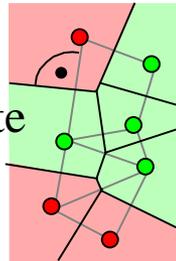
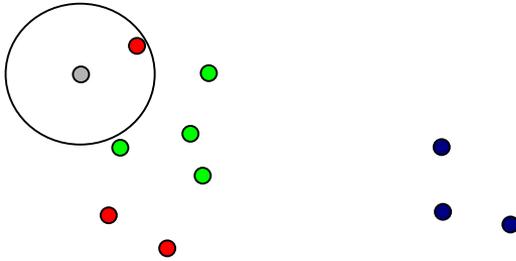


Instanzbasiertes Lernen (*instance based learning*)

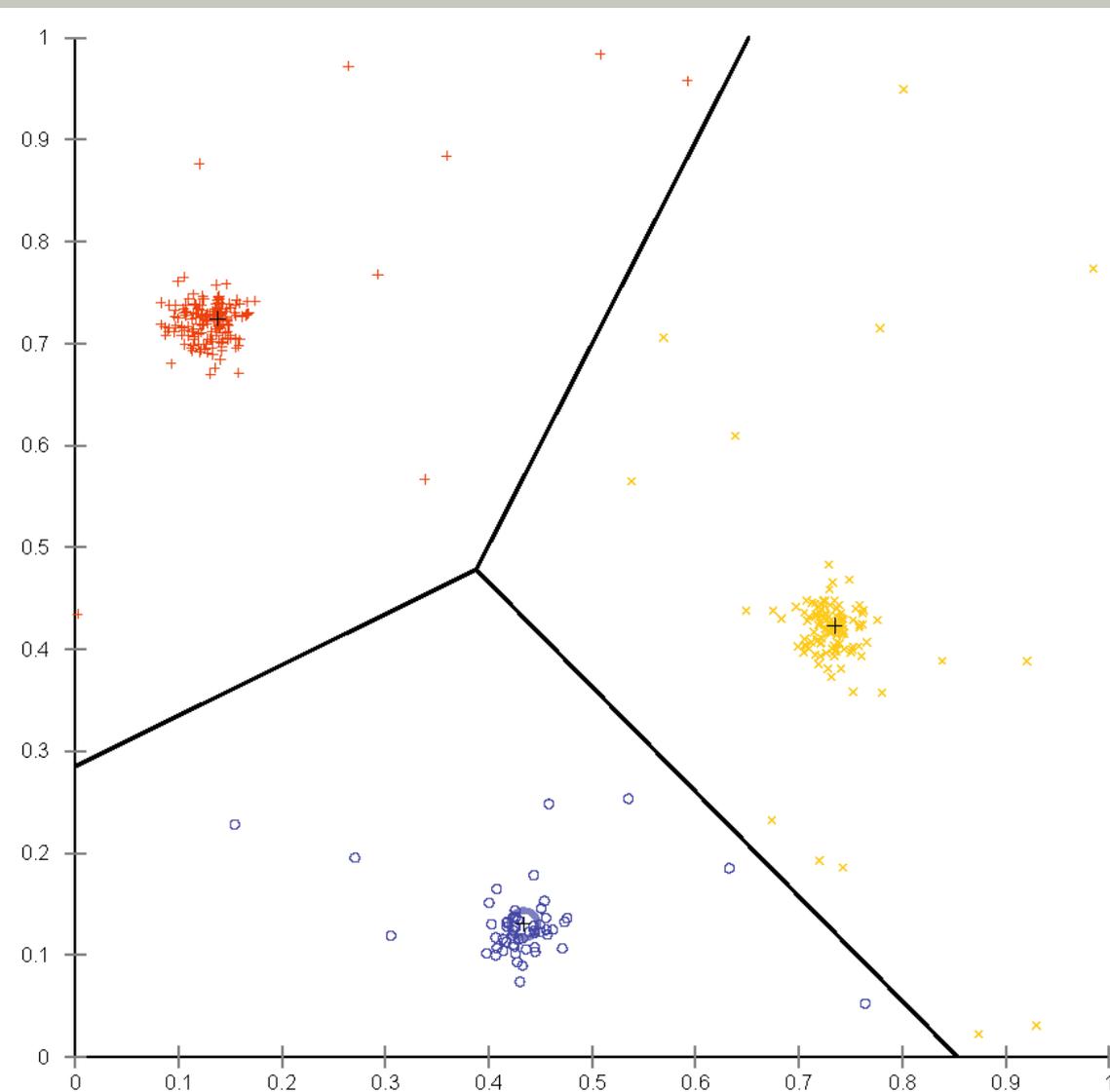
- Einfacher Nächster-Nachbar-Klassifikator:
Zuordnung zu der Klasse des nächsten Nachbarpunkts
- Im Beispiel: Nächster Nachbar ist eine Schraube
- Regionen der Klassenzuordnung können als *Voronoi-Diagramme* dargestellt werden:

Mittel-
senkrechte

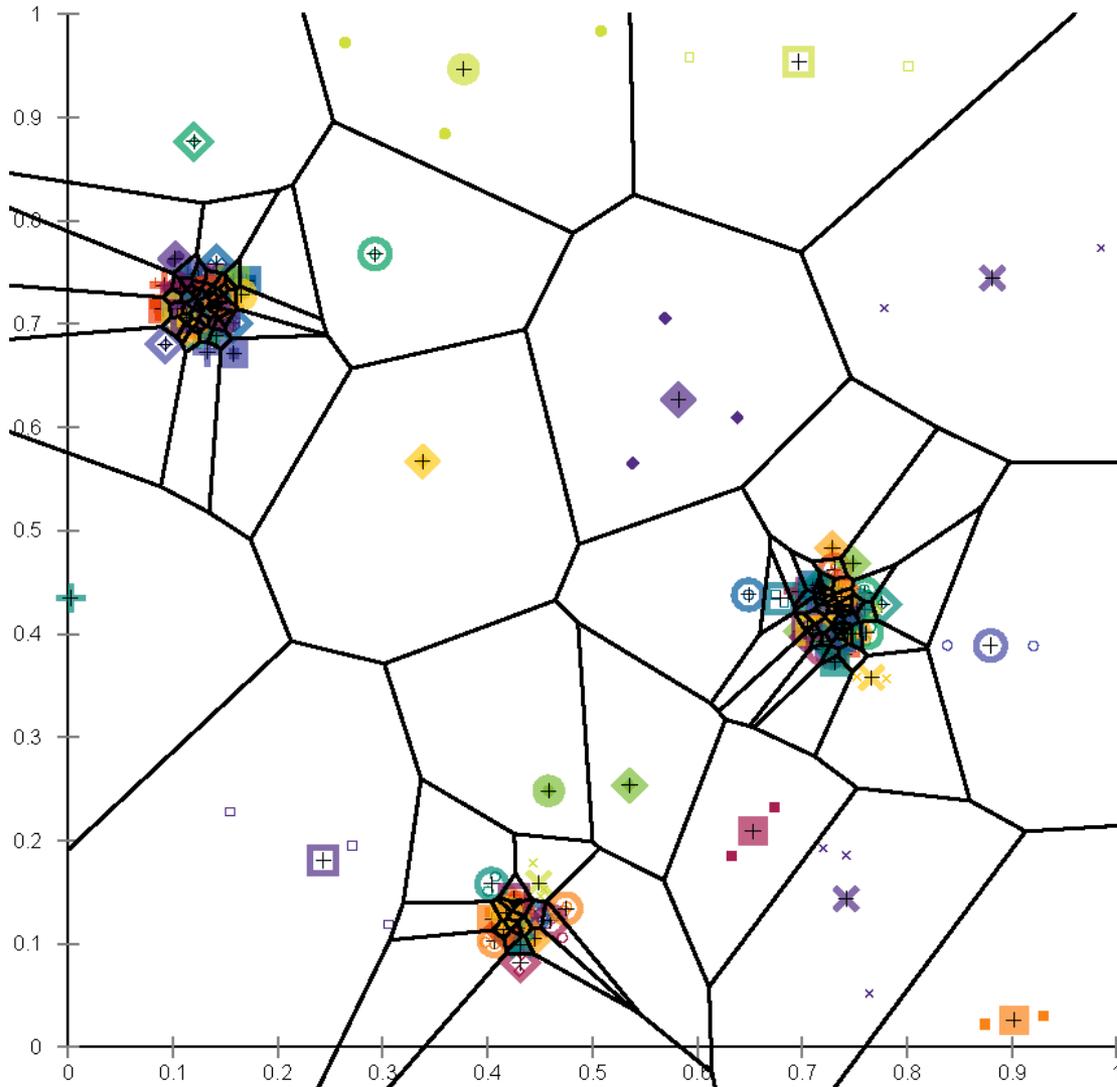




- Problem: Punkt links oben wahrscheinlich nur Ausreißer
- Besser: Betrachte mehr als nur einen Nachbarn
⇒ k -Nächste-Nachbarn-Klassifikator
- *Entscheidungsmenge*
die Menge der zur Klassifikation betrachteten k -nächsten Nachbarn
- *Entscheidungsregel*
Wie bestimmt man aus den Klassen der Entscheidungsmenge die Klasse des zu klassifizierenden Objekts?
 - Interpretiere Häufigkeit einer Klasse in der Entscheidungsmenge als Wahrscheinlichkeit der Klassenzugehörigkeit
 - Maximum-Likelihood-Prinzip: Mehrheitsentscheidung
 - Ggf. Gewichtung



- $k=1$: Voronoi-Zellen separieren „Trainingsobjekte“
- Zuordnung neuer Punkte zum nächstgelegenen „Trainings“-Objekt

