

# INTELLIGENTE DATENANALYSE IN MATLAB

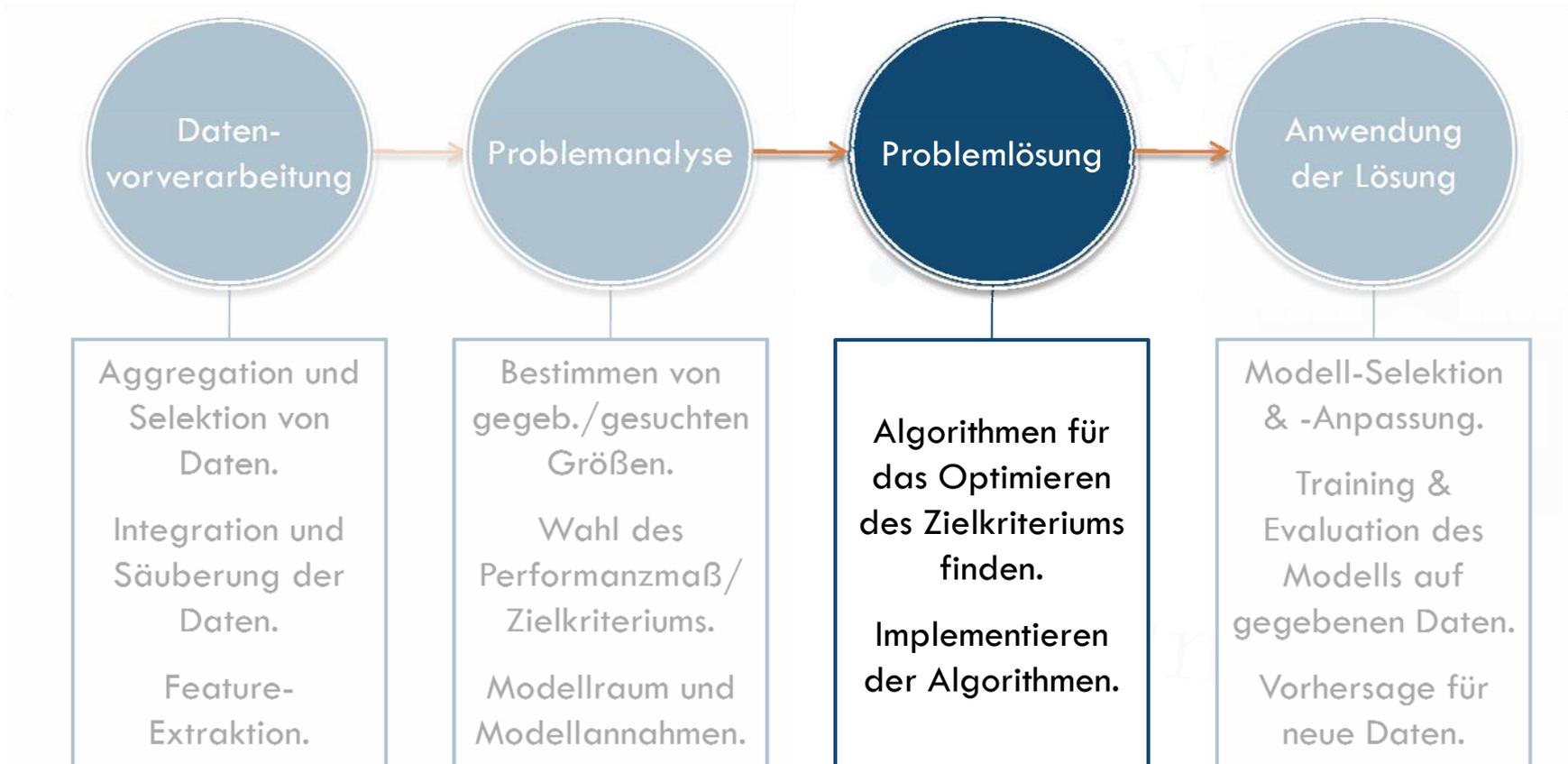
Unüberwachtes Lernen

# Literatur

- Chris Bishop: Pattern Recognition and Machine Learning.
- Jiawei Han und Micheline Kamber: Data Mining – Concepts and Techniques.
- Ulrike von Luxburg: A Tutorial on Spectral Clustering.  
[http://www.kyb.mpg.de/publications/attachments/Luxburg06\\_TR\\_%5B0%5D.pdf](http://www.kyb.mpg.de/publications/attachments/Luxburg06_TR_%5B0%5D.pdf)
- Matteo Matteucci: A Tutorial on Clustering Algorithms.  
[http://home.dei.polimi.it/matteucc/Clustering/tutorial\\_html](http://home.dei.polimi.it/matteucc/Clustering/tutorial_html)

# Überblick

## □ Schritte der Datenanalyse:



# Unüberwachtes Lernen:

## Problemstellung

- Gegeben: Trainingsdaten mit unbekannten Zielattributen (ungelabelte Daten).
- Eingabe: Instanzen.
- Ausgabe: Belegung der Zielattribute.
  - Clustern: Nominaler Wertebereich des Zielattributs (z.B. Jahreszeiten, {Cluster1, Cluster2, ...}), d.h. jeder Datenpunkt wird zu genau einem Cluster zugeordnet.
  - Dichteschätzung: Numerischer Wertebereich des Zielattributs, d.h. jedem Datenpunkt wird Wahrscheinlichkeitsverteilung über alle Cluster zugeordnet.

# Unüberwachtes Lernen:

## Beispiel

### □ Beispiel *Clustern* (Tabellendarstellung):

Monat	Bewölkung	Temperatur	Luftfeuchtigkeit	Wind	Jahreszeit
Juli	sonnig	warm	hoch	wenig	?
September	sonnig	warm	hoch	stark	?
August	bedeckt	warm	hoch	wenig	?
April	Regen	mild	hoch	wenig	?
Oktober	Regen	kühl	normal	wenig	?
Dezember	Regen	kühl	normal	stark	?
Januar	bedeckt	kühl	normal	stark	?
Juli	sonnig	mild	hoch	wenig	?
Februar	sonnig	kühl	normal	wenig	?
März	Regen	mild	normal	wenig	?
November	sonnig	mild	normal	stark	?
August	bedeckt	mild	hoch	stark	?
Juni	bedeckt	warm	normal	wenig	?
April	Regen	mild	hoch	stark	?

