w+1)} - \boldsymbol{\vartheta}^\circ\|}{\|\boldsymbol{\vartheta}^{(w)} - \boldsymbol{\vartheta}^\circ\|}, \quad (2.20)$$

gdje je $\|\cdot\|$ proizvoljna norma na \mathbb{R}^d . Lako se pokaže da, uz određene uvjete regularnosti vrijedi da mjeru konvergencije možemo poistovjetiti s najvećom svojstvenom vrijednosti matrice $J(\boldsymbol{\vartheta}^\circ)$. U praksi ipak brzinu konvergencije računamo formulom

$$r = \lim_{w \rightarrow \infty} \frac{\|\boldsymbol{\vartheta}^{(w+1)} - \boldsymbol{\vartheta}^{(w)}\|}{\|\boldsymbol{\vartheta}^{(w)} - \boldsymbol{\vartheta}^{(w-1)}\|}. \quad (2.21)$$

U potpoglavlju 2.2 prezentirat ćemo pristup modeliranju Gaussovih mješavina stvarne populacije. Podatke ćemo analizirati pomoću metoda objašnjenih u poglavlju 1 i algoritma

⁷Ovaj pristup računanja brzine konvergencije objašnjen je u 2. poglavlju knjige Finite Mixture Models.

maksimizacije očekivanja pa ovdje navodimo njegov pseudo kod.

	unos: Učitaj početnu vrijednost vektora parametara $\boldsymbol{\vartheta} \leftarrow (\eta_1, \dots, \eta_K, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K)$
	ispis: Spremi vektor procijenjenih parametara distribucije $\hat{\boldsymbol{\vartheta}} \leftarrow \boldsymbol{\vartheta}$
1	ponavljaj
2	E–korak: Izračunaj trenutnu odgovornost $h_{ik}, \forall i = 1, \dots, N, \forall k \in \mathcal{I}$.
3	$h_{ik} \leftarrow \frac{\mathbb{P}(\mathbf{Y} = \mathbf{y}_i \mid s_i = k, \boldsymbol{\vartheta})\eta_k}{\sum_{j=1}^K \mathbb{P}(\mathbf{Y} = \mathbf{y}_i \mid s_i = j, \boldsymbol{\vartheta})\eta_j}$
4	M–korak: Izračunaj nove vrijednosti vektora parametara u ovisnosti o trenutnoj vrijednosti $h_{ik}, \forall k \in \mathcal{I}$.
	$\eta_k \leftarrow \frac{1}{N} \sum_{i=1}^n h_{ik}$
	$\boldsymbol{\mu}_k \leftarrow \frac{\sum_{i=1}^N h_{ik} \mathbf{y}_i}{\sum_{i=1}^N h_{ik}}$
	$\boldsymbol{\Sigma}_k \leftarrow \frac{\sum_{i=1}^N h_{ik} (\mathbf{y}_i - \boldsymbol{\mu}_i) (\mathbf{y}_i - \boldsymbol{\mu}_i)^\tau}{\sum_{i=1}^N h_{ik}}$
5	Izračunaj novu vrijednost funkcije log–vjerodostojnosti.
	$\log L(\mathbf{y} \mid \boldsymbol{\vartheta}) \leftarrow \sum_{i=1}^N \log \left(\sum_{k=1}^K \eta_k f(\mathbf{y}_i \mid \boldsymbol{\theta}_k) \right)$
6	ako je zadovoljen uvjet konvergencije onda
7	Izadi iz petlje.
8	inače
9	Prijeđi na sljedeću iteraciju.
10	$w \leftarrow w + 1$
11	kraj
12	dok ne bude <u>zadovoljen uvjet konvergencije</u>

Algoritam 1: EM algoritam

2.2 Praktična primjena algoritma

U ovom potpoglavlju ćemo primijeniti metode procjene parametara na stvarnim podacima te obraditi rezultate.⁸ Modeliramo populaciju pingvina otočja Palmer na Antarktici prema podacima prikupljenima u sklopu istraživanja značajki triju vrsta pingvina Adélie, Chinstrap i Gentoo. Otočje se sastoji od otoka Torgersen, Biscoe i Dream, a podaci su prikupljeni od 2007. do 2009. godine. Istraživači, na čelu s dr. Kristen Gorman, su bilježili mjere veličine, podatke o leglu, udio relevantnih izotopa u krvi, otok na kojemu je jedinka pronađena te pripadnost vrsti. Takvi podaci su primjer potpunog skupa podataka, ali ćemo ih modelirati pod pretpostavkom da unaprijed ne znamo kojoj vrsti pingvina određena jedinka pripada. Cilj nam je na temelju mjera veličine promatrane jedinice odrediti vrstu pingvina. Budući da se radi o životinjskoj populaciji, smisleno je da je distribucija svake komponente višedimenzionalna normalna distribucija. Za svaku ćemo vrstu pingvina odrediti parametre distribucije i pripadne proporcije miješanja. Nakon završetka modeliranja, usporedit ćemo pridruživanja dobivena EM algoritmom sa stvarnim podacima i izračunati pogrešku. Ovo je potpoglavlje podijeljeno u cjeline prema fazi modeliranja.

