

## 6.034. Design. Assignment. 2.

April, 5,, 2005,

**Weka Script Due:** Friday April 8, in recitation

**Paper Due:** Wednesday April 13, in class

**Oral reports:** Friday April 15, by appointment

The goal of this assignment is for you to gain some practice in the application of machine learning algorithms to real data. We give you two data sets and a framework that will allow you to experiment with different learning algorithms on that data.

### 1- Data-Sets-

We ask you to build classifiers for two data sets:

1. **credit-g-500.arff:** This is a two-class data set related to credit rating in Germany. Information on the data set can be found at the top of the file.
2. **digits-2-4-5-9.arff:** This is a collection of 4x14 binary images of hand-written digits (2,4,5,9); see figure. There are 250 samples of each digit. Each image is converted into a feature vector by listing the content of the array in row-major order.



### 2- Experiments-

We would like you to find an effective learning algorithm for each of these data sets. In order to do so, you should think about the general strengths and weaknesses of the different learning algorithms, as well as of experiment with them on the data.

Using the algorithm you choose, generate a classifier, and its prediction of how well it will perform on new data. We will fund that classifier on some additional data and compare its performance to your predicted performance.

### 3- Write-up-

In your write-up, describe in detail the process by which you developed the classifiers you did. You should answer the following questions, in detail, including supporting data and graphs or tables.

- What algorithms are generally expected to be appropriate for these data sets?

- How did you choose among the different algorithms? Report your chosen algorithm, as well as fit feast three others that you tried.
- How did you choose parameter settings for each algorithm? Report the parameters that gave you the best results.
- How did you come up with the prediction or how well the classifiers you delivered would perform on previously unseen data? Report your prediction.
- Compare the best performance you got on each data set with the performance you had picked the class(a) that random (unbiased coin flip) for (b) by always predicting the most prevalent class in the training data.
- What classifier would you use in the credit data if it were twice as expensive to say that a person with bad credit was going to have good credit, as to say that a person with good credit would have bad credit?
- What two attributes seem to be most relevant in each data set? Or is it the case that they're all just about equally significant? Explain how you determined this, and why you think you obtained the answer you did.
- In the multiple-class digits problem, which two digits are most frequently confused by your classifier. Does that make sense to you?

## 4- Grading-

There will be a late penalty of 20% per day assessed, with no credit given for assignments turned in after the final report.

Grading will be broken down as follows:

- 30: Good plan for choosing and validating algorithm, parameters, and classifiers
- 15: How effective are the classifiers on the new data
- 10: How good is the supplied performance prediction
- 5: Completing the Weka script given that the end of this handout
- 20: Clarity and organization of written report
- 20: Clarity and understanding in final report

## 5- Software-

We ask you to use the Weka environment for machine learning. You can download the software from:

<http://www.cs.waikato.ac.nz/~ml/weka/>

The software is written in Java and should run under Windows, Linux and Mac. A word of warning: Weka will often run but consume memory and need to be re-started, so have results as you go.

Within this system, you can find the major algorithms that we've studied:

- K-Nearest Neighbor (called K in Weka)
- Decision trees (called J48 in Weka)
- Naive Bayes (called Naive Bayes in Weka)
- SVM (called SMO in Weka)

There are many other algorithms that you can experiment with if you'd like to, but we expect you to consider these methods. Note that the SMO implementation is relatively slow compared to others. You might want to do cross-validation sparingly with SMO since you find that it takes a long time to run.

You should try at least one additional interesting processing step on at least one of the flat sets. You might, for example:

- Run feature selection or dimensionality reduction algorithm on the digits flat
- Try normalizing the attributes before applying the nearest neighbor

Don't just try something at random, though. Think about what is likely to help classification performance, and justify your choice in your write-up.

## 6- Using Weka-

To make sure that you don't (immediately) run out of memory when running the program. Under Linux and MacOSX, you should start Weka by connecting to the Weka directory (weka-3-4-4) and calling java with the following arguments:

```
java -mx100000000 -oss100000000 -jar weka.jar
```

If you are running Windows, Weka will install under

```
c:\Program Files\Weka-3-4
```

In that folder, you should see a file called `runWeka.bat`, edit that file to add the `-mx` and `-oss` arguments to the `java` call. When you start Weka from the Start menu, you should see a console window with the appropriate call to `java`.

There are also a couple of documentation PDFs in the Weka directory: `tutorial.pdf` and `ExplorerGuide.pdf`. This page also has some useful information and links:

<http://www.comp.leeds.ac.uk/andyr/teaching/db32/weka-db32.pdf>

Here's a script to follow that will expose you to the basics of operating Weka. Report the performance values that we ask for below. Hand them in at recitation on April 8. This will be worth 5 points on your grade for the assignment. If you have any questions, please bring them up at recitation.

```
java -mx100000000 -oss100000000 -jar weka.jar
Go to "Weka" GUI "Chooser" window
Click "Explorer"
Choose "Preprocess" Tab at the "top" of the "new" window
Open "File"
<pick "breast-cancer.arff">
```

Clicking on the name of the different "Attributes" shows a histogram of the values on the bottom right, colored by the class variable (or whatever attribute is chosen in the pull-down above the graph). If you click on the class attribute, you'll see how many of each class there are in the flat set.