

Kapitel ML:IX

IX. Clusteranalyse

- ❑ Einordnung Data Mining
- ❑ Einführung in die Clusteranalyse
- ❑ Hierarchische Verfahren
- ❑ Iterative Verfahren
- ❑ Dichtebaserte Verfahren
- ❑ Cluster Evaluation
- ❑ Constrained Cluster Analysis

ML:IX-1 Cluster Analysis

©STEIN 2002-2012

Einordnung Data Mining

Definition 1 (Data Mining)

Unter Data Mining versteht man das systematische, in der Regel automatisierte oder halbautomatische Entdecken und Extrahieren bislang unbekannter Zusammenhänge aus großen Mengen von Daten.

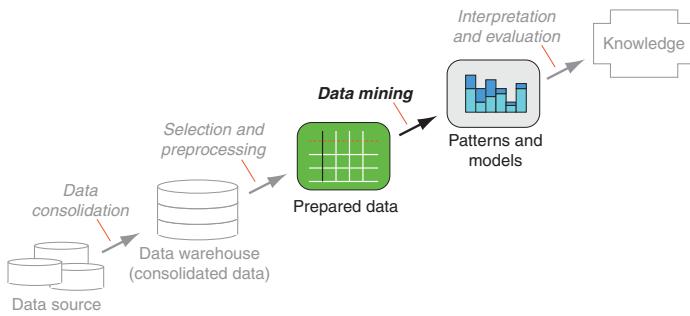
Data Mining umfasst folgende Schritte:

1. Aufgabendefinition
2. Datenselektion
3. Datenvorbereitung und -transformation
4. Mustererkennung
5. Kommunikation, Präsentation

ML:IX-2 Cluster Analysis

©STEIN 2002-2012

Einordnung Data Mining



Definition 2 (Knowledge Discovery in Databases, KDD)

Wissensentdeckung in Datenbanken (*Knowledge Discovery in Databases*) ist der nichttriviale Prozess der Identifikation gültiger, neuer, potentiell nützlicher und schlussendlich verständlicher Muster in großen Datenbeständen.

[vgl. Fayyad 1996, Wrobel 1998]

ML:IX-3 Cluster Analysis

©STEIN 2002-2012

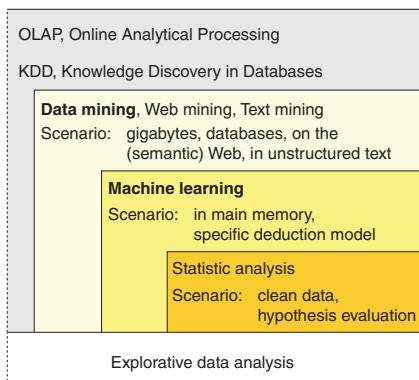
Bemerkungen:

- ❑ Data-Mining-Techniken werden der *explorativen Datenanalyse* zugeordnet. Ziel der explorativen Datenanalyse ist – über die Darstellung der Daten hinaus – die Suche nach Strukturen und Besonderheiten. Die explorative Datenanalyse wird eingesetzt, wenn die Fragestellung oder die Wahl eines geeigneten statistischen Modells unklar ist.
- ❑ In der Data-Mining-Definition wird auf den Begriff der Information verzichtet: Data Mining wird der signifikanten Ebene der Semiotik zugeordnet. Die *Interpretation* der entdeckten Muster, also die im Rahmen der explorativen Datenanalyse stattfindende Auseinandersetzung mit Informationen im Sinne eines subjektiven Wissenszuwachses, die auf der pragmatischen Ebene abläuft, gehört in das Gebiet des Knowledge Discovery in Databases, KDD.
- ❑ Vor allem im kommerziellen Bereich wird der Begriff des Data Mining synonym zu Wissensentdeckung in Datenbanken (KDD) verwendet. Data Mining ist aber nur ein Teilschritt innerhalb des KDD-Prozesses, nämlich der Analyseschritt zur Mustererkennung.
- ❑ Unter Web Mining versteht man die Übertragung von Techniken des Data Mining zur (teil)automatischen Extraktion von Informationen aus dem Internet, speziell dem World Wide Web.
- ❑ Mit Text Mining wird die Entdeckung neuer und für den Benutzer relevanter Informationen aus Textdaten bezeichnet.

ML:IX-4 Cluster Analysis

©STEIN 2002-2012

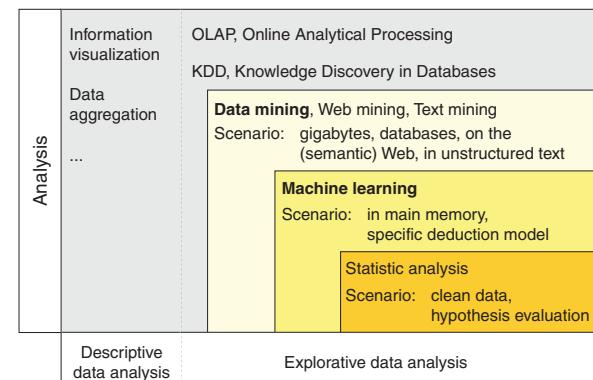
Einordnung Data Mining



ML:IX-5 Cluster Analysis

©STEIN 2002-2012

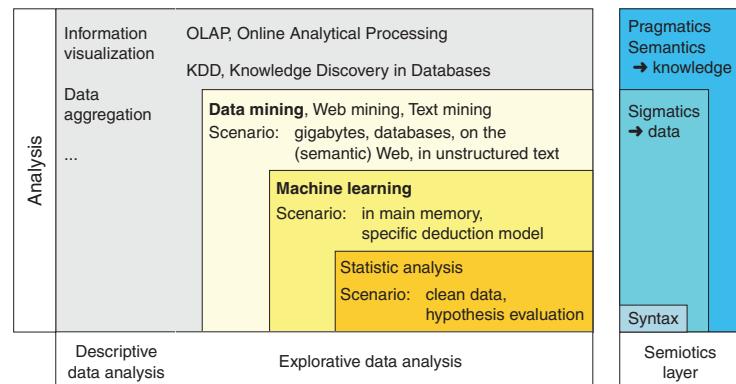
Einordnung Data Mining



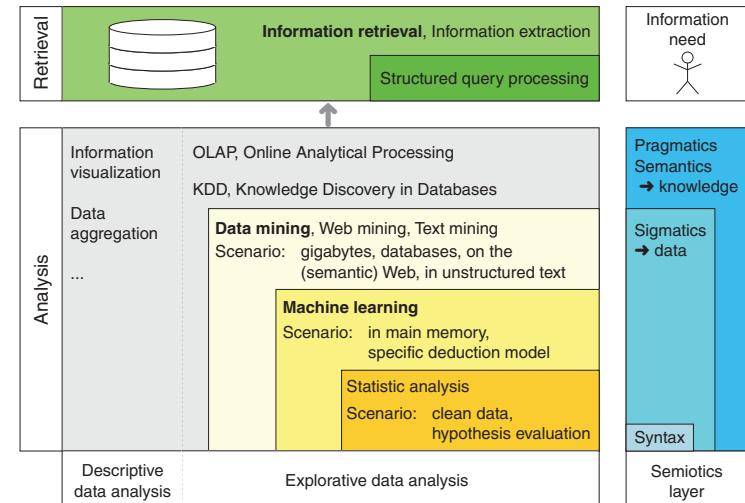
ML:IX-6 Cluster Analysis

©STEIN 2002-2012

Einordnung Data Mining



Einordnung Data Mining



Bemerkungen:

- Eine klare Abgrenzung von Machine Learning und Data Mining ist nicht immer möglich; ein wichtiger Unterschied resultiert aus der Größe der behandelten Datenmengen: Anwendungen des Machine Learning laufen typischerweise im Hauptspeicher ab; die Disziplin des Data Mining entstand aus der Notwendigkeit, maschinelle Analyseverfahren auf große Datenbanken anzuwenden.
- Ein Schwerpunkt des Machine Learning ist der eigentliche Lern- bzw. Deduktionsprozess wie die Theorie des analogen Schließens, das Lernen aus Beispielen, oder der Einfluss der Umwelt auf das Lernen. Hingegen ist die treibende Kraft hinter Data Mining die Industrie- und Geschäftswelt mit ihren großen Datenbanken.
- Zu den bekannten Aufgabenstellung des Data Mining gehören: ungerichtete Abhängigkeitsanalyse zur Identifikation signifikanter Dependenzen zwischen den Attributen eines Informationsobjektes (Beispiel: Warenkorbanalyse), Gruppenbildung und Klassifikationsprobleme, Filtern von Prozessdaten, Prognoseaufgaben.

ML: IX-9 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Einordnung Data Mining

Methoden und Techniken

- Clusteranalyse
- propositionale (oder relative) Regellernverfahren
- assoziative Regellernverfahren
- Hauptkomponenten- und Faktorenanalyse

Kapitel ML:IX (Fortsetzung)

IX. Clusteranalyse

- Einordnung Data Mining
- Einführung in die Clusteranalyse
- Hierarchische Verfahren
- Iterative Verfahren
- Dichtebasierter Verfahren
- Cluster Evaluation
- Constrained Cluster Analysis

ML: IX-11 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Einführung in die Clusteranalyse

Clusteranalyse ist die **unüberwachte** Klassifikation einer Menge von Objekten in Gruppen; dabei wird folgendes Ziel verfolgt:

1. Ähnlichkeit innerhalb der Gruppen maximieren
2. Ähnlichkeit zwischen den Gruppen minimieren

ML: IX-12 Cluster Analysis

©STEIN 2002-2012

Einführung in die Clusteranalyse

Clusteranalyse ist die **unüberwachte** Klassifikation einer Menge von Objekten in Gruppen; dabei wird folgendes Ziel verfolgt:

1. Ähnlichkeit innerhalb der Gruppen maximieren
2. Ähnlichkeit zwischen den Gruppen minimieren

Anwendungen

- ❑ Identifikation gleichartiger Käufergruppen
- ❑ „höhere“ Bildverarbeitung, im Sinne von Objekterkennung
- ❑ Suche nach ähnlichen Genprofilen
- ❑ Spezifikation von Syndromen
- ❑ Analyse von Verkehrsdaten in Computernetzen
- ❑ Visualisierung komplexer Graphen
- ❑ Textkategorisierung im Information-Retrieval

Bemerkungen:

- ❑ Die Problemstellung der Clusteranalyse ist umgekehrt zur Problemstellung der Varianzanalyse. Ziel der Varianzanalyse ist die Überprüfung, ob eine gegebene nominalskalierte Variable Gruppen definiert, deren Mitglieder sich in abhängigen (intervallskalierten) Variablen unterscheiden. Ziel der Clusteranalyse ist die Erzeugung einer solchen nominalskalierten Variable durch die Entdeckung der Ausprägungen dieser Variable. Jedes Cluster korrespondiert zu einer Variablenausprägung.
- ❑ Die Clusteranalyse ist ein Verfahren zur Strukturerzeugung: Man weiß quasi nichts über die zu entdeckende Variable, insbesondere nichts über die Anzahl ihrer Werteausprägungen. Die Varianzanalyse ist ein Verfahren zur Strukturprüfung.

Einführung in die Clusteranalyse

x_1, \dots, x_n sind die zu n Objekten gehörenden p -dimensionalen Merkmalsvektoren:

	Merkmal 1	Merkmal 2	...	Merkmal p
x_1	x_{11}	x_{12}	...	x_{1p}
x_2	x_{21}	x_{22}	...	x_{2p}
:				
x_n	x_{n1}	x_{n2}	...	x_{np}

Einführung in die Clusteranalyse

x_1, \dots, x_n sind die zu n Objekten gehörenden p -dimensionalen Merkmalsvektoren:

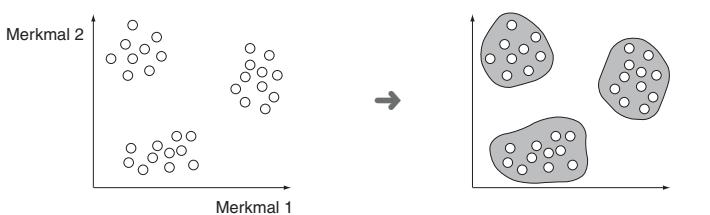
	Merkmal 1	Merkmal 2	...	Merkmal p	kein Zielkonzept
x_1	x_{11}	x_{12}	...	x_{1p}	c_1
x_2	x_{21}	x_{22}	...	x_{2p}	c_2
:					:
x_n	x_{n1}	x_{n2}	...	x_{np}	c_n

Einführung in die Clusteranalyse

x_1, \dots, x_n sind die zu n Objekten gehörenden p -dimensionalen Merkmalsvektoren:

	Merkmal 1	Merkmal 2	...	Merkmal p	kein Zielkonzept
x_1	x_{11}	x_{12}	...	x_{1p}	c_1
x_2	x_{21}	x_{22}	...	x_{2p}	c_2
:					:
x_n	x_{n1}	x_{n2}	...	x_{np}	c_n

30 zweidimensionale Merkmalsvektoren ($n = 30, p = 2$):



Einführung in die Clusteranalyse

Definition 2 (exklusives Clustering) [vgl. Splitting]

Sei X eine Menge von Merkmalsvektoren. Ein exklusives Clustering \mathcal{C} von X , $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$, $C_i \subseteq X$, ist eine Aufteilung von X in nichtleere und paarweise disjunkte Teilmengen C_i mit $\bigcup_{C_i \in \mathcal{C}} C_i = X$.

Einführung in die Clusteranalyse

Definition 2 (exklusives Clustering) [vgl. Splitting]

Sei X eine Menge von Merkmalsvektoren. Ein exklusives Clustering \mathcal{C} von X , $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$, $C_i \subseteq X$, ist eine Aufteilung von X in nichtleere und paarweise disjunkte Teilmengen C_i mit $\bigcup_{C_i \in \mathcal{C}} C_i = X$.

Algorithmen zur Clusteranalyse sind unüberwachte Lernverfahren:

- der Lernprozess ist selbstorganisiert
- es gibt keinen externen Lehrer
- es gibt ein aufgabenunabhängiges Optimierungskriterium

Einführung in die Clusteranalyse

Definition 2 (exklusives Clustering) [vgl. Splitting]

Sei X eine Menge von Merkmalsvektoren. Ein exklusives Clustering \mathcal{C} von X , $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$, $C_i \subseteq X$, ist eine Aufteilung von X in nichtleere und paarweise disjunkte Teilmengen C_i mit $\bigcup_{C_i \in \mathcal{C}} C_i = X$.

Algorithmen zur Clusteranalyse sind unüberwachte Lernverfahren:

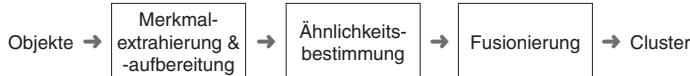
- der Lernprozess ist selbstorganisiert
- es gibt keinen externen Lehrer
- es gibt ein aufgabenunabhängiges Optimierungskriterium

Überwachtes Lernen:

- es gibt Lernziele: das Zielkonzept, gewünschte Reaktionen, etc.
- es gibt ein aufgabenabhängiges Optimierungskriterium
- es gibt Information darüber, wie eine Verbesserung im Optimierungskriterium zu erzielen ist. Stichwort: Instructive Feedback

Einführung in die Clusteranalyse

Grundsritte einer Clusteranalyse



Einführung in die Clusteranalyse

Merkmalextrahierung und -aufbereitung

Gesucht sind evtl. neue Merkmale mit hoher Varianz. Techniken:

- Analyse von Streuungsparametern
- Dimensionsreduktion: Faktorenanalyse, multidimensionale Skalierung
- Visuelle Analyse: Scatter-Plots, Box-plots

Einführung in die Clusteranalyse

Merkmalextrahierung und -aufbereitung

Gesucht sind evtl. neue Merkmale mit hoher Varianz. Techniken:

- Analyse von Streuungsparametern
- Dimensionsreduktion: Faktorenanalyse, multidimensionale Skalierung
- Visuelle Analyse: Scatter-Plots, Box-plots

Standardisierung von Variablen (Merkmälen) ist problematisch:



Einführung in die Clusteranalyse

Berechnung von Distanzen oder Ähnlichkeiten

	Merkmal 1	Merkmal 2	...	Merkmal p
x_1	x_{11}	x_{12}	...	x_{1p}
x_2	x_{21}	x_{22}	...	x_{2p}
:				
x_n	x_{n1}	x_{n2}	...	x_{np}

	x_1	x_2	...	x_n
x_1	0	$d(x_1, x_2)$...	$d(x_1, x_n)$
x_2	-	0	...	$d(x_2, x_n)$
:				
x_n	-	-	...	0

Bemerkungen:

- Die Distanzmatrix ist oft implizit durch eine Metrik auf dem Merkmalsraum definiert.
- Die Distanzmatrix kann als Adjazenzmatrix eines gewichteten, ungerichteten Graphen G , $G = \langle V, E, w \rangle$, interpretiert werden: Die Menge X der Merkmalsvektoren wird bijektiv auf eine Knotenmenge V abgebildet; eine Distanz $d(x_i, x_j)$ entspricht dem Gewicht $w(\{u, v\})$ der Kante $\{u, v\} \in E$ zwischen den mit x_i und x_j assoziierten Knoten u und v .

ML-IX-25 Cluster Analysis

©STEIN 2002-2012

Einführung in die Clusteranalyse

Berechnung von Distanzen oder Ähnlichkeiten (Fortsetzung)

Anforderungen an eine Distanzfunktion:

1. $d(x_1, x_2) \geq 0$
2. $d(x_1, x_1) = 0$
3. $d(x_1, x_2) = d(x_2, x_1)$
4. $d(x_1, x_3) \leq d(x_1, x_2) + d(x_2, x_3)$

ML-IX-26 Cluster Analysis

©STEIN 2002-2012

Einführung in die Clusteranalyse

Berechnung von Distanzen oder Ähnlichkeiten (Fortsetzung)

Anforderungen an eine Distanzfunktion:

1. $d(x_1, x_2) \geq 0$
2. $d(x_1, x_1) = 0$
3. $d(x_1, x_2) = d(x_2, x_1)$
4. $d(x_1, x_3) \leq d(x_1, x_2) + d(x_2, x_3)$

Bei intervallskalierten Variablen Verwendung der Minkowsky-Metrik:

$$d(x_1, x_2) = \left(\sum_{i=1}^p |x_{1i} - x_{2i}|^r \right)^{1/r}$$

mit

- $r = 1$. Manhattan- oder Hamming-Distanz, L_1 -Norm
- $r = 2$. Euklidische Distanz, L_2 -Norm
- $r = \infty$. Maximum-Distanz, L_∞ -Norm bzw. L_{\max} -Norm

ML-IX-27 Cluster Analysis

©STEIN 2002-2012

Einführung in die Clusteranalyse

Berechnung von Distanzen oder Ähnlichkeiten (Fortsetzung)

Eine Clusteranalyse verlangt kein spezielles Skalenniveau der Merkmale.

- Verallgemeinerung der Distanzfunktion zur (Un)Ähnlichkeitsfunktion durch Verzicht auf die Dreiecksungleichung. (Un)Ähnlichkeiten lassen sich zwischen binären, nominalen und ordinalen Variablen quantifizieren.

ML-IX-28 Cluster Analysis

©STEIN 2002-2012

Einführung in die Clusteranalyse

Berechnung von Distanzen oder Ähnlichkeiten (Fortsetzung)

Eine Clusteranalyse verlangt kein spezielles Skalenniveau der Merkmale.

- Verallgemeinerung der Distanzfunktion zur (Un)Ähnlichkeitsfunktion durch Verzicht auf die Dreiecksungleichung. (Un)Ähnlichkeiten lassen sich zwischen binären, nominalen und ordinalen Variablen quantifizieren.

Ähnlichkeitskoeffizienten für zwei Vektoren, x_1 , x_2 , mit binären Merkmalen:

$$\text{Simple Matching Coefficient (SMC)} = \frac{f_{11} + f_{00}}{f_{11} + f_{00} + f_{01} + f_{10}}$$

$$\text{Jaccard-Koeffizient (J)} = \frac{f_{11}}{f_{11} + f_{01} + f_{10}}$$

mit

f_{11} = Anzahl der Merkmale mit Ausprägung 1 sowohl in x_1 als auch in x_2

f_{00} = Anzahl der Merkmale mit Ausprägung 0 sowohl in x_1 als auch in x_2

f_{01} = Anzahl der Merkmale mit Ausprägung 0 in x_1 und Ausprägung 1 in x_2

f_{10} = Anzahl der Merkmale mit Ausprägung 1 in x_1 und Ausprägung 0 in x_2

ML-IX-29 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

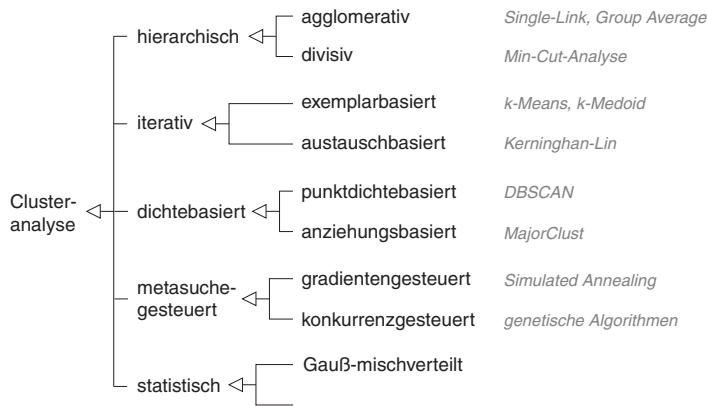
- Die Definition der Ähnlichkeitskoeffizienten lässt sich auf nominale Variablen erweitern.
- Heterogene Metriken wie HEOM und HVDM ermöglichen die kombinierte Verrechnung verschiedener Skalenniveaus.
- Die Berechnung von Korrelationskoeffizienten zwischen zwei Vektoren über alle Merkmale (also nicht zwischen zwei Merkmalen über alle Vektoren) ermöglicht den Vergleich von Profilen. Beispiel: Q-Korrelationskoeffizient
- Die Konstruktion eines geeigneten Ähnlichkeitsmaßes stellt oft die größte Herausforderung bei der Cluster-Analyse dar. Typische Problemfelder:
 - Normalisierung
 - Empfindlichkeit bei Ausreißern
 - Korrelationen zwischen Merkmalen
 - unterschiedliche Wichtigkeit der Merkmale
- Ähnlichkeitsmaße lassen sich kanonisch in Unähnlichkeitsmaße umrechnen – und umgekehrt.

ML-IX-30 Cluster Analysis

©STEIN 2002-2012

Einführung in die Clusteranalyse

Prinzipien der Fusionierung



ML-IX-31 Cluster Analysis

©STEIN 2002-2012

Kapitel ML: IX (Fortsetzung)

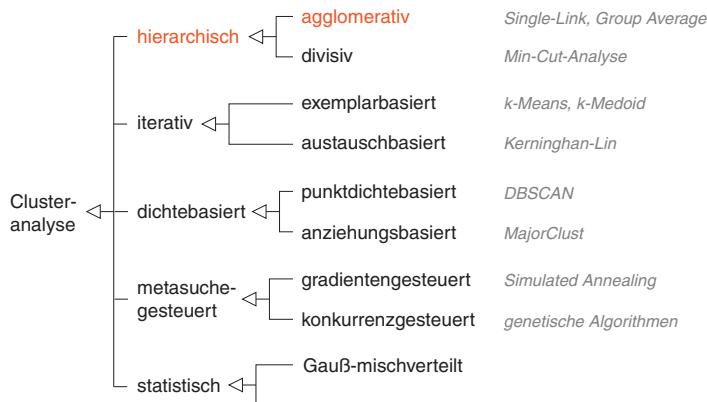
IX. Clusteranalyse

- Einordnung Data Mining
- Einführung in die Clusteranalyse
- Hierarchische Verfahren
- Iterative Verfahren
- Dichtebaserte Verfahren
- Cluster Evaluation
- Constrained Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Prinzipien der Fusionierung



ML-IX-33 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Algorithmus zur hierarchisch-agglomerativen Clusteranalyse

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d_c . Distance measure between two clusters.

