

# Mathe III: Statistische Methoden

## Teil 2

*Garance PARIS*

*Sommersemester 2010*

# Data Mining

*1. Juli 2010*

- Künstliche Generierung von Wissen aus Erfahrung
- Erkennung komplexer Muster und Regelmäßigkeiten in vorhandenen Daten
- Ziel: Verallgemeinerung
  - Über das Nachschlagen bereits gesehener Beispiele hinausgehen
  - Beurteilung unbekannter Daten
- Beispiele:
  - Die Gleichung einer Geraden an Hand zweier Punkten bestimmen
  - Einen unbekanntem Text mit Wortkategorien annotieren

- Künstliche Generierung von Wissen aus Erfahrung
- Erkennung komplexer Muster und Regelmäßigkeiten in vorhandenen Daten
- Ziel: Verallgemeinerung
  - Über das Nachschlagen bereits gesehener Beispiele hinausgehen
  - Beurteilung unbekannter Daten
- Beispiele:
  - Die Gleichung einer Geraden an Hand zweier Punkten bestimmen
  - Einen unbekanntem Text mit Wortkategorien annotieren

Analyse im Voraus **strukturierter** Daten

## Beispiel

Der Manager eines Golf-Clubs möchte wissen, wann er viele Kunden zu erwarten hat, damit er Studenten als Aushilfe einstellen kann, und wann keiner spielen will, damit er seinen Angestellten freigeben kann.

Zwei Wochen lang führt er Buch darüber, wie das Wetter ist und ob er viele oder wenige Kunden an dem Tag hat.

Er schreibt sich auf:

- ob das Wetter heiter, bewölkt oder regnerisch ist,
- wie warm es ist,
- wieviel Luftfeuchtigkeit es gibt,
- ob der Wind stark weht oder nicht,
- und wieviele Kunden er an dem Tag hat

## Beispiel: Der Wetterdatensatz

Das Ergebnis ist für einen Menschen nicht besonders gut durchschaubar. . .

Outlook	Temp.	Humidity	Windy	Play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

## *Was kann gelernt werden?*

Einige Aufgabetypen:

**Klassifikation:**

Instanzen vordefinierten Kategorien zuweisen

**Clustering:**

Instanzen ihrer Ähnlichkeit nach in Gruppen aufteilen (ohne vordefinierte Kategorien)

**Numerische Vorhersage:**

Der numerische Wert eines Merkmals auf Grund der Werte aller anderen Merkmale bestimmen

## Beispiel: Ein numerischer Datensatz

Outlook	Temp.	Humidity	Windy	# Cust.
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

## Beispiel: Ein numerischer Datensatz

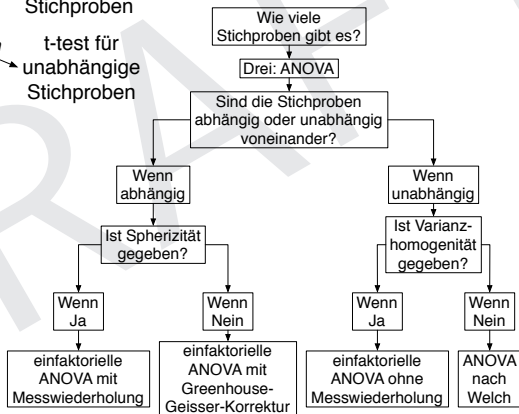
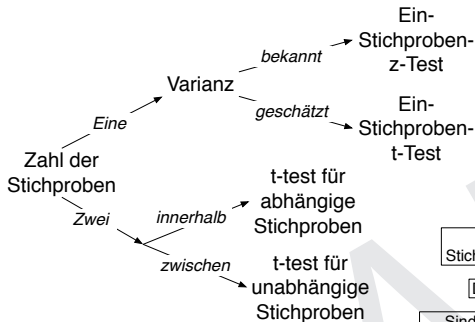
Outlook	Temp.	Humidity	Windy	# Cust.
sunny	85	85	false	12
sunny	80	90	true	10
overcast	83	86	false	30
rainy	70	96	false	15
rainy	68	80	false	16
rainy	65	70	true	5
overcast	64	65	true	24
sunny	72	95	false	10
sunny	69	70	false	18
rainy	75	80	false	20
sunny	75	70	true	25
overcast	72	90	true	18
overcast	81	75	false	32
rainy	71	91	true	8

