

Tutorial Exercises for the Weka Explorer

17

The best way to learn about the Explorer interface is simply to use it. This chapter presents a series of tutorial exercises that will help you learn about Explorer and also about practical data mining in general. The first section is introductory, but we think you will find the exercises in the later sections quite thought-provoking.

We begin with a quick, guided tour of the Explorer interface, examining each of the panels and what they can do, which largely parallels the introduction given in Chapter 11. Our screenshots are from Weka 3.6, although almost everything is the same with other versions.

17.1 INTRODUCTION TO THE EXPLORER INTERFACE

Invoke Weka from the Windows Start menu (on Linux or the Mac, double-click *weka.jar* or *weka.app*, respectively). This starts up the Weka GUI Chooser (shown in Figure 11.3(a)). Click the *Explorer* button to enter the Weka Explorer. The Pre-process panel (shown in Figure 11.3(b)) opens up when the Explorer interface is started.

Loading a Dataset

Load a dataset by clicking the *Open file* button in the top left corner of the panel. Inside the *data* folder, which is supplied when Weka is installed, you will find a file named *weather.nominal.arff*. This contains the nominal version of the standard “weather” dataset in Table 1.2. Open this file (the screen will look like Figure 11.3(b)).

As the result shows, the weather data has 14 instances, and 5 attributes called *outlook*, *temperature*, *humidity*, *windy*, and *play*. Click on the name of an attribute in the left subpanel to see information about the selected attribute on the right, such as its values and how many times an instance in the dataset has a particular value. This information is also shown in the form of a histogram. All attributes in this dataset are “nominal”—that is, they have a predefined finite set of values. The last attribute, *play*, is the “class” attribute; its value can be *yes* or *no*.

Familiarize yourself with the Preprocess panel by doing the following exercises. The solutions to these and other exercises in this section are given at the end of the section.

Exercise 17.1.1. What are the values that the attribute *temperature* can have?

Exercise 17.1.2. Load a new dataset. Click the *Open file* button and select the file *iris.arff*, which corresponds to the iris dataset in Table 1.4. How many instances does this dataset have? How many attributes? What is the range of possible values of the attribute *petallength*?

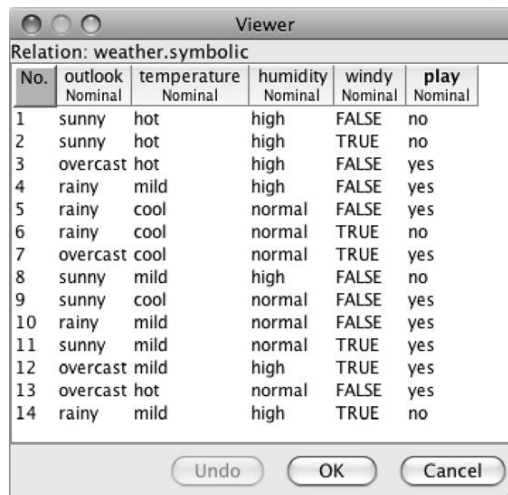
The Dataset Editor

It is possible to view and edit an entire dataset from within Weka. To do this, load the *weather.nominal.arff* file again. Click the *Edit* button from the row of buttons at the top of the Preprocess panel. This opens a new window called Viewer, which lists all instances of the weather data (see Figure 17.1).

Exercise 17.1.3. What is the function of the first column in the Viewer window?

Exercise 17.1.4. What is the class value of instance number 8 in the weather data?

Exercise 17.1.5. Load the iris data and open it in the editor. How many numeric and how many nominal attributes does this dataset have?



| No. | outlook Nominal | temperature Nominal | humidity Nominal | windy Nominal | play Nominal |
|-----|--------------------|------------------------|---------------------|------------------|-----------------|
| 1 | sunny | hot | high | FALSE | no |
| 2 | sunny | hot | high | TRUE | no |
| 3 | overcast | hot | high | FALSE | yes |
| 4 | rainy | mild | high | FALSE | yes |
| 5 | rainy | cool | normal | FALSE | yes |
| 6 | rainy | cool | normal | TRUE | no |
| 7 | overcast | cool | normal | TRUE | yes |
| 8 | sunny | mild | high | FALSE | no |
| 9 | sunny | cool | normal | FALSE | yes |
| 10 | rainy | mild | normal | FALSE | yes |
| 11 | sunny | mild | normal | TRUE | yes |
| 12 | overcast | mild | high | TRUE | yes |
| 13 | overcast | hot | normal | FALSE | yes |
| 14 | rainy | mild | high | TRUE | no |

FIGURE 17.1

The data viewer.

Applying a Filter

As you know, Weka “filters” can be used to modify datasets in a systematic fashion—that is, they are data preprocessing tools. Reload the *weather.nominal* dataset, and let’s remove an attribute from it. The appropriate filter is called *Remove*; its full name is

```
weka.filters.unsupervised.attribute.Remove
```

Examine this name carefully. Filters are organized into a hierarchical structure of which the root is *weka*. Those in the *unsupervised* category don’t require a class attribute to be set; those in the *supervised* category do. Filters are further divided into ones that operate primarily on attributes (the *attribute* category) and ones that operate primarily on instances (the *instance* category).

Click the *Choose* button in the Preprocess panel to open a hierarchical menu (shown in Figure 11.9(a)) from which you select a filter by following the path corresponding to its full name. Use the path given in the full name above to select the *Remove* filter. The text “Remove” will appear in the field next to the *Choose* button.

Click on the field containing this text. The Generic Object Editor window, which is used throughout Weka to set parameter values for all of the tools, opens. In this case it contains a short explanation of the *Remove* filter (shown in Figure 11.9(b))—click *More* to get a fuller description (Figure 11.9(c)). Enter 3 into the *attribute-Indices* field and click the *OK* button. The window with the filter options closes. Now click the *Apply* button on the right, which runs the data through the filter. The filter removes the attribute with index 3 from the dataset, and you can see that this has happened. This change does not affect the dataset in the file; it only applies to the data held in memory. The changed dataset can be saved to a new ARFF file by pressing the *Save* button and entering a file name. The action of the filter can be undone by pressing the *Undo* button. Again, this applies to the version of the data held in memory.

What we have described illustrates how filters are applied to data. However, in the particular case of *Remove*, there is a simpler way of achieving the same effect. Instead of invoking a filter, attributes can be selected using the small boxes in the Attributes subpanel and removed using the *Remove* button that appears at the bottom, below the list of attributes.

Exercise 17.1.6. Load the *weather.nominal* dataset. Use the filter *weka.unsupervised.instance.RemoveWithValues* to remove all instances in which the *humidity* attribute has the value *high*. To do this, first make the field next to the *Choose* button show the text *RemoveWithValues*. Then click on it to get the Generic Object Editor window, and figure out how to change the filter settings appropriately.

Exercise 17.1.7. Undo the change to the dataset that you just performed, and verify that the data has reverted to its original state.