Pripremna obrada podataka

Koristimo tablicu `penguins.csv` koja sadrži kategoričke varijable (*vrsta*, *otok*, *spol*, *godina*) i numeričke varijable (*duljina kljuna*, *visina kljuna*, *duljina peraje*, *masa tijela*). Strukturu podataka možemo vidjeti u tablici 2.1.

r.br.	vrsta	otok	duljina kljuna	visina kljuna	duljina peraje	masa tijela	spol	godina
20	Adelie	Torgersen	46	21.5	194	4200	m	2007
21	Adelie	Biscoe	37.8	18.3	174	3400	f	2007
66	Adelie	Biscoe	41.6	18	192	3950	m	2008
74	Adelie	Torgersen	45.8	18.9	197	4150	m	2008
155	Gentoo	Biscoe	48.7	14.1	210	4450	f	2007
189	Gentoo	Biscoe	42.6	13.7	213	4950	f	2008
277	Chinstrap	Dream	50	19.5	196	3900	m	2007
335	Chinstrap	Dream	45.6	19.4	194	3525	f	2009

Tablica 2.1: Tablični prikaz podataka o pingvinima

Prikazan je dio originalnih podataka pri čemu je zaglavlje tablice prevedeno na hrvatski jezik. Također, muške jedinice označavamo slovom **m**, a ženske slovom **f**.

⁸Podaci su preuzeti s GitHub stranice <https://github.com/allisonhorst/palmerpenguins> dr. Allison Horst 17. kolovoza 2020. te su javno dostupni prema *Creative Commons Zero v1.0 Universal* licenci.

Jedinke prikazane u tablici 2.1 su izabirane tako da reprezentiraju vrijednosti kategoričkih varijabli. Izmjerene duljine su navedene u milimetrima, a masa u gramima. Populaciju pingvina modeliramo prema numeričkim varijablama. Svaku jedinku stoga karakteriziramo 4–dimenzionalnim slučajnim vektorom $Y = (Y_1, Y_2, Y_3, Y_4)$ čije komponente redom odgovaraju navedenim numeričkim varijablama. Prema notaciji iz definicije (1.2.1) stavimo $r = 4$. U tablici su zabilježene značajke o 344 jedinke pingvina, ali u analizu ne uvrštavamo one jedinke koje nemaju zabilježene sve numeričke varijable. Modeliranje stoga provodimo na populaciji koja se sastoji od 342 jedinke pa je $N = 342$. Također, iz informacija o istraživanju, znamo da se populacija pingvina otočja Palmer sastoji od tri vrste, stoga je $K = 3$. Da bi model bio dovoljno dobar, potrebno je prilagoditi broj opažanja promatranom broju značajki kako bi se izbjegla prenaučenosť i podnaučenosť modela. Često korištena smjernica o broju potrebnih opažanja u ovisnosti od broja promatranih značajki zahtijeva da su uređeni relacijom $2^r < N$. Ovaj zahtjev je u našem slučaju zadovoljen pa smatramo da imamo dovoljan broj podataka da model bude dobar.

Procjena parametara EM algoritmom

Računalni kôd je pisan u jeziku *Python* [15] i možemo ga pronaći u dodatku A. Svi grafički prikazi dani u ovom radu napravljeni su korištenjem Pythonovog paketa *matplotlib* i softvera *Orange* [2]. Prvi korak u modeliranju mješavine je inicijalizacija parametara, a ona je provedena algoritmom *k*–srednjih vrijednosti ++.⁹ čiji je kôd također napisan u dodatku A.

Algoritam *k*–srednjih vrijednosti ++

Algoritam *K++* je implementiran u skripti *k_means_pp.py*, gdje učitavamo originalnu tablicu *penguins.csv*. Nakon izvođenja algoritma, spremamo novu tablicu *potpuni_podaci.csv* dobivenu spajanjem originalne tablice i vektora pridruživanja. Ovaj se algoritam razlikuje od običnog algoritma *k*–srednjih vrijednosti samo u načinu izbora prvog *centroida* svake komponente. Centroid je središte mase komponente, a računamo ga kao srednju vrijednost svih podataka pridruženih toj komponenti. Prvi centroid biramo nasumično, dok ostalih $K - 1$ biramo tako da iz skupa preostalih jedinki izaberemo one koje su, po nekoj metrici prostora \mathbb{R}^r , najudaljenije od već izabranih centroida. Grupiranje sad nastavljamo običnim algoritmom *k*–srednjih vrijednosti koji podatke grupira izmjenjujući *korak pridruživanja* i *korak određivanja novih centroida*. U koraku pridruživanja svakoj jedinki dodijelimo oznaku one komponente čiji je centroid njoj najbliži. Nakon što se sve jedinke pridruže nekoj komponenti, nove centroide računamo kao aritmetičku sredinu podataka po grupama, a algoritam se ponavlja do postizanja konvergencije. Po završetku provođenja algoritma, jedinke su pridružene komponentama i tada procjenjujemo distribucijske parametre pojedine

⁹U nastavku teksta koristimo skraćeni naziva algoritma — *algoritam K++*.

vrste metodom momenata. Time je završen postupak inicijalizacije parametara mješavine i inicijalnog pridruživanja jedinki komponentama. Ovaj način grupiranja podataka spada u klasu tvrdih grupiranja jer je svaku jedinku moguće pridružiti točno jednoj komponenti. Iako su podaci pridruženi komponentama, a distribucijski parametri vrsta procijenjeni, algoritam K++ nije optimalan izbor za grupiranje konačnih mješavina. Razlog tomu je što komponente ne možemo striktno odvojiti što je upravo slučaj s populacijom pingvina koju modeliramo. Kad u istom biomu živi više genetski sličnih podvrsta iste vrste životinje ili biljke, često dolazi do miješanja među populacijama. Može se dogoditi da neku jedinku krivo pridružimo komponenti kojoj ne pripada, samo zato što joj je centroid te komponente najbliži. Zbog toga je bolje koristiti pristup koji za svaku komponentu računa vjerojatnost da joj jedinka pripada. Upravo je na tom principu koncipiran EM algoritam s kojim i nastavljamo modeliranje.

