

# Clusteranalyse

## Aufgabe 1: Clusteranalyse mit k-Means

In der folgenden Tabelle seien der *weight index* und der *pH index* von drei Medikamenten gegeben:

Objekt	weight index	pH index
Medizin A	1	1
Medizin B	2	1
Medizin C	4	3

- a) Gruppieren Sie die Medikamente mithilfe des k-means Clusterverfahrens in zwei Cluster. Medikamente A und B sollen dabei die Startpartition für je eine Klasse sein. Als Distanzmaß soll die euklidische Distanz verwendet werden.
- b) Gehen Sie nun davon aus, dass drei zusätzliche Medikamente in der Clusteranalyse mitbezogen werden:

Objekt	weight index	pH index
Medizin D	5	4
Medizin E	3	2
Medizin F	6	4

Dabei ergeben sich für das k-means Verfahren sowie für das Single Linkage Verfahren folgende Clusterungen

Medikament	A	B	C	D	E	F
k-means	2	2	1	1	2	1
Single Linkage	1	1	1	2	1	2

Geben Sie eine geeignete Maßzahl für die Ähnlichkeit der beiden Clusterungen an und berechnen Sie diese.

- c) **Zum Selbststudium:** Bestimmen Sie die Cluster, die in Teilaufgabe b) für das k-means Verfahren sowie das Single Linkage Verfahren angegeben sind, sowohl von Hand als auch mit R.

**Lösung:**

a) Gesucht: 2 Klassen / Cluster mit optimaler Zuordnung der Medikamente.

Distanz: euklidische Distanz

(1) Iteration 0:  $\mathbf{x}_1^{(0)} = (1, 1)^T$ ,  $\mathbf{x}_2^{(0)} = (2, 1)^T$

$$\mathbf{D}^{(0)} = \begin{pmatrix} \textcircled{0} & 1 & 3,61 \\ 1 & \textcircled{0} & \textcircled{2,83} \end{pmatrix} \begin{array}{l} \rightarrow \text{Distanz zu } \mathbf{x}_1^{(0)} \\ \rightarrow \text{Distanz zu } \mathbf{x}_2^{(0)} \end{array}$$

$A \quad B \quad C$  (umkringelt  $\hat{=}$  minimales Element pro Spalte)

$$\text{z.B. } D(C, \mathbf{x}_1^{(0)}) = \sqrt{(4-1)^2 + (3-1)^2} = 3,61$$

$$\Rightarrow \mathcal{C}^{(0)} = \{\{A\}, \{B, C\}\}$$

(2) Iteration 1:

$$\begin{aligned} \mathbf{x}_1^{(1)} &= (1, 1)^T & \{A\} \\ \mathbf{x}_2^{(1)} &= \left(\frac{2+4}{2}, \frac{1+3}{2}\right)^T = (3, 2)^T & \{B, C\} \end{aligned}$$

$$\mathbf{D}^{(1)} = \begin{pmatrix} \textcircled{0} & \textcircled{1} & 3,61 \\ 2,24 & 1,41 & \textcircled{1,41} \end{pmatrix} \begin{array}{l} \rightarrow \text{Distanz zu } \mathbf{x}_1^{(1)} \\ \rightarrow \text{Distanz zu } \mathbf{x}_2^{(1)} \end{array}$$

$A \quad B \quad C$  (umkringelt  $\hat{=}$  minimales Element pro Spalte)

$$\text{z.B. } D(C, \mathbf{x}_2^{(1)}) = \sqrt{(4-3)^2 + (3-2)^2} = 1,41$$

$$\Rightarrow \mathcal{C}^{(1)} = \{\{A, B\}, \{C\}\}$$

(3) Iteration 2:

$$\begin{aligned} \mathbf{x}_1^{(2)} &= \left(\frac{1+2}{2}, \frac{1+1}{2}\right)^T = \left(\frac{3}{2}, 1\right)^T & \{A, B\} \\ \mathbf{x}_2^{(2)} &= (4, 3)^T & \{C\} \end{aligned}$$

$$\mathbf{D}^{(2)} = \begin{pmatrix} \textcircled{0,5} & \textcircled{0,5} & 3,20 \\ 3,61 & 2,83 & \textcircled{0} \end{pmatrix} \begin{array}{l} \rightarrow \text{Distanz zu } \mathbf{x}_1^{(2)} \\ \rightarrow \text{Distanz zu } \mathbf{x}_2^{(2)} \end{array}$$