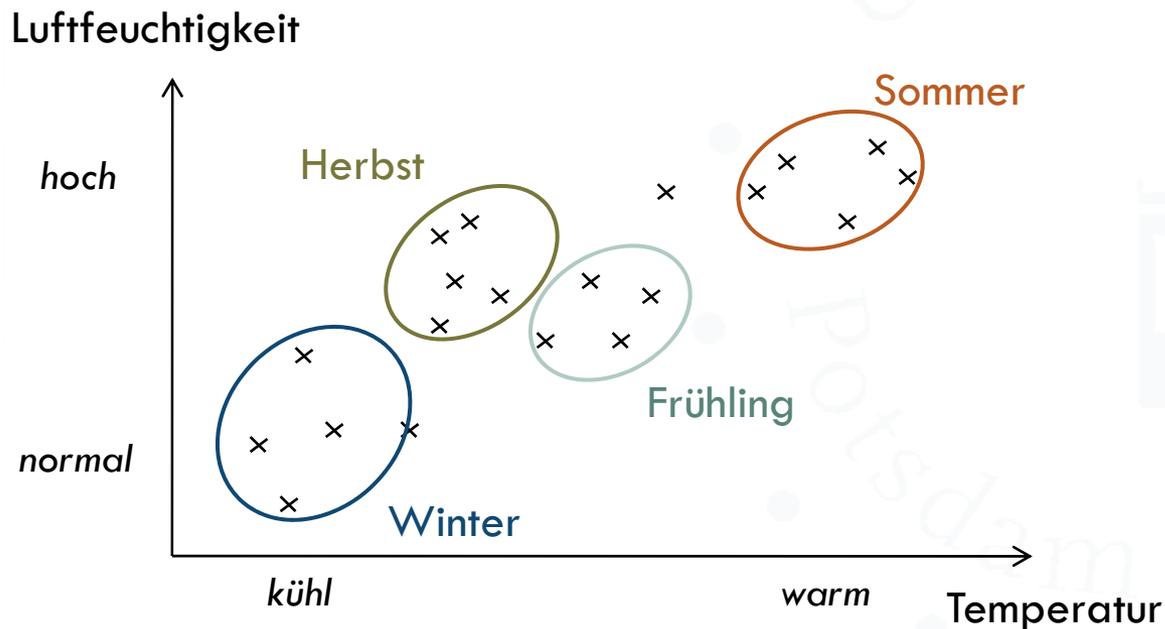
Trainingsdaten

Zielgröße

# Unüberwachtes Lernen:

## Beispiel

- Beispiel *Clustern* (Diagramm bzgl. der Attribute Luftfeuchtigkeit und Temperatur):



# Motivation

- Besseres Verständnis/Beschreibung der Daten:
  - Kundensegmentierung: Zielgerichtetes Marketing, Produktentwicklung usw. für einzelne Zielgruppen.
  - Landnutzung/Stadtplanung: Identifizieren von Regionen mit ähnlichen geographischen, klimatischen, städtebaulichen Eigenschaften.
  - Risikoanalyse: Klassen von Kunden mit unterschiedlichen Versicherungsrisiken erkennen.
- Teil der Datenvorverarbeitung für weitere Analyse:
  - Diskretisierung von numerischen Attributen.
  - Outlier-Detection.

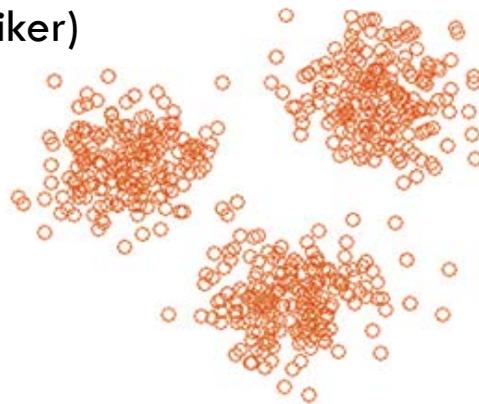
# Motivation:

## Anwendung

- Überblick über eine Dokumentenkollektion.
  - Suche nach Stichwort „Kohl“ liefert viele Dokumente.
  - Idee: Zeige dem Nutzer Cluster um genauere Auswahl des Themas zu ermöglichen.

Helmut Kohl (Politiker)

Kohl's (US Kaufhaus)



Kohl (Gemüse)

# Motivation:

## Anwendung

- Spam-Kampagnen identifizieren.
  - Spam-Kampagne ist große Menge ähnlicher (aber nicht gleicher) Emails.
  - Idee: Email-Attribute (z.B. enthaltene Wörter) clustern und Spam-Cluster durch nachgeschalteten Klassifikator erkennen.

Hello. This is Terry Hagan. We are accepting your mortgage application.

Our company confirms you are eligible for a \$250,000

loan for a \$380.00/month. Approval minute, so please fill out the form of

Best Regards, Terry Hagan; Senior Trades/Finance Department North

Dear Mr/Mrs, This is Brenda Dunn. We are accepting your mortgage application.

Our office confirms you can get a \$228,000 loan for a \$371.00 per month payment. Follow the link to our website and submit your contact information.

