

**Knowledge Discovery in Databases**  
 WS 2017/18

**Übungsblatt 10: Classification**

Besprechung: 25. und 26.01.2018

**Aufgabe 10-1 Bewertung von Klassifikatoren**

Gegeben sei ein Datensatz mit bekannter Klassenzugehörigkeit der Objekte. Um die Qualität eines Klassifikators  $K$  zu ermitteln wurden die Objekte mittels  $K$  klassifiziert. Die Klassifikationsergebnisse sind in der folgenden Tabelle dargestellt.

ID	Objektklasse	$K(o)$	ID	Objektklasse	$K(o)$
$O_1$	A	A	$O_2$	B	A
$O_3$	A	C	$O_4$	C	C
$O_5$	C	B	$O_6$	B	B
$O_7$	A	A	$O_8$	A	A
$O_9$	A	A	$O_{10}$	B	C
$O_{11}$	B	A	$O_{12}$	C	A
$O_{13}$	C	C	$O_{14}$	C	C
$O_{15}$	B	B			

- Berechnen Sie anhand der tabellierten Ergebnisse Precision und Recall jeder Klasse.
- Um ein vollständiges Maß für die Güte der Klassifikation bezüglich einer Klasse zu haben, wird häufig auch das sogenannte  $F_1$ -Measure (harmonisches Mittel zwischen Precision und Recall) verwendet. Das  $F_1$ -Measure für Klasse  $i$  ist wie folgt definiert:

$$F_1(K, i) = \frac{2 \cdot \text{Recall}(K, i) \cdot \text{Precision}(K, i)}{\text{Recall}(K, i) + \text{Precision}(K, i)}$$

Berechnen Sie das  $F_1$ -Measure für alle Klassen.

- Berechnen Sie die durchschnittliche Precision, den durchschnittlichen Recall und daraus das  $F_1$ -Measure.

**Aufgabe 10-2 Naive Bayes**

Die Ski-Saison ist eröffnet. Um zuverlässig zu entscheiden, wann Sie Skifahren gehen können und wann nicht, können Sie einen Klassifikator (z.B. Naive Bayes) benutzen. Der Klassifikator wird mit Ihren Erfahrungswerten aus dem letzten Jahr trainiert. Berücksichtigt werden dabei folgende Attribute:

Das Wetter: Das Attribut `Wetter` kann die folgenden drei Werte annehmen: Sonne, Regen und Schnee.

Die Schneehöhe: Das Attribut `Schneehöhe` kann die folgenden zwei Werte annehmen:  $\geq 50$  (Es liegen mindestens 50 cm Schnee) und  $< 50$  (Es liegen weniger als 50 cm Schnee).

Angenommen, Sie wollten letztes Jahr 8-mal zum Skifahren gehen. Die folgende Tabelle gibt Ihre jeweiligen Entscheidungen wieder:

Wetter	Schneehöhe	Skifahren ?
Sonne	< 50	nein
Regen	< 50	nein
Regen	≥ 50	nein
Schnee	≥ 50	ja
Schnee	< 50	nein
Sonne	≥ 50	ja
Schnee	≥ 50	ja
Regen	< 50	ja

- (a) Berechnen Sie die *a priori* Wahrscheinlichkeiten für die beiden Klassen Skifahren = ja und Skifahren = nein (auf den Trainingsdaten)!
- (b) Berechnen Sie für die Klassen die Werteverteilungen aller Attribute.
- (c) Entscheiden Sie, ob Sie bei den folgenden Wetter- und Schneebedingungen Skifahren gehen oder nicht! Verwenden Sie dazu den naiven Bayes-Klassifikator.

	Wetter	Schneehöhe
Tag A	Sonne	≥ 50
Tag B	Regen	< 50
Tag C	Schnee	< 50

### Aufgabe 10-3 Entscheidungs bäume

Sie wollen die Risikoklasse eines Autofahrers anhand der folgenden Merkmale vorhersagen:

- Zeit seit Bestehen der Fahrprüfung(1-2 Jahre, 2-7 Jahre, >7 Jahre)
- Geschlecht (männlich, weiblich)
- Wohnort(Stadt, Land)

Für Ihre Analyse stehen Ihnen folgende manuell eingeteilte Testbeispiele zu Verfügung:

Person	Zeit seit der Fahrprüfung	Geschlecht	Wohnort	Risikoklasse
1	1-2	m	Stadt	niedrig
2	2-7	m	Land	hoch
3	>7	w	Land	niedrig
4	1-2	w	Land	hoch
5	>7	m	Land	hoch
6	1-2	m	Land	hoch
7	2-7	w	Stadt	niedrig
8	2-7	m	Stadt	niedrig

- (a) Konstruieren Sie anhand dieser Trainingsdaten einen Entscheidungsbaum. Benutzen Sie beim Split den Informationsgewinn als Maß für die Unreinheit. Erzeugen Sie dabei für jeden Attributwert einen eigenen Ast. Der Entscheidungsbaum soll terminieren, wenn alle Instanzen im Blatt die gleiche Klasse haben. Die Anwendung eines Pruning-Algorithmus ist nicht erforderlich!
- (b) Wenden Sie Ihren Entscheidungsbaum auf folgende Autofahrer an:  
 Person A: 1-2, w, Land  
 Person B: 2-7, m, Stadt  
 Person C: 1-2, w, Stadt