The Visualize Panel

Now take a look at Weka's data visualization facilities. These work best with numeric data, so we use the iris data. Load *iris.arff*, which contains the iris dataset of Table 1.4 containing 50 examples of three types of Iris: *Iris setosa*, *Iris versicolor*, and *Iris virginica*.

Click the *Visualize* tab to bring up the Visualize panel (shown in Figure 11.17). Click the first plot in the second row to open up a window showing an enlarged plot using the selected axes. Instances are shown as little crosses, the color of which depends on the instance's class. The *x*-axis shows the *sepalength* attribute, and the *y*-axis shows *petalwidth*.

Clicking on one of the crosses opens up an Instance Info window, which lists the values of all attributes for the selected instance. Close the Instance Info window again.

The selection fields at the top of the window containing the scatter plot determine which attributes are used for the *x*- and *y*-axes. Change the *x*-axis to *petalwidth* and the *y*-axis to *petallength*. The field showing *Color: class (Num)* can be used to change the color coding.

Each of the barlike plots to the right of the scatter plot window represents a single attribute. In each bar, instances are placed at the appropriate horizontal position and scattered randomly in the vertical direction. Clicking a bar uses that attribute for the *x*-axis of the scatter plot. Right-clicking a bar does the same for the *y*-axis. Use these bars to change the *x*- and *y*-axes back to *sepalength* and *petalwidth*.

The *Jitter* slider displaces the cross for each instance randomly from its true position, and can reveal situations where instances lie on top of one another. Experiment a little by moving the slider.

The *Select Instance* button and the *Reset*, *Clear*, and *Save* buttons let you modify the dataset. Certain instances can be selected and the others removed. Try the *Rectangle* option: Select an area by left-clicking and dragging the mouse. The *Reset* button changes into a *Submit* button. Click it, and all instances outside the rectangle are deleted. You could use *Save* to save the modified dataset to a file. *Reset* restores the original dataset.

The Classify Panel

Now we apply a classifier to the weather data. Load the weather data again. Go to the Preprocess panel, click the *Open file* button, and select *weather.nominal.arff* from the data directory. Then switch to the Classify panel (shown in Figure 11.4(b)) by clicking the *Classify* tab at the top of the window.

Using the C4.5 Classifier

As you learned in Chapter 11 (page 410), the C4.5 algorithm for building decision trees is implemented in Weka as a classifier called *J48*. Select it by clicking the *Choose*

button near the top of the *Classify* tab. A dialog window appears showing various types of classifier. Click the *trees* entry to reveal its subentries, and click *J48* to choose that classifier. Classifiers, like filters, are organized in a hierarchy: *J48* has the full name *weka.classifiers.trees.J48*.

The classifier is shown in the text box next to the *Choose* button: It now reads *J48 -C 0.25 -M 2*. This text gives the default parameter settings for this classifier, which in this case rarely require changing to obtain good performance.

For illustrative purposes we evaluate the performance using the training data, which has been loaded in the Preprocess panel—this is *not* generally a good idea because it leads to unrealistically optimistic performance estimates. Choose *Use training set* from the *Test options* part of the Classify panel. Once the test strategy has been set, the classifier is built and evaluated by pressing the *Start* button. This processes the training set using the currently selected learning algorithm, C4.5 in this case. Then it classifies all the instances in the training data and outputs performance statistics. These are shown in Figure 17.2(a).

Interpreting the Output

The outcome of training and testing appears in the Classifier Output box on the right. Scroll through the text and examine it. First, look at the part that describes the decision tree, reproduced in Figure 17.2(b). This represents the decision tree that was built, including the number of instances that fall under each leaf. The textual representation is clumsy to interpret, but Weka can generate an equivalent graphical version.

Here's how to get the graphical tree. Each time the *Start* button is pressed and a new classifier is built and evaluated, a new entry appears in the Result List panel in the lower left corner of Figure 17.2(a). To see the tree, right-click on the entry *trees.J48* that has just been added to the result list and choose *Visualize tree*. A window pops up that shows the decision tree in the form illustrated in Figure 17.3. Right-click a blank spot in this window to bring up a new menu enabling you to auto-scale the view. You can pan around by dragging the mouse.

Now look at the rest of the information in the Classifier Output area. The next two parts of the output report on the quality of the classification model based on the chosen test option.

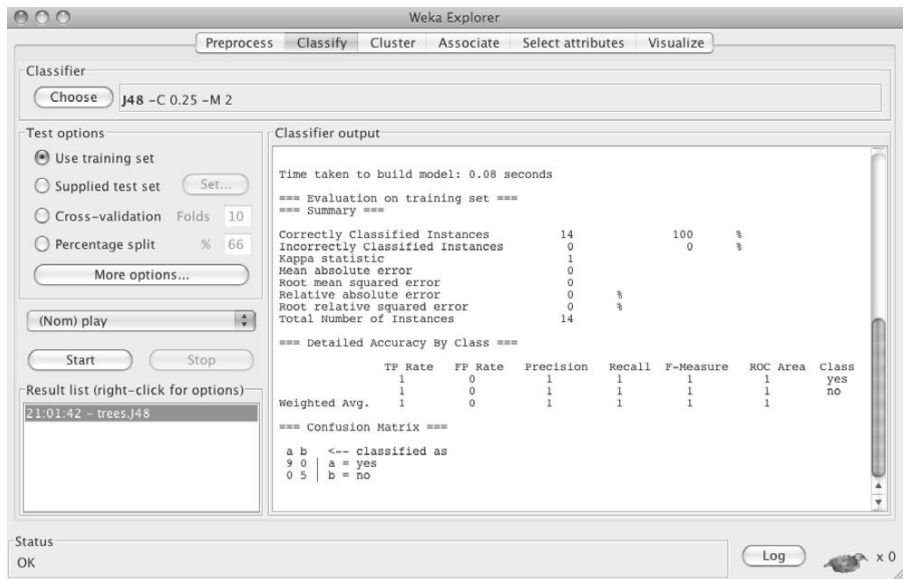
This text states how many and what proportion of test instances have been correctly classified:

```
Correctly Classified Instances   14    100%
```

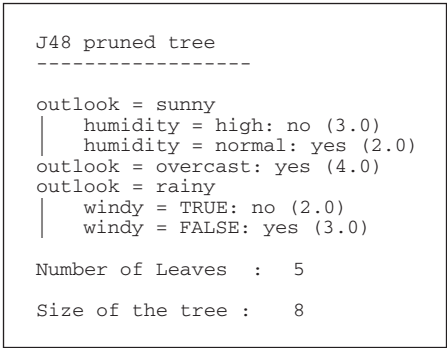
This is the accuracy of the model on the data used for testing. In this case it is completely accurate (100%), which is often the case when the training set is used for testing.

At the bottom of the output is the confusion matrix:

```
=== Confusion Matrix ===
 a b    <- classified as
 9 0 | a = yes
 0 5 | b = no
```



(a)



(b)

FIGURE 17.2 Output after building and testing the classifier: (a) screenshot and (b) decision tree.

Each element in the matrix is a count of instances. Rows represent the true classes, and columns represent the predicted classes. As you can see, all 9 *yes* instances have been predicted as *yes*, and all 5 *no* instances as *no*.