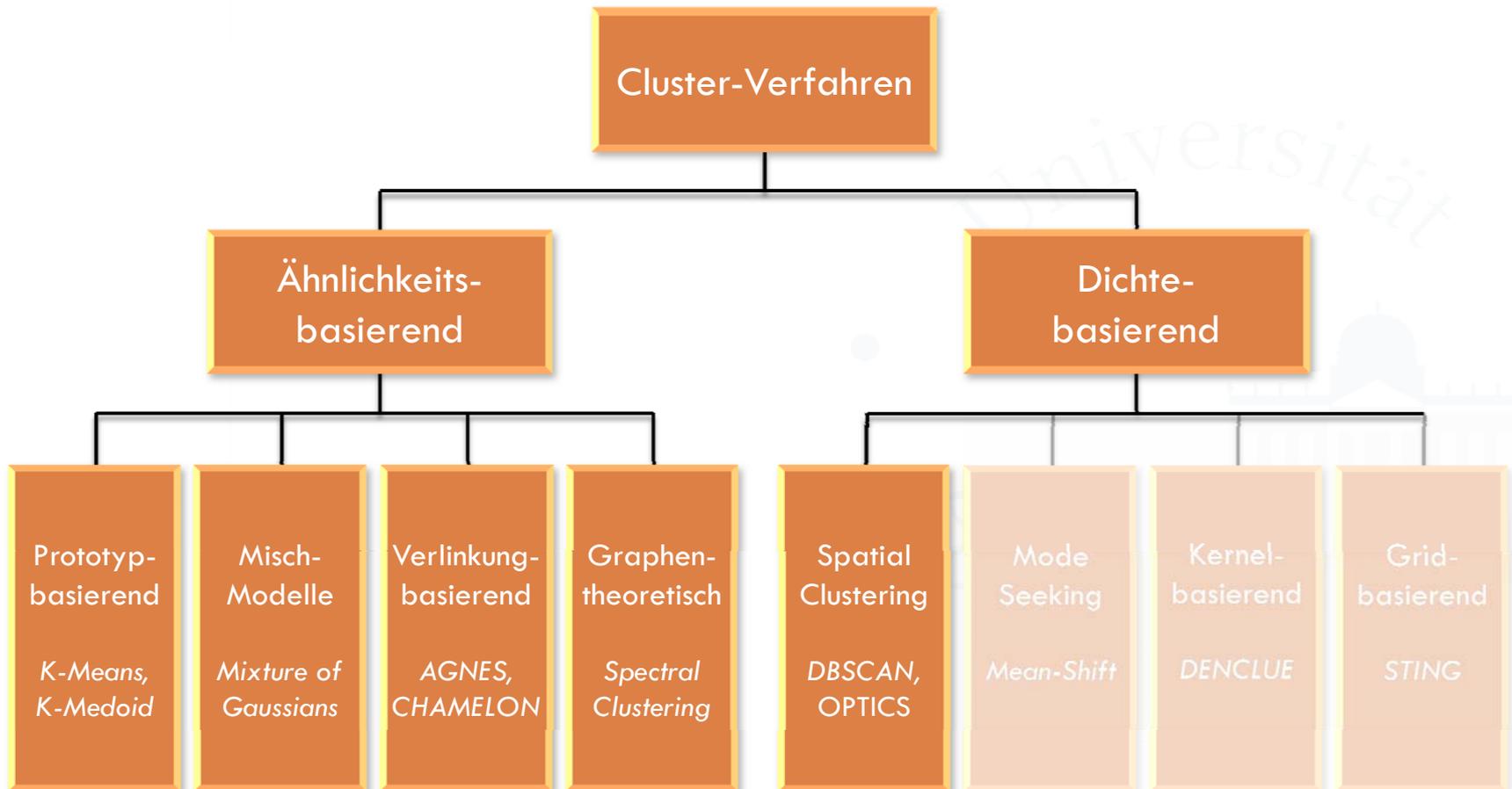
Best Regards, Brenda Dunn; Accounts Manager Trades/Finance Department East Office

# Evaluierung

- Qualitätsmerkmale eines Clusterings:
  - Hohe Ähnlichkeit zweier Datenpunkte eines Clusters (*intra-cluster similarity*).
  - Geringe Ähnlichkeit zwischen Datenpunkten verschiedener Cluster (*inter-cluster similarity*).
  - Anzahl, Form und Größenvarianz der Cluster.
  - Interpretierbarkeit, d.h. gefundene Cluster entsprechen echten (versteckten) Clustern.

Stichprobenartig durch Experten prüfen.

# Cluster-Verfahren



# Ähnlichkeitsbasierendes Clustern:

## Prototyp-Verfahren



- Gegeben:
  - ▣ Ungelabelte Daten  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ .
  - ▣ Anzahl vermuteter Cluster  $k$  mit  $1 < k < n$ .
  - ▣ Ähnlichkeitsmaß zwischen Datenpunkten.
- Gesucht: Partitionierung der Daten in  $k$  Cluster.
- Ziel: Kleiner Abstand zw. Punkten im selben Cluster und großer Abstand zw. Punkten verschiedener Cluster.
  - ▣ Exponentiell viele Partitionierungen  $\Rightarrow$  Suche NP-hart!
  - ▣ Heuristische Suche (lokal optimal): *K-Means* und *K-Medoids*.

# Ähnlichkeitsbasierendes Clustern:

## Prototyp-Verfahren: K-Means

- Gesucht ist eine Zuweisung  $\mathbf{R} = \{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_n\}$  der Daten  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  zu den Clustern mit

$$\mathbf{r}_i \in \{0, 1\}^k \quad r_{ij} = \begin{cases} 1 & \mathbf{x}_i \text{ in Cluster } j \\ 0 & \text{sonst} \end{cases}$$

z.B.  $\mathbf{r}_5 = [0 \ 0 \ 1 \ 0]^T$  falls das 5. Beispiel in Cluster 3 liegt.

- Seien  $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k\}$  Cluster-Mittelpunkte (Prototypen).
- Ziel: Minimaler quadratischer Abstand zu den Cluster-Mittelpunkten:

$$\min \sum_{i=1}^n \sum_{j=1}^k r_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$$

# Ähnlichkeitsbasierendes Clustern:

## Prototyp-Verfahren: K-Means

- Gleichzeitiges Minimieren über  $\{\mu_1, \mu_2, \dots, \mu_k\}$  und  $\{r_1, r_2, \dots, r_n\}$  sehr schwierig.
- Idee: Abwechselnde Minimierung.
- Algorithmus:

K-Means(*Instanzen*  $\mathbf{x}_i$ , *Clusteranzahl*  $k$ )

Setze  $l=0$  und wähle zufällig  $\forall i \mu_i^0 = \mathbf{x}_i$

DO

$$\{\mathbf{r}_1^{l+1}, \dots, \mathbf{r}_n^{l+1}\} = \arg \min_{\mathbf{r}_1, \dots, \mathbf{r}_n} \sum_{i=1}^n \sum_{j=1}^k r_{ij} \|\mathbf{x}_i - \mu_j^l\|^2$$

$$\{\mu_1^{l+1}, \dots, \mu_n^{l+1}\} = \arg \min_{\mu_1, \dots, \mu_n} \sum_{i=1}^n \sum_{j=1}^k r_{ij}^{l+1} \|\mathbf{x}_i - \mu_j\|^2$$

$l = l + 1$

WHILE  $\{\mathbf{r}_1^l, \dots, \mathbf{r}_n^l\} \neq \{\mathbf{r}_1^{l-1}, \dots, \mathbf{r}_n^{l-1}\}$