EM algoritam

Nakon provedene inicijalizacije parametara i pridruživanja, možemo krenuti s algoritmom maksimizacije očekivanja koji je detaljno objašnjen u potpoglavlju 2.1. Pokrećemo skriptu `em_algoritam.py` koja učitava tablicu `potpuni_podaci.csv` dobivenu korištenjem algoritma `k-srednjih` vrijednosti `++`. Zadajemo uvjet zaustavljanja na način da se mora ispuniti barem jedna od dvije tvrdnje uvjeta.

- Razlika vektora parametara procijenjenih u uzastopnim iteracijama je dovoljno mala, tj. manja od $\delta = 10^{-7}$

$$|\boldsymbol{\vartheta}^{(w+1)} - \boldsymbol{\vartheta}^{(w)}| < 10^{-7}, \quad \text{za neki } w \in \mathbb{N}_0$$

- Razlika vrijednosti funkcija log-vjerodostojnosti u uzastopnim iteracijama je dovoljno mala, tj. manja od $\varepsilon = 10^{-5}$

$$|L(\boldsymbol{\vartheta}^{(w+1)}) - L(\boldsymbol{\vartheta}^{(w)})| < 10^{-5}, \quad \text{za neki } w \in \mathbb{N}_0$$

Ako je za neki $w \in \mathbb{N}_0$ ispunjena jedna od navedenih tvrdnji onda smatramo da je algoritam iskonvergirao. Dobiveni vektor pridruživanja upisujemo u novu tablicu `konacna_tablica.csv` u kojoj su navedene vrijednosti pojedinih karakteristika jedinki, kao i pridruživanje ranije određeno algoritmom K++.

Analiza rezultata

Uz ovako definiran uvjet zaustavljanja, konvergencija se postiže nakon nekoliko desetaka iteracija EM algoritma. Uspoređujući rezultate dobivene pokretanjem računalnog kôda u

više navrata, primijetili smo da se procijenjeni parametri ne razlikuju značajno, tj. da dolaze iz $\varepsilon = 10^{-4}$ okoline stvarnog parametra. Razlika nastaje zbog slučajnog odabira prvog centroida u inicijalizaciji parametara algoritmom K++. Prikazat ćemo rezultate dobivene jednim pokretanjem algoritma te ih interpretirati. Konačna tablicu za jedinke predstavljene u tablici 2.1 je

r.br.	vrsta	duljina kljuna	visina kljuna	duljina peraje	masa tijela	spol	godina	K++	EM
20	Adelie	46	21.5	194	4200	m	2007	0	0
21	Adelie	37.8	18.3	174	3400	f	2007	1	0
66	Adelie	41.6	18	192	3950	m	2008	1	0
74	Adelie	45.8	18.9	197	4150	m	2008	0	1
155	Gentoo	48.7	14.1	210	4450	f	2007	0	2
189	Gentoo	42.6	13.7	213	4950	f	2008	2	2
277	Chinstrap	50	19.5	196	3900	m	2007	1	1
335	Chinstrap	45.6	19.4	194	3525	f	2009	1	1

Tablica 2.2: Tablični prikaz pridruživanja jedinki pingvina

U stupcu **K++** dane su oznake komponenata određene algoritmom K++, dok su u stupcu **EM** navedeni identifikatori komponenata određenih EM algoritmom.

Iz originalnih podataka o istraživanju znamo za svaku jedinku kojoj vrsti ona pripada pa možemo provjeriti točnost pridruživanja dobivenih EM algoritmom. Iz populacije 342 jedinke pingvina tek su četiri krivo pridružene, stoga je postotak točnog pridruživanja 98.83%. Radi se o dvije jedinke vrste Adelie koje su pridružene vrsti Chinstrap i o dvije jedinke vrste Chinstrap pridružene vrsti Adelie. To nam daje naslutiti da su diferencijacijska obilježja vrste Gentoo najizraženija. U tablici 2.2 vidimo da je jedinka vrste Adelie pod rednim brojem 74 pogrešno pridružena vrsti Chinstrap. Promjenom uvjeta zaustavljanja algoritma na način da smanjimo dozvoljenu razliku među parametrima, nismo izbjegli krivo pridruživanje za ove četiri jedinke. Potpuno točno pridruživanje dobili smo samo u slučaju kada za uvjet zaustavljanja zahtjevamo ispunjenje obje tvrdnje uvjeta. Tada je potrebno čak 30% više iteracija da se osigura konvergencija. Navodimo procijenjene distribucijske i težinske parametre komponenata miješane populacije pingvina te prikazujemo dijagram rasipanja po vrstama:

Adelie: $\eta_A = 0.446$, $\mu_A = (38.81, 18.32, 189.71, 3691.57)$

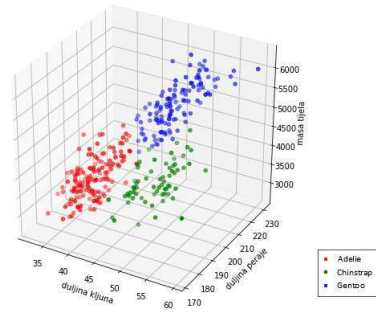
$$\Sigma_A = \begin{bmatrix} 7.00 & 1.17 & 4.46 & 622.21 \\ 1.17 & 1.49 & 2.53 & 324.89 \\ 4.46 & 2.53 & 39.94 & 1351.78 \\ 622.21 & 324.89 & 1351.78 & 208064.76 \end{bmatrix} \quad (2.22)$$

Chinstrap: $\eta_C = 0.194$, $\mu_C = (49.00, 18.48, 196.52, 3754.61)$

$$\Sigma_C = \begin{bmatrix} 9.97 & 2.10 & 7.12 & 531.50 \\ 2.10 & 1.21 & 4.00 & 240.67 \\ 7.12 & 4.00 & 48.09 & 1673.55 \\ 531.50 & 240.67 & 1673.55 & 144047.84 \end{bmatrix} \quad (2.23)$$

Gentoo: $\eta_G = 0.360$, $\mu_G = (47.50, 14.98, 217.19, 5076.02)$

$$\Sigma_G = \begin{bmatrix} 9.42 & 1.93 & 13.11 & 1031.17 \\ 1.93 & 0.95 & 4.46 & 352.80 \\ 13.11 & 4.46 & 41.71 & 2278.47 \\ 1031.17 & 352.80 & 2278.47 & 252067.11 \end{bmatrix} \quad (2.24)$$



Slika 2.2: Dijagram rasipanja po komponentama

Prikazan je dijagram rasipanja populacije pingvina po vrstama za vrijednosti numeričkih varijabli *duljina_kljuna*, *duljina_peraje* i *masa_tijela*. Možemo primijetiti da se razmještaji pingvina vrsta Adelle i Chinstrap preklapaju.