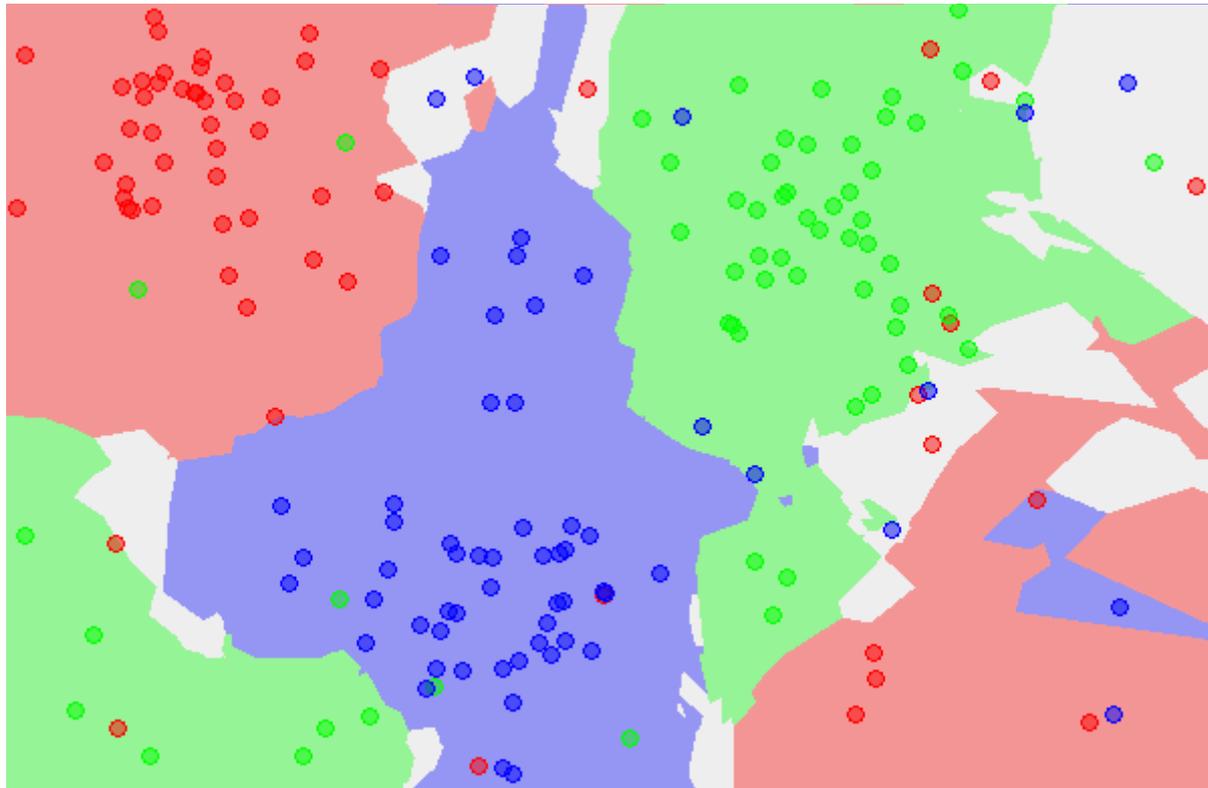


viele Trainingsobjekte



sehr kleinteilige, flexible
Klassengrenzen

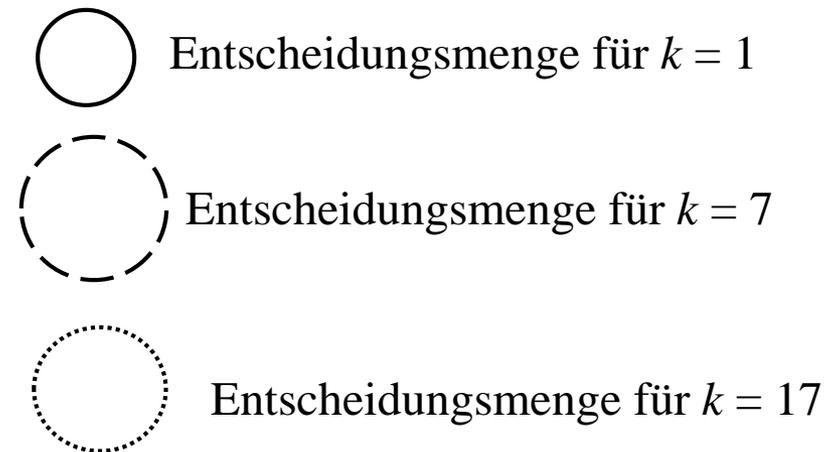
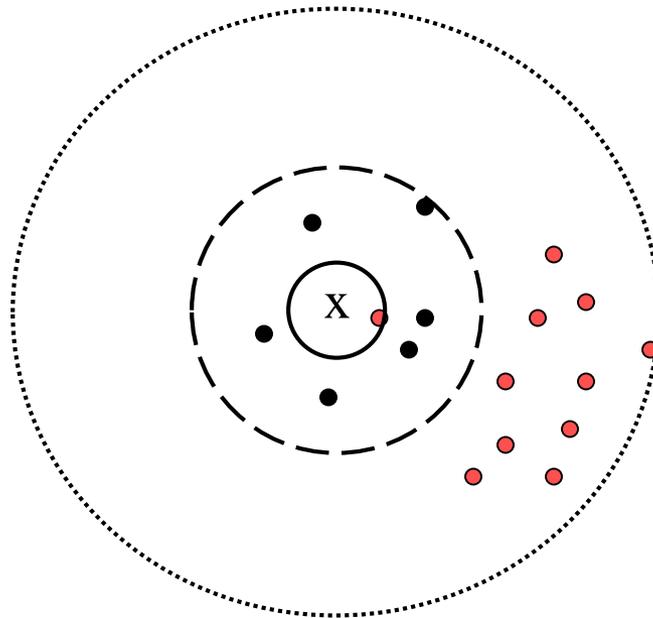
Nächste-Nachbarn-Klassifikatoren



- 1-NN
- 5-NN

Wahl des Parameters k

- „zu kleines“ k : hohe Sensitivität gegenüber Ausreißern
- „zu großes“ k : viele Objekte aus anderen Clustern (Klassen) in der Entscheidungsmenge.
- mittleres k : höchste Klassifikationsgüte, oft $1 \ll k < 10$



x: zu klassifizieren

- Standardregel

⇒ wähle die Mehrheitsklasse der Entscheidungsmenge

- Gewichtete Entscheidungsregel

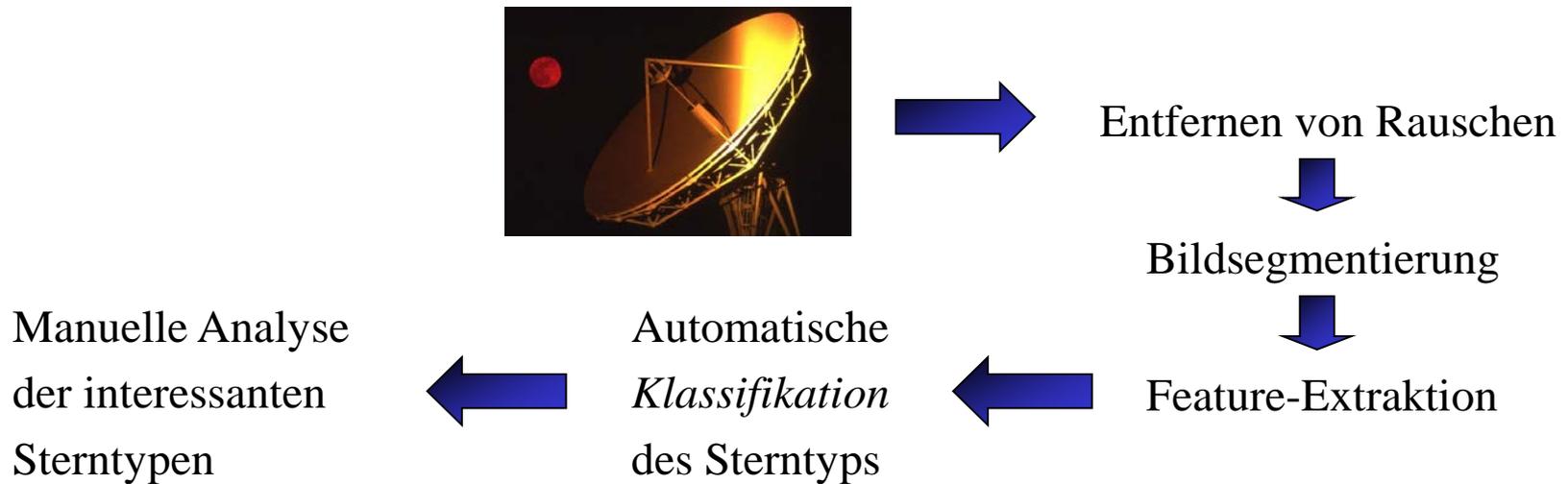
gewichte die Klassen der Entscheidungsmenge

- nach Distanz, meist invers quadriert: $weight(dist) = 1/dist^2$
- nach Verteilung der Klassen (oft sehr ungleich!)