RETURN  $\{\mathbf{r}_1^l, \dots, \mathbf{r}_n^l\}$

Expectation

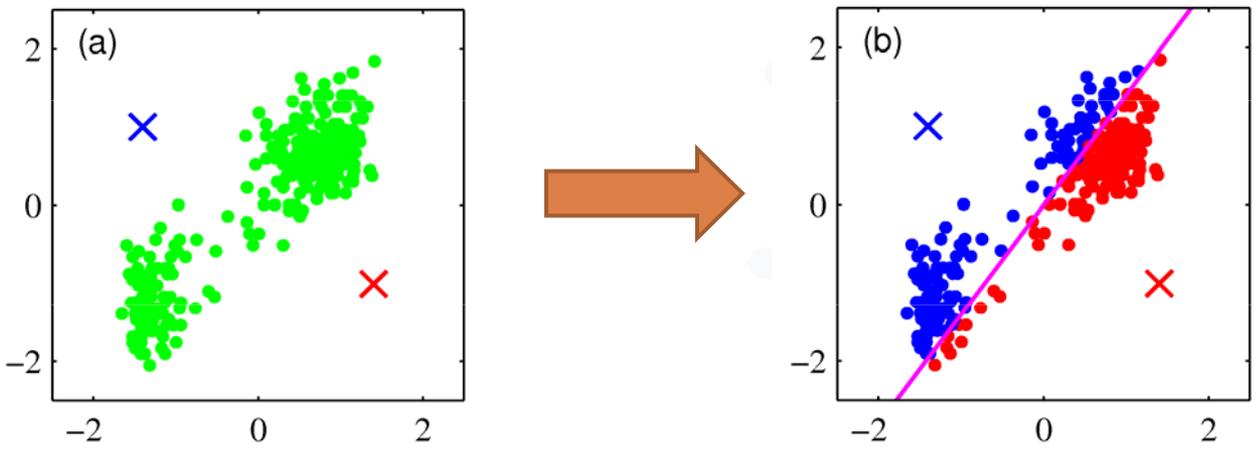
Maximization

# Ähnlichkeitsbasierendes Clustern:

## Prototyp-Verfahren: K-Means

□ Expectation-Schritt:  $\{\mathbf{r}_1^{l+1}, \dots, \mathbf{r}_n^{l+1}\} = \arg \min_{\mathbf{r}_1, \dots, \mathbf{r}_n} \sum_{i=1}^n \sum_{j=1}^k r_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j^l\|^2$

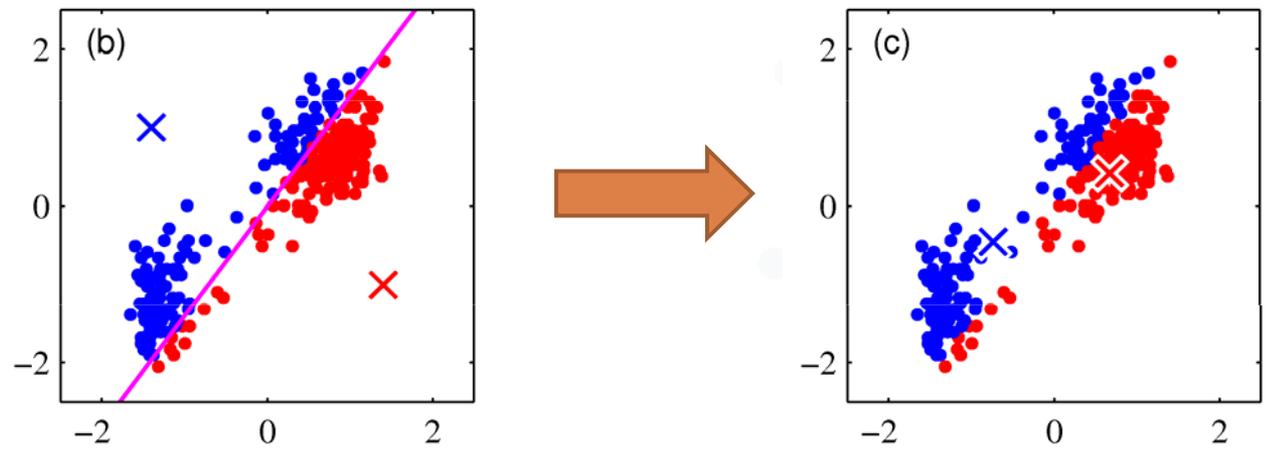
▣ Ordne jeden Punkt dem ihm nächsten Cluster-Mittelpunkt zu.



# Ähnlichkeitsbasierendes Clustern:

## Prototyp-Verfahren: K-Means

- Maximization-Schritt:  $\{\mu_1^{l+1}, \dots, \mu_n^{l+1}\} = \arg \min_{\mu_1, \dots, \mu_n} \sum_{i=1}^n \sum_{j=1}^k r_{ij}^{l+1} \|\mathbf{x}_i - \mu_j\|^2$
- Bestimme neue Cluster-Mittelpunkte  $\mu_j^{l+1} = \frac{\sum_{i=1}^n r_{ij}^{l+1} \mathbf{x}_i}{\sum_{i=1}^n r_{ij}^{l+1}}$ .



# Ähnlichkeitsbasierendes Clustern:

## Prototyp-Verfahren: K-Means

### □ Vorteile:

- Einfach zu implementieren
- Relativ schnell.
  - In  $O(n \cdot k)$  pro Iteration.
  - Effizienten Berechnung neuer Cluster-Mittelpunkte möglich.

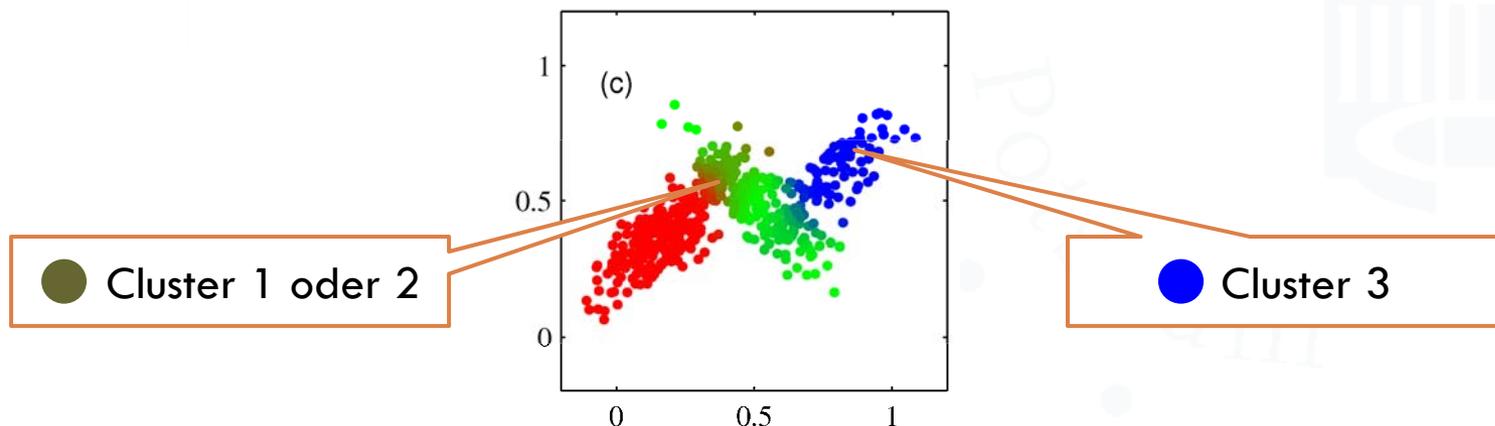
### □ Nachteile:

- Nur lokales Optimum garantiert (unterschiedliche Startwerte = unterschiedliche Lösungen).
- Harte Zuweisung zu Clustern (nicht-probabilistisch).
- Anzahl Cluster muss vorgeben werden.
- Nur für numerische Attribute geeignet.