Algoritmu je trebalo 58 iteracija do postizanja konvergencije što vremenski odgovara trajanju od 5 minuta. Brzinu konvergencije računali smo formulom (2.20) te ona iznosi $r = 0.620$. Kada bi red veličine promatrane populacije bio veći 10^m , $m \in \mathbb{N}$ puta, izvjesno

je da bi izvođenje algoritma trajalo znatno dulje. Sljedeći korak modeliranja populacije bila bi optimizacija računalnog algoritma. Smanjenje broja potrebnih iteracija možemo postići korištenjem neopadajućeg preslikavanja M u koraku maksimizacije. Budući da promatramo mješavinu triju višedimenzionalnih normalnih distribucija, trebalo je procijeniti $3K$ parametara. Spomenimo još Akaikeov i Bayesov informacijski kriterij koji penaliziraju modele s velikim brojem parametara, a oni za ovaj model iznose $AIC = 10319.376$ te $BIC = 10353.889$. Uz pretpostavku da populacija pingvina ima četiri ili više komponenata, računali smo pripadne informacijske kriterije i primijetili da se oni povećavaju kako broj komponenata raste. Time izvodimo zaključak da se optimalan model promatrane mješavine sastoji od tri komponente što i odgovara stvarnoj populaciji koju modeliramo. Populacija pingvina na kojoj smo implementirali algoritam pokazala se kao idealna za demonstraciju modeliranja Gaussovih mješavina.

Dodatak A

Računalni kôdovi u Pythonu

Ovdje su navedeni kôdovi korišteni za implementiranje EM algoritma za modeliranje populacije pingvina otočja Palmer. Za pisanje kôdova korišten je jezik *Python 3.7.4* i njegovi dodatni paketi *numpy* [10], *pandas* [6], *time*, *math*, *random*, *sys* [14], *scipy.stats* [16] i *matplotlib* [5].