## Klassifikation: Terminologie

- Konzept:** Das, was gelernt werden soll  
Bsp.: Ob unter angegebenen Bedingungen gerne gespielt wird oder nicht
- Instanz:** Ein einziges Beispiel im Datensatz  
Bsp.: Das Wetter an einem der 14 Tagen
- Attribut:** Ein Merkmal einer Instanz  
Bsp.: *outlook, temperature, humidity, windy*
- Wert:** Wert eines Attributs  
Bsp.: Für *outlook*: sunny, overcast, rainy
- Kategorie oder Klasse:**  
Ziel der Klassifikation, Bsp.: yes, no

# Entscheidungsbäume

*1. Juli 2010*



## *Entscheidungsbäume*

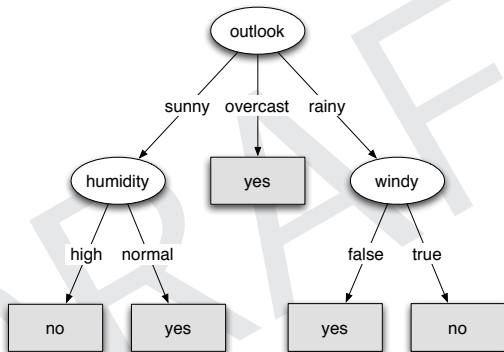
### *Allgemein*

Ein von einem Experten erstellter Baum,  
um aufeinanderfolgende, hierarchische Entscheidungen  
zu veranschaulichen

### *Data Mining*

Ein auf Basis verfügbarer Daten erstelltes Modell,  
das benutzt werden kann, um Voraussagen über  
neue Daten zu machen

## Baum für den Wetter-Datensatz



## *ID3: Intuition*

- Bestimmen, welches Attribut die Daten am besten klassifiziert
- Dieses Attribut als Wurzel des Baums verwenden
- Einen Zweig für jeden Wert erstellen, den das Attribut annehmen kann
- Dieses Prozess für jeden Baumzweig wiederholen, jeweils mit der Untermenge der zu klassifizierenden Daten
- Rekursiv weitermachen, bis die Klassifikation unter allen Blättern eindeutig ist

## *Auswahl des klassifizierenden Attribut*

**Frage:** Wie wird das beste Attribut bestimmt?

**Antwort:** Es wird das Attribut gewählt, das  
... die meiste Information beiträgt

## Auswahl des klassifizierenden Attribut

**Frage:** Wie wird das beste Attribut bestimmt?

**Antwort:** Es wird das Attribut gewählt, das  
... die meiste Information beiträgt  
... die Entropie im Datensatz am meisten reduziert

**Algorithmus:** *ID3* (Kommerzieller Nachfolger: *C4.5*)

## Informationsgewinn

Für jeden Attribut  $A$ :

$$\text{Gain}(S, A) = E(S) - \sum_{v \in A} \frac{|S_v|}{|S|} E(S_v)$$

$S$  : Datensatz

$E(S)$  : Entropie im Datensatz oder Untermenge davon

$S_v$  : Untermenge von  $S$  für die  $A$  den Wert  $v$  hat

$|S|$  : Mächtigkeit von  $S$

$|S_v|$  : Mächtigkeit von  $S_v$

Der Informationsgewinn  $\text{Gain}(S, A)$  ist die Reduzierung der Entropie, die man dadurch erwartet, dass man den Wert von Attribut  $A$  kennt