$A \quad B \quad C$  (umkringelt  $\hat{=}$  minimales Element pro Spalte)

$$\Rightarrow \mathcal{C}^{(2)} = \mathcal{C}^{(1)} \Rightarrow \text{Ende}$$

---

## Wiederholung:

(Adjustierter) Rand-Index:

Ziel: Quantifizierung der Ähnlichkeit der Ergebnisse zweier Clusterungen

Grundlage: Aufstellen einer Vierfeldertafel bzgl. der Einordnung aller Objektpaare

		Clusterung B	
		1	0
Clusterung A	1	a	b
	0	c	d

- a entspricht der Anzahl aller Objektpaare, die sowohl für Clusterung A als auch für Clusterung B im gleichen Cluster sind.
- b entspricht der Anzahl aller Objektpaare, die für Clusterung A im gleichen Cluster aber für Clusterung B in verschiedenen Clustern sind.
- c entspricht der Anzahl aller Objektpaare, die für Clusterung A in verschiedenen Clustern sind aber für Clusterung B gleichen Cluster sind.
- d entspricht der Anzahl aller Objektpaare, die sowohl für Clusterung A als auch für Clusterung B in zwei verschiedenen Clustern sind

Dieses Verfahren kann auf beliebige Clusteranzahlen verallgemeinert werden. Dabei brauchen die Clusterungen A und B *nicht* die gleiche Clusteranzahl aufweisen.

Rand-Index:  $R = \frac{a+d}{a+b+c+d} = \frac{a+d}{\binom{n}{2}}$

Problematik: Zufällige Übereinstimmungen

⇒ **Adjustierter Rand-Index** korrigiert dafür:  $R_{adj} = \frac{\binom{n}{2}(a+d) - [(a+b)(a+c) + (c+d)(b+d)]}{\binom{n}{2}^2 - [(a+b)(a+c) + (c+d)(b+d)]}$

(Quelle: Santos & Embrechts (2009))

---

Invertiere zunächst eine der Clusterungen → größere Übereinstimmung der Cluster → bessere Übersichtlichkeit:

Medikament	A	B	C	D	E	F
k-means	1	1	2	2	1	2
Single Linkage	1	1	1	2	1	2

Vierfeldertafeln:

Vergleiche zunächst Medikament A mit Medikament B – F:		1	0
	1	2	–
	0	1	2

Medikament B verglichen mit Medikament C – F:		1	0
	1	1	–
	0	1	2

...

		Single Linkage	
Insgesamt ergibt sich:		1	0
	k-means	4	2
		3	6

⇒ Adj. Rand-Index:

$$\begin{aligned}
 R_{adj} &= \frac{\binom{6}{2}(4+6) - [(4+2)(4+3) + (3+6)(2+6)]}{\binom{6}{2}^2 - [(4+2)(4+3) + (3+6)(2+6)]} \\
 &= \frac{15 \cdot 10 - [6 \cdot 7 + 9 \cdot 8]}{225 - [6 \cdot 7 + 9 \cdot 8]} \\
 &= \frac{150 - [42 + 72]}{225 - [42 + 72]} = \frac{36}{111} = 0,3243
 \end{aligned}$$

Zusätzliche Ergänzung: nicht-adj. Rand-Index:

$$R = \frac{4+6}{\binom{6}{2}} = \frac{10}{15} \approx 0,67$$

→ Die Clusterungen stimmen nur mäßig überein. Zwar wurde nur Medikament C anders eingeteilt, bei insgesamt lediglich 6 Medikamenten scheint dies aber bereits einen großen Unterschied zu machen (bereits viele Übereinstimmungen durch Zufälligkeiten zu erwarten).