A.1 Algoritam maksimizacije očekivanja

```
1 import numpy as np
2 import pandas as pd
3 import time
4 from time import sleep
5 import math
6 from scipy.stats import multivariate_normal as mvn
7 import linalg
8 import sys
9
10 start = time.time()
11
12 # ----- FUNKCIJE -----#
13 def postignuta_konvergencija(eta_s, eta_n, mu_s, mu_n, sigma_s, sigma_n,
14     delta_ograda, epsilon_ograda):
15     zaustavljanje = 0
16     log_vj_s = log_vjerodostojnost (eta_s, mu_s, sigma_s)
17     log_vj_n = log_vjerodostojnost (eta_n, mu_n, sigma_n)
18     if ( ( np.linalg.norm(np.asarray(eta_novi) - np.asarray(eta_stari)
19         ,2) < delta_ograda ) and ( np.linalg.norm(np.asarray(mu_novi) - np.
20         asarray(mu_stari),2) < delta_ograda ) and ( np.linalg.norm(
21         sigma_stari[0].as_matrix() - sigma_novi[0], 2) < delta_ograda ) and
22         ( np.linalg.norm(sigma_stari[1].as_matrix() - sigma_novi[1], 2) <
23         delta_ograda ) and ( np.linalg.norm(sigma_stari[2].as_matrix() -
```

```
sigma_novi[2], 2) < delta_ograda ) ) or (np.abs(log_vj_n - log_vj_s
) < epsilon_ograda) ):
18     zaustavljanje = 1 #uvjet zaustavljanja je zadovoljen i
    postignuta je konvergencija
19
20     return(zaustavljanje)
21
22 def log_vjerodostojnost (trenutna_tezina, trenutno_ocekivanje,
    trenutna_kovarijacijska):
23     log_vj = 0
24     for i in range (N):
25         pomocna_suma = 0
26         vektor = znacajke.loc[i][: -1]
27         for k in range(K):
28             pomocna_suma += trenutna_tezina[k] * mvn.pdf(vektor, mean=
    trenutno_ocekivanje[k], cov=trenutna_kovarijacijska[k])
29         logaritam = np.log(pomocna_suma)
30         log_vj += logaritam
31     return(log_vj)
32
33 def brzina_konvergencije (eta_ps, eta_s, eta_n, mu_ps, mu_s, mu_n,
    sigma_ps, sigma_s, sigma_n):
34     brzina_konvergencije_eta = (np.linalg.norm(np.asarray(eta_n) - np.
    asarray(eta_s), 2))/(np.linalg.norm(np.asarray(eta_ps) - np.asarray(
    eta_s), 2))
35     brzina_konvergencije_mu = (np.linalg.norm(np.asarray(mu_n) - np.
    asarray(mu_s), 2))/(np.linalg.norm(np.asarray(mu_ps) - np.asarray(
    mu_s), 2))
36     brzina_konvergencije_sigma0 = (np.linalg.norm(np.matrix(sigma_n[0])
    - np.matrix(sigma_s[0]), 2))/(np.linalg.norm(np.matrix(sigma_ps[0])
    - np.matrix(sigma_s[0]), 2))
37     brzina_konvergencije_sigma1 = (np.linalg.norm(np.matrix(sigma_n[1])
    - np.matrix(sigma_s[1]), 2))/(np.linalg.norm(np.matrix(sigma_ps[1])
    - np.matrix(sigma_s[1]), 2))
38     brzina_konvergencije_sigma2 = (np.linalg.norm(np.matrix(sigma_n[2])
    - np.matrix(sigma_s[2]), 2))/(np.linalg.norm(np.matrix(sigma_ps[2])
    - np.matrix(sigma_s[2]), 2))
39     brzina = max(brzina_konvergencije_eta, brzina_konvergencije_mu,
    brzina_konvergencije_sigma0, brzina_konvergencije_sigma1,
    brzina_konvergencije_sigma2)
40     return(brzina)
41
42 def aic (broj_komponenti, log_v):
43     br_procijenjenih = 3*broj_komponenti
44     broj = 2* (br_procijenjenih - log_v)
45     return(broj)
46
```

```

47 def bic (broj_komponenti, broj_opazanja, log_v):
48     n = broj_opazanja
49     br_procijenjenih = 3*broj_komponenti
50     k = br_procijenjenih
51     broj = k * np.log(n) - 2 * (log_v)
52     return(broj)
53 # ----- FUNKCIJE -----#
54
55 podaci = pd.read_csv('potpuni_podaci.csv',header=0)
56 znacajke = podaci.copy()
57 z = znacajke.drop(['vrsta', 'otok', 'spol', 'godina'], axis=1)
58 znacajke = z
59
60 # ----- INICIJALIZACIJA ----- #
61 z0 = znacajke[znacajke["kMeans_pridruzivanje"].isin([0])]
62 znacajke_0 = z0.drop(['kMeans_pridruzivanje'], axis=1)
63
64 z1 = znacajke.loc[znacajke['kMeans_pridruzivanje'] == 1]
65 znacajke_1 = z1.drop(['kMeans_pridruzivanje'], axis=1)
66
67 z2 = znacajke.loc[znacajke['kMeans_pridruzivanje'] == 2]
68 znacajke_2 = z2.drop(['kMeans_pridruzivanje'], axis=1)
69
70 znacajke_num = znacajke.drop(['kMeans_pridruzivanje'], axis =1)
71
72 mu0 = znacajke_0.mean().astype(np.int)
73 mu1 = znacajke_1.mean().astype(np.int)
74 mu2 = znacajke_2.mean().astype(np.int)
75
76 mu_stari = (mu0, mu1, mu2)
77
78 sigma0 = znacajke_0.cov()
79 sigma1 = znacajke_1.cov()
80 sigma2 = znacajke_2.cov()
81
82 sigma_stari = (sigma0, sigma1, sigma2)
83
84 velicina_komponente0 = np.shape(znacajke_0)[0]
85 velicina_komponente1 = np.shape(znacajke_1)[0]
86 velicina_komponente2 = np.shape(znacajke_2)[0]
87
88 N = np.shape(znacajke)[0]
89 K = 3
90
91 eta0 = 1/N*velicina_komponente0
92 eta1 = 1/N*velicina_komponente1
93 eta2 = 1/N*velicina_komponente2