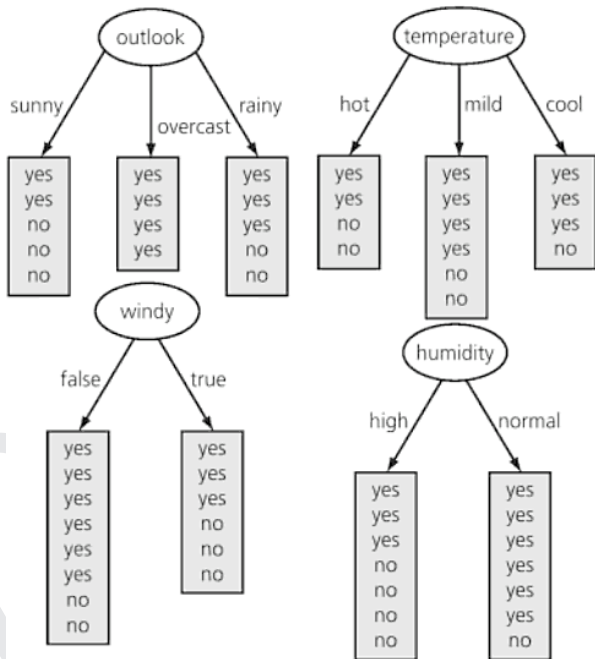
$$E(S) = - \sum_{k=1}^c p_k \log_2 p_k$$

$c$ : Anzahl Kategorien

$p_k$ : Anteil der Instanzen in  $S$ , die Kategorie  $k$  angehören

### Informationstheoretische Interpretation

- $E(S)$ : Anzahl Bits, die benötigt werden, um ein beliebiges Element aus  $S$  zu kodieren
- $E(S)$  ist gleich null, wenn alle Instanzen von  $S$  der selben Kategorie angehören
- Bei einer binären Klassifikation ist  $E(S)$  gleich eins, wenn  $S$  gleich viele Instanzen aus jeder Klasse enthält



- Erst die Gesamtentropie berechnen:  
Von 14 Instanzen werden 5 als “*play=yes*”  
klassifiziert, 9 als “*play=no*”

$$\begin{aligned} E(S) &= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \\ &= 0.41 + 0.53 = 0.94 \end{aligned}$$

- Dann für jedes Attribut und für jeden Wert  
die Entropie berechnen, z. B. für “*windy = false*”:

$$E(S_{false}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

Dies kann wie folgt umgeformt werden:

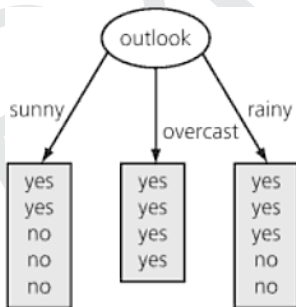
$$E(S_{false}) = \frac{8 \log_2 8 - 6 \log_2 6 - 2 \log_2 2}{8} = 0.81$$

## Beispiel mit dem Wetterdatensatz

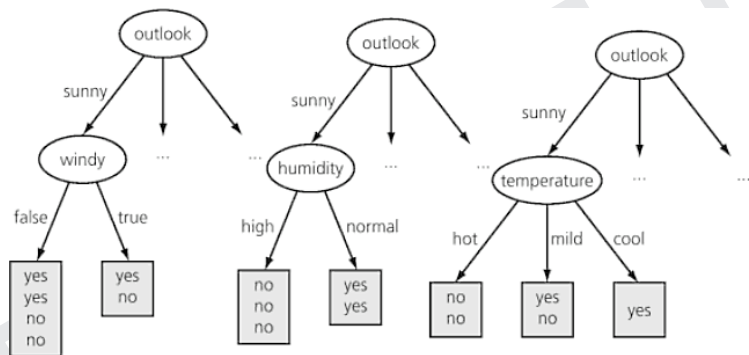
- Bei "*windy = true*" kommen dreimal "*play = yes*" vor und dreimal "*play = no*", daher ist die Entropie gleich eins
- Daraus und aus der Gesamtentropie setzt sich  $\text{Gain}(S, \text{windy})$  zusammen:

$$\begin{aligned}\text{Gain}(S, \text{windy}) &= E(S) - \frac{|S_{\text{false}}|}{|S|} E(S_{\text{false}}) - \frac{|S_{\text{true}}|}{|S|} E(S_{\text{true}}) \\ &= 0.94 - \frac{8}{14} * 0.81 - \frac{6}{14} * 1 = 0.048\end{aligned}$$

