

Grundlagen

Ziel

Konstruktion einer Hierarchie von Clustern (meist repräsentiert durch ein sog. *Dendrogramm*), so dass immer die Cluster mit minimaler Distanz verschmolzen werden

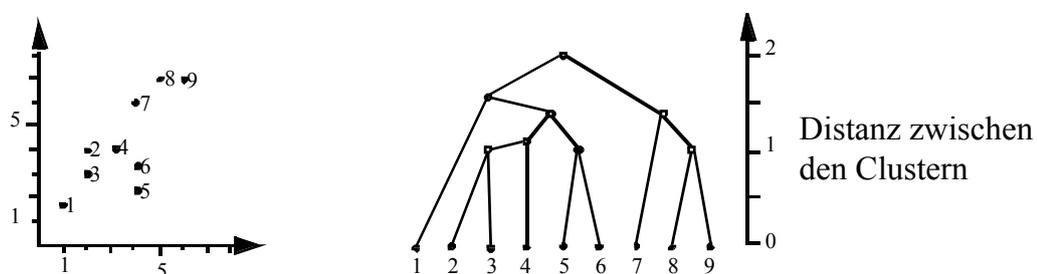
Dendrogramm

ein Baum, dessen Knoten jeweils einen Cluster repräsentieren, mit folgenden Eigenschaften:

- die Wurzel repräsentiert die ganze DB
- die Blätter repräsentieren einzelne Objekte
- ein innerer Knoten repräsentiert einen Cluster bestehend aus allen Objekten des darunter liegenden Teilbaums

Grundlagen

Beispiel eines Dendrogramms



Typen von hierarchischen Verfahren (meistens best-first Heuristiken)

Bottom-Up Konstruktion des Dendrogramms (*agglomerative*)

Top-Down Konstruktion des Dendrogramms (*divisive*)

Agglomeratives hierarchisches Clustering

1. Bilde initiale Cluster, die jeweils aus einem Objekt bestehen, und bestimme die Distanzen zwischen allen Paaren dieser Cluster.
2. Bilde einen neuen Cluster aus den zwei Clustern, welche die geringste Distanz zueinander haben.
3. Bestimme die Distanz zwischen dem neuen Cluster und allen anderen Clustern.
4. Wenn sich alle Objekte in einem einzigen Cluster befinden: Fertig, andernfalls wiederhole ab Schritt 2.

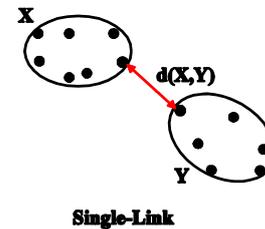
Distanzfunktionen für Cluster

- Die Verfahren unterscheiden sich anhand ihrer Distanzfunktionen für Cluster
- Beliebte Distanzfunktionen (Clustermodell: Kompaktheit):
 - Single Link
 - Complete Link
 - Average Link
- Sei eine Distanzfunktion $dist(x,y)$ für Paare von Objekten gegeben
- Seien X, Y Cluster, d.h. Mengen von Objekten.

Single-Link Distanz

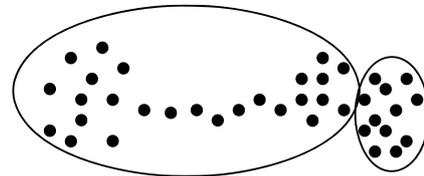
Definition:

$$distSL(X, Y) = \min_{x \in X, y \in Y} dist(x, y)$$



Eigenschaften:

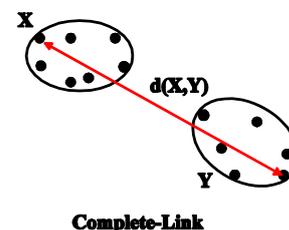
- Effiziente Implementierung (z.B. SLINK): $O(n^2)$
- Single-Link Effekt: „kettenförmige“ Cluster, Cluster werden durch wenige, kettenförmig verteilte Objekte vereinigt
 - Cluster mit starker Streuung
 - Cluster mit langgezogener Struktur



Complete-Link Distanz

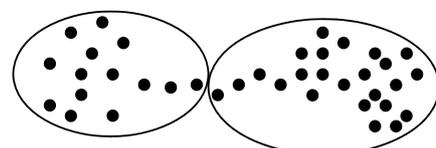
Definition:

$$distCL(X, Y) = \max_{x \in X, y \in Y} dist(x, y)$$



Eigenschaften:

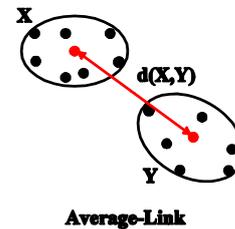
- Effiziente Implementierung (z.B. CLINK): $O(n^2)$
- Complete-Link Effekt
 - Kleine, stark abgegrenzte Cluster
 - Gleichgroße, konvexe Cluster



Average-Link Distanz

Definition:

$$distAL(X, Y) = \frac{1}{|X| \cdot |Y|} \cdot \sum_{x \in X, y \in Y} dist(x, y)$$



Eigenschaften:

- Keine effiziente Implementierung
- Kompromiss zwischen Single- und Complete-Link Ansatz

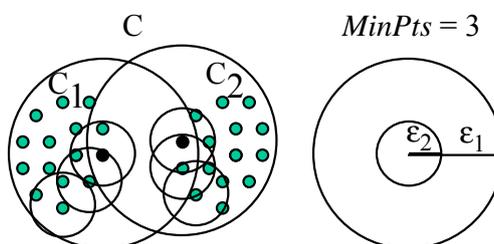
Diskussion

- + erfordert keine Kenntnis der Anzahl k der Cluster
- + findet nicht nur ein flaches Clustering, sondern eine ganze Hierarchie
- + ein einzelnes Clustering kann aus dem Dendrogramm gewonnen werden, z.B. mit Hilfe eines horizontalen Schnitts durch das Dendrogramm (erfordert aber wieder Anwendungswissen: wo ist ein sinnvoller Schnitt?)
- Entscheidungen können nicht zurückgenommen werden (best first search)
- Single-Link-Effekte, Complete-Link-Effekte
- Ineffizienz: Laufzeitkomplexität von mindestens $O(n^2)$ für n Objekte

Dichtebasiertes hierarchisches Clustering

[Ankerst, Breunig, Kriegel & Sander 1999]

- für einen konstanten *MinPts*-Wert sind dichte-basierte Cluster bzgl. eines kleineren ϵ vollständig in Clustern bzgl. eines größeren ϵ enthalten



- in einem DBSCAN-ähnlichen Durchlauf gleichzeitig das Clustering für verschiedene Dichte-Parameter bestimmen
- zuerst die dichteren Teil-Cluster, dann den dünneren Rest-Cluster
- kein Dendrogramm, sondern eine auch noch bei sehr großen Datenmengen übersichtliche Darstellung der Cluster-Hierarchie
- Clustermodell: Dichte

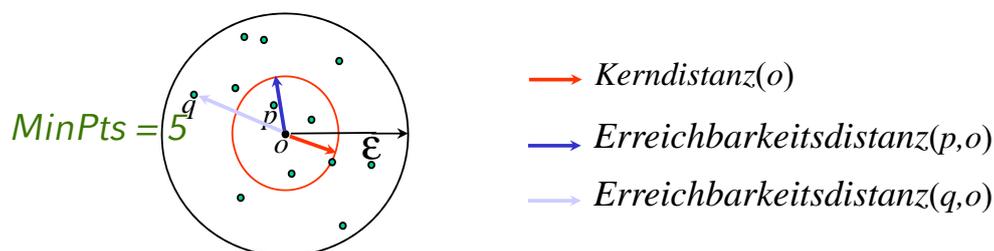
Grundbegriffe

Kerndistanz eines Objekts o bzgl. ϵ und *MinPts*

$$\text{Kerndistanz}_{\epsilon, \text{MinPts}}(o) = \begin{cases} \text{UNDEFINIERT}, & \text{wenn } |RQ(o, \epsilon)| < \text{MinPts} \\ \text{MinPtsDistanz}(o), & \text{sonst} \end{cases}$$

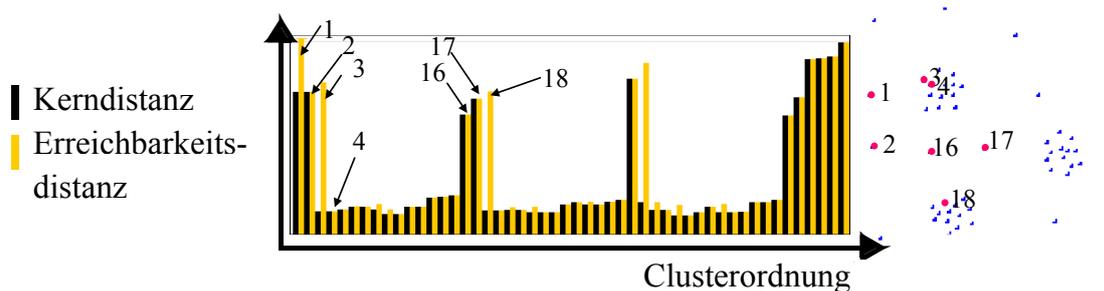
Erreichbarkeitsdistanz eines Objekts p relativ zu einem Objekt o

$$\text{Erreichbarkeitsdistanz}_{\epsilon, \text{MinPts}}(p, o) = \begin{cases} \text{UNDEFINIERT} & \text{wenn } |RQ(o, \epsilon)| < \text{MinPts} \\ \max\{\text{Kerndistanz}(o), \text{dist}(o, p)\}, & \text{sonst} \end{cases}$$



Clusterordnung

- OPTICS liefert nicht direkt ein (hierarchisches) Clustering, sondern eine „Clusterordnung“ bzgl. ϵ und $MinPts$
- Clusterordnung bzgl. ϵ und $MinPts$
 - beginnt mit einem beliebigen Objekt
 - als nächstes wird das Objekt besucht, das zur Menge der bisher besuchten Objekte die minimale Erreichbarkeitsdistanz besitzt