Problem: Klasse mit zu wenig Instanzen ($< k/2$) in der Trainingsmenge bekommt keine Chance, ausgewählt zu werden, selbst bei optimaler Distanzfunktion

- Klasse A: 95 %, Klasse B 5 %
- Entscheidungsmenge = {A, A, A, A, B, B, B}
- Standardregel ⇒ A, gewichtete Regel ⇒ B

Analyse astronomischer Daten



Klassifikation des Sterntyps mit Nächste-Nachbarn Klassifikator
 basierend auf dem Hipparcos-Katalog

Hipparcos-Katalog [ESA 1998]

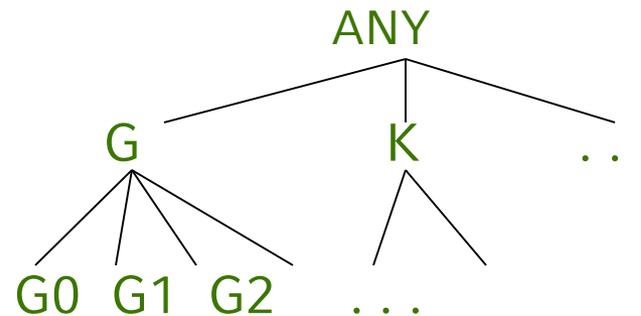
- enthält ca. 118 000 Sterne
- mit 78 Attributen (Helligkeit, Entfernung, Farbe, . . .)
- Klassenattribut: Spektraltyp (Attribut H76)

z.B.

H76: G0

H76: G7.2

H76: KIII/IV



- Werte des Spektraltyps sind vage
 ➔ Hierarchie von Klassen: benutze die erste Ebene der Klassenhierarchie

Verteilung der Klassen

Klasse	#Instanzen	Anteil Instanzen	
K	32 036	27.0	} häufige Klassen
F	25 607	21.7	
G	22 701	19.3	
A	18 704	15.8	
B	10 421	8.8	
M	4 862	4.1	
O	265	0.22	} seltene Klassen
C	165	0.14	
R	89	0.07	
W	75	0.06	
N	63	0.05	
S	25	0.02	
D	27	0.02	

Experimentelle Untersuchung [Poschenrieder 1998]

- Distanzfunktion
 - mit 6 Attributen (Farbe, Helligkeit und Entfernung)
 - mit 5 Attributen (ohne Entfernung)
 - ⇒ beste Klassifikationsgenauigkeit mit 6 Attributen
- Anzahl k der Nachbarn
 - ⇒ beste Klassifikationsgenauigkeit für $k = 15$
- Entscheidungsregel
 - Gewichtung nach Distanz
 - Gewichtung nach Klassenverteilung
 - ⇒ beste Klassifikationsgenauigkeit bei Gewichtung nach Distanz aber nicht nach Klassenverteilung

Klassifikation von Sternen

Klasse	Falsch klassifiziert	Korrekt klassifiziert	Klassifikationsgenauigkeit
K	408	2338	85.1%
F	350	2110	85.8%
G	784	1405	64.2%
A	312	975	75.8%
B	308	241	43.9%
M	88	349	79.9%
C	4	5	55.6%
R	5	0	0%
W	4	0	0%
O	9	0	0%
N	4	1	20%
D	3	0	0%
S	1	0	0%
Total	2461	7529	75.3%



hohe Klassifikationsgenauigkeit für die häufigen Klassen, schlechte Genauigkeit für die seltenen Klassen

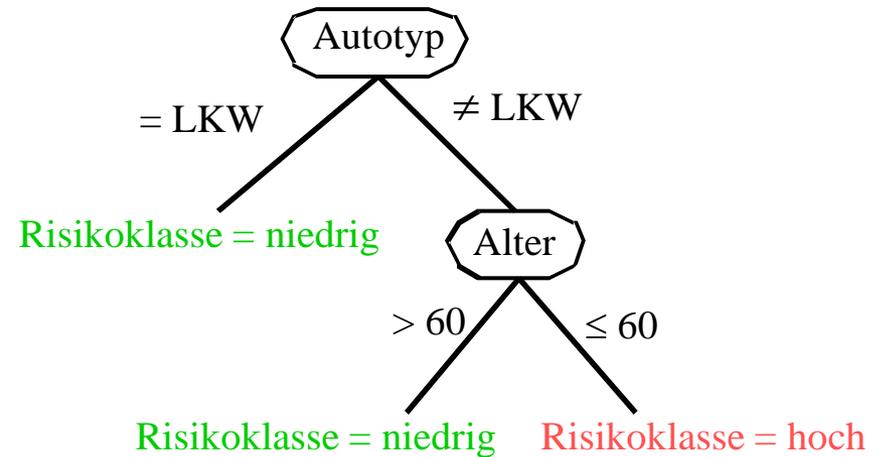
die meisten seltenen Klassen besitzen weniger als $k / 2 = 8$ Instanzen!

Diskussion

- + Anwendbarkeit: erfordert nur die Trainingsdaten
- + hohe Klassifikationsgenauigkeit in vielen Anwendungen
- + inkrementell: Klassifikator kann sehr einfach an neue Trainingsobjekte adaptiert werden
- + keine Trainingsphase erforderlich („lazy learner“)
- Ineffizienz: Auswertung des “Modells” erfordert k -nächste-Nachbarn Anfrage an die Datenbank
- liefert kein explizites Wissen über die Klassen
- Wahl von k problematisch bei sehr unterschiedlichen Klassengrößen

Motivation

ID	Alter	Autotyp	Risiko
1	23	Familie	hoch
2	17	Sport	hoch
3	43	Sport	hoch
4	68	Familie	niedrig
5	32	LKW	niedrig

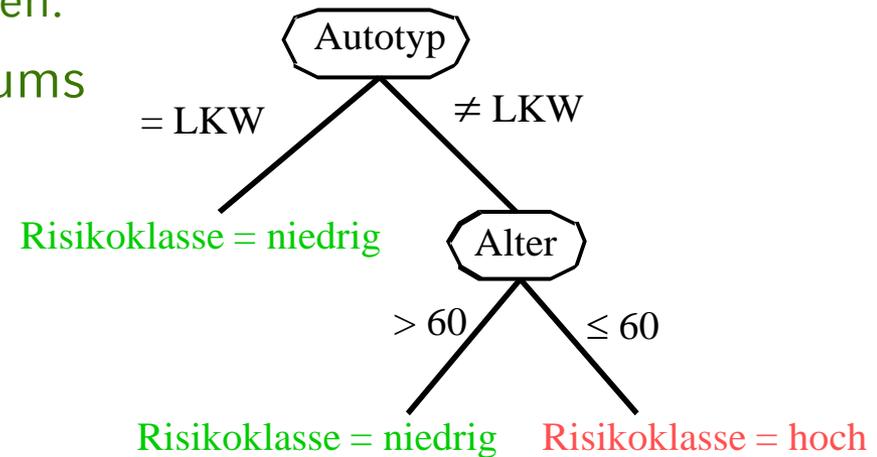


- finden explizites Wissen
- Entscheidungsbäume liefern ein für den Benutzer verständliches Modell (Hierarchie von Regeln)

- Ein *Entscheidungsbaum* ist ein Baum mit folgenden Eigenschaften:
 - ein innerer Knoten repräsentiert ein Attribut,
 - eine Kante repräsentiert einen Test auf dem Attribut des Vaterknotens,
 - ein Blatt repräsentiert eine der Klassen.