- Ähnlich werden berechnet:  
 $Gain(S, windy) = 0.048$   
 $Gain(S, outlook) = 0.247$   
 $Gain(S, temperature) = 0.029$   
 $Gain(S, humidity) = 0.152$
- *Outlook* erzielt den größten Informationsgewinn, daher wird es als Baumwurzel gewählt



## Informationsgewinn: Beispiel

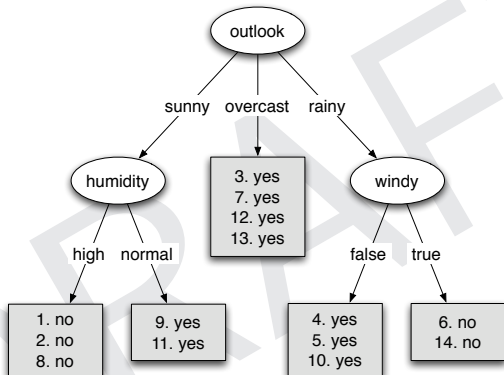


$$\text{Gain}(S, \text{windy}) = 0.020$$

$$\text{Gain}(S, \text{temperature}) = 0.571$$

$$\text{Gain}(S, \text{humidity}) = 0.971$$

*Humidity* erzielt den größten Informationsgewinn, daher wird es als nächsten Knoten unter "*outlook = sunny*" eingesetzt



- Alle Instanzen unter einem Blatt werden gleich klassifiziert (Entropie = 0)
- Das Attribut *temperature* wird nicht verwendet

## *Interpretation*

- Wenn das Wetter bewölkt ist, spielen Leute immer
- Es gibt Leute, die sogar bei Regen spielen, aber nicht wenn der Wind weht
- Wenn die Sonne scheint, spielen Leute gerne, außer wenn die Luftfeuchtigkeit hoch ist

### Schlußfolgerungen:

- Der Manager kann an Tagen, an denen Regen und starken Wind vorausgesagt wird, seinem Personal frei geben
- An sonnigen Tagen, an denen es Luft gibt, kann er sich überlegen, ob er zusätzlich ein paar Studenten einstellt

## Repräsentation

- Nicht-terminale Knoten testen ein Attribut
- Baumzweige entsprechen den Werten, die ein Attribut annehmen kann
- Blätter weisen Instanzen einer Kategorie zu
- Der Pfad zwischen Wurzel und Blatt stellt eine Konjunktion dar

## Umwandlung in Entscheidungsregeln

```
IF (outlook==sunny)  $\wedge$  (humidity==high)
THEN play=no
IF (outlook==sunny)  $\wedge$  (humidity==normal)
THEN play=yes
...
```

## *Eigenschaften des Algorithmus*

- Allgemein werden flache Bäume bevorzugt
- Es findet eine Suche durch den Hypothesenraum statt
- Die Suche ist aber unvollständig, d. h. nicht alle möglichen Bäume werden in Betracht gezogen
- Es ist möglich, dass der “beste” (kleinste) Baum nicht gefunden wird  
(z.B. wenn es einen insgesamt kleineren Baum gibt, der eine andere Wurzel hat, als die, die anfangs den höchsten Informationsgewinn bietet)

# Naiv-Bayes-Klassifikation

*5. Juli 2010*

## Die Intuition

Die Wahrscheinlichkeit einer Zielkategorie für eine Instanz hängt von zwei Faktoren ab:

- die Gesamtwahrscheinlichkeit der Kategorie überhaupt, Engl. *prior probability (a-priori-Wahrscheinlichkeit)*  
Bsp.: Wie oft kommen *play=yes* bzw. *play=no* insgesamt im Datensatz vor?
- die Wahrscheinlichkeit der Kategorie bei der von der Instanz beigetragenen Evidenz oder Information, Engl. *posterior probability*  
Hier: Wahrscheinlichkeit der einzelnen Attributwerte der Instanz

## Das Maximum-a-posteriori

Wir wollen für neue Instanzen die beste Kategorie gegeben der Evidenz finden, auf Englisch das *“maximum a posteriori”*:

$$\begin{aligned}c_{map} &= \operatorname{argmax}_c P(c|e) \\ &= \operatorname{argmax}_c \frac{P(c)P(e|c)}{P(e)} \\ &= \operatorname{argmax}_c P(c)P(e|c)\end{aligned}$$