Exercise 17.1.8. How would this instance be classified using the decision tree?

outlook = sunny, temperature = cool, humidity = high, windy = TRUE

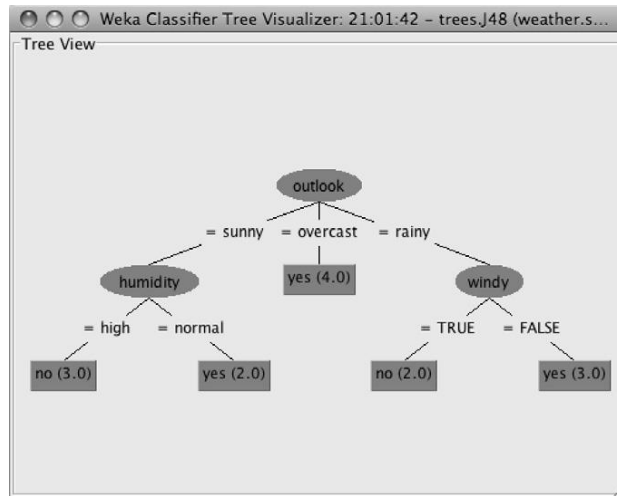


FIGURE 17.3

The decision tree that has been built.

Setting the Test Method

When the *Start* button is pressed, the selected learning algorithm is run and the dataset that was loaded in the Preprocess panel is used with the selected test protocol. For example, in the case of tenfold cross-validation this involves running the learning algorithm 10 times to build and evaluate 10 classifiers. A model built from the *full* training set is then printed into the Classifier Output area: This may involve running the learning algorithm one final time. The remainder of the output depends on the test protocol that was chosen using test options; these options were discussed in Section 11.1.

Exercise 17.1.9. Load the iris data using the Preprocess panel. Evaluate C4.5 on this data using (a) the training set and (b) cross-validation. What is the estimated percentage of correct classifications for (a) and (b)? Which estimate is more realistic?

Visualizing Classification Errors

Right-click the *trees.J48* entry in the result list and choose *Visualize classifier errors*. A scatter plot window pops up. Instances that have been classified correctly are marked by little crosses; ones that are incorrect are marked by little squares.

Exercise 17.1.10. Use the *Visualize classifier errors* function to find the wrongly classified test instances for the cross-validation performed in Exercise 17.1.9. What can you say about the location of the errors?

17.2 NEAREST-NEIGHBOR LEARNING AND DECISION TREES

In this section you will experiment with nearest-neighbor classification and decision tree learning. For most of it, a real-world forensic glass classification dataset is used.

We begin by taking a preliminary look at the dataset. Then we examine the effect of selecting different attributes for nearest-neighbor classification. Next we study class noise and its impact on predictive performance for the nearest-neighbor method. Following that we vary the training set size, both for nearest-neighbor classification and for decision tree learning. Finally, you are asked to interactively construct a decision tree for an image segmentation dataset.

Before continuing you should review in your mind some aspects of the classification task:

- How is the accuracy of a classifier measured?
- To make a good classifier, are all the attributes necessary?
- What is class noise, and how would you measure its effect on learning?
- What is a learning curve?
- If you, personally, had to invent a decision tree classifier for a particular dataset, how would you go about it?

The Glass Dataset

The glass dataset *glass.arff* from the U.S. Forensic Science Service contains data on six types of glass. Glass is described by its refractive index and the chemical elements that it contains; the aim is to classify different types of glass based on these features. This dataset is taken from the UCI datasets, which have been collected by the University of California at Irvine and are freely available on the Web. They are often used as a benchmark for comparing data mining algorithms.

Find the dataset *glass.arff* and load it into the Explorer interface. For your own information, answer the following exercises, which review material covered in the previous section.

Exercise 17.2.1. How many attributes are there in the dataset? What are their names? What is the class attribute? Run the classification algorithm *IBk* (*weka.classifiers.lazy.IBk*). Use cross-validation to test its performance, leaving the number of folds at the default value of 10. Recall that you can examine the classifier options in the Generic Object Editor window that pops up when you click the text beside the *Choose* button. The default value of the KNN field is 1: This sets the number of neighboring instances to use when classifying.

Exercise 17.2.2. What is the accuracy of *IBk* (given in the Classifier Output box)? Run *IBk* again, but increase the number of neighboring instances to $k = 5$ by entering this value in the KNN field. Here and throughout this section, continue to use cross-validation as the evaluation method.

Exercise 17.2.3. What is the accuracy of *IBk* with five neighboring instances ($k = 5$)?

Attribute Selection

Now we investigate which subset of attributes produces the best cross-validated classification accuracy for the *IBk* algorithm on the glass dataset. Weka contains automated attribute selection facilities, which are examined in a later section, but it is instructive to do this manually.

Performing an exhaustive search over all possible subsets of the attributes is infeasible (why?), so we apply the backward elimination procedure described in Section 7.1 (page 311). To do this, first consider dropping each attribute individually from the full dataset, and run a cross-validation for each reduced version. Once you have determined the best eight-attribute dataset, repeat the procedure with this reduced dataset to find the best seven-attribute dataset, and so on.

Exercise 17.2.4. Record in Table 17.1 the best attribute set and the greatest accuracy obtained in each iteration. The best accuracy obtained in this process is quite a bit higher than the accuracy obtained on the full dataset.

Exercise 17.2.5. Is this best accuracy an unbiased estimate of accuracy on future data? Be sure to explain your answer. (*Hint:* To obtain an unbiased estimate of accuracy on future data, we must not look at the test data *at all*

Table 17.1 Accuracy Obtained Using *IBk*, for Different Attribute Subsets

| Subset Size (No. of Attributes) | Attributes in “Best” Subset | Classification Accuracy |
|------------------------------------|--------------------------------|----------------------------|
| 9 | | |
| 8 | | |
| 7 | | |
| 6 | | |
| 5 | | |
| 4 | | |
| 3 | | |
| 2 | | |
| 1 | | |
| 0 | | |

when producing the classification model for which the estimate is being obtained.)

Class Noise and Nearest-Neighbor Learning

Nearest-neighbor learning, like other techniques, is sensitive to noise in the training data. In this section we inject varying amounts of class noise into the data and observe the effect on classification performance.

You can flip a certain percentage of class labels in the data to a randomly chosen other value using an unsupervised attribute filter called *AddNoise*, in *weka.filters.unsupervised.attribute*. However, for this experiment it is important that the test data remains unaffected by class noise. Filtering the training data without filtering the test data is a common requirement, and is achieved using a metalearner called *FilteredClassifier*, in *weka.classifiers.meta*, as described near the end of Section 11.3 (page 444). This metalearner should be configured to use *IBk* as the classifier and *AddNoise* as the filter. *FilteredClassifier* applies the filter to the data before running the learning algorithm. This is done in two batches: first the training data and then the test data. The *AddNoise* filter only adds noise to the first batch of data it encounters, which means that the test data passes through unchanged.

