





Skript zur Vorlesung:

Einführung in die Informatik: Systeme und Anwendungen Sommersemester 2012

Kapitel 4: Data Mining

Vorlesung: Prof. Dr. Christian Böhm Übungen: Sebastian Goebl

Skript © 2010 Christian Böhm, Peer Kröger, Arthur Zimek

http://www.dbs.ifi.lmu.de/cms/ Einführung in die Informatik Systeme und Anwendungen





Überblick



4.1 Einleitung

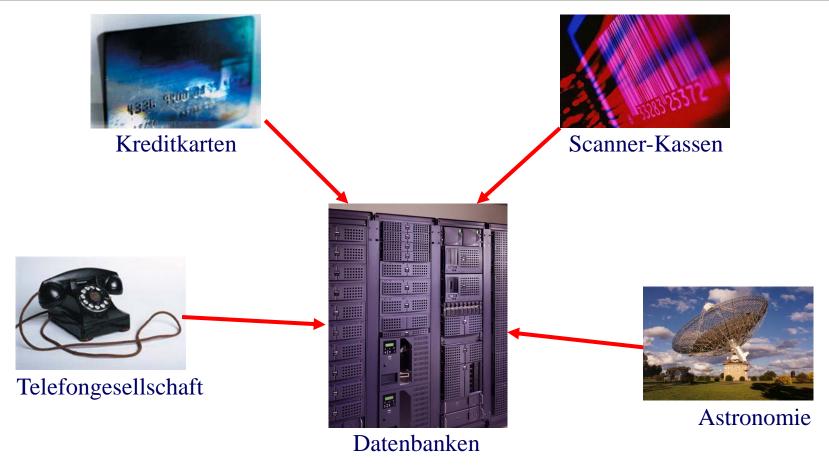
4.2 Clustering

4.3 Klassifikation



Motivation





- Riesige Datenmengen werden in Datenbanken gesammelt
- Analysen können nicht mehr manuell durchgeführt werden



Definition KDD



- Knowledge Discovery in Databases (KDD) ist der Prozess der (semi-) automatischen Extraktion von Wissen aus Datenbanken, das
 - gültig
 - bisher unbekannt
 - und potentiell nützlich ist.

Bemerkungen:

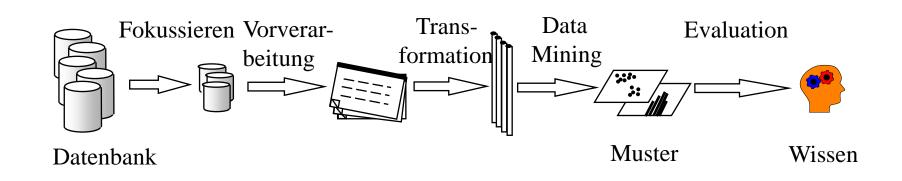
- (*semi-*) *automatisch*: im Unterschied zu manueller Analyse. Häufig ist trotzdem Interaktion mit dem Benutzer nötig.
- *gültig*: im statistischen Sinn.
- bisher unbekannt: bisher nicht explizit, kein "Allgemeinwissen".
- potentiell nützlich: für eine gegebene Anwendung.



Der KDD-Prozess (Modell)



Prozessmodell nach Fayyad, Piatetsky-Shapiro & Smyth



Verwaltung (File/DB) Selektion relevanter Daten

Beschaffung der Daten

Fokussieren:

Vorverarbeitung:

Integration von Daten aus unterschiedlichen QuellenVervollständigung

Transformation

Konsistenzprüfung

 Diskretisierung numerischer Merkmale Ableitung neuer MerkmaleSelektion relevanter Merkm

Data MiningGenerierung der Musterbzw. Modelle

Evaluation

Bewertung der Interessantheit durch den Benutzer

 Validierung: Statistische Prüfung der Modelle

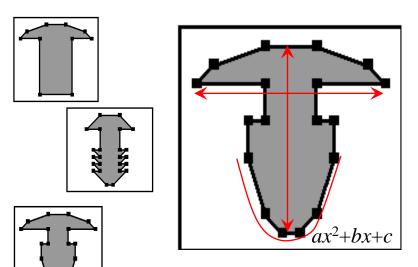


Objekt-Merkmale (Feature)



- Oft sind die betrachteten Objekte komplex
- Eine Aufgabe des KDD-Experten ist dann, geeignete Merkmale (*Features*) zu definieren bzw. auszuwählen, die für die Unterscheidung (Klassifikation, Ähnlichkeit) der Objekte relevant sind.