# Ähnlichkeitsbasierendes Clustern:

## Mischmodelle

- Gegeben:
  - Ungelabelte Daten  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ .
  - Anzahl vermuteter Cluster  $k$  muss nicht bekannt sein.
  - Verteilungsannahme der Datenpunkte.
- Gesucht: Probabilistische Partitionierung der Daten.



# Ähnlichkeitsbasierendes Clustern:

## Mischmodelle

### □ Idee:

- Generatives (Misch-)Modell welches Daten  $\mathbf{X}$  erzeugt hat mit Modell-Parameter  $\Theta$ .
- Cluster-Zuordnungen  $\mathbf{R} = \{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_n\}$  sind versteckte Variablen des Modells.

### □ Ziel:

- Parameter  $\Theta$  mit maximaler A-Posteriori-Wahrscheinlichkeit (MAP):

$$\Theta^* = \arg \max_{\Theta} p(\Theta | \mathbf{X}) = \arg \max_{\Theta} p(\mathbf{X} | \Theta) p(\Theta)$$

# Ähnlichkeitsbasierendes Clustern:

## Mischmodelle: Mixture of Gaussians

- Gesucht ist eine Zuweisung  $\{\pi_1, \pi_2, \dots, \pi_n\}$  der Daten  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  zu den Clustern mit

$$\pi_i \in [0, 1]^k \quad \sum_{j=1}^k \pi_{ij} = 1$$

- Annahme:
  - ▣ Daten-erzeugendes Modell ist Kombination von Gauß-Verteilungen mit unterschiedlichen Mittelwerten und Kovarianzen

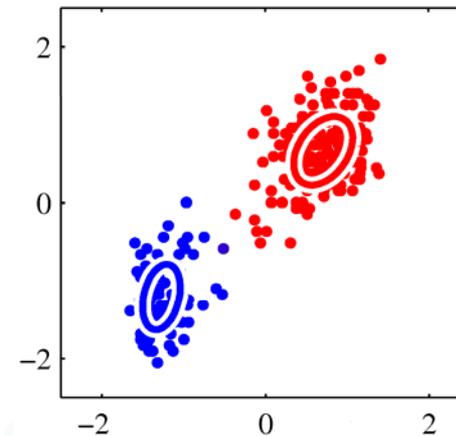
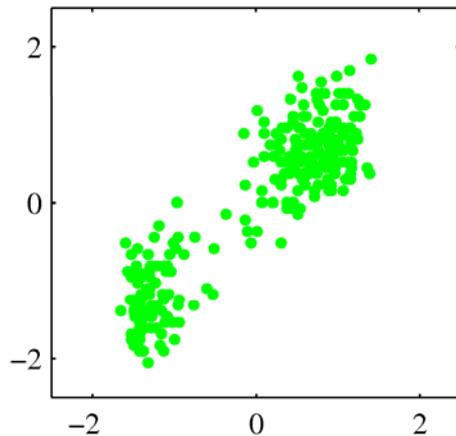
$$p(\mathbf{X} | \Theta) = \prod_{i=1}^n p(\mathbf{x}_i | \Theta) = \prod_{i=1}^n \left( \sum_{j=1}^k \pi_{ij} N(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \right)$$

mit Parametern  $\Theta = (\{\pi_{ij}\}, \{\boldsymbol{\mu}_j\}, \{\boldsymbol{\Sigma}_j\})$ .

# Ähnlichkeitsbasierendes Clustern:

## Mischmodelle: Mixture of Gaussians

- Schätzen der Parameter durch abwechselnde Optimierung analog zu K-Means (EM-Algorithmus).



# Ähnlichkeitsbasierendes Clustern:

## Mischmodelle: Mixture of Gaussians

### □ Vorteile:

- Probabilistische Zuweisung zu Clustern.
- Anzahl Cluster muss nicht vorgeben werden.
  - Automatischer Trade-off zwischen Anzahl Clustern und Anpassung an Daten.

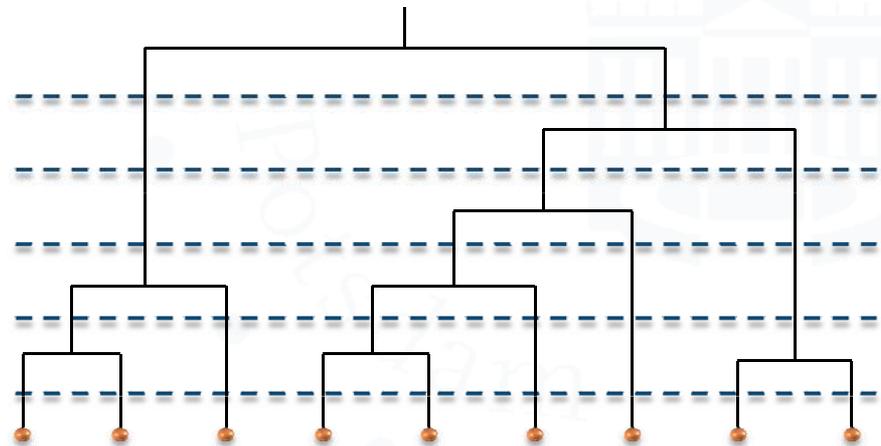
### □ Nachteile:

- Langsamer und komplexer als K-Means.
- Nur für numerische Attribute geeignet.

# Ähnlichkeitsbasierendes Clustern:

## Verlinkungsbasierte Verfahren

- Gegeben:
  - ▣ Ungelabelte Daten  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ .
  - ▣ Anzahl vermuteter Cluster  $k$  muss nicht bekannt sein.
  - ▣ Abstandsmaß *dist* zwischen Datenpunkten.
- Gesucht:
  - ▣ Darstellung der Daten in Form eines *Dendrogramms*.