Algorithmus OPTICS

Datenstrukturen

- SeedList
 - speichert Punkte mit „aktueller“ Erreichbarkeitsdistanz aufsteigend sortiert
- ClusterOrder
 - resultierende Clusterordnung wird schrittweise aufgebaut

Hauptschleife:

```
SeedList = ∅;
WHILE es gibt noch unmarkierte Objekte in DB DO
  IF SeedList = ∅
    THEN
      füge beliebiges noch unmarkiertes Objekt in ClusterOrder ein mit
      Erreichbarkeitsdistanz ∞;
    ELSE
      füge erstes Objekt aus der SeedList mit aktueller Erreichbarkeitsdistanz in
      ClusterOrder ein;
  // sei obj das zuletzt in ClusterOrder eingefügte Objekt
  markiere obj als bearbeitet;
  FOR ALL neighbor ∈ RQ(obj, ε) DO
    // insert/update neighbor in SeedList mit referenzobjekt obj
    SeedList.update(neighbor, obj);
```

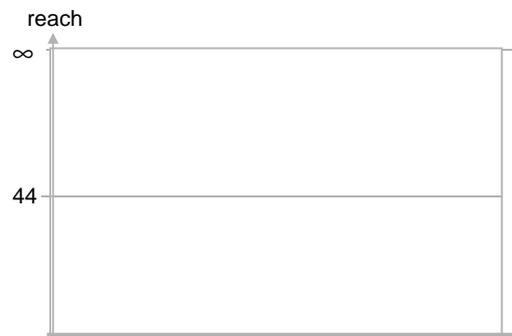
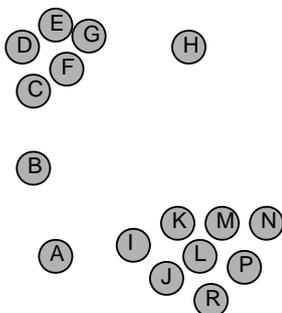
Algorithmus OPTICS

Einfügen/Updaten eines Objekts o in SeedList

- Beachte: Für alle Objekte p in SeedList ist die „aktuelle“ Erreichbarkeitsdistanz $p.rdist$ gespeichert.
- SeedList ist nach $p.rdist$ aufsteigend sortiert (als Heap organisiert)
- Referenzobjekt: obj

```
SeedList :: update( $o$ ,  $obj$ )
Berechne Erreichbarkeitsdistanz $_{\epsilon, MinPts}(o, obj) := rdistneu_o$ ;
IF  $o$  ist bereits in SeedList THEN
    IF  $rdistneu_o \leq o.rdist$  THEN
         $o.rdist := rdistneu_o$ ;
        verschiebe  $o$  in SeedList (nach vorne); // aufsteigen im Heap
    ELSE
        //  $o$  ist noch nicht in SeedList => normales Einfügen in Heap
        füge  $o$  mit  $o.rdist := rdistneu_o$  in SeedList ein;
```

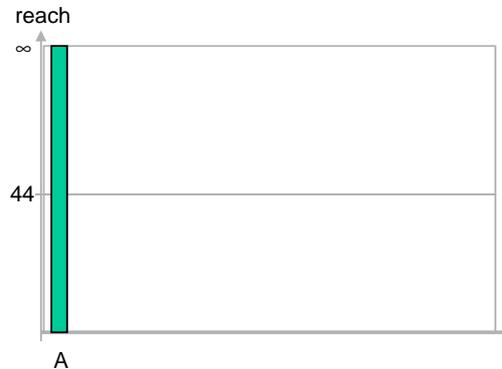
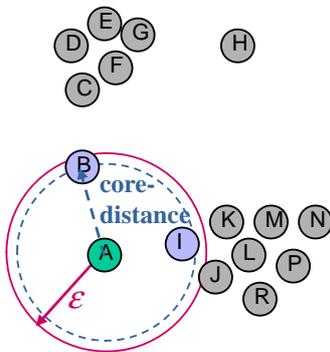
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list:

5.4 Hierarchische Verfahren

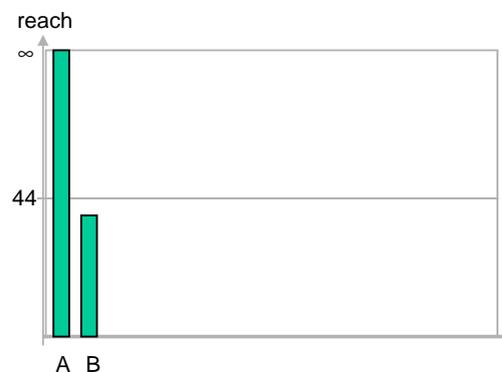
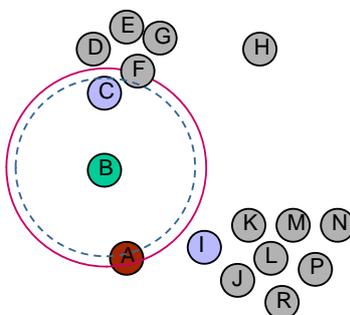
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (B, 40) (I, 40)