Exercise 17.2.6. Record in Table 17.2 the cross-validated accuracy estimate of *IBk* for 10 different percentages of class noise and neighborhood sizes $k = 1$, $k = 3$, $k = 5$ (determined by the value of k in the k -nearest-neighbor classifier).

Exercise 17.2.7. What is the effect of increasing the amount of class noise?

Exercise 17.2.8. What is the effect of altering the value of k ?

| Table 17.2 Effect of Class Noise on <i>IBk</i> , for Different Neighborhood Sizes | | | |
|---|---------|---------|---------|
| Percentage Noise | $k = 1$ | $k = 3$ | $k = 5$ |
| 0% | | | |
| 10% | | | |
| 20% | | | |
| 30% | | | |
| 40% | | | |
| 50% | | | |
| 60% | | | |
| 70% | | | |
| 80% | | | |
| 90% | | | |
| 100% | | | |

Table 17.3 Effect of Training Set Size on *IBk* and *J48*

| Percentage of Training Set | <i>IBk</i> | <i>J48</i> |
|----------------------------|------------|------------|
| 10% | | |
| 20% | | |
| 30% | | |
| 40% | | |
| 50% | | |
| 60% | | |
| 70% | | |
| 80% | | |
| 90% | | |
| 100% | | |

Varying the Amount of Training Data

This section examines learning curves, which show the effect of gradually increasing the amount of training data. Again, we use the glass data, but this time with both *IBk* and the C4.5 decision tree learners, implemented in Weka as *J48*.

To obtain learning curves, use *FilteredClassifier* again, this time in conjunction with *weka.filters.unsupervised.instance.Resample*, which extracts a certain specified percentage of a given dataset and returns the reduced dataset.¹ Again, this is done only for the first batch to which the filter is applied, so the test data passes unmodified through the *FilteredClassifier* before it reaches the classifier.

Exercise 17.2.9. Record in Table 17.3 the data for learning curves for both the one-nearest-neighbor classifier (i.e., *IBk* with $k = 1$) and *J48*.

Exercise 17.2.10. What is the effect of increasing the amount of training data?

Exercise 17.2.11. Is this effect more pronounced for *IBk* or *J48*?

Interactive Decision Tree Construction

One of Weka's classifiers is interactive: It lets the user—you!—construct your own classifier. Here's a competition: Who can build a classifier with the highest predictive accuracy?

Follow the procedure described in Section 11.2 (page 424). Load the file *segment-challenge.arff* (in the data folder that comes with the Weka distribution). This dataset

¹This filter performs sampling with replacement, rather than sampling without replacement, but the effect is minor and we will ignore it here.

has 20 attributes and 7 classes. It is an image segmentation problem, and the task is to classify images into seven different groups based on properties of the pixels.

Set the classifier to *UserClassifier*, in the *weka.classifiers.trees* package. We use a separate test set (performing cross-validation with *UserClassifier* is incredibly tedious!), so in the *Test options* box choose the *Supplied test set* option and click the *Set* button. A small window appears in which you choose the test set. Click *Open file* and browse to the file *segment-test.arff* (also in the Weka distribution's data folder). On clicking *Open*, the small window updates to show the number of attributes (20) in the data. The number of instances is not displayed because test instances are read incrementally (so that the Explorer interface can process larger test files than can be accommodated in main memory).

Click *Start*. *UserClassifier* differs from all other classifiers: It opens a special window and waits for you to build your own classifier in it. The tabs at the top of the window switch between two views of the classifier. The *Tree visualizer* shows the current state of your tree, and the nodes give the number of class values there. The aim is to come up with a tree of which the leaf nodes are as pure as possible. To begin with, the tree has just one node—the root node—containing all the data. More nodes will appear when you proceed to split the data in the *Data visualizer*.

Click the *Data visualizer* tab to see a two-dimensional plot in which the data points are color-coded by class, with the same facilities as the Visualize panel discussed in Section 17.1. Try different combinations of *x*- and *y*-axes to get the clearest separation you can find between the colors. Having found a good separation, you then need to select a region in the plot: This will create a branch in the tree. Here's a hint to get you started: Plot *region-centroid-row* on the *x*-axis and *intensity-mean* on the *y*-axis (the display is shown in Figure 11.14(a)); you can see that the red class (*sky*) is nicely separated from the rest of the classes at the top of the plot.

There are four tools for selecting regions in the graph, chosen using the dropdown menu below the *y*-axis selector. *Select Instance* identifies a particular instance. *Rectangle* (shown in Figure 11.14(a)) allows you to drag out a rectangle on the graph. With *Polygon* and *Polyline* you build a free-form polygon or draw a free-form polyline (left-click to add a vertex and right-click to complete the operation).

When you have selected an area using any of these tools, it turns gray. (In Figure 11.14(a) the user has defined a rectangle.) Clicking the *Clear* button cancels the selection without affecting the classifier. When you are happy with the selection, click *Submit*. This creates two new nodes in the tree, one holding all the instances covered by the selection and the other holding all remaining instances. These nodes correspond to a binary split that performs the chosen geometric test.

Switch back to the *Tree visualizer* view to examine the change in the tree. Clicking on different nodes alters the subset of data that is shown in the *Data visualizer* section. Continue adding nodes until you obtain a good separation of the classes—that is, the leaf nodes in the tree are mostly pure. Remember, however, that you should not overfit the data because your tree will be evaluated on a separate test set.

When you are satisfied with the tree, right-click any blank space in the *Tree visualizer* view and choose *Accept The Tree*. Weka evaluates the tree against the test set and outputs statistics that show how well you did.

Exercise 17.2.12. You are competing for the best accuracy score of a hand-built *UserClassifier* produced on the *segment-challenge* dataset and tested on the *segment-test* set. Try as many times as you like. When you have a good score (anything close to 90% correct or better), right-click the corresponding entry in the Result list, save the output using *Save result buffer*, and copy it into your answer for this exercise. Then run *J48* on the data to see how well an automatic decision tree learner performs on the task.

17.3 CLASSIFICATION BOUNDARIES

In this section we examine the classification boundaries that are produced by different types of models. To do this, we use Weka's *Boundary Visualizer*, which is not part of the Explorer interface. To find it, start up the Weka GUI Chooser as usual from the Windows Start menu (on Linux or the Mac, double-click *weka.jar* or *weka.app*, respectively) and select *BoundaryVisualizer* from the *Visualization* menu at the top.

The boundary visualizer shows a two-dimensional plot of the data and is most appropriate for datasets with two numeric attributes. We will use a version of the iris data without the first two attributes. To create this, start up the Explorer interface, load *iris.arff* using the *Open file* button, and remove the first two attributes (*sepal-length* and *sepalwidth*) by selecting them and clicking the *Remove* button that appears at the bottom. Then save the modified dataset to a file (using *Save*) called, say, *iris.2D.arff*.

Now leave the Explorer interface and open this file for visualization using the boundary visualizer's *Open file* button. Initially, the plot just shows the data in the dataset.