# Ähnlichkeitsbasierendes Clustern:

## Verlinkungsbasierte Verfahren

### □ Idee:

#### ■ *Agglomerative Hierarchical Clustering.*

- Zu Beginn bildet jeder Datenpunkt einen eigenen Cluster.
- Benachbarte Cluster werden sukzessive verschmolzen (bottom-up).

#### ■ *Divisive Hierarchical Clustering.*

- Zu Beginn bilden alle Daten einen gemeinsamen Cluster.
- Cluster werden sukzessive gesplittet (top-down).

### □ Ziel:

- Iteratives Aufbauen des Clusterings bis vorgegebene Qualität erreicht ist.

# Ähnlichkeitsbasierendes Clustern:

## Verlinkungsbasierte Verfahren: Agglomerative Nesting (AGNES)

### □ Algorithmus:

AGNES(*Instanzen*  $\mathbf{x}_i$ )

Setze  $C_i = \{\mathbf{x}_i\} \forall i$

DO

$D_{ij} = \text{dist}(C_i, C_j) \forall i, j$

$(i^*, j^*) = \arg \min_{i, j} (D_{ij})$

Verschmelze  $C_{i^*}$  und  $C_{j^*}$

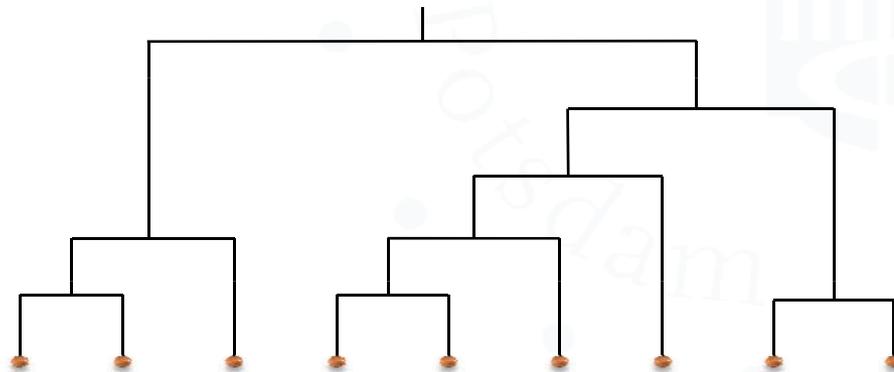
WHILE  $D_{i^*j^*} < \varepsilon$

RETURN  $\{C_i\}$

Jeder Datenpunkt ein eigener Cluster

Distanz zwischen zwei Clustern?

Mindest-Qualität (Abbruchbedingung)



# Ähnlichkeitsbasierendes Clustern:

## Verlinkungsbasierte Verfahren: Agglomerative Nesting (AGNES)

### □ Distanz zwischen zwei Clustern:

□ Single Linkage: 
$$dist(C_i, C_j) = \min_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} (dist(\mathbf{x}, \mathbf{y}))$$

□ Complete Linkage: 
$$dist(C_i, C_j) = \max_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} (dist(\mathbf{x}, \mathbf{y}))$$

□ Average Linkage: 
$$dist(C_i, C_j) = \frac{1}{|C_i| |C_j|} \sum_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} dist(\mathbf{x}, \mathbf{y})$$

□ Average Group Linkage: 
$$dist(C_i, C_j) = \frac{1}{|C_i \cup C_j|} \sum_{\mathbf{x}, \mathbf{y} \in C_i \cup C_j} dist(\mathbf{x}, \mathbf{y})$$

□ Centroid: 
$$dist(C_i, C_j) = dist \left( \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}, \frac{1}{|C_j|} \sum_{\mathbf{y} \in C_j} \mathbf{y} \right)$$

# Ähnlichkeitsbasierendes Clustern:

## Verlinkungsbasierte Verfahren: Agglomerative Nesting (AGNES)

### □ Vorteile:

- Einfach zu implementieren
- Relativ schnell.
- Iteratives Verfahren, für Online-Clustering geeignet.
- Für nominale, ordinale und numerische Attribute geeignet.
- Anzahl Cluster muss nicht vorgegeben werden.

### □ Nachteile:

- Nur lokales Optimum garantiert.
- Harte Zuweisung zu Clustern (nicht-probabilistisch).
- Maß für Cluster-Distanz & Abbruchbedingung muss vorgegeben werden.

# Ähnlichkeitsbasierendes Clustern:

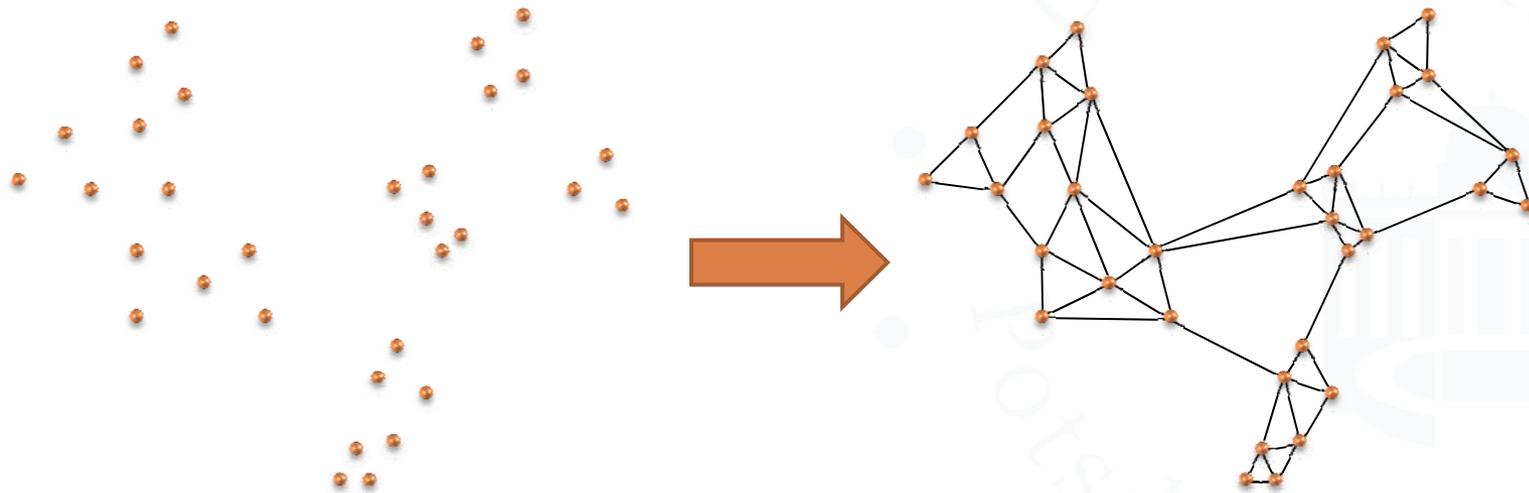
## Graphentheoretische Verfahren

- Gegeben:
  - Ungelabelte Daten  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ .
  - Anzahl vermuteter Cluster  $k$  mit  $1 < k < n$ .
  - Ähnlichkeitsmaß (Kernel) *sim* zwischen Datenpunkten.
- Gesucht: Partitionierung der Daten in  $k$  Cluster.
- Ziel: Hohe Ähnlichkeit zw. Punkten im selben Cluster und geringe Ähnlichkeit zw. Punkten verschiedener Cluster.
  - Repräsentation der Daten als Graph und Partitionieren des Graphs.