Output: $T = \langle V_T, E_T \rangle$. Cluster hierarchy or dendrogram.

1. $\mathcal{C} = \{ \{v\} \mid v \in V \}$ // initial clustering
- 2.
3. **WHILE** $|\mathcal{C}| > 1$ **DO**
4. *update_distance_matrix*(\mathcal{C}, G, d_c)
5. $\{C, C'\} = \underset{\{C_i, C_j\} \in \mathcal{C}: C_i \neq C_j}{\operatorname{argmin}} d_c(C_i, C_j)$
6. $\mathcal{C} = (\mathcal{C} \setminus \{C, C'\}) \cup \{C \cup C'\}$ // merging
- 7.
8. **ENDDO**
9. **RETURN**(T)

Vergleiche hierzu den Algorithmus zur hierarchisch-divisiven Clusteranalyse.

©STEIN 2002-2012

Hierarchische Verfahren

Algorithmus zur hierarchisch-agglomerativen Clusteranalyse

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d_c . Distance measure between two clusters.

Output: $T = \langle V_T, E_T \rangle$. Cluster hierarchy or dendrogram.

1. $\mathcal{C} = \{ \{v\} \mid v \in V \}$ // initial clustering
2. $V_T = \{v_C \mid C \in \mathcal{C}\}, E_T = \emptyset$ // initial dendrogram
3. **WHILE** $|\mathcal{C}| > 1$ **DO**
4. *update_distance_matrix*(\mathcal{C}, G, d_c)
5. $\{C, C'\} = \underset{\{C_i, C_j\} \in \mathcal{C}: C_i \neq C_j}{\operatorname{argmin}} d_c(C_i, C_j)$
6. $\mathcal{C} = (\mathcal{C} \setminus \{C, C'\}) \cup \{C \cup C'\}$ // merging
7. $V_T = V_T \cup \{v_{C,C'}\}, E_T = E_T \cup \{\{v_{C,C'}, v_C\}, \{v_{C,C'}, v_{C'}\}\}$ // dendrogram
8. **ENDDO**
9. **RETURN**(T)

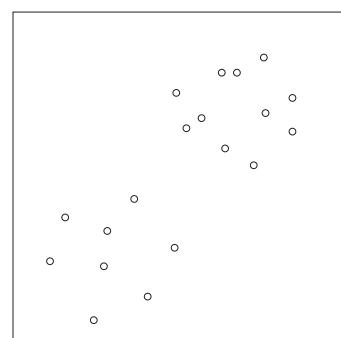
Vergleiche hierzu den Algorithmus zur hierarchisch-divisiven Clusteranalyse.

ML-IX-35 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

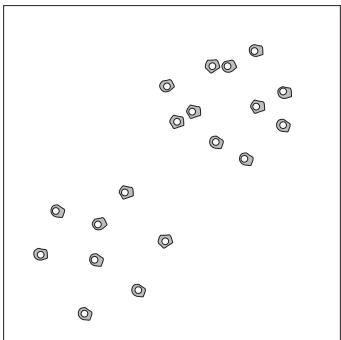
Single-Link: Cluster-Distanzmaß d_c = Nearest-Neighbor



©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor

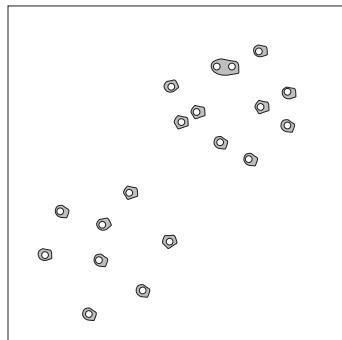


ML-IX-37 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor



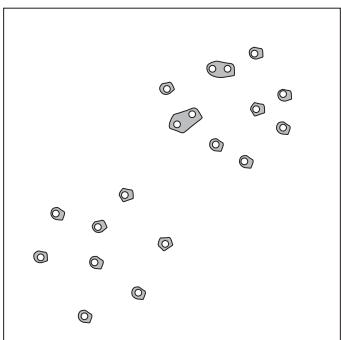
→ Distanz

ML-IX-38 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor



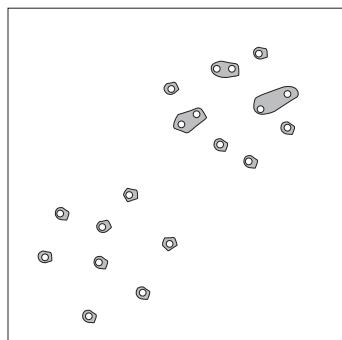
→ Distanz

ML-IX-39 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor



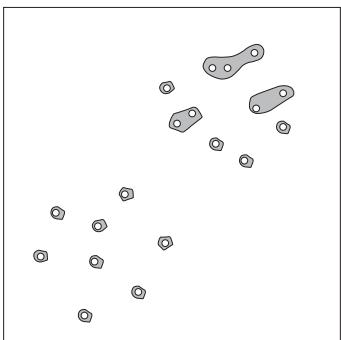
→ Distanz

ML-IX-40 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor



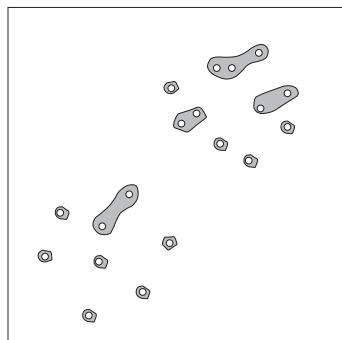
→ Distanz

ML-IX-41 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor



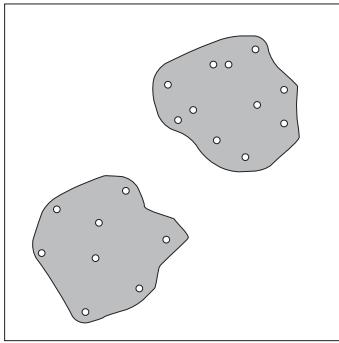
→ Distanz

ML-IX-42 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor

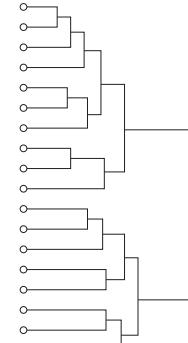
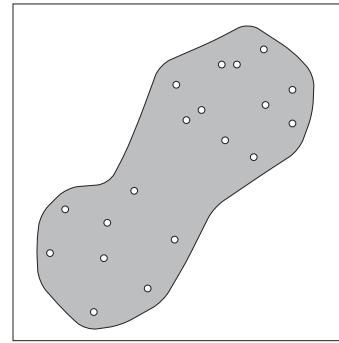


ML-IX-43 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Single-Link: Cluster-Distanzmaß d_C = Nearest-Neighbor



ML-IX-44 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Distanzmaße hierarchisch-agglomerativer Verfahren [\[Eigenschaften\]](#)

$$d_C(C, C') = \min_{v \in C} d(u, v)$$

Single-Link
(Nearest-Neighbor)

$$d_C(C, C') = \max_{u \in C, v \in C'} d(u, v)$$

Complete-Link
(Furthest-Neighbor)

$$d_C(C, C') = \frac{1}{|C| \cdot |C'|} \sum_{u \in C, v \in C'} d(u, v)$$

(Group-)Average-Link

$$d_C(C, C') = \sqrt{\frac{2 \cdot |C| \cdot |C'|}{|C| + |C'|} \cdot ||\bar{u} - \bar{v}||}$$

Ward (Varianz)

Siehe Verwendung im Algorithmus:

- [hierarchisch-agglomerative Clusteranalyse](#)
- [hierarchisch-divisive Clusteranalyse](#)

ML-IX-45 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Ward-Kriterium

Ward ist ein Varianzkriterium; es entspricht der doppelten Zunahme der Wurzel aus der Fehlerquadratsumme, SSE , in dem neuen Cluster, der durch die Vereinigung der beiden Cluster C und C' entsteht. Herleitung:

$$SSE(C) = \sum_{u \in C} ||\bar{u} - u||^2$$

©STEIN 2002-2012

Hierarchische Verfahren

Ward-Kriterium

Ward ist ein Varianzkriterium; es entspricht der doppelten Zunahme der Wurzel aus der Fehlerquadratsumme, SSE , in dem neuen Cluster, der durch die Vereinigung der beiden Cluster C und C' entsteht. Herleitung:

$$\begin{aligned} SSE(C) &= \sum_{u \in C} ||\bar{u} - u||^2 = \sum_{u \in C} (||\bar{u}||^2 - 2 \cdot \langle u, \bar{u} \rangle + ||u||^2) \\ &= |C| \cdot ||\bar{u}||^2 - 2|C| \cdot ||\bar{u}||^2 + \sum_{u \in C} ||u||^2 = \sum_{u \in C} ||u||^2 - |C| \cdot ||\bar{u}||^2 \end{aligned}$$

Hierarchische Verfahren

Ward-Kriterium

Ward ist ein Varianzkriterium; es entspricht der doppelten Zunahme der Wurzel aus der Fehlerquadratsumme, SSE , in dem neuen Cluster, der durch die Vereinigung der beiden Cluster C und C' entsteht. Herleitung:

$$\begin{aligned} SSE(C) &= \sum_{u \in C} ||\bar{u} - u||^2 = \sum_{u \in C} (||\bar{u}||^2 - 2 \cdot \langle u, \bar{u} \rangle + ||u||^2) \\ &= |C| \cdot ||\bar{u}||^2 - 2|C| \cdot ||\bar{u}||^2 + \sum_{u \in C} ||u||^2 = \sum_{u \in C} ||u||^2 - |C| \cdot ||\bar{u}||^2 \end{aligned}$$

$$SSE(C') = \sum_{v \in C'} ||\bar{v} - v||^2 - |C'| \cdot ||\bar{v}||^2$$

$$SSE(C \cup C') = \sum_{w \in (C \cup C')} ||w||^2 - |C \cup C'| \cdot ||\bar{w}||^2, \quad \text{mit } \bar{w} = \frac{|C| \cdot \bar{u} + |C'| \cdot \bar{v}}{|C| + |C'|}$$

$$SSE(C \cup C') - SSE(C) - SSE(C') = \dots = \frac{|C| \cdot |C'|}{|C| + |C'|} \cdot ||\bar{u} - \bar{v}||^2$$

\bar{u} bzw. \bar{v} bezeichnen den Mittelwert der Punkte $u \in C$ bzw. $v \in C'$.

ML-IX-47 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Hierarchische Verfahren

Update-Formel für Cluster-Distanzen

Nach Vereinigung der Cluster C und C' zu einem neuen Cluster sind dessen Distanzen, in Zeichen: $d_C(C \cup C', C_i)$, zu den anderen Clustern C_i zu berechnen.

Die folgende Update-Formel (Lance-Williams-Formel) ermöglicht eine effiziente Berechnung dieser Distanzen:

$$d_C(C \cup C', C_i) = \alpha \cdot d_C(C, C_i) + \\ \beta \cdot d_C(C', C_i) + \\ \gamma \cdot d_C(C, C') + \\ \delta \cdot |d_C(C, C_i) - d_C(C', C_i)|$$

Die Konstanten $\alpha, \beta, \gamma, \delta$ sind spezifisch für Single-Link, Complete-Link, Average-Link und Ward, und lassen sich auf Basis der entsprechenden Berechnungsvorschriften für d_C herleiten.

ML-IX-49 Cluster Analysis

©STEIN 2002-2012

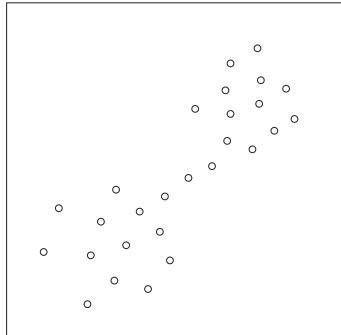
©STEIN 2002-2012

Bemerkungen:

- ❑ Link-basierte Verfahren arbeiten sowohl mit beliebigen Distanz- als auch Ähnlichkeitsmaßen.
- ❑ Single-Link kann unmittelbar mit einem Minimum-Spanning-Tree-Algorithmus realisiert werden.
- ❑ Varianzbasierte Verfahren sind nur sinnvoll, falls alle Merkmale Intervallskalenniveau besitzen.
- ❑ Idee der Lance-Williams-Formel: Anstatt immer wieder alle Elemente von zwei Clustern zu betrachten, nimmt die Update-Formel Bezug auf die im vorherigen Schritt berechneten Distanzen.

Hierarchische Verfahren

Chaining-Problematik bei Single-Link (d_C = Nearest-Neighbor)



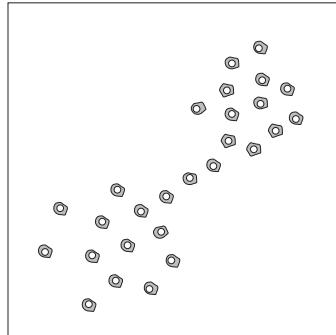
ML-IX-51 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link (d_C = Nearest-Neighbor)

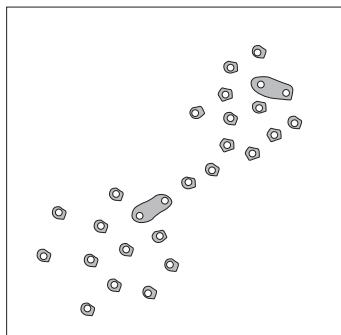


ML-IX-52 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link (d_C = Nearest-Neighbor)



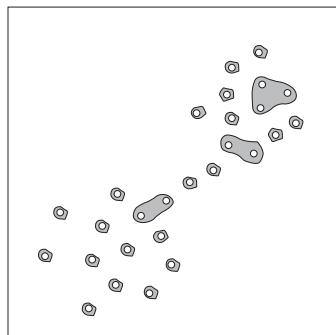
ML-IX-53 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link (d_C = Nearest-Neighbor)

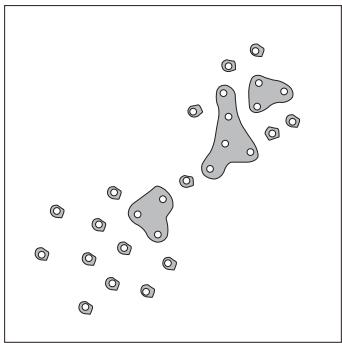


ML-IX-54 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_c = \text{Nearest-Neighbor}$)

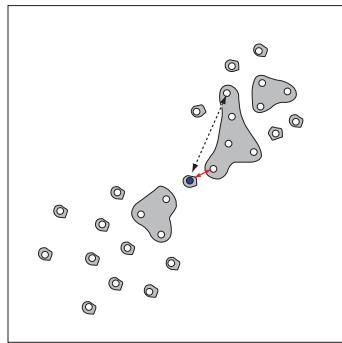


ML-IX-55 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_c = \text{Nearest-Neighbor}$)

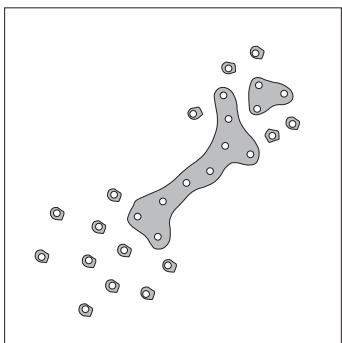


ML-IX-56 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_c = \text{Nearest-Neighbor}$)

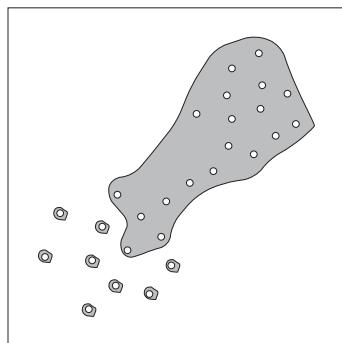


ML-IX-57 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_c = \text{Nearest-Neighbor}$)

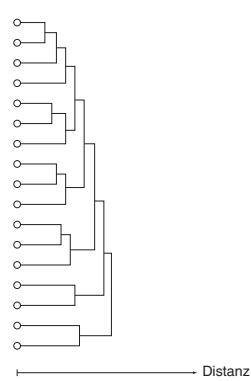
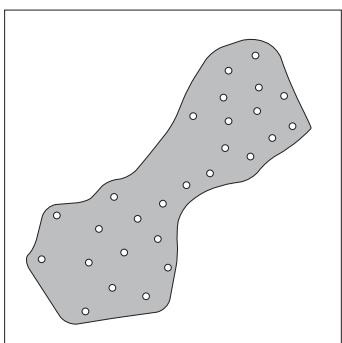


ML-IX-58 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_c = \text{Nearest-Neighbor}$)

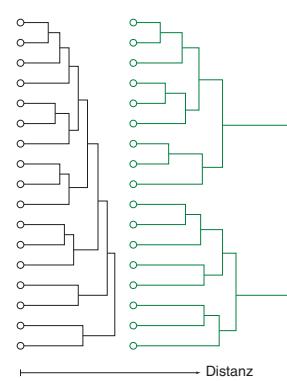
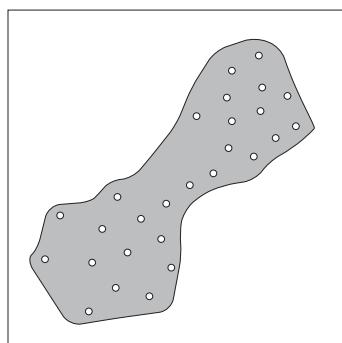


ML-IX-59 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_c = \text{Nearest-Neighbor}$)



ML-IX-60 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

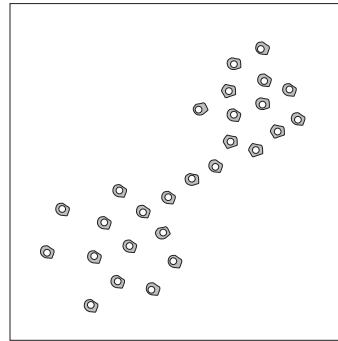
- ❑ Mit einer k -Nearest-Neighbor-Variante könnte man das Problem entschärfen.
- ❑ Bei k -Nearest-Neighbor werden größere Cluster bei der Agglomeration bevorzugt, da sie mehr und damit – statistisch gesehen – auch mehr nähere Nachbarn besitzen.

ML-IX-61 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

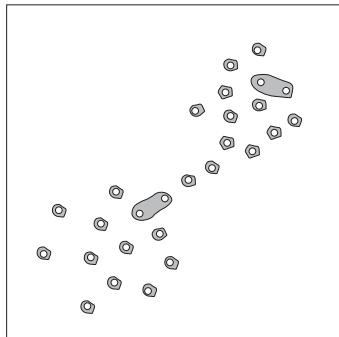
Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)



©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)

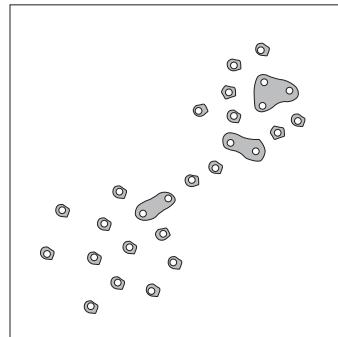


ML-IX-63 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

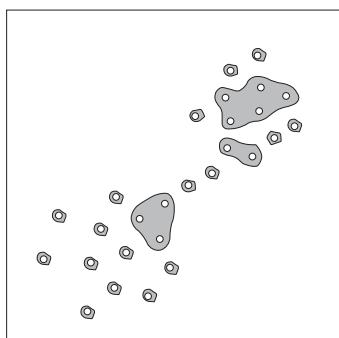
Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)



©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)

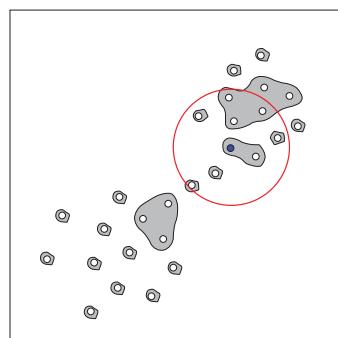


ML-IX-65 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)

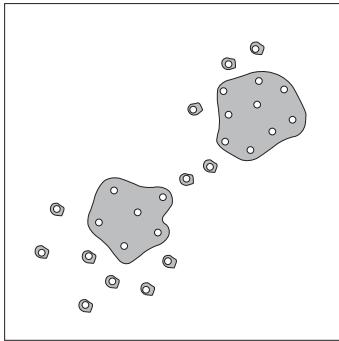


ML-IX-66 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)

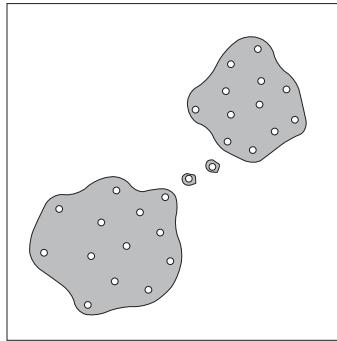


ML-IX-67 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)

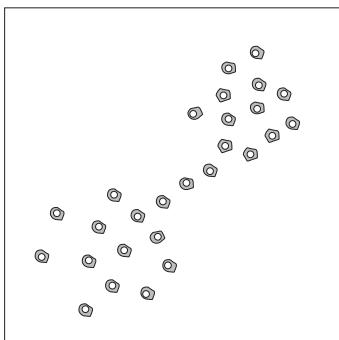


ML-IX-68 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)



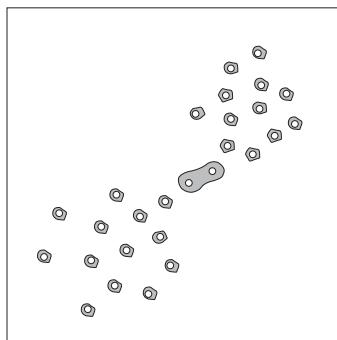
In speziellen Situationen kann auch k -Nearest-Neighbor versagen.

ML-IX-69 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)



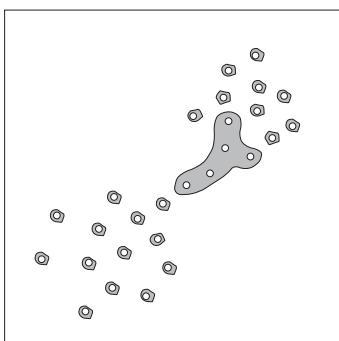
In speziellen Situationen kann auch k -Nearest-Neighbor versagen.

ML-IX-70 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)



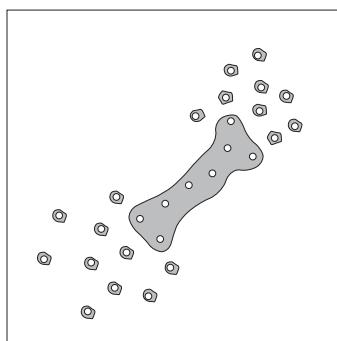
In speziellen Situationen kann auch k -Nearest-Neighbor versagen.

ML-IX-71 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Chaining-Problematik bei Single-Link ($d_C = k$ -Nearest-Neighbor)



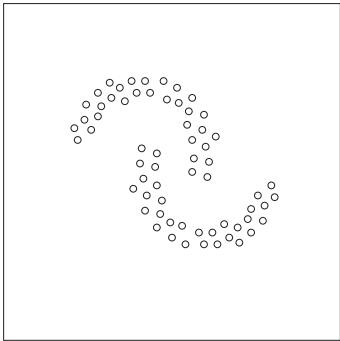
In speziellen Situationen kann auch k -Nearest-Neighbor versagen.

