

MASCHINELLES LERNEN & DATAMINING

Vorlesung im Wintersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 25. August 2017

Teil V

Gruppierung von Objekten

Überwachungsszenarien

Etikettierung der Lerndatenobjekte nach Klassenzugehörigkeit ?

Überwachtes Lernen

Der Lehrer stellt Zielwert **aller** Lernobjekte bereit.

{ Klassifikation
Vorhersage }

Halbüberwachtes Lernen

Der Lehrer verrät Zielwert **weniger** Lernobjekte.

{ Bootstrap
Transduktion }

Reinforcement Lernen

Der Lehrer übt **Erfolgskontrolle** („*feed-back*“).

{ Spielstrategie
Aktionsplanung }

Unüberwachtes Lernen

Der Lehrer stellt **keinerlei** Zielwerte bereit.

{ Gruppierung
Assoziation }

Gruppierung a.k.a. Clusteranalyse

Partitionierung der Datenobjekte in Ballungs- oder Häufungsgebiete

Objektrepräsentation

Vektorraum · Attribute · Metrik

Zielgröße

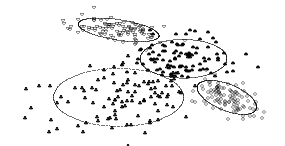
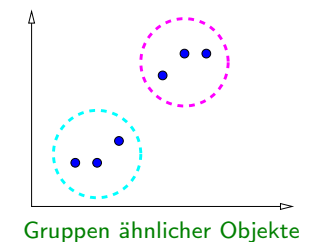
global · lokal · ad hoc

Zerlegungsstrategie

- top-down
- bottom-up
- Austausch (K-means, EM)
- split & merge

Gruppenrepräsentation

Mengen · Prototypen · Verteilungen
Formeln · Regeln



Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K-means

Mischungsidentifikation

Relationale Gruppierung

Konzeptuelle Gruppierung

Spektrale Gruppierung

Clustergütemaße

Zusammenfassung

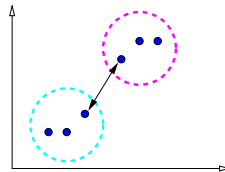
Mengendistanzfunktionen

$$d : \mathfrak{P}\Omega \times \mathfrak{P}\Omega \rightarrow \mathbb{R}_0^+$$

Single-Linkage

Kürzeste Brücke zwischen zwei Gruppen

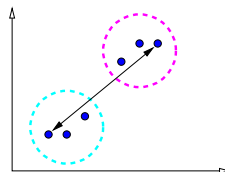
$$d_{\text{SL}}(A, B) \stackrel{\text{def}}{=} \min_{x \in A} \min_{y \in B} d(x, y)$$



Complete-Linkage

Durchmesser nach Vereinigung zweier Gruppen

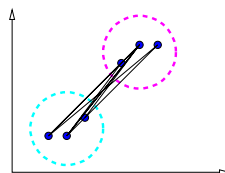
$$d_{\text{CL}}(A, B) \stackrel{\text{def}}{=} \max_{x \in A} \max_{y \in B} d(x, y)$$



Average-Linkage

Mittlerer bipartiter Punkteabstand

$$d_{\text{AL}}(A, B) \stackrel{\text{def}}{=} \frac{1}{|A| \cdot |B|} \sum_{x \in A} \sum_{y \in B} d(x, y)$$



Agglomerative Gruppierung

Generischer Bottom-up-Algorithmus

(Algorithmus)

Gegeben sind die Datenobjekte $\mathbf{x}_1, \dots, \mathbf{x}_T \in \Omega$

1 INITIALISIERUNG

Starte mit $K = T$ Gruppen $\omega_t = \{\mathbf{x}_t\}$, $t = 1..T$.

2 DISTANZBERECHNUNG

Berechne für alle $1 \leq \kappa < \lambda \leq K$:

$$D_{\kappa\lambda} \stackrel{\text{def}}{=} d(\omega_{\kappa}, \omega_{\lambda})$$

3 VEREINIGUNG

Vereinige die beiden Gruppen ω_{κ^*} , ω_{λ^*} mit

$$(\kappa^*, \lambda^*) = \underset{\lambda, \kappa}{\operatorname{argmin}} D_{\kappa\lambda}$$

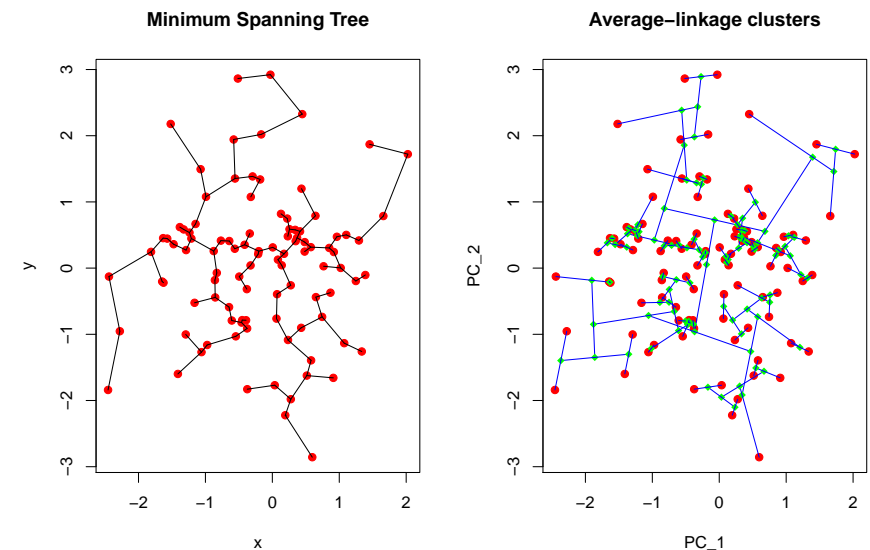
4 TERMINIERUNG

Wenn $K = 1$ dann **ENDE**, sonst \rightsquigarrow 2.

(zumFolienA)

Kettenbildung und Lassoefekt

Beispiel mit $T = 100$ Objekten im \mathbb{R}^2



Mengendistanzfunktionen

Welche ist die beste ?

Single-Linkage

„Ketteneffekt“

- erzeugt minimalen Spannbaum
- ⊕ schwach monoton
inv. monot. d -Transf.

Complete-Linkage

„Lassoefekt“

- extrem anfällig gegen Ausreißer
- ⊕ schwach monoton
inv. monot. d -Transf.

Average-Linkage

weder Ketten- noch Lassoefekt

- bevorzugt sphärische Ballungsgebiete
- ⊕ schwach monoton
 Δ -invariant

Getrimmte Distanzen

Diese Effekte lassen sich abmildern, wenn in der Distanzformel jeweils $q > 1$ kleinste bzw. größte Distanzen eliminiert werden, wodurch die Einflußdramatik eventueller Ausreißer eingedämmt wird.

Dreiecksungleichung und Monotonie

Satz (Lance & Williams, 1967)

Es sei eine rekursive Form der Mengendistanzfunktion vorausgesetzt.

1. Gilt $\alpha_1 + \alpha_2 \geq 1$, $\beta \geq 0$ und $\gamma = 0$ und gilt die Dreiecksungleichung für alle Gruppendifferenzen, so gilt sie auch noch nach der $d(\cdot, \cdot)$ -optimalen Vereinigung:

$$d(A_1 \uplus A_2, B) + d(A_1 \uplus A_2, C) \geq d(B, C)$$

2. Gilt $\alpha_1 + \alpha_2 + \beta \geq 1$ und $\gamma = 0$, so steigen die Gruppendifferenzen schwach monoton an:

$$d(A_1 \uplus A_2, B) \geq d(A_1, A_2)$$

3. Erfüllt das Agglomerationsverfahren sogar die strenge Monotonie, so gelten für jede intermediäre Gruppenstruktur (mit Gruppe A) die **ultrametrischen** Ungleichungen:

$$\forall \mathbf{x}, \mathbf{y} \in A, \forall \mathbf{z} \notin A: \quad d(\mathbf{x}, \mathbf{y}) < d(\mathbf{x}, \mathbf{z})$$

Lance-Williams-Rekursion

Effizient auswertbare Mengendistanzfunktionen

Definition

Gehorcht eine Distanzfunktion $d : \mathfrak{P}\Omega \times \mathfrak{P}\Omega \rightarrow \mathbb{R}$ in eindeutiger Weise dem Schema

$$\begin{aligned} d(\{\mathbf{x}\}, \{\mathbf{y}\}) &= d(\mathbf{x}, \mathbf{y}) \\ d(A_1 \uplus A_2, B) &= \alpha_1 \cdot d(A_1, B) + \alpha_2 \cdot d(A_2, B) + \beta \cdot d(A_1, A_2) \\ &\quad + \gamma \cdot |d(A_1, B) - d(A_2, B)| \end{aligned}$$

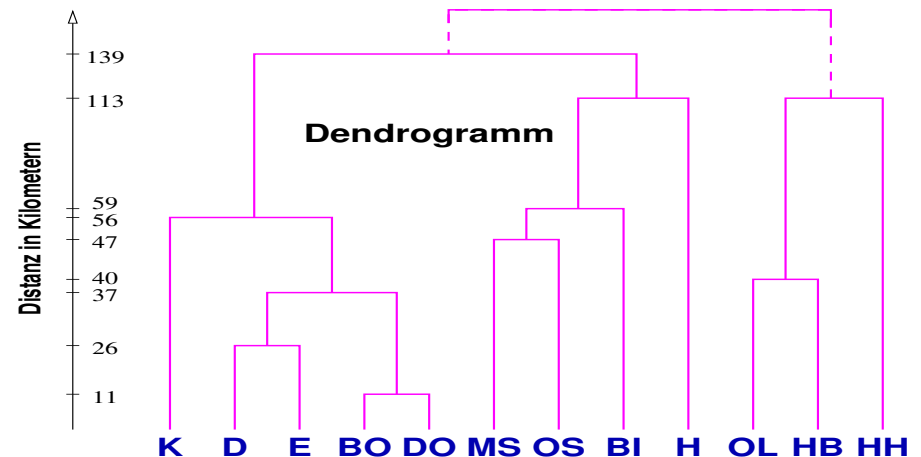
so heißt diese Vorschrift **Lance-Williams-Rekursion** mit den reellwertigen Parametern $\alpha_1 \geq 0$, $\alpha_2 \geq 0$, β und γ .

Bemerkung

Die drei X-Linkage-Funktionen besitzen alle die Lance-Williams-Gestalt:

1. Single-Linkage: $\alpha_1 = \alpha_2 = \frac{1}{2}$, $\beta = 0$, $\gamma = -\frac{1}{2}$
2. Complete-Linkage: $\alpha_1 = \alpha_2 = \frac{1}{2}$, $\beta = 0$, $\gamma = +\frac{1}{2}$
3. Average-Linkage: $\alpha_1 = \frac{|A_1|}{|A_1| + |A_2|}$, $\alpha_2 = \frac{|A_2|}{|A_1| + |A_2|}$, $\beta = 0$, $\gamma = 0$

Beispiel — Dendrogramm für Städtedistanzen



Strenge Monotonie

Je später zwei Gruppen im agglomerativen Clusteralgorithmus vereinigt werden, desto größer ist ihre Mengendistanz.

(Nichtmonotonie \rightsquigarrow Inversionen des Dendrogramms)

Weitere Distanzfunktionen

Simple-Average

Keine globale Semantik, aber schwach monoton und Δ -invariant:

$$d_{SA}(A_1 \uplus A_2, B) \stackrel{\text{def}}{=} \frac{1}{2} \cdot d_{SA}(A_1, B) + \frac{1}{2} \cdot d_{SA}(A_2, B)$$

Lance-Williams-Parameter: $\alpha_1 = \alpha_2 = \frac{1}{2}$, $\beta = 0$, $\gamma = 0$

Zentroid-Verfahren

Für numerische Attribute; weder schwach monoton noch Δ -invariant:

$$d_{ZEN}(A, B) \stackrel{\text{def}}{=} \|\mu(A) - \mu(B)\|^2$$

Lance-Williams-Parameter: $\alpha_1 = \frac{|A_1|}{|A_1| + |A_2|}$, $\alpha_2 = \frac{|A_2|}{|A_1| + |A_2|}$, $\beta = -\alpha_1 \alpha_2$, $\gamma = 0$

Median/Gower-Verfahren

Wie Zentroid; ignoriert aber die relativen Größen vereinigter Gruppen:

$$\alpha_1 = \alpha_2 = \frac{1}{2}, \quad \beta = -\frac{1}{4}, \quad \gamma = 0$$

Divisive Gruppierung

Generischer Top-down-Algorithmus

(Algorithmus)

Gegeben sind die Datenobjekte $\mathbf{x}_1, \dots, \mathbf{x}_T \in \Omega$.

1 INITIALISIERUNG

Starte mit $K = 1$ Gruppe(n) $\omega_1 = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$.

2 HETEROGENITÄTSKRITERIUM

Berechne für alle $1 \leq \kappa \leq K$ die Gruppenheterogenität

$$H_\kappa \stackrel{\text{def}}{=} d(\omega_\kappa).$$

3 AUFSPALTUNG

Zerlege diejenige Gruppe ω_{κ^*} mit

$$\kappa^* = \underset{\kappa}{\operatorname{argmax}} H_\kappa$$

in zwei disjunkte Teilgruppen (z.B. via Austauschverfahren).

4 TERMINIERUNG

Wenn $K = T$, dann **Ende**, sonst \rightsquigarrow 2.

(zumf3hlog(A))

Ward-Verfahren

Ähneln der Zentroiddistanz · Garantiert aber Distanzmonotonie

Ward-Zielgröße

Das globale Clusterverzerrungsmaß

$$\varepsilon_{\text{WARD}}(\{\omega_1, \dots, \omega_K\}) \stackrel{\text{def}}{=} \sum_{\lambda=1}^K \sum_{\mathbf{x} \in \omega_\lambda} \|\mathbf{x} - \mu_\lambda\|^2, \quad \mu_\lambda = \mu(\omega_\lambda)$$

führt auf den **Heterogenitätszuwachs**

$$d_{\text{WARD}}(A, B) = \varepsilon' - \varepsilon = \frac{|A| \cdot |B|}{|A| + |B|} \cdot \|\mu(A) - \mu(B)\|^2$$

bei Vereinigung der Gruppen $A = \omega_\kappa$ und $B = \omega_\lambda$ und diese Formel wiederum auf eine Lance-Williams-Darstellung:

$$d_{\text{WARD}}(A_1 + A_2, B) = \frac{(|A_1| + |B|) \cdot d(A_1, B) + (|A_2| + |B|) \cdot d(A_2, B) - |B| \cdot d(A_1, A_2)}{|A_1| + |A_2| + |B|}$$

Heterogenitätskriterien

Gruppendurchmesser

$$d_{\text{DIAM}}(\omega) \stackrel{\text{def}}{=} \max_{\mathbf{x}, \mathbf{y} \in \omega} d(\mathbf{x}, \mathbf{y})$$

Mittlere Innergruppenspanne

$$d_{\text{AD}}(\omega) \stackrel{\text{def}}{=} \frac{1}{|\omega|^2 - |\omega|} \sum_{\mathbf{x}, \mathbf{y} \in \omega} d(\mathbf{x}, \mathbf{y})$$

Empirische Gruppenvarianz

$$d_{\text{VAR}}(\omega) \stackrel{\text{def}}{=} \frac{1}{|\omega|} \sum_{\mathbf{x} \in \omega} \|\mathbf{x} - \mu(\omega)\|^2 = \text{spur}(\mathbf{S}(\omega))$$

Gaußäquivalente Entropie

$$d_{\text{CE}}(\omega) \stackrel{\text{def}}{=} -\frac{2}{|\omega|} \cdot \log \mathcal{N}(\omega \mid \mu(\omega), \mathbf{S}(\omega)) = \text{const} + \log \det \mathbf{S}(\omega)$$

Hierarchische Gruppierung

Divisive Gruppierung $\hat{=}$ Top-down-Induktion

Kontrollflußregelung durch Heterogenitätsmaß; $O(T \cdot n_{\text{split}})$
Polythetische Verzweigungsfragen · extensionale Zerlegung

Blinde Gruppierung

Keine Heterogenitätsprüfung
Balancierte Aufspaltung in 2^b Gruppen

Agglomerative Gruppierung $\hat{=}$ Bottom-up-Iteration

Gierige Verschmelzung mit Aufwand $O(T \cdot T^2)$

ISODATA-Algorithmus

„Split+merge“-Strategie
Pulsierende Folge von Teilungen & Verschmelzungen

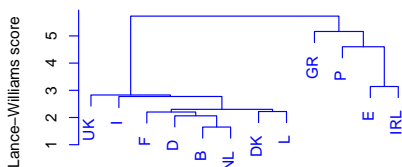
Welches ist die beste Gruppierungsstufe ?

