

# Maschinelles Lernen und Data Mining

Projekt

WS 05/06

Auf diesem Blatt finden Sie die Problemstellungen für den Angewandten Teil der Übungen zur Vorlesung Maschinelles Lernen und Data Mining. Das Ziel dieser Aufgaben ist es, Ihnen zu ermöglichen praktische Erfahrungen mit einem Data Mining Werkzeug zu sammeln, und dabei Kenntnisse zu erwerben, die letztendlich auch für die Klausur relevant sein werden. Beachten Sie, daß es oft keine eindeutigen Lösungen der Aufgaben gibt. Was zählt ist die praktische Erprobung der Verfahren, die Antworten können sich abhängig von den eingesetzten Verfahren bzw. abhängig von den untersuchten Datensätzen unterscheiden.

Es wird erwartet, daß Sie die Aufgaben selbständig bzw. in kleinen Gruppen lösen, und die gefundenen Lösungen abgeben. Die Lösungen werden in der letzten Übungsstunde (16. 2. 2006) kurz durchbesprochen. Eventuell auftretende Probleme können Sie in den laufenden Übungsstunden oder im Vorlesungs-Forum ansprechen.

Sie finden die Homepage der Weka Machine Learning Library unter <http://www.cs.waikato.ac.nz/ml/weka/>. Dort können Sie sich die Java-Software, einige Kurzanleitungen, sowie Beispieldatenbanken herunterladen. Wir empfehlen, daß Sie die Experimente mit dem *Explorer* GUI (`weka.gui.explorer.Explorer`) durchführen, Sie können sie aber auch von der Command-line oder einem anderen GUI (`weka.gui.GUIChooser`) durchführen.

Dokumentation zu Weka finden Sie auf der Weka-Homepage. Beachten Sie insbesondere auch die Tips & Tricks ([http://www.cs.waikato.ac.nz/~ml/weka/tips\\_and\\_tricks.html](http://www.cs.waikato.ac.nz/~ml/weka/tips_and_tricks.html)). Dort wird z.B. erklärt, wie Sie den Speicher der Java VM vergrößern können, falls Probleme auftreten.

Als Lösung einer Aufgabe wird erwartet, daß Sie die wesentlichen Resultate (und die Schritte zu Ihrer Lösung, d.h. z.B. die aufgerufenen Routinen, die verwendeten Parameter, etc.) schriftlich zusammenfassen bzw. ggf. Grafiken zur Illustration verwenden. Ein Aneinanderhängen der Outputs von Weka ist nicht zielführend. Wir werden uns alle erhaltenen Lösungen ansehen, eine direkte Benotung der Abgaben findet jedoch nicht statt.

Abgabe der Aufgaben bis **7. 2. 2006**.

Viel Spaß bei der Durchführung dieser Aufgaben!

## 1 Aufgabe: Auswahl der Daten

Mit der Installation von Weka kommen einige Beispieldatensätze, ein jar-File mit vielen weiteren finden Sie auf der Homepage (<http://prdownloads.sourceforge.net/weka/datasets-UCI.jar>). Eine kurze Beschreibung der Daten befindet sich oft als Kommentar in den ersten Zeilen des arff-Files.

## 2 Aufgabe: Evaluierung

Wählen Sie fünf Datensets aus und testen Sie den Entscheidungsbaumlerner J48, eine Implementation von C4.5, an Ihren Datensets.

- Benutzen Sie die gesamten Daten als Trainingsmenge und bestimmen Sie die Genauigkeit auf dem Trainings-Set
- Bestimmen Sie die Genauigkeit mittels 2-fold, 5-fold, 10-fold und 20-fold-Crossvalidation.
- Iterieren Sie die Crossvalidation, z.B. mittels 10-facher 10-fold-Crossvalidation.
- Bestimmen Sie den Fehler mittels Leave-one-out cross-validation

Diskutieren Sie die unterschiedlichen Genauigkeitsabschätzungen.

In der Folge wird, so nicht anders angegeben, mit "Genauigkeit" immer die mit 10-facher Cross-validation geschätzte Genauigkeit betrachtet.

## 3 Aufgabe: Noise und Pruning

Verwenden Sie das Datenset, das in der vorigen Aufgabe die größte Genauigkeit auf dem Trainings-Set erzielt hat.

Stören Sie die Klasseninformation in diesem Datenset durch Hinzufügen von verschiedenen Levels von Noise (z.B., 5%, 10%, 25%, 50%, 75%, 100%; `weka.filters.unsupervised.attribute.AddNoise`).