5.4 Hierarchische Verfahren

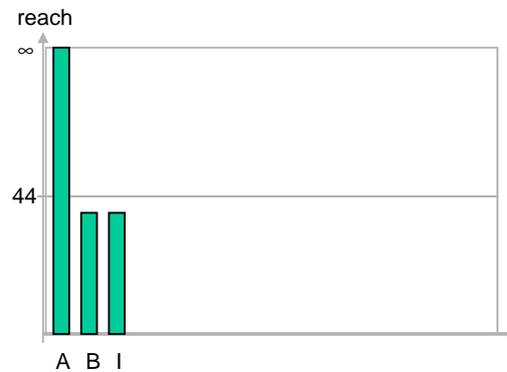
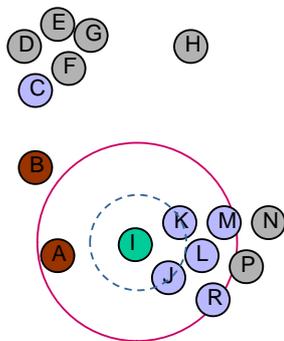
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (I, 40) (C, 40)

5.4 Hierarchische Verfahren

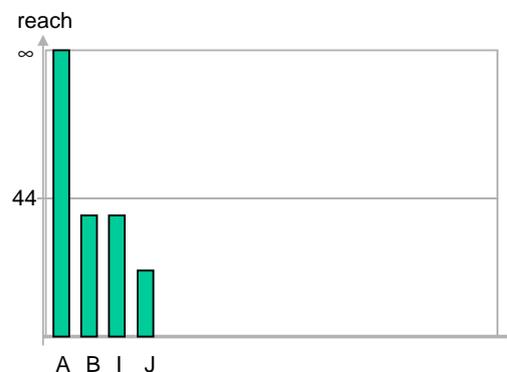
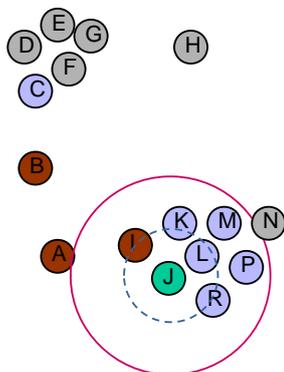
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (J, 20) (K, 20) (L, 31) (C, 40) (M, 40) (R, 43)

5.4 Hierarchische Verfahren

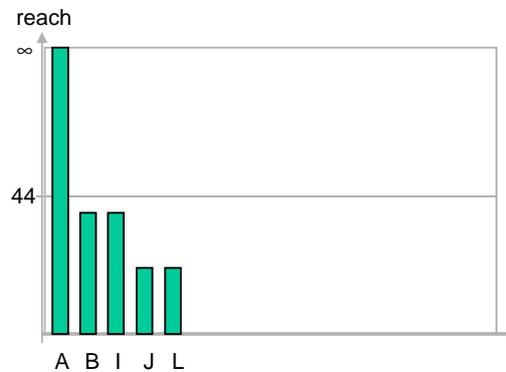
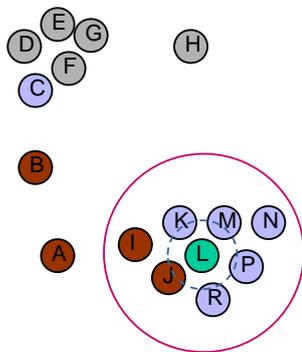
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (L, 19) (K, 20) (R, 21) (M, 30) (P, 31) (C, 40)

5.4 Hierarchische Verfahren

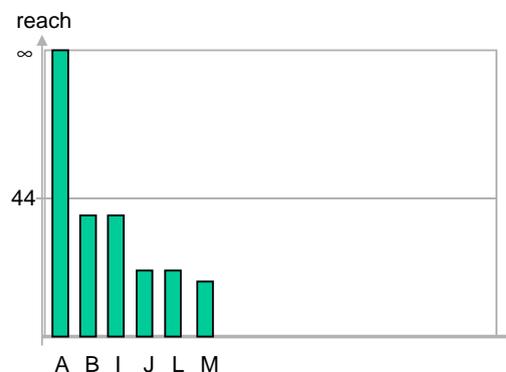
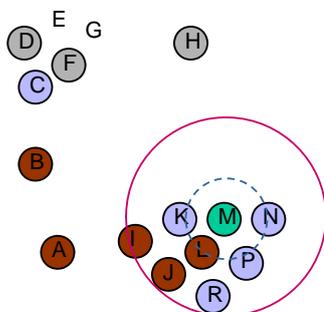
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (M, 18) (K, 18) (R, 20) (P, 21) (N, 35) (C, 40)