- Konstruktion eines Entscheidungsbaums

- anhand der Trainingsmenge
- Top-Down



- Anwendung eines Entscheidungsbaums

Durchlauf des Entscheidungsbaum von der Wurzel zu einem der Blätter

⇒ eindeutiger Pfad

Zuordnung des Objekts zur Klasse des erreichten Blatts

Basis-Algorithmus

- Anfangs gehören alle Trainingsdatensätze zur Wurzel.
- Das nächste Attribut wird ausgewählt (Splitstrategie).
- Die Trainingsdatensätze werden unter Nutzung des Splitattributs partitioniert.
- Das Verfahren wird rekursiv für die Partitionen fortgesetzt.
⇒ lokal optimierender Algorithmus

Abbruchbedingungen

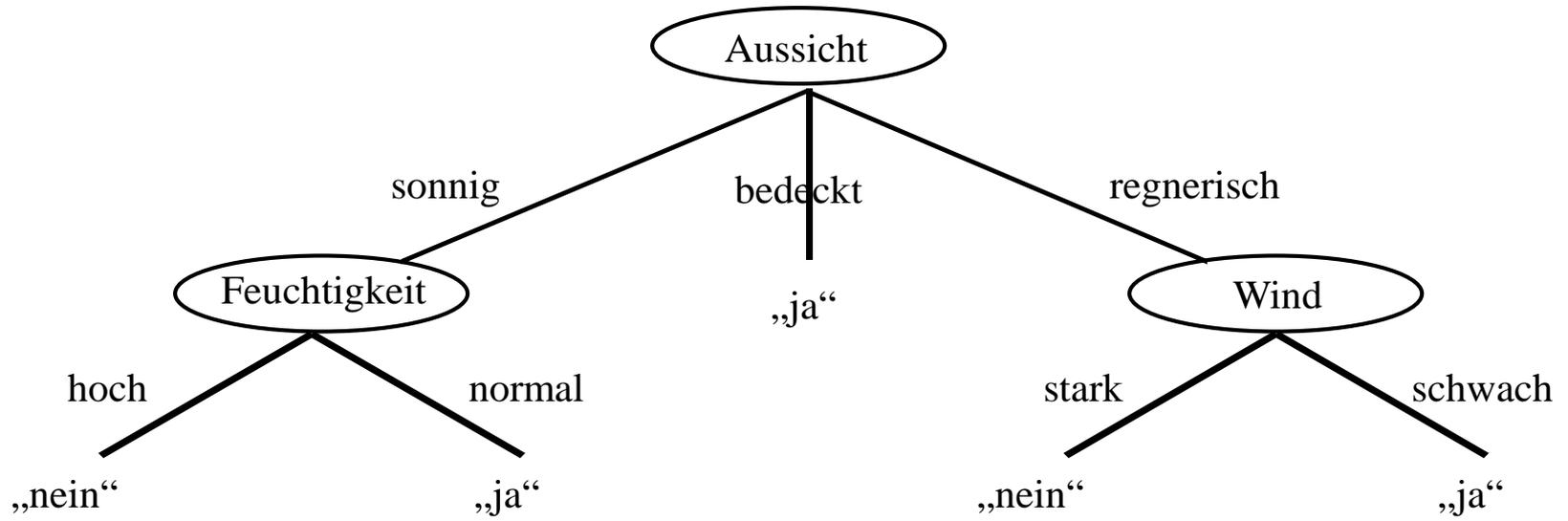
- keine weiteren Splitattribute
- alle Trainingsdatensätze eines Knotens gehören zur selben Klasse

Beispiel

Tag	Aussicht	Temperatur	Feuchtigkeit	Wind	Tennispielen
1	sonnig	heiß	hoch	schwach	nein
2	sonnig	heiß	hoch	stark	nein
3	bedeckt	heiß	hoch	schwach	ja
4	regnerisch	mild	hoch	schwach	ja
5	regnerisch	kühl	normal	schwach	ja
6	regnerisch	kühl	normal	stark	nein
7	bedeckt	kühl	normal	stark	ja
8	sonnig	mild	hoch	schwach	nein
9	sonnig	kühl	normal	schwach	ja
10	regnerisch	mild	normal	schwach	ja
11	sonnig	mild	normal	stark	ja
12	bedeckt	mild	hoch	stark	ja
13	bedeckt	heiß	normal	schwach	ja
14	regnerisch	mild	hoch	stark	nein

Ist heute ein Tag zum Tennispielen?

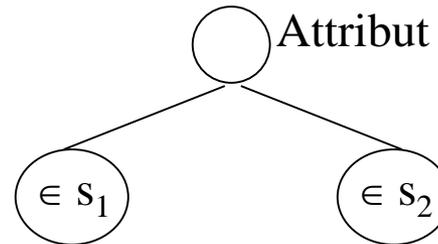
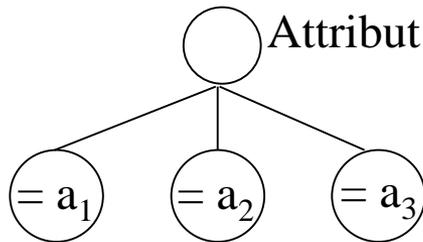
Beispiel



Typen von Splits

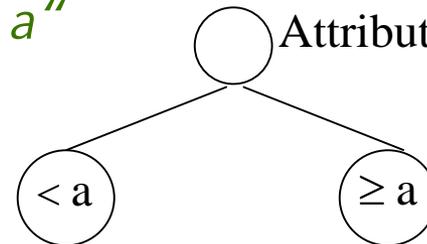
- **Kategorische Attribute**

Splitbedingungen der Form „Attribut = a“ or „Attribut ∈ set“
viele mögliche Teilmengen



- **Numerische Attribute**

Splitbedingungen der Form „Attribut < a“
viele mögliche Splitpunkte



Wo sollen diskrete Attribute gesplittet werden?