ML-IX-72 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)

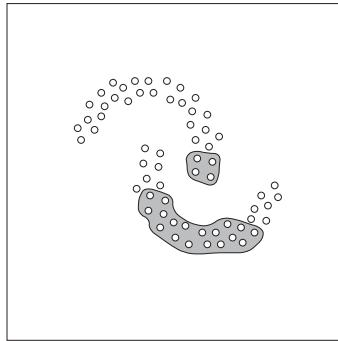


ML-IX-73 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)

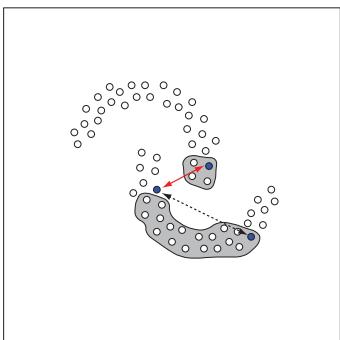


ML-IX-74 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)

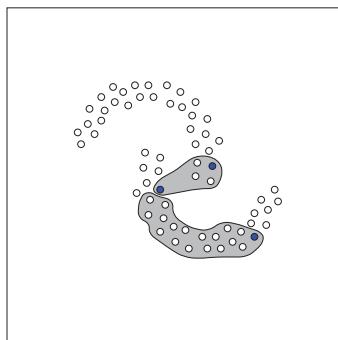


ML-IX-75 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)

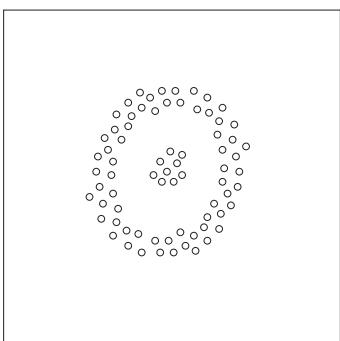


ML-IX-76 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)

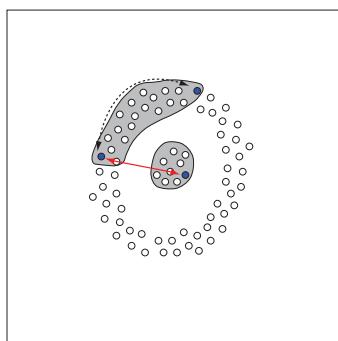


ML-IX-77 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)

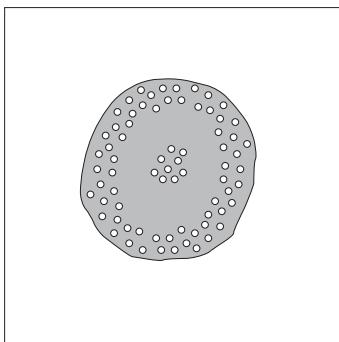


ML-IX-78 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)

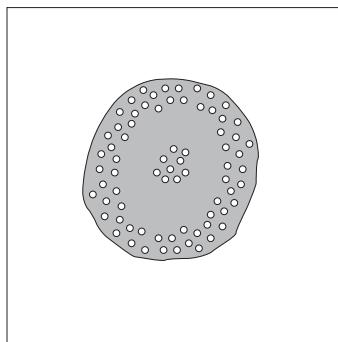


ML-IX-79 Cluster Analysis

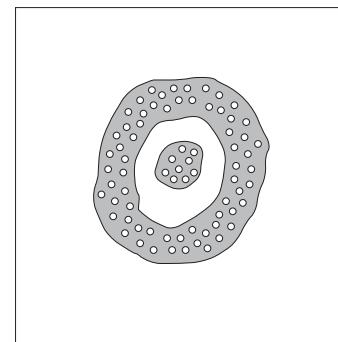
©STEIN 2002-2012

Hierarchische Verfahren

Überlappungsproblematik bei Complete-Link (d_c = Furthest-Neighbor)



Wirklichkeit



Wunsch

©STEIN 2002-2012

Hierarchische Verfahren

Eigenschaften hierarchisch-agglomerativer Verfahren [Distanzmaße]

Geometrische Eigenschaften:

	Single-Link	Complete-Link	Average-Link	Ward
Charakteristik	kontrahierend:	dilatierend:	konservativ:	konservativ:
Cluster-Zahl	niedrig	hoch	mittel	mittel
Cluster-Form	ausgedehnt	klein	kompakt	sphärisch
Verkettungstendenz	stark	niedrig	niedrig	niedrig
Ausreißerentdeckung	sehr gut	schlecht	mittel	mittel

ML-IX-81 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Eigenschaften hierarchisch-agglomerativer Verfahren [Distanzmaße]

Geometrische Eigenschaften:

	Single-Link	Complete-Link	Average-Link	Ward
Charakteristik	kontrahierend:	dilatierend:	konservativ:	konservativ:
Cluster-Zahl	niedrig	hoch	mittel	mittel
Cluster-Form	ausgedehnt	klein	kompakt	sphärisch
Verkettungstendenz	stark	niedrig	niedrig	niedrig
Ausreißerentdeckung	sehr gut	schlecht	mittel	mittel

Datenbezogene Eigenschaften:

verrauschte Daten	empfindlich	empfindlich	unbeeinflusst	unbeeinflusst
Merkmaltransformation	invariant	invariant	–	–

ML-IX-82 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Eigenschaften hierarchisch-agglomerativer Verfahren [Distanzmaße]

Geometrische Eigenschaften:

	Single-Link	Complete-Link	Average-Link	Ward
Charakteristik	kontrahierend:	dilatierend:	konservativ:	konservativ:
Cluster-Zahl	niedrig	hoch	mittel	mittel
Cluster-Form	ausgedehnt	klein	kompakt	sphärisch
Verkettungstendenz	stark	niedrig	niedrig	niedrig
Ausreißerentdeckung	sehr gut	schlecht	mittel	mittel

Datenbezogene Eigenschaften:

verrauschte Daten	empfindlich	empfindlich	unbeeinflusst	unbeeinflusst
Merkmaltransformation	invariant	invariant	–	–

Eigenschaften des Cluster-Distanzmaßes:

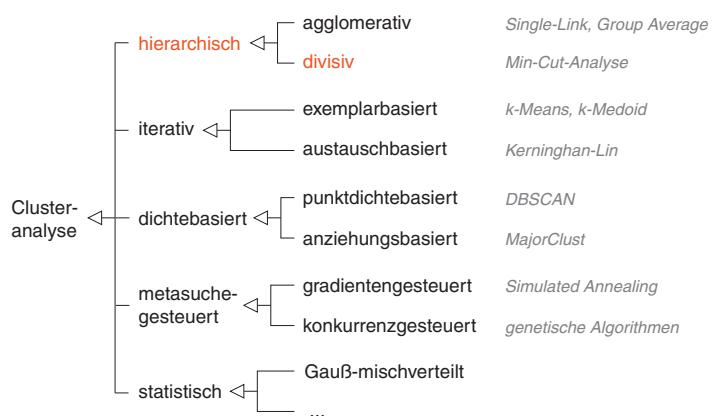
d_c monoton	✓	✓	✓	✓
d_c reihenfolgeunabh.	✓	✓	✓	✓
d_c konsistent	→ 0	→ ∞	✓	→ ∞

ML-IX-83 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Prinzipien der Fusionierung



ML-IX-84 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Algorithmus zur hierarchisch-divisiven Clusteranalyse

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d_C . Distance measure between two clusters.

Output: $T = \langle V_T, E_T \rangle$. Cluster hierarchy or dendrogram.

```

1.  $\mathcal{C} = \{V\}$  // initial clustering
2.
3. WHILE  $\exists C_x : (C_x \in \mathcal{C} \wedge |C_x| > 1)$  DO
4.    $\{C, C'\} = \underset{\substack{\{C_i, C_j\}: \\ C_i \cup C_j = C_x \wedge C_i \cap C_j = \emptyset}}{\operatorname{argmax}} d_C(C_i, C_j)$ 
5.    $\mathcal{C} = (\mathcal{C} \setminus \{C_x\}) \cup \{C, C'\}$  // splitting
6.
7. ENDDO
8. RETURN( $T$ )

```

Vergleiche hierzu den Algorithmus zur hierarchisch-agglomerativen Clusteranalyse.

ML-IX-85 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

- Im Prinzip kann d_C wie bei den hierarchisch-agglomerativen Verfahren gewählt werden. Die Worst-Case-Komplexität ist exponentiell statt quadratisch.
- Hierarchisch-divisive Verfahren werden oft als *monothetische* Variante konstruiert: in jedem Entscheidungsschritt wird nur *eine* Variable betrachtet.
- Im Gegensatz zu hierarchisch-agglomerativen Verfahren darf ein hierarchisch-divisives Verfahren keinen „Fehler“ bei den ersten Schritten machen.
- Ein mächtiges hierarchisch-divisives Clusteranalyse-Verfahren entsteht mit

$$\text{sim}_C(C, C') = \sum_{e \in \text{cut}(\{C, C'\})} w(e) \quad \text{bzw.} \quad d_C(C, C') = \frac{1}{\text{sim}_C(C, C')}$$

ML-IX-87 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

MinCut-Clusteranalyse

Definition 3 (Cut, minimaler Cut)

Sei $G = \langle V, E, w \rangle$ ein Graph mit nicht-negativer Gewichtsfunktion w ; weiterhin sei $U \subset V$ eine nichtleere Teilmenge der Knotenmenge V und $\bar{U} = V \setminus U$. Der Cut zwischen U und \bar{U} ist wie folgt definiert:

$$\text{cut}(\{U, \bar{U}\}) = \{\{u, v\} \mid \{u, v\} \in E, u \in U, v \in \bar{U}\}$$

Weiterhin bezeichne $w(\{U, \bar{U}\})$ das Gewicht oder die Kapazität von $\text{cut}(\{U, \bar{U}\})$:

$$w(\{U, \bar{U}\}) = \sum_{e \in \text{cut}(\{U, \bar{U}\})} w(e)$$

$\text{cut}(\{U, \bar{U}\})$ heißt minimaler Cut von G (*Minimum Capacity Cut*), wenn für alle Zerlegungen $\{W, \bar{W}\}$, $W, \bar{W} \neq \emptyset$ gilt:

$$w(\{U, \bar{U}\}) \leq w(\{W, \bar{W}\})$$

ML-IX-89 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Algorithmus zur hierarchisch-divisiven Clusteranalyse

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d_C . Distance measure between two clusters.

Output: $T = \langle V_T, E_T \rangle$. Cluster hierarchy or dendrogram.

```

1.  $\mathcal{C} = \{V\}$  // initial clustering
2.  $V_T = \{v_C \mid C \in \mathcal{C}\}$ ,  $E_T = \emptyset$  // initial dendrogram
3. WHILE  $\exists C_x : (C_x \in \mathcal{C} \wedge |C_x| > 1)$  DO
4.    $\{C, C'\} = \underset{\substack{\{C_i, C_j\}: \\ C_i \cup C_j = C_x \wedge C_i \cap C_j = \emptyset}}{\operatorname{argmax}} d_C(C_i, C_j)$ 
5.    $\mathcal{C} = (\mathcal{C} \setminus \{C_x\}) \cup \{C, C'\}$  // splitting
6.    $V_T = V_T \cup \{v_C, v_{C'}\}$ ,  $E_T = E_T \cup \{\{v_{C_x}, v_C\}, \{v_{C_x}, v_{C'}\}\}$  // dendrogram
7. ENDDO
8. RETURN( $T$ )

```

Vergleiche hierzu den Algorithmus zur hierarchisch-agglomerativen Clusteranalyse.

ML-IX-86 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

MinCut-Clusteranalyse

Definition 3 (Cut, minimaler Cut)

Sei $G = \langle V, E, w \rangle$ ein Graph mit nicht-negativer Gewichtsfunktion w ; weiterhin sei $U \subset V$ eine nichtleere Teilmenge der Knotenmenge V und $\bar{U} = V \setminus U$. Der Cut zwischen U und \bar{U} ist wie folgt definiert:

$$\text{cut}(\{U, \bar{U}\}) = \{\{u, v\} \mid \{u, v\} \in E, u \in U, v \in \bar{U}\}$$

ML-IX-88 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

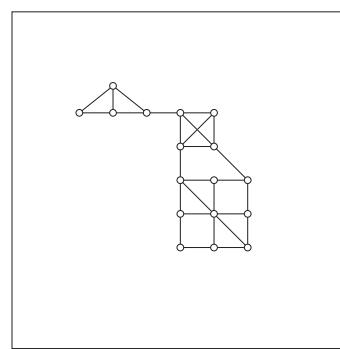
MinCut-Clusteranalyse

Definition 3 (Cut, minimaler Cut)

Sei $G = \langle V, E, w \rangle$ ein Graph mit nicht-negativer Gewichtsfunktion w ; weiterhin sei $U \subset V$ eine nichtleere Teilmenge der Knotenmenge V und $\bar{U} = V \setminus U$. Der Cut zwischen U und \bar{U} ist wie folgt definiert:

Hierarchische Verfahren

MinCut-Clusteranalyse



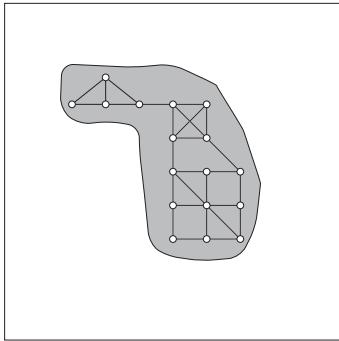
ML-IX-89 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Hierarchische Verfahren

MinCut-Clusteranalyse

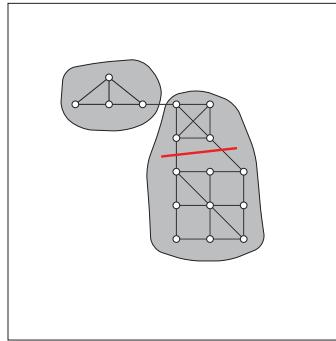


ML-IX-91 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

MinCut-Clusteranalyse

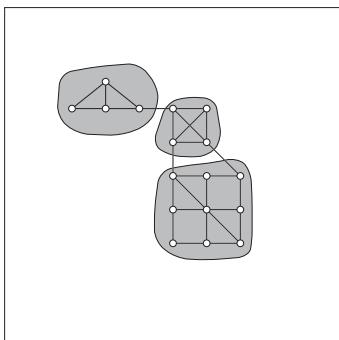


ML-IX-92 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

MinCut-Clusteranalyse



ML-IX-93 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

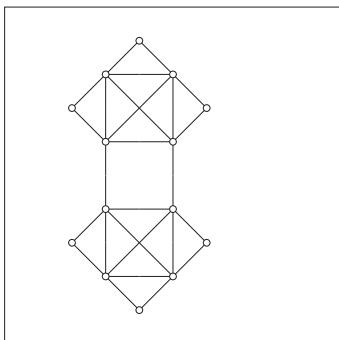
- Jede Aufteilung erfordert die Berechnung eines Cuts minimaler Kapazität (*Minimum-Capacity-Cut*). Beachte, dass kein Knoten als Quelle oder Senke vorgegeben ist.
- Die Laufzeitkomplexität des besten, bekannten Algorithmus zur Berechnung eines minimalen Cuts ist in $O(|V| \cdot |E| + |V|^2 \cdot \log |V|)$. [Nagamochi/Ono/Ibaraki 1994]
- Es sind $|V| - 1$ Berechnungen eines Minimum-Capacity-Cuts notwendig, um eine vollständige Zerlegung (= ein Knoten pro Cluster) herzustellen.
- Der Aufwand zur Berechnung eines Minimum- $s-t$ -Cuts, also eines Cuts bei vorgegebener Quelle s und Senke t , ist in $O(|V|^2 \log(|E|))$.
- Der Aufwand zur Berechnung eines Balanced Min-Cut (k -way, $k \geq 2$) ist NP-vollständig.

ML-IX-94 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Splitting-Problematik bei MinCut-Clusteranalyse

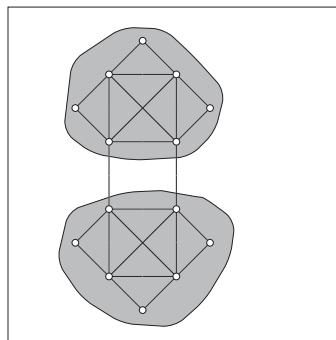


ML-IX-95 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Splitting-Problematik bei MinCut-Clusteranalyse

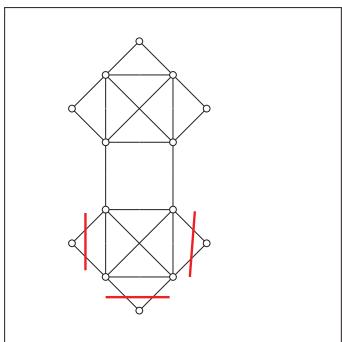


ML-IX-96 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Splitting-Problematik bei MinCut-Clusteranalyse



Ausweg: Normalisierung der Cut-Kapazität bzgl. der Größe der Knotenmengen.

ML-IX-97 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

- Die Bestimmung eines Cuts minimaler, normalisierter Kapazität ist NP-vollständig.
- Es gibt effiziente Näherungslösungen zur Berechnung von $\bar{w}(\{U, \bar{U}\})$. Das Verfahren wird angewandt im Bereich der Bildsegmentierung und der Gene-Expression-Cluster-Analyse. [Shi/Malik 2000]

ML-IX-99 Cluster Analysis

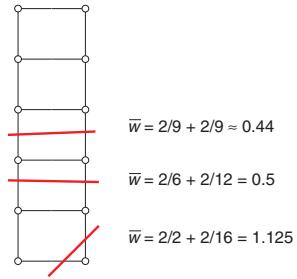
©STEIN 2002-2012

Hierarchische Verfahren

Splitting-Problematik bei MinCut-Clusteranalyse

$$\text{Normalisierte Cut-Kapazität: } \bar{w}(\{U, \bar{U}\}) = \frac{w(\{U, \bar{U}\})}{w(\{U, V\})} + \frac{w(\{U, \bar{U}\})}{w(\{\bar{U}, V\})}$$

Illustration von \bar{w} :



$$\bar{w} = 2/9 + 2/9 \approx 0.44$$

$$\bar{w} = 2/6 + 2/12 = 0.5$$

$$\bar{w} = 2/2 + 2/16 = 1.125$$

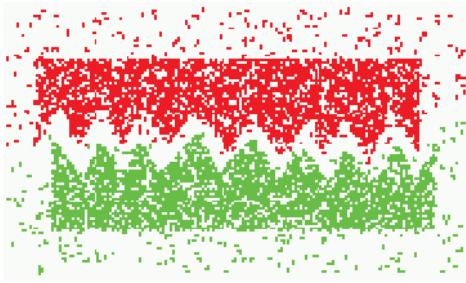
ML-IX-98 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Kombination hierarchischer Verfahren

Das System Chameleon kombiniert die Schritte Graphausdünnung, Graphpartitionierung und hierarchische Cluster-Analyse [Karypis/Han/Kumar 2000] :



Die Clusterdistanz $d_C(C, C')$ ist definiert als $d_C = \frac{1}{R_I(C, C') \cdot (R_C(C, C'))^\alpha}$

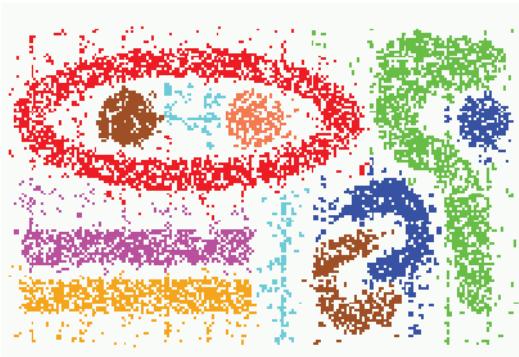
ML-IX-101 Cluster Analysis

©STEIN 2002-2012

Hierarchische Verfahren

Kombination hierarchischer Verfahren

Chameleon [Karypis/Han/Kumar 2000] :



Der Parameter α in d_C ist vom Anwender problemabhängig zu bestimmen.

ML-IX-102 Cluster Analysis

©STEIN 2002-2012

Kapitel ML: IX (Fortsetzung)

IX. Clusteranalyse

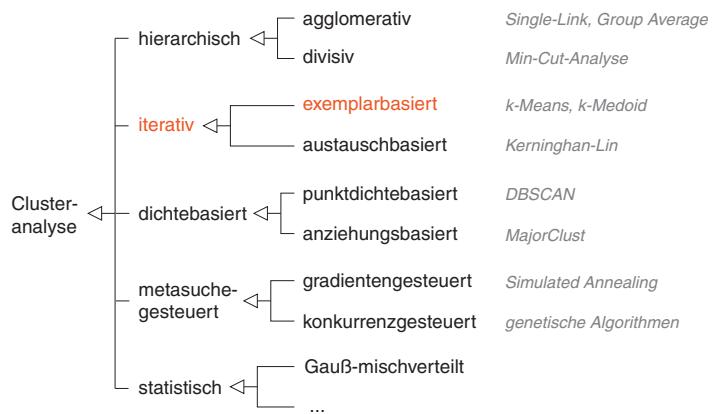
- Einordnung Data Mining
- Einführung in die Clusteranalyse
- Hierarchische Verfahren
- Iterative Verfahren
- Dichtebaserte Verfahren
- Cluster Evaluation
- Constrained Cluster Analysis

ML-IX-103 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

Prinzipien der Fusionierung



Single-Link, Group Average
Min-Cut-Analyse

k-Means, k-Medoid
Kernighan-Lin

DBSCAN
MajorClust

Simulated Annealing
genetische Algorithmen

ML-IX-104 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

Algorithmus zur exemplarbasierten Clusteranalyse

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d . Distance function for nodes in V .
 e . Minimization criterion for cluster representatives, based on d .
 k . Number of desired clusters.

Output: r_1, \dots, r_k . Cluster representatives.

```

1.
2. FOR i = 1 to k DO  $r_i(t) = \text{choose}(V)$  // init representatives
3.
4.
5.

6. FOREACH  $v \in V$  DO // find nearest representative (cluster)
7.    $i = \underset{j:j \in [1, \dots, k]}{\text{argmin}} d(r_j(t), v)$ ,  $C_i = C_i \cup \{v\}$ 
8. ENDDO
9. FOR i = 1 to k DO  $r_i(t) = \text{minimize}(e(C_i))$  // update

10.
11.

```

ML-IX-105 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

Algorithmus zur exemplarbasierten Clusteranalyse

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d . Distance function for nodes in V .
 e . Minimization criterion for cluster representatives, based on d .
 k . Number of desired clusters.

Output: r_1, \dots, r_k . Cluster representatives.

```

1.  $t = 0$ 
2. FOR i = 1 to k DO  $r_i(t) = \text{choose}(V)$  // init representatives
3. REPEAT
4.    $t = t + 1$ 
5.   FOR i = 1 to k DO  $C_i = \emptyset$ 

6.   FOREACH  $v \in V$  DO // find nearest representative (cluster)
7.      $i = \underset{j:j \in [1, \dots, k]}{\text{argmin}} d(r_j(t), v)$ ,  $C_i = C_i \cup \{v\}$ 
8.   ENDDO
9.   FOR i = 1 to k DO  $r_i(t) = \text{minimize}(e(C_i))$  // update

10. UNTIL (convergence( $r_1(t), \dots, r_k(t)$ ) OR  $t > t_{\max}$ )
11. RETURN( $\{r_1(t), \dots, r_k(t)\}$ )

```

ML-IX-106 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

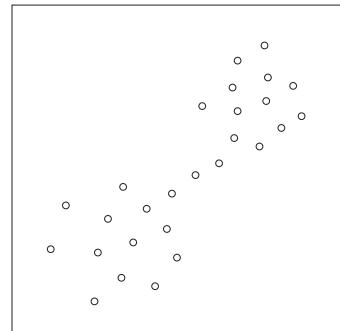
- Die Cluster-Repräsentanten werden Centroide bzw. allgemein Medoide genannt.
- Die Funktion $\text{choose}(V)$ realisiert eine zufällige Auswahl ohne Zurücklegen.
- Als Distanzfunktion d wird bei metrischen Daten meistens der euklidische Abstand zwischen zwei Punkten gewählt. Ein alternativer und allgemeiner Ansatz ist die Verwendung des kürzesten Weges in G .
- Als Minimierungskriterium e wird bei metrischen Daten meistens die Summe der quadrierten Abweichungen zum Cluster-Repräsentanten gewählt (Varianzkriterium). Sind die Punkte $v \in V$ Vektoren aus dem \mathbb{R}^p , lässt sich der optimale Cluster-Repräsentant durch komponentenweise Berechnung des arithmetischen Mittels ermitteln.

