

# Data Mining Lab 5: More Programming Concepts and Model Evaluation

## 1 Introduction

In this lab we are going to expand our repertoire of R programming skills and then look at model evaluation for the so-called ‘churn’ data set which was introduced in the last lab.

**Note** that §2 might feel rather disconnected from what we are doing in the course, but this is to keep the initial examples simple. We will apply the new material from §2 to the course material in §3. This lab contains fewer exercises, but each will require significantly more thought than previous labs.

If you need to refer to previous labs or to download the data set, they will be on the course labs website:

<http://www.maths.tcd.ie/~louis/DataMining/>

## 2 Further Programming

In the very first lab we refreshed the bare basics of R, but didn’t delve into any of the programmatic power it has. Thus far, we have basically given one-line-at-a-time instructions to be run (with some exceptions which you were not expected to follow at the time). Now we will look at control structures and functions which will enable you to build more powerful solutions.

### 2.1 Control Structures: ‘if’

The `if` statement provides a simple way to perform different actions depending upon whether some condition you specify is true or not. The basic structure is as follows:

```
if(<condition>) {  
    <do this if condition evaluates to TRUE>  
} else {  
    <do this if condition evaluates to FALSE>  
}
```

For example, the following generates a  $\text{Normal}(\mu = 0, \sigma^2 = 1)$  random number and prints to the screen which side of 0.25 it lies:

```
i = rnorm(1)  
i  
if(i > 0.25) {  
    cat("i is greater than a quarter")  
} else {  
    cat("i is less than or equal to a quarter")  
}
```

**Exercise:** Use R to simulate a fair coin toss. Your code should simply print “Heads” or “Tails” to the screen once every time you run it, each time with a probability  $p = \frac{1}{2}$  for a head or a tail.

## 2.2 Control Structures: ‘for’

A frequent issue you will face when doing any statistical coding is the desire to perform the same action repeatedly many times. To avoid having to copy-and-paste the section of code to repeat it (eg if we want to do something 1000 times!), we can simply use loops which will perform an action as often as we specify. The `for` loop is the simplest of these (there are also `while` and `repeat` loops). The structure is:

```
for(<variable> in <vector>) {  
  <do this, one time for each element in the vector, setting variable  
    to the value of that element each time>  
}
```

For example, the following repeatedly calculates the log of a value and prints it. In this case, we are evaluating the log of the numbers 1 to a set value (10 here). Notice we also keep track of the sum of all the calculations to date and print this after the loop.

```
total = 0  
upto = 10  
for(j in 1:upto) {  
  total = total + log(j)  
  cat("The natural log of", j, "is", log(j), "\n")  
}  
cat("The sum of the natural logs of 1 to", upto, "is", total)
```

Be sure to input and run this and ensure you understand it completely. Note:

- that the `"\n"` simply creates a new line so that the next output is on the next line.
- we could have typed `1:10` instead of `1:upto`, but this way we can change one thing and all the results adjust.
- here we used the variable `j` which had the value from each element of the vector we looped over. However, you don't have to ... if you ignore it you just get an identical section of code run as many times as there are elements in your vector.

**Exercise:** Using a `for` loop, simulate 100 tosses of a fair coin and count how many heads you get in a variable called `total` which you print out after the loop.

## 2.3 Functions

Functions allow us to bundle up code into easy to use commands that we can use without having to think about what's going on inside. We've used many many functions on the course so far, like `rpart()`, which were written by the creators of R. However, things which are specific to our own work can be made much easier if we also create functions rather than having reams of code in a file. Combining loops and functions can also give us powerful abstractions that make statistical work very easy. The structure is:

```
<name of function> = function(<arg1>, <arg2>) {  
  <do things here, with the variables arg1 and arg2 available for use>  
  <this last line is what the function will return as the result>  
}
```

and it is called like so

```
<result> = <name of function>(<variable to pass 1>, <variable to pass 2>)
```

This kind of meta-definition is getting horrible looking, but an example should clarify! The following draws the specified number of cards from a deck of 52 with replacement:

```
drawCards = function(number) {  
  faceValues = c(1:10, "Jack", "Queen", "King")  
  suits = c("Hearts", "Diamonds", "Spades", "Clubs")  
  f = sample(faceValues, number, replace=TRUE)  
  s = sample(suits, number, replace=TRUE)  
  paste(f, "of", s)  
}  
  
drawCards(5)
```

Note:

- `sample(v, n, replace=TRUE)` draws `n` elements with replacement from a vector `v` with equally likely probabilities.
- `paste()` combines all of the elements of its input into a single string, element by element.