## Aufgabe 2: Modellbasiertes Clustering

Gehen Sie davon aus, dass Sie eine modellbasierte Clusteranalyse auf bivariaten Daten durchführen. Dazu liegt Ihnen das Ergebnis eines entsprechenden Ansatzes in Form der resultierenden Mischverteilung vor. Konkret sei diese für eine Clusteranzahl von 2 gegeben durch

$$\hat{f}(\mathbf{x}) = \sum_{k=1}^2 \hat{p}(k) f(\mathbf{x} | \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}}_k),$$

wobei  $f$  der Dichte einer bivariaten Normalverteilung entspricht mit den Parametern

$$\hat{\boldsymbol{\mu}}_1 = \begin{pmatrix} -2 \\ 2 \end{pmatrix}, \quad \hat{\boldsymbol{\Sigma}}_1 = \begin{pmatrix} \sigma_{1;1}^2 & \sigma_{12;1} \\ \sigma_{12;1} & \sigma_{2;1}^2 \end{pmatrix}, \quad \hat{\boldsymbol{\mu}}_2 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad \hat{\boldsymbol{\Sigma}}_2 = \begin{pmatrix} \sigma_{1;2}^2 & \sigma_{12;2} \\ \sigma_{12;2} & \sigma_{2;2}^2 \end{pmatrix}$$

und mit  $\hat{p}(1) = \hat{p}(2) = 0,5$ . Außerdem seien drei in dem Datensatz enthaltene Datenpunkte konkret angegeben mit  $\mathbf{x}_1 = (0, 0)^T$ ,  $\mathbf{x}_2 = (-1, 1)^T$  und  $\mathbf{x}_3 = (1, 3)^T$  (welche nach Annahme ebenfalls für die Berechnung der Mischverteilung berücksichtigt wurden).

- Machen Sie sich mit dem Konzept des (computergestützten) modellbasierten Clusters vertraut. Betrachten Sie hierbei auch die R-Funktion `Mclust()` aus dem Paket `mclust`. Lesen Sie sich dazu insbesondere die Hilfe zu `mclustModelNames` durch und machen Sie sich die einzelnen Modelle verständlich. Wie erhält man generell aus der Mischverteilung eine Partitionierung der Daten?
- Welchen der beiden Cluster werden die Datenpunkte  $\mathbf{x}_1$ ,  $\mathbf{x}_2$  und  $\mathbf{x}_3$  für das Modell *spherical*, *equal volume* zugeordnet? Veranschaulichen Sie dies in einer Skizze.
- Überlegen Sie sich jeweils konkrete Kovarianzmatrizen  $\hat{\boldsymbol{\Sigma}}_1$  und  $\hat{\boldsymbol{\Sigma}}_2$ , für die unter den folgenden Modellen die jeweils vorgegebene Zuordnung getroffen wird:

	Modell	Zuordnung
i)	<i>spherical, unequal volume</i>	$\mathbf{x}_1, \mathbf{x}_2 \mapsto C_1, \mathbf{x}_3 \mapsto C_2$
ii)	<i>diagonal, equal volume and shape</i>	$\mathbf{x}_2, \mathbf{x}_3 \mapsto C_1, \mathbf{x}_1 \mapsto C_2$
iii)	<i>ellipsoidal, varying volume, shape and orientation</i>	$\mathbf{x}_3 \mapsto C_1, \mathbf{x}_1, \mathbf{x}_2 \mapsto C_2$

*Hinweise:*

- Die Dichte der multivariaten Normalverteilung lässt sich mittels der Funktion `dmvnorm()` aus dem Paket `mvtnorm` berechnen.
- In dem R-Code zu Aufgabe 3 (b) auf dem 2. Übungsblatt können Sie sich noch einmal die Zusammenhänge zwischen den Parametern und der Form der multivariaten Normalverteilung veranschaulichen.