ML-IX-107 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

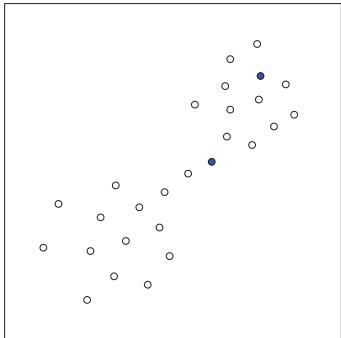
k -Means mit Minimierungskriterium e = Varianz



©STEIN 2002-2012

Iterative Verfahren

k -Means mit Minimierungskriterium $e = \text{Varianz}$

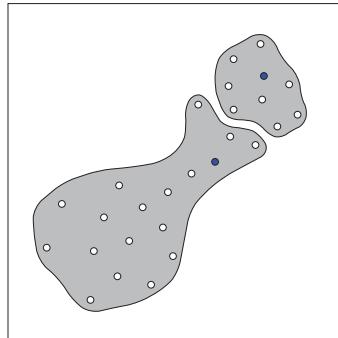


ML-IX-109 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means mit Minimierungskriterium $e = \text{Varianz}$

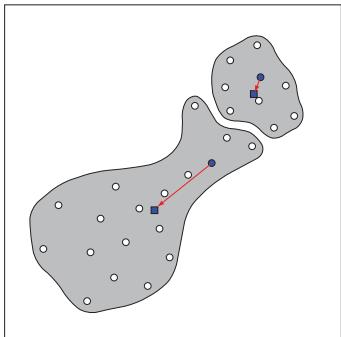


ML-IX-110 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means mit Minimierungskriterium $e = \text{Varianz}$

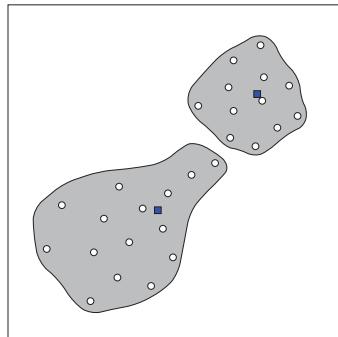


ML-IX-111 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means mit Minimierungskriterium $e = \text{Varianz}$

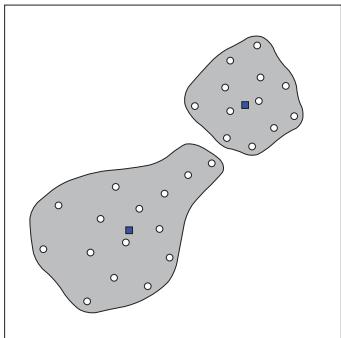


ML-IX-112 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means mit Minimierungskriterium $e = \text{Varianz}$

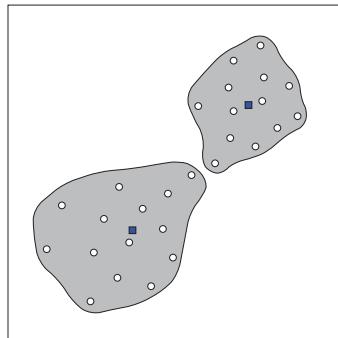


ML-IX-113 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means mit Minimierungskriterium $e = \text{Varianz}$

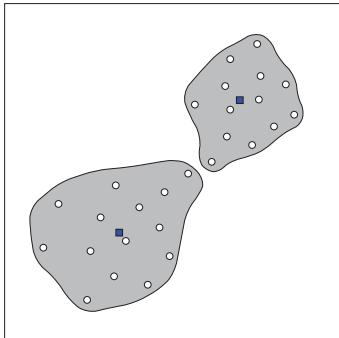


ML-IX-114 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means mit Minimierungskriterium e = Varianz



ML-IX-115 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

- \bar{v}_{C_i} bezeichne den Mittelwert der Punkte $v \in C_i$.
- Der Cluster-Repräsentant ist abkürzend mit r_i anstatt mit $r_i(t)$ bezeichnet.
- Die Centroid-Berechnung bei k -Means als Mittelwert der Cluster-Elemente entspricht einer lokalen, also Cluster-spezifischen Varianzminimierung.
- Der Medoid (Zentralelement) eines Cluster ist dasjenige Element, für das die Summe aller Distanzen zu diesem Element minimal ist. Ein Vorteil bei Verwendung von Medoiden ist das robustere Verhalten gegenüber Ausreißern und damit eventuell eine in der Anzahl der Iterationen schnellere Konvergenz.
- Bei Fuzzy- k -Means bezeichnet $\mu_i(v)$ den Zugehörigkeitswert von $v \in V$ zu Cluster C_i .
- k -Medoid und k -Center arbeiten sowohl mit beliebigen Distanz- als auch mit Ähnlichkeitsmaßen.
- k -Means und Fuzzy- k -Means setzen intervallskalierte Merkmale voraus.
- k -Means kann unmittelbar als Kohonen Self-Organizing-Map, SOM, einem speziellen neuronalen Netz, umgesetzt werden:
 - Die Netztopologie besteht aus einer Eingangsschicht mit Knoten in der Anzahl der Merkmale und einer Verarbeitungsschicht (*Competitive Layer*) mit k Knoten.
 - Der Lernalgorithmus bestimmt für einen Merkmalsvektor auf Basis der Kantengewichte das *Winning Neuron*, dessen Gewichte gemäß des Lernparameters η erhöht werden.

ML-IX-117 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

Minimierungskriterien exemplarbasierter Verfahren

$$e(C_i) = \sum_{v \in C_i} (v - r_i)^2$$

$$r_i = \bar{v}_{C_i}$$

Centroid-Berechnung
mittels Varianzminimierung
(k -Means)

$$e(C_i) = \sum_{v \in C_i} |v - r_i|$$

$$r_i \in C_i$$

Medoid-Berechnung
(k -Medoid)

$$e(C_i) = \max_{v \in C_i} |v - r_i|$$

$$r_i \in C_i$$

k -Center

$$e(C_i) = \sum_{v \in V} (\mu_i(v))^2 \cdot (v - r_i)^2 \quad r_i = \frac{\sum_{v \in V} (\mu_i(v))^2 \cdot v}{\sum_{v \in V} (\mu_i(v))^2}$$

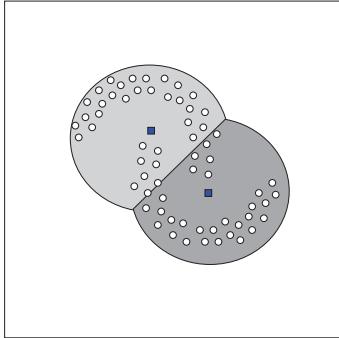
Fuzzy- k -Means

ML-IX-116 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means versus Single-Link



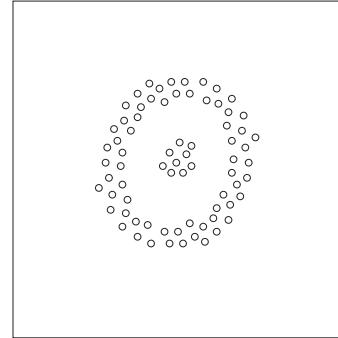
Exemplarbasierte Verfahren versagen bei verschrankt liegenden Clustern.

ML-IX-119 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k -Means versus Single-Link



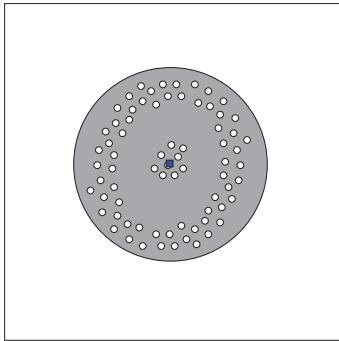
Exemplarbasierte Verfahren versagen bei verschrankt liegenden Clustern.

ML-IX-120 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k-Means versus Single-Link



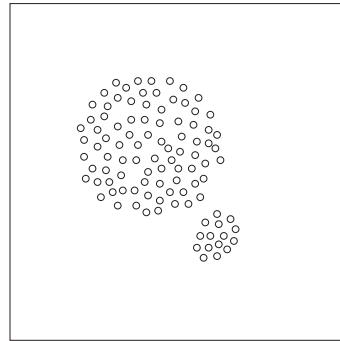
Exemplarbasierte Verfahren versagen bei verschränkt liegenden Clustern.

ML-IX-121 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k-Means versus Single-Link



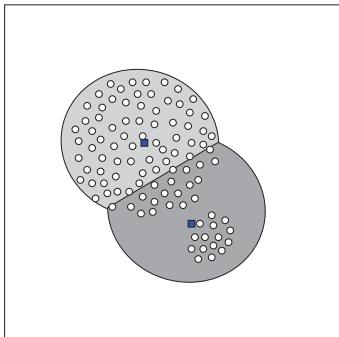
Exemplarbasierte Verfahren versagen bei extremen Größendifferenzen.

ML-IX-122 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

k-Means versus Single-Link



Exemplarbasierte Verfahren versagen bei extremen Größendifferenzen.

ML-IX-123 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

Exklusive versus nicht-exklusive Clusteranalyse

Sei $\mathcal{C} = \{C_1, \dots, C_k\}$ eine Partitionierung einer Menge V mit $\bigcup_{i=1 \dots k} C_i = V$.

- exklusive Clusteranalyse: $\forall i, j \in \{1, \dots, k\} : i \neq j$ impliziert $C_i \cap C_j = \emptyset$
- nicht-exklusive Clusteranalyse erlaubt mehrfache Cluster-Zugehörigkeit

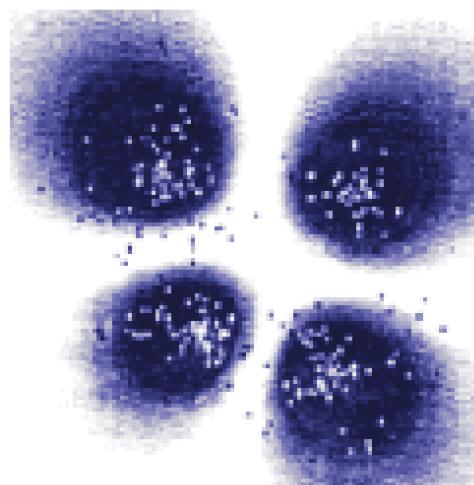
©STEIN 2002-2012

Iterative Verfahren

Exklusive versus nicht-exklusive Clusteranalyse

Sei $\mathcal{C} = \{C_1, \dots, C_k\}$ eine Partitionierung einer Menge V mit $\bigcup_{i=1 \dots k} C_i = V$.

- exklusive Clusteranalyse: $\forall i, j \in \{1, \dots, k\} : i \neq j$ impliziert $C_i \cap C_j = \emptyset$
- nicht-exklusive Clusteranalyse erlaubt mehrfache Cluster-Zugehörigkeit
- Fuzzy-Clusteranalyse quantifiziert die Cluster-Zugehörigkeit der $v \in V$ mit Zugehörigkeitsfunktionen $\mu_i(v)$, $i \in \{1, \dots, k\}$.



[Höppner/Klawonn/Kruse 1997]

ML-IX-125 Cluster Analysis

©STEIN 2002-2012

ML-IX-126 Cluster Analysis

©STEIN 2002-2012

Iterative Verfahren

Exklusive versus nicht-exklusive Clusteranalyse

Anwendung der Fuzzy-Clusteranalyse zur Darstellung von Gehirngewebe:



[Pham/Prince/Dagher/Xn 1996]

ML-IX-127 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Bemerkungen:

- Die linguistische Variable bei der Fuzzy-Modellierung hat k Ausprägungen, entsprechend den Clustern C_1, \dots, C_k .
- Üblich ist eine Normalisierungsrestriktion für Zugehörigkeitsfunktionen: $\sum_{i=1\dots k} \mu_i(v) = 1$
- Varianten von Fuzzy- k -Means ohne die Normalisierungsrestriktion haben den Nachteil, dass Punkte mit kleinen Zugehörigkeitswerten zu einem Cluster wie Ausreißer behandelt werden, anstatt den Cluster in ihre Richtung zu bewegen. Deshalb ist es sinnvoll, das Iterationsverfahren in einer Initialisierungsphase zunächst mit der Restriktion anzuwenden.
- Eine Klassifikation durch eine unscharfe Cluster-Analyse ist vorteilhaft, wenn keine klar ausgebildete Klassenstruktur vorliegt bzw. wenn sich nicht alle Vektoren eindeutig einer Klasse zuordnen lassen.
- Ein Nachteil der Fuzzy-Clusteranalyse ist, dass keine Repräsentanten für Cluster ermittelt werden.

Kapitel ML:IX (Fortsetzung)

IX. Clusteranalyse

- Einordnung Data Mining
- Einführung in die Clusteranalyse
- Hierarchische Verfahren
- Iterative Verfahren
- Dichtebaserte Verfahren
- Cluster Evaluation
- Constrained Cluster Analysis

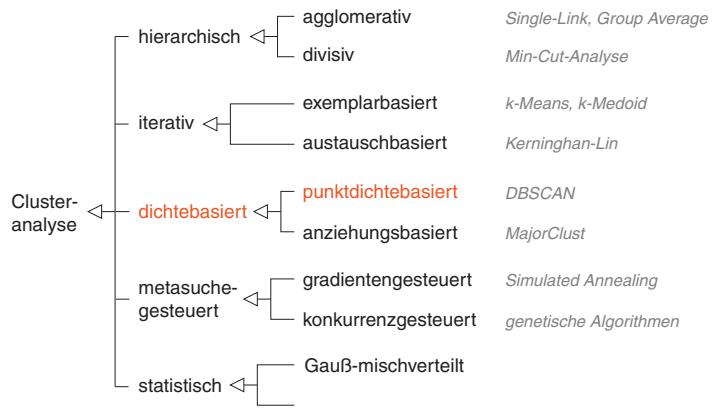
ML-IX-129 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Dichtebaserte Verfahren

Prinzipien der Fusionierung



ML-IX-130 Cluster Analysis

©STEIN 2002-2012

Dichtebaserte Verfahren

Dichtebaserte Verfahren versuchen, den Graphen $G = \langle V, E, w \rangle$ bzw. die Punktmenge V in Bereiche gleicher Dichte aufzuteilen.

Ansätze zur Dichteschätzung:

- parameterbasiert: die unterliegende Verteilung ist bekannt
- parameterlos: Histogramme (Konstruktion von Bar-Charts), Kerndichteschätzer (Überlagerung kontinuierlicher Funktionen)

ML-IX-131 Cluster Analysis

©STEIN 2002-2012

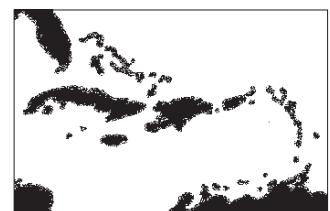
Dichtebaserte Verfahren

Dichtebaserte Verfahren versuchen, den Graphen $G = \langle V, E, w \rangle$ bzw. die Punktmenge V in Bereiche gleicher Dichte aufzuteilen.

Ansätze zur Dichteschätzung:

- parameterbasiert: die unterliegende Verteilung ist bekannt
- parameterlos: Histogramme (Konstruktion von Bar-Charts), Kerndichteschätzer (Überlagerung kontinuierlicher Funktionen)

Beispiel (karibische Inseln):

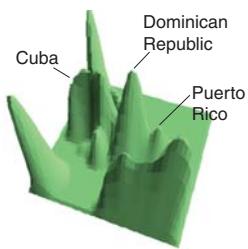


ML-IX-132 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

Dichteschätzung mit Gauß'schem Kern für das Beispiel

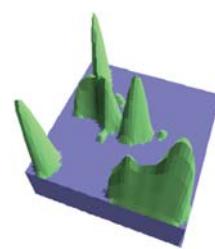


ML-IX-133 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

Dichteschätzung mit Gauß'schem Kern für das Beispiel

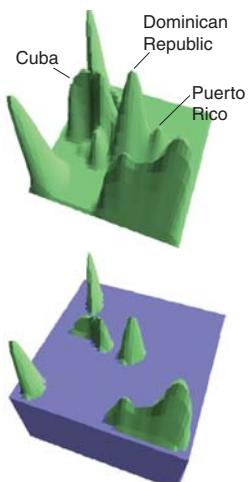


ML-IX-134 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

Dichteschätzung mit Gauß'schem Kern für das Beispiel

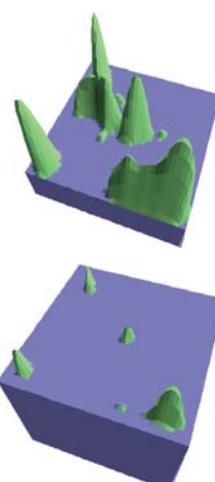
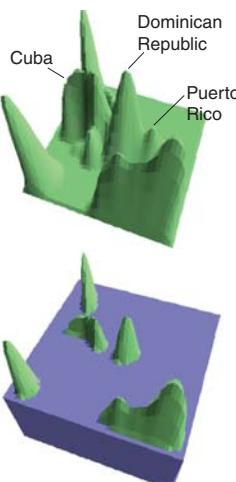


ML-IX-135 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

Dichteschätzung mit Gauß'schem Kern für das Beispiel



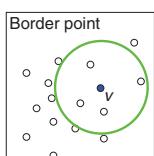
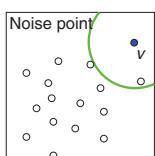
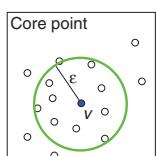
ML-IX-136 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN: Prinzip der Dichteschätzung [Ester et al. 1996]

Sei $N_\varepsilon(v)$ die ε -Nachbarschaft eines Punktes v . Unterscheidung von drei Punkttypen:



1. v ist Kernpunkt (core point) $\Leftrightarrow |N_\varepsilon(v)| \geq \text{MinPts}$
2. v ist Rauschen (noise point) $\Leftrightarrow v$ ist von keinem Kernpunkt aus dichteerreichbar (density-reachable)
3. v ist Randpunkt (border point) in allen anderen Fällen

ML-IX-137 Cluster Analysis

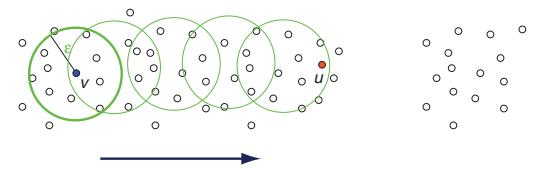
©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN: Prinzip der Dichteschätzung

Ein Punkt u ist dichteerreichbar von einem Punkt v , falls gilt:

- (a) $u \in N_\varepsilon(v)$, wobei v ein Kernpunkt ist, oder
- (b) es gibt eine Menge von Punkten $\{v_1, \dots, v_l\}$:
 $v_{i+1} \in N_\varepsilon(v_i)$ und v_i ist Kernpunkt, $i = 1, \dots, l-1$, mit $v_1 = v$, $v_l = u$.



Bedingung (b) lässt sich als transitive Anwendung von Bedingung (a) auffassen.

ML-IX-138 Cluster Analysis

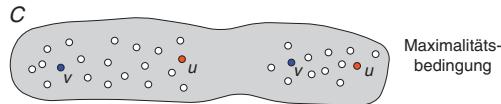
©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN: Cluster-Interpretation

Ein Cluster $C \subseteq V$ erfüllt folgende Bedingungen:

1. $\forall u, v : \text{Falls } v \in C \text{ und } u \text{ dichterreichbar von } v, \text{ dann ist } u \in C.$



ML-IX-139 Cluster Analysis

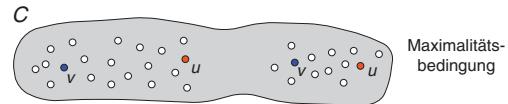
©STEIN 2002-2012

Dichtebasierter Verfahren

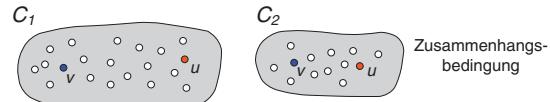
DBSCAN: Cluster-Interpretation

Ein Cluster $C \subseteq V$ erfüllt folgende Bedingungen:

1. $\forall u, v : \text{Falls } v \in C \text{ und } u \text{ dichterreichbar von } v, \text{ dann ist } u \in C.$



2. $\forall u, v : u \text{ ist dichteverbunden mit } v, \text{ d.h., es existiert ein Punkt } t \text{ von dem } u \text{ und } v \text{ dichterreichbar sind.}$



ML-IX-140 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN: Algorithmus

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d . Distance function for nodes in V .
 ε . Neighborhood radius.
 MinPts . Lower bound for point number in ε -neighborhood.

Output: $\gamma : V \rightarrow \mathbb{Z}$. Cluster assignment function.

- 1.
- 2.
- 3.
4. $N_\varepsilon(v) = \text{neighborhood}(G, d, v, \varepsilon)$
5. **IF** $|N_\varepsilon(v)| \geq \text{MinPts}$ **THEN** // v is core point
6. $i = i + 1$
7. $C_i = \text{density_reachable_hull}(G, d, N_\varepsilon(v))$ // form a new cluster
8. **FOREACH** $v \in C_i$ **DO** $\gamma(v) = i$
9. **ELSE** $\gamma(v) = -1$ // v is _preliminarily_ classified as noise
- 10.
- 11.

ML-IX-141 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN: Algorithmus

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d . Distance function for nodes in V .
 ε . Neighborhood radius.
 MinPts . Lower bound for point number in ε -neighborhood.

Output: $\gamma : V \rightarrow \mathbb{Z}$. Cluster assignment function.