Wähle die „richtige“ Clusteranzahl $K \in \{1, 2, \dots, T\}$

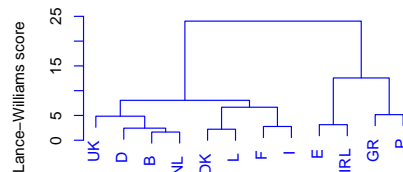
Beispiel — Agrarnationen der EU (1993)

Vergleich unterschiedlicher Lance-Williams-Distanzen

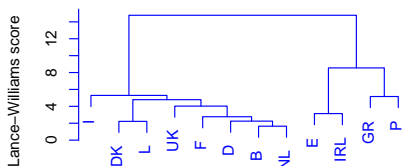
single-Distanz



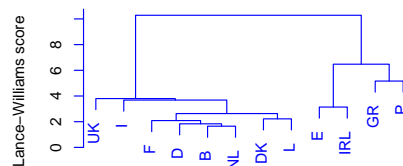
complete-Distanz



average-Distanz



centroid-Distanz



Beispiel — Agrarnationen der EU (1993)

Datensatz 'agriculture' (cluster)

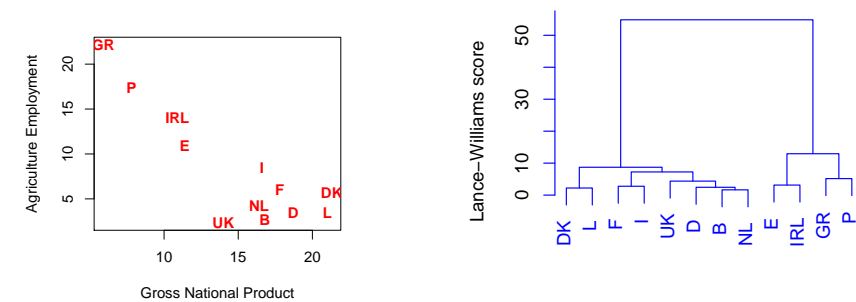
12 europäische Länder

Attribut x_1 = Bruttosozialprodukt der Hauptstadt

Attribut x_2 = Bevölkerungsanteil (%) in landwirtschaftlicher Anstellung

	B	DK	D	GR	E	F	IRL	I	L	NL	P	UK
x_1	16.8	21.3	18.7	5.9	11.4	17.8	10.9	16.6	21.0	16.4	7.8	14.0
x_2	2.7	5.7	3.5	22.2	10.9	6.0	14.0	8.5	3.5	4.3	17.4	2.3

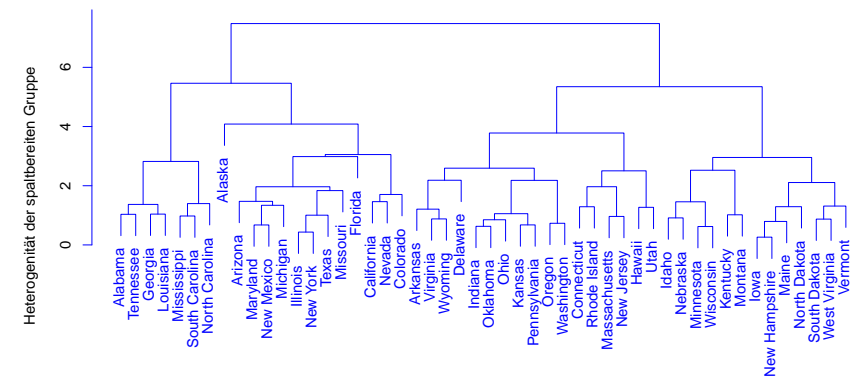
Ward-Distanz



Beispiel — Verbrechenstatistik

Divisive Gruppierung mit 'diana'/R

diana (USArrests, metric='euclidean', stand=TRUE)



Datensatz 'USArrests' (datasets)

50 Objekte: Kriminalstatistiken aller US-Bundesstaaten (1973)

3 Attribute: „Murder“, „Assault“, „Rape“ (Anzahl je 10^5 Einwohner) und

1 Attribut: „UrbanPop“ (Prozentsatz Stadtbevölkerung)

Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K-means

Mischungsidentifikation

Relationale Gruppierung

Konzeptuelle Gruppierung

Spektrale Gruppierung

Clustergütemaße

Zusammenfassung

Permutationsverfahren

Gieriges Suchverfahren · extensional · alle Metriken

(Algorithmus)

1 INITIALISIERUNG

Wähle eine Startpartition $\omega_1, \dots, \omega_K$ mit vorgegebenen $|\omega_\kappa| = T_\kappa$.

2 VERZERRUNGSDIFFERENZEN

Berechne für alle $\mathbf{x} \in \omega_\kappa$ und $\mathbf{y} \in \omega_\lambda$ mit $\kappa \neq \lambda$

$$\Delta\varepsilon(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \varepsilon(\{\dots, \underbrace{\omega'_\kappa, \omega'_\lambda}_{\mathbf{x} \leftrightarrow \mathbf{y}}, \dots\}) - \varepsilon(\{\omega_1, \dots, \omega_K\}) .$$

3 VERTAUSCHUNG

Vertausche innerhalb der aktuellen Partition das Datenvektorpaar

$$(\mathbf{x}^*, \mathbf{y}^*) = \operatorname{argmin} \Delta\varepsilon(\mathbf{x}, \mathbf{y}) .$$

4 TERMINIERUNG

Wenn $\varepsilon \leq \theta$ dann **Ende** sonst \rightsquigarrow 2.

(zumf3h0g1A)

Scharfe Gruppierung

bei vorgegebener Gruppenanzahl $K \in \mathbb{N}$

GESUCHT

ist eine K -Partition des Datensatzes $\omega \subset \Omega$.

- **extensional:** Teilmengensystem $\omega_1 \uplus \omega_2 \uplus \dots \uplus \omega_K = \omega$
- **intensional:** Gruppenprototypen $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K \in \Omega$

Verzerrung einer Gruppe

hinsichtlich einer Objektraummetrik $d : \Omega \times \Omega \rightarrow \mathbb{R}$:

$$\varepsilon(\omega_\kappa) \stackrel{\text{def}}{=} \sum_{\mathbf{x} \in \omega_\kappa} d(\mathbf{x}, \boldsymbol{\mu}(\omega_\kappa)) , \quad \kappa = 1, \dots, K$$

Verzerrung einer Partition

$$\varepsilon(\{\omega_1, \dots, \omega_K\}) \stackrel{\text{def}}{=} \sum_{\kappa=1}^K \varepsilon(\omega_\kappa)$$

⇒ Kombinatorische Optimierungsaufgabe

Intensionale Gruppierung

Gruppierung mit Prototypen — „Vektorquantisierung“

Lemma

Es sei $\omega_1, \dots, \omega_K$ eine Gruppierung der Elemente $\mathbf{x}_1, \dots, \mathbf{x}_T$ des metrischen Raumes (Ω, d) , welche die globale Verzerrung minimiert. Dann gibt es **Gruppenprototypen** $\mathbf{z}_1, \dots, \mathbf{z}_K \in \Omega$ mit

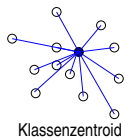
1. Jedes \mathbf{z}_κ ist Zentroid seiner Gruppe ω_κ :

$$\mathbf{z}_\kappa = \operatorname{argmin}_{\mathbf{y} \in \Omega} \sum_{\mathbf{x} \in \omega_\kappa} d(\mathbf{x}, \mathbf{y})$$

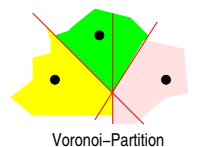
2. Jeder Datenvektor \mathbf{x}_t , $t = 1, \dots, T$ gehört zu der Gruppe des nächstliegenden Prototypen:

$$\mathbf{x}_t \in \omega_\kappa \iff d(\mathbf{x}_t, \mathbf{z}_\kappa) = \min_{\lambda} d(\mathbf{x}_t, \mathbf{z}_\lambda)$$

Für die euklidische Distanz gilt natürlich $\mathbf{z}_\kappa = \boldsymbol{\mu}(\omega_\kappa)$ für alle $\kappa = 1, \dots, K$.



Klassenzentroid



Voronoi-Partition

Stapelweiser K-means-Algorithmus

Lloyd 1957 · Forgy 1965

(Algorithmus)

1 INITIALISIERUNG

Wähle eine zufällige Startpartition

$$\omega_1 \uplus \omega_2 \uplus \dots \uplus \omega_K = \omega$$

2 REPRÄSENTATION

Berechne alle neuen Prototypen

$$\mathbf{z}_\kappa = \mu_{\text{ZEN}}(\omega_\kappa) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{y} \in \Omega} \sum_{\mathbf{x} \in \omega_\kappa} d(\mathbf{x}, \mathbf{y})$$

3 REKLASSIFIKATION

Berechne alle neuen Gruppen

$$\omega_\kappa = \left\{ \mathbf{x}_t \in \omega \mid \operatorname{argmin}_\lambda d(\mathbf{x}_t, \mathbf{z}_\lambda) = \kappa \right\}$$

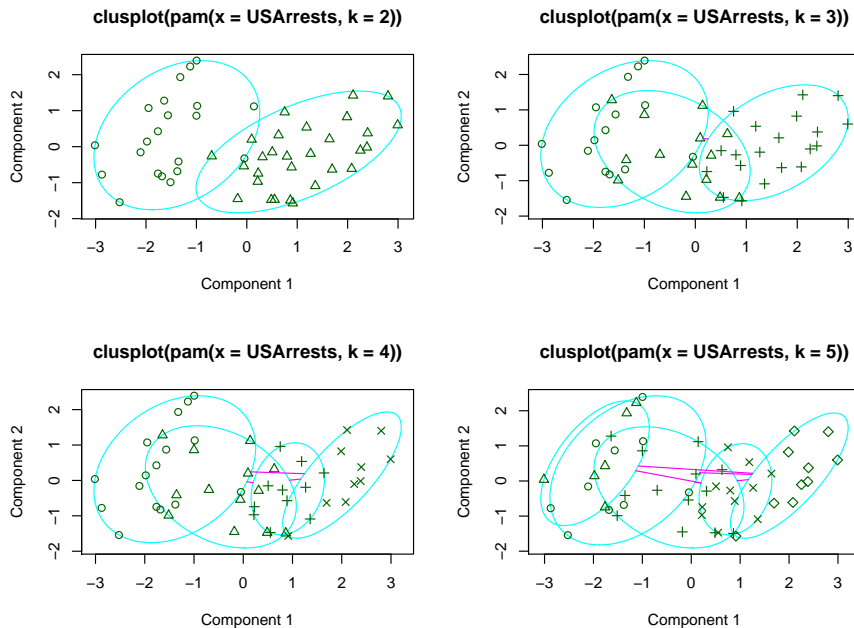
4 TERMINIERUNG

Wenn $\varepsilon(\{\omega_1, \dots, \omega_K\}) \leq \theta$ dann **Ende** sonst \rightsquigarrow 2.

(sumftihogIA)

Beispiel — 'USArrests'-Datensatz

K-medoids-Algorithmus minimiert $\|\cdot\|^1$ -Summe · robuster als K-means



Inkrementeller K-means-Algorithmus

MacQueen 1967

(Algorithmus)

1 INITIALISIERUNG

Wähle zufällige Startprototypen $\{\mathbf{z}_1, \dots, \mathbf{z}_K\}$, setze $t \leftarrow 1$.

2 REKLASSIFIKATION

Wähle $\mathbf{y} = \mathbf{x}_{t \bmod T}$ und setze $\kappa = \operatorname{argmin}_\lambda d(\mathbf{y}, \mathbf{z}_\lambda)$.

3 REPRÄSENTATION

Verschiebe $\mathbf{z}_\kappa \leftarrow \alpha_t \cdot \mathbf{y} + (1 - \alpha_t) \cdot \mathbf{z}_\kappa$.

4 TERMINIERUNG

Wenn $\varepsilon(\cdot) \leq \theta$ dann **Ende** sonst $t \leftarrow t + 1$ und \rightsquigarrow 2.

(sumftihogIA)

Bemerkungen

1. Die Gewinnerprototypen \mathbf{z}_κ werden nach jedem Einzelschritt aktualisiert.
2. Die Datenprobe wird zyklisch oder randomisiert durchlaufen.
3. Distanz $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|^2 \Rightarrow$ Zentroid \triangleq Mittelwert.
4. Mittelungsgewichte *exponentiell* ($\alpha_t \equiv \alpha_0$) oder *kumulativ* ($\alpha_t = 1/|\omega_\kappa|$).
5. Schnellerer Abstieg — aber Oszillationsgefahr!

Unschärfe Gruppierung

GESUCHT

ist eine **Zugehörigkeitsfunktion** für den Datensatz $\omega \subset \Omega$:

$$\mathbf{u} : \begin{cases} \omega & \rightarrow [0, 1]^K \\ \mathbf{x}_t & \mapsto \{u_{\kappa,t}\}_{\kappa=1}^K \end{cases}, \quad \sum_{\kappa} u_{\kappa,t} = 1 \quad (\forall t)$$

Fuzzy K-means Zielgröße

Distanzfunktion ist (hier) der quadrierte euklidische Abstand:

$$\varepsilon(\{\mathbf{u}_\kappa\}, \{\mathbf{z}_\kappa\}) = \sum_{\kappa=1}^K \sum_{t=1}^T (u_{\kappa}(\mathbf{x}_t))^{\alpha} \cdot \|\mathbf{x}_t - \mathbf{z}_\kappa\|^2, \quad \alpha \geq 1$$

Opt. Prototypen/Zugehörigk.

Normierung \rightsquigarrow Lagrangemultiplikatoren

$$\sum_{\mathbf{x} \in \omega} \beta_{\mathbf{x}} \cdot \left(\sum_{\kappa=1}^K u_{\kappa}(\mathbf{x}) - 1 \right)$$

Spezialfälle

- $\alpha = 1$: scharfe Datenmengen
- $\alpha = 2$: unscharfe Datenmengen
- $\alpha = \infty$: identische Gruppen

Fuzzy K-means-Algorithmus

(Algorithmus)

1 INITIALISIERUNG

Wähle zufällige Startzugehörigkeiten $u_{\kappa,t} \in [0, 1]$.

2 PROTOTYPEN

Für alle $1 \leq \kappa \leq K$ berechne

$$\mathbf{z}_{\kappa} = \frac{\sum_{\mathbf{x} \in \omega} (u_{\kappa}(\mathbf{x}))^{\alpha} \cdot \mathbf{x}}{\sum_{\mathbf{x} \in \omega} (u_{\kappa}(\mathbf{x}))^{\alpha}}$$

3 ZUGEHÖRIGKEITEN

Für alle $1 \leq \kappa \leq K$ und $\mathbf{x} \in \omega$ berechne neue unscharfe Gruppen:

$$u_{\kappa}(\mathbf{x}) = \frac{1}{\sum_{\lambda=1}^K \left(\frac{\|\mathbf{x} - \mathbf{z}_{\kappa}\|^2}{\|\mathbf{x} - \mathbf{z}_{\lambda}\|^2} \right)^{\frac{1}{\alpha-1}}}$$

4 TERMINIERUNG

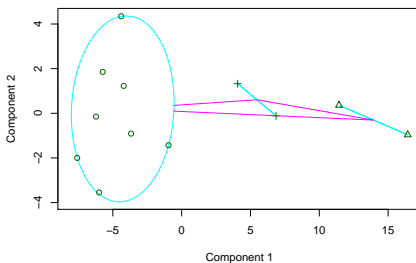
Wenn $\varepsilon \leq \theta$ dann **Ende** sonst \rightsquigarrow 2.

(zumfitting(A))

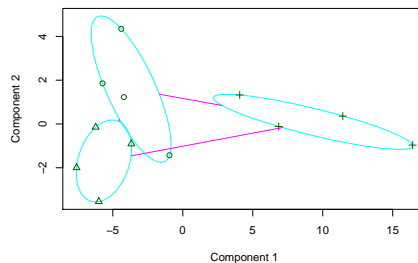
Beispiel — 'agriculture'-Datensatz

Fuzzy K-means-Algorithmus ($\alpha \in \{\sqrt{2}^i \mid i = 1, 2, 3, 4\}$)

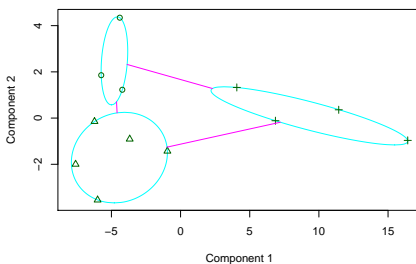
clusplot(fanny(x = agriculture, k = 3, memb.exp = 1.4))



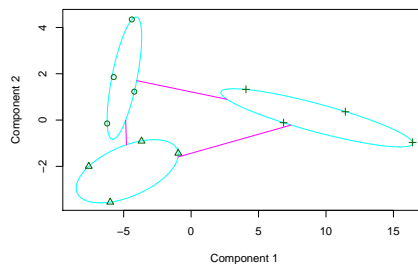
clusplot(fanny(x = agriculture, k = 3, memb.exp = 2))



clusplot(fanny(x = agriculture, k = 3, memb.exp = 2.8))



clusplot(fanny(x = agriculture, k = 3, memb.exp = 4))



Zugehörigkeitsfunktionen

Fuzzy Clustering mit steileren Abklingraten

Harmonische Zugehörigkeitsfunktion

(Standardverfahren: Fuzzy K-means)

Cauchy Zugehörigkeitsfunktion

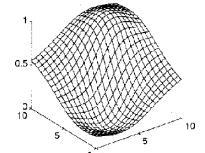
mit den Halbwertsbreiten $\eta_{\kappa} > 0$:

$$u_{\kappa}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{1 + \left(\frac{\|\mathbf{x} - \mathbf{z}_{\kappa}\|^2}{\eta_{\kappa}} \right)^{\frac{1}{\alpha-1}}}$$

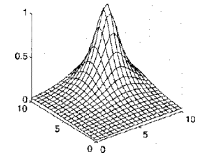
Hyperkonische Zugehörigkeitsfunktion

mit den Radien $r_{\kappa} > 0$:

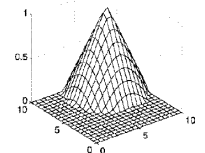
$$u_{\kappa}(\mathbf{x}) \stackrel{\text{def}}{=} \begin{cases} 1 - \|\mathbf{x} - \mathbf{z}_{\kappa}\|/r_{\kappa} & \text{falls } \|\mathbf{x} - \mathbf{z}_{\kappa}\| \leq r_{\kappa} \\ 0 & \text{sonst} \end{cases}$$



K-means



Cauchy



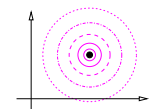
Hyperkonisch

Geometrische Clusterformen

Pendikulare Linien im \mathbb{R}^N

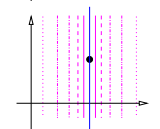
Punktförmiges Zentrum im \mathbb{R}^2

$$d^2(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|^2 = (x_1 - z_1)^2 + (x_2 - z_2)^2$$



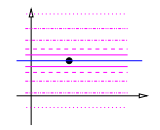
Vertikal linienförmig im \mathbb{R}^2

$$d^2(\mathbf{x}, \mathbf{z}) = (x_1 - z_1)^2$$



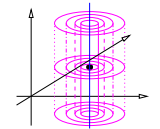
Horizontal linienförmig im \mathbb{R}^2

$$d^2(\mathbf{x}, \mathbf{z}) = (x_2 - z_2)^2$$



Vertikal linienförmig im \mathbb{R}^3

$$\begin{aligned} d^2(\mathbf{x}, \mathbf{z}) &= (x_1 - z_1)^2 + (x_2 - z_2)^2 \\ &= \|\mathbf{x} - \mathbf{z}\|^2 - (x_3 - z_3)^2 \end{aligned}$$



Perpendikulare Linienzentren

Unendlich lange, koordinatenachsenparallele Cluster

Linienförmiges Klassenzentrum

des \mathbb{R}^N in Richtung der x_n -Achse, $n \in \{1, \dots, N\}$:

$$\begin{aligned} d^2(\mathbf{x} \mid \mathbf{z}, n) &= \|\mathbf{x} - \mathbf{z}\|^2 - (x_n - z_n)^2 \\ &= \|\mathbf{x} - \mathbf{z}\|^2 - (\mathbf{e}_n^\top \cdot (\mathbf{x} - \mathbf{z}))^2 \\ &= \|\mathbf{x} - \mathbf{z}\|^2 - \|\mathbf{e}_n^\top \cdot (\mathbf{x} - \mathbf{z})\|^2 = d^2(\mathbf{x} \mid \mathbf{z}, \mathbf{e}_n) \end{aligned}$$

Vom euklidischen Abstand wird also die Norm einer Achsenprojektion subtrahiert.