# Ähnlichkeitsbasierendes Clustern:

## Graphentheoretische Verfahren

- Ähnlichkeit zwischen Datenpunkten (Knoten) bilden gewichtete Kanten:



# Ähnlichkeitsbasierendes Clustern:

## Graphentheoretische Verfahren



### □ Konstruktion des Graphen:

- $\varepsilon$ -Neighborhood-Graph: Verbinde Knoten  $\mathbf{x}_i$  mit  $\mathbf{x}_j$  falls  $sim(\mathbf{x}_i, \mathbf{x}_j) > \varepsilon$ , und setze Kantengewicht auf 1.
- $k$ -Nearest-Neighbor-Graph: Verbinde Knoten  $\mathbf{x}_i$  mit  $\mathbf{x}_j$  falls  $\mathbf{x}_i$  einer der  $k$ -nächsten Nachbarn von  $\mathbf{x}_j$  ist oder/und  $\mathbf{x}_j$  einer der  $k$ -nächsten Nachbarn von  $\mathbf{x}_i$  ist, und setze Kantengewicht auf  $sim(\mathbf{x}_i, \mathbf{x}_j)$ .
- Vollständiger Graph: Verbinde alle Knoten  $\mathbf{x}_i$  mit  $\mathbf{x}_j$ , und setze Kantengewicht auf  $sim(\mathbf{x}_i, \mathbf{x}_j)$ .

### □ Partitionierung des Graphen:

- Kanten zwischen Clustern (Teilgraphen) haben geringe Gewichte (geringe inter-cluster similarity).
- Kanten innerhalb eines Clusters haben hohe Gewichte (hohe intra-cluster similarity).

# Ähnlichkeitsbasierendes Clustern:

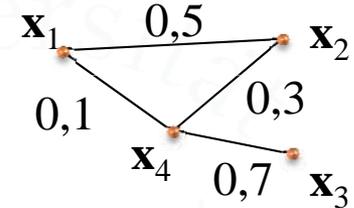
## Graphentheoretische Verfahren

### □ Repräsentation eines (ungerichteten) Graphen:

#### □ Adjazenzmatrix:

- Enthält Kantengewichte  $A_{ij}$ .

$$A = \begin{bmatrix} 1 & 0,5 & 0 & 0,1 \\ 0,5 & 1 & 0 & 0,3 \\ 0 & 0 & 1 & 0,7 \\ 0,1 & 0,3 & 0,7 & 1 \end{bmatrix}$$



#### □ Knotengrad-Matrix:

$$D = \begin{bmatrix} d_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & d_n \end{bmatrix}$$

$$d_i = \sum_{j=1}^n A_{ij}$$

#### □ Laplace-Matrix:

- Unnormalisiert:
- Normalisiert (Random Walk):
- Symmetrisch normalisiert:

$$L_{un} = D - A$$

$$L_{rw} = I - D^{-1}A$$

$$L_{sym} = I - D^{-1/2}AD^{-1/2}$$

# Ähnlichkeitsbasierendes Clustern:

## Graphentheoretische Verfahren: Spectral Clustering

### □ Algorithmus:

SpecClust(*Instanzen*  $\mathbf{x}_i$ , *Clusteranzahl*  $k$ )

Konstruiere Graph aus *Instanzen*  $\mathbf{x}_i$

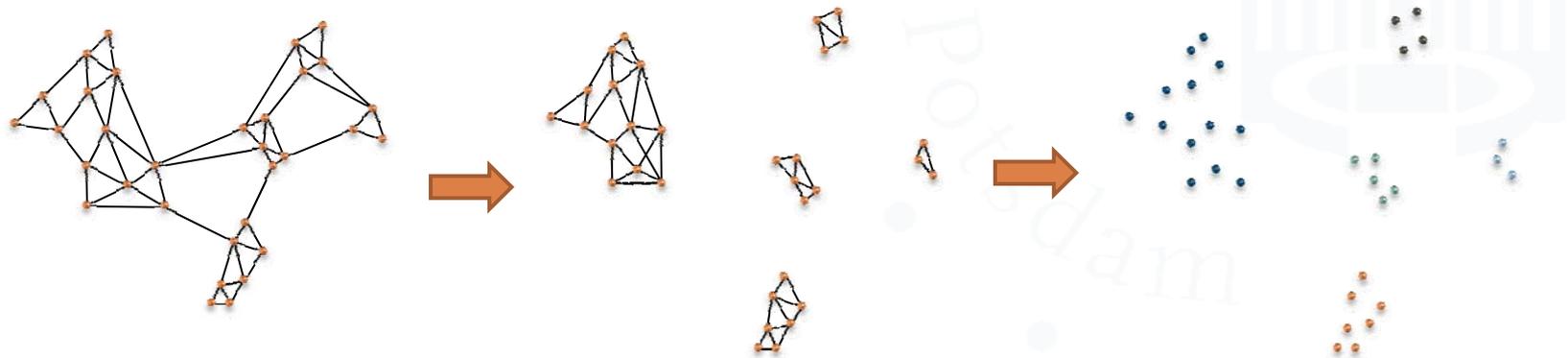
Berechne zugehörige Laplace-Matrix  $\mathbf{L}$

Berechne die  $k$  Eigenvektoren  $\mathbf{v}_i \in \mathbb{R}^n$  mit den  $k$  kleinsten Eigenwerten

Setze  $[\mathbf{x}'_1 \ \mathbf{x}'_2 \ \dots \ \mathbf{x}'_n] = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_k]^T$

$\{\mathbf{r}_1, \dots, \mathbf{r}_n\} = \text{K-Means}(\text{Instanzen } \mathbf{x}'_i, \text{Clusteranzahl } k)$

RETURN  $\{\mathbf{r}_1, \dots, \mathbf{r}_n\}$



# Ähnlichkeitsbasierendes Clustern:

## Graphentheoretische Verfahren: Spectral Clustering

### □ Vorteile:

- Cluster können beliebige Form haben.
- Einfach zu implementieren.
- Relativ schnell (falls Laplace-Matrix sparse).
- Meist hohe Qualität.