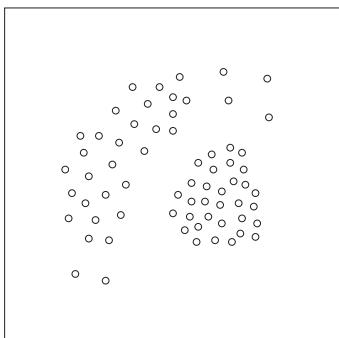
1. $i = 0$
2. **WHILE** $\exists v : (v \in V \text{ AND } \gamma(v) = \perp)$ **DO** // \perp = unclassified
3. $v = \text{choose_unclassified_point}(V)$
4. $N_\varepsilon(v) = \text{neighborhood}(G, d, v, \varepsilon)$
5. **IF** $|N_\varepsilon(v)| \geq \text{MinPts}$ **THEN** // v is core point
6. $i = i + 1$
7. $C_i = \text{density_reachable_hull}(G, d, N_\varepsilon(v))$ // form a new cluster
8. **FOREACH** $v \in C_i$ **DO** $\gamma(v) = i$
9. **ELSE** $\gamma(v) = -1$ // v is _preliminarily_ classified as noise
10. **ENDDO**
11. **RETURN**(γ)

ML-IX-142 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN

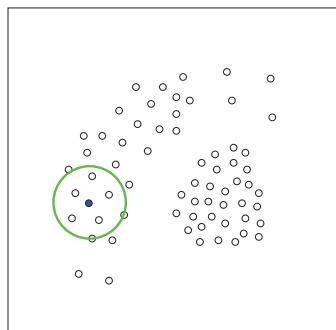


ML-IX-143 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN

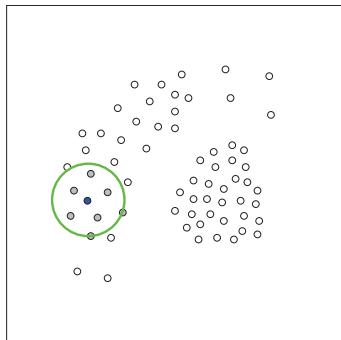


ML-IX-144 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN

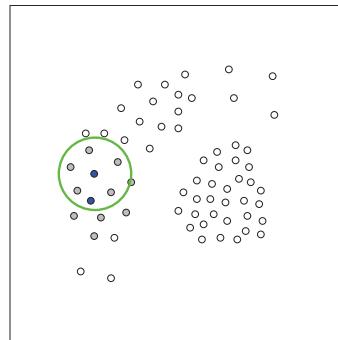


ML-IX-145 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN

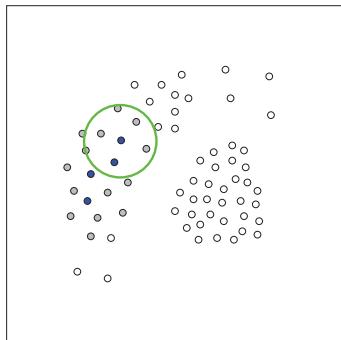


ML-IX-146 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN

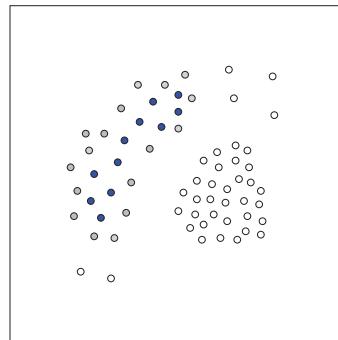


ML-IX-147 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN



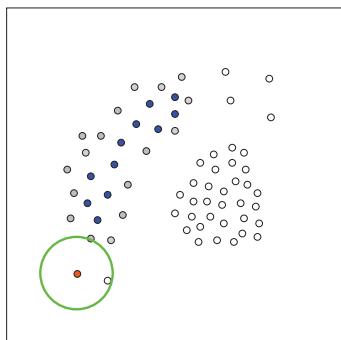
- Core point
- Border point

ML-IX-148 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN



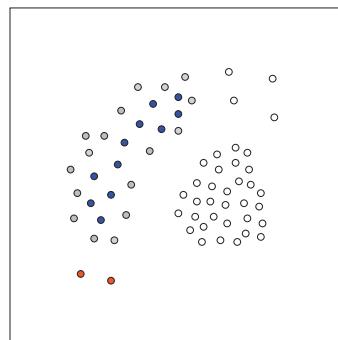
- Core point
- Border point

ML-IX-149 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN



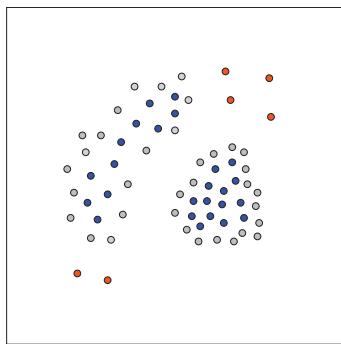
- Core point
- Border point
- Noise point

ML-IX-150 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN

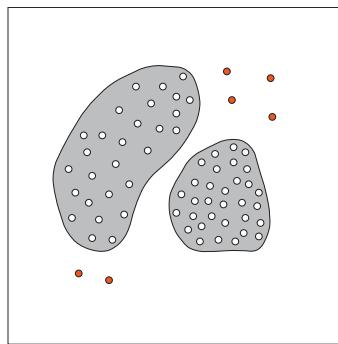


- Core point
- Border point
- Noise point

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN

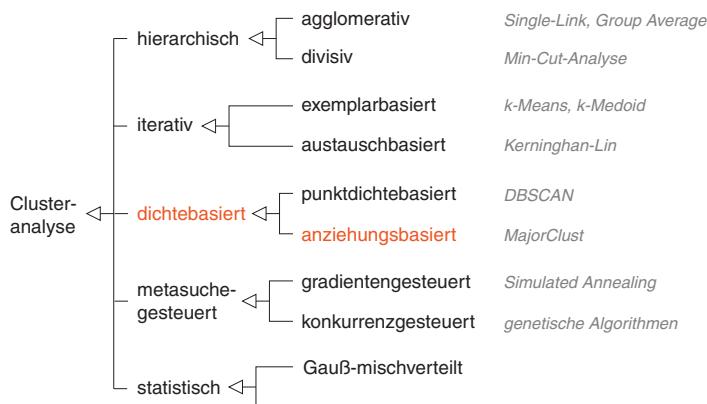


- Noise point

©STEIN 2002-2012

Dichtebasierter Verfahren

Prinzipien der Fusionierung



ML-IX-153 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung [Stein/Niggemann 1999]

Die gewichteten Kanten im Graph $G = \langle V, E, w \rangle$ werden als Kräfte interpretiert; Knoten desselben Clusters bündeln ihrer Kräfte. Illustration:

eindeutige Zugehörigkeitsentscheidung (mit Agglomeration):



©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung [Stein/Niggemann 1999]

Die gewichteten Kanten im Graph $G = \langle V, E, w \rangle$ werden als Kräfte interpretiert; Knoten desselben Clusters bündeln ihrer Kräfte. Illustration:

eindeutige Zugehörigkeitsentscheidung (mit Agglomeration):



eindeutige Zugehörigkeitsentscheidung
(mit Cluster-Wechsel):



Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung [Stein/Niggemann 1999]

Die gewichteten Kanten im Graph $G = \langle V, E, w \rangle$ werden als Kräfte interpretiert; Knoten desselben Clusters bündeln ihrer Kräfte. Illustration:

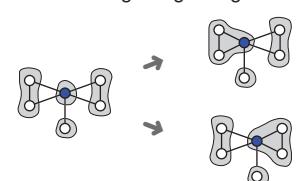
eindeutige Zugehörigkeitsentscheidung (mit Agglomeration):



eindeutige Zugehörigkeitsentscheidung
(mit Cluster-Wechsel):



mehrdeutige Zugehörigkeit:



ML-IX-155 Cluster Analysis

©STEIN 2002-2012

ML-IX-156 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Algorithmus

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d . Distance function for nodes in V .
Output: $\gamma : V \rightarrow \mathbb{N}$. Cluster assignment function.

```

1.
2.
3.
4.

5. FOREACH  $v \in V$  DO
6.    $\gamma^* = \operatorname{argmax}_{i: i \in \{1, \dots, |V|\}} \sum_{\{u,v\}: \{u,v\} \in E \wedge \gamma(u)=i} w(u,v)$ 
7.   IF  $\gamma(v) \neq \gamma^*$  THEN  $\gamma(v) = \gamma^*$ ,  $t = \text{False}$  ENDIF // relabeling
8. ENDDO

9.
10.

```

ML-IX-157 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Algorithmus

Input: $G = \langle V, E, w \rangle$. Weighted graph.
 d . Distance function for nodes in V .
Output: $\gamma : V \rightarrow \mathbb{N}$. Cluster assignment function.

```

1.  $i = 0$ ,  $t = \text{False}$ 
2. FOREACH  $v \in V$  DO  $i = i + 1$ ,  $\gamma(v) = i$  ENDDO
3. UNLESS  $t$  DO
4.    $t = \text{True}$ 

5. FOREACH  $v \in V$  DO
6.    $\gamma^* = \operatorname{argmax}_{i: i \in \{1, \dots, |V|\}} \sum_{\{u,v\}: \{u,v\} \in E \wedge \gamma(u)=i} w(u,v)$ 
7.   IF  $\gamma(v) \neq \gamma^*$  THEN  $\gamma(v) = \gamma^*$ ,  $t = \text{False}$  ENDIF // relabeling
8. ENDDO

9. ENDDO
10. RETURN( $\gamma$ )

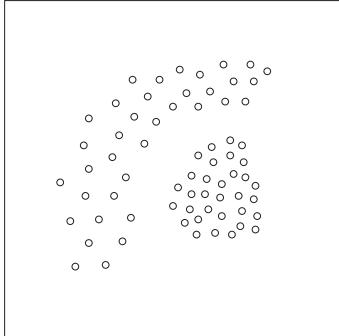
```

ML-IX-158 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

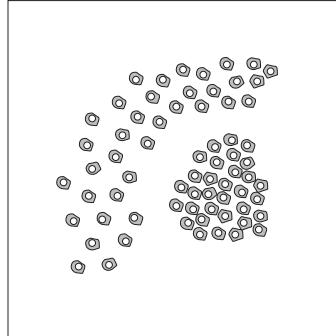


ML-IX-159 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

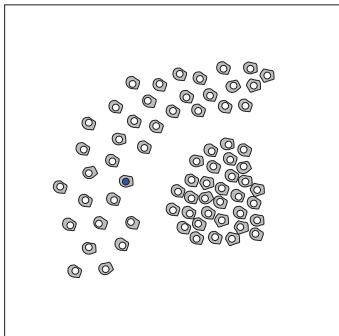


ML-IX-160 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

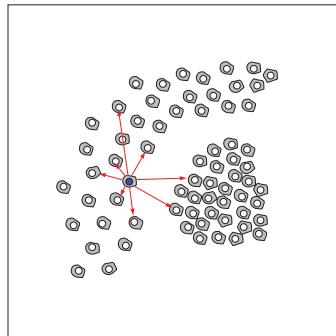


ML-IX-161 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

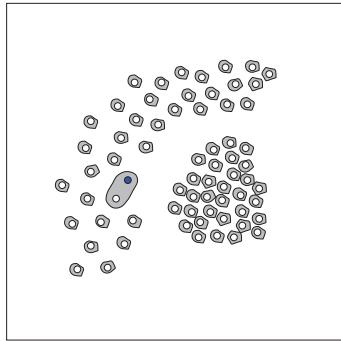


ML-IX-162 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

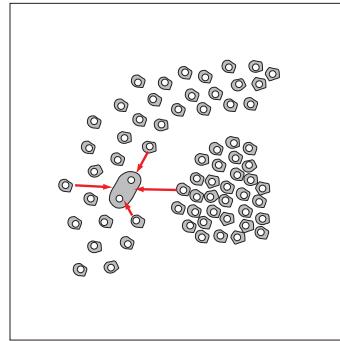


ML-IX-163 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

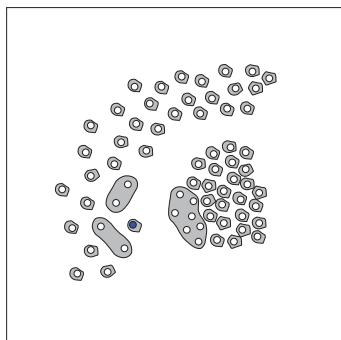


ML-IX-164 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

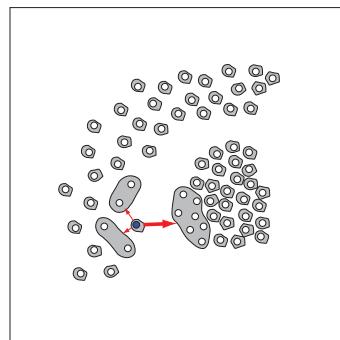


ML-IX-165 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

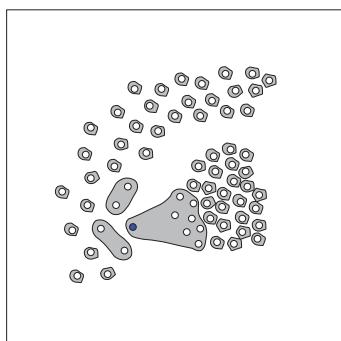


ML-IX-166 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

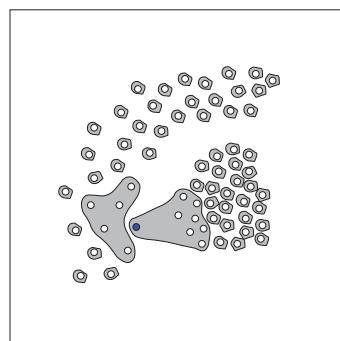


ML-IX-167 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

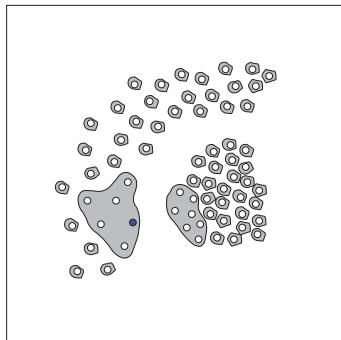


ML-IX-168 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust

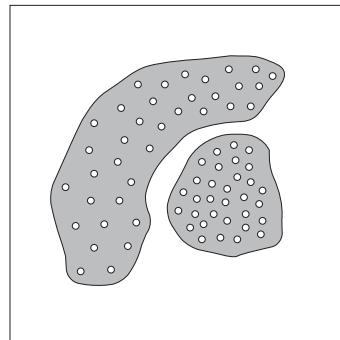


ML-IX-169 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust



ML-IX-170 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung (Fortsetzung)

Jedes $\mathcal{C} = \{C_1, \dots, C_k\}$ induziert k Teilgraphen. MajorClust ist eine Heuristik zur Maximierung des gewichteten partiellen Kantenzusammenhangs, $\Lambda(\mathcal{C})$.

$$\Lambda(\mathcal{C}) = \sum_{i=1}^k |C_i| \cdot \lambda_i$$

ML-IX-171 Cluster Analysis

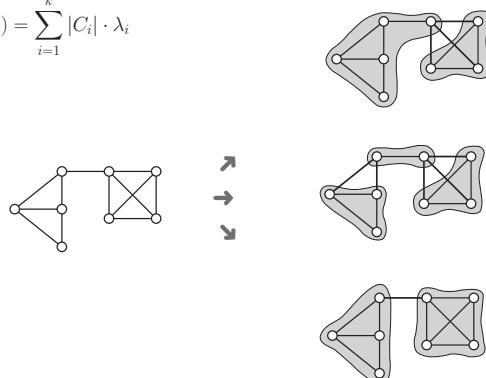
©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung (Fortsetzung)

Jedes $\mathcal{C} = \{C_1, \dots, C_k\}$ induziert k Teilgraphen. MajorClust ist eine Heuristik zur Maximierung des gewichteten partiellen Kantenzusammenhangs, $\Lambda(\mathcal{C})$.

$$\Lambda(\mathcal{C}) = \sum_{i=1}^k |C_i| \cdot \lambda_i$$



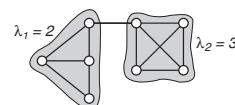
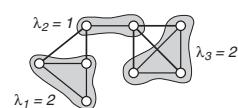
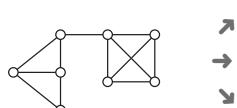
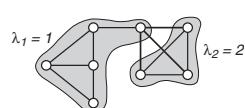
©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung (Fortsetzung)

Jedes $\mathcal{C} = \{C_1, \dots, C_k\}$ induziert k Teilgraphen. MajorClust ist eine Heuristik zur Maximierung des gewichteten partiellen Kantenzusammenhangs, $\Lambda(\mathcal{C})$.

$$\Lambda(\mathcal{C}) = \sum_{i=1}^k |C_i| \cdot \lambda_i$$



ML-IX-173 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung (Fortsetzung)

Jedes $\mathcal{C} = \{C_1, \dots, C_k\}$ induziert k Teilgraphen. MajorClust ist eine Heuristik zur Maximierung des gewichteten partiellen Kantenzusammenhangs, $\Lambda(\mathcal{C})$.

$$\Lambda(\mathcal{C}) = \sum_{i=1}^k |C_i| \cdot \lambda_i$$

$$\lambda_1 = 1 \quad \lambda_2 = 2 \\ \Lambda = 5 \cdot 1 + 3 \cdot 2 = 11$$

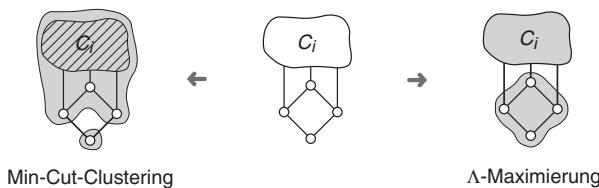
$$\lambda_2 = 1 \quad \lambda_3 = 2 \\ \lambda_1 = 2 \\ \Lambda = 3 \cdot 2 + 2 \cdot 1 + 3 \cdot 2 = 14$$

$$\lambda_1 = 2 \quad \lambda_2 = 3 \\ \Lambda = \Lambda^* = 4 \cdot 2 + 4 \cdot 3 = 20$$

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung (Fortsetzung)

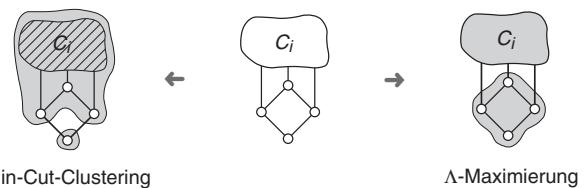


ML-IX-175 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

MajorClust: Prinzip der Dichteschätzung (Fortsetzung)



Satz 1 (Strong Splitting Condition) [Stein/Niggemann 1999]

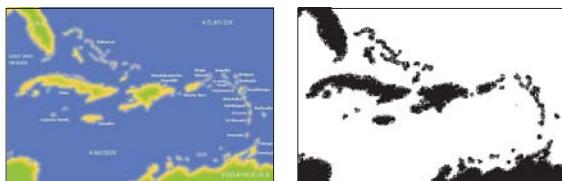
Sei $\mathcal{C} = \{C_1, \dots, C_k\}$ eine Partitionierung eines Graphen $G = \langle V, E, w \rangle$; weiterhin bezeichne $\lambda(G)$ den Kantenzusammenhang von G und $\lambda_1, \dots, \lambda_k$ die Kantenzusammenhänge der von den C_1, \dots, C_k induzierten Subgraphen.

Gilt $\lambda(G) < \min\{\lambda_1, \dots, \lambda_k\}$ (Strong Splitting Condition), so liefert Λ -Maximierung eine Aufteilung am minimalen Cut.

Dichtebasierter Verfahren

DBSCAN versus MajorClust: niedrigdimensionale Daten

Karte der karibischen Inseln, etwa 20.000 Punkte:



ML-IX-177 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN versus MajorClust: niedrigdimensionale Daten

Karte der karibischen Inseln, etwa 20.000 Punkte:



DBSCAN:



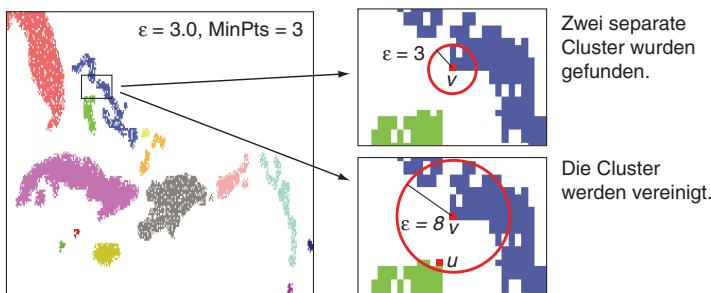
ML-IX-178 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN versus MajorClust: niedrigdimensionale Daten

Problematik geeigneter ε -Werte bei DBSCAN:



ML-IX-179 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN versus MajorClust: niedrigdimensionale Daten

Karte der karibischen Inseln, etwa 20.000 Punkte:



MajorClust:



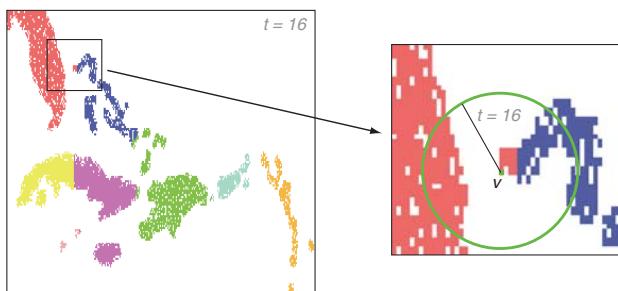
ML-IX-180 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN versus MajorClust: niedrigdimensionale Daten

Problematik der globalen Analyse (keine Beschränkung auf eine ε -Nachbarschaft) bei MajorClust:



ML:IX-181 Cluster Analysis

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN versus MajorClust: hochdimensionale Daten

Dokumentenkategorisierung mit dem Reuters-Korpus:

- ❑ 1000 Dokumente
- ❑ 10 Kategorien: Politik, Kultur, Wirtschaft, etc.
- ❑ die Dokumente sind gleichverteilt, gehören genau zu einer Kategorie
- ❑ Dimension des Merkmalsraums: > 10.000

DBSCAN:

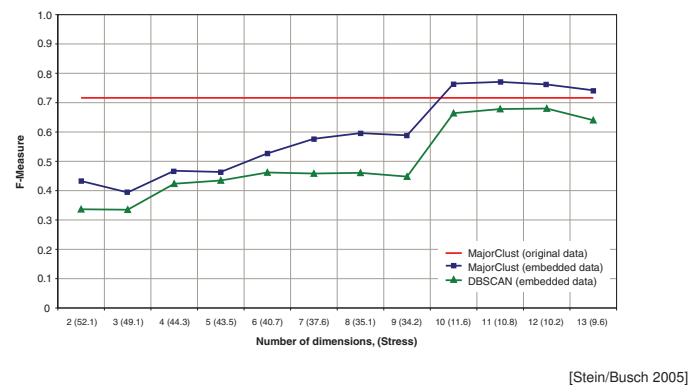
- ❑ degeneriert mit steigender Zahl der Dimensionen
- ❑ Ursache ist die Bestimmung der ε -Nachbarschaft
- ❑ Ausweg ist eine Dimensionreduktion, z. B. eine Einbettung der Daten mittels multidimensionaler Skalierung (MDS)

©STEIN 2002-2012

Dichtebasierter Verfahren

DBSCAN versus MajorClust: hochdimensionale Daten

Klassifikationsergebnisse (F -Measure), aufgetragen über Dimensionalität:



ML:IX-183 Cluster Analysis

©STEIN 2002-2012

Bemerkungen:

- ❑ Das Problem der Nachbarschaftssuche in hochdimensionalen Räumen ist meistens nicht effizient lösbar: Ab Dimensionen größer als 10-20 ist das lineare Durchsuchen aller Merkmalsvektoren effizienter als die Verwendung von hochentwickelten, raumpartitionierenden Datenstrukturen wie R -Tree, X -Tree, Quadtree, KD-Tree, etc. Einen Ausweg bieten die Ansätze wie Locality-Sensitive-Hashing oder Fuzzy-Fingerprinting. Siehe auch: [Weber 99] [Gionis/Indyk/Motwani 99-04] [Stein 05] [Stein/SMZE 05]
- ❑ DBSCAN verwendet zur Bestimmung der ε -Nachbarschaft die R -Tree-Datenstruktur. Diese Datenstruktur leistet einen wesentlichen (wenn nicht den größten) Teil der Cluster-Analyse innerhalb von DBSCAN.
- ❑ Möchte man DBSCAN für hochdimensionale Daten verwenden, ist eine Einbettung der Daten in einen niedrigdimensionalen Raum unvermeidbar. Dabei ist zu bedenken, dass eine Dimensionreduktion durch Einbettung rechenintensiv ist und das gute Laufzeitverhalten von DBSCAN zunicht macht.

©STEIN 2002-2012

Chapter ML:IX (continued)

IX. Clusteranalyse

- ❑ Einordnung Data Mining
- ❑ Einführung in die Clusteranalyse
- ❑ Hierarchische Verfahren
- ❑ Iterative Verfahren
- ❑ Dichtebasierter Verfahren
- ❑ Cluster Evaluation
- ❑ Constrained Cluster Analysis

Cluster Evaluation

Overview

"The validation of clustering structures is the most difficult and frustrating part of cluster analysis. Without a strong effort in this direction, cluster analysis will remain a black art accessible only to those true believers who have experience and great courage."

[Jain/Dubes 1990]

ML:IX-185 Cluster Analysis

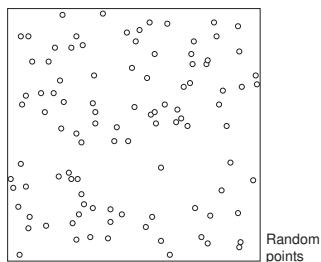
©STEIN 2002-2012

ML:IX-186 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

[from Tan/Steinbach/Kumar 05]



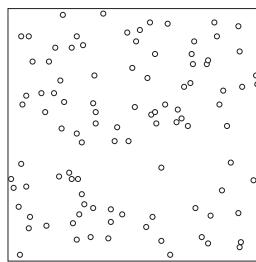
Random points

ML-IX-187 Cluster Analysis

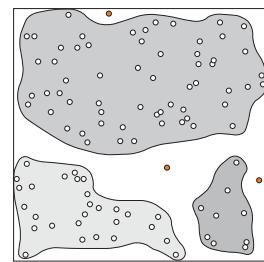
©STEIN 2002-2012

Cluster Evaluation

[from Tan/Steinbach/Kumar 05]



Random points

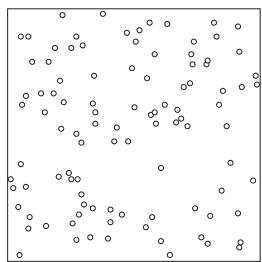


DBSCAN

©STEIN 2002-2012

Cluster Evaluation

[from Tan/Steinbach/Kumar 05]



Random points

ML-IX-189 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Overview

Cluster evaluation can address different issues:

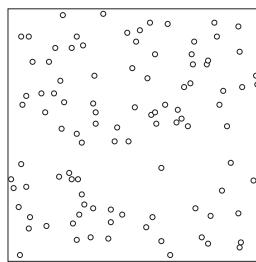
- ❑ Provide evidence whether data contains non-random structures.
- ❑ Relate found structures in the data to externally provided class information.
- ❑ Rank alternative clusterings with regard to their quality.
- ❑ Determine the ideal number of clusters.
- ❑ Provide information to choose a suited clustering approach.