=> An den Stellen, die die Qualität maximieren.

Idee: Ordnen der numerischen Attributwerte

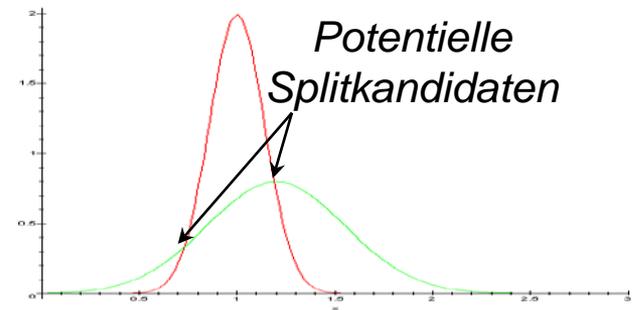
Wert	0.9	0.8	0.65	0.5	0.45	0.3	0.15	0.01
Klasse	A	A	B	B	B	A	A	A

Potentielle Splitkandidaten

Teste die Kombination, die den höchsten Information Gain erzielen.

Schnellere Methode:

- Bilde Gauß-Kurve über alle Klassen
- Wähle Schnittpunkte der Gauß-Kurven als Kandidaten.



Qualitätsmaße für Splits

Gegeben

- eine Menge T von Trainingsobjekten
- eine disjunkte, vollständige Partitionierung T_1, T_2, \dots, T_m von T
- p_i die relative Häufigkeit der Klasse c_i in T

Gesucht

- ein Maß der *Unreinheit* einer Menge S von Trainingsobjekten in Bezug auf die Klassenzugehörigkeit
- ein Split von T in T_1, T_2, \dots, T_m , der dieses Maß der Unreinheit *minimiert*
 - ⇒ Informationsgewinn, Gini-Index

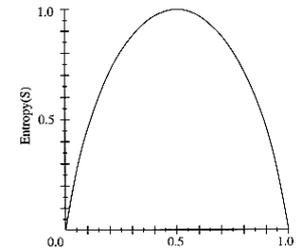
Informationsgewinn

Entropie: minimale Anzahl von Bits zum Codieren der Nachricht, mit der man die Klasse eines zufälligen Trainingsobjekts mitteilen möchte. Die *Entropie* für eine Menge T von Trainingsobjekten ist definiert als

$$\text{entropie}(T) = -\sum_{i=1}^k p_i \cdot \log p_i$$

$$\text{entropie}(T) = 0, \text{ falls } p_i = 1 \text{ für ein } i$$

$$\text{entropie}(T) = 1 \text{ für } k = 2 \text{ Klassen mit } p_i = 1/2$$



Das Attribut A habe die Partitionierung T_1, T_2, \dots, T_m erzeugt.

Der *Informationsgewinn* des Attributs A in Bezug auf T ist definiert als

$$\text{Informationsgewinn}(T, A) = \text{entropie}(T) - \sum_{i=1}^m \frac{|T_i|}{|T|} \cdot \text{entropie}(T_i)$$

Beispiel

9 „ja“ 5 „nein“ Entropie = 0,940



hoch normal

3 „ja“ 4 „nein“

Entropie = 0,985

6 „ja“ 1 „nein“

Entropie = 0,592

9 „ja“ 5 „nein“ Entropie = 0,940



schwach stark

6 „ja“ 2 „nein“

Entropie = 0,811

3 „ja“ 3 „nein“

Entropie = 1,0

$$\text{Informationsgewinn}(T, \text{Feuchtigkeit}) = 0,94 - \frac{7}{14} \cdot 0,985 - \frac{7}{14} \cdot 0,592 = 0,151$$

$$\text{Informationsgewinn}(T, \text{Wind}) = 0,94 - \frac{8}{14} \cdot 0,811 - \frac{6}{14} \cdot 1,0 = 0,048$$

⇒ Feuchtigkeit liefert den höheren Informationsgewinn

Gini-Index

Gini-Index für eine Menge T von Trainingsobjekten

$$gini(T) = 1 - \sum_{j=1}^k p_j^2$$

- kleiner Gini-Index \Leftrightarrow geringe Unreinheit,
- großer Gini-Index \Leftrightarrow hohe Unreinheit

Das Attribut A habe die Partitionierung T_1, T_2, \dots, T_m erzeugt.

Gini-Index des Attributs A in Bezug auf T ist definiert als

$$gini_A(T) = \sum_{i=1}^m \frac{|T_i|}{|T|} \cdot gini(T_i)$$

Beispiel

9 „ja“ 5 „nein“, |T|=14



hoch normal

3 „ja“ 4 „nein“

6 „ja“ 1 „nein“

$$gini(T) = 1 - \left(\frac{3^2}{7^2} + \frac{4^2}{7^2} \right) = \frac{24}{49}$$

$$gini(T) = 1 - \left(\frac{6^2}{7^2} + \frac{1^2}{7^2} \right) = \frac{12}{49}$$

$$gini_{Feuchtigkeits}(T) = \frac{7}{14} \cdot \frac{24}{49} + \frac{7}{14} \cdot \frac{12}{49} = \frac{18}{49}$$

