

Data Science and R Programming

Rob Hayward

June 20, 2014

Introduction

This is an introduction to data science and R programming. There are two Coursera courses. The main classes can be found [Data science toolkit](#) and [R Programming](#).

0.1 Git

Some rules for adding files

- `git add .` adds all new files
- `git add -u` updates tracking for files that were changed or deleted
- `git add -A` all new files and all changes.

If you try to mix different classes in a vector, R will coerce to the lowest common denominator. For example, if you mix numeric and character, you will get two characters; if you try to mix logical and numeric, you will get numeric, if you try to mix logical and character, you will get character.

A list is a special object. It can change class.

```
x <- list(1, "a", TRUE, 1 + 4)
x
## [[1]]
## [1] 1
##
## [[2]]
## [1] "a"
##
```

```
## [[3]]
## [1] TRUE
##
## [[4]]
## [1] 5
```

The elements are indexed by double brackets.

0.2 Factors

Factors are integer vectors with a label. The level is determined by alphabetical order. Table can be used on factors. "Unclass" will strip the class and take the factor down to an integer. The levels can be set using "level" argument. Elements of an object can usually be named with "name". Lists and matrices can have names.

0.3 Subsetting

There are a number of ways to sub-set

- [always returns an object of the same class as the original. Can select more than one element.
- [[is used to extract elements of a list or data frame. Can only extract a single element and the class may not be an element or a dataframe.
- \$ used to extract named elements of a list or data frame. Same attributes as above.

Subsetting a matrix will usually return a vector. However, a matrix can be returned by setting drop = FALSE. This can be important if taking a column of a matrix. You will get a vector rather than a matrix unless you set drop = FALSE

```
x <- list(foo = 1:4, bar = 0.6)
x[1]

## $foo
## [1] 1 2 3 4

x[[1]]
```

```
## [1] 1 2 3 4

x$bar

## [1] 0.6

x[["bar"]]

## [1] 0.6
```

The `[[]]` can be used with computed indices.

```
x <- list(foo = 1:4, bar = 0.6, baz = "hello")
name <- "foo"
x[[name]]

## [1] 1 2 3 4

x$name

## NULL

x$foo

## [1] 1 2 3 4
```

```
x <- list(a = list(10, 12, 14), b = c(3.14, 2.18))
x[[c(1, 2)]]

## [1] 12

x[[1]][[3]]

## [1] 14

x[[c(2, 1)]]

## [1] 3.14
```

0.4 reading data

One way to speed up reading of large files is to read a small amount and to then create a vector with the column classes to input into the `colClasses` section.

```
initial <- read.table("datatable.txt", nrows = 100)
classes <- sapply(initial, class)
tabAll <- readtable("datatable.txt", colClasses = classes)
```

Interfaces with the outside world

- file connection to a file
- gzfile connection to a gzip compression
- bzfile connection to a bzip2 compression
- url connection to a web page

0.5 Control Structures

- if, else: testing a condition
- for: executing a loop a fixed number of times
- while: executing a loop *while* a condition is in place
- repeat execute an infinite loop
- break: break the execution of a loop
- next: skip an iteration of a loop
- return: exit function

Nested loops are inside each other. These are usually going to be used with things like matrices.

```
x <- matrix(1:6, 2, 3)
for(i in seq_len(nrow(x))) {
  for(j in seq_len(ncol(x))) {
    print(x[i,j])
  }
}
```

As another example that will keep flipping a coin until the score goes above or below a number.

```
z <- 5
while(z >= 3 && z <=10) {
  print(z)
  coin <- rbinm(1, 1, 0.5)

  if(coin ==1) { ## randomm walk
    z <- z + 1
  } else {
    z <- z -1
  }
}
```

Repeat could be used for running a loop until a certain tolerance is reached. For example,

```
x0 <- 1
tol <- 1e-8

repeat {
  x1 computeEstimate()

  if(abs(x1 - x0) < tol) {
    break
  } else {
    x0 <- x1
  }
}
```

This is used in optimisation algorithms. Usually good to add a limit on the number of repeats. For example a loop.

1 Functions

Create an object of a class function *Lazy Evaluation* means that the arguments of the function are only evaluated as they are needed. If an argument is not used or evaluated, there is no problem. However, if the function tries to perform some activity with an argument that does not exist (print for example), it will throw an error.

The `...` argument is used to repeated already existing argument defaults. For example, to extend a copy of the `plot` function,

```
myplot <- function(x, y, type = 'l', ...) {  
  plot(x, y, type = type, ...)  
}
```

They can also be used with generic functions to pass additional arguments to the function. All elements that are named after the `...` have to be named explicitly.

1.1 Scoping

How does R assign values to symbols?