5.4 Hierarchische Verfahren

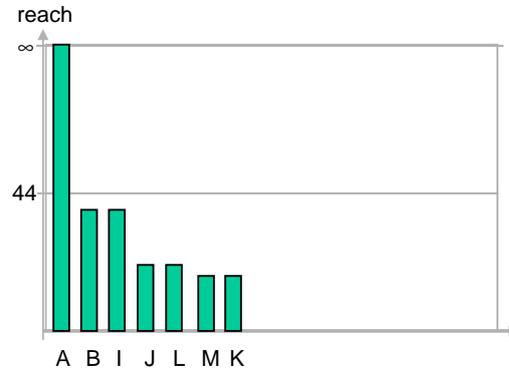
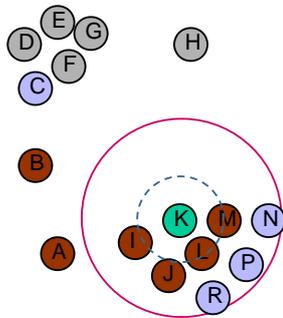
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (K, 18) (N, 19) (R, 20) (P, 21) (C, 40)

5.4 Hierarchische Verfahren

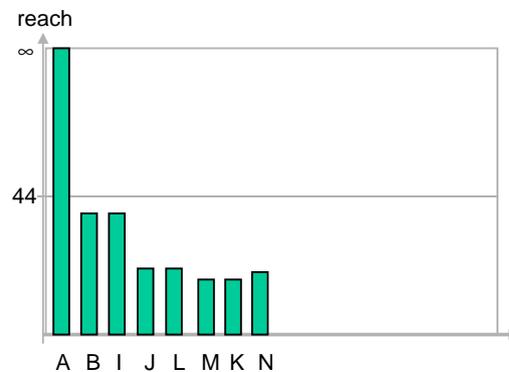
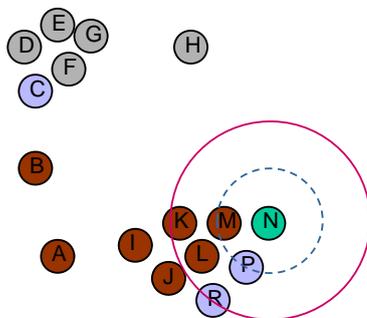
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (N, 19) (R, 20) (P, 21) (C, 40)

5.4 Hierarchische Verfahren

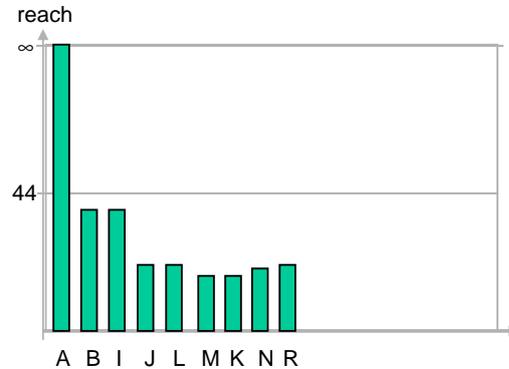
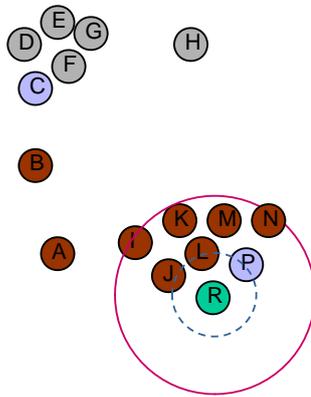
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (R, 20) (P, 21) (C, 40)

5.4 Hierarchische Verfahren

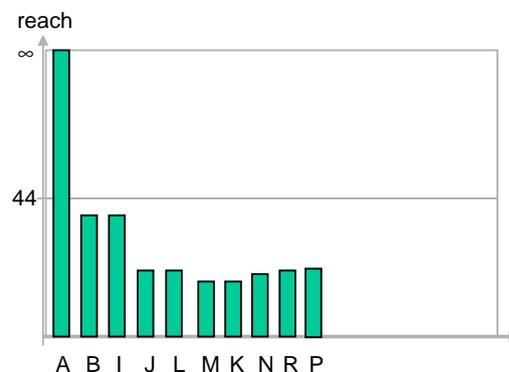
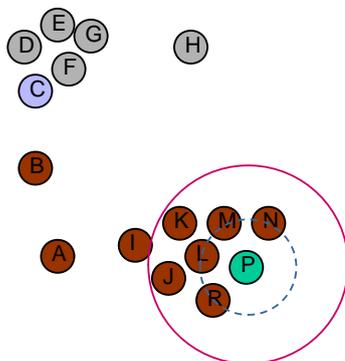
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (P, 21) (C, 40)

5.4 Hierarchische Verfahren

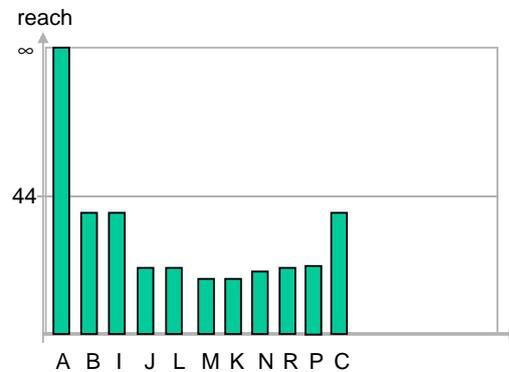
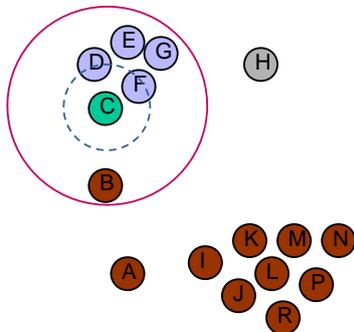
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (C, 40)