```

```
94
95 eta_stari = (eta0, eta1, eta2)
96
97
98 log_vjerodostojnost_stara = log_vjerodostojnost (eta_stari, mu_stari,
    sigma_stari)
99 # ----- INICIJALIZACIJA ----- #
100
101 # ----- 1. ITERACIJA ----- #
102 vjerojatnost = znacajke_num.copy()
103
104 vjerojatnost.rename(columns = {'duljina_kljuna':'komponenta0', '
    visina_kljuna' : 'komponenta1', 'duljina_peraje':'komponenta2'},
    inplace=True)
105 vjerojatnost.drop(['masa_tijela'], inplace = True, axis=1)
106
107 for index, row in znacajke.iterrows():
108     i = index
109     vektor = znacajke.loc[i][: -1]
110     for j in range(3):
111         vjerojatnost.iloc[i,j] = mvn.pdf(vektor, mean=mu_stari[j], cov =
            sigma_stari[j])
112
113 odgovornost = vjerojatnost.copy()
114 odgovornost.rename(columns = {'komponenta0':0, 'komponenta1':1, '
    komponenta2':2}, inplace=True)
115 novo_pridruzivanje = odgovornost.idxmax(axis=1)
116
117 znacajke["kMeans_pridruzivanje"] = novo_pridruzivanje
118
119 z0 = znacajke[znacajke["kMeans_pridruzivanje"].isin([0])]
120 znacajke_0 = z0.drop(['kMeans_pridruzivanje'], axis=1)
121
122 z1 = znacajke.loc[znacajke['kMeans_pridruzivanje'] == 1]
123 znacajke_1 = z1.drop(['kMeans_pridruzivanje'], axis=1)
124
125 z2 = znacajke.loc[znacajke['kMeans_pridruzivanje'] == 2]
126 znacajke_2 = z2.drop(['kMeans_pridruzivanje'], axis=1)
127
128 znacajke_num = znacajke.drop(['kMeans_pridruzivanje'], axis =1)
129
130 for index, row in vjerojatnost.iterrows():
131     i = index
132     for j in range(K):
133         odgovornost.iloc[i,j] = vjerojatnost.iloc[i][j]*eta_stari[j]/np.
            dot(vjerojatnost.loc[i], eta_stari)
134
```

```
135 eta_novi = [1/N*odgovornost.sum(axis=0)[0], 1/N*odgovornost.sum(axis=0)
136           [1], 1/N*odgovornost.sum(axis=0)[2]]
137 mu_novi = []
138 for k in range(K):
139     hy = 0
140     for i in range(N):
141         hy += odgovornost.iloc[i][k]*znacajke_num.loc[i]
142     hy = hy/odgovornost.sum(axis=0)[k]
143     mu_novi.append(hy)
144
145 sigma_n = []
146 for k in range(K):
147     sig = 0
148     for i in range(N):
149         sig += odgovornost.iloc[i][k] * np.outer(((znacajke_num.loc[i]-
150           mu_novi[k]).to_numpy()), np.transpose((znacajke_num.loc[i]-mu_novi[k]
151           ]).to_numpy())) / odgovornost.sum(axis=0)[k]
152     sigma_n.append(sig)
153
154 sigma_novi = [pd.DataFrame(np.matrix(sigma_n[0])), pd.DataFrame(np.
155           matrix(sigma_n[1])), pd.DataFrame(np.matrix(sigma_n[2]))]
156 log_vjerodostojnost_nova = log_vjerodostojnost (eta_novi, mu_novi,
157           sigma_novi)
158
159 # ----- 1. ITERACIJA ----- #
160
161 br_iteracija = 1
162
163 delta = 1e-7
164 epsilon = 1e-5
165
166 if np.abs(log_vjerodostojnost_nova - log_vjerodostojnost_stara) < delta
167 :
168     uvjet = 1
169 else:
170     uvjet = 0
171
172 # ----- PETLJA ----- #
173
174 while uvjet == 0 :
175     br_iteracija += 1
176
177     trenutno_vrijeme = time.time()
178     elapsed = trenutno_vrijeme - start
179     minute = elapsed//60
```

```
176     sekunde = round(elapsed - minute*60)
177
178     sys.stdout.write('\r')
179     sys.stdout.write("Vrijeme izvodjenja koda: [%d min : %d s] ---
Broj trenutne iteracije: %d"%(minute,sekunde, br_iteracija))
180
181     sys.stdout.flush()
182     sleep(0.1)
183
184     mu_prastari = mu_stari
185     eta_prastari = eta_stari
186     sigma_prastari = sigma_stari
187
188     mu_stari = mu_novi.copy()
189     eta_stari = eta_novi.copy()
190     sigma_stari = sigma_novi.copy()
191
192     log_vjerodostojnost_stara = log_vjerodostojnost_nova.copy()
193
194     for index, row in znacajke.iterrows():
195         i = index
196         vektor = znacajke.loc[i][:-1]
197         for j in range(K):
198             vjerojatnost.iloc[i,j] = mvn.pdf(vektor, mean=mu_stari[j],
cov=sigma_stari[j])
199
200     odgovornost = vjerojatnost.copy()
201     odgovornost.rename(columns = {'komponenta0':0, 'komponenta1':1, '
komponenta2':2}, inplace=True)
202
203     novo_pridruzivanje = odgovornost.idxmax(axis=1)
204
205     znacajke["kMeans_pridruzivanje"] = novo_pridruzivanje
206
207     z0 = znacajke[znacajke["kMeans_pridruzivanje"].isin([0])]
208     znacajke_0 = z0.drop(['kMeans_pridruzivanje'], axis=1)
209
210     z1 = znacajke.loc[znacajke['kMeans_pridruzivanje'] == 1]
211     znacajke_1 = z1.drop(['kMeans_pridruzivanje'], axis=1)
212
213     z2 = znacajke.loc[znacajke['kMeans_pridruzivanje'] == 2]
214     znacajke_2 = z2.drop(['kMeans_pridruzivanje'], axis=1)
215
216     znacajke_num = znacajke.drop(['kMeans_pridruzivanje'], axis =1)
217
218     for index, row in vjerojatnost.iterrows():
219         i = index
```

```
220     for j in range(3):
221         odgovornost.iloc[i,j] = vjerojatnost.iloc[i][j]*eta_stari[j]
222         ]/np.dot(vjerojatnost.loc[i], eta_stari)
223
224     eta_novi = [1/N*odgovornost.sum(axis=0)[0], 1/N*odgovornost.sum(axis
225 =0)[1], 1/N*odgovornost.sum(axis=0)[2]]
226
227     mu_novi = []
228     for k in range(K):
229         hy = 0
230         for i in range (N):
231             hy += odgovornost.iloc[i][k]*znacajke_num.loc[i]
232         hy = hy/odgovornost.sum(axis=0)[k]
233         mu_novi.append(hy)
234
235     sigma_n = []
236     for k in range (K):
237         sig = 0
238         for i in range(N):
239             sig += odgovornost.iloc[i][k] * np.outer(((znacajke_num.loc[
240 i]-mu_novi[k]).to_numpy()), np.transpose((znacajke_num.loc[i]-
241 mu_novi[k]).to_numpy())) / odgovornost.sum(axis=0)[k]
242         sigma_n.append(sig)
243
244     sigma_novi = [pd.DataFrame(np.matrix(sigma_n[0])), pd.DataFrame(np.
245 matrix(sigma_n[1])), pd.DataFrame(np.matrix(sigma_n[2]))]
246
247     log_vjerodostojnost_nova = log_vjerodostojnost (eta_novi, mu_novi,
248 sigma_novi)
249
250     je_li_iskonvergiralo = postignuta_konvergencija(eta_stari, eta_novi,
251 mu_stari, mu_novi, sigma_stari, sigma_novi, delta, epsilon)
252     uvjet = je_li_iskonvergiralo
253
254     brzina = brzina_konvergencije(eta_prastari, eta_stari, eta_novi,
255 mu_prastari, mu_stari, mu_novi, sigma_prastari, sigma_stari,
256 sigma_novi)
257
258     print()
259     print()
260     print("Brzina konvergencije je ", brzina)
261     print()
262
263     iznos_a = aic (K, log_vjerodostojnost_nova)
264     iznos_b = bic (K, N, log_vjerodostojnost_nova)
```



```
258
259 print()
260 print("Akaikeov informacijski kriterij iznosi AIC = ", iznos_a)
261 print()
262 print()
263 print("Bayesov informacijski kriterij iznosi BIC = ", iznos_b)
264 print()
265 # ----- PETLJA ----- #
266
267 # ----- KRAJ ALGORITMA ----- #
268
269 with open(r'sredine.txt', 'w') as f1:
270     f1.write("ocekivanje\n")
271     for item in mu_novi:
272         f1.write("%s\n" % item)
273
274 with open(r'kovarijacijske.txt', 'w') as f2:
275     for item in sigma_novi:
276         f2.write("%s\n" % item)
277
278 with open(r'tezine.txt', 'w') as f3:
279     f3.write("tezina\n")
280     for item in eta_novi:
281         f3.write("%s\n" % item)
282
283 print("Nakon ", br_iteracija, ". iteracije, završavamo s algoritmom i
    analiziramo rezultate.")
284
285 podaci['EM_pridruzivanje'] = novo_pridruzivanje
286 podaci.to_csv(r'konacna_tablica.csv', header=True, index=False)
287
288 trenutno_vrijeme = time.time()
289 end = time.time()
290 elapsed = end - start
291 minute = elapsed//60
292 sekunde = round(elapsed - minute*60)
293
294 print ("Ukupno vrijeme za postizanje konvergencije pod ovim uvjetima je
    ", minute,"minuta i ",sekunde, "sekundi.")
295 print("vektor tezina:", eta_novi)
296 print("vektor ocekivanja:", mu_novi)
297 print("kovarijacijske matrice:", sigma_novi)
298
299 print("----- KRAJ -----")
```