Visualizing 1R

The purpose of the boundary visualizer is to show the predictions of a given model for every possible combination of attribute values—that is, for every point in the two-dimensional space. The points are color-coded according to the prediction the model generates. We will use this to investigate the decision boundaries that different classifiers generate for the reduced iris dataset.

Start with the 1R rule learner. Use the *Choose* button of the boundary visualizer to select *weka.classifiers.rules.OneR*. Make sure you tick *Plot training data*; otherwise, only the predictions will be plotted. Then click the *Start* button. The program starts plotting predictions in successive scan lines. Click the *Stop* button once the

plot has stabilized—as soon as you like, in this case—and the training data will be superimposed on the boundary visualization.

Exercise 17.3.1. Explain the plot based on what you know about 1R. (*Hint:* Use the Explorer interface to look at the rule set that 1R generates for this data.)

Exercise 17.3.2. Study the effect of the *minBucketSize* parameter on the classifier by regenerating the plot with values of 1, and then 20, and then some critical values in between. Describe what you see, and explain it. (*Hint:* You could speed things up by using the Explorer interface to look at the rule sets.)

Now answer the following questions by thinking about the internal workings of 1R. (*Hint:* It will probably be fastest to use the Explorer interface to look at the rule sets.)

Exercise 17.3.3. You saw earlier that when visualizing 1R the plot always has three regions. But why aren't there more for small bucket sizes (e.g., 1)? Use what you know about 1R to explain this apparent anomaly.

Exercise 17.3.4. Can you set *minBucketSize* to a value that results in less than three regions? What is the smallest possible number of regions? What is the smallest value for *minBucketSize* that gives this number of regions? Explain the result based on what you know about the iris data.

Visualizing Nearest-Neighbor Learning

Now let's examine the classification boundaries created by the nearest-neighbor method. Use the boundary visualizer's *Choose* button to select the *IBk* classifier (*weka.classifiers.lazy.IBk*) and plot its decision boundaries for the reduced iris data.

OneR's predictions are categorical: For each instance, they predict one of the three classes. In contrast, *IBk* outputs probability estimates for each class, and the boundary visualizer uses them to mix the colors red, green, and blue that correspond to the three classes. *IBk* estimates class probabilities by looking at the set of *k*-nearest neighbors of a test instance and counting the number in each class.

Exercise 17.3.5. With $k = 1$, which is the default value, it seems that the set of *k*-nearest neighbors could have only one member and therefore the color will always be pure red, green, or blue. Looking at the plot, this is indeed almost always the case: There is no mixing of colors because one class gets a probability of 1 and the others a probability of 0. Nevertheless, there is a small area in the plot where two colors are in fact mixed. Explain this. (*Hint:* Examine the data carefully using the Explorer interface's Visualize panel.)

Exercise 17.3.6. Experiment with different values of *k*, say 5 and 10. Describe what happens as *k* increases.

Visualizing *NaïveBayes*

Turn now to the *NaïveBayes* classifier. Its assumption that attributes are conditionally independent given a particular class value means that the overall class probability is obtained by simply multiplying the per-attribute conditional probabilities together (and taking into account the class prior probabilities as well). In other words, with two attributes, if you know the class probabilities along the x - and y -axes (and the class prior), you can calculate the value for any point in space by multiplying them together (and then normalizing). This is easy to understand if you visualize it as a boundary plot.

Plot the predictions of *NaïveBayes*. But first discretize the attribute values. By default, Weka's *NaïveBayes* classifier assumes that the attributes are normally distributed given the class. You should override this by setting *useSupervisedDiscretization* to *true* using the Generic Object Editor window. This will cause *NaïveBayes* to discretize the numeric attributes in the data with a supervised discretization technique. In most practical applications of *NaïveBayes*, supervised discretization works better than the default method. It also produces a more comprehensible visualization, which is why we use it here.

Exercise 17.3.7. The plot that is generated by visualizing the predicted class probabilities of *NaïveBayes* for each pixel location is quite different from anything we have seen so far. Explain the patterns in it.

Visualizing Decision Trees and Rule Sets

Decision trees and rule sets are similar to nearest-neighbor learning in the sense that they are quasi-universal: In principle, they can approximate any decision boundary arbitrarily closely. In this section, we look at the boundaries generated by *JRip* and *J48*.

Generate a plot for *JRip*, with default options.

Exercise 17.3.8. What do you see? Relate the plot to the output of the rules that you get by processing the data in the Explorer.

Exercise 17.3.9. The *JRip* output assumes that the rules will be executed in the correct sequence. Write down an equivalent set of rules that achieves the same effect regardless of the order in which they are executed. Generate a plot for *J48*, with default options.

Exercise 17.3.10. What do you see? Again, relate the plot to the output that you get by processing the data in the Explorer interface. One way to control how much pruning *J48* performs is to adjust the minimum number of instances required in a leaf, *minNumObj*.

Exercise 17.3.11. Suppose you want to generate trees with 3, 2, and 1 leaf node, respectively. What are the exact ranges of values for *minNumObj* that achieve this, given default values for the other parameters?

Messing with the Data

With the Boundary Visualizer you can modify the data by adding or removing points.

Exercise 17.3.12. Introduce some noise into the data and study the effect on the learning algorithms we looked at above. What kind of behavior do you observe for each algorithm as you introduce more noise?

17.4 PREPROCESSING AND PARAMETER TUNING

Now we look at some useful preprocessing techniques, which are implemented as filters, as well as a method for automatic parameter tuning.

Discretization

As we know, there are two types of discretization techniques: *unsupervised* ones, which are “class blind,” and *supervised* ones, which take the class value of the instances into account when creating intervals. Weka’s main unsupervised method for discretizing numeric attributes is *weka.filters.unsupervised.attribute.Discretize*. It implements these two methods: equal-width (the default) and equal-frequency discretization.

Find the glass dataset *glass.arff* and load it into the Explorer interface. Apply the unsupervised discretization filter in the two different modes explained previously.

Exercise 17.4.1. What do you observe when you compare the histograms obtained? The one for equal-frequency discretization is quite skewed for some attributes. Why?

The main *supervised* technique for discretizing numeric attributes is *weka.filters.supervised.attribute.Discretize*. Locate the iris data, load it, apply the supervised discretization scheme, and look at the histograms obtained. Supervised discretization strives to create intervals within which the class distribution is consistent, although the distributions vary from one interval to the next.

Exercise 17.4.2. Based on the histograms obtained, which of the discretized attributes would you consider to be most predictive? Reload the glass data and apply supervised discretization to it.

Exercise 17.4.3. For some attributes there is only a single bar in the histogram. What does that mean?

Discretized attributes are normally coded as nominal attributes, with one value per range. However, because the ranges are ordered, a discretized attribute is actually on an ordinal scale. Both filters have the ability to create binary attributes rather than multivalued ones, by setting the option *makeBinary* to *true*.

Exercise 17.4.4. Choose one of the filters and use it to create binary attributes. Compare the result with the output generated when *makeBinary* is *false*. What do the binary attributes represent?