5.4 Hierarchische Verfahren

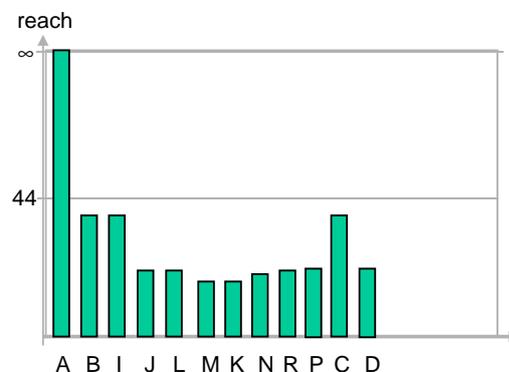
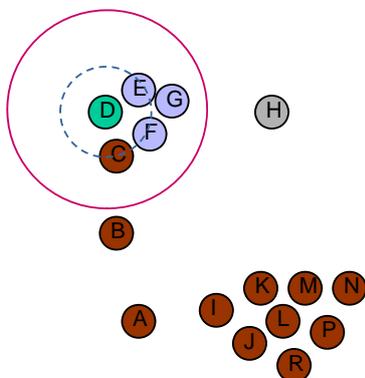
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (D, 22) (F, 22) (E, 30) (G, 35)

5.4 Hierarchische Verfahren

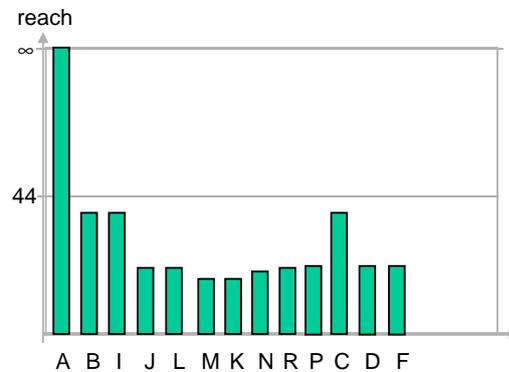
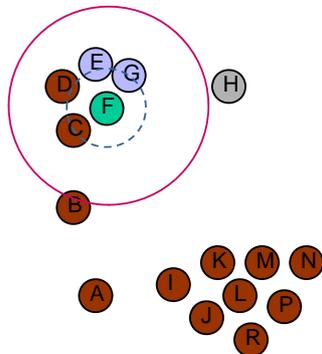
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (F, 22) (E, 22) (G, 32)

5.4 Hierarchische Verfahren

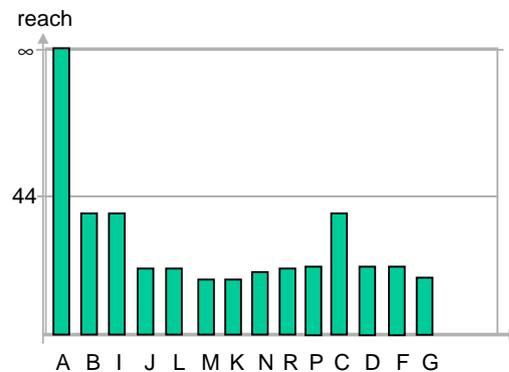
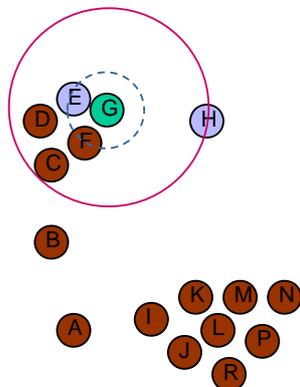
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (G, 17) (E, 22)

5.4 Hierarchische Verfahren

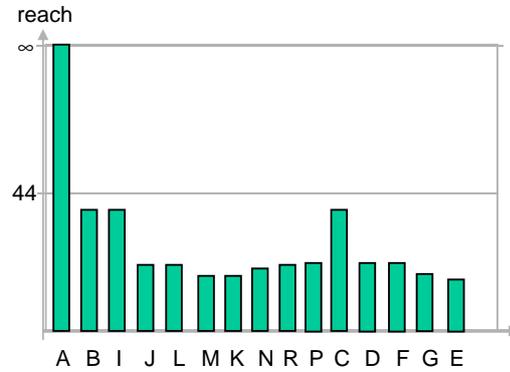
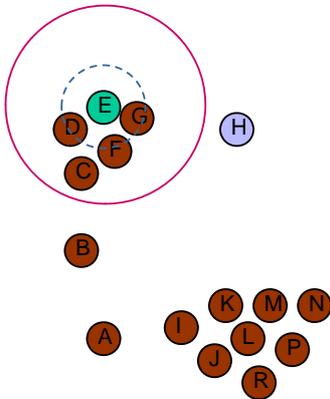
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (E, 15) (H, 43)