Beobachten sie die Genauigkeit und Größe der gelernten Bäume auf dem Original bzw. den gestörten Datensets für J48

- mit den Default-Parametern
- ohne Pruning (-U und -M 1)

Experimentieren Sie ein wenig mit den Parametern -C und -M und versuchen Sie, die Kombination zu finden, die die höchste Genauigkeit liefert (auf den mit 10% Noise gestörten Daten).

**Anmerkung:** Ein Noise-Level von  $x\%$  wird erzeugt, indem bei  $x\%$  aller Beispiele das Label des Beispiels durch ein zufällig ausgewähltes Label eines der anderen Klassen ersetzt wird. Bei Zwei-Klassen-Problemen werden Sie feststellen, daß die Performanz bei 100% Noise identisch ist mit der Performanz bei 0% Noise (Warum?). Adaptieren Sie in diesem Fall die Schranken in geeigneter Weise (hier entspricht 50% Noise zufälligen Daten).

## 4 Aufgabe: Regel-Lernen

In Weka finden Sie u.a. JRip, eine Nachimplementierung des bekanntesten Regellerners Ripper, und `ConjunctiveRule`, ein Lerner, der nur eine einzige Regel lernt (ähnlich dem in der Vorlesung besprochenen `Batch-FindG`).

Vergleichen Sie die Genauigkeit und die Größe (Anzahl der Regeln, Anzahl der Bedingungen) der von `ConjunctiveRule`, JRip, und J48 gelernten Konzepte auf zehn Datensets.

Führen Sie auf den Ergebnissen paarweise einen Signifikanz-Test (Vorzeichen-Test) durch, um festzustellen welche Methode am besten funktioniert.

## 5 Aufgabe: ROC-Kurven

Vergleichen Sie für einen ausgewählten Datensatz die ROC-Kurven bzw. die Fläche unter diesen Kurven für die Klassifizierer J48 und `NaiveBayes`. Sie können die ROC-Kurven betrachten, indem Sie mit der rechten Maustaste im Fenster "Result List" den Menü-Punkt "Threshold List" auswählen.

Interpretieren Sie die Resultate. Sie können die Werte, die zum Zeichnen der Kurve verwendet wurden, auch mit "Save" in ein ARFF-File exportieren, und dieses (nach Löschen des Headers) in Grafik-Programme importieren. So können Sie z.B. beide Kurven (für J48 und `NaiveBayes`) übereinander legen.

## 6 Aufgabe: Pre-Processing

Wählen Sie ein Datenset mit vielen numerischen Attributen aus. Erstellen Sie eine diskretisierte Version (`weka.filters.supervised.attribute.Discretize`).

Schätzen Sie die Genauigkeit von J48 mittels Cross-validation auf den ursprünglichen Daten und auf den diskretisierten Daten ab.

Der Meta-Classifer `FilteredClassifier` erlaubt, eine Kombination einer Pre-processing Methode und eines Classifiers zu einem neuen Classifier zu machen. Erzeugen Sie die Kombination `Discretize` und J48 und schätzen Sie deren Genauigkeit auf den ursprünglichen Daten ab.

Wie interpretieren Sie den Vergleich der Genauigkeiten und der Größe der gelernten Bäume dieser drei Experimente?

## 7 Aufgabe: Entdecken von Assoziationsregeln

Das Datenset `adult` (<http://www.ke.informatik.tu-darmstadt.de/lehre/ws06/mldm/projekt/adult.arff>) enthält Daten von 48842 US Bürgern über Geschlecht, Ausbildung, Familienstand, Beruf, Einkommen (class Variable), etc. Versuchen Sie, mit dem Apriori-Algorithmus aus Weka in diesem Datenset *interessante* Regeln zu finden. Sie können dabei sowohl die Optionen von Weka ausprobieren (z.B. -T das Maß, nach dem die Regeln sortiert werden) als auch das Datenset verändern (z.B. durch Entfernen einzelner Attribute). Beachten

Sie, daß in der Version zum Download zwei numerische Attribute enthalten sind, die Sie diskretisieren oder einfach entfernen können. Falls die Laufzeiten zu lange werden (mehrere Minuten), können Sie auch auf einer Teilmenge der Daten arbeiten.

## **8 Aufgabe: Ensemble-Lernen**

Vergleichen Sie die Genauigkeit von J48, Bagging mit J48, AdaBoost mit J48 und RandomForest auf fünf Datensets. Erhöhen Sie die Anzahl der Iterationen bei den drei Ensemble-Verfahren und beobachten Sie die Entwicklung der erzielten Genauigkeiten.