More on Discretization

Here we examine the effect of discretization when building a *J48* decision tree for the data in *ionosphere.arff*. This dataset contains information about radar signals returned from the ionosphere. “Good” samples are those showing evidence of some type of structure in the ionosphere, while for “bad” ones the signals pass directly through the ionosphere. For more details, take a look at the comments in the ARFF file. Begin with unsupervised discretization.

Exercise 17.4.5. For *J48*, compare cross-validated accuracy and the size of the trees generated for (1) the raw data, (2) data discretized by the unsupervised discretization method in default mode, and (3) data discretized by the same method with binary attributes.

Now turn to supervised discretization. Here a subtle issue arises, discussed near the end of Section 11.3 (page 432). If Exercise 17.4.5 were simply repeated using a supervised discretization method, the result would be overoptimistic. In effect, because cross-validation is used for evaluation, *the data in the test set has been taken into account when determining the discretization intervals*. This does not give a fair estimate of performance on fresh data.

To evaluate supervised discretization fairly, use *FilteredClassifier* from Weka’s metalearners. This builds the filter using the training data only, and then evaluates it on the test data using the discretization intervals computed for the training data. After all, that is how you would have to process fresh data in practice.

Exercise 17.4.6. Using *FilteredClassifier* and *J48*, compare cross-validated accuracy and the size of the trees generated for (4) supervised discretization in default mode, and (5) supervised discretization with binary attributes.

Exercise 17.4.7. Compare these with the results for the raw data from Exercise 17.4.5. How can decision trees generated from discretized data possibly be better predictors than ones built from raw numeric data?

Automatic Attribute Selection

In most practical applications of supervised learning not all attributes are equally useful for predicting the target. For some learning schemes, redundant and/or irrelevant attributes can result in less accurate models. As you found in Section 17.2, it is tedious to identify useful attributes in a dataset manually; automatic attribute selection methods are usually more appropriate.

Attribute selection methods can be divided into filter and wrapper methods (see Section 7.1, page 308). The former apply a computationally efficient heuristic to measure the quality of a subset of attributes; the latter measure the quality of an attribute subset by building and evaluating an actual classification model, which is more expensive but often delivers superior performance.

The Explorer interface's *Select attributes* panel applies attribute selection methods to datasets. The default is to use *CfsSubsetEval*, described in Section 11.8 (page 488), which evaluates subsets of attributes. An alternative is to evaluate attributes individually using an evaluator like *InfoGainAttributeEval* (see Section 11.8, page 491) and then rank them by applying a special “search” method, namely the *Ranker*, as described Section 11.8 (page 490).

Exercise 17.4.8. Apply the ranking technique to the labor negotiations data in *labor.arff* to determine the four most important attributes based on information gain.²

CfsSubsetEval aims to identify a subset of attributes that are highly correlated with the target while not being strongly correlated with one another. It searches through the space of possible attribute subsets for the “best” one using the *BestFirst* search method by default, although other methods can be chosen. In fact, choosing *GreedyStepwise* and setting *searchBackwards* to *true* gives backward elimination, the search method you used manually in Section 17.2.

To use the wrapper method rather than a filter method, such as *CfsSubsetEval*, first select *WrapperSubsetEval* and then configure it by choosing a learning algorithm to apply and setting the number of cross-validation folds to use when evaluating it on each attribute subset.

Exercise 17.4.9. On the same data, run *CfsSubsetEval* for correlation-based selection, using the *BestFirst* search. Then run the wrapper method with *J48* as the base learner, again using the *BestFirst* search. Examine the attribute subsets that are output. Which attributes are selected by both methods? How do they relate to the output generated by ranking using information gain?

More on Automatic Attribute Selection

The *Select attributes* panel allows us to gain insight into a dataset by applying attribute selection methods to it. However, as with supervised discretization, using this information to reduce a dataset becomes problematic if some of the reduced data is used for testing the model (as in cross-validation). Again, the reason is that we have

²Note that most evaluators, including *InfoGainAttributeEval* and *CfsSubsetEval*, discretize numeric attributes using Weka's supervised discretization method before evaluating them.

looked at the class labels in the test data while selecting attributes, and using the test data to influence the construction of a model biases the accuracy estimates obtained.

This can be avoided by dividing the data into training and test sets and applying attribute selection to the training set only. However, it is usually more convenient to use *AttributeSelectedClassifier*, one of Weka's metalearners, which allows an attribute selection method and a learning algorithm to be specified as part of a classification scheme. *AttributeSelectedClassifier* ensures that the chosen set of attributes is selected based on the training data only.

Now we test the three attribute selection methods from above in conjunction with *NaïveBayes*. *NaïveBayes* assumes independence of attributes, so attribute selection can be very helpful. You can see the effect of redundant attributes by adding multiple copies of an attribute using the filter *weka.filters.unsupervised.attribute.Copy* in the Preprocess panel. Each copy is obviously perfectly correlated with the original.

Exercise 17.4.10. Load the diabetes classification data in *diabetes.arff* and add copies of the first attribute. Measure the performance of *NaïveBayes* (with *useSupervisedDiscretization* turned on) using cross-validation after you have added each one. What do you observe?

Do the above three attribute selection methods, used in conjunction with *AttributeSelectedClassifier* and *NaïveBayes*, successfully eliminate the redundant attributes? Run each method from within *AttributeSelectedClassifier* to see the effect on cross-validated accuracy and check the attribute subset selected by each method. Note that you need to specify the number of ranked attributes to use for the *Ranker* method. Set this to 8 because the original diabetes data contains 8 attributes (excluding the class). Specify *NaïveBayes* as the classifier to be used inside the wrapper method because this is the classifier for which we want to select a subset.

Exercise 17.4.11. What can you say regarding the performance of the three attribute selection methods? Do they succeed in eliminating redundant copies? If not, why?

Automatic Parameter Tuning

Many learning algorithms have parameters that can affect the outcome of learning. For example, the decision tree learner C4.5 has two parameters that influence the amount of pruning (we saw one, the minimum number of instances required in a leaf, in Section 17.3). The *k*-nearest-neighbor classifier *IBk* has a parameter (*k*) that sets the neighborhood size. But manually tweaking parameter settings is tedious, just like manually selecting attributes, and presents the same problem: The test data must not be used when selecting parameters; otherwise, the performance estimate will be biased.

Weka's metalearner *CVParameterSelection* searches for the best parameter settings by optimizing cross-validated accuracy on the training data. By default, each setting is evaluated using tenfold cross-validation. The parameters to optimize are specified using the *CVParameters* field in the Generic Object Editor window. For each parameter, three pieces of information must be supplied: (1) a string that names it using its letter code (which can be found in the Javadoc for the corresponding classifier—see Section 14.2, page 525); (2) a numeric range of values to evaluate; and (3) the number of steps to try in this range (note that the parameter is assumed to be numeric). Click on the *More* button in the Generic Object Editor window for more information and an example.

For the diabetes data used in the previous section, use *CVParameterSelection* in conjunction with *IBk* to select the best value for the neighborhood size, ranging from 1 to 10 in 10 steps. The letter code for the neighborhood size is *K*. The cross-validated accuracy of the parameter-tuned version of *IBk* is directly comparable with its accuracy using default settings because tuning is performed by applying inner cross-validation runs to find the best parameter value for each training set occurring in the outer cross-validation—and the latter yields the final performance estimate.