A.2 Algoritam k-srednjih vrijednosti ++

```
1 import numpy as np
2 import pandas as pd
3 import time
4 import random
5 import math
6
7 penguins_data = pd.read_csv("penguins.csv")
8 penguins = penguins_data.drop(columns=['Unnamed: 0'])
9 podaci = penguins.rename(columns = {'species':'vrsta', 'island':'otok',
   'bill_length_mm':'duljina_kljuna', 'bill_depth_mm':'visina_kljuna',
   'flipper_length_mm':'duljina_peraje', 'body_mass_g':'masa_tijela',
   'sex':'spol', 'year':'godina'})
10 potpuni_podaci = podaci.dropna(subset=['duljina_kljuna', 'visina_kljuna',
   'duljina_peraje', 'masa_tijela'])
11 opazanja = potpuni_podaci.reset_index()[['duljina_kljuna', '
   visina_kljuna', 'duljina_peraje', 'masa_tijela']].values.tolist()
12
13 br_tocaka = len(opazanja)
14 tocke = np.array(opazanja)
15 dimenzija_tocaka = np.size(tocke[0])
16 K = 3
17 br_centara = K
18
19 def euklidska_metrika(x,y):
20     dist = 0
21     for i in range(len(x)):
22         dist = dist + (x[i]-y[i])*(x[i]-y[i])
23     return dist
24
25 # inicijalizacija centroida
26 centri = []
27 centri.append(tocke[random.randint(0,br_tocaka-1)][:])
28 centri
29
30 for i in range(br_centara-1):
31     maxi = 0
32     pom = 0;
33
34     for j in range(br_tocaka):
35         for z in range(np.size(centri, axis=0)):
36             if euklidska_metrika(tocke[j][:],centri[z]) > maxi:
37                 maxi = euklidska_metrika(tocke[j][:],centri[z])
38                 pom = j
```

```
39     centri.append(tocke[pom][:])
40
41 group_id = []
42 for i in range(br_tocaka):
43     group_id.append(0)
44
45 br_iteracija = 100
46
47 for temp in range(br_iteracija):
48     # korak pridruzivanja
49     for i in range(br_tocaka):
50         udaljenost_trenutna = euklidska_metrika(tocke[i], centri[0])
51         k = 0
52
53         for j in range(br_centara):
54             udaljenost_pom = euklidska_metrika(tocke[i], centri[j])
55             if udaljenost_pom < udaljenost_trenutna:
56                 udaljenost_trenutna = udaljenost_pom
57                 k = j
58         group_id[i] = k
59
60     # korak racunanja novih centroida
61     for i in range(br_centara):
62         niz = []
63         for j in range(dimenzija_tocaka):
64             niz.append(0)
65         cluster_count = 0
66         for j in range(br_tocaka):
67             if group_id[j] == i:
68                 cluster_count += 1
69                 niz += tocke[j]
70
71         centil95 = np.percentile(centri[i], 95)
72         if cluster_count > 0:
73             for j in range(dimenzija_tocaka):
74                 if centri[i][j] > centil95:
75                     centri[i][j] = niz[j]*1.0/cluster_count
76
77 indeksi = list(potpuni_podaci.index.values)
78 pridruzivanja = pd.Series(group_id, index = indeksi)
79 potpuni_podaci['kMeans_pridruzivanje'] = pridruzivanja
80 potpuni_podaci.to_csv(r'potpuni_podaci.csv', header=True, index=False)
```