ML-IX-191 Cluster Analysis

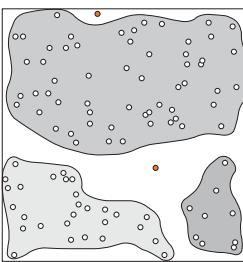
©STEIN 2002-2012

Cluster Evaluation

[from Tan/Steinbach/Kumar 05]



Random points



DBSCAN

©STEIN 2002-2012

Cluster Evaluation

Overview

Cluster evaluation can address different issues:

- ❑ Provide evidence whether data contains non-random structures.
- ❑ Relate found structures in the data to externally provided class information.
- ❑ Rank alternative clusterings with regard to their quality.
- ❑ Determine the ideal number of clusters.
- ❑ Provide information to choose a suited clustering approach.

1. External validity measures:

Analyze how close is a clustering to a reference.

2. Internal validity measures:

Analyze intrinsic characteristics of a clustering.

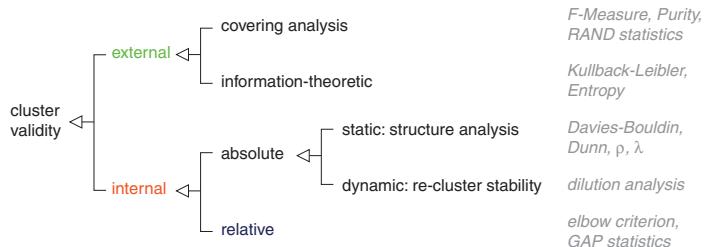
3. Relative validity measures:

Analyze the sensitivity (of internal measures) during clustering generation.

©STEIN 2002-2012

Cluster Evaluation

Overview



1. External validity measures:

Analyze how close is a clustering to a reference.

2. Internal validity measures:

Analyze intrinsic characteristics of a clustering.

3. Relative validity measures:

Analyze the sensitivity (of internal measures) during clustering generation.

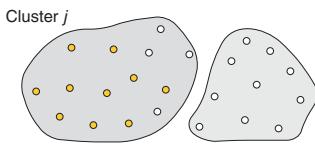
Cluster Evaluation

External Validity Measures: *F*-Measure



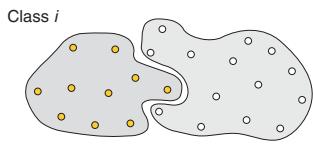
Cluster Evaluation

External Validity Measures: *F*-Measure



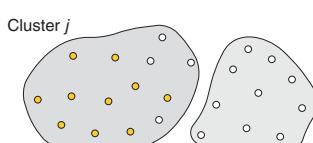
Cluster Evaluation

External Validity Measures: *F*-Measure



Cluster Evaluation

External Validity Measures: *F*-Measure

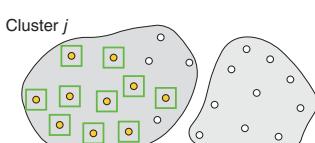


(node-based analysis)

		Truth	
		P	N
Hypothesis	P		
	N		

Cluster Evaluation

External Validity Measures: *F*-Measure

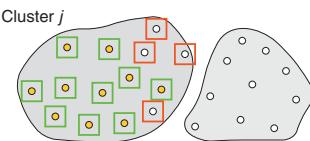


(node-based analysis)

		Truth	
		P	N
Hypothesis	P	TP (a)	
	N		

Cluster Evaluation

External Validity Measures: *F*-Measure

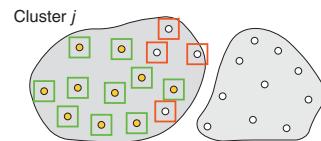


(node-based analysis)

		Truth	
		P	N
Hypothesis	P	TP (a)	FP (b)
	N		

ML-IX-199 Cluster Analysis

©STEIN 2002-2012



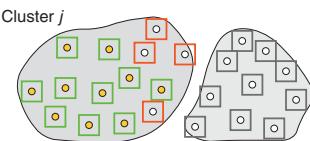
(node-based analysis)

		Truth	
		P	N
Hypothesis	P	TP (a)	FP (b)
	N	FN (c)	

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



(node-based analysis)

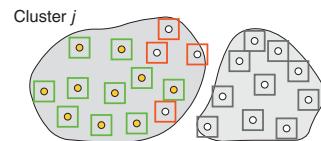
		Truth	
		P	N
Hypothesis	P	TP (a)	FP (b)
	N	FN (c)	TN (d)

ML-IX-201 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



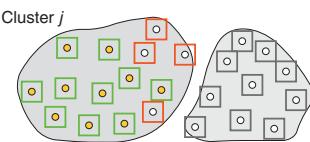
(node-based analysis)

		Truth	
		P	N
Hypothesis	P	TP (a)	FP (b)
	N	FN (c)	TN (d)

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



(node-based analysis)

		Truth	
		P	N
Hypothesis	P	TP (a)	FP (b)
	N	FN (c)	TN (d)

Precision: Recall: *F*-measure:

$$\frac{a}{a+b}$$

$$\frac{a}{a+c}$$

F-measure:

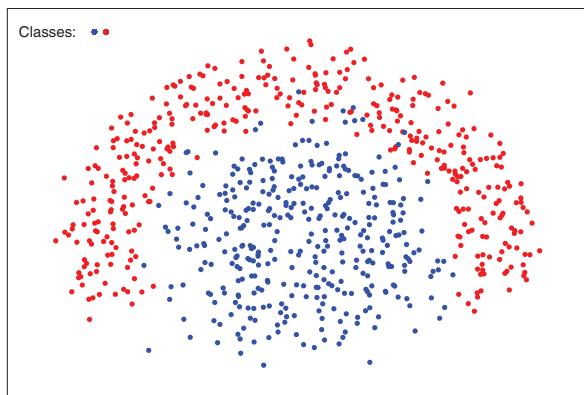
$$F_\alpha = \frac{1 + \alpha}{\frac{1}{precision} + \frac{\alpha}{recall}}$$

$\alpha = 1$
 $\alpha \in (0; 1)$
 $\alpha > 1$

harmonic mean
favor precision over recall
favor recall over precision

Cluster Evaluation

External Validity Measures: *F*-Measure



©STEIN 2002-2012

ML-IX-203 Cluster Analysis

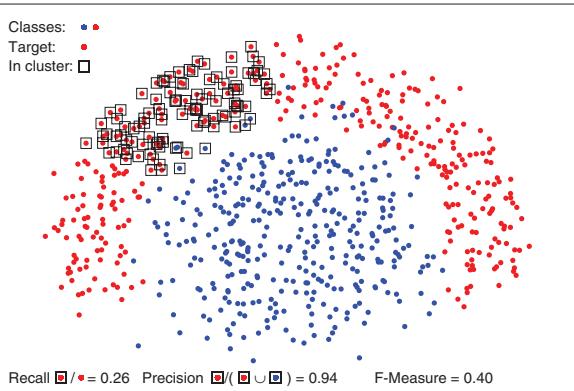
©STEIN 2002-2012

ML-IX-204 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



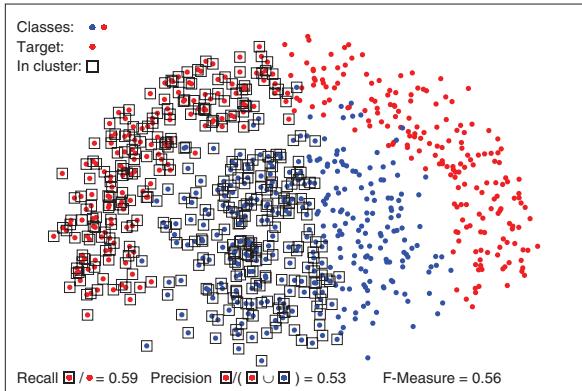
High precision, low recall \Rightarrow low *F*-measure.

ML-IX-205 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



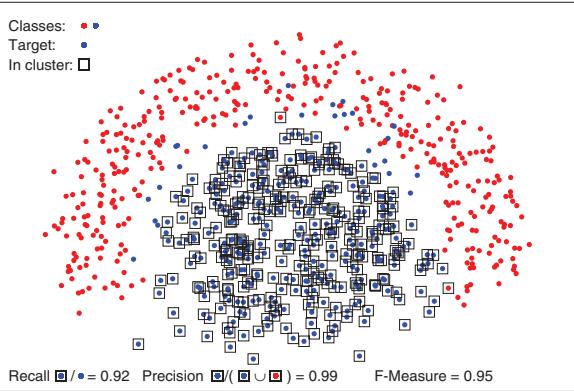
Low precision, low recall \Rightarrow low *F*-measure.

ML-IX-206 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



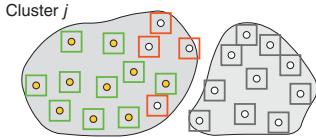
High precision, high recall \Rightarrow high *F*-measure.

ML-IX-207 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



(node-based analysis)

□ Clustering $\mathcal{C} = \{C_1, \dots, C_k\}$ and classification $\mathcal{C}^* = \{C_1^*, \dots, C_l^*\}$ of D .

□ $F_{i,j}$ is the *F*-measure of a cluster j computed *with respect to a class i*.

Recall of cluster j with respect to class i is $|C_j \cap C_i^*| / |C_i^*|$

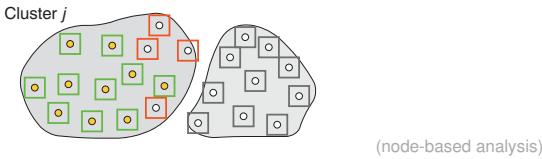
Precision of cluster j with respect to class i is $|C_j \cap C_i^*| / |C_j|$

ML-IX-208 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



(node-based analysis)

□ Clustering $\mathcal{C} = \{C_1, \dots, C_k\}$ and classification $\mathcal{C}^* = \{C_1^*, \dots, C_l^*\}$ of D .

□ $F_{i,j}$ is the *F*-measure of a cluster j computed *with respect to a class i*.

Recall of cluster j with respect to class i is $|C_j \cap C_i^*| / |C_i^*|$

Precision of cluster j with respect to class i is $|C_j \cap C_i^*| / |C_j|$

□ Micro-averaged *F*-measure for $\langle D, \mathcal{C}, \mathcal{C}^* \rangle$:

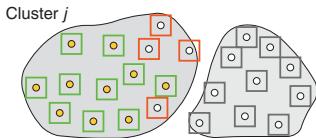
$$F = \sum_{i=1}^l \frac{|C_i^*|}{|D|} \cdot \max_{j=1,\dots,k} \{F_{i,j}\}$$

ML-IX-209 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: *F*-Measure



(node-based analysis)

□ Clustering $\mathcal{C} = \{C_1, \dots, C_k\}$ and classification $\mathcal{C}^* = \{C_1^*, \dots, C_l^*\}$ of D .

□ $F_{i,j}$ is the *F*-measure of a cluster j computed *with respect to a class i*.

Recall of cluster j with respect to class i is $|C_j \cap C_i^*| / |C_i^*|$

Precision of cluster j with respect to class i is $|C_j \cap C_i^*| / |C_j|$

□ Macro-averaged *F*-measure for $\langle D, \mathcal{C}, \mathcal{C}^* \rangle$:

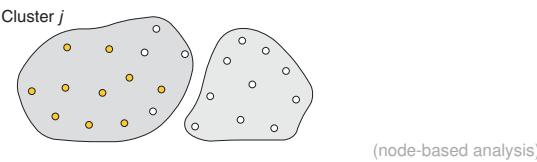
$$F = \frac{1}{l} \sum_{i=1}^l \max_{j=1,\dots,k} \{F_{i,j}\}$$

ML-IX-210 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: Entropy

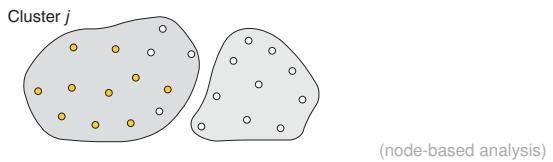


ML-IX-211 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: Entropy

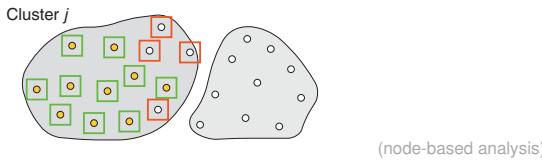


- A cluster C acts as information source \mathcal{L} .

\mathcal{L} emits cluster labels L_1, \dots, L_l with probabilities $P(L_1), \dots, P(L_l)$.

Cluster Evaluation

External Validity Measures: Entropy



- A cluster C acts as information source \mathcal{L} .

\mathcal{L} emits cluster labels L_1, \dots, L_l with probabilities $P(L_1), \dots, P(L_l)$.

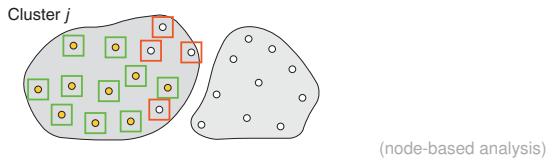
$$\hat{P}(\square) = 10/14, \quad \hat{P}(\square) = 4/14$$

ML-IX-213 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: Entropy



- A cluster C acts as information source \mathcal{L} .

\mathcal{L} emits cluster labels L_1, \dots, L_l with probabilities $P(L_1), \dots, P(L_l)$.

$$\hat{P}(\square) = 10/14, \quad \hat{P}(\square) = 4/14$$

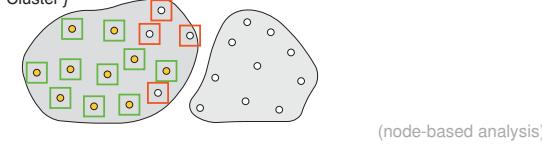
- Entropy of \mathcal{L} :

$$H(\mathcal{L}) = -\sum_{i=1}^l P(L_i) \cdot \log_2(P(L_i))$$

$$\text{Entropy of } C_j \text{ wrt. } \mathcal{C}^* : H(C_j) = -\sum_{C_j \cap C_i^* \neq \emptyset} |C_j \cap C_i^*| / |C_j| \cdot \log_2(|C_j \cap C_i^*| / |C_j|)$$

Cluster Evaluation

External Validity Measures: Entropy



- A cluster C acts as information source \mathcal{L} .

\mathcal{L} emits cluster labels L_1, \dots, L_l with probabilities $P(L_1), \dots, P(L_l)$.

$$\hat{P}(\square) = 10/14, \quad \hat{P}(\square) = 4/14$$

- Entropy of \mathcal{L} :

$$H(\mathcal{L}) = -\sum_{i=1}^l P(L_i) \cdot \log_2(P(L_i))$$

$$\text{Entropy of } C_j \text{ wrt. } \mathcal{C}^* : H(C_j) = -\sum_{C_j \cap C_i^* \neq \emptyset} |C_j \cap C_i^*| / |C_j| \cdot \log_2(|C_j \cap C_i^*| / |C_j|)$$

- Entropy of \mathcal{C} wrt. \mathcal{C}^* :

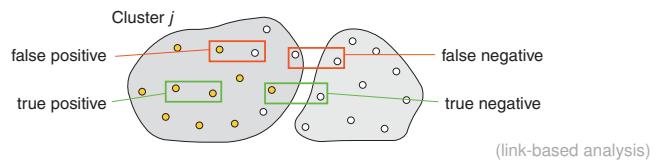
$$H(\mathcal{C}) = \sum_{C_j \in \mathcal{C}} |C_j| / |D| \cdot H(C_j)$$

ML-IX-215 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

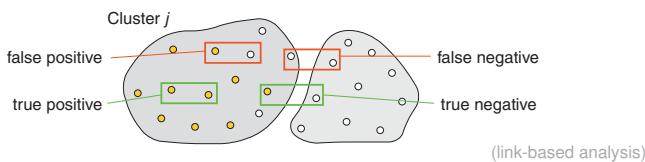
External Validity Measures: Rand, Jaccard



©STEIN 2002-2012

Cluster Evaluation

External Validity Measures: Rand, Jaccard



$$\square R(\mathcal{C}) = \frac{|TP| + |TN|}{|TP| + |TN| + |FP| + |FN|} = \frac{|TP| + |TN|}{n(n-1)/2}, \quad \text{with } n = |D|$$

$$\square J(\mathcal{C}) = \frac{|TP|}{|TP| + |FP| + |FN|}$$

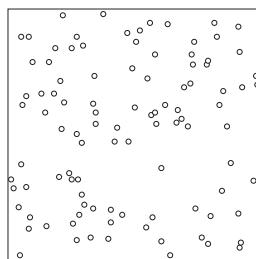
ML-IX-217 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]

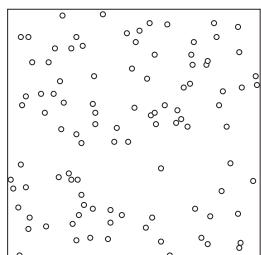


ML-IX-218 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



$$\begin{pmatrix} 1.0 & 0.2 & 0.1 & 0.3 & \dots & 0.1 & 0.0 \\ - & 1.0 & 0.1 & 0.0 & \dots & 0.0 & 0.2 \\ & - & - & - & - & \vdots & \\ & - & - & - & - & - & 1.0 & 0.6 \\ & - & - & - & - & - & - & 1.0 \end{pmatrix}$$

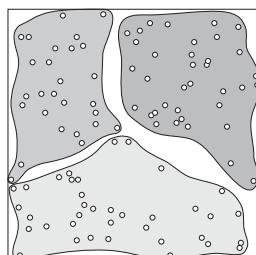
ML-IX-219 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



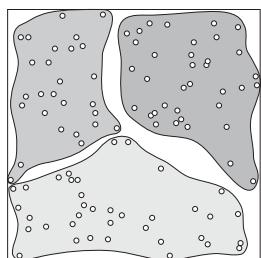
$$\begin{pmatrix} 1.0 & 0.2 & 0.1 & 0.3 & \dots & 0.1 & 0.0 \\ - & 1.0 & 0.1 & 0.0 & \dots & 0.0 & 0.2 \\ & - & - & - & - & \vdots & \\ & - & - & - & - & - & 1.0 & 0.6 \\ & - & - & - & - & - & - & 1.0 \end{pmatrix}$$

ML-IX-220 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



$$\begin{pmatrix} 1.0 & 0.2 & 0.1 & 0.3 & \dots & 0.1 & 0.0 \\ - & 1.0 & 0.1 & 0.0 & \dots & 0.0 & 0.2 \\ & - & - & - & - & \vdots & \\ & - & - & - & - & - & 1.0 & 0.6 \\ & - & - & - & - & - & - & 1.0 \end{pmatrix} \sim \begin{pmatrix} 1 & 0 & 0 & 1 & \dots & 0 & 0 \\ - & 1 & 0 & 0 & \dots & 0 & 1 \\ & - & - & - & - & \vdots & \\ & - & - & - & - & - & 1 & 1 \\ & - & - & - & - & - & - & 1 \end{pmatrix}$$

- ❑ Construct occurrence matrix based on cluster analysis.
- ❑ Compare similarity matrix to occurrence matrix: correlation τ

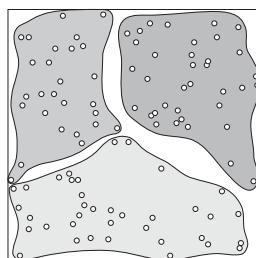
ML-IX-221 Cluster Analysis

©STEIN 2002-2012

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



$$\begin{pmatrix} 1.0 & 0.2 & 0.1 & 0.3 & \dots & 0.1 & 0.0 \\ - & 1.0 & 0.1 & 0.0 & \dots & 0.0 & 0.2 \\ & - & - & - & - & \vdots & \\ & - & - & - & - & - & 1.0 & 0.6 \\ & - & - & - & - & - & - & 1.0 \end{pmatrix} \sim \begin{pmatrix} 1 & 0 & 0 & 1 & \dots & 0 & 0 \\ - & 1 & 0 & 0 & \dots & 0 & 1 \\ & - & - & - & - & \vdots & \\ & - & - & - & - & - & 1 & 1 \\ & - & - & - & - & - & - & 1 \end{pmatrix}$$

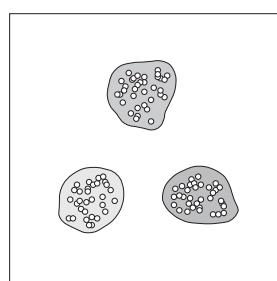
- ❑ Construct occurrence matrix based on cluster analysis.
- ❑ Compare similarity matrix to occurrence matrix: correlation τ

ML-IX-222 Cluster Analysis

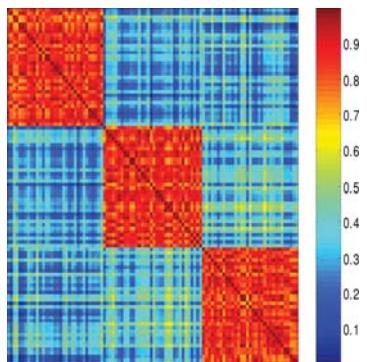
©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



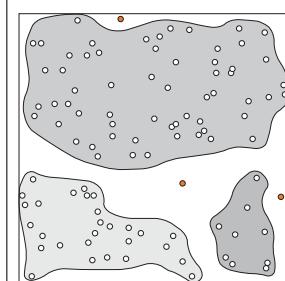
k-means at structured data.



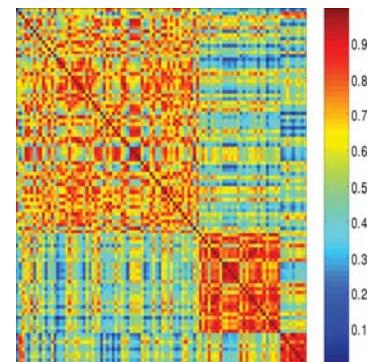
Similarity matrix sorted by cluster label.

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



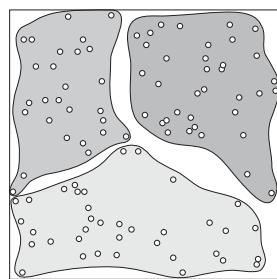
DBSCAN at random data.



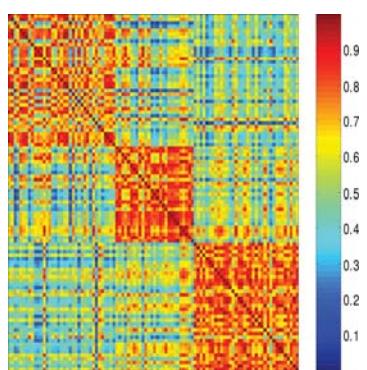
Similarity matrix sorted by cluster label.

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



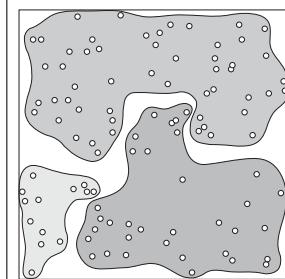
k-means at random data.