### □ Nachteile:

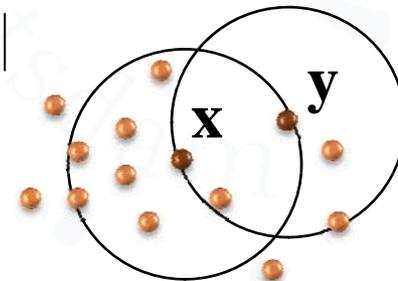
- Nur lokales Optimum garantiert.
- Harte Zuweisung zu Clustern (nicht-probabilistisch).
- Anzahl Cluster muss vorgeben werden.

# Dichtebasierendes Clustern:

## Spatial Clustering

- Idee: Cluster-Grenzen sollen durch Regionen mit geringer Datendichte verlaufen.
  - ▣ Zwei Punkte  $x$  und  $y$  gehören zum gleichen Cluster falls es einen „Pfad“ zwischen beiden Punkten gibt.
  
- *Direkte Erreichbarkeit* des Punktes  $y$  von  $x$ :
  - ▣ Nachbarschaftsbedingung:  $y \in N_\varepsilon(\mathbf{x}) = \{\mathbf{z} \mid \text{dist}(\mathbf{x}, \mathbf{z}) < \varepsilon\}$
  - ▣ *Core-Point*-Bedingung:  $\#_{\min} \leq |N_\varepsilon(\mathbf{x})|$

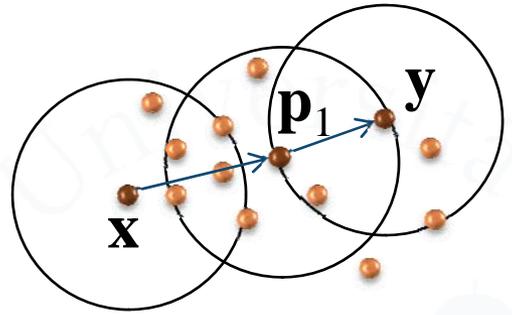
Mindestanzahl von  
benachbarten Punkten



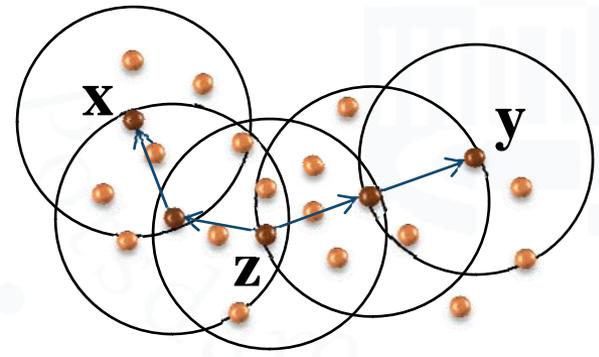
# Dichtebasierendes Clustern:

## Spatial Clustering

- Punkt  $y$  ist von  $x$  erreichbar:
  - Es gibt einen Pfad  $p_1, p_2, \dots, p_m$  mit  $p_1 = x$  und  $p_m = y$ , so dass  $p_{i+1}$  direkt erreichbar von  $p_i$  ist für  $i = 1 \dots m$ .



- Punkt  $y$  ist mit  $x$  verbunden:
  - Es gibt einen Punkt  $z$ , so dass  $x$  und  $y$  erreichbar von  $z$  sind.



# Dichtebasierendes Clustern:

## Spatial Clustering: DBSCAN

### □ Algorithmus:

DBSCAN(*Instanzen*  $\mathbf{x}_i$ , *Radius*  $\varepsilon$ , *Mindest-Clustergröße*  $g$ )

Setze  $l=0, D = \{\mathbf{x}_i\}$

DO

Wähle zufällig  $\mathbf{x} \in D$

$P = \{\mathbf{y} \mid \mathbf{y} \text{ ist von } \mathbf{x} \text{ erreichbar}\}$

IF  $|P| < g$

$D = D \setminus \mathbf{x}$

ELSE

$l = l + 1$

$C_l = P$

$D = D \setminus P$

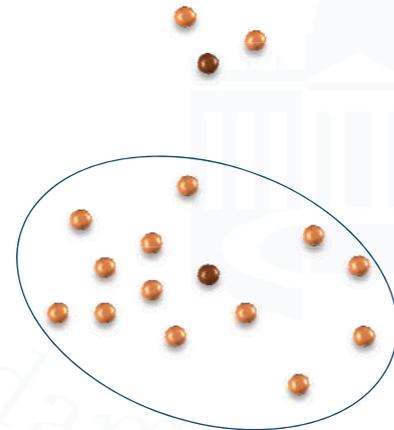
WHILE  $D \neq \emptyset$

RETURN  $\{C_1, \dots, C_l\}$

$P$  enthält  $\mathbf{x}$

$\mathbf{x}$  ist Outlier  
(nicht Teil eines Clusters)

$\mathbf{x}$  ist Core-Point eines  
Clusters



# Dichtebasierendes Clustern:

## Spatial Clustering: DBSCAN

### □ Vorteile:

- Cluster können beliebige Form haben.
- Für nominale, ordinale und numerische Attribute geeignet.
- Anzahl Cluster muss nicht vorgegeben werden.

### □ Nachteile:

- Nur lokales Optimum garantiert.
- Harte Zuweisung zu Clustern (nicht-probabilistisch).
- Radius & Mindest-Clustergröße müssen vorgegeben werden.

# Zusammenfassung

- Prototyp-Verfahren (z.B. K-Means):
  - Schnell, numerische Attribute, bekannte Clusteranzahl.
- Mischmodelle (z.B. Mixture of Gaussians):
  - Probabilistische Cluster-Zuweisung, numerische Attribute, unbekannte Clusteranzahl.
- Verlinkungsbasierte Verfahren (z.B. AGNES):
  - Iterativ, beliebige Attribute, unbekannte Clusteranzahl.
- Graphbasierte Verfahren (z.B. Spectral Clustering):
  - Beliebige Clusterform, hohe Qualität.
- Spatial Clustering (z.B. DBSCAN):
  - Beliebige Clusterform, unbekannte Clusteranzahl.