### a) Details

*The following models are available in package mclust:*

#### univariate mixture

“E” = equal variance (one-dimensional)

“V” = variable variance (one-dimensional)

**multivariate mixture**

“EII” = spherical, equal volume

“VII” = spherical, unequal volume

“EEI” = diagonal, equal volume and shape

“VEI” = diagonal, varying volume, equal shape

“EVI” = diagonal, equal volume, varying shape

“VVI” = diagonal, varying volume and shape

“EEE” = ellipsoidal, equal volume, shape, and orientation

“EVE” = ellipsoidal, equal volume and orientation

“VEE” = ellipsoidal, equal shape and orientation

“VVE” = ellipsoidal, equal orientation

“EEV” = ellipsoidal, equal volume and equal shape

“VEV” = ellipsoidal, equal shape

“EVV” = ellipsoidal, equal volume

“VVV” = ellipsoidal, varying volume, shape, and orientation

**single component**

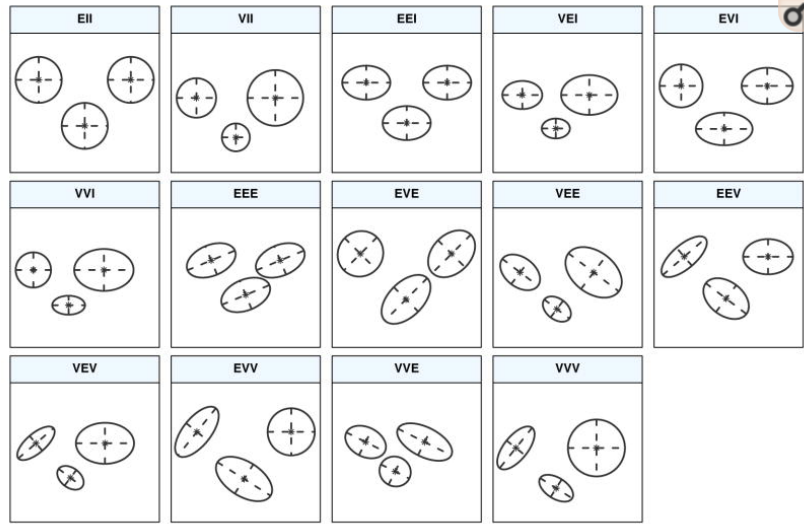
“X” = univariate normal

“XII” = spherical multivariate normal

“XXI” = diagonal multivariate normal

“XXX” = ellipsoidal multivariate normal

- Die Gruppeneinteilung erfolgt anhand dessen, wie plausibel es für das jeweilige Objekt ist einer bestimmten Gruppe anzugehören.
- Kennzahl: Posteriori-Dichte
- Diese ergibt sich für das Objekt  $i$  mit Merkmalsvektor  $x_i$  aus dem Produkt der Priori-Wahrscheinlichkeit und der Dichte der jeweiligen Gruppenverteilung im Punkt  $x_i$ .
- Ein Objekt wird derjenigen Gruppe zugeordnet, für die es die höchste Posteriori-Dichte aufweist.



Ellipses of isodensity for each of the 14 Gaussian models obtained by eigen-decomposition in case of three groups in two dimensions.

**Quelle:** mclust 5: Clustering, Classification and Density Estimation Using Gaussian Finite Mixture Models. Scrucca, L. et al.