Bayes'scher Satz

$P(e)$  fällt weg, weil es konstant und unabhängig von der Hypothese ist

## Parameter-Schätzung

- $P(c)$  schätzen ist einfach: Man berechnet die relative Häufigkeit jeder Hypothese in den Trainingsdaten
- $P(e|c)$  schätzen ist wegen Sparse-Data schwieriger, da man nie genug Daten hat, damit alle mögliche Attribut-Kombinationen vorkommen
- Es wird daher angenommen, dass die Attribute unabhängig voneinander sind

$$P(e|c) = P(a_1, a_2, \dots, a_n|c) = \prod_{i=1}^n P(a_i|c)$$

- $P(a|c)$ : Relative Häufigkeit des Attributs  $a$  unter den Instanzen, die Kategorie  $c$  angehören

## Beispiel

- Klassifizierung von Instanz  $outlook = sunny$ ,  
 $temperature = cool$ ,  $humidity = high$ ,  $windy = true$ :

$$\begin{aligned}c_{map} &= \operatorname{argmax}_{c \in \{yes, no\}} P(c) \cdot \prod_i P(a_i|c) \\ &= \operatorname{argmax}_{c \in \{yes, no\}} P(c) \cdot P(sunny|c) \cdot P(cool|c) \cdot \\ &\quad P(high|c) \cdot P(true|c)\end{aligned}$$

- Berechnung der relativen Häufigkeit von  $P(sunny|yes)$ :  
Insgesamt gibt es 9 Instanzen mit  $play = yes$ ,  
davon 2 bei denen  $outlook = sunny$ ,  
also ist  $P(sunny|yes) = \frac{2}{9}$

- Zuerst die Wahrscheinlichkeit jeder Kategorie in den Trainingsdaten berechnen

Gesamthäufigkeiten für "play":  $yes = 9$ ,  $no = 5$

Wahrscheinlichkeiten:  $yes = \frac{9}{14}$ ,  $no = \frac{5}{14}$

- Für jede Kategorie die Wahrscheinlichkeit der einzelnen Attribute berechnen

		freq.		P	
		yes	no	yes	no
Outlook	sunny	2	3	$\frac{2}{9}$	$\frac{3}{5}$
	overcast	4	0	$\frac{4}{9}$	0
	rainy	3	2	$\frac{3}{9}$	$\frac{2}{5}$
	$\Sigma$	9	5	1	1
Temperature	...				
Humidity	...				
Windy	...				

## Vorgehen, Teil 2

- Wahrscheinlichkeit jeder Kategorie für die vorliegende Instanz berechnen

$$P(\text{yes}) \cdot P(\text{sunny}|\text{yes}) \cdot P(\text{cool}|\text{yes}) \cdot P(\text{high}|\text{yes}) \cdot \\ P(\text{true}|\text{yes}) = \frac{9}{14} * \frac{2}{9} * \frac{3}{9} * \frac{3}{9} * \frac{3}{9} = 0.0053$$

$$P(\text{no}) \cdot P(\text{sunny}|\text{no}) \cdot P(\text{cool}|\text{no}) \cdot P(\text{high}|\text{no}) \cdot \\ P(\text{true}|\text{no}) = \frac{5}{14} * \frac{3}{5} * \frac{1}{5} * \frac{4}{5} * \frac{3}{5} = 0.0206$$

- Kategorie zuweisen

$$c_{\text{map}} = \operatorname{argmax}_{c \in \text{yes, no}} P(c) \cdot \prod_i P(a_i|c) \\ = \text{no}$$

## Textklassifikation

Bei der Klassifikation von Dokumenten werden die Wörter des Textes als Attribute eingesetzt:

$$\begin{aligned}c_{map} &= \operatorname{argmax}_c P(c)P(c|d) \\ &= \operatorname{argmax}_c P(c) \prod_i P(w_i|c)\end{aligned}$$

$P(w|c)$  : Wahrscheinlichkeit, dass Wort  $w$  in einem  
: Dokument mit Klasse  $c$  vorkommt