Bibliografija

- [1] A. Dempster, N. Laird i D. Rubin, Maximum Likelihood from Incomplete Data via the EM Algorithm, *Journal of the Royal Statistical Society. Series B (Methodological)* **39** (1977), br. 1, 1–38.
- [2] J. Demšar et al., Orange: Data Mining Toolbox in Python, *Journal of Machine Learning Research* **14** (2013), 2349–2353, <http://jmlr.org/papers/v14/demsar13a.html>, version 3.23.1.
- [3] W. Feller, On a General Class of "Contagious" Distributions, *Annals of Mathematical Statistics* **14** (1943), br. 4, 389–400.
- [4] S. Frühwirth-Schnatter, Finite Mixture and Markov Switching Models, Springer Science + Business Media, LLC, New York, 2006.
- [5] J. D. Hunter, Matplotlib: A 2D graphics environment, *Computing in science & engineering* **9** (2007), br. 3, 90–95.
- [6] W. McKinney et al., Data structures for statistical computing in python, *Proceedings of the 9th Python in Science Conference*, sv. 445, Austin, TX, 2010, str. 51–56.
- [7] G. J. McLachlan, S. X. Lee i S. I. Rathnayake, Finite Mixture Models, *Annual Review of Statistics and Its Application* **6** (2019), br. 1, 355–378.
- [8] G. J. McLachlan i D. Peel, Finite Mixture Models, *Wiley Series in Probability and Statistics*, John Wiley & Sons, Inc, New York, 2000.
- [9] S. K. Ng, T. Krishnan i G. J. McLachlan, The EM algorithm, *Papers/Humboldt-Universität Berlin, Center for Applied Statistics and Economics (CASE)* **24** (2004).
- [10] T. E. Oliphant, A guide to NumPy, **1**, Trelgol Publishing USA, 2006.
- [11] K. Pearson, Contributions to the Mathematical Theory of Evolution, *Philosophical Transactions of the Royal Society of London A*, **185** (1894), 71–110.

- [12] P. E. Rossi, G. M. Allenby i R. McCulloch, Bayesian Statistics and Marketing, Wiley Series in Probability and Statistics, John Wiley & Sons, Ltd, West Sussex, 2005.
- [13] N. Sarapa, Teorija vjerojatnosti, Školska knjiga, Zagreb, 2002.
- [14] G. Van Rossum, The Python Library Reference, release 3.8.2, Python Software Foundation, 2020.
- [15] G. Van Rossum i F. L. Drake, Python 3 Reference Manual, CreateSpace, Scotts Valley, CA, 2009, version 3.7.4.
- [16] P. Virtanen et al., SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python, 2020.
- [17] C. F. J. Wu, On the Convergence Properties of the EM Algorithm, The Annals of Statistics **11** (1983), br. 1, 95–103.
- [18] J. Šnajder i B. Dalbelo Bašić, Strojno učenje, 2014, https://www.fer.unizg.hr/_download/repository/StrojnoUcenje.pdf, datum pristupa (11.10.2019.).

Sažetak

Model Gaussovih mješavina često se pojavljuje u raznim znanstvenim istraživanjima. Mješavina je populacija koja se sastoji od komponenata čije vjerojatnosne distribucije mogu, ali ne moraju nužno, pripadati istoj familiji distribucija. U ovom smo radu proučavali Gaussove mješavine čiji je broj komponenata konačan i unaprijed poznat. Za takve mješavine konstruirali smo model tako što smo procijenili distribucijske parametre svake komponente pa komponentama pridružili opažanja. Koristili smo metode procjene parametara iz teorije frekvencionističke i Bayesovske statistike koje su na kraju i objedinjene algoritmom maksimizacije očekivanja (*EM algoritam*). Svaki model mješavine potrebno je prilagoditi dostupnim informacijama o populaciji. Najjednostavniji slučaj koji modeliramo je kada su poznate gustoće komponenata pa im još preostaje pridružiti jedinke korištenjem Bayesovog naivnog klasifikatora. Ako su gustoće komponenata nepoznate, ali je poznata realizacija pridruživanja svakog opažanja, potrebno je samo procijeniti parametre gustoća. U tom smo slučaju procjenu proveli metodom maksimalne vjerodostojnosti kao i Bayesovom metodom. Za zadnji slučaj smo ostavili primjer populacije za koju unaprijed ne znamo niti parametre gustoća komponenata niti realizaciju pridruživanja jedinki. Taj je problem jednostavno riješen pomoću EM algoritma koji je u radu teorijski objašnjen, ali i realiziran na stvarnoj populaciji. Proces modeliranja Gaussovih mješavina ilustrirali smo na populaciji tri vrste pingvina i pritom za svaku jedinku odredili pridruživanje pomoću EM algoritma. Model smo izrazili kao linearnu kombinaciju višedimenzionalnih normalnih gustoća čije smo parametre procijenili.

Summary

Gaussian mixture model occurs commonly in various scientific research. A mixture is a population that consists of components whose probability distributions could, but should not necessarily, belong to the same family of distributions. This thesis studies Gaussian mixtures whose number of components is finite and known in advance. For those mixtures, we constructed a model in a way that the distribution parameters were estimated and observations were assigned to the components. Parameter estimation methods that were used come both from frequentist and Bayesian statistics and are consolidated with the Expectation–Maximization algorithm (*the EM algorithm*). Each mixture model should be adjusted according to the available information on the population. The simplest case that is modelled is the one where densities of the components are known and only allocation of the observation remains. Allocation is conducted by using the naïve Bayes classifier. In the case where component densities are unknown and the allocation of the observations is known, only parameter estimation is required. Unknown parameters are estimated by using the method of maximum likelihood estimation along with the Bayesian method of estimation. The last studied case is the example of a population for which both the component parameters and the allocation is unknown. That problem is easily solved by using the EM algorithm which is in this thesis theoretically explained and implemented on the real–life population. The process of modelling Gaussian mixtures is illustrated on the population of three penguin species for which all the units are allocated to the components by using the EM algorithm. The model is given as a linear combination of multivariate normal distributions whose parameters were estimated.

Životopis

Rođena sam 14. kolovoza 1994. u Zagrebu. Srednju školu, zagrebačku XV. gimnaziju, upisujem 2009. godine, a preddiplomski studij Matematika na Prirodoslovno matematičkom fakultetu – Matematičkom odsjeku 2013. godine. Diplomski studij Matematička statistika upisujem na istom fakultetu 2018. godine.