So now, anytime I ever want to do something in R which involves drawing cards from a deck, all I need to remember and use is `drawCards()` and I don't clutter my code with the rest!

**Exercise:** Create a function called `tossCoin` which simulates coin tosses. It should take a single argument specifying how many coins to toss, and should return just the total number of heads from the simulation.

## 3 Telecom Churn Model Evaluation

We return now to the world of data mining and will put the above control structures and functions theory to use in evaluating the classification tree model of the telecom churn data.

### 3.1 Getting Setup

**Exercise:** Load the telecom churn data into a variable called `churn` in your workspace. (Hint: see lab 2, §2.1 if you've forgotten)

**Exercise:** Load the `rpart` package, which contains the tree building functions. (Hint: see lab 2, §3.1 if you've forgotten)

Now run the following code which will create our train/test split; fit a tree on the train data; and then calculate the “Yes” prediction probabilities for the test data.

```
> test_rows = sample.int(nrow(churn), nrow(churn)/3)  
> test = churn[test_rows,]  
> train = churn[-test_rows,]  
  
> fit = rpart(churn ~ ., data=train, parms=list(split="gini"))  
> fit  
  
> probs = predict(fit, test)[,2]  
> probs
```

### 3.2 Goal

The goal today is to evaluate the model which has been fit by R. Remember the cut-off probability ( $\alpha$ ) for splitting Yes/No is the primary thing over which we have control *after* the model is fit (we have CP and loss matrices specified *before*) and so that plays a key role in the two things we look at today: Confusion Charts (at different  $\alpha$ ) and Receiver Operating Curves (a curve created by varying  $\alpha$ ). R has functions to create both these, but as an exercise we will create them ourselves.

### 3.3 Confusion Charts

Confusion charts can help us see whether our classifier is making false positive or false negative errors more frequently for a given  $\alpha$ . We will need a function which takes (i) our prediction probabilities for yes ( $\mathbb{P}[\text{churn} = \text{Yes} | \mathbf{x}]$ ); (ii) our chosen cut-off  $\alpha$ ; and which then returns actual Yes/No predictions. Below is a template with one omission.

```
makePrediction = function(probYes, cutoff) {  
  prediction = vector(length=length(probYes))  
  for(i in 1:length(prediction)) {  
    if(???) {  
      prediction[i] = "Yes"  
    } else {  
      prediction[i] = "No"  
    }  
  }  
  factor(prediction, levels=c("No", "Yes"))  
}
```

So, `probYes` will be a vector of probabilities and `cutoff` will be a number between 0 and 1 which is  $\alpha$ .

**Exercise:** Understand what is going on and so figure out what ??? on the fourth line should be.

Now we can compute the confusion table for any  $\alpha$  we choose. Here's for  $\alpha = 0.4$  as an example:

```
> pred = makePrediction(probs, 0.4)  
> table(pred, test$churn, dnn=c("Predicted", "Actual"))
```

### 3.4 ROC

Having created a function to calculate the predictions for any given  $\alpha$ , we can now use this from within *another* function which calculates it over a range in order to calculate and draw a ROC. We will design our function to take the model we fitted and the test set, then draw the ROC for it by computing the true and false positive rates over a range of 100 values between 0 and 1.

```
drawROC = function(fit, test) {  
  alpha = seq(0, 1, length.out=100)  
  
  tp = vector(length=100)  
  fp = vector(length=100)  
  
  probs = predict(fit, test)[,2]  
  for(i in 1:100) {  
    pred = makePrediction(probs, ???)  
    confusion = table(pred, test$churn, dnn=c("Predicted", "Actual"))  
    tp[i] = confusion[2,2]/(confusion[1,2]+confusion[2,2])  
    fp[i] = confusion[2,1]/(confusion[1,1]+confusion[2,1])  
  }  
  
  plot(???, ???, type="l")  
}
```

**Exercise:** Understand what is going on and so figure out what the three ???'s should be.

Once you've done that, run the following to see the ROC for the model we fitted to the train data in §3.1

```
> drawROC(fit, test)
```