For example,

```
lm <- function{x * x}  
lm
```

There is already a function called `lm` in R. R will search for each of the packages in a list. This can be found using the `search` function.

```
search()  
  
## [1] ".GlobalEnv"      "package:knitr"     "package:stats"  
## [4] "package:graphics" "package:grDevices" "package:utils"  
## [7] "package:datasets" "package:methods"   "Autoloads"  
## [10] "package:base"
```

This will search the environments. The first search is the Global Environment is the workspace. If you have something in the Global Environment, it will take that first. In this case, it will find the recently defined `lm` function in the Global Environment and will use that instead of the one in `stats` package. The base package is always the last on the list.

When a package is loaded with the `library` function, that package gets put in the second place. R has separate namespaces for functions and other objects. However, in the global environment, it is only possible to have one object for each name.

Scoping rules deal with the way that values are assigned to variables. Remember that there are variables that are the arguments of the function and other variables. R uses *lexical* or *static* scoping rather than *dynamic* scoping. This is particularly useful for statistical calculations.

For example,

```
f <- function(x, y) {
  x^2 + y/z
}
```

There are two formal arguments x and y but it is not clear where z comes from. This is a free variable that was not defined. Scoping rules define how values are assigned to these *free variables*.

Scoping rules

- The values of the free variables are searched for in the environment in which the variable was defined.
- An environment is a collection of (symbol, value) pairs. For example, x is a symbol and 3.1 may be the value.
- Every environment has a parent. It is possible for an environment to have multiple children.
- The only environment without a parent is the *empty environment*
- A function and an environment creates a closure.
- Therefore, if the function is defined in the global environment, this will be the first place that a value for the free variable will be sought.
- If this does not work, the next stop is the parent environment. This is the next environment on the search list. If there is no value in the package, the search continues down until the empty environment is reached.
- If there is no value in all the environments, an error is returned.

This becomes important when functions are defined within other functions. Therefore, a function can return a function. In this case the function is defined within another function (not the global environment).

For example,

```
make.power <- function(n) {
  pow <- function(x) {
    x^n
  }
  pow
}
cube <- make.power(3)
```

```

square <- make.power(2)
cube(3)

## [1] 27

square(3)

## [1] 9

```

Here `n` is a free variable as it is not defined in the `pow` function. However, it is defined in the `make.power` function.

To explain the difference between lexical and dynamic scoping, look at the following.

```

y <- 10
f <- function(x) {
  y <- 2
  y^2 + g(x)
}
g(x) <- function(x) {
}

## Error: could not find function "g<-"

f(3)

## Error: could not find function "g"

```

In Lexical scoping, the value of `y` in the function `g` is looked up in the environment in which the function was defined, in this case the global environment, so the value of `y` is 10.

With dynamic scoping, the value of `y` is looked up in the environment from which the function was called (sometimes referred to as the *calling environment*). In that case the value of `y` would be 2.

if a function is called from the global environment, the global environment and the calling environment are the same so that Lexical scoping can appear to be dynamic scoping.

The consequences of Lexical scoping

- All objects in R must be stored in memory.
- All functions must carry a pointer to their respective defining environment.

1.2 Vectorisation

This can be a way to avoid looping. Vectorising makes things simple. It applies the function to each element of a vector. For example,

```
x <- 1:4
y <- 6:9
x + y

## [1]  7  9 11 13

x > 9

## [1] FALSE FALSE FALSE FALSE

y == 8

## [1] FALSE FALSE  TRUE FALSE

x * y

## [1]  6 14 24 36

x/y

## [1] 0.1667 0.2857 0.3750 0.4444
```

For matrix operation

```
x <- 1:4
y <- 6:9
x * y # element wise multiplication

## [1]  6 14 24 36

x/y # element wise multiplication

## [1] 0.1667 0.2857 0.3750 0.4444

# True matrix multiplication
x %*% y

##           [,1]
## [1,]      80
```

1.3 Dates and times

Times have two formats

- POSIXct large integer with the date and time
- POSIXlt a list with information on the day of the week, year etc

Generic functions `weekday`, `months`, `quarters` give the day, month name or quarter when applied to a date, `POSIXct` or `POSIXlt` function. Things can be moved back and forward between the `POSIXct` and `POSIXlt` functions.

```
x <- Sys.time()
x

## [1] "2014-06-20 14:28:01 BST"

p <- as.POSIXlt(x)
names(unclass(p))

## [1] "sec"    "min"    "hour"    "mday"    "mon"    "year"    "yday"    "isdst"

p$sec

## [1] 1.317
```