Beispiel: CAD-Zeichnungen:



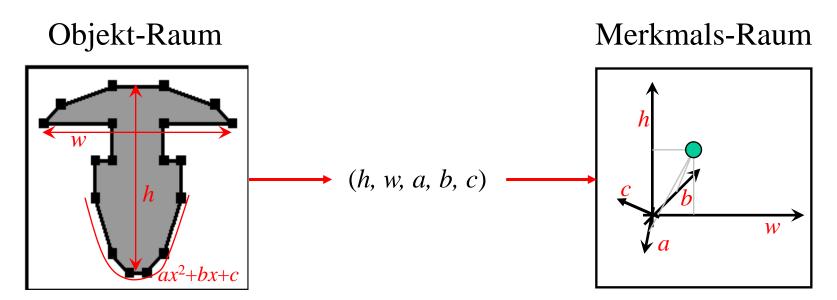
Mögliche Merkmale:

- Höhe h
- Breite w
- Kurvatur-Parameter (*a*,*b*,*c*)



Feature-Vektoren





- Im Kontext von statistischen Betrachtungen werden die Merkmale häufig auch als *Variablen* bezeichnet
- Die ausgewählten Merkmale werden zu Merkmals-Vektoren (*Feature Vector*) zusammengefasst
- Der Merkmalsraum ist häufig hochdimensional (im Beispiel 5-dim.)

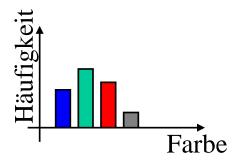


Feature-Vektoren (weitere Beispiele)



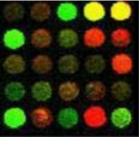
Bilddatenbanken: Farbhistogramme



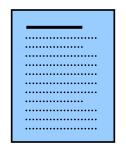


Gen-Datenbanken: Expressionslevel





Text-Datenbanken: Begriffs-Häufigkeiten



Data 25
Mining 15
Feature 12
Object 7
...

Der Feature-Ansatz ermöglicht einheitliche Behandlung von Objekten verschiedenster Anwendungsklassen



Feature: verschiedene Kategorien



Nominal (kategorisch)

Charakteristik:

Nur feststellbar, ob der Wert gleich oder verschieden ist. Keine Richtung (besser, schlechter) und kein Abstand.

Merkmale mit nur zwei Werten nennt man *dichotom*.

Beispiele:

Geschlecht (dichotom) Augenfarbe Gesund/krank (dichotom)

Ordinal

Charakteristik:

Es existiert eine Ordnungsrelation (besser/schlechter) zwischen den Kategorien, aber kein einheitlicher Abstand.

Beispiele:

Schulnote (metrisch?) Güteklasse Altersklasse

Metrisch

Charakteristik:

Sowohl Differenzen als auch Verhältnisse zwischen den Werten sind aussagekräftig. Die Werte können diskret oder stetig sein.

Beispiele:

Gewicht (stetig) Verkaufszahl (diskret) Alter (stetig oder diskret)



Ähnlichkeit von Objekten



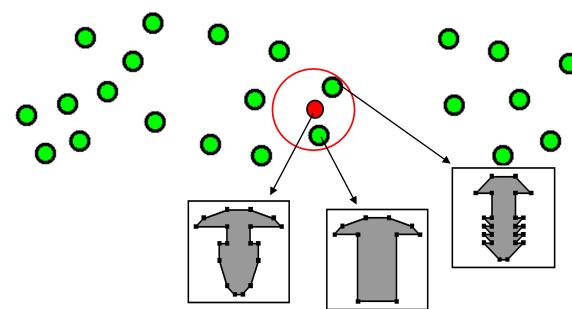
- Spezifiziere Anfrage-Objekt $q \in DB$ und...
 - ... suche ähnliche Objekte Range-Query (Radius ε)

$$RQ(q,\varepsilon) = \{ o \in DB \mid \delta(q,o) \le \varepsilon \}$$

– ... suche die *k* ähnlichsten Objekte – Nearest Neighbor Query

 $NN(q,k) \subseteq DB$ mit k Objekten, sodass

$$\forall o \in NN(q,k), p \in DB \setminus NN(q,k) : \delta(q,o) \leq \delta(q,p)$$



alternative Schreibweise für Mengendifferenz:

$$A \setminus B = A - B$$



Ähnlichkeitsmaße im Feature-Raum



Euklidische Norm (L₂):

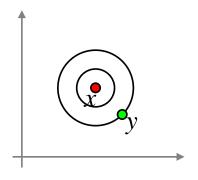
$$\delta_1(x,y) = ((x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots)^{1/2}$$

Manhattan-Norm (L_1) :

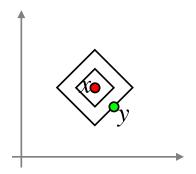
$$\delta_2(x,y) = |x_1 - y_1| + |x_2 - y_2| + \dots$$