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5096736/>

b) *spherical, equal volume*  $\Rightarrow \hat{\Sigma}_1 = \hat{\Sigma}_2 = \begin{pmatrix} \alpha & 0 \\ 0 & \alpha \end{pmatrix}$

Damit gilt:

$$\det(\hat{\Sigma}_k) = \alpha^2 \quad \text{und} \quad \hat{\Sigma}_k^{-1} = \begin{pmatrix} 1/\alpha & 0 \\ 0 & 1/\alpha \end{pmatrix}$$

Dichtefunktion:

$$\begin{aligned} f(\mathbf{x}|\hat{\mu}_k, \hat{\Sigma}_k) &= \frac{1}{\sqrt{(2\pi)^2 \det(\hat{\Sigma}_k)}} \exp \left\{ -\frac{1}{2} \left( (\mathbf{x} - \hat{\mu}_k)^\top \hat{\Sigma}_k^{-1} (\mathbf{x} - \hat{\mu}_k) \right) \right\} \\ &= \frac{1}{2\pi\alpha} \exp \left\{ -\frac{1}{2\alpha} \left( \underbrace{(\mathbf{x} - \hat{\mu}_k)^\top (\mathbf{x} - \hat{\mu}_k)}_{=||\mathbf{x} - \hat{\mu}_k||^2} \right) \right\} \end{aligned}$$

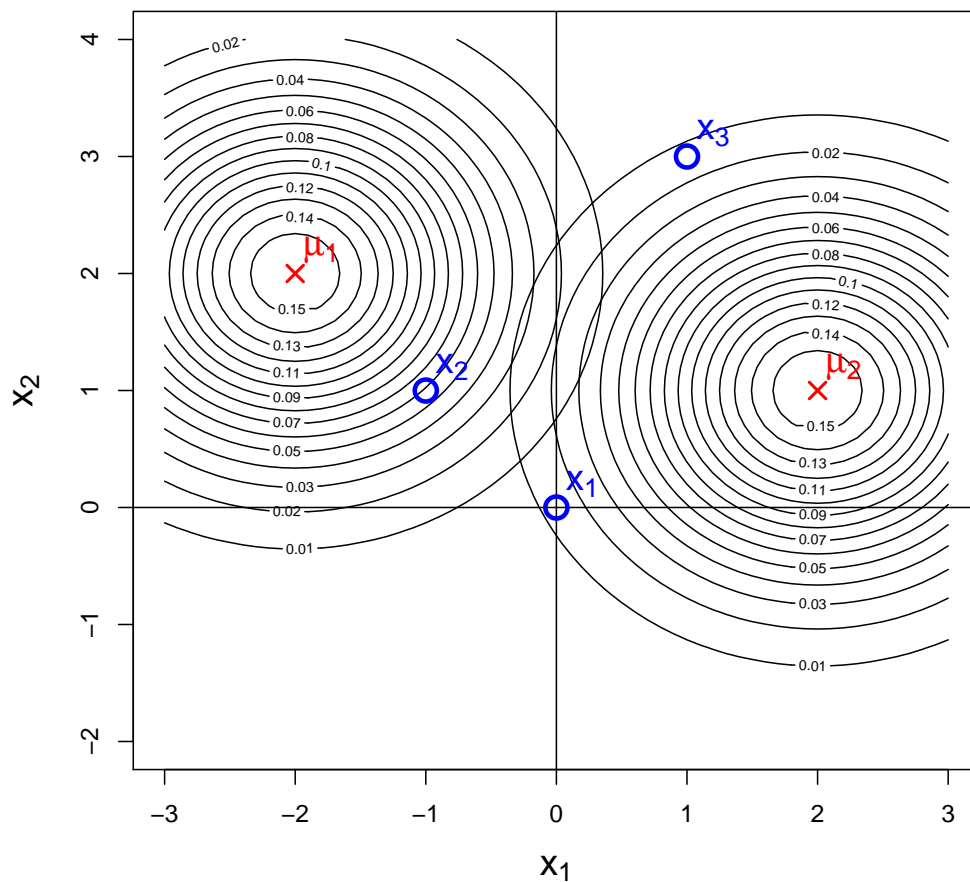
Somit:

$$\begin{aligned} \hat{p}(1) \cdot f(\mathbf{x}|\hat{\mu}_1, \hat{\Sigma}_1) &\leq \hat{p}(2) \cdot f(\mathbf{x}|\hat{\mu}_2, \hat{\Sigma}_2) \\ \hat{p}^{(1)} &\stackrel{\hat{p}^{(2)}}{\iff} f(\mathbf{x}|\hat{\mu}_1, \hat{\Sigma}_1) \leq f(\mathbf{x}|\hat{\mu}_2, \hat{\Sigma}_2) \\ &\iff ||\mathbf{x} - \hat{\mu}_1||^2 \geq ||\mathbf{x} - \hat{\mu}_2||^2 \end{aligned}$$