Exercise 17.4.12. What accuracy is obtained in each case? What value is selected for the parameter-tuned version based on cross-validation on the full data set? (*Note:* This value is output in the Classifier Output text area because, as mentioned earlier, the model that is output is the one built from the full dataset.)

Now consider parameter tuning for *J48*. If there is more than one parameter string in the *CVParameters* field, *CVParameterSelection* performs a grid search on the parameters simultaneously. The letter code for the pruning confidence parameter is *C*, and you should evaluate values from 0.1 to 0.5 in five steps. The letter code for the minimum leaf size parameter is *M*, and you should evaluate values from 1 to 10 in 10 steps.

Exercise 17.4.13. Run *CVParameterSelection* to find the best parameter value setting. Compare the output you get to that obtained from *J48* with default parameters. Has accuracy changed? What about tree size? What parameter values were selected by *CVParameterSelection* for the model built from the full training set?

17.5 DOCUMENT CLASSIFICATION

Next we perform some experiments in document classification. The raw data is text, and this is first converted into a form suitable for learning by creating a dictionary of terms from all the documents in the training corpus and making a numeric

attribute for each term using Weka's unsupervised attribute filter *StringToWordVector*. There is also the class attribute, which gives the document's label.

Data with String Attributes

The *StringToWordVector* filter assumes that the document text is stored in an attribute of type *String*—a nominal attribute without a prespecified set of values. In the filtered data, this is replaced by a fixed set of numeric attributes, and the class attribute is put at the beginning, as the first attribute.

To perform document classification, first create an ARFF file with a string attribute that holds the document's text—declared in the header of the ARFF file using *@attribute document string*, where *document* is the name of the attribute. A nominal attribute is also needed to hold the document's classification.

Exercise 17.5.1. Make an ARFF file from the labeled mini-documents in Table 17.4 and run *StringToWordVector* with default options on this data. How many attributes are generated? Now change the value of the option *minTermFreq* to 2. What attributes are generated now?

Exercise 17.5.2. Build a *J48* decision tree from the last version of the data you generated.

Exercise 17.5.3. Classify the new documents in Table 17.5 based on the decision tree generated from the documents in Table 17.4. To apply the same

Table 17.4 Training Documents

| Document Text | Classification |
|--|----------------|
| The price of crude oil has increased significantly | yes |
| Demand for crude oil outstrips supply | yes |
| Some people do not like the flavor of olive oil | no |
| The food was very oily | no |
| Crude oil is in short supply | yes |
| Use a bit of cooking oil in the frying pan | no |

Table 17.5 Test Documents

| Document Text | Classification |
|--|----------------|
| Oil platforms extract crude oil | unknown |
| Canola oil is supposed to be healthy | unknown |
| Iraq has significant oil reserves | unknown |
| There are different types of cooking oil | unknown |

filter to both training and test documents, use *FilteredClassifier*, specifying the *StringToWordVector* filter and *J48* as the base classifier. Create an ARFF file from Table 17.5, using question marks for the missing class labels. Configure *FilteredClassifier* using default options for *StringToWordVector* and *J48*, and specify your new ARFF file as the test set. Make sure that you select *Output predictions* under *More options* in the Classify panel. Look at the model and the predictions it generates, and verify that they are consistent. What are the predictions?

Classifying Actual Documents

A standard collection of newswire articles is widely used for evaluating document classifiers. *ReutersCorn-train.arff* and *ReutersGrain-train.arff* are training sets derived from this collection; *ReutersCorn-test.arff* and *ReutersGrain-test.arff* are corresponding test sets. The actual documents in the corn and grain data are the same; only the labels differ. In the first dataset, articles concerning corn-related issues have a class value of 1 and the others have 0; the aim is to build a classifier that identifies “corny” articles. In the second, the labeling is performed with respect to grain-related issues; the aim is to identify “grainy” articles.

Exercise 17.5.4. Build classifiers for the two training sets by applying *FilteredClassifier* with *StringToWordVector* using (1) *J48* and (2) *NaiveBayesMultinomial*, evaluating them on the corresponding test set in each case. What percentage of correct classifications is obtained in the four scenarios? Based on the results, which classifier would you choose?

Other evaluation metrics are used for document classification besides the percentage of correct classifications: They are tabulated under *Detailed Accuracy By Class* in the Classifier Output area—the number of true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN). The statistics output by Weka are computed as specified in Table 5.7; the *F*-measure is mentioned in Section 5.7 (page 175).

Exercise 17.5.5. Based on the formulas in Table 5.7, what are the best possible values for each of the output statistics? Describe when these values are attained.

The Classifier Output also gives the ROC area (also known as AUC), which, as explained in Section 5.7 (page 177), is the probability that a randomly chosen positive instance in the test data is ranked above a randomly chosen negative instance, based on the ranking produced by the classifier. The best outcome is that all positive examples are ranked above all negative examples, in which case the AUC is 1. In the worst case it is 0. In the case where the ranking is essentially random, the AUC is 0.5, and if it is significantly less than this the classifier has performed anti-learning!

Exercise 17.5.6. Which of the two classifiers used above produces the best AUC for the two Reuters datasets? Compare this to the outcome for percent correct. What do the different outcomes mean?

The ROC curves discussed in Section 5.7 (page 172) can be generated by right-clicking on an entry in the result list and selecting *Visualize threshold curve*. This gives a plot with FP Rate on the x -axis and TP Rate on the y -axis. Depending on the classifier used, this plot can be quite smooth or it can be fairly irregular.

Exercise 17.5.7. For the Reuters dataset that produced the most extreme difference in Exercise 17.5.6, look at the ROC curves for class 1. Make a very rough estimate of the area under each curve, and explain it in words.

Exercise 17.5.8. What does the ideal ROC curve corresponding to perfect performance look like?

Other types of threshold curves can be plotted, such as a precision–recall curve with Recall on the x -axis and Precision on the y -axis.

Exercise 17.5.9. Change the axes to obtain a precision–recall curve. What is the shape of the ideal precision–recall curve, corresponding to perfect performance?

Exploring the *StringToWordVector* Filter

By default, the *StringToWordVector* filter simply makes the attribute value in the transformed dataset 1 or 0 for all single-word terms, depending on whether the word appears in the document or not. However, as mentioned in Section 11.3 (page 439), there are many options:

- *outputWordCounts* causes actual word counts to be output.
- *IDFTransform* and *TFTransform*: When both are set to *true*, term frequencies are transformed into $TF \times IDF$ values.
- *stemmer* gives a choice of different word-stemming algorithms.
- *useStopList* lets you determine whether or not stopwords are deleted.
- *tokenizer* allows different tokenizers for generating terms, such as one that produces word n -grams instead of single words.

There are several other useful options. For more information, click on *More* in the Generic Object Editor window.

Exercise 17.5.10. Experiment with the options that are available. What options give a good AUC value for the two datasets above, using *NaiveBayesMultinomial* as the classifier?