Maximums-Norm (L_{∞}):

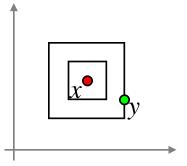
$$\delta_{\infty}(x,y) = \max\{|x_1-y_1|, |x_2-y_2|,...\}$$



Abstand in Euklidischen Raum (natürliche Distanz)



Die Unähnlichkeiten der einzelnen Merkmale werden direkt addiert



Die Unähnlichkeit des am wenigsten ähnlichen Merkmals zählt

Verallgemeinerung L_p -Abstandsmaß:

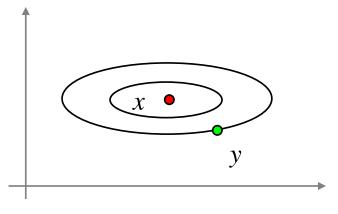
$$\delta_{p}(x,y) = (|x_1 - y_1|^p + |x_2 - y_2|^p + ...)^{1/p}$$



Gewichtete Ähnlichkeitsmaße

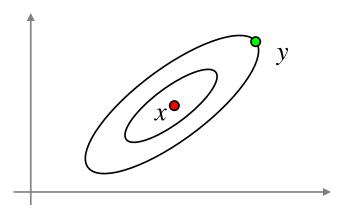


• Viele Varianten gewichten verschiedene Merkmale unterschiedlich stark.



Gewichtete Euklidische Distanz

$$\delta_{p,w}(x,y) = \sqrt[p]{\sum_{i=1}^d w_i \cdot |x_i - y_i|^p}$$



Mahalanobis Distanz

$$\delta_{\Sigma}(x, y) = \sqrt{(x - y)^T \cdot \Sigma^{-1} \cdot (x - y)}$$

 $\Sigma = \text{Kovarianz-Matrix}$



Kategorien von Data Mining

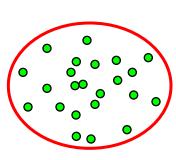


- Wichtigste Data-Mining-Verfahren auf Merkmals-Vektoren:
 - Clustering
 Outlier Detection
 Klassifikation
 Regression
 normalerweise unsupervised
 normalerweise supervised
- <u>Supervised</u>: In Trainingsphase wird eine Funktion gelernt, die in der Testphase angewandt wird.
- <u>Unsupervised</u>: Es gibt keine Trainingsphase. Die Methode findet Muster, die einem bestimmten Modell entsprechen.
- Darüber hinaus gibt es zahlreiche Verfahren, die nicht auf Merkmalsvektoren, sondern direkt auf Texten, Mengen, Graphen usw. arbeiten.



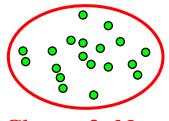
Clustering







0



Cluster 2: Nägel

Ein Grundmodell des Clustering ist:

Zerlegung (Partitionierung) einer Menge von Objekten (bzw. Feature-Vektoren) in Teilmengen (Cluster), so dass

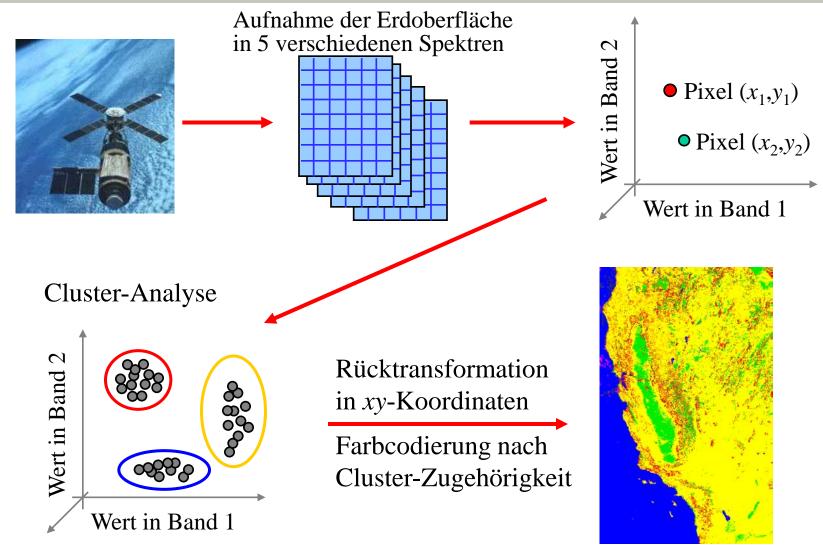
- die Ähnlichkeit der Objekte innerhalb eines Clusters maximiert
- die Ähnlichkeit der Objekte verschiedener Cluster minimiert wird

Idee: Die verschiedenen Cluster repräsentieren meist unterschiedliche Klassen von Objekten; bei evtl. unbek. Anzahl und Bedeutung der Klassen



Anwendung: Thematische Karten





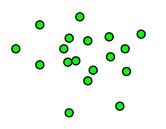


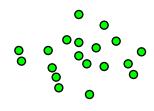
Outlier Detection





Datenfehler? Betrug?





