

Tutorial 8

DSA1101

Introduction to Data Science

October 26, 2018

Exercise 1. Receiver Operating Characteristic (ROC) Curve in R

In last week's lecture, we introduced the ROC curve which is a common tool to evaluate classifiers in terms of trade-off between True Positive Rate (TPR) and False Positive Rate (FPR) when the classification threshold varies. We have also studied area under the curve (AUC), which is calculated by measuring the area under the ROC curve. Higher AUC scores mean the classifier performs better. In this week's tutorial, we will illustrate how AUC can be computed with the R package 'ROCR'.

Consider the dataset 'bank-sample.csv' we discussed in the lectures. For this exercise, we will predict the binary outcome **subscribed** using the naïve Bayes classifier, and generate the ROC curve of the naïve Bayes classifier built on a training set of 2,000 instances and tested on a testing set of 100 instances. All the datasets for this exercise are posted under the folder 'Tutorial 8' in IVLE.

- (a) Load the training and testing bank sample datasets.

```
1 # training set
2 banktrain <- read.table("bank-sample.csv",header=TRUE,sep=",")
3 # drop a few columns
4 drops <- c("balance", "day", "campaign", "pdays", "previous", "month"
5           )
6 banktrain <- banktrain[,!(names(banktrain) %in% drops)]
7 # testing set
8 banktest <- read.table("bank-sample-test.csv",header=TRUE,sep=",")
9 banktest <- banktest[,!(names(banktest) %in% drops)]
```

- (b) Build the naïve Bayes classifier based on the training dataset, and perform prediction for the test dataset.

```
1 library(e1071)
2
3 # build the naive Bayes classifier
4 nb_model <- naiveBayes(subscribed~.,
5 data=banktrain)
6 # perform on the testing set
7 nb_prediction <- predict(nb_model,
8 # remove column "subscribed"
9 banktest[, -ncol(banktest)],
10 type='raw')
```

- (c) Plot the ROC curve for the naïve Bayes classifier.

```
1 library(ROCR)
2
3 score <- nb_prediction[, c("yes")]
4
5 actual_class <- banktest$subscribed == 'yes'
6 pred <- prediction(score, actual_class)
7
8 perf <- performance(pred, "tpr", "fpr")
9 plot(perf, lwd=2, xlab="False Positive Rate (FPR)",
10 ylab="True Positive Rate (TPR)")
11 abline(a=0, b=1, col="gray50", lty=3)
```

- (d) Compute AUC for the naïve Bayes classifier.

```
1 auc <- performance(pred, "auc")
2 auc <- unlist(slot(auc, "y.values"))
3 auc
```