

FIT3152 Data analytics. Tutorial 08:

Classification models

Pre-tutorial Activity

The “diamonds” data set comes packaged with ggplot2 and contains data about the price of diamonds as well as information on size and the 4 Cs affecting diamond price: carat (size), cut, colour and clarity.

1. Re-create the decision tree from week 07 pre-tutorial activity to predict the price level of diamonds using a 70%/30% split for training and testing datasets. Check prediction accuracy.
2. Perform a cross validation test and prune the decision tree based on the misclassification rate. Check the accuracy of prediction and comment on any improvements.
3. Using the pruned decision tree do prediction as probabilities and draw a ROC curve.
4. Create a lift chart using the pruned decision tree output.

Tutorial Activities

Topics Covered:

- Improving the decision tree by Cross Validation and Pruning
- Naïve Bayes classification
- Classifier performance evaluation using ROC charts and Lift

References: Introduction to Data Mining, Tan, Steinbach and Kumar (available from Hargrave-Andrew library); An Introduction to Statistical Learning with applications in R, 2nd Ed. (Springer Texts in Statistics), James, Witten, Hastie and Tibshirani, Chapter 8 (available on-line from Monash Library); Reference Manuals to each of the packages used (listed below), available from CRAN.

- 1 Work through the examples in the lecture slides. For the examples using R you will need to download and install the packages in the script below.

```
# set working directory to desktop
# setwd("~/Desktop")
# clean up the environment before starting
rm(list = ls())
install.packages("tree")
library(tree)
install.packages("e1071")
library(e1071)
install.packages("ROCR")
library(ROCR)
```

2 Naïve Bayes Classification

ID	Colour	Size	Teeth	IsFriendly
1	red	medium	yes	no
2	blue	big	no	yes
3	green	medium	no	no
4	green	small	yes	no
5	blue	big	yes	yes
6	blue	small	yes	yes
7	red	small	no	yes
8	red	medium	no	yes
9	blue	medium	yes	yes
10	green	small	no	no

Given the Aliens data in the table above, use Naïve Bayes classification to predict whether or not the following aliens are friendly.

ID	Colour	Size	Teeth	IsFriendly
11	red	big	yes	?
12	green	big	yes	?
13	blue	small	no	?

Calculate the classification confidence (probability) for each of the 3 cases to be classified as friendly.

3 Creating ROC Charts

The Aliens table below provides information about several aliens including a class attribute, IsFriendly. A decision tree has been used to classify the data and obtained confidence values for IsFriendly='yes'.

a) Using this information, create a ROC chart for the Alien data.

Remember ROC charts graph TPR vs FPR for varying confidence thresholds. You can use the tables below the data table to assist you with the calculations.

b) What does this ROC chart tell you about the classifier?

c) What is the AUC value and what does this tell you about the classifier?

ID	Colour	Size	Teeth	IsFriendly	Confidence (IsFriendly 'yes')
1	red	medium	yes	no	0.7
2	blue	big	no	yes	0.9
3	green	medium	no	no	0.4
4	green	small	yes	no	0.1
5	blue	big	yes	yes	0.9
6	blue	small	yes	yes	0.8
7	red	small	no	yes	0.3
8	red	medium	no	yes	0.6
9	blue	medium	yes	yes	0.7
10	green	small	no	no	0.6