Verallgemeinerung auf „schräge“ Cluster ?

Es ist naheliegend, daß dieser Zusammenhang auch für den nichtperpendikularen Fall gilt.

M-dimensionale Hyperflächenzentren

Satz (Pythagoras)

Sei $0 \leq M \leq N$. Für alle $\mathbf{x} \in \mathbb{R}^N$ berechnet sich der lotrechte Abstand zwischen \mathbf{x} und der M-dimensionalen Hyperfläche mit dem

Aufpunktvektor \mathbf{z} und den orthonormalen **Richtungsvektoren** $\mathbf{u}_1, \dots, \mathbf{u}_M \in \mathbb{R}^N$ gemäß

$$\min_{a_1, \dots, a_M} \left\| \mathbf{x} - \left(\mathbf{z} + \sum_{m=1}^M a_m \mathbf{u}_m \right) \right\|^2 = \|\mathbf{x} - \mathbf{z}\|^2 - \left\| \mathbf{U}^\top (\mathbf{x} - \mathbf{z}) \right\|^2,$$

wenn $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_M)$ ist.

Beweis.

Der Abstand für ein M-dimensionales Zentrum, das durch die orthonormalen Vektoren $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_M) \in \mathbb{R}^{N \times M}$ aufgespannt wird:

$$d^2(\mathbf{x} \mid \mathbf{z}, \mathbf{U}) = \|\mathbf{x} - \mathbf{z}\|^2 - \sum_{m=1}^M (\mathbf{u}_m^\top (\mathbf{x} - \mathbf{z}))^2 = \|\mathbf{x} - \mathbf{z}\|^2 - \left\| \mathbf{U}^\top (\mathbf{x} - \mathbf{z}) \right\|^2$$

□

Achsenrotation

Datentransformation $\phi: \mathbf{x} \mapsto \mathbf{U}^\top \mathbf{x}$, $\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{E}$

Rotationen sind distanzinvariant

$$\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} = \mathbf{x}^\top \mathbf{E} \mathbf{x} = \mathbf{x}^\top \mathbf{U} \mathbf{U}^\top \mathbf{x} = \left\| \mathbf{U}^\top \mathbf{x} \right\|^2 = \sum_{n=1}^N (\mathbf{u}_n^\top \mathbf{x})^2$$

Summendarstellung mit den Spaltenvektoren $\mathbf{u}_1, \dots, \mathbf{u}_N$ von \mathbf{U} .

Linienbezogener Abstand in \mathbf{u}_n -Richtung

$$\begin{aligned} d^2(\mathbf{x} \mid \mathbf{z}, \mathbf{u}_n) &= \|\phi \mathbf{x} - \phi \mathbf{z}\|^2 - ((\phi \mathbf{x})_n - (\phi \mathbf{z})_n)^2 \\ &= \|\phi(\mathbf{x} - \mathbf{z})\|^2 - (\mathbf{u}_n^\top \mathbf{x} - \mathbf{u}_n^\top \mathbf{z})^2 \\ &= \|\mathbf{x} - \mathbf{z}\|^2 - (\mathbf{u}_n^\top (\mathbf{x} - \mathbf{z}))^2 \end{aligned}$$

Flächenbezogener Abstand in \mathbf{u}, \mathbf{v} -Richtung

(Richtungsvektoren \mathbf{u}, \mathbf{v} normiert und senkrecht zueinander)

$$d^2(\mathbf{x} \mid \mathbf{z}, \mathbf{u}, \mathbf{v}) = \|\mathbf{x} - \mathbf{z}\|^2 - (\mathbf{u}^\top (\mathbf{x} - \mathbf{z}))^2 - (\mathbf{v}^\top (\mathbf{x} - \mathbf{z}))^2$$

Fuzzy K-Varieties

Definition

Das (unscharfe) Gruppierungsverfahren mit der Zielgröße

$$\varepsilon(\{\omega_\kappa\}) = \sum_{\kappa=1}^K \sum_{\mathbf{x} \in \omega} u_\kappa(\mathbf{x})^\alpha \cdot d^2(\mathbf{x} \mid \mathbf{z}_\kappa, \mathbf{U}_\kappa)$$

heißt **fuzzy K-varieties**-Algorithmus; im Spezialfall $M = 1$ heißt es **fuzzy K-lines**-Algorithmus.

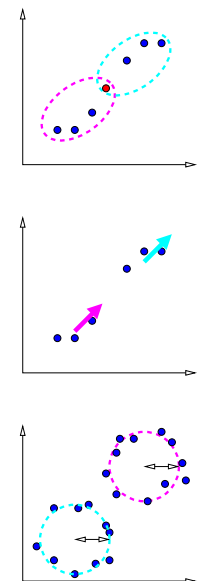
Elliptotypzentren („fuzzy K-elliptotypes“)

$$d^2(\mathbf{x} \mid \mathbf{z}, \mathbf{U}) = \|\mathbf{x} - \mathbf{z}\|^2 - \rho \cdot \left\| \mathbf{U}^\top (\mathbf{x} - \mathbf{z}) \right\|^2$$

Spezialfälle: $\rho = 0$ Punktzentrum · $\rho = 1$ Hyperflächenzentrum

Hyperkugelschalen („fuzzy K-shells“)

$$d^2(\mathbf{x} \mid \mathbf{z}, r) = (\|\mathbf{x} - \mathbf{z}\| - r)^2$$



Gradientenabstieg für Fuzzy K -Varieties

Lemma

Die Minimierung der Zielgröße mit Lagrangemultiplikatoren für die Normierungsbedingungen liefert die Bestimmungsgleichungen

$$\mathbf{z}_\kappa = \frac{\sum_{\mathbf{x} \in \omega} (u_\kappa(\mathbf{x}))^\alpha \cdot \mathbf{x}}{\sum_{\mathbf{x} \in \omega} (u_\kappa(\mathbf{x}))^\alpha}$$

für die **Aufpunktvektoren**,

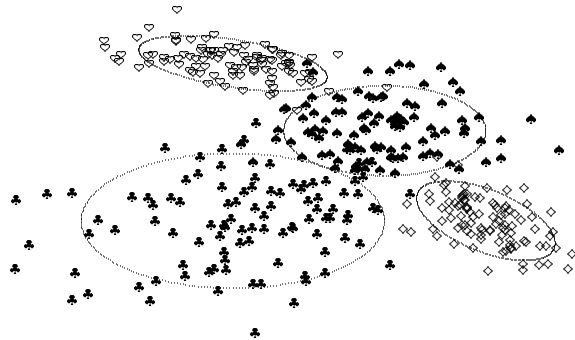
$$u_\kappa(\mathbf{x}) = 1 / \sum_{\lambda=1}^K \left(\frac{d^2(\mathbf{x} | \mathbf{z}_\kappa, \mathbf{U}_\kappa)}{d^2(\mathbf{x} | \mathbf{z}_\lambda, \mathbf{U}_\lambda)} \right)^{\frac{1}{\alpha-1}}$$

für die **Gruppenzugehörigkeiten** und für die **Gruppenkovarianzen**

$$\mathbf{S}_\kappa = \sum_{\mathbf{x} \in \omega} u_\kappa(\mathbf{x})^\alpha (\mathbf{x} - \mathbf{z}_\kappa)(\mathbf{x} - \mathbf{z}_\kappa)^\top.$$

Die m -te Spalte $\mathbf{u}_{\kappa,m}$ von \mathbf{U}_κ schließlich ergibt sich als Eigenvektor zum m -größten Eigenwert von \mathbf{S}_κ .

Identifikation von Mischverteilungen



Problem

Angenommen, obige Daten sind gemäß

$f(\mathbf{x}) = \sum_{\kappa=1}^K \pi_\kappa \cdot g(\mathbf{x} | \theta_\kappa)$ mischverteilt. Wie lauten die **bestpassenden** (ML)

Verteilungsparameter $\hat{\pi}_\kappa, \hat{\theta}_\kappa, \kappa = 1, \dots, K$ des Modells?

Lösung

Im Normalverteilungsfall $g(\mathbf{x} | \theta_\kappa) = \mathcal{N}(\mathbf{x} | \mu_\kappa, \mathbf{S}_\kappa)$ existiert eine **asymptotisch** eindeutige Lösung sowie ein **lokales** Optimierungsverfahren (EM).

Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K -means

Mischungsidentifikation

Relationale Gruppierung

Konzeptuelle Gruppierung

Spektrale Gruppierung

Clustergütemaße

Zusammenfassung

EM-Algorithmus

zur Identifikation gaußscher Mischverteilungen

(Algorithmus)

1 INITIALISIERUNG

Wähle zufällige Startparameter $(\pi_\kappa, \mu_\kappa, \mathbf{S}_\kappa), \kappa = 1, \dots, K$.

2 ERWARTUNGSWERTE

Berechne für $\kappa = 1..K$ und $t = 1..T$ die a posteriori Wahrscheinlichkeiten

$$\gamma_{\kappa,t} \stackrel{\text{def}}{=} P(\Omega_\kappa | \mathbf{x}_t) = \frac{P(\Omega_\kappa) \cdot P(\mathbf{x}_t | \Omega_\kappa)}{P(\mathbf{x}_t)} \propto \pi_\kappa \cdot \mathcal{N}(\mathbf{x}_t | \mu_\kappa, \mathbf{S}_\kappa)$$

3 MAXIMIERUNG

$$\pi_\kappa \leftarrow \frac{\sum_t \gamma_{\kappa,t}}{\sum_\lambda \sum_t \gamma_{\lambda,t}}, \quad \mu_\kappa \leftarrow \frac{\sum_t \gamma_{\kappa,t} \mathbf{x}_t}{\sum_t \gamma_{\kappa,t}}, \quad \mathbf{S}_\kappa \leftarrow \frac{\sum_t \gamma_{\kappa,t} \mathbf{x}_t \mathbf{x}_t^\top}{\sum_t \gamma_{\kappa,t}} - \mu_\kappa \mu_\kappa^\top$$

4 TERMINIERUNG

Wenn die ML-Zielgröße $\ell(\dots)$ stagniert dann **Ende** sonst \rightsquigarrow **2**.

(summiertlog(A))

Konvergenzeigenschaften

des EM-Algorithmus für Gaußsche Mischverteilungsmodelle (GMM)

1. Schwache Monotonie

Verfahren erreicht stationären Punkt

$$\ell(\theta_0) \leq \ell(\theta_1) \leq \ell(\theta_2) \leq \ell(\theta_3) \leq \dots \leq \ell(\theta_j) \leq \dots \leq \dots$$

2. Beschränktheit

pathologische Aufgabenstellung („ill-posed problem“)

$$\mathcal{N}(\mu_\lambda, \Sigma_\lambda) = \mathcal{N}(\mathbf{x}_t, \mathbf{0})$$

3. Lokale Maxima

viele relative Maxima mit $\ell(\theta) < \infty$ und großem Einzugsbereich

4. Zyklischer Iterationsverlauf

Kraterrandphänomen

$$\theta_1 \neq \theta_2 \neq \dots \neq \theta_m \quad \text{mit} \quad \ell(\theta_1) = \ell(\theta_2) = \dots = \ell(\theta_m)$$

Problematik des Rangdefizits

$$\text{rg}(\mathbf{S}_\kappa) < N \Rightarrow \det(\mathbf{S}_\kappa) = 0 \Rightarrow \mathbf{S}_\kappa^{-1} = ?$$

Gratregularisierung

Anisotropes Aufblasen der Konzentrationsellipse („Speckschicht“)

$$\mathbf{S}^{(\delta)} \stackrel{\text{def}}{=} \mathbf{S} + \delta \cdot \mathbf{E} = \begin{pmatrix} s_{11} + \delta & s_{12} & \dots & s_{1N} \\ s_{21} & s_{22} + \delta & \dots & s_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ s_{N1} & s_{N2} & \dots & s_{NN} + \delta \end{pmatrix}$$

Fixierung & Verklebung

Alle Eigenschaften des EM-Algorithmus bleiben erhalten:

- Kovarianzmatrizen fixieren ($\forall \kappa : \mathbf{S}_\kappa \stackrel{!}{=} \mathbf{S}^*$) \rightsquigarrow keine pathologische Lösung
- Kovarianzmatrizen verkleben ($\forall \kappa, \lambda : \mathbf{S}_\kappa \stackrel{!}{=} \mathbf{S}_\lambda$) \rightsquigarrow mehr Robustheit

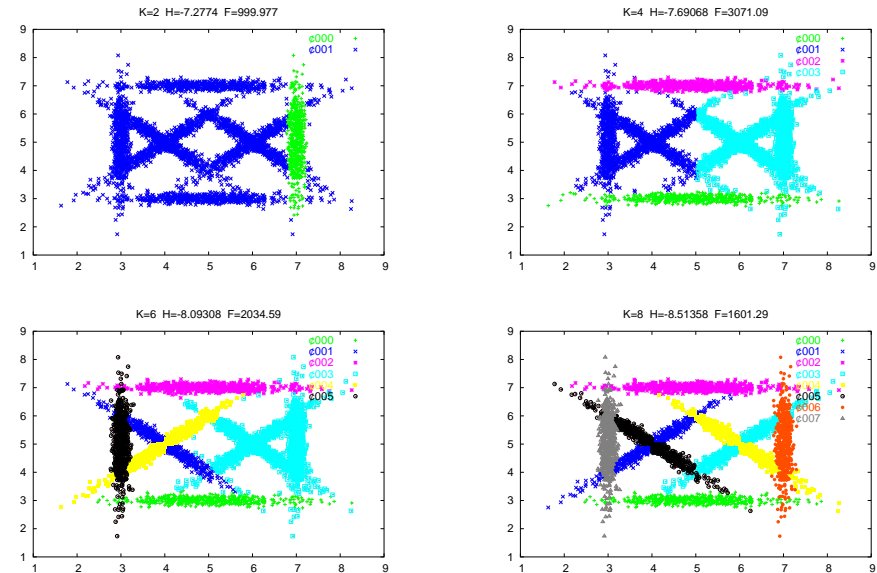
Hintergrundkomponente

Streuungsintensives Rückweisungscluster zur Ausreißerbehandlung

$$f_0(\cdot) = \mathcal{N}(\cdot | \mu(\omega), \mathbf{S}_0) \quad \text{mit} \quad \mathbf{S}_0 = \mathbf{S}(\omega) \text{ oder } \mathbf{S}_0 = \sigma_0^2 \cdot \mathbf{E}$$

Verhalten in unkritischen Fällen

Was kann EM, das K-mean nicht kann ?



Probabilistische PCA

Zerlegung des \mathbb{R}^N in systematisch und in zufällig streuende Komponenten

Normalverteilungsmodelle für rangdefizite Daten

Das homogene Faktoranalysemodell

$$\mathbb{X} = \mu + \mathbb{E} + \mathbf{A} \cdot \mathbb{V} \quad \text{mit} \quad \begin{cases} \mu \in \mathbb{R}^N \\ \mathbf{A} \in \mathbb{R}^{N \times M} \\ \mathbb{E} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \cdot \mathbf{E}_N) \\ \mathbb{V} \sim \mathcal{N}(\mathbf{0}, \mathbf{E}_M) \end{cases}$$

(Dimension $M \leq N$; PCA-Annahme identischer Störvarianzen)

- besitzt als Ladungsvektoren die bereits hinlänglich bekannten M führenden Hauptachsen der Verteilungsellipse,
- definiert aber gleichzeitig eine explizite Wahrscheinlichkeitsverteilung für die Daten.