5.4 Hierarchische Verfahren

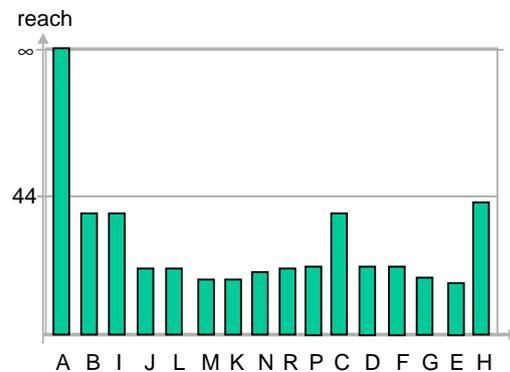
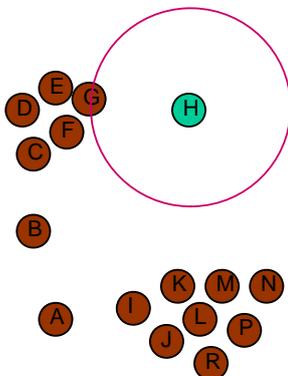
- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$



seed list: (H, 43)

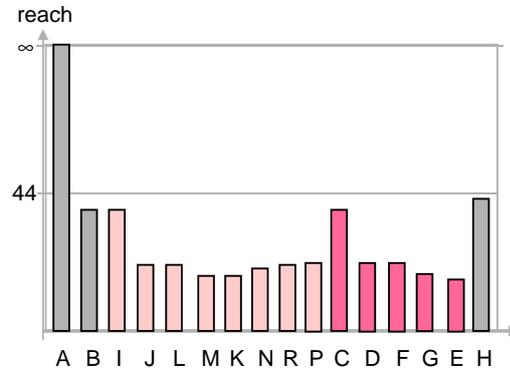
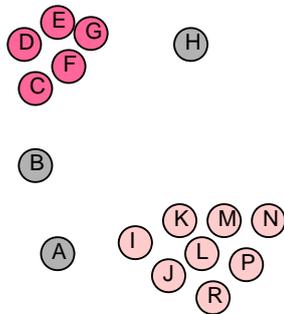
5.4 Hierarchische Verfahren

- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$

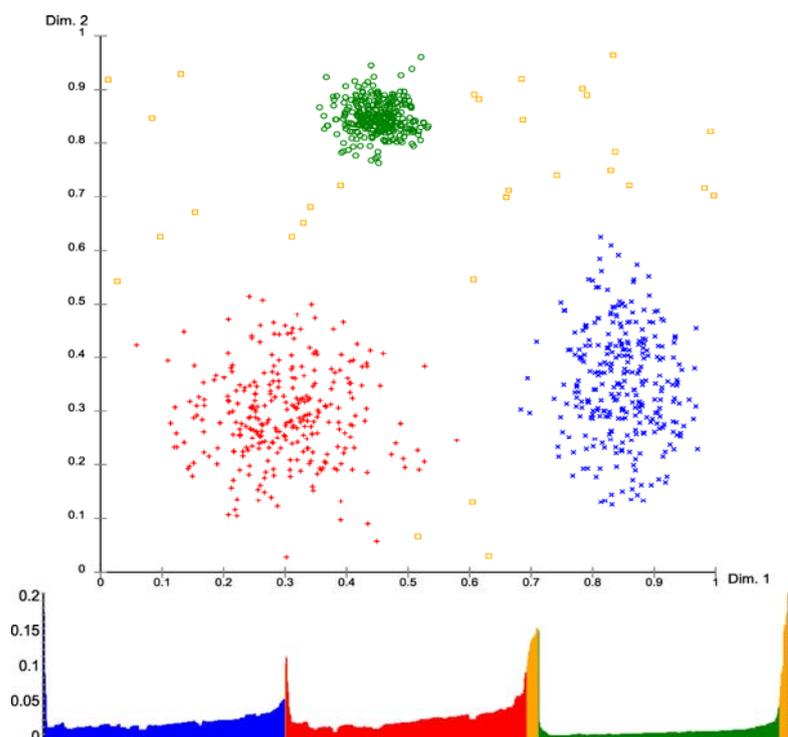


seed list: -

- Example Database (2-dimensional, 16 points)
- $\epsilon = 44$, $MinPts = 3$

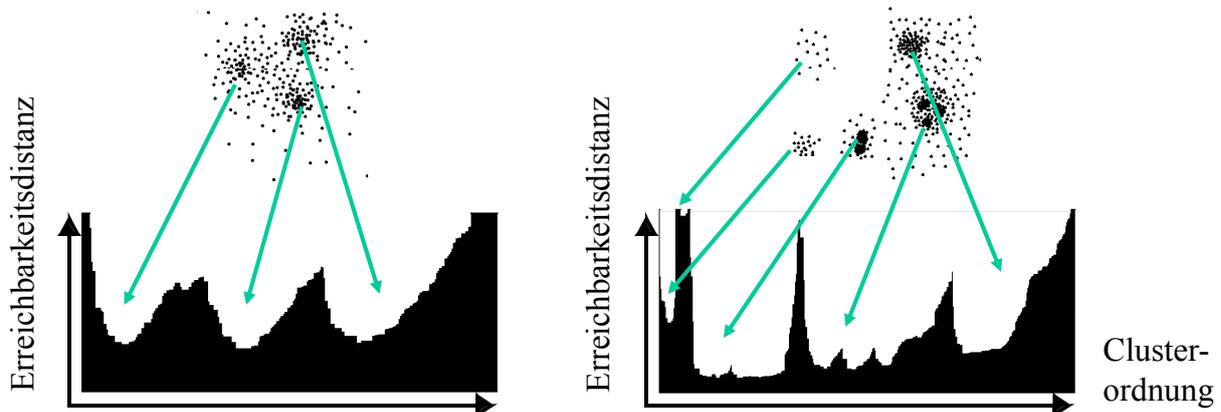


Erreichbarkeits-Diagramm:



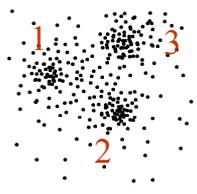
Erreichbarkeits-Diagramm

Zeigt die Erreichbarkeitsdistanzen (bzgl. ϵ und $MinPts$) der Objekte als senkrechte, nebeneinanderliegende Balken in der durch die Clusterordnung der Objekte gegebenen Reihenfolge

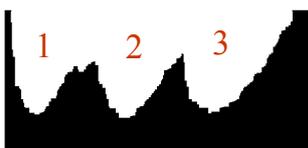


Schnitt: ϵ für DBSCAN: Dichte-Niveau

Parameter-Sensitivität



$MinPts = 10, \epsilon = 10$



optimale Parameter

$MinPts = 10, \epsilon = 5$



kleineres ϵ

$MinPts = 2, \epsilon = 10$



kleineres $MinPts$



Clusterordnung ist robust gegenüber den Parameterwerten
gute Resultate wenn Parameterwerte „groß genug“

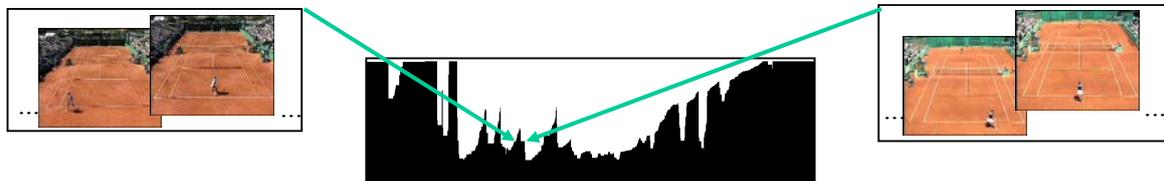
Heuristische Parameter-Bestimmung

ϵ

- wähle größte *MinPts*-Distanz aus einem Sample oder
- berechne durchschnittliche *MinPts*-Distanz für gleichverteilte Daten

MinPts

- glätte Erreichbarkeits-Diagramm
- vermeide "single-" bzw. "*MinPts*-link" Effekt

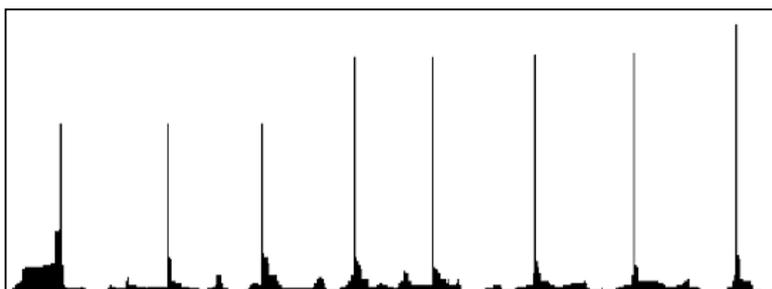


105

Manuelle Analyse der Cluster

Mit Erreichbarkeits-Diagramm

- gibt es Cluster?
- wieviele Cluster?
- sind die Cluster hierarchisch geschachtelt?
- wie groß sind die Cluster?



Erreichbarkeits-Diagramm

106

Automatisches Entdecken von Clustern

ξ -Cluster [Ankerst, Breunig, Kriegel, Sander 99]

Cluster = Teilsequenz der Clusterordnung mit mindestens *MinPts* Punkte

1. beginnt in einem Gebiet *x*-steil *abfallender* Erreichbarkeitsdistanzen
2. endet in einem Gebiet *x*-steil *steigender* Erreichbarkeitsdistanzen bei etwa demselben absoluten Wert

ClusterTree [Sander, Qin, Lu, Niu, Kovarsky 02]

Cluster sind geteilt durch "signifikante" lokale Maxima

GradientClustering [Brecheisen, Kriegel, Kröger, Pfeifle 04]

Kombination aus ξ -Cluster und ClusterTree