Not all of the attributes (i.e., terms) are important when classifying documents. The reason is that many words are irrelevant for determining an article's topic. Weka's *AttributeSelectedClassifier*, using ranking with *InfoGainAttributeEval* and the *Ranker* search, can eliminate less useful attributes. As before, *FilteredClassifier* should be used to transform the data before passing it to *AttributeSelectedClassifier*.

Exercise 17.5.11. Experiment with this, using default options for *StringToWordVector* and *NaiveBayesMultinomial* as the classifier. Vary the number of the most informative attributes that are selected from the information gain–based ranking by changing the value of the *numToSelect* field in the *Ranker*. Record the AUC values you obtain. How many attributes give the best AUC for the two datasets discussed before? What are the best AUC values you managed to obtain?

17.6 MINING ASSOCIATION RULES

In order to get some experience with association rules, we work with *Apriori*, the algorithm described in Section 4.5 (page 144). As you will discover, it can be challenging to extract useful information using this algorithm.

Association-Rule Mining

To get a feel for how to apply *Apriori*, start by mining rules from the *weather.nominal.arff* data that was used in Section 17.1. Note that this algorithm expects data that is purely nominal: If present, numeric attributes must be discretized first. After loading the data in the Preprocess panel, click the *Start* button in the Associate panel to run *Apriori* with default options. It outputs 10 rules, ranked according to the confidence measure given in parentheses after each one (they are listed in Figure 11.16). As we explained in Chapter 11 (page 430), the number following a rule's antecedent shows how many instances satisfy the antecedent; the number following the conclusion shows how many instances satisfy the entire rule (this is the rule's "support"). Because both numbers are equal for all 10 rules, the confidence of every rule is exactly 1.

In practice, it can be tedious to find minimum support and confidence values that give satisfactory results. Consequently, as explained in Chapter 11, Weka's *Apriori* runs the basic algorithm several times. It uses the same user-specified minimum confidence value throughout, given by the *minMetric* parameter. The support level is expressed as a proportion of the total number of instances (14 in the case of the weather data), as a ratio between 0 and 1. The minimum support level starts at a certain value (*upperBoundMinSupport*, default 1.0). In each iteration the support is decreased by a fixed amount (*delta*, default 0.05, 5% of the instances) until either a certain number of rules has been generated (*numRules*, default 10 rules) or the support reaches a certain "minimum minimum" level (*lowerBoundMinSupport*,

default 0.1)—because rules are generally uninteresting if they apply to less than 10% of the dataset. These four values can all be specified by the user.

This sounds pretty complicated, so we will examine what happens on the weather data. The Associator output text area shows that the algorithm managed to generate 10 rules. This is based on a minimum confidence level of 0.9, which is the default and is also shown in the output. The *Number of cycles performed*, which is shown as 17, indicates that *Apriori* was actually run 17 times to generate these rules, with 17 different values for the minimum support. The final value, which corresponds to the output that was generated, is 0.15 (corresponding to $0.15 \times 14 \approx 2$ instances).

By looking at the options in the Generic Object Editor window, you can see that the initial value for the minimum support (*upperBoundMinSupport*) is 1 by default, and that delta is 0.05. Now, $1 - 17 \times 0.05 = 0.15$, so this explains why a minimum support value of 0.15 is reached after 17 iterations. Note that *upperBoundMinSupport* is decreased by delta *before* the basic *Apriori* algorithm is run for the first time.

The Associator output text area also shows how many frequent item sets were found, based on the last value of the minimum support that was tried (0.15 in this example). In this case, given a minimum support of two instances, there are 12 item sets of size 1, 47 item sets of size 2, 39 item sets of size 3, and six item sets of size 4. By setting *outputItemSets* to *true* before running the algorithm, all those different item sets and the number of instances that support them are shown. Try it out!

Exercise 17.6.1. Based on the output, what is the support for this item set?

outlook = rainy humidity = normal windy = FALSE play = yes

Exercise 17.6.2. Suppose you want to generate all rules with a certain confidence and minimum support. This can be done by choosing appropriate values for *minMetric*, *lowerBoundMinSupport*, and *numRules*. What is the total number of possible rules for the weather data for each combination of values in Table 17.6?

Table 17.6 Number of Rules for Different Values of Minimum Confidence and Support

| Minimum Confidence | Minimum Support | Number of Rules |
|--------------------|-----------------|-----------------|
| 0.9 | 0.3 | |
| 0.9 | 0.2 | |
| 0.9 | 0.1 | |
| 0.8 | 0.3 | |
| 0.8 | 0.2 | |
| 0.8 | 0.1 | |
| 0.7 | 0.3 | |
| 0.7 | 0.2 | |
| 0.7 | 0.1 | |

Apriori has some further parameters. If *significanceLevel* is set to a value between 0 and 1, the association rules are filtered based on a χ^2 test with the chosen significance level. However, applying a significance test in this context is problematic because of the multiple comparison problem: If a test is performed hundreds of times for hundreds of association rules, it is likely that significant effects will be found just by chance—that is, an association seems to be statistically significant when really it is not. Also, the χ^2 test is inaccurate for small sample sizes (in this context, small support values).

There are alternative measures for ranking rules. As well as confidence, *Apriori* supports *lift*, *leverage*, and *conviction*, which can be selected using *metricType*. More information is available by clicking *More* in the Generic Object Editor window.

Exercise 17.6.3. Run *Apriori* on the weather data with each of the four rule-ranking metrics, and default settings otherwise. What is the top-ranked rule that is output for each metric?

Mining a Real-World Dataset

Now consider a real-world dataset, *vote.arff*, which gives the votes of 435 U.S. congressmen on 16 key issues gathered in the mid-1980s, and also includes their party affiliation as a binary attribute. This is a purely nominal dataset with some missing values (corresponding to abstentions). It is normally treated as a classification problem, the task being to predict party affiliation based on voting patterns. However, association-rule mining can also be applied to this data to seek interesting associations. More information on the data appears in the comments in the ARFF file.

Exercise 17.6.4. Run *Apriori* on this data with default settings. Comment on the rules that are generated. Several of them are quite similar. How are their support and confidence values related?

Exercise 17.6.5. It is interesting to see that none of the rules in the default output involve *Class = republican*. Why do you think that is?

Market Basket Analysis

In Section 1.3 (page 26) we introduced market basket analysis—analyzing customer purchasing habits by seeking associations in the items they buy when visiting a store. To do market basket analysis in Weka, each transaction is coded as an instance of which the attributes represent the items in the store. Each attribute has only one value: If a particular transaction does not contain it (i.e., the customer did not buy that item), this is coded as a missing value.

Your job is to mine supermarket checkout data for associations. The data in *supermarket.arff* was collected from an actual New Zealand supermarket. Take a look at this file using a text editor to verify that you understand the structure. The

main point of this exercise is to show you how difficult it is to find any interesting patterns in this type of data!

Exercise 17.6.6. Experiment with *Apriori* and investigate the effect of the various parameters described before. Write a brief report on the main findings of your investigation.