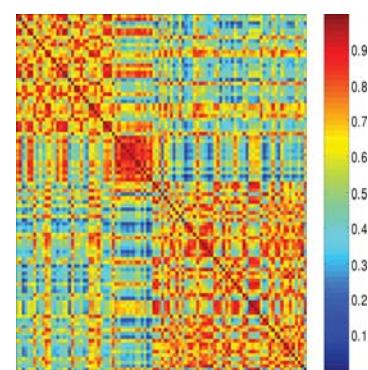
Similarity matrix sorted by cluster label.

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



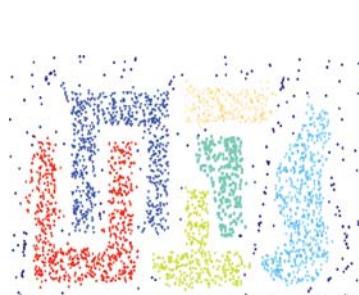
Complete link at random data.



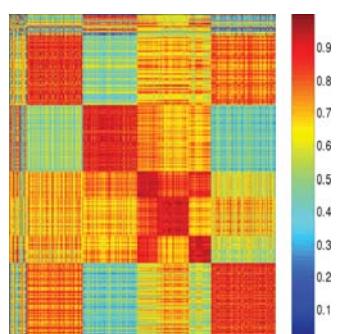
Similarity matrix sorted by cluster label.

Cluster Evaluation

Internal Validity Measures: Link Correlation [from Tan/Steinbach/Kumar 05]



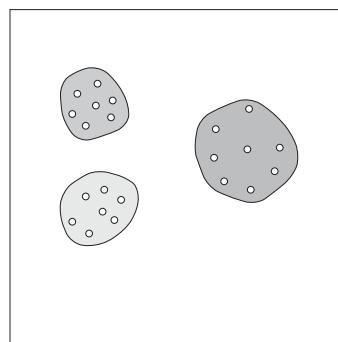
DBSCAN at structured data.



Similarity matrix sorted by cluster label.

Cluster Evaluation

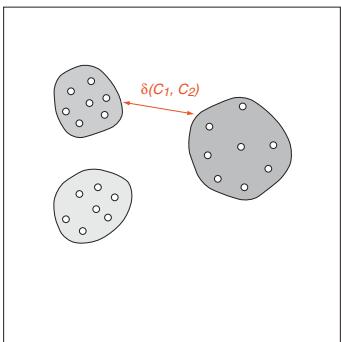
Internal Validity Measures: Structural Analysis



- Distance between two clusters, $\delta(C_1, C_2)$.
- Diameter of a cluster, $\Delta(C)$.
- Scatter within a cluster, $\sigma^2(C)$, SSE.

Cluster Evaluation

Internal Validity Measures: Structural Analysis



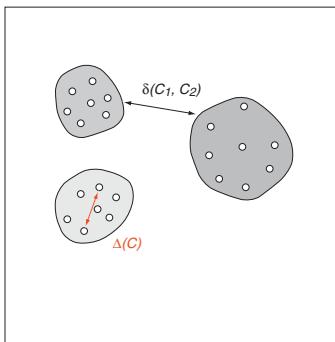
- Distance between two clusters, $\delta(C_1, C_2)$.
- Diameter of a cluster, $\Delta(C)$.
- Scatter within a cluster, $\sigma^2(C)$, SSE.

ML-IX-229 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Structural Analysis



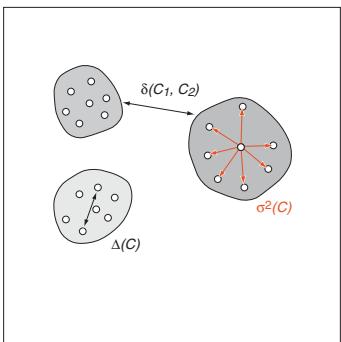
- Distance between two clusters, $\delta(C_1, C_2)$.
- Diameter of a cluster, $\Delta(C)$.
- Scatter within a cluster, $\sigma^2(C)$, SSE.

ML-IX-230 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Structural Analysis



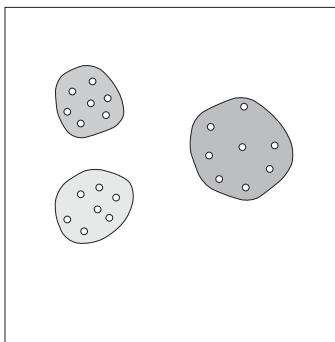
- Distance between two clusters, $\delta(C_1, C_2)$.
- Diameter of a cluster, $\Delta(C)$.
- Scatter within a cluster, $\sigma^2(C)$, SSE.

ML-IX-231 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Dunn Index



$$I(\mathcal{C}) = \frac{\min_{i \neq j} \{\delta(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}},$$

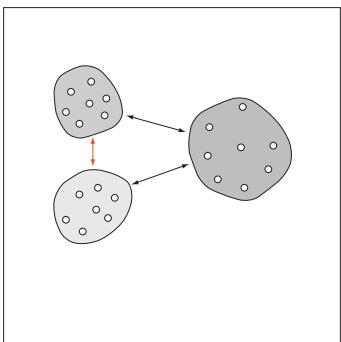
$$I(\mathcal{C}) \rightarrow \max$$

ML-IX-232 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Dunn Index



$$I(\mathcal{C}) = \frac{\min_{i \neq j} \{\delta(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}},$$

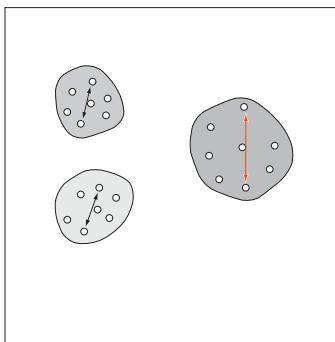
$$I(\mathcal{C}) \rightarrow \max$$

ML-IX-233 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Dunn Index



$$I(\mathcal{C}) = \frac{\min_{i \neq j} \{\delta(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}},$$

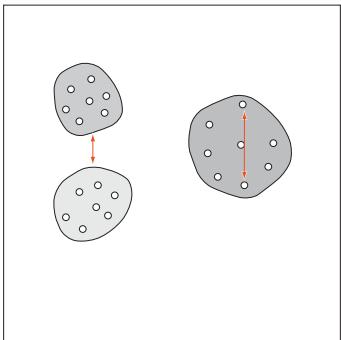
$$I(\mathcal{C}) \rightarrow \max$$

ML-IX-234 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Dunn Index



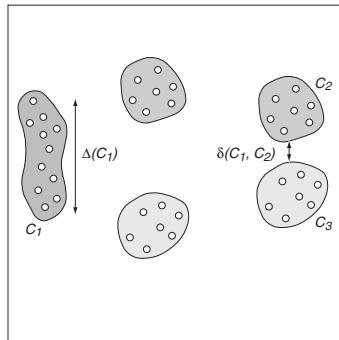
$$I(\mathcal{C}) = \frac{\min_{i \neq j} \{\delta(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}},$$

$$I(\mathcal{C}) \rightarrow \max$$

- Dunn is susceptible to noise.
- Dunn is biased towards the worst substructure in a clustering (cf. min)
- Dunn cannot put distances and diameters into relation.

Cluster Evaluation

Internal Validity Measures: Dunn Index



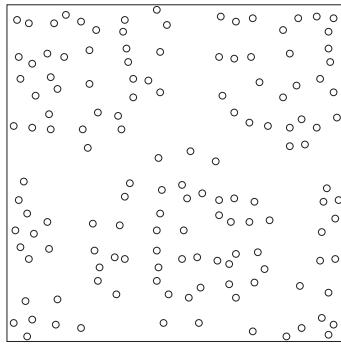
$$I(\mathcal{C}) = \frac{\min_{i \neq j} \{\delta(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}},$$

$$I(\mathcal{C}) \rightarrow \max$$

- Dunn is susceptible to noise.
- Dunn is biased towards the worst substructure in a clustering (cf. min)
- Dunn cannot put distances and diameters into relation.

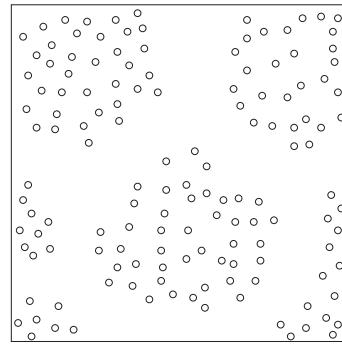
Cluster Evaluation

Internal Validity Measures: Expected Density ρ [Stein/Meyer zu Eissen 2007]



Cluster Evaluation

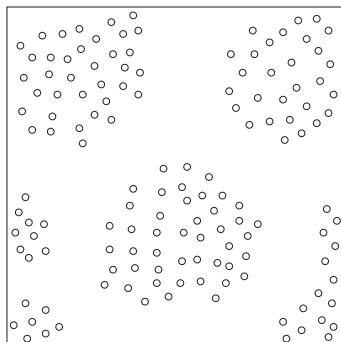
Internal Validity Measures: Expected Density ρ [Stein/Meyer zu Eissen 2007]



Different retrieval models yield different similarity graphs.

Cluster Evaluation

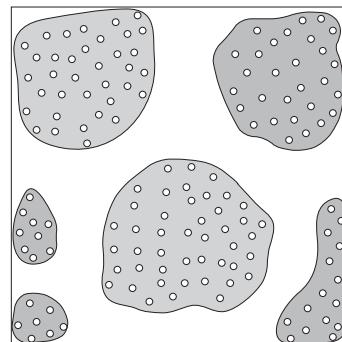
Internal Validity Measures: Expected Density ρ [Stein/Meyer zu Eissen 2007]



Different retrieval models yield different similarity graphs.

Cluster Evaluation

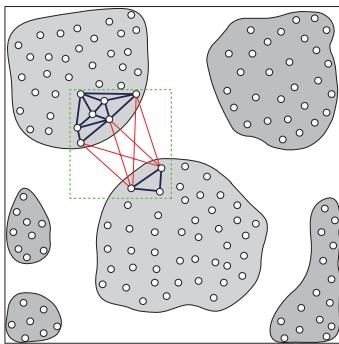
Internal Validity Measures: Expected Density ρ



Compare (for alternative clusterings) the similarity density within the clusters to the average similarity of the entire graph.

Cluster Evaluation

Internal Validity Measures: Expected Density ρ



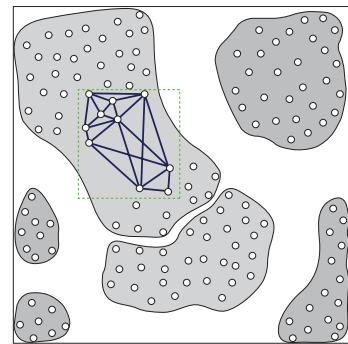
Compare (for alternative clusterings) the similarity density within the clusters to the average similarity of the entire graph.

ML-IX-241 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Expected Density ρ



Compare (for alternative clusterings) the similarity density within the clusters to the average similarity of the entire graph.

ML-IX-242 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Expected Density ρ

Graph $G = \langle V, E \rangle$

- G is called sparse [dense] if $|E| = O(|V|)$ [$O(|V|^2)$]
- the density θ computes from the equation $|E| = |V|^\theta$

ML-IX-243 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Expected Density ρ

Graph $G = \langle V, E \rangle$

- G is called sparse [dense] if $|E| = O(|V|)$ [$O(|V|^2)$]
- the density θ computes from the equation $|E| = |V|^\theta$

Similarity graph $G = \langle V, E, w \rangle$, $|E| \sim w(G) := \sum_{e \in E} w(e)$

- the density θ computes from the equation $w(G) = |V|^\theta$

©STEIN 2002-2012

Cluster Evaluation

Internal Validity Measures: Expected Density ρ

Graph $G = \langle V, E \rangle$

- G is called sparse [dense] if $|E| = O(|V|)$ [$O(|V|^2)$]
- the density θ computes from the equation $|E| = |V|^\theta$

Similarity graph $G = \langle V, E, w \rangle$, $|E| \sim w(G) := \sum_{e \in E} w(e)$

- the density θ computes from the equation $w(G) = |V|^\theta$

Induced subgraph G_i for class C_i

- the expected density ρ compares class C_i to the density average in D

$$\rho(G_i) = \frac{w(G_i)}{|V_i|^\theta}$$

ML-IX-245 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Relative Validity Measures: Elbow Criterion

1. Hyperparameters of a clustering algorithm: p_1, \dots, p_m
 - number of centroids for k -means
 - stopping level for hierarchical algorithms
 - neighborhood size for DBSCAN
2. Clusterings $\mathcal{C} = \{C_{p_1}, \dots, C_{p_m}\}$ associated with p_1, \dots, p_m .
3. Points of an error curve $\{(p_i, e(C_{p_i})) \mid i = 1, \dots, m\}$.

©STEIN 2002-2012

Cluster Evaluation

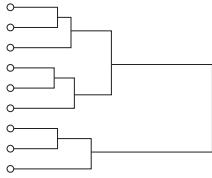
Relative Validity Measures: Elbow Criterion

- Hyperparameters of a clustering algorithm: p_1, \dots, p_m

- number of centroids for k -means
- stopping level for hierarchical algorithms
- neighborhood size for DBSCAN

- Clusterings $\mathcal{C} = \{\mathcal{C}_{p_1}, \dots, \mathcal{C}_{p_m}\}$ associated with p_1, \dots, p_m .

- Points of an error curve $\{(p_i, e(\mathcal{C}_{p_i})) \mid i = 1, \dots, m\}$.

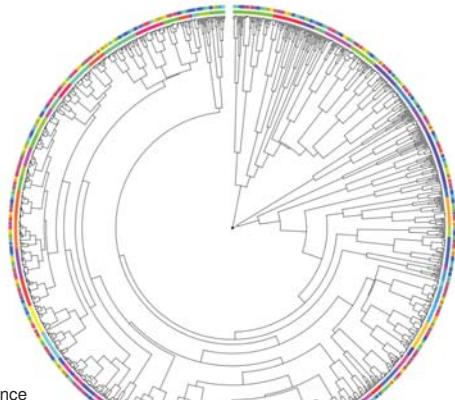


ML-IX-247 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Relative Validity Measures: Elbow Criterion



d_C : Hamming distance
Merging: complete link

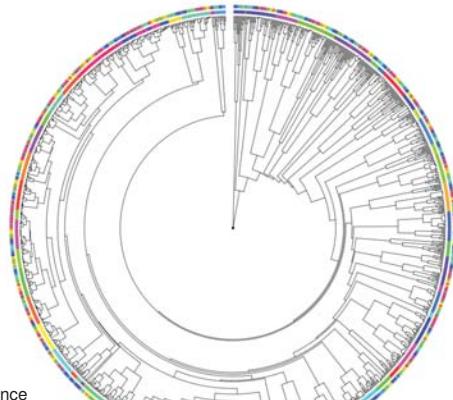
<http://cs.jhu.edu/~razvanm/fs-expedition/2.6.x.html>
Relations between 1377 file systems for Linux Kernel 2.6.0. [Razvan Musaloiu-E. 2009]

ML-IX-249 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Relative Validity Measures: Elbow Criterion



d_C : Hamming distance
Merging: group average link

<http://cs.jhu.edu/~razvanm/fs-expedition/2.6.x.html>
Relations between 1377 file systems for Linux Kernel 2.6.0. [Razvan Musaloiu-E. 2009]

ML-IX-250 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Correlation between External and Internal Measures

In the wild, we are not given a reference classification.

- An external evaluation is not possible.
(though many papers report on such experiments)
- Resort to an internal evaluation.
(connectivity, squared error sums, distance-diameter heuristics, etc.)

Cluster Evaluation

Correlation between External and Internal Measures

In the wild, we are not given a reference classification.

- An external evaluation is not possible.
(though many papers report on such experiments)
- Resort to an internal evaluation.
(connectivity, squared error sums, distance-diameter heuristics, etc.)

"To which extent can an internal evaluation be used to predict for a clustering its distance from the best reference classification—say, to predict the F-measure?"

$$\operatorname{argmax}_{\phi} \{ \tau(X, Y) \mid x = F(\mathcal{C}), y = \phi(\mathcal{C}), \mathcal{C} \in \mathcal{C} \}$$

[Stein/Meyer zu Eissen 2007]

ML-IX-251 Cluster Analysis

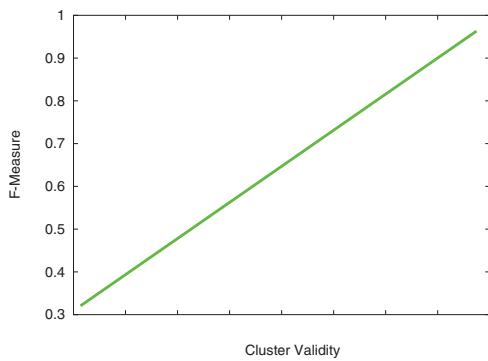
©STEIN 2002-2012

ML-IX-252 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Correlation between External and Internal Measures



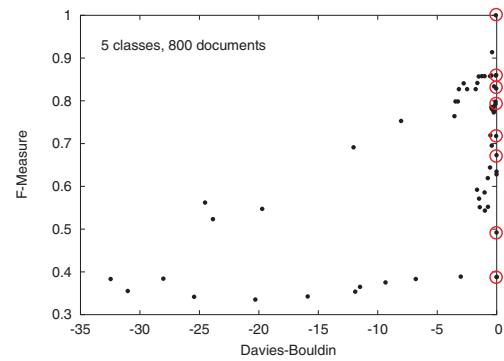
Perfect correlation (desired).

ML-IX-253 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Correlation between External and Internal Measures



$$\text{Davies-Bouldin: } \frac{1}{k} \cdot \sum_{i=1}^k \max_j \frac{s(C_i) + s(C_j)}{\delta(C_i, C_j)}$$

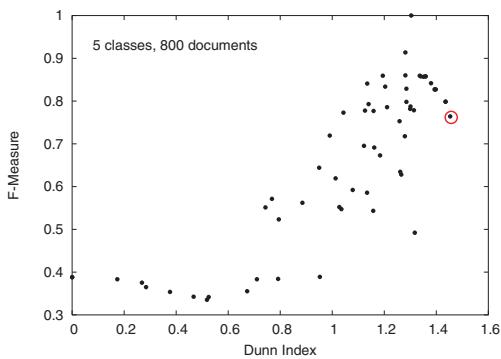
Prefers spherical clusters.

ML-IX-254 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Correlation between External and Internal Measures



$$\text{Dunn Index: } \frac{\min_{i \neq j} \{\delta(C_i, C_j)\}}{\max_{1 \leq l \leq k} \{\Delta(C_l)\}}$$

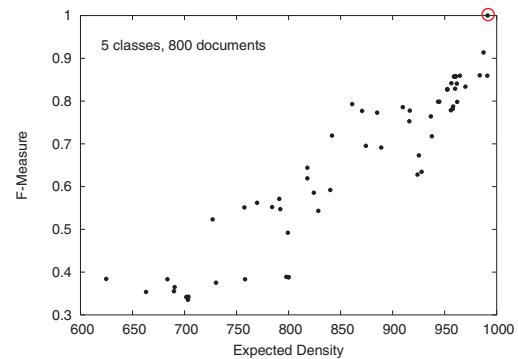
Maximizes dilatation = inter/intra-cluster-diameter.

ML-IX-255 Cluster Analysis

©STEIN 2002-2012

Cluster Evaluation

Correlation between External and Internal Measures



$$\text{Expected Density: } \bar{\rho} = \sum_{i=1}^k \frac{|V_i|}{|V|} \cdot \frac{w(G_i)}{|V_i|^\theta}$$

Independent of cluster forms and sizes.

ML-IX-256 Cluster Analysis

©STEIN 2002-2012

Chapter ML:IX (continued)

IX. Clusteranalyse

- ❑ Einordnung Data Mining
- ❑ Einführung in die Clusteranalyse
- ❑ Hierarchische Verfahren
- ❑ Iterative Verfahren
- ❑ Dichtebasierte Verfahren
- ❑ Cluster Evaluation
- ❑ Constrained Cluster Analysis

Constrained Cluster Analysis

Person Resolution Task



ML-IX-257 Cluster Analysis

©STEIN 2002-2012

ML-IX-258 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

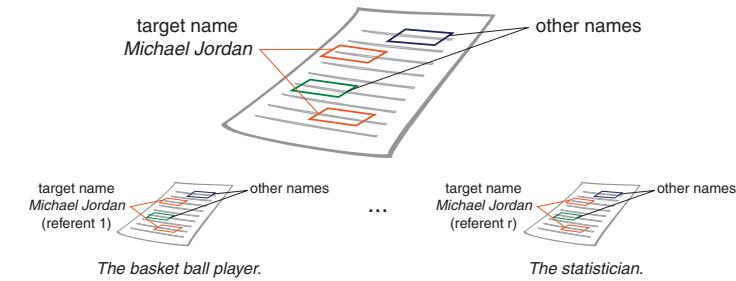
Person Resolution Task



ML-IX-259 Cluster Analysis

Constrained Cluster Analysis

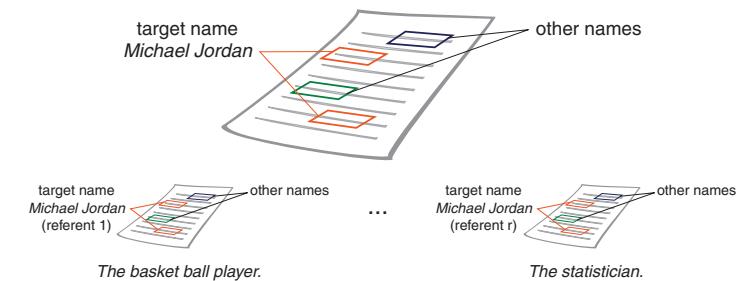
Person Resolution Task



ML-IX-261 Cluster Analysis

Constrained Cluster Analysis

Person Resolution Task



❑ Facts about the Spock data mining challenge:

Target names: $|T| = 44$
 Refers: $|R| = 1101$
 Documents: $|D_{train}| = 27\,000$ (labeled $\approx 2.3\text{GB}$)
 $|D_{test}| = 75\,000$ (unlabeled $\approx 7.8\text{GB}$)

ML-IX-263 Cluster Analysis

Constrained Cluster Analysis

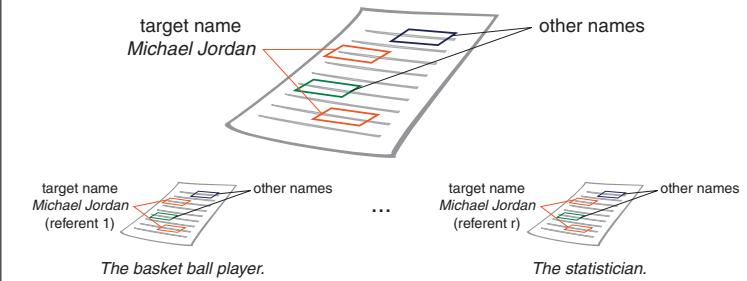
Person Resolution Task



©STEIN 2002-2012

Constrained Cluster Analysis

Person Resolution Task

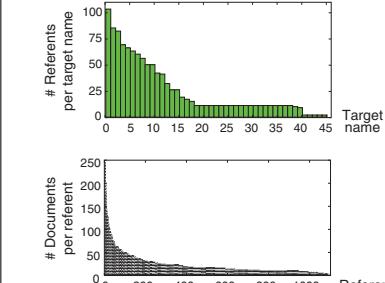


ML-IX-262 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

Person Resolution Task



- ❑ up to 105 referents for a single target name
- ❑ about 25 referents on average per target name

- ❑ about 23 documents on average per referent

❑ Facts about the Spock data mining challenge:

Target names: $|T| = 44$
 Refers: $|R| = 1101$
 Documents: $|D_{train}| = 27\,000$ (labeled $\approx 2.3\text{GB}$)
 $|D_{test}| = 75\,000$ (unlabeled $\approx 7.8\text{GB}$)

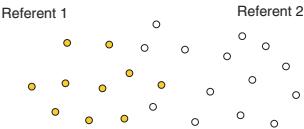
©STEIN 2002-2012

ML-IX-264 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

Applied to Multi-Document Resolution



1. Model similarities → new and established retrieval models:

- global and context-based vector space models
- explicit semantic analysis
- ontology alignment