Setze  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2$  ein: z.B.:  $\left\| \begin{pmatrix} 0 \\ 0 \end{pmatrix} - \begin{pmatrix} -2 \\ 2 \end{pmatrix} \right\|^2 = 2^2 + (-2)^2 = 8.$

Distanzen				
	$\boldsymbol{\mu}_1$		$\boldsymbol{\mu}_2$	$\Rightarrow$ Cluster
$\mathbf{x}_1$	8	>	5	2
$\mathbf{x}_2$	2	<	9	1
$\mathbf{x}_3$	10	>	5	2

Abbildung des Modells *spherical, equal volume*:



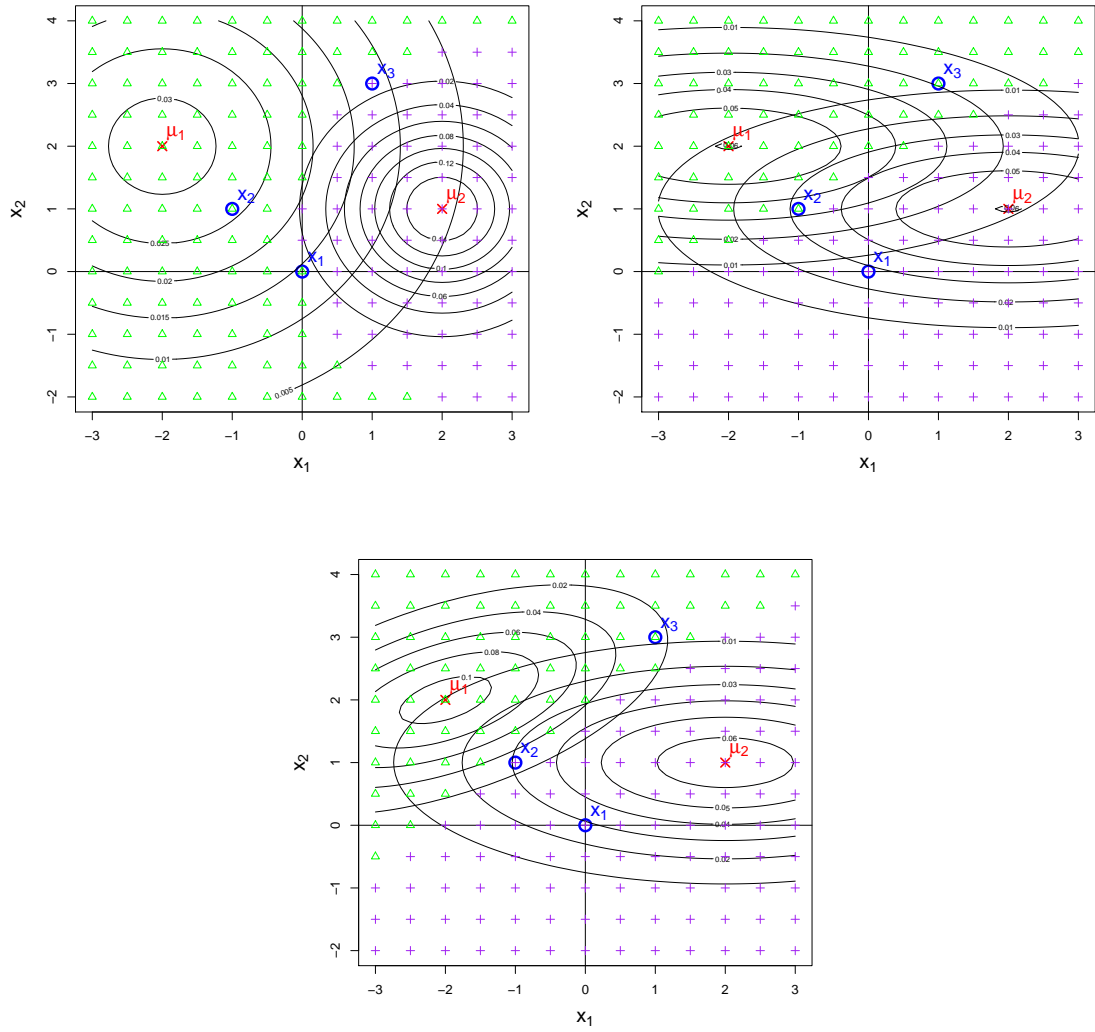
c) Mögliche Matrizen:

(i)  $\hat{\boldsymbol{\Sigma}}_1 = \begin{pmatrix} 5 & 0 \\ 0 & 5 \end{pmatrix}$  und  $\hat{\boldsymbol{\Sigma}}_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

(ii)  $\hat{\boldsymbol{\Sigma}}_1 = \begin{pmatrix} 7 & 0 \\ 0 & 1 \end{pmatrix}$  und  $\hat{\boldsymbol{\Sigma}}_2 = \begin{pmatrix} 7 & 0 \\ 0 & 1 \end{pmatrix}$

(iii)  $\hat{\boldsymbol{\Sigma}}_1 = \begin{pmatrix} 3 & 0.9 \\ 0.9 & 1 \end{pmatrix}$  und  $\hat{\boldsymbol{\Sigma}}_2 = \begin{pmatrix} 6 & 0 \\ 0 & 1 \end{pmatrix}$





- Grüne Dreiecke = Posteriori-Dichte des Clusters 1 ist größer
- Violette Kreuze = Posteriori-Dichte des Clusters 2 ist größer

### Aufgabe 3: Clusteranalyse in R II

Der Datensatz `geyser` aus dem R-Paket `MASS` beinhaltet für 299 Eruptionen des berühmten *Old Faithful* Geysirs im Yellowstone Nationalpark die Wartezeit seit der vorangegangenen Eruption (in Minuten) sowie die Eruptionsdauer (in Minuten). Im Folgenden soll dieser Datensatz mit Hilfe eines Mischmodellansatzes untersucht werden. Für die Verteilung in den Klassen wird eine bivariate Normalverteilung angenommen.

- a) Skizzieren Sie kurz die verwendeten Modellannahmen des hier betrachteten Mischmodellansatzes. Wie kann das Mischmodell geschätzt werden?
- b) Plotten Sie die Daten und bestimmen Sie visuell eine geeignete Anzahl an Klassen für den Mischmodellansatz.
- c) Führen Sie das modellbasierte Clustering durch:
  - i) Clustern Sie die Geysir-Daten mit Hilfe der Funktion `Mclust()` aus dem Package `mclust`. Verwenden Sie unterschiedliche Annahmen (z.B. „EII“, „VVI“, „VVV“) für die Kovarianzstruktur in den Klassen. Vergleichen Sie die sich daraus ergebenden Modelle und evaluieren Sie die Güte der finalen Partitionierung.
  - ii) Nutzen Sie die `Mclust()` Funktion, um unter allen möglichen modellbasierten Clusterverfahren für die aktuellsten Daten die bzgl. dem BIC optimalsten Annahmen und Anzahl der Klassen zu bestimmen.
- d) Möglichkeiten zur Evaluierung der Güte einzelner Clusterungen sowie zum Vergleich verschiedener Clustering-Ergebnisse:
  - i) Schränken Sie den `geyser`-Datensatz für eine übersichtlichere Darstellung auf die ersten 50 Zeilen ein und führen Sie ein k-Means Clustering mit 4 Clustern auf dem reduzierten Datensatz durch.
  - ii) Evaluieren Sie die Ergebnisse des k-Means Clusterings auf Basis verschiedener graphischer Darstellungen:
    - Silhouette-Plot (Funktion `factoextra::fviz_silhouette`)
    - Stripes-Plot (`flexclust::stripes`)
    - Heatmap der Distanzmatrix (`lattice::levelplot` oder `factoextra::fviz_dist`)
  - iii) Nutzen Sie den Rand-Index (Funktion `mclust::adjustedRandIndex`), um das k-Means Clustering mit dem in c) ii) erhaltenen BIC-optimalen modellbasierten Clustering zu vergleichen. Schätzen Sie hierfür das modellbasierte Clustering ebenfalls nur auf den ersten 50 Zeilen der Daten.