PPCA-Schätzung bei bekanntem Modellrang M

Lemma (Tipping & Bishop 1999)

Der Zufallsvektor des homogenen FA-Modells ist gemäß $\mathbb{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{S})$ normalverteilt mit der Kovarianzmatrix

$$\mathbf{S} = \mathbf{A}\mathbf{A}^\top + \sigma^2 \cdot \mathbf{E}_N.$$

Der **ML-Schätzer** für \mathbf{S} ergibt sich durch Einsetzen der Schätzwerte

$$\begin{aligned}\hat{\mathbf{A}} &= \mathbf{U}_M \cdot (\mathbf{D}_M - \sigma^2 \cdot \mathbf{E}_M)^{1/2} \\ \hat{\sigma}^2 &= \frac{1}{N-M} \cdot \sum_{j=M+1}^N \lambda_j\end{aligned}$$

mit der $(M : N)$ -eigenzerlegten Datenkovarianzmatrix

$$\hat{\mathbf{S}}(\omega) \stackrel{!}{=} (\mathbf{U}_M, \mathbf{U}'_M) \cdot \begin{pmatrix} \mathbf{D}_M & \mathbf{0} \\ \mathbf{0}^\top & \mathbf{D}'_M \end{pmatrix} \cdot (\mathbf{U}_M, \mathbf{U}'_M)^\top, \quad \begin{pmatrix} \mathbf{D}_M & \mathbf{0} \\ \mathbf{0}^\top & \mathbf{D}'_M \end{pmatrix} = \text{diag}(\lambda_1, \dots, \lambda_N).$$

Invertierung der PPCA-Kovarianzmatrix

Kovarianzstruktur

Die ursprüngliche Darstellung zeigt eine rangverminderte Darstellung mit einem additiv aufgeprägtem Fehlergrad von σ^2 auf der Diagonalen.

$$\mathbf{S} = \mathbf{A}\mathbf{A}^\top + \sigma^2 \cdot \mathbf{E}_N$$

Die gleichwertige alternative Darstellung präsentiert eine vollständige Eigenzerlegung mit den kanonischen Eigenvektoren und -werten; nur die letzten $(N - M)$ Eigenwerte wurden gemittelt.

$$\mathbf{S} = \mathbf{U}_M \mathbf{D}_M \mathbf{U}_M^\top + \sigma^2 \cdot \mathbf{U}'_M \mathbf{U}'_M^\top$$

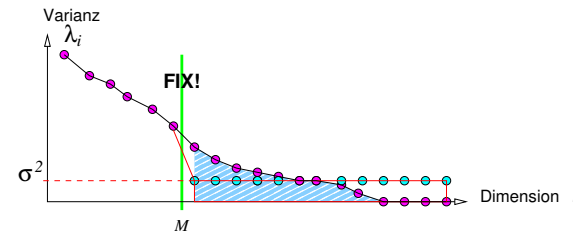
Inverse Kovarianz

Diese Inversionsformel verwendet ausschließlich die M führenden Eigenvektoren sowie die führenden reziproken Eigenwerte.

$$\mathbf{S}^{-1} = \frac{1}{\sigma^2} \cdot \left\{ \mathbf{E}_N - \mathbf{U}_M \cdot (\mathbf{E}_M - \sigma^2 \cdot \mathbf{D}_M^{-1}) \cdot \mathbf{U}_M^\top \right\}$$

Schätzung der mittleren Reststreuung unter ausschließlicher Verwendung der $M \ll N$ Hauptachsen

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{N-M} \cdot \text{spur}(\mathbf{D}'_M) \\ &= \frac{1}{N-M} \left\{ \text{spur} \begin{pmatrix} \mathbf{D}_M & \mathbf{0} \\ \mathbf{0}^\top & \mathbf{D}'_M \end{pmatrix} - \text{spur}(\mathbf{D}_M) \right\} \\ &= \frac{1}{N-M} \left\{ \text{spur}(\hat{\mathbf{S}}(\omega)) - \text{spur}(\mathbf{D}_M) \right\} \\ &= \frac{1}{N-M} \left\{ \frac{1}{T} \sum_{t=1}^T \|\mathbf{x}_t - \boldsymbol{\mu}\|^2 - \sum_{j=1}^M \lambda_j \right\}\end{aligned}$$



PPCA-Schätzung bei bekannter Störvarianz σ^2

Lemma (Meincke & Ritter 2000)

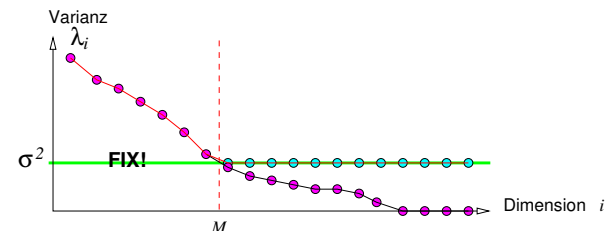
Ein Zufallsvektor sei gemäß $\mathbb{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{S})$ normalverteilt mit der Kovarianzmatrix

$$\mathbf{S} = \boldsymbol{\Psi} + \sigma^2 \cdot \mathbf{E}_N$$

mit **bekannter** Störvarianz σ^2 und positiv-semidefinitem $\boldsymbol{\Psi}$ mit **unbekanntem** Rang $\nu = \text{ran } \boldsymbol{\Psi}$, $\nu \leq N$.

Mit der empirischen Datenkovarianz $\hat{\mathbf{S}}(\omega)$ und ihrer Eigenzerlegung $\hat{\mathbf{S}}(\omega) = \mathbf{U} \cdot \mathbf{D} \cdot \mathbf{U}^\top$, $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_N)$ ergeben sich die ML-Schätzer

$$\hat{\nu} = |\{\lambda_i \mid \lambda_i > \sigma^2\}| \quad \text{und} \quad \hat{\boldsymbol{\Psi}} = \mathbf{U}_\nu \cdot (\mathbf{D}_\nu - \sigma^2 \cdot \mathbf{E}_\nu) \cdot \mathbf{U}_\nu^\top.$$



Effiziente Berechnung der PPCA-Dichtewerte

auch in extrem hochdimensionalen ($N \gg M$) Vektorräumen

Determinante $\det(\mathbf{S})$

Determinanten sind rotationsinvariant ($\det(\mathbf{S}) = \det(\mathbf{U}^\top \mathbf{S} \mathbf{U})$); also gilt:

$$\det(\mathbf{S}) = \prod_{i=1}^N \tilde{\lambda}_i = \sigma^{2 \cdot (N-M)} \cdot \prod_{i=1}^M \lambda_i.$$

Quadratische Form $(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu})$

Wegen der Darstellung von \mathbf{S}^{-1} gilt für die quadratische Form

$$\begin{aligned} \mathbf{y}^\top \mathbf{S}^{-1} \mathbf{y} &= \frac{1}{\sigma^2} \cdot \left\{ \mathbf{y}^\top \mathbf{E}_N \mathbf{y} - \mathbf{y}^\top \mathbf{U}_M \cdot (\mathbf{E}_M - \sigma^2 \mathbf{D}_M^{-1}) \cdot \mathbf{U}_M^\top \mathbf{y} \right\} \\ &= \frac{1}{\sigma^2} \cdot \left\{ \|\mathbf{y}\|^2 - \|\tilde{\mathbf{y}}\|^2 + \sigma^2 \cdot \tilde{\mathbf{y}}^\top \mathbf{D}_M^{-1} \tilde{\mathbf{y}} \right\} \\ &= \frac{\|\mathbf{y}\|^2 - \|\tilde{\mathbf{y}}\|^2}{\sigma^2} + \sum_{i=1}^M \frac{\tilde{y}_i^2}{\lambda_i} \end{aligned}$$

unter Verwendung des Hauptachsenprojektionsvektors

$$\tilde{\mathbf{y}} = \mathbf{U}_M^\top \mathbf{y} = \mathbf{U}_M^\top \cdot (\mathbf{x} - \boldsymbol{\mu}).$$

Zweistufiges EM-Abkühlverfahren

(Algorithmus)

- 1 Vorwahl von $\sigma_{\max}^2 > 0$, $\sigma_{\min}^2 > 0$ und $\alpha \in (0, 1)$.
- 2 Setze $\theta \leftarrow (\boldsymbol{\mu}, \dots, \boldsymbol{\mu})$, $m \leftarrow 0$ und $\sigma_m^2 \leftarrow \sigma_{\max}^2$.
- 3 SPHÄRISCHE GRUPPIERUNG (EM)

$$\ell(\theta \mid \omega, \sigma_m^2) = \sum_t \sum_{\kappa} \gamma_{\kappa, t} \cdot \log(\pi_{\kappa} \cdot \mathcal{N}(\mathbf{x}_t \mid \boldsymbol{\mu}_{\kappa}, \sigma_m^2 \mathbf{E})) \xrightarrow{!} \max$$

- 4 Setze $m \leftarrow m + 1$ und $\sigma_m^2 \leftarrow \alpha \cdot \sigma_{m-1}^2$.
- 5 Wenn $K_{\text{eff}} < K$ dann \rightsquigarrow 3.
- 6 LOKALADAPTIVE PPCA-GRUPPIERUNG (EM)

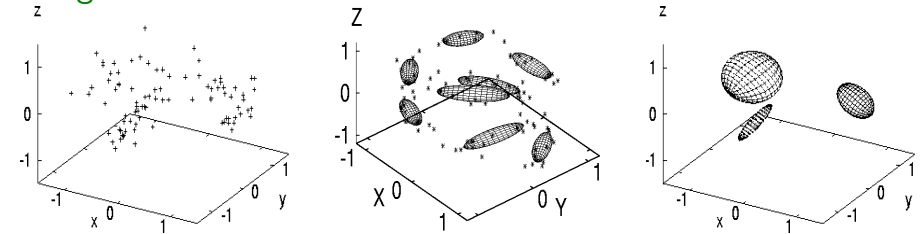
$$\mathcal{L}(\Theta \mid \omega, \sigma_m^2) = \sum_t \sum_{\kappa} \gamma_{\kappa, t} \cdot \log(\pi_{\kappa} \cdot \mathcal{N}(\mathbf{x}_t \mid \boldsymbol{\mu}_{\kappa}, \boldsymbol{\Psi}_{\kappa} + \sigma_m^2 \mathbf{E})) \xrightarrow{!} \max$$

- 7 Setze $m \leftarrow m + 1$ und $\sigma_m^2 \leftarrow \alpha \cdot \sigma_{m-1}^2$.
- 8 Wenn $\sigma_m^2 > \sigma_{\min}^2$ dann \rightsquigarrow 6 sonst Ende.

(zumflieg!A)

PPCA-Mischverteilungsidentifikation

Rangdefizite Cluster



Punktmenge im \mathbb{R}^3

unterschiedliche Richtungen

unterschiedliche Ränge

Sphärische Gruppierung

Alle Gaußkomponenten sphärisch
(kugelförmige Konzentration;
konstante Streuung σ_m^2)

Klassenmittelwertvektoren $\boldsymbol{\mu}_{\kappa}$

Ende sobald Anzahl K_{eff}
Gruppenprototypen gleich
Sollgruppenzahl K ist.

Lokaladaptive Gruppierung

PPCA-Gaußkomponenten
(zeppelförmige Konzentration;
variable Effektivdimension ν_{κ})

Rangdefiziente Matrizen $\boldsymbol{\Psi}_{\kappa}$

⬇ Reststreuung

⬆ Ränge (Parameterkomplexität)

Beispiel — Handgeschriebene Ziffern

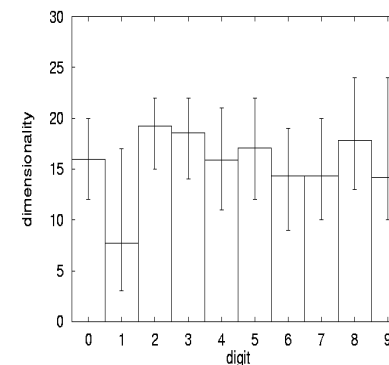
MNIST Datensammlung, LeCun 1998

Datensatz

60 000 (10 000) Lern- und Testmuster

Originalziffern 28×28 Pixel zu 8 Bit \rightsquigarrow Umrasterung 8×8

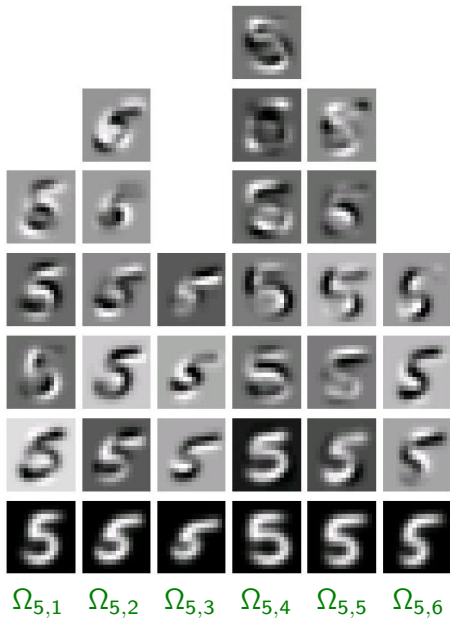
Verschiedene Gruppenstärken $K \in \{1, 2, 4, 8, 16\}$ getestet



PPCA-Dimensionen

alle zehn Ziffernklassen
 ν -Durchschnitt und min/max
 $K = 16$ Mischungskomponenten
im „Gefrierpunkt“ σ_{\min}^2

Beispiel — die Ziffernklasse „5“



PPCA-Mischung

für Ziffernklasse Ω_5

Gruppenanzahl

$K = 6$ gewählt

Modellrang

$M \in \{4, 5, 6, 7\} \Rightarrow M \ll 64$

Hauptachsen

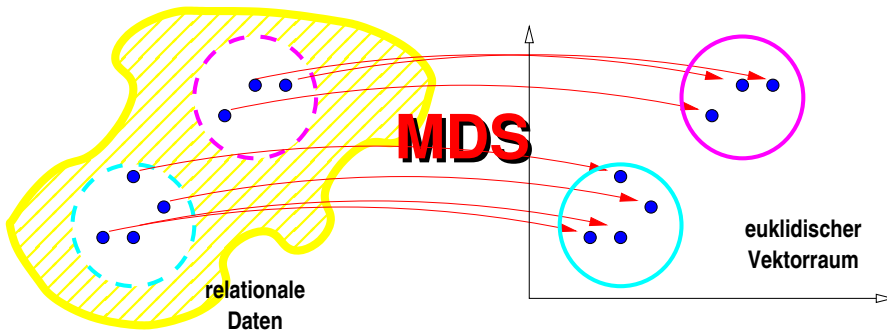
von unten nach oben

aufgetürmt

als 8×8 -Grauwertbild

Relationale Gruppierung

Datenobjekte mit wechselseitiger Distanz — ohne Attribute



MDS-Gruppierung

1. Mehrdimensionale Skalierung von ω nach $\omega' = \{\mathbf{x}_1, \dots, \mathbf{x}_T\} \subseteq \mathbb{R}^N$
2. K-means Gruppierung des Datensatzes ω'
3. Aufprägung der ω' -Gruppierung auf die Urbilder ω

Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K-means

Mischungsidentifikation

Relationale Gruppierung

Konzeptuelle Gruppierung

Spektrale Gruppierung

Clustergütemaße

Zusammenfassung

Fuzzy K-medoids für Metriken

Datensatz

Objektmenge $\omega = \{o_1, \dots, o_T\}$ mit der Abstandsmatrix

$$R = [r_{s,t}] \in \mathbb{R}^{T \times T}, \quad r_{s,t} = d(o_s, o_t)$$

Insbesondere gelte die **Symmetrie** $R^T = R$ und die **Definitheit** $\text{diag}(R) = 0$.

Zugehörigkeitsfunktionen

- **Harmonisch:**
$$u_{\kappa}(o_t) = 1 / \sum_{\lambda=1}^K \left(\frac{r_{t_{\kappa},t}}{r_{t_{\lambda},t}} \right)^{\frac{2}{\alpha-1}}$$
- **Cauchy / possibilistisch:**
$$u_{\kappa}(o_t) = \left(1 + \left(\frac{r_{t_{\kappa},t}}{\sqrt{\eta_{\kappa}}} \right)^{\frac{2}{\alpha-1}} \right)^{-1}$$
- **Hyperkonisch:**
$$u_{\kappa}(o_t) = \begin{cases} 1 - r_{t_{\kappa},t}/\rho_{\kappa} & \text{falls } r_{t_{\kappa},t} \leq \rho_{\kappa} \\ 0 & \text{sonst} \end{cases}$$

Für *harmonische* Zugehörigkeiten gilt die Normierungseigenschaft $\sum_{\lambda} u_{\lambda}(o_t) = 1$.

RACE — Relationaler Austauschalgorithmus

(Algorithmus)

GEGEBEN

Datenrelation $R \in \mathbb{R}^{T \times T}$, Gruppenzahl $K \in \mathbb{N}$, Iterationen $I \in \mathbb{N}$.

1 INITIALISIERUNG

Setze $i \leftarrow 1$.

Wähle zufällige Prototypenindizes $\{t_1, \dots, t_K\} \subseteq \{1, \dots, T\}$.

2 ITERATIONSSCHRITT

- Bestimme alle Zugehörigkeiten $u_\kappa(o_t)$
- Bestimme die „Restenergien“

$$e_{\kappa,t} = \sum_{\lambda \neq \kappa} u_\lambda(o_t)$$

- Bestimme die neuen Prototypen

$$t_\kappa \leftarrow \underset{t=1..T}{\operatorname{argmin}} e_{\kappa,t}$$

3 TERMINIERUNG

Wenn $i = I$ dann \rightsquigarrow Ende sonst $i \leftarrow i + 1, \rightsquigarrow$ 2.

(sumrtiniogIA)

Beispiel — Text Mining

Automatische Erstellung eines Stichwortinventars

Datensammlung

Kapitel 2 aus dem Buch „Information Mining“ (Th. Runkler)

Alle Formeln und Sonderzeichen wurden entfernt.

Großbuchstaben \mapsto Kleinbuchstaben

Objektmenge und Metrik

1605 Wortvorkommen, davon $T = 564$ verschieden

Matrix $R \in \mathbb{R}^{564 \times 564}$ der Levenshteinabstände

Verarbeitung

RACE-Algorithmus mit $K = 20$ Gruppen und $I = KT = 11280$ Schritten

ESS-Defuzzifizierung auf 28 (bzw. 29) Wörter/Gruppe

Die 20 häufigsten Wörter des Textes

die der und für in als werden ist sich mit
oder den sind ein auch daten läßt können abstand wird

Defuzzifizierung

Finales Schärfen (Aushärten) der Gruppenzugehörigkeiten

$$u_\kappa(o_t) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{falls } u_\kappa(o_t) = \max_\lambda u_\lambda(o_t) \\ 0 & \text{sonst} \end{cases}, \quad \kappa = 1, \dots, K$$

Diese Aushärtungsregel ergibt Partitionen mit variablen Gruppenstärken.