T =	Predicted Class Labels			T =	Predicted Class Labels			T =	Predicted Class Labels		
Actual		Class = 1	Class = 0	Actual		Class = 1	Class = 0	Actual		Class = 1	Class = 0
Class	Class = 1			Class	Class = 1			Class	Class = 1		
Labels	Class = 0			Labels	Class = 0			Labels	Class = 0		
TPR =		FPR =		TPR =		FPR =		TPR =		FPR =	
T =	Predicted Class Labels			T =	Predicted Class Labels			T =	Predicted Class Labels		
Actual		Class = 1	Class = 0	Actual		Class = 1	Class = 0	Actual		Class = 1	Class = 0
Class	Class = 1			Class	Class = 1			Class	Class = 1		
Labels	Class = 0			Labels	Class = 0			Labels	Class = 0		
TPR =		FPR =		TPR =		FPR =		TPR =		FPR =	
T =	Predicted Class Labels			T =	Predicted Class Labels			T =	Predicted Class Labels		
Actual		Class = 1	Class = 0	Actual		Class = 1	Class = 0	Actual		Class = 1	Class = 0
Class	Class = 1			Class	Class = 1			Class	Class = 1		
Labels	Class = 0			Labels	Class = 0			Labels	Class = 0		
TPR =		FPR =		TPR =		FPR =		TPR =		FPR =	

IsFriendly	Confidence (IsFriendly 'yes')	Threshold	X Axis FPR = FP/(FP + TN)	Y Axis TPR = TP/(TP + FN)
no	0.1	0.1		
yes	0.3	0.2		
no	0.4	0.3		
yes	0.6	0.4		
no	0.6	0.5		
no	0.7	0.6		
yes	0.7	0.7		
yes	0.8	0.8		
yes	0.9	0.9		
yes	0.9	1.0		

4 Creating a Lift chart using the Aliens data

If you select an alien at random, what is the chance it will be friendly?

- Using the data provided, what is the value of the Lift if you choose an alien from the subset which the classifier is at least 80% confident of?
- Now sketch a Lift chart, using steps of 20% (i.e. 80%, 60%, 40%, 20%)

5 Fitting a tree, cross validation and pruning in R

The Zoo data set “zoo.data.csv” contains data relating to seven classes of animals. Construct a decision tree using 70% training and 30% test data. Note that you will need to make sure the model knows that “type” is a factor. Make a table showing the actual vs predicted classifications using your test data. Calculate the accuracy of your model (by calculating correct classifications/all classifications). Are there any classes that are particularly difficult to predict?

Using cross validation and the “cv.tree” function see whether or not you are able to create a more accurate tree by pruning. Make a prediction with the new tree and compare its accuracy with your original tree.

6 Implementing Naïve Bayes classification in R

Again, use the Zoo data set split into a 70% training set and 30% test set. Fit a Naïve Bayes classifier. Compare the difference in accuracy between using the decision tree classifier (in the previous question) and the Naïve Bayes classifier. Does Naïve Bayes perform better on some classes than others?

7 The Japanese credit data “JapaneseCredit.csv” has 690 records relating to credit card applications, and whether they were approved or not. All attribute names and values have been anonymised. Ref. <http://archive.ics.uci.edu/ml/datasets/Japanese+Credit+Screening>

There are a variety of attribute types: continuous, and nominal (some with a small number of different values, and some with large). There are also some missing values that have been indicated as NA. The class value indicates “+” or “-” indicating whether or not a credit applicant was approved.

(a) Split the data into a 70% training and a 30% test set and create a classification model using each of the following techniques in turn:

- Decision Tree
- Naïve Bayes

(b) Using the test data, classify each of the test cases into “+” or “-”. Create a confusion matrix for each and report the accuracy of each model.

(c) Now, calculate the confidence of predicting a “+” outcome for each of the test cases and construct an ROC curve for each classifier. You should be able to plot all the curves on the same axis. Use a different colour for each classifier.

What does the ROC curve tell you about each of the classifiers? Is there a single “best” classifier?

(d) Examining each of the models, determine the most important variables in predicting whether or not an applicant would be granted a loan.

8 Some of the performance measures mentioned or covered in the lecture and tutorial are: accuracy, precision, recall, true positive rate (TPR), false positive rate (FPR), sensitivity, specificity, AUC, lift.

(a) Write the formula or a simple explanation for how you calculate each.

(b) Write a simple explanation (or as simple as possible) of what each of these measures indicates about the classifier.