### Lösung:

- a)
  - Beobachtungen  $\mathbf{x}_1, \dots, \mathbf{x}_n$  mit  $\mathbf{x}_i \in \mathbb{R}^p$ , hier:  $p = 2$
  - Item-Werte bekannt, aber unbekannte Gruppenzugehörigkeit  $r_1, \dots, r_n$
  - bei gegebener Gruppenzugehörigkeit ist  $\mathbf{x}_i$  normalverteilt:

- $\mathbf{x}_i|r \sim \mathcal{N}(\boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r), r \in \{1, \dots, g\}$
- $f_r(\mathbf{x}_i) = f(\mathbf{x}_i|r) = f(\mathbf{x}_i|\boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r)$

- Priori-Wahrscheinlichkeit der Gruppenzugehörigkeit:

$$p(r), r \in \{1, \dots, g\}$$

- Annahme: Mischverteilung:

$$f(\mathbf{x}) = \sum_{r=1}^g p(r) f(\mathbf{x}|r)$$

- Posterior-Wahrscheinlichkeit für die Gruppenzugehörigkeit:

$$\hat{p}(r|\mathbf{x}_i) = \frac{\hat{p}(r)\hat{f}(\mathbf{x}_i|r)}{\hat{f}(\mathbf{x}_i)} =: \hat{p}_{ir} \quad (1)$$

- Gruppeneinteilung über marginale, geschätzte Posteriori-Wahrscheinlichkeit:

$$\mathcal{C}_r = \{\mathbf{x}_i | \hat{p}_{ir} \geq \hat{p}_{is}, r \neq s\}, r \in \{1, \dots, g\}$$

- Schätzung mit *EM-Algorithmus*, Iteration bis zur Konvergenz:

(1) E-Schritt:

- \* Gegeben  $\hat{p}(r), \hat{f}(\mathbf{x}|r), \hat{f}(\mathbf{x})$
- \* Berechne  $\hat{p}_{ir}$  gemäß (??)

(2) M-Schritt:

- \* Gegeben  $\hat{p}_{ir}$ , update  $\hat{p}(r) = \frac{1}{n} \sum_{i=1}^n \hat{p}_{ir}$
- \*  $(\hat{\boldsymbol{\mu}}_r, \hat{\boldsymbol{\Sigma}}_r) = \arg \max_{\boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r} \sum_{i=1}^n \hat{p}_{ir} \cdot \log(f_r(\mathbf{x}_i|\boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r)), r \in \{1, \dots, g\}$  (gewichteter ML-Schätzer)

- häufig zwecks Parameterökonomie: weniger flexible Annahmen für die Kovarianzen  $\boldsymbol{\Sigma}_r$ , z.B.

- \*  $\boldsymbol{\Sigma}_r = \sigma^2 \cdot I, r \in \{1, \dots, g\}$
- \*  $\boldsymbol{\Sigma}_r = \boldsymbol{\Sigma}$
- \*  $\boldsymbol{\Sigma}_r = \sigma_r^2 \cdot I, r \in \{1, \dots, g\}$

b) R-Code

c) R-Code

d) R-Code