ESS-Defuzzifizierung

(„equal size subset“)

1 INITIALISIERUNG

Setze $\mathcal{I} \leftarrow \{1, 2, \dots, T\}$ (Indexmenge)

Setze $\mathcal{C}_\kappa \leftarrow \emptyset$ (Gruppenindexmenge) für alle $\kappa = 1..K$

2 ITERATION (für alle $t = 1, \dots, T$)

- Setze Gruppenindex $\lambda = t \bmod K + 1$.
- Bestimme bestpassenden Restindex $t^* = \operatorname{argmax}_{t \in \mathcal{I}} u_\lambda(o_t)$.
- Verschiebe Index t^* von \mathcal{I} nach \mathcal{C}_λ .

3 TERMINIERUNG

Jede Gruppe enthält entweder $\lfloor T/K \rfloor$ oder $\lceil T/K \rceil$ Datenelemente.

(sumrtiniogIA)

Beispiel — Text Mining

Gruppenprototypen $o_{t_1}, o_{t_2}, \dots, o_{t_{20}}$

originalsignal	mengenschreibweise	bzw
inkompatibilität	quantisierungsschritte	wertkontinuierlich
intervallskalierten	unterschiedlichen	übereinstimmungen
matrixdarstellung	mindestabtastrate	abtastzeitpunkten
datencharakteristika	ordnungsrelation	objektdatensatz
quantisierungsfehler	polygonzug	kovarianzmatrix
speicherplatzes	kaufmännisches	

Gruppen Ω_1, Ω_2 und Ω_6

(die Wörter mit den höchsten Zugehörigkeitsbewertungen)

Gruppe 1	Gruppe 2	Gruppe 6
originalsignal	mengenschreibweise	wertkontinuierlich
zeitsignal	matrixschreibweise	zeitkontinuierlich
ordinal	beschreiben	kontinuierliche
signal	schriftweise	wertebereich
signals	schreiben	rekonstruieren
digitalen	beschreibt	nichtnumerisch
zeitsignalen	geschrieben	willkürlich
proportional	beschrieben	konstruierten

Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K-means

Mischungsidentifikation

Relationale Gruppierung

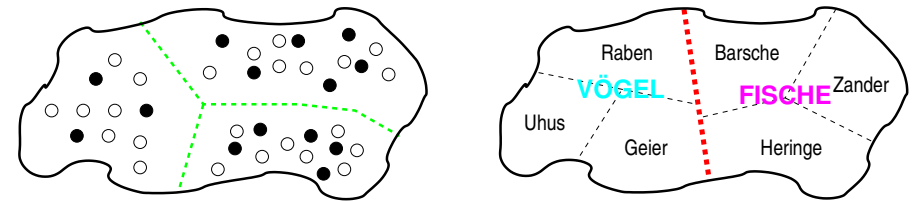
Konzeptuelle Gruppierung

Spektrale Gruppierung

Clustergütemaße

Zusammenfassung

Konzeptuelle Gruppierung



Gegeben

Objektbereich Ω · Hypothesenraum \mathcal{H} · Beispielmenge $\omega \subseteq \Omega$

Gesucht

eine **intensionale Partition** von ω , d.h.:

*Eine Folge von Hypothesen h_1, \dots, h_K ,
welche die beobachteten Beispiele aus ω
sowie auch neue Objekte aus $\Omega \setminus \omega$
überschneidungsfrei gegeneinander abgrenzen.*

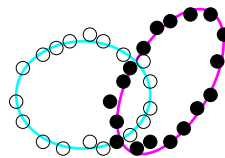
Und welches **Gütekriterium** optimiert die Partition ? ... ?

Lokale versus globale Objektähnlichkeit

Traditionelle Clusteranalyse

Lokaler Ähnlichkeitsbegriff (Hamming-Distanz)

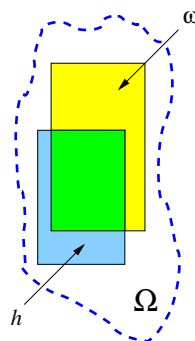
$$d(\mathbf{x}, \mathbf{y}) = \#\{\ell \mid x_\ell \neq y_\ell\}$$



Konzeptuelle Gruppierung

Globale Ähnlichkeit (engste \mathcal{H} -Umfassung)

$$d(\mathbf{x}, \mathbf{y}) = \min\{\sigma(h, \omega) \mid h \models \mathbf{x}, \mathbf{y}\}$$



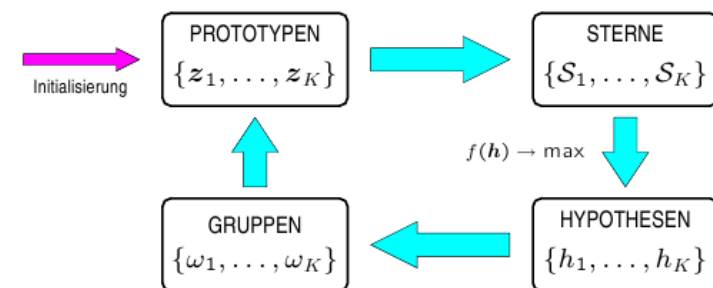
Überdeckungsgrad

der Hypothese h durch die Daten ω :

$$\sigma(h, \omega) = \frac{\#\{\mathbf{x} \in \omega \mid h \models \mathbf{x}\}}{\#\{\mathbf{x} \in \Omega \mid h \models \mathbf{x}\}} = \frac{\text{green}}{\text{blue}}$$

Wiederholtes Austauschen von Gruppenelementen

„conceptual K-means“



Gütekriterium

Kumulativer Überdeckungsgrad

$$f(\mathbf{h}) = \frac{1}{K} \cdot \sum_{k=1}^K \sigma(h_k, \omega)$$

Gruppenprototypen

müssen unbedingt ω angehören!
(Sterne, Überdeckung)

- Medoide
- Pseudomediane

Wahl der Gruppenzentren

Medoid

Dasjenige Gruppenelement mit minimaler **Exzentrizität**:

$$\mu_{\text{med}}(\omega_{\kappa}) = \underset{\mathbf{x} \in \omega_{\kappa}}{\operatorname{argmin}} \varepsilon(\mathbf{x}, \omega_{\kappa}) = \underset{\mathbf{x} \in \omega_{\kappa}}{\operatorname{argmin}} \sum_{\mathbf{y} \in \omega_{\kappa}} d(\mathbf{x}, \mathbf{y})$$

Pseudomedian

Kombiniere je nach Skalentyp der \mathcal{X}_n , $n = 1, \dots, N$ komponentenweise Mittelwerte, Zentroide, Mediane:

$$\mu_{\text{pseudo}}(\omega_{\kappa}) \stackrel{\text{def}}{=} (\mu_{\kappa,1}, \dots, \mu_{\kappa,N})^{\top}, \quad \mu_{\kappa,n} \stackrel{\text{def}}{=} \mu(\{x_n \mid \mathbf{x} \in \omega_{\kappa}\})$$

Ergibt eine hocheffizient berechenbare Näherung $\mu_{\text{pseudo}} \approx \mu_{\text{med}}$.

Reintegration

Wegen $\mu_{\text{pseudo}}(\omega_{\kappa}) \notin \omega_{\kappa}$ verwende den nächsten ω_{κ} -Nachbarn:

$$\tilde{\mu}_{\text{pseudo}}(\omega_{\kappa}) = \underset{\mathbf{x} \in \omega_{\kappa}}{\operatorname{argmin}} d(\mu_{\text{pseudo}}(\omega_{\kappa}), \mathbf{x})$$

Beispiel — Gebrauchtwagenhandel

Datensammlung

Objekt:	\mathbf{o}_1	\mathbf{o}_2	\mathbf{o}_3	\mathbf{o}_4	\mathbf{o}_5	\mathbf{o}_6	\mathbf{o}_7	\mathbf{o}_8
x_1 Geschwindigkeit	<i>h</i>	<i>m</i>	<i>h</i>	<i>l</i>	<i>m</i>	<i>l</i>	<i>h</i>	<i>m</i>
x_2 Farbe	<i>r</i>	<i>r</i>	<i>g</i>	<i>b</i>	<i>b</i>	<i>g</i>	<i>b</i>	<i>r</i>
x_3 Preis	<i>h</i>	<i>l</i>	<i>h</i>	<i>rh</i>	<i>rl</i>	<i>l</i>	<i>rh</i>	<i>rh</i>

$$\begin{aligned} x_1 \in \mathcal{X}_1 &= \{\text{high, medium, low}\} \\ x_2 \in \mathcal{X}_2 &= \{\text{red, blue, green}\} \\ x_3 \in \mathcal{X}_3 &= \{\text{high, rel_high, rel_low, low}\} \end{aligned}$$

Hypothesen als Attributkomplexe

zum Beispiel: $x_1 \in \{h\} \wedge x_2 \in \{b, g\} \wedge x_3 \in \{h, rh, rl\}$

<i>r</i>	1			
<i>b</i>		7		
<i>g</i>	3			
	<i>h</i>	<i>rh</i>	<i>rl</i>	<i>l</i>

Ebene $x_1 = h$

<i>r</i>		8		2
<i>b</i>			5	
<i>g</i>				
	<i>h</i>	<i>rh</i>	<i>rl</i>	<i>l</i>

Ebene $x_1 = m$

<i>r</i>				
<i>b</i>		4		
<i>g</i>				6
	<i>h</i>	<i>rh</i>	<i>rl</i>	<i>l</i>

Ebene $x_1 = l$

Überdeckung

$$\sigma(h, \omega) = \frac{2}{6} = 0.\bar{3}$$

Konzeptueller Austauschalgorithmus

(Algorithmus)

1 INITIALISIERUNG

Wähle Klassenzahl K und wähle $\mathbf{z}_1, \dots, \mathbf{z}_K$ zufällig aus ω .

2 PROTOTYPEN \Rightarrow STERNE

$$\mathcal{S}_{\kappa} = \mathcal{S}(\mathbf{z}_{\kappa} \mid \{\mathbf{z}_1, \dots, \mathbf{z}_K\} \setminus \{\mathbf{z}_{\kappa}\})$$

3 STERNE \Rightarrow HYPOTHESEN

Bestimme ω -**eindeutigen** Hypothesensatz $\mathbf{h} \in \mathcal{S}_1 \times \dots \times \mathcal{S}_K$ mit

$$f(\mathbf{h}) = \text{MAX}$$

4 HYPOTHESEN \Rightarrow GRUPPEN

$$\omega_{\kappa} = \{x \in \omega \mid h_{\kappa} \models x\}$$

5 GRUPPEN \Rightarrow PROTOTYPEN

$$\mathbf{z}_{\kappa} = \mu_{\text{med}}(\omega_{\kappa})$$

6 TERMINIERUNG Wenn $f(\mathbf{h}) \geq \theta$ dann \rightsquigarrow Ende sonst \rightsquigarrow 2.

Beispiel — Gebrauchtwagenhandel

Erster Iterationsschritt

P1 Wähle als initiale Prototypen $\mathbf{z}_1 = \mathbf{o}_1$ und $\mathbf{z}_2 = \mathbf{o}_2$ aus

S1 Stern von \mathbf{z}_1 :

$$h_1^1 = (x_2 = r, b) \wedge (x_3 = h, rh, rl) \quad \text{und} \quad h_1^2 = (x_1 = h, l)$$

Stern von \mathbf{z}_2 :

$$h_2^1 = (x_1 = h) \wedge (x_2 = g) \vee (x_3 = l) \quad \text{und} \quad h_2^2 = (x_1 = m)$$

G1 Dann gilt

$$\begin{aligned} h_1^1 &\models \mathbf{o}_1, \mathbf{o}_4, \mathbf{o}_5, \mathbf{o}_7, \mathbf{o}_8 \\ h_1^2 &\models \mathbf{o}_1, \mathbf{o}_3, \mathbf{o}_4, \mathbf{o}_6, \mathbf{o}_7 \end{aligned}$$

$$\begin{aligned} h_2^1 &\models \mathbf{o}_2, \mathbf{o}_3, \mathbf{o}_6 \\ h_2^2 &\models \mathbf{o}_2, \mathbf{o}_5, \mathbf{o}_8 \end{aligned}$$

H1 Nur die Kombinationen (h_1^1, h_2^1) und (h_1^2, h_2^2) bilden konzeptuelle Partitionen

$$f(h_1^1) + f(h_2^1) = 5/18 + 3/12 = 38/72$$

$$f(h_1^2) + f(h_2^2) = 5/24 + 3/12 = 33/72$$

Beispiel — Gebrauchtwagenhandel

Zweiter Iterationsschritt

P2 Objekte und ihr Median in h_1^1

Attribut	\mathbf{o}_1	\mathbf{o}_4	\mathbf{o}_5	\mathbf{o}_7	\mathbf{o}_8	Modus
x_1	h	l	m	h	m	h, m
x_2	r	b	b	b	r	b
x_3	h	rh	rl	rh	rh	rh

⇒ Zentroid ist \mathbf{o}_7

Die Hypothese h_2^1 hat Pseudomedian (m, g, l) und Median \mathbf{o}_6

⇒ neue Gruppenprototypen sind $\mathbf{z}_1 = \mathbf{o}_7, \mathbf{z}_2 = \mathbf{o}_6$

S2 Stern von \mathbf{o}_7 :

$$h_1^1 = (x_3 = h, rh)$$

$$h_1^2 = (x_1 = h, m) \wedge (x_2 = r, b) \wedge (x_3 = r, rh)$$

Stern von \mathbf{o}_6 :

$$h_2^1 = (x_1 = m, l) \wedge (x_3 = rl, l)$$

G2 Dann gilt

$$h_1^1 \models \mathbf{o}_1, \mathbf{o}_3, \mathbf{o}_4, \mathbf{o}_7, \mathbf{o}_8$$

$$h_1^2 \models \mathbf{o}_1, \mathbf{o}_7, \mathbf{o}_8$$

$$h_2^1 \models \mathbf{o}_2, \mathbf{o}_5, \mathbf{o}_6$$

H2 Nur die Kombination (h_1^1, h_2^1) bildet eine konzeptuelle Partition

$$f(h_1^1) + f(h_2^1) = 5/18 + 3/12 = 38/72$$

Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K-means

Mischungsidentifikation

Relationale Gruppierung

Konzeptuelle Gruppierung

Spektrale Gruppierung

Clustergütemaße

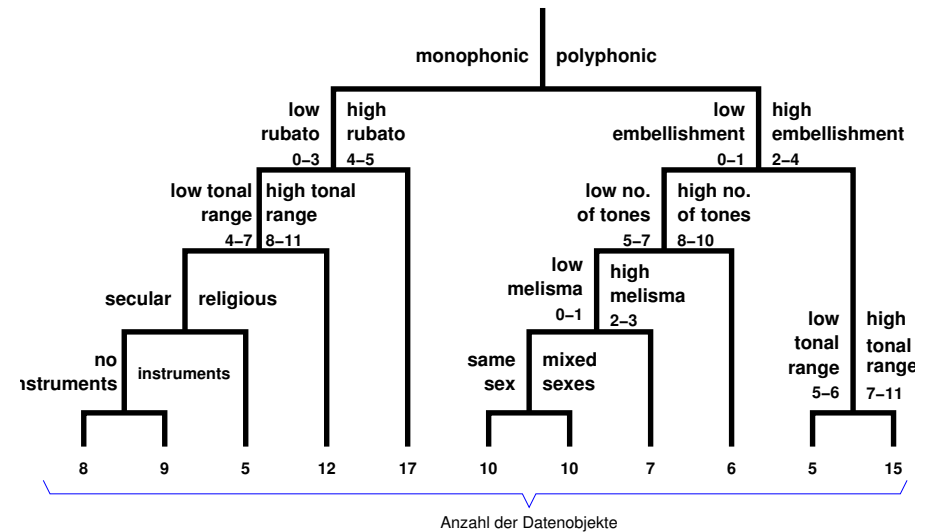
Zusammenfassung

Beispiel — Taxonomie spanischer Volkslieder

Gattungsdendrogramm nach konzeptueller Division

100 Lernbeispiele spanischer Volkslieder

22 Attribute mit nominalen / ordinalen Wertebereichen



Google Page Rank

„Gute Webseiten werden von guten Webseiten erwähnt.“

Relevanz und Qualität

Seitenbewertung = Anfragepassung + Seriositätsmaß

$$\text{score}_q^{\text{Google}}(\text{doc}) = \text{Rel}_q(\text{doc}) + \text{rank}(\text{doc})$$

Worldwide Web als gerichteter Graph

Adjazenzmatrix $\mathbf{A} \in \{1, 0\}^{T \times T}$ mit $a_{st} = 1$ ⇔ $\text{doc}_i \mapsto \text{doc}_j$

Irrfahrtmodell

Der „Random Surfer“ besucht Webseiten mit W'keit p_j und der Politik

$$p_j = (1 - \beta) \cdot \frac{1}{T} + \beta \cdot \sum_i p_i \cdot a_{ij} \cdot \frac{1}{\sum_k a_{ik}}$$

Die **Gleichgewichtsverteilung** gehorcht einer Eigenwertaufgabe ($\lambda = 1$):

$$\mathbf{B} \cdot \mathbf{p} = \left((1 - \beta) \cdot \frac{1}{T} + \beta \cdot \tilde{\mathbf{A}} \right) \cdot \mathbf{p} = \mathbf{p} = \lambda \cdot \mathbf{p}, \quad \tilde{a}_{ij} \stackrel{\text{def}}{=} a_{ij} / \sum_k a_{ik}$$

Schnitte in gewichteten Graphen

Definition

Sei $(\mathcal{K}, \mathcal{E}, \mathbf{A})$ ein ungerichteter, gewichteter Graph mit nicht-negativer, symmetrischer **Affinitätsmatrix** \mathbf{A} . Für zwei Knotenmengen \mathcal{B}, \mathcal{C} sei

$$\ell_{\text{aff}}(\mathcal{B}, \mathcal{C}) = \sum_{s \in \mathcal{B}} \sum_{t \in \mathcal{C}} A_{st}, \quad \tilde{\ell}_{\text{aff}}(\mathcal{B}, \mathcal{C}) = \frac{\ell_{\text{aff}}(\mathcal{B}, \mathcal{C})}{\ell_{\text{aff}}(\mathcal{B}, \mathcal{K})}$$

definiert. Eine Menge $\mathcal{C} \subset \mathcal{K}$ mit minimalem $\ell_{\text{aff}}(\mathcal{C}, \mathcal{K} \setminus \mathcal{C})$ bzw. mit minimalem $\tilde{\ell}_{\text{aff}}(\mathcal{C}, \mathcal{K} \setminus \mathcal{C})$ heißt **Schnitt** oder **normierter Schnitt**. Eine Partition $\mathcal{C}_1, \dots, \mathcal{C}_K$ von \mathcal{K} mit minimalem

$$\tilde{\ell}_{\text{aff}}(\{\mathcal{C}_\kappa\}_{\kappa=1}^K) \stackrel{\text{def}}{=} \frac{1}{K} \cdot \sum_{\kappa=1}^K \tilde{\ell}_{\text{aff}}(\mathcal{C}_\kappa, \mathcal{K} \setminus \mathcal{C}_\kappa)$$

heißt **normierter K-Schnitt**.