Outlier Detection bedeutet: Ermittlung von untypischen Daten

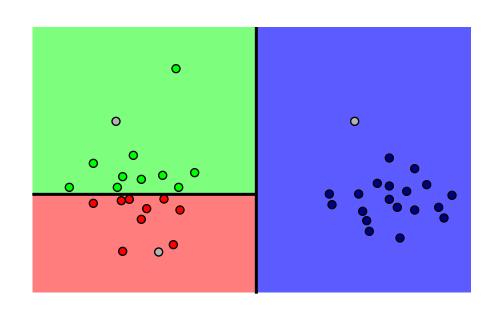
Anwendungen:

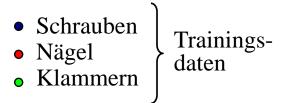
- Entdeckung von Missbrauch etwa bei
 - Kreditkarten
 - Telekommunikation
- Datenbereinigung (Messfehler)



Klassifikation







Neue Objekte

Aufgabe:

Lerne aus den bereits klassifizierten *Trainingsdaten* die *Regeln*, um neue Objekte nur aufgrund der Merkmale zu klassifizieren

Das Ergebnismerkmal (Klassenvariable) ist nominal (kategorisch)

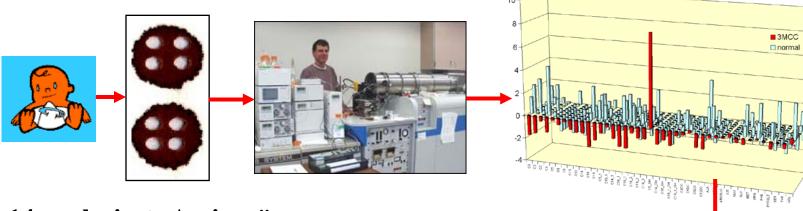


Anwendung: Neugeborenen-Screening LMU



Blutprobe des Neugeborenen Massenspektrometrie

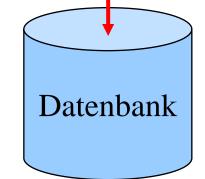
Metabolitenspektrum



14 analysierte Aminosäuren:

alanine
arginine
argininosuccinate
citrulline
glutamate
glycine
methionine

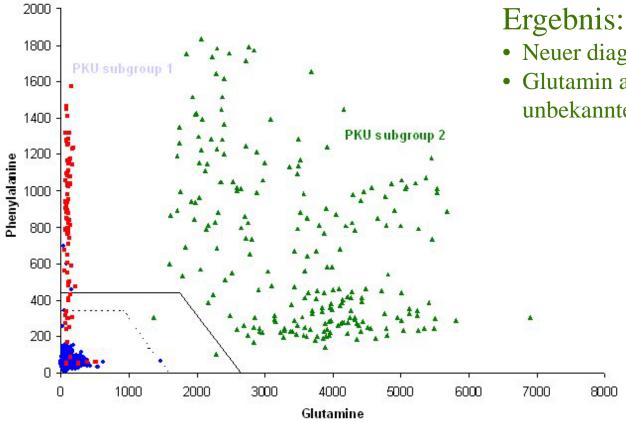
phenylalanine
pyroglutamate
serine
tyrosine
valine
leuzine+isoleuzine
ornitine





Anwendung: Neugeborenen-Screening



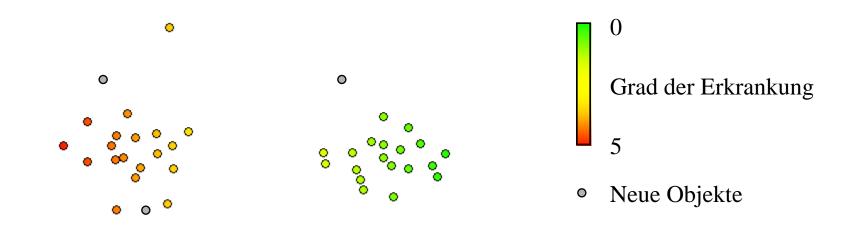


- Neuer diagnostischer Test
- Glutamin als bisher unbekannter Marker



Regression





Aufgabe:

Ähnlich zur Klassifikation, aber das Ergebnis-Merkmal, das gelernt bzw. geschätzt werden soll, ist *metrisch*.