9 „ja“ 5 „nein“, |T|=14



schwach stark

6 „ja“ 2 „nein“

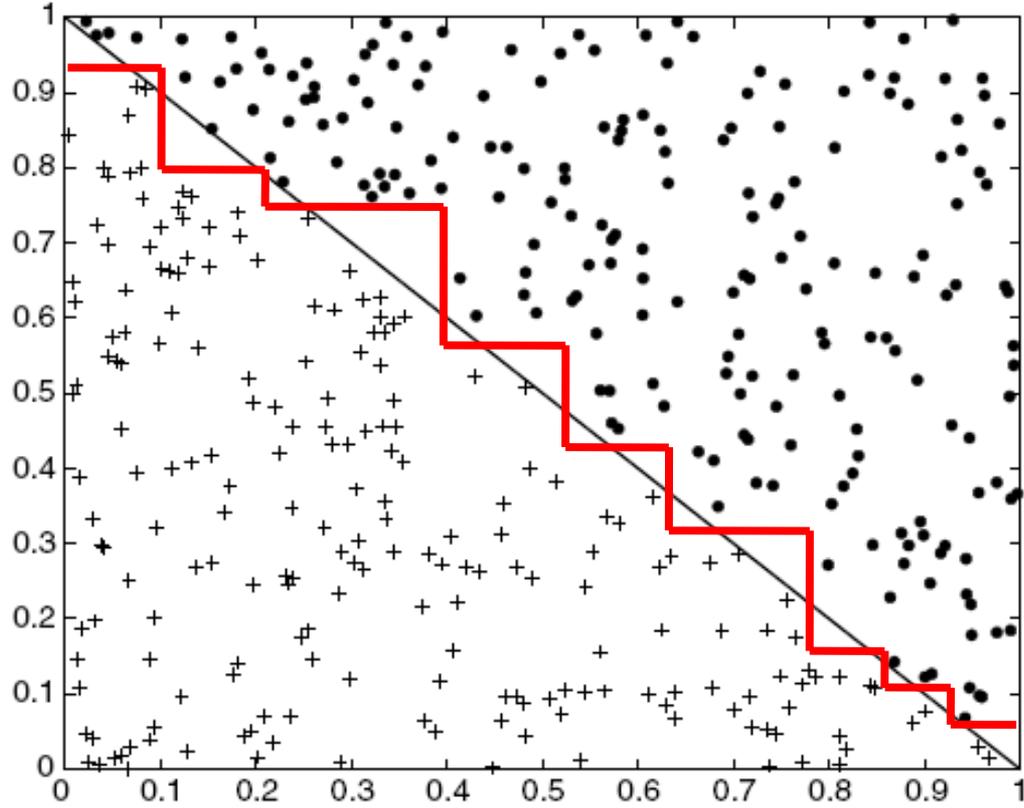
3 „ja“ 3 „nein“

$$gini(T) = 1 - \left(\frac{6^2}{8^2} + \frac{2^2}{8^2} \right) = \frac{24}{64} = \frac{3}{8}$$

$$gini(T) = 1 - \left(\frac{3^2}{6^2} + \frac{3^2}{6^2} \right) = \frac{18}{36} = \frac{1}{2}$$

$$gini_{Wind}(T) = \frac{8}{14} \cdot \frac{3}{8} + \frac{6}{14} \cdot \frac{1}{2} = \frac{3}{7} = \frac{21}{49}$$

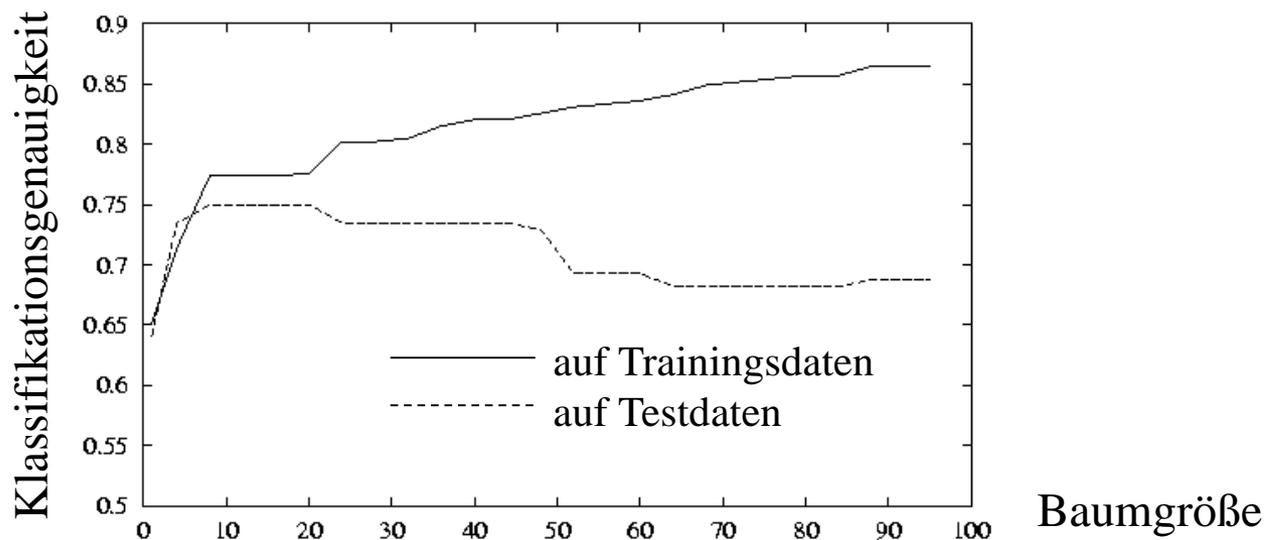
⇒ Feuchtigkeits liefert kleineren Gini-Index (geringere Unreinheit)



Einführung

Overfitting bei der Konstruktion eines Entscheidungsbaums, wenn es zwei Entscheidungsbäume E und E' gibt mit

- E hat auf der Trainingsmenge eine kleinere Fehlerrate als E' ,
- E' hat auf der Grundgesamtheit der Daten eine kleinere Fehlerrate als E .



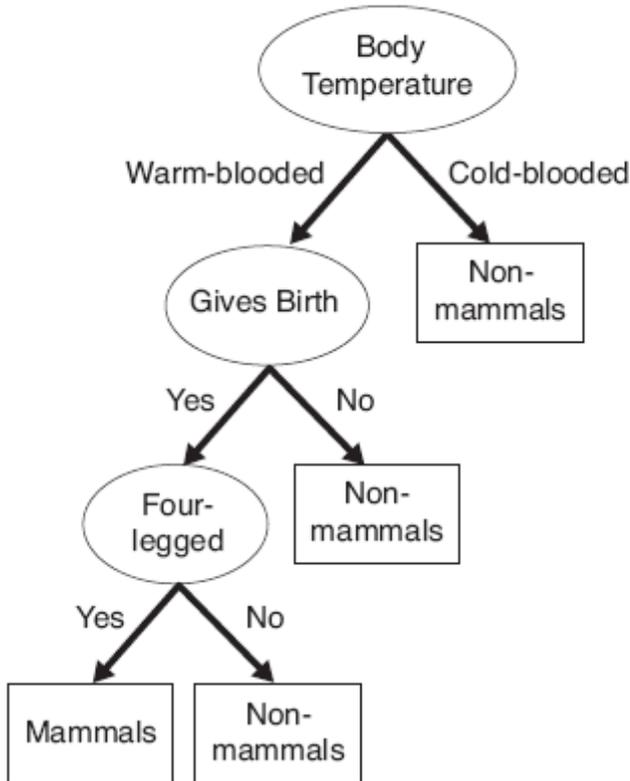
Trainingsdaten

Name	Body Temp	Gives Birth	Four-legged	Hibernates	Mammal
porcupine	warm-blooded	yes	yes	yes	yes
cat	warm-blooded	yes	yes	no	yes
bat	warm-blooded	yes	no	yes	no*
whale	warm-blooded	yes	no	no	no*
salamander	cold-blooded	no	yes	yes	no
komodo dragon	cold-blooded	no	yes	no	no
python	cold-blooded	no	no	yes	no
salmon	cold-blooded	no	no	no	no
eagle	warm-blooded	no	no	no	no
guppy	cold-blooded	yes	no	no	no