Bemerkung

Für die Affinitätsmatrix \mathbf{A} gilt $A_{ss} = 0$ und $A_{st} = A_{ts}$ für alle $s, t \in \{1, \dots, T\}$.

K-NC als Spurmaximierung

„K-way normalized cut“

Matrixalgebraische Formulierung

Die **Indikatormatrix** $\mathbf{C} \in \{1, 0\}^{T \times K}$ beschreibt die Zugehörigkeit der Knoten v_t zu den Gruppen \mathcal{C}_κ :

$$C_{t\kappa} \stackrel{\text{def}}{=} \begin{cases} 1 & t \in \mathcal{C}_\kappa \\ 0 & t \notin \mathcal{C}_\kappa \end{cases}$$

Die **Diagonalmatrix** $\mathbf{D} \in \mathbb{R}^{T \times T}$ enthält je Knoten die Summe seiner ausgehenden (einlaufenden) Kantengewichte:

$$\mathbf{D} = \text{diag}(\{d_s\}), \quad d_s \stackrel{\text{def}}{=} \sum_{t=1}^T A_{st}$$

Lemma

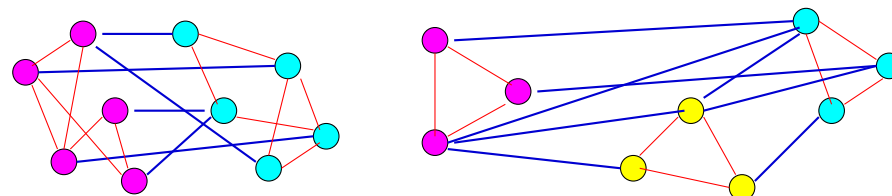
Das K-NC Kriterium ist äquivalent zur Maximierung der Größe

$$\frac{1}{K} \cdot \text{spur}(\mathbf{Z}^\top \mathbf{A} \mathbf{Z}) \quad \text{mit} \quad \mathbf{Z} \stackrel{\text{def}}{=} \mathbf{C} \cdot (\mathbf{C}^\top \mathbf{D} \mathbf{C})^{-1/2}.$$

Gewöhnliche & normierte Schnitte

Fakt

Die Berechnung eines (normierten) 2-Schnittes ist beweisbar NP-hart.



Gewöhnlicher Schnitt

Dichotome Partition von \mathcal{K} mit **minimaler Mengenaffinität**:

Summe der Querverbindungsgewichte zwischen \mathcal{C} und $\mathcal{K} \setminus \mathcal{C}$

Normierter Schnitt

Minimale relative Mengenaffinität:

Proportion der Querverbindungsgewichte zu den Gewichten aller \mathcal{C} verlassenden Kanten

Normierter K-Schnitt

Minimale Summe aller relativen Affinitäten zwischen den Schnittmengen \mathcal{C}_κ und ihren **Komplementen** $\mathcal{K} \setminus \mathcal{C}_\kappa$

Relaxationslösung

F. Chung: *Spectral Graph Theory*, AMS 1997

Skalierung

Die umskalierte Matrix $\tilde{\mathbf{Z}} := \mathbf{D}^{1/2} \mathbf{Z} \in \mathbb{R}^{T \times K}$ besitzt offenbar orthonormale Spalten:

$$\tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}} = \mathbf{E} \in \mathbb{R}^{K \times K}$$

➡ Spurmaximierung durch Berechnung der K ersten Eigenvektoren

Lemma

Die Matrix $\tilde{\mathbf{Z}} \in \mathbb{R}^{T \times K}$ mit den K oberen Eigenvektoren von $\tilde{\mathbf{A}} := \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ als Spalten maximiert die Spur

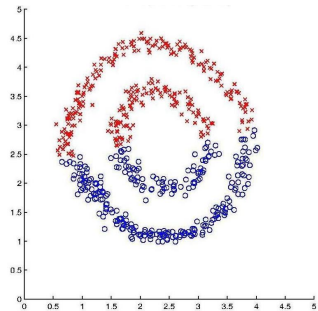
$$\text{spur} \left(\underbrace{\tilde{\mathbf{Z}}^\top \cdot \mathbf{D}^{-1/2}}_{\mathbf{Z}^\top} \cdot \mathbf{A} \cdot \underbrace{\mathbf{D}^{-1/2} \cdot \tilde{\mathbf{Z}}}_{\mathbf{Z}} \right)$$

unter der Bedingung $\tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}} = \mathbf{E}$.

Gelockerte Bedingung

Ersetze die komplizierte Strukturforderung für \mathbf{Z} durch Orthonormalitätsbedingung (links).

Spektrale Gruppierung statt K-means



K-means Algorithmus

modelliert ausschließlich **konvexe** Ballungsgebiete und findet nur **lokale** Verzerrungsminima.

2 Ringwolken — 2-means-Gruppierung

1. Bestimme $\tilde{\mathbf{A}}$
2. Berechne $\tilde{\mathbf{Z}}$ (EWP)
3. Ermittle $\mathbf{Z} = \mathbf{D}^{-1/2} \tilde{\mathbf{Z}}$
4. Errate (!) \mathbf{C} aus \mathbf{Z}

Spektrale Gruppierung

Die metrische Struktur der Datenobjekte wird in einen gewichteten Graphen transformiert; anschließend wird der Graph in Polynomialzeit durch Berechnung eines Semi-Schnittes in K Teilgraphen partitioniert.

Ng-Jordan-Weiss Algorithmus

1 AFFINITÄTSMATRIX

$$\mathbf{A} \in \mathbb{R}^{T \times T} \text{ mit } A_{st} \stackrel{\text{def}}{=} \begin{cases} 0 & s = t \\ \exp \{-\| \mathbf{x}_s - \mathbf{x}_t \|^2 / 2\sigma^2\} & s \neq t \end{cases}$$

2 LAPLACEMATRIX

$$\mathbf{L} \in \mathbb{R}^{T \times T} \text{ via Normierung } L_{st} \stackrel{\text{def}}{=} \frac{A_{st}}{\sqrt{\rho_s \cdot \rho_t}}, \quad \rho_s \stackrel{\text{def}}{=} \sum_{r=1}^T A_{rs}$$

3 K FÜHRENDE EIGENVEKTOREN

$$\mathbf{L} = \sum_{t=1}^T d_t^2 \cdot \mathbf{u}_t \mathbf{u}_t^\top, \quad \mathbf{U} \stackrel{\text{def}}{=} (\mathbf{u}_1, \dots, \mathbf{u}_K) \in \mathbb{R}^{T \times K}$$

4 ZEILENWEISE NORMIERUNG

$$\tilde{\mathbf{U}} \stackrel{\text{def}}{=} \mathbf{C}^{-1/2} \cdot \mathbf{U}, \quad C_{st} = \begin{cases} \sum_{\kappa=1}^K U_{t\kappa}^2 & s = t \\ 0 & s \neq t \end{cases}$$

5 K-MEANS GRUPPIERUNG

der Matrixzeilen $\tilde{\pi} : \{\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_T\} \rightarrow \{1, 2, \dots, K\}$

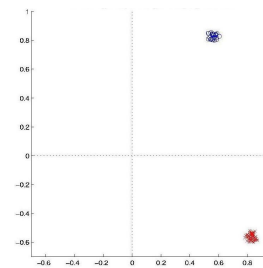
6 PARTITIONIERUNG

$$\text{der Originaldaten } \pi : \begin{cases} \{\mathbf{x}_1, \dots, \mathbf{x}_T\} & \rightarrow \{1, 2, \dots, k\} \\ \mathbf{x}_t & \mapsto \tilde{\pi}(\tilde{\mathbf{v}}_t) \end{cases}$$

Gruppieren im Spektralraum

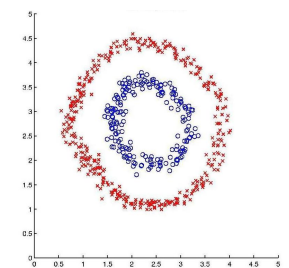
des Ähnlichkeitsgraphen

$$\mathbb{R}^{T \times N} \xrightarrow{\text{Affinität}} \mathbb{R}^{T \times T} \xrightarrow{\text{Eigenraum}} \mathbb{R}^{T \times K} \xrightarrow{\text{Gruppen}} \mathbb{R}^{K \times K}$$



Zeilen der $\mathbb{R}^{T \times 2}$ -Matrix

Gruppen-ID



Spektrale Gruppen ($K = 2$)

Eigenraummatrix $\mathbf{U} \in \mathbb{R}^{T \times K}$

Die K -dimensionalen Zeilenvektoren weisen eine hochdiskriminante Gruppenstruktur (innerer Ring — äußerer Ring) auf.

Warum funktioniert der NJW-Algorithmus ?

Beobachtung (geodätische Gruppenbildung)

Die Zeilen der Matrix \mathbf{U} bilden eine K -dimensionale Repräsentation der Daten, in der Objekte mit kurzem Verbindungsweg — geodätisch, nicht Luftlinie — nahe beieinander liegen. warum?

Idealtypisches Szenarium ($K = 3$)

Objekte unterschiedlicher Gruppe besitzen den euklidischen Abstand ∞ .

- Affinitätsmatrix und Laplacematrix sind von **Blockdiagonalform** (geeignete Nummerierung der \mathbf{x}_t vorausgesetzt)

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_1 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{A}_3 \end{pmatrix}, \quad \mathbf{L} = \begin{pmatrix} \mathbf{L}_1 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{L}_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{L}_3 \end{pmatrix}$$

- Die Eigenvektoren von \mathbf{L} sind die *mit Nullen aufgefüllten* Eigenvektoren der Blockmatrizen \mathbf{L}_κ .
- Dank der Doppelnormierung von \mathbf{L} besitzt jeder Block *genau einen* maximalen **Eigenwert Eins**.

Warum funktioniert der NJW-Algorithmus ?

Idealtypisches Szenarium ($K = 3$)

Die K Haupteigenvektoren von \mathbf{L} rekrutieren sich aus den K **Blockgewinnern**.

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}'_1 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{u}'_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{u}'_3 \end{pmatrix} \in \mathbb{R}^{T \times 3}, \quad \tilde{\mathbf{U}} = \begin{pmatrix} \mathbf{100} \\ \mathbf{010} \\ \mathbf{001} \end{pmatrix} \in \mathbb{R}^{T \times 3}$$

- Die Matrix \mathbf{U} besitzt die K Eigenvektoren als Spalten.
- Alle Zeilen von \mathbf{U} enthalten **genau einen** Eintrag ungleich Null.
- In der *zeilensummennormierten* Matrix $\tilde{\mathbf{U}}$ wird der Eintrag zur Eins.
- Das Gruppieren der Zeilen (\approx Einheitsvektoren) ergibt zwangsläufig genau die richtigen Cluster!

Details zum NJW-Algorithmus

Lineare Störungstheorie

Analyse des NJW **ohne** Intercluster-Distanzen = ∞

Stewart & Sun: *Matrix Perturbation Theory*, 1990

- Eigengap $\lambda_K - \lambda_{K+1}$ als untere Schranke der Gruppierungsstabilität

Abklingparameter $\sigma^2 > 0$

Minimale Endverzerrung nach dem K -means Clustering

- Skalarer Optimierungsablauf für σ^2

Startpartition für K -means

Die (idealen) Gruppenzentren liegen auf der Einheitssphäre. Sie stehen paarweise senkrecht aufeinander.

- Sukzessive Auswahl derjenigen \mathbf{v}_t als Saatpunkte, die zu allen bereits selektierten Kandidaten maximal orthogonal sind.

Warum funktioniert der NJW-Algorithmus ?

Denkfehler

Die K Haupteigenvektoren besitzen den gemeinsamen Eigenwert Eins.

- Sie sind also keineswegs eindeutig bestimmt und voller Nullen, sondern spannen lediglich einen eindeutig bestimmten K -dimensionalen Unterraum auf.

Rettung der Argumentation

Statt \mathbf{U} erhalten wir $\mathbf{U}' = \mathbf{U}\mathbf{R}$ mit einer Rotationsmatrix $\mathbf{R} \in \mathbb{R}^{K \times K}$.

Bezeichne \mathbf{v}_t^\top die t -te Zeile von \mathbf{U} .

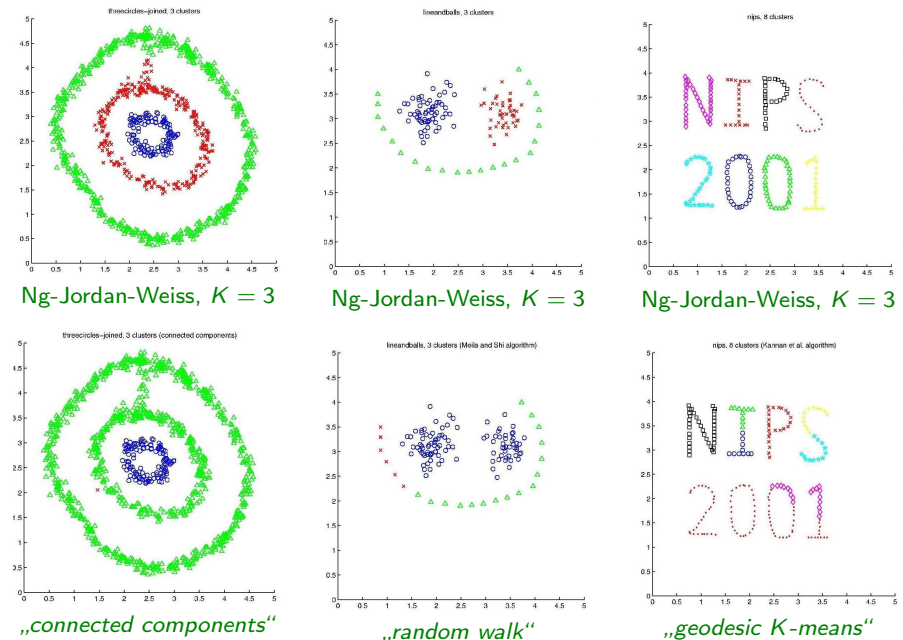
- $\mathbf{v}_t^\top \mathbf{R}$ ist die t -te Zeile von \mathbf{U}' und für die Quadratnorm gilt:

$$(\mathbf{v}_t^\top \mathbf{R}) \cdot (\mathbf{R}^\top \mathbf{v}_t) = \mathbf{v}_t^\top \cdot (\mathbf{R}\mathbf{R}^\top) \cdot \mathbf{v}_t = \mathbf{v}_t^\top \cdot \mathbf{v}_t = \|\mathbf{v}_t\|^2$$

Die Zeilennormierung macht $\mathbf{U}\mathbf{R}$ zu $\tilde{\mathbf{U}}\mathbf{R}$.

Zu clustern sind nicht mehr die **Einheitsvektoren** des \mathbb{R}^K , aber immerhin noch die Vektoren einer Orthonormalbasis \mathbf{R} des Raumes.

Beispiele — NJW & Wettbewerber



Ng-Jordan-Weiss

Die Konkurrenz

Sonstige spektral orientierte Verfahren

Gruppierung nach Zusammenhangskomponenten

Bilde den ungerichteten Graphen mit den Kanten

$$(s, t) \in \mathcal{E} \quad \Leftrightarrow \quad \|\mathbf{x}_s - \mathbf{x}_t\| \leq \theta .$$

Wähle die Schranke θ so, daß $\#(\text{ZSH-Komponenten}) = K$ ist.

Gruppierung nach dem Irrfahrtprinzip (Meila & Shi)

Unterm Strich derselbe Ablauf wie beim NJW-Algorithmus, aber:

- Nur die *Zeilen* der Affinitätsmatrix \mathbf{A} werden normiert.
- Die Zeilen der Eigenvektormatrix \mathbf{U} werden *nicht* normiert.

Geodätischer K-means (Kannan & Vempala & Vetta)

Wiederum derselbe Ablauf wie beim NJW-Algorithmus, aber:

- Nur die *Zeilen* der Affinitätsmatrix \mathbf{A} werden normiert.
- Die *Urbilder* der Clusterzentroide werden als *Repräsentanten* genutzt.

Dualisierte Berechnungen für K-means

Lemma

Sei $\omega = \{\mathbf{x}_1, \dots, \mathbf{x}_T\} \subset \Omega$ und $\phi : \Omega \rightarrow \mathbb{H}$ eine Expansion mit dem zugehörigen Kernoperator $K(\cdot, \cdot)$.