## Anwendung: Spam herausfiltern

- Datensatz: Ein Korpus von E-Mail-Nachrichten, annotiert als "Spam" oder "Ham"
- Aufgabe: Eine neue Nachricht als "Spam" oder "Ham" klassifizieren
- Attribute: Vokabular der E-Mail-Nachrichten im Trainingskorpus

Bsp.: Eine E-Mail mit dem Inhalt "Get rich fast!!!"

$$\begin{aligned}c_{map} &= \operatorname{argmax}_{c \in \text{spam, ham}} P(c) \cdot \prod_i P(a_i|c) \\ &= \operatorname{argmax}_{c \in \text{spam, ham}} P(c) \cdot \\ &\quad P(\text{get}|c) \cdot P(\text{rich}|c) \cdot P(\text{fast}|c) \cdot P(!!!|c)\end{aligned}$$

## Anwendung: Spam Herausfiltern

- Problem: “rich” und “!!!” sind in Ham-E-Mails selten
- Null-Werte führen dazu, dass Kategorien ununterscheidbar werden

$$\begin{aligned}P(\text{ham}) &= P(c) \cdot P(\text{get}|c) \cdot P(\text{rich}|c) \cdot P(\text{fast}|c) \cdot P(\text{!!!}|c) \\ &= P(c) \cdot P(\text{get}|c) \cdot 0 \cdot P(\text{fast}|c) \cdot 0 = \boxed{0} !\end{aligned}$$

- Mögliche Lösung: Kleine Werte zu Zähler und Nenner hinzufügen, damit beide ungleich null werden

$$P(a|c) = \frac{n_a + \frac{1}{k}}{n + 1}$$

$n$  : Anzahl Trainingsinstanzen in Kategorie  $c$

$n_a$  : Anzahl Trainingsinstanzen mit Attribut  $a$  in Kategorie  $c$

$k$  : Anzahl Werte, die Attribut  $a$  annehmen kann

# Clusteranalyse

*6. Juli 2010*

## Clustering vs. Klassifikation

- Bei der Klassifikation werden Instanzen vordefinierten Klassen zugeordnet
- Beim Clustering entdeckt der Algorithmus “natürliche” Klassen, die die Instanzen in Gruppen mit ähnlichen Eigenschaften teilen
- Deutscher Begriff: *Ballungsanalyse*

## *Betreutes und unbetreutes Lernen*

Betreutes Lernen (Engl. “supervised learning”):

Der Algorithmus lernt, indem er Eingabe-Ausgabe-Beispiele analysiert, die die Rolle eines “Lehrers” übernehmen

Unbetreutes Lernen (Engl. “unsupervised learning”):

Vorkategorisierte Beispiele sind nicht verfügbar

### Halb-betreutes Lernen (Engl. "semi-supervised learning")

- Nur ein Teil der Trainingsbeispiele ist bereits mit einer Kategorie annotiert
- Kombination aus betreutes und unbetreutes Lernen

### Bestärkendes Lernen (Engl. "reinforcement learning")

- Die Lernkomponente ist Teil eines Systems oder Agents, das lernen soll, wie in potenziell auftretenden Situationen zu handeln ist
- Der Algorithmus erfährt aber nicht explizit, was die richtige Kategorisierung wäre
- Stattdessen lernt er durch Feedback von der Umwelt, das ihm in der Form von Belohnung und Bestrafung gegeben wird, den Erfolg des Agenten zu maximieren

## Wozu Clustering verwenden?

### Explorative Datenanalyse

Um ein Gefühl für die vorhandenen Daten und ihre Eigenschaften zu gewinnen

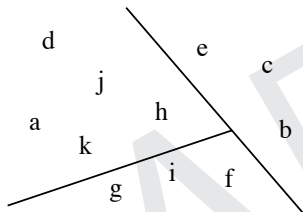
### Binning

Instanzen entdecken, die sich ähnlich verhalten und daher ähnlich behandelt werden können, um Abhilfe bei Sparse-Data-Problemen zu schaffen