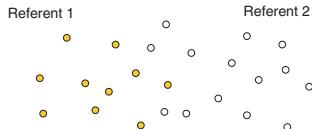
2. Learn class memberships (supervised) → logistic regression

3. Find equivalence classes (unsupervised) → cluster analysis:

- (a) adaptive graph thinning
- (b) multiple, density-based cluster analysis
- (c) clustering selection by expected density maximization

Constrained Cluster Analysis

Applied to Multi-Document Resolution



1. Model similarities → new and established retrieval models:

- global and context-based vector space models
- explicit semantic analysis
- ontology alignment

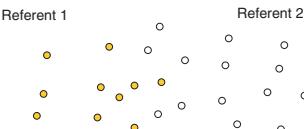
2. Learn class memberships (supervised) → logistic regression

3. Find equivalence classes (unsupervised) → cluster analysis:

- (a) adaptive graph thinning
- (b) multiple, density-based cluster analysis
- (c) clustering selection by expected density maximization

Constrained Cluster Analysis

Applied to Multi-Document Resolution



1. Model similarities → new and established retrieval models:

- global and context-based vector space models
- explicit semantic analysis
- ontology alignment

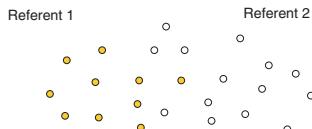
2. Learn class memberships (supervised) → logistic regression

3. Find equivalence classes (unsupervised) → cluster analysis:

- (a) adaptive graph thinning
- (b) multiple, density-based cluster analysis
- (c) clustering selection by expected density maximization

Constrained Cluster Analysis

Applied to Multi-Document Resolution



1. Model similarities → new and established retrieval models:

- global and context-based vector space models
- explicit semantic analysis
- ontology alignment

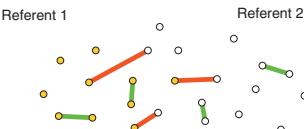
2. Learn class memberships (supervised) → logistic regression

3. Find equivalence classes (unsupervised) → cluster analysis:

- (a) adaptive graph thinning
- (b) multiple, density-based cluster analysis
- (c) clustering selection by expected density maximization

Constrained Cluster Analysis

Applied to Multi-Document Resolution



1. Model similarities → new and established retrieval models:

- global and context-based vector space models
- explicit semantic analysis
- ontology alignment

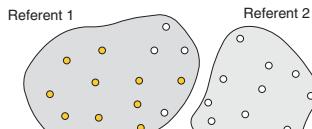
2. Learn class memberships (supervised) → logistic regression

3. Find equivalence classes (unsupervised) → cluster analysis:

- (a) adaptive graph thinning
- (b) multiple, density-based cluster analysis
- (c) clustering selection by expected density maximization

Constrained Cluster Analysis

Applied to Multi-Document Resolution



1. Model similarities → new and established retrieval models:

- global and context-based vector space models
- explicit semantic analysis
- ontology alignment

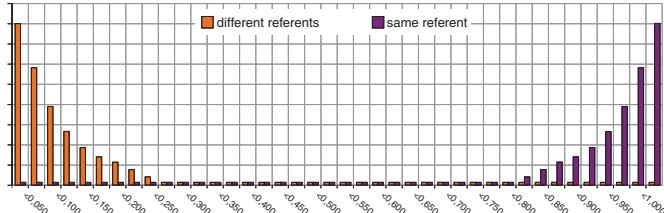
2. Learn class memberships (supervised) → logistic regression

3. Find equivalence classes (unsupervised) → cluster analysis:

- (a) adaptive graph thinning
- (b) multiple, density-based cluster analysis
- (c) clustering selection by expected density maximization

Constrained Cluster Analysis

Idealized Class Membership Distribution over Similarities



Similarity distributions for document pairs from **different referents** and **same referent**.

Logistic regression task:

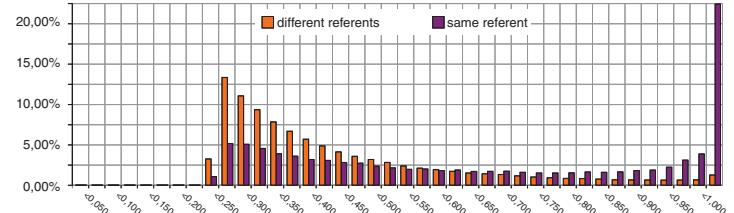
- sample size: 400 000
- classes imbalance: **non-target class : target class** $\approx 25:1$
- items are drawn uniformly distributed wrt. non-targets and targets
- items are uniformly distributed over the groups of target names

ML-IX-271 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

Membership Distribution under *tf-idf* Vector Space Model



Model details:

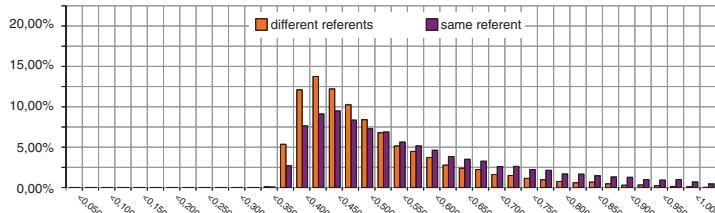
- corpus size: 25 000 documents
- dictionary size: 1,2 Mio terms
- stopwords number: 850
- stopword volume: 36%

ML-IX-272 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

Membership Distribution under Context-Based Vector Space Model



Model details:

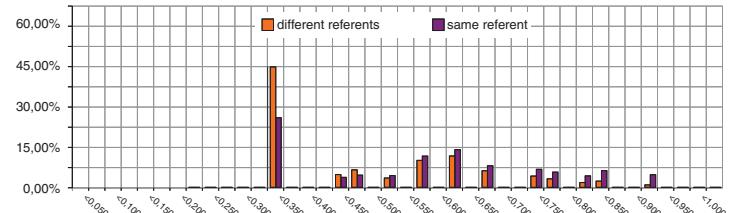
- corpus size: 25 000 documents
- dictionary size: 1,2 Mio terms
- stopwords number: 850
- stopword volume: 36%

ML-IX-273 Cluster Analysis

©STEIN 2002-2012

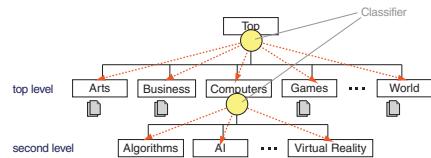
Constrained Cluster Analysis

Membership Distribution under Ontology Alignment Model



Model details:

- DMOZ open directory project
- > 5 million documents
- 12 top-level categories
- 31 second level categories
- ML: hierarchical Bayes
- training set: 100 000 pages



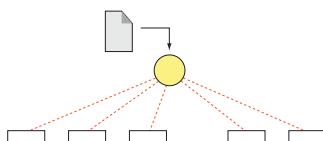
ML-IX-274 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

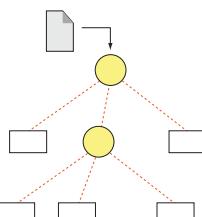
In-Depth: Multi-Class Hierarchical Classification

Flat (big-bang) classification



- + simple realization
- loss of discriminative power with increasing number of categories

Hierarchical (top-down) classification



- + specialized classifiers (divide and conquer)
- misclassification at higher levels can never become repaired

ML-IX-275 Cluster Analysis

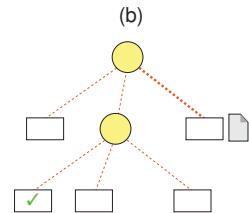
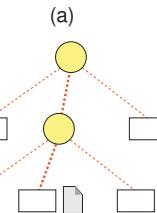
©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Multi-Class Hierarchical Classification

State of the art of effectiveness analyses:

1. independence assumption between categories
2. neglection of both hierarchical structure and degree of misclassification



Improvements:

- Consider similarity $\varphi(C_i, C_j)$ between correct and wrong category.
- Consider graph distance $d(C_i, C_j)$ between correct and wrong category.

ML-IX-276 Cluster Analysis

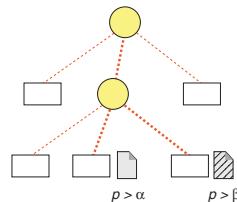
©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Multi-Class Hierarchical Classification

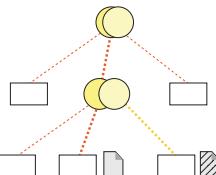
Improvements continued:

Multi-label (multi path) classification



- ❑ traverse more than one path and return all labels
- ❑ employ probabilistic classifiers with a threshold: split a path or not

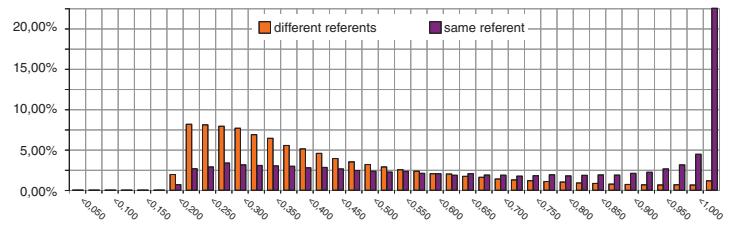
Multi-classifier (ensemble) classification



- ❑ classification result is a majority decision
- ❑ employ different classifier (different types or differently parameterized)

Constrained Cluster Analysis

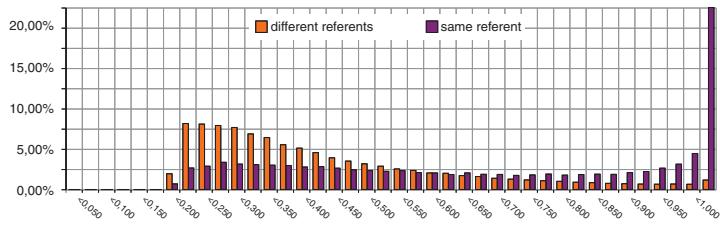
Membership Distribution under Optimized Retrieval Model Combination



Retrieval Model	$F_{1/3}$ -Measure
tfidf vector space	0.39
context-based vector space	0.32
ESA Wikipedia persons	0.30
phrase structure grammar	0.17
ontology alignment	0.15
optimized combination	0.42

Constrained Cluster Analysis

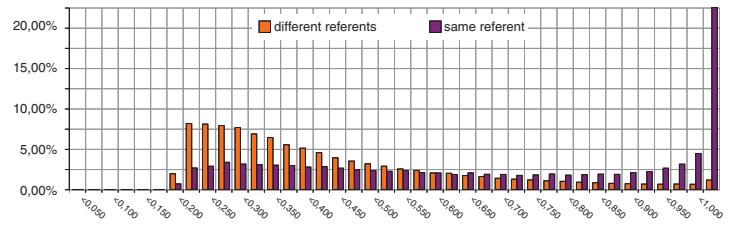
Membership Distribution under Optimized Retrieval Model Combination



Retrieval Model	$F_{1/3}$ -Measure
tfidf vector space	0.39
context-based vector space	0.32
ESA Wikipedia persons	0.30
phrase structure grammar	0.17
ontology alignment	0.15
optimized combination	0.42

Constrained Cluster Analysis

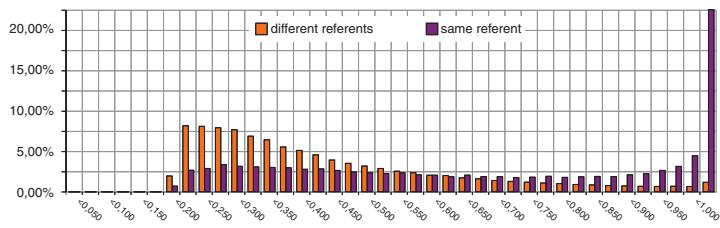
Membership Distribution under Optimized Retrieval Model Combination



Retrieval Model	$F_{1/3}$ -Measure
tfidf vector space	0.39
context-based vector space	0.32
ESA Wikipedia persons	0.30
phrase structure grammar	0.17
ontology alignment	0.15
optimized combination	0.42

Constrained Cluster Analysis

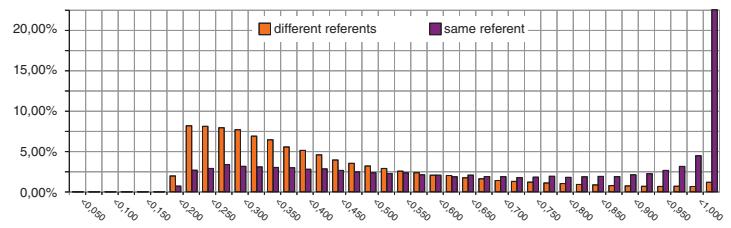
Membership Distribution under Optimized Retrieval Model Combination



Retrieval Model	$F_{1/3}$ -Measure
tfidf vector space	0.39
context-based vector space	0.32
ESA Wikipedia persons	0.30
phrase structure grammar	0.17
ontology alignment	0.15
optimized combination	0.42

Constrained Cluster Analysis

Membership Distribution under Optimized Retrieval Model Combination



Referent 1	Referent 2	Referent m
○ ○	○ ○	○ ○
○ ○	○ ○	○ ○
...
○ ○	○ ○	○ ○

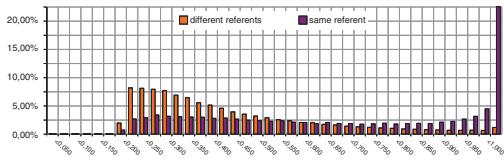
In the example:

- ❑ precision = 0.4
- ❑ recall = 0.43
- ❑ $F_{1/3} = 0.41$

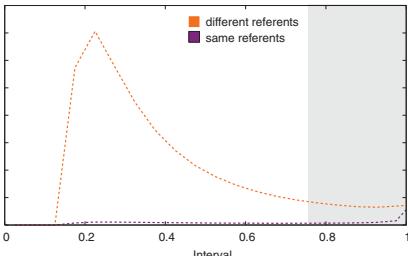
(if false negatives are uniformly distributed)

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Consideration of imbalance:



ML-IX-283 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

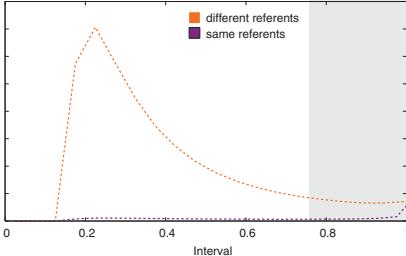
In-Depth: Analysis of Classifier Effectiveness

□ class imbalance factor (*CIF*) of 25

⇒ precision in interval [0.725; 1] for edges between same referents: ≈ 0.17

How can $F_{1/3} = 0.42$ be achieved via cluster analysis?

Consideration of imbalance:

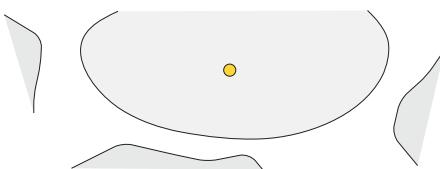


ML-IX-284 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Assumption: uniform distribution of referents over documents (here: 25 clusters with $|C| = 23$)

- ⇒ $|TP|$ true 1-similarities per cluster (here: 130 @ threshold 0.725)
- ⇒ $\frac{|TP|}{|C|}$ degree of true positives per node (here: 11)
- ⇒ $|TP|(\frac{1}{precision} - 1)$ false 1-similarities per cluster (here: 760)

Density-based cluster analysis: effective false positives, FP^* , connect to same cluster

- ⇒ analyze $P(|FP^*| > k \mid D, R_{iid})$ (here: $E(|FP^*|) = 2.7$)
- ⇒ edge tie factor (*ETF*) specifies the excess of true positives until tie (here: 3...5)

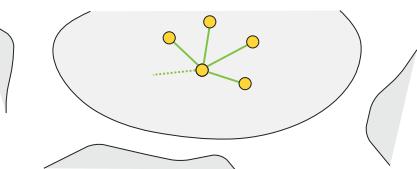
$$ETF = \frac{|TP|}{|C| \cdot E(|FP^*|)}, \quad \text{effective precision} = \text{precision} \cdot \frac{CIF}{ETF}$$

ML-IX-285 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Assumption: uniform distribution of referents over documents (here: 25 clusters with $|C| = 23$)

- ⇒ $|TP|$ true 1-similarities per cluster (here: 130 @ threshold 0.725)
- ⇒ $\frac{|TP|}{|C|}$ degree of true positives per node (here: 11)
- ⇒ $|TP|(\frac{1}{precision} - 1)$ false 1-similarities per cluster (here: 760)

Density-based cluster analysis: effective false positives, FP^* , connect to same cluster

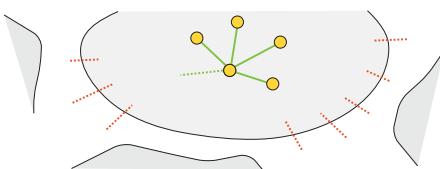
- ⇒ analyze $P(|FP^*| > k \mid D, R_{iid})$ (here: $E(|FP^*|) = 2.7$)
- ⇒ edge tie factor (*ETF*) specifies the excess of true positives until tie (here: 3...5)

$$ETF = \frac{|TP|}{|C| \cdot E(|FP^*|)}, \quad \text{effective precision} = \text{precision} \cdot \frac{CIF}{ETF}$$

©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Assumption: uniform distribution of referents over documents (here: 25 clusters with $|C| = 23$)

- ⇒ $|TP|$ true 1-similarities per cluster (here: 130 @ threshold 0.725)
- ⇒ $\frac{|TP|}{|C|}$ degree of true positives per node (here: 11)
- ⇒ $|TP|(\frac{1}{precision} - 1)$ false 1-similarities per cluster (here: 760)

Density-based cluster analysis: effective false positives, FP^* , connect to same cluster

- ⇒ analyze $P(|FP^*| > k \mid D, R_{iid})$ (here: $E(|FP^*|) = 2.7$)
- ⇒ edge tie factor (*ETF*) specifies the excess of true positives until tie (here: 3...5)

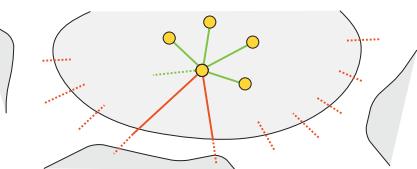
$$ETF = \frac{|TP|}{|C| \cdot E(|FP^*|)}, \quad \text{effective precision} = \text{precision} \cdot \frac{CIF}{ETF}$$

ML-IX-287 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Assumption: uniform distribution of referents over documents (here: 25 clusters with $|C| = 23$)

- ⇒ $|TP|$ true 1-similarities per cluster (here: 130 @ threshold 0.725)
- ⇒ $\frac{|TP|}{|C|}$ degree of true positives per node (here: 11)
- ⇒ $|TP|(\frac{1}{precision} - 1)$ false 1-similarities per cluster (here: 760)

Density-based cluster analysis: effective false positives, FP^* , connect to same cluster

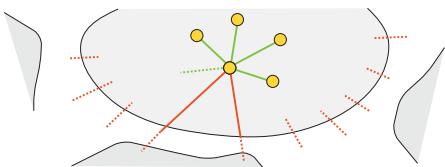
- ⇒ analyze $P(|FP^*| > k \mid D, R_{iid})$ (here: $E(|FP^*|) = 2.7$)
- ⇒ edge tie factor (*ETF*) specifies the excess of true positives until tie (here: 3...5)

$$ETF = \frac{|TP|}{|C| \cdot E(|FP^*|)}, \quad \text{effective precision} = \text{precision} \cdot \frac{CIF}{ETF}$$

©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Assumption: uniform distribution of referents over documents (here: 25 clusters with $|C| = 23$)

$$\Rightarrow |TP| \text{ true 1-similarities per cluster (here: 130 @ threshold 0.725)}$$

$$\Rightarrow \frac{|TP|}{|C|} \text{ degree of true positives per node (here: 11)}$$

$$\Rightarrow |TP| \left(\frac{1}{precision} - 1 \right) \text{ false 1-similarities per cluster (here: 760)}$$

Density-based cluster analysis: **effective false positives**, FP^* , connect to same cluster

$$\Rightarrow \text{analyze } P(|FP^*| > k \mid D, R_{iid}) \text{ (here: } E(|FP^*|) = 2.7)$$

\Rightarrow edge tie factor (ETF) specifies the excess of true positives until tie (here: 3...5)

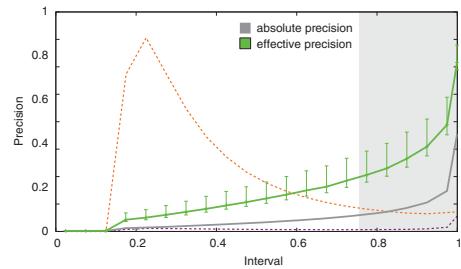
$$ETF = \frac{|TP|}{|C| \cdot E(|FP^*|)}, \quad \text{effective precision} = precision \cdot \frac{CIF}{ETF}$$

ML-IX-289 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Density-based cluster analysis: **effective false positives**, FP^* , connect to same cluster

$$\Rightarrow \text{analyze } P(|FP^*| > k \mid D, R_{iid}) \text{ (here: } E(|FP^*|) = 2.7)$$

\Rightarrow edge tie factor (ETF) specifies the excess of true positives until tie (here: 3...5)

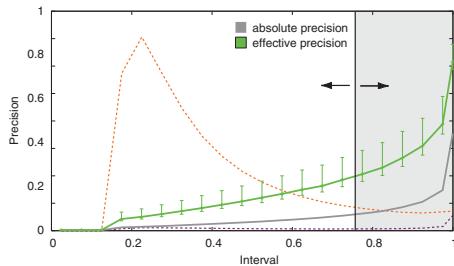
$$ETF = \frac{|TP|}{|C| \cdot E(|FP^*|)}, \quad \text{effective precision} = precision \cdot \frac{CIF}{ETF}$$

ML-IX-290 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

In-Depth: Analysis of Classifier Effectiveness



Determine optimum similarity threshold for class-membership function:

$$\theta^* = \underset{\theta \in [0;1]}{\operatorname{argmax}} \left\{ \frac{1 + \alpha}{\frac{ETF}{precision_\theta \cdot CIF} + \frac{\alpha}{recall_\theta}} \right\}$$

θ^* considers co-variate shift, introduces model formation bias and sample selection bias.

ML-IX-291 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

Model Selection: Our Risk Minimization Strategy

Retrieval Model	$F_{1/3}$ -Measure
<i>tfidf</i> vector space	0.39
context-based vector space	0.32
ESA Wikipedia persons	0.30
phrase structure grammar	0.17
ontology alignment	0.15
optimized combination	0.42
Ensemble cluster analysis	0.40

Ensemble cluster analysis: higher bias, better generalization.

(1) Do we speculate on a better fit?

(2) Do we expect a significant covariate shift, more noise, etc.?

Constrained Cluster Analysis

Recap

1. Multi-document resolution can be tackled with constrained cluster analysis.
2. Constraints are derived from labeled examples.
3. Class membership function ties constraints to multiple retrieval models.
4. Advanced density-based clustering technology is key.

ML-IX-293 Cluster Analysis

©STEIN 2002-2012

Constrained Cluster Analysis

References

- Disambiguating Web Appearances of People in a Social Network.
[R. Bekkerman, A. McCallum. WWW 2005]
- A Bayesian Model for Supervised Clustering with the Dirichlet Process Prior.
[H. Daumé III, D. Marcu. Journal MLR 2005]
- Computing Semantic Relatedness Using Wikipedia-based Explicit Semantic Analysis.
[E. Gabrilovich, S. Markovitch. IJCAI 2007]
- Unsupervised Discrimination of Person Names in Web Contexts.
[T. Pedersen, A. Kulkarni. CICLing 2007]
- On Information Need and Categorizing Search.
[S. Meyer zu Eissen. Dissertation, Paderborn University, 2007]
- Weighted Experts: A Solution for the Spock Data Mining Challenge.
[B. Stein, S. Meyer zu Eissen. I-KNOW 2008]
- GRAPE: A System for Disambiguating and Tagging People Names in Web Search.
[L. Jiang, W. Shen, J. Wang, N. An]

ML-IX-294 Cluster Analysis

©STEIN 2002-2012