1. Das Zentroidelement der termexpandierten Daten $\phi(\omega)$ bezüglich des quadratischen euklidischen Abstandes $\|\cdot\|_{\mathbb{H}}^2$ ist der Mittelwertvektor

$$\mu = \frac{1}{T} \cdot \sum_{t=1}^T \phi(\mathbf{x}_t) .$$

2. Der Abstand zwischen μ und einem expandierten Objekt $\phi(\mathbf{y})$, $\mathbf{y} \in \Omega$, läßt sich mit $O(T^2)$ Kernoperatorauswertungen berechnen.
3. Die Berechnung der Abstände von μ zu allen $\phi(\mathbf{x}_t)$, $t = 1, \dots, T$, erfordert i.a. den Aufwand $O(T^2N)$, wenn $\Omega = \mathbb{R}^N$ ist.

K-means mit Termexpansion

Der Kernel Trick

Gruppieren in einem (impliziten) Expansionsraum $\phi(\Omega)$:

$$\phi : \Omega \rightarrow \mathbb{H} , \quad \langle \phi \mathbf{x}, \phi \mathbf{y} \rangle_{\mathbb{H}} = K(\mathbf{x}, \mathbf{y})$$

Der **Kernoperator** $K(\cdot, \cdot)$ simuliert das „Rechnen“ im RKHS \mathbb{H} .

Optimale K-Gruppierung in $\phi(\Omega)$

Ein Kodebuch $\{\mu_1, \dots, \mu_K\}$ mit minimaler **Verzerrung**

$$\varepsilon(\{\omega_\kappa\}_\kappa) = \sum_{\kappa=1}^K \sum_{\mathbf{x} \in \omega_\kappa} \|\phi(\mathbf{x}) - \mu_\kappa\|_{\mathbb{H}}^2$$

Berechnung von
Gruppenzentr(oid)en μ_κ

$\mu_\kappa \in \mathbb{H}$ mittelt (endlich
viele) expandierte Objekte.

Berechnung von
Prototypdistanzen $\|f - g\|_{\mathbb{H}}^2$

f ist ein expandiertes Objekt.
 g ist ein Gruppenzentrum.

Beweis.

1. \mathbb{H} als Hilbertraum ist insbesondere auch ein Vektorraum.
2. Kernbasierte Distanzberechnung (speziell $\mathbf{y} = \mathbf{x}_r \in \omega$):

$$\begin{aligned} \|\phi \mathbf{y} - \mu\|_{\mathbb{H}}^2 &= \left\| \phi \mathbf{y} - \frac{1}{T} \sum_{\mathbf{x}_t} \phi \mathbf{x}_t \right\|_{\mathbb{H}}^2 \\ &= \langle \phi \mathbf{y}, \phi \mathbf{y} \rangle - 2 \cdot \frac{1}{T} \cdot \sum_{\mathbf{x}_t} \langle \phi \mathbf{y}, \phi \mathbf{x}_t \rangle + \frac{1}{T^2} \cdot \sum_{\mathbf{x}_s, \mathbf{x}_t} \langle \phi \mathbf{x}_s, \phi \mathbf{x}_t \rangle \\ &= K(\mathbf{y}, \mathbf{y}) - 2 \cdot \frac{1}{T} \cdot \sum_{\mathbf{x}_t} K(\mathbf{y}, \mathbf{x}_t) + \frac{1}{T^2} \cdot \sum_{\mathbf{x}_s, \mathbf{x}_t} K(\mathbf{x}_s, \mathbf{x}_t) \\ &= \frac{1}{T^2} \cdot \left\{ T^2 \cdot G_{rr} - 2T \cdot \sum_{\mathbf{x}_t} G_{rt} + \sum_{\mathbf{x}_s, \mathbf{x}_t} G_{st} \right\} \\ &= G_{rr} - 2\bar{g}_r + \bar{G} \end{aligned}$$

3. Jede Kernoperatorauswertung kostet $O(N)$, die Berechnung der Gramschen Matrix kostet $O(T^2N)$, und die Mittelungen über \mathbf{G} und ihre Zeilen \mathbf{g}_r bleiben bei $O(T^2)$.

Kernel K-means Algorithmus

Kostenpunkt: $O(T^2N + T^2I)$

(Algorithmus)

1 STARTWERTE

Startgruppierung $\{\omega_\kappa^{(0)}\}$ und Kernmatrix \mathbf{G} , $G_{st} = K(\mathbf{x}_s, \mathbf{x}_t)$.

2 NEUE OBJEKTZUGEHÖRIGKEIT

$$\kappa^*(\mathbf{x}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\lambda=1..K} \|\phi\mathbf{x} - \mu_\lambda^{(i)}\|$$

3 IMPLIZITE ZENTROIDBERECHNUNG

$$\omega_\lambda^{(i+1)} \leftarrow \{\mathbf{x}_t \mid \kappa^*(\mathbf{x}_t) = \lambda\}$$

(keine *explizite* Berechnung der Mittelwertvektoren $\mu_\lambda^{(i+1)} \in \mathbb{H}$)

4 TERMINIERUNG

Wenn $\varepsilon^{(i)}(\{\omega_\kappa\}) \leq \theta$ dann **Ende** sonst $i \leftarrow i + 1$ und \rightsquigarrow 2.

(sumfnhogIA)

WKMM vs. normierter Schnitt

Minimale Verzerrung \leftrightarrow maximale NC-Matrixspur

Lemma

Die Gruppenverzerrung des WKMM ist gleichwertig zum Ausdruck

$$\varepsilon(\{\omega_\kappa\}, \mathbf{w}) = \operatorname{spur}(\mathbf{W}^{1/2} \mathbf{G} \mathbf{W}^{1/2}) - \operatorname{spur}(\mathbf{U}^\top \cdot \mathbf{W}^{1/2} \mathbf{G} \mathbf{W}^{1/2} \cdot \mathbf{U})$$

mit

$$\Phi = (\phi\mathbf{x}_1, \dots, \phi\mathbf{x}_T)$$

$$\mathbf{G} = \Phi^\top \Phi$$

$$\mathbf{W} = \operatorname{diag}(\{w(\mathbf{x}_t)\}_t)$$

$$\mathbf{U} = \operatorname{diag}(\{s_\lambda^{-1/2} \cdot \mathbf{W}_\lambda^{1/2} \cdot \mathbf{1}_\lambda\}), \quad s_\lambda = \mathbf{1}_\lambda^\top \mathbf{W}_\lambda \mathbf{1}_\lambda.$$

Die ersten K Eigenvektoren der Matrix $\mathbf{W}^{1/2} \mathbf{G} \mathbf{W}^{1/2}$ liefern eine verzerrungsminimale Lösung unter der **Relaxationsbedingung** $\mathbf{U}^\top \mathbf{U} = \mathbf{E}$ für die blockstrukturierte $(T \times K)$ -Matrix \mathbf{U} .

(Beweis durch exzessives Nachrechnen)

Gewichteter Kernel K-means

Gewichtetes Verzerrungskriterium

Objektabhängige Gewichtung der Zentrumsabstände:

$$\varepsilon(\{\omega_\kappa\}, \mathbf{w}) = \sum_{\kappa=1}^K \sum_{\mathbf{x} \in \omega_\kappa} w(\mathbf{x}) \cdot \|\phi(\mathbf{x}) - \mu_\kappa\|^2$$

\rightarrow Gruppenzentroide \triangleq **gewichtete** Mittelwertvektoren.

Optimale Gruppenstruktur

Minimale Verzerrung \Leftrightarrow maximale Spur (NC):

- Kernmatrix $\mathbf{W}^{1/2} \mathbf{G} \mathbf{W}^{1/2}$ korrespondiert mit normierter Affinität $\tilde{\mathbf{A}}$.
- Mantelmatrix \mathbf{U} korrespondiert mit $\tilde{\mathbf{Z}}$.

Jede Zeile von \mathbf{U} ist ein skaliertes K -Einheitsvektor.

$\rightarrow \mathbf{U}^\top \mathbf{U} = \mathbf{E}_{(K)}$ (\mathbf{U} besitzt orthonormale Spalten)

WKMM kann (im Relaxationssinne) auch durch NJW gelöst werden!

Spektrale Gruppierung

Vorzüge und nachteilige Eigenschaften

NJW-Algorithmus

- einstufiges Verfahren
- metrische Distanzen
- K -means über \mathbb{R}^K
- Eigenvektoren $O(T^2K)$
- Gruppen raten $O(TK^2I)$
- **Relaxationslösung** !?!

WKMM-Algorithmus

- iteratives Verfahren
- Mercer-reskalierbar
- K -means über \mathbb{H} /dual
- Gram-Matrix $O(T^2N)$
- dualer K -means $O(T^2I)$
- **Startpartition** ??

Hybride spektrale Gruppierung

(Algorithmus)

- 1 Berechne ggf. die Affinitäten \mathbf{A} , Zeilensummen \mathbf{D} und $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$.
- 2 Startgruppierung $\{\omega_\kappa^{(0)}\}$ via **NJW-Algorithmus** auf $\tilde{\mathbf{A}}$.
- 3 Berechne die (virtuelle) Gram-Matrix $\mathbf{G} = \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1}$ und setze $\mathbf{W} = \mathbf{D}$.
- 4 Führe den i -ten **Weighted-Kernel-K-Means**-Doppelschritt durch.
- 5 Wenn $\varepsilon^{(i)}(\{\omega_\kappa\}) \leq \theta$ dann **Ende** sonst $i \leftarrow i + 1$ und \rightsquigarrow 4.

(sumfnhogIA)

Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K-means

Mischungsidentifikation

Relationale Gruppierung

Konzeptuelle Gruppierung

Spektrale Gruppierung

Clustergütemaße

Zusammenfassung

Cluster Recovery Index

Externe Qualitätskriterien — Vergleich mit „Goldstandard“

Ausgangssituation

Objektmenge $\{o_1, \dots, o_T\}$ mit **wahrer** und **hypothetischer** Gruppierung:

$$\omega_1^* \uplus \omega_2^* \uplus \dots \uplus \omega_{K^*}^* \quad \text{versus} \quad \hat{\omega}_1 \uplus \hat{\omega}_2 \uplus \dots \uplus \hat{\omega}_K$$

Aufgabenstellung

Berechnen eines Übereinstimmungsmaßes zwischen $\{\omega_\kappa^*\}$ und $\{\hat{\omega}_\kappa\}$.

- Vergleich einer Ist-Lösung mit der Soll-Lösung
- Vergleich zweier Lösungen zweier Methoden

Problem

Gruppierungen sind nur eindeutig bis auf **Indexpermutation**.

Gütemaße für Gruppen & Partitionen

Fragen über Fragen

- Mit welcher Vorgabeordnung $K \in \mathbb{N}$ **starte** ich K-means ?
- Welche Zerlegungsebene **wähle** ich als Resultat aus?
(agnes, diana, pam & Co.)
- Zerfällt ω_κ in **noch kleinere** Gruppen ?
- Trifft gelernte Partition $\{\hat{\omega}_\kappa\}_\kappa$ die **wahre** Gruppenstruktur ?

Problem

Weder **Gruppenverzerrungen** $\varepsilon(\omega_\kappa)$

noch **Gesamtverzerrung** $\varepsilon(\{\omega_\kappa\}_\kappa)$

beantworten auch nur eine dieser Fragen !

Der Rand-Index

Überschneidungsfreie, scharfe Gruppen (W. M. Rand, 1971)

Kreuzadjazenzstatistiken

Die Objekte o_1, \dots, o_T bilden $M = \binom{T}{2}$ ungeordnete Paaren $\{o_s, o_t\}$.

	gleiche $\{\hat{\omega}_\kappa\}$ -Gruppe	verschiedene $\{\hat{\omega}_\kappa\}$ -Gruppen	
gleiche $\{\omega_\kappa^*\}$ -Gruppe	M_{11}	M_{10}	$M_{1.}$
verschiedene $\{\omega_\kappa^*\}$ -Gruppen	M_{01}	M_{00}	$M_{0.}$
	$M_{.1}$	$M_{.0}$	M

Definition

Unter dem **Rand-Index** zweier scharfer Objektpartitionen verstehen wir den relativen Anteil

$$C_{\text{rand}} \stackrel{\text{def}}{=} \frac{M_{11} + M_{00}}{M}$$

der kohärent gruppierten Punktepaare $\{o_s, o_t\}$.

Der bereinigte Rand-Index

„adjusted Rand index“ (Hubert & Arabie, 1985)

Problem

- ⊕ Maximum $C_{rand} = 1$ wird für äquivalente Partitionen angenommen.
- ⊖ Hohe C_{rand} -Werte entstehen auf Grund zufälliger Korrespondenzen.

Definition

Unter dem **bereinigten Rand-Index** zweier scharfer Objektpartitionen verstehen wir den Quotienten

$$C_{ari} \stackrel{\text{def}}{=} \frac{\text{observed} - \text{expected}}{\text{maximum} - \text{expected}} = \frac{C_{rand} - \{M_1 \cdot M_{.1} + M_0 \cdot M_{.0}\} / M^2}{1 - \{M_1 \cdot M_{.1} + M_0 \cdot M_{.0}\} / M^2}$$

aus **beobachtetem** und **größtmöglichem** Übertreffen der allein zufallsbedingten Gruppenkohärenz.

Bemerkung

Ein **störbereinigtes** und **permutationsinvariantes** Vergleichsmaß ist auch die **Transinformation** $\mathcal{H}(\mathbb{K}_1) + \mathcal{H}(\mathbb{K}_2) - \mathcal{H}(\mathbb{K}_1, \mathbb{K}_2)$ zwischen dem wahren und dem hypothetischen Gruppenindex der Objekte.

Bestimmung des Hopkins-Index

(Algorithmus)

- 1 RESAMPLING ω und $\bar{\omega}$
Ziehe $S \in \mathbb{N}$ Vektoren $\mathbf{z}_1, \dots, \mathbf{z}_S$ aus ω .
Ziehe $S \in \mathbb{N}$ Vektoren $\mathbf{y}_1, \dots, \mathbf{y}_S$ aus der konvexen Hülle

$$\bar{\omega} \stackrel{\text{def}}{=} \left\{ \sum_{t=1}^T a_t \mathbf{x}_t \mid \sum_{t=1}^T a_t = 1, a_t \geq 0 \right\}$$

- 2 PUNKTDICHTE IN $\bar{\omega}$
Berechne die kumulative Punkt-Mengen-Distanz

$$\mathcal{E}_{\bar{\omega}}[d^*(\mathbb{X}, \omega)] \approx D_{\bar{\omega}} = \sum_{s=1}^S \min_{\mathbf{x} \in \omega} d(\mathbf{y}_s, \mathbf{x})$$

- 3 PUNKTDICHTE IN ω
Berechne die „leave-one-out“ kumulative Punkt-Mengen-Distanz

$$\mathcal{E}_{\omega}[d^*(\mathbb{X}, \omega^{(\mathbb{X})})] \approx D_{\omega} = \sum_{s=1}^S \min_{\mathbf{x} \in \omega \setminus \{\mathbf{z}_s\}} d(\mathbf{z}_s, \mathbf{x})$$

- 4 AUSGABE
Berechne den Quotienten $C_{hop} = D_{\bar{\omega}} / (D_{\bar{\omega}} + D_{\omega})$.

(zumr3hog1A)

Die Heterogenität einer Punktmenge

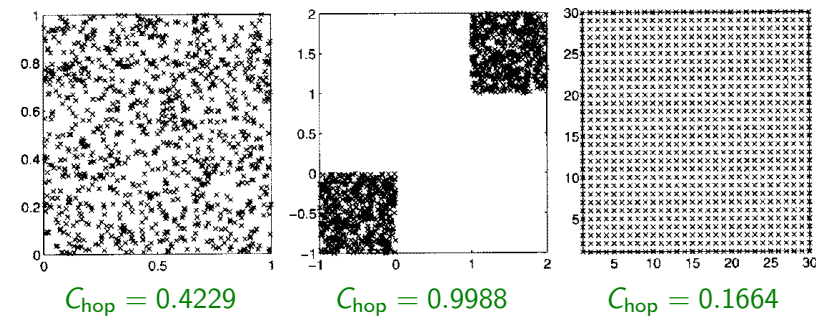
Zerfällt $\omega \subset \Omega$ in noch kleinere Gruppen ?

Definition

Für die Punktmenge ω mit der konvexen Hülle $\bar{\omega} \supset \omega$ ist der **Hopkins-Index** durch die relative Punktdichte

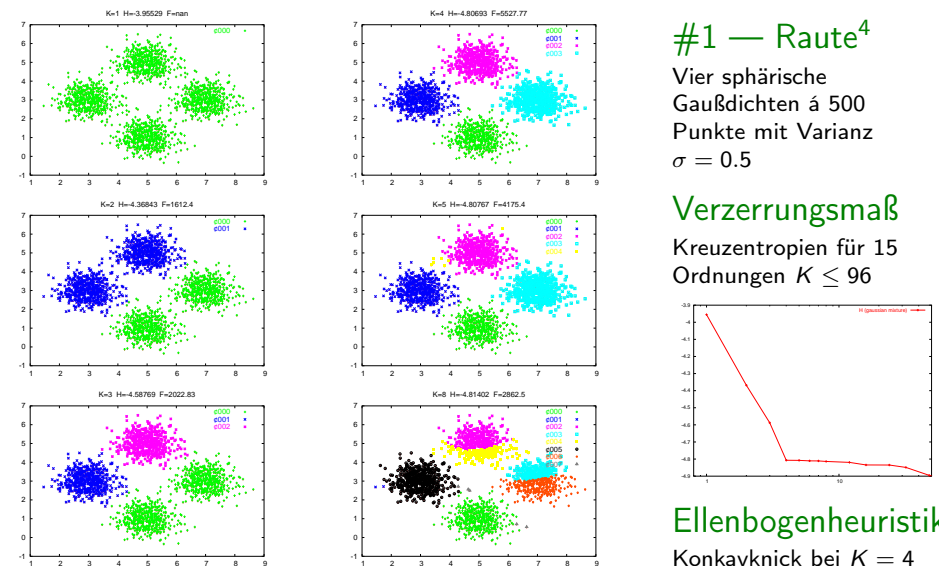
$$C_{hop} \stackrel{\text{def}}{=} \frac{\mathcal{E}_{\bar{\omega}}[d^*(\mathbb{X}, \omega)]}{\mathcal{E}_{\bar{\omega}}[d^*(\mathbb{X}, \omega)] + \mathcal{E}_{\bar{\omega}}[d^*(\mathbb{X}, \omega^{(\mathbb{X})})]}$$

im Gesamtbereich der ω -Hülle definiert.