*Fehlerhafte Daten

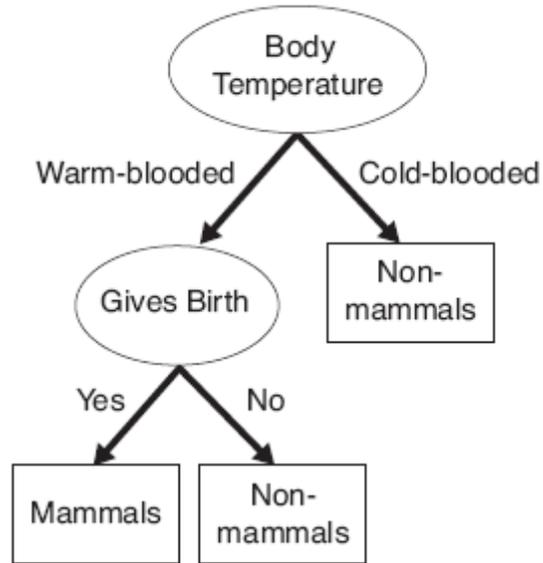
Testdaten

Name	Body Temp	Gives Birth	Four-legged	Hibernates	Mammal
human	warm-blooded	yes	no	no	yes
pigeon	warm-blooded	no	no	no	no
elephant	warm-blooded	yes	yes	no	yes
leopard shark	cold-blooded	yes	no	no	no
turtle	cold-blooded	no	yes	no	no
penguin	warm-blooded	no	no	no	no
eel	cold-blooded	no	no	no	no
dolphin	warm-blooded	yes	no	no	yes
spiny anteater	warm-blooded	no	yes	yes	yes
gila monster	cold-blooded	no	yes	yes	no



(a) Model M1

Trainingerror: 0%
Testerror: 30%



(b) Model M2

Trainingerror: 20%
Testerror: 10%

M1:

- Fehlklassifikation von human und dolphin aufgrund von fehlerhaften Trainingsdaten
- spiny anteater: ungewöhnlicher Fall

M2:

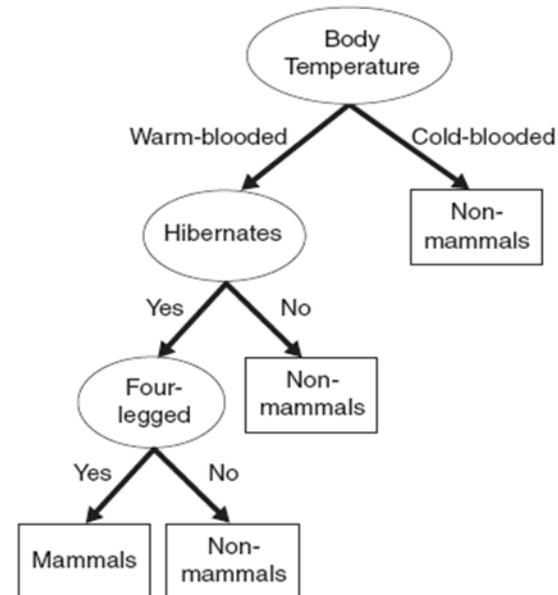
ungewöhnlicher Fall kann nicht vermieden werden

Name	Body Temp	Gives Birth	Four-legged	Hibernates	Mammal
salamander	cold-blooded	no	yes	yes	no
poorwill	warm-blooded	no	no	yes	no
platypus	warm-blooded	no	yes	yes	yes
eagle	warm-blooded	no	no	no	no
guppy	cold-blooded	yes	no	no	no

Problem:

Mangel an repräsentativen Daten

⇒ Fehlerrate auf Testset: 30%



Ansätze zum Vermeiden von Overfitting

- Entfernen von fehlerhaften Trainingsdaten
insbesondere widersprüchliche Trainingsdaten
- Wahl einer geeigneten Größe der Trainingsmenge
nicht zu klein, nicht zu groß
- Wahl einer geeigneten Größe des minimum support
minimum support:

Anzahl der Datensätze, die mindestens zu einem Blattknoten
des Baums gehören müssen



minimum support $\gg 1$

Ansätze zum Vermeiden von Overfitting

- Wahl einer geeigneten Größe der minimum confidence
minimum confidence: Anteil, den die Mehrheitsklasse eines
 Blattknotens mindestens besitzen muß.
 *minimum confidence* \ll 100%
 Blätter können auch fehlerhafte Datensätze oder Rauschen
 „absorbieren“
- nachträgliches Pruning des Entscheidungsbaums
 Abschneiden der überspezialisierten Äste
 kleinere Bäume sind weniger anfällig für Overfitting (Occam's Razor)

Fehlerreduktions-Pruning [Mitchell 1997]

- Aufteilung der klassifizierten Daten in Trainingsmenge und Testmenge
- Konstruktion eines Entscheidungsbaums E für die Trainingsmenge
- Pruning von E mit Hilfe der Testmenge T
 - bestimme denjenigen Teilbaum von E , dessen Abschneiden den Klassifikationsfehler auf T am stärksten reduziert
 - entferne diesen Teilbaum
 - fertig, falls kein solcher Teilbaum mehr existiert



nur anwendbar, wenn genügend viele klassifizierte Daten

Diskussion

- + Interpretation des gefundenen Baumes relativ einfach
- + Implizite Gewichtung der Attribute
- + Leistungsfähiger Klassifikator, häufig in der Praxis verwendet
- + Effiziente Auswertung des gefundenen Modells
- Finden eines optimalen Entscheidungsbaums ist exponentiell
- Heuristische Methoden können nur lokales Optimum finden
- Anfällig für Overfitting (besondere Methoden zur Vermeidung von Overfitting für Entscheidungsbäume)