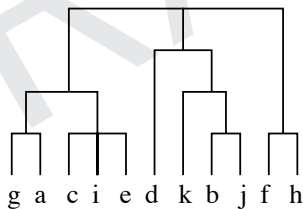
- In einem Korpus findet man die Sequenzen “*am Donnerstag*” und “*am Freitag*” sowie *donnertags* und *freitags*
- Außerdem hat man “*am Montag*”, aber *montags* kommt nicht vor
- Wenn wir wissen, dass *Donnerstag*, *Freitag* und *Montag* sich syntaktisch ähnlich verhalten, können wir *montags* inferieren

## Flach

Garance PARIS

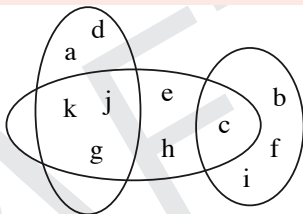


## Hierarchisch (Dendrogramme)



### Disjunktiv

Instanzen können mehreren Cluster angehören



### Probabilistisch oder "soft"

Für jeden Cluster wird eine Wahrscheinlichkeit angegeben, dass eine Instanz ihm zugeordnet wird

	1	2	3
a	0.4	0.1	0.5
b	0.1	0.8	0.1
c	0.3	0.3	0.4
d	0.1	0.1	0.8
e	0.4	0.2	0.4
f	0.1	0.4	0.5
g	0.7	0.2	0.1
h	0.5	0.4	0.1

## *Intuition bei k-means*

- Bestimme  $k$ , die Anzahl der gewünschten Cluster
- Wähle  $k$  beliebige Punkte als Cluster-Zentren aus
- Weise jede Instanz dem nächsten Cluster-Zentrum zu
- Berechne den Mittelpunkt für jeden Cluster und verwende ihn als neues Zentrum
- Weise alle Instanzen wieder dem nächsten Cluster-Zentrum zu
- Iteriere, bis alle Cluster stabil sind

## Der Algorithmus

- Jede Instanz  $\vec{x}$  im Training Set wird als Vektor mit einem Wert pro Attribut repräsentiert

$$\vec{x} = (x_1, x_2, \dots, x_n)$$

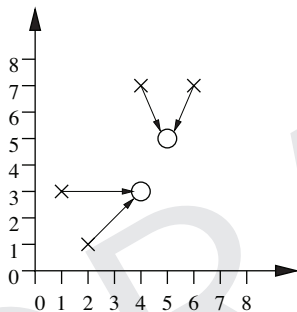
- Die Distanz zwischen zwei Vektors  $\vec{x}$  and  $\vec{y}$  ist definiert als (euklidische Distanz):

$$|\vec{x} - \vec{y}| = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

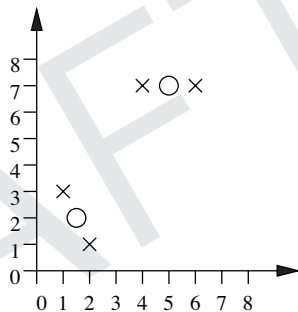
- Der Mittelpunkt  $\vec{\mu}$  einer Menge Vektoren  $c_j$  ist definiert als:

$$\vec{\mu} = \frac{1}{|c_j|} \sum_{\vec{x} \in c_j} \vec{x}$$

## Beispiel



Die Instanzen  
(Kreuzchen) werden  
anfangs zum nächsten  
Cluster-Zentrum  
(Kreise) zugewiesen



Der Mittelpunkt jedes  
Cluster wird dann  
berechnet und als neuen  
Zentrum verwendet

## *Eigenschaften von k-means*

- Flaches Clustering-Verfahren
- Konzeptuell einfach
- Effizient bei großen Datenmengen
- Nicht geeignet für Nominaldaten
- Findet nur ein lokales Maximum, keine globales
- Die Cluster hängen sehr von der initialen Wahl der Cluster-Zentren
- Kann man für hierarchisches Clustering verwenden
- Andere Distanzmaße können verwendet werden (z. B. der Cosinus)
- Alternative: der EM-Algorithmus
  - Legt iterativ die Parameter eines Modells so fest, dass es die gesehenen Daten optimal erklärt
  - HMMs sind eine Anwendung des EM-Algorithmus