Die beste Anzahl von Gruppen

Gretchenfrage $K \in \{1, 2, 3, 4, 5, 8\}$ bei der GMM-Identifikation



#1 — Raute⁴

Vier sphärische Gaußdichten á 500 Punkte mit Varianz $\sigma = 0.5$

Verzerrungsmaß

Kreuzentropien für 15 Ordnungen $K \leq 96$

Ellenbogenheuristik

Konkavknick bei $K = 4$ nebst Sättigungsauslauf

Ordnung und Güte einer Gruppierung

Je mehr Gruppen — desto paßgenauer das Datenmodell

Lemma

Für einen Datensatz $\omega = \{\mathbf{x}_1, \dots, \mathbf{x}_T\} \subset \mathbb{R}^N$ und die Ordnung $K \in \mathbb{N}$ bezeichne

$$\varepsilon_{VQ}^{(K)}(\omega) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{z}_1, \dots, \mathbf{z}_K} \sum_{t=1}^T \min_{\kappa} \|\mathbf{x}_t - \mathbf{z}_{\kappa}\|^2$$

die minimale **quadratisch-euklidische Verzerrung** der K -Partition und

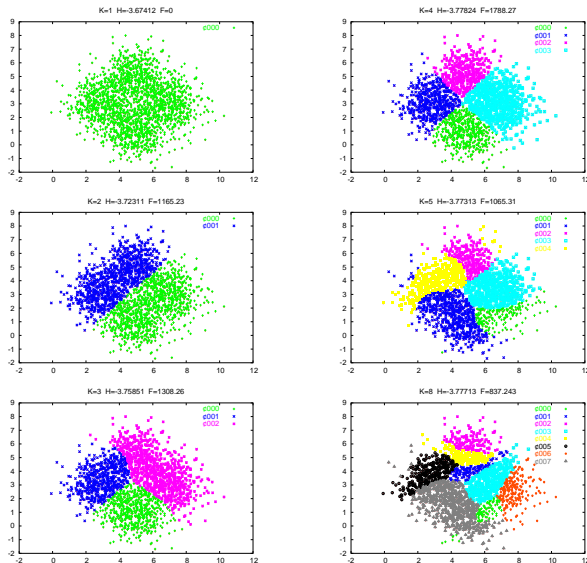
$$\ell_{GMM}^{(K)}(\omega) \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \mathcal{M}_K} \sum_{t=1}^T \log \sum_{\kappa=1}^K \left\{ \pi_{\kappa}^{(\theta)} \cdot \mathcal{N}(\mathbf{x}_t \mid \mu_{\kappa}^{(\theta)}, \mathbf{S}_{\kappa}^{(\theta)}) \right\}$$

die maximale Güte eines K -**Mischverteilungsmodells** für den Datensatz. Dann gelten die Antitonie und die Monotonie

$$K \leq K' \Rightarrow \begin{cases} \varepsilon_{VQ}^{(K)}(\omega) & \geq \varepsilon_{VQ}^{(K')}(\omega) \\ \ell_{GMM}^{(K)}(\omega) & \leq \ell_{GMM}^{(K')}(\omega) \end{cases}.$$

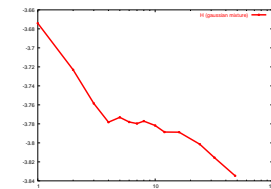
Die beste Anzahl von Gruppen

kann oft durch den Pseudo- F -Wert ermittelt werden

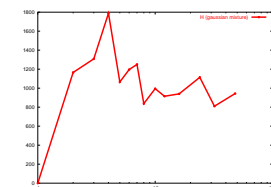


#2 — Raute⁴

Vier sphärische Gaußdichten mit Varianz $\sigma = 1.0$



Negatives $\ell_{GMM}^{(K)}(\text{EM})$



Pseudo- F -Wert

Der Pseudo- F -Wert

belohnt gute Gruppentrennung & bestraft große Gruppenanzahl

Problem

Die allermeisten Gruppierungskriterien **verbessern sich systematisch** mit wachsender Gruppenzahl K , sind also zur Auswahl der Gruppenzahl völlig ungeeignet.

Bemerkung

Die Monotonie ist in den meisten Kurven verletzt, denn K -means und GMM-Identifikation sind **lokale** Optimierungsverfahren.

Definition (Calinski-Harabasz)

Die Vergleichsgröße

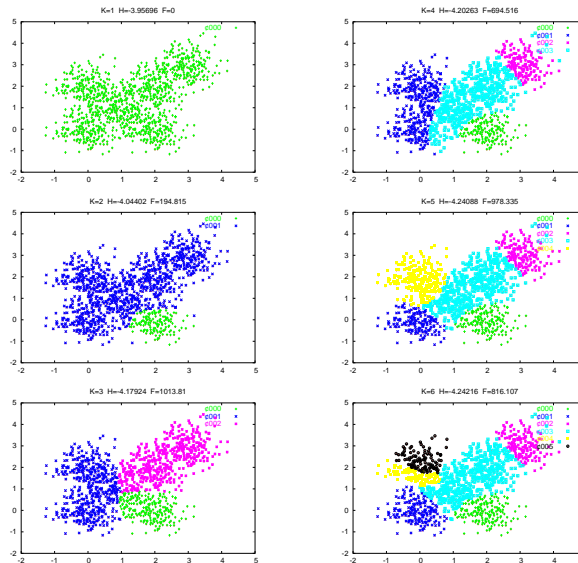
$$C_{\text{pseudo}}(\{\omega_1, \dots, \omega_K\}) \stackrel{\text{def}}{=} \frac{\text{spur}(\mathbf{S}_B) / (K - 1)}{\text{spur}(\mathbf{S}_W) / (T - K)}$$

heißt **Pseudo- F -Wert** der Gruppierung $\{\omega_1, \dots, \omega_K\}$.

Es bezeichnen \mathbf{S}_B die **Zwischengruppenstreuung** und \mathbf{S}_W die **Innergruppenstreuung** der Datenpartition.

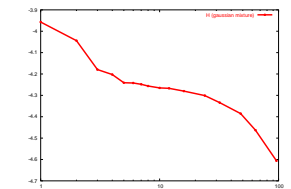
Hier versagt der Pseudo- F -Wert !

Der EM-Algorithmus zur Mischungsidentifikation findet nur suboptimale GMM

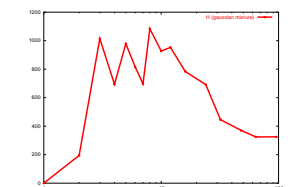


#3 - Kreuz⁶

Sechs Dichten in schrägliegender Kreuzform



Negatives $\ell_{GMM}^{(K)}(\text{EM})$

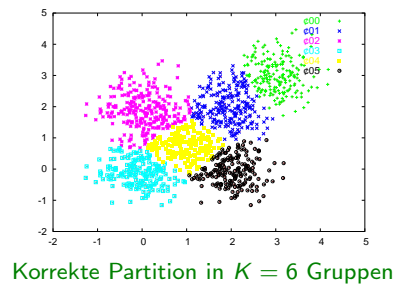
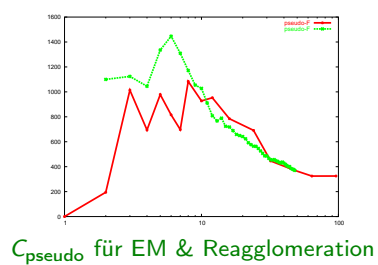


Pseudo- F -Wert

Dreiphasige Gruppierung

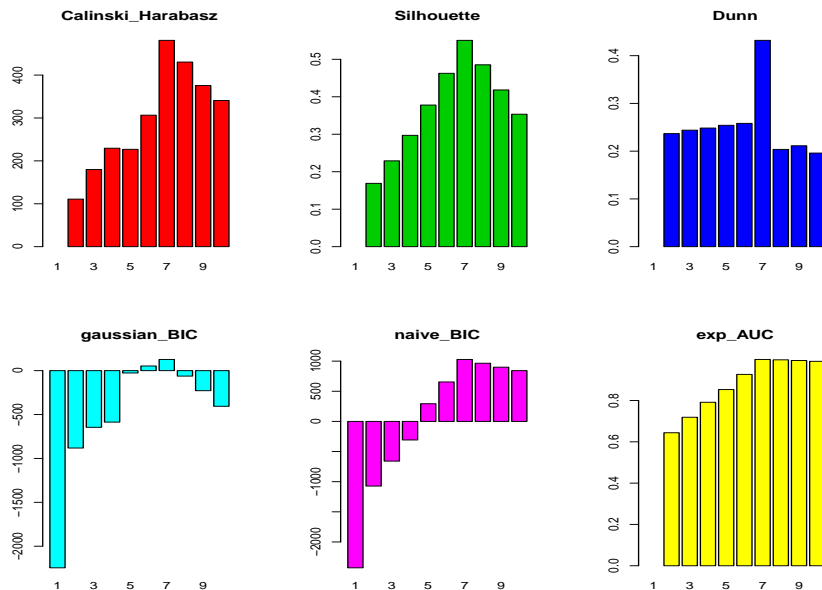
(Algorithmus)

- 1 DIVISIVE GRUPPIERUNG
Sukzessive Zerlegung mit 2-means in $K_{\max} = 2^b$ Gruppen.
- 2 GMM-IDENTIFIKATION
Austauschiteration mit Gaußschem Mischverteilungsmodell (EM-Schritte).
- 3 REAGGLOMERATION
Bottom-up Gruppierung durch sukzessive ℓ_{GMM} -Maximierung.



Beispiel: regularisierte Gütemaße im Vergleich

Sieben Cluster im \mathbb{R}^7 · Hierarchische Gruppierung für $K = 1, 2, \dots, 10$



Regularisierte Gütemaße für scharfe Gruppierungen

„cluster validity index“

Dunn's ISODATA

$$C_{\text{iso}} \stackrel{\text{def}}{=} \frac{\min_{\kappa \neq \lambda} \min_{x \in \omega_{\kappa}} \min_{y \in \omega_{\lambda}} d(x, y)}{\max_{\kappa} \max_{x \in \omega_{\kappa}} \max_{y \in \omega_{\kappa}} d(x, y)}$$

Rousseeuw's Silhouette

$$C_{\text{sil}} \stackrel{\text{def}}{=} \frac{1}{T} \sum_{t=1}^T \frac{b_t - a_t}{\max(a_t, b_t)} \quad \text{mit } D_{\kappa, t} = \overline{d(\omega_{\kappa}, x_t)} \text{ und } \begin{cases} a_t = D_{\kappa(t), t} \\ b_t = \min_{\lambda \neq \kappa(t)} D_{\lambda, t} \end{cases}$$

Expected Area Under Curve

$$C_{\text{auc}} \stackrel{\text{def}}{=} \frac{1}{T} \sum_{t=1}^T \text{AUC}([d(x_{\bullet}, x_t)], [x_{\bullet} \in \omega_{\kappa(t)}])$$

Uni-/multivariat gaußsches BIC

$$C_{\text{bic}} \stackrel{\text{def}}{=} -\log \ell_{\text{GMM}}^{(K)} + \log(T) \cdot \text{df}(K, N) \quad \text{mit } \text{df}(K, N) = \begin{cases} K + 2NK \\ K + (N + 3) \frac{NK}{2} \end{cases} \quad \begin{matrix} \text{(naiv)} \\ \text{(sonst)} \end{matrix}$$

Gütemaße für unscharfe Gruppierungen

Tendenz zur monotonen Verbesserung mit der Gruppenanzahl K

Partitionskoeffizient

$$C_{\text{part}} \stackrel{\text{def}}{=} \frac{1}{T} \sum_{t=1}^T \sum_{\kappa=1}^K u_{\kappa}^2(x_t) \xrightarrow{!} \max$$

Proportionsexponent

$$C_{\text{prop}} \stackrel{\text{def}}{=} \frac{1}{T} \sum_{t=1}^T \max_{\kappa} u_{\kappa}(x_t) \xrightarrow{!} \max$$

Klassifikationsentropie

$$C_{\text{entro}} \stackrel{\text{def}}{=} -\frac{1}{T} \sum_{t=1}^T \sum_{\kappa=1}^K u_{\kappa}(x_t) \cdot \log_2 u_{\kappa}(x_t) \xrightarrow{!} \min$$

Bemerkung

Partitionskoeffizient C_{part} , Proportionsexponent C_{prop} und Klassifikationsentropie C_{entro} nehmen ihre Optimalwerte ($1/1/0$) für die **scharfen** Gruppierungen $\{u_{\kappa}(\cdot)\}$ an.

Mischungsidentifikation in \mathbb{R}

Ein Zoo konkurrierender Modelle — Kovarianzauslegung $\mathbf{S}_\kappa := \sigma_\kappa^2 \cdot \mathbf{U}_\kappa \mathbf{D}_\kappa \mathbf{U}_\kappa^\top$

Sphärische Modelle

EII $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma^2 \cdot \mathbf{E})$

VII $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma_\kappa^2 \cdot \mathbf{E})$

Global oder κ -variabel?

Volumen, Gesamtstreuung

Gestalt, Dynamik

Orientierung, Hauptachsen

σ_κ^2

\mathbf{D}_κ

\mathbf{U}_κ

Diagonale Modelle

EEI $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma^2 \cdot \mathbf{D})$

VEI $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma_\kappa^2 \cdot \mathbf{D})$

EVI $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma^2 \cdot \mathbf{D}_\kappa)$

VVI $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma_\kappa^2 \cdot \mathbf{D}_\kappa)$

Ellipsoidale Modelle

EEE $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma^2 \cdot \mathbf{U} \mathbf{D} \mathbf{U}^\top)$

EEV $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma^2 \cdot \mathbf{U}_\kappa \mathbf{D} \mathbf{U}_\kappa^\top)$

VEV $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma_\kappa^2 \cdot \mathbf{U}_\kappa \mathbf{D} \mathbf{U}_\kappa^\top)$

VVV $\mathcal{N}(\boldsymbol{\mu}_\kappa, \sigma_\kappa^2 \cdot \mathbf{U}_\kappa \mathbf{D}_\kappa \mathbf{U}_\kappa^\top)$

Hierarchische Gruppierung: agglomerativ/divisiv

Austauschverfahren: (fuzzy) K-means

Mischungsidentifikation

Relationale Gruppierung

Konzeptuelle Gruppierung

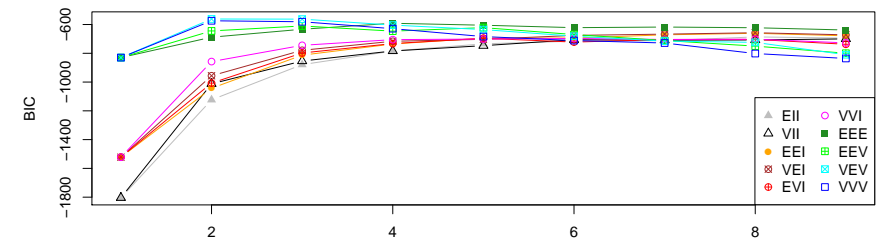
Spektrale Gruppierung

Clustergütemaße

Zusammenfassung

Mischungsidentifikation in \mathbb{R}

BIC-Entscheidung: optimales GM-Modell & optimale Gruppenzahl K



Agglomeratives Gruppieren

`hc (modelName, data, ...)`

'R'-Paket: *mclust*

EM-Iteration

`em (modelName, data, parameters, ...)`

Start mit E-Schritt

`me (modelName, data, z, ...)`

Start mit M-Schritt

Bayes-Informationsmaß

`Mclust (data, G=1:9, modelNames)` alle $K \in \{1, \dots, 9\}$, alle GMM-Typen

Zusammenfassung (5)

1. Das Ziel der **Gruppierung (Clusteranalyse)** ist die **unüberwachte** Zerlegung eines Datensatzes in **explizit** oder **implizit** charakterisierte Teilmengen von Objekten.
2. Die **hierarchischen** Verfahren arbeiten *bottom-up* (**agglomerativ**) oder *top-down* (**divisiv**); das Resultat ist ein **Grupp dendrogramm**.
3. Die **Austauschverfahren** geben eine **Anzahl** $K \in \mathbb{N}$ vor und erzeugen **scharfe** (*K-means*) oder **unscharfe** (*fuzzy K-means*) K -Partitionen.
4. Die **EM-Gruppierung** modelliert die Daten durch Identifikation einer gaußschen **Mischverteilung**.
5. Neben **sphärischen** Gruppen lassen sich auch **rangdefiziente** Ballungsgebiete ermitteln (*K-Elliptotypes* oder **Probabilistic PCA**).
6. **Relationale** Datensätze werden entweder *agglomerativ* gruppiert oder — wie auch **nominal** skalierte Objekte — mit einem **K-medoids**-Austauschverfahren (*RACE*, *K-Sterne*).
7. Die **spektrale** Methode löst eine **Minimalschnittaufgabe** im Affinitätsgraphen der Datenobjekte und läuft auf eine Art gaußschen **Kernel-K-means**-Algorithmus hinaus.
8. Ermittlung der **Clusteranzahl** durch **Ellenbogenheuristik** oder **Pseudo-F-Wert**.