The `strptime` function will convert dates in character format into date or time objects.

```
datestring <- c("January 10, 2012 10:40", "December 9, 2011 09:10")
x <- strptime(datestring, "%B %d, %Y %H:%M")
x

## [1] "2012-01-10 10:40:00" "2011-12-09 09:10:00"

class(x)

## [1] "POSIXlt" "POSIXt"
```

It is possible to use regular operations on dates and times. However, different classes cannot be mixed. For example,

```
x <- as.Date("2012-01-01")
y <- strptime("9 Jan 2011 11:34.21", "%d %b %Y %H:%S")
x - y

## Warning: Incompatible methods ("-.Date", "-.POSIXt") for "-"
## Error: non-numeric argument to binary operator

x <- as.POSIXlt(x)
x - y

## Time difference of 356.5 days
```

These operators keep track of very difficult things like leap years, leap seconds daylight savings and time zones. For example,

```
x <- as.Date("2012-03-01")
y <- as.Date("2012-02-28")
x - y

## Time difference of 2 days

x <- as.POSIXct(x)
y <- as.POSIXct(y, tz = "GMT")
y - x

## Time difference of -2 days
```

The first shows the leap year and the second shows the difference in timezones.

1.4 Loop functions

There are a number of loop functions in R

- lapply Look over list and evaluate a function for each element.
- sapply Same but simplify the result
- apply Margins of an array
- {tapply Apply over a subset of a vector
- mapply Multivariate version of tapply

Also `split` is useful.

These make a lot of use of *anonymous functions*. These are functions that do not have names. For example, write a function to extract the first column.

```
x <- list(a = matrix(1:4, 2, 2), b = matrix(1:6, 3, 2))
lapply(x, function(elt) elt[, 1])

## $a
## [1] 1 2
##
## $b
## [1] 1 2 3
```

`sapply` will simplify. For example, it may return a vector if there is a list of single objects or it may return a matrix.

Special functions

- `rowSums` – `apply(x, 1, sum)`
- `rowMean` = `apply(x, 1, mean)`
- `colSums` = `apply(x, 2, sum)`
- `colMean` = `apply(x, 2, mean)`

This also works on an array.

```
a <- array(rnorm(2 * 2 * 10), c(2, 2, 10))
apply(a, c(1, 2), mean)

##           [,1]      [,2]
## [1,] -0.5248 -0.48346
## [2,]  0.2254 -0.03591

rowMeans(a, dim = 2)

##           [,1]      [,2]
## [1,] -0.5248 -0.48346
## [2,]  0.2254 -0.03591
```

There are 10 2×2 matrices and we want to take the mean of the matrices. It will keep the first two dimensions and average over the third dimension. Hence `c(1, 2)`

For tapply a function can be applied across a subset of a vector. For example, here there are three different types of random variable and they are associated with factors one to three with the gl function.

```
x <- c(rnorm(10), runif(10), rnorm(10, 1))
f <- gl(3, 10)
f

## [1] 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 3 3 3
## Levels: 1 2 3

tapply(x, f, mean)

##      1      2      3
## 0.5283 0.4080 0.4635
```

split can be used to split a vector into parts that are identified by a factor. It will return a list. For example,

```
s <- split(airquality, airquality$Month)
sapply(s, function(x) colMeans(x[, c("Ozone", "Solar.R", "Wind")], na.rm = TRUE))

##           5           6           7           8           9
## Ozone      23.62    29.44    59.115    59.962    31.45
## Solar.R  181.30   190.17   216.484   171.857   167.43
## Wind       11.62    10.27     8.942     8.794    10.18
```

Split can also be done over multiple factors. For example,

```
x <- rnorm(10)
f1 <- gl(2, 5)
f2 <- gl(5, 2)
f1

## [1] 1 1 1 1 1 2 2 2 2 2
## Levels: 1 2

f2

## [1] 1 1 2 2 3 3 4 4 5 5
## Levels: 1 2 3 4 5
```

```

interaction(f1, f2)

## [1] 1.1 1.1 1.2 1.2 1.3 2.3 2.4 2.4 2.5 2.5
## Levels: 1.1 2.1 1.2 2.2 1.3 2.3 1.4 2.4 1.5 2.5

str(split(x, list(f1, f2)))

## List of 10
## $ 1.1: num [1:2] -0.689 0.483
## $ 2.1: num(0)
## $ 1.2: num [1:2] 0.192 0.191
## $ 2.2: num(0)
## $ 1.3: num -0.402
## $ 2.3: num 1.55
## $ 1.4: num(0)
## $ 2.4: num [1:2] -1.376 -0.408
## $ 1.5: num(0)
## $ 2.5: num [1:2] 0.694 0.359

str(split(x, list(f1, f2), drop = TRUE)) # drop empty levels

## List of 6
## $ 1.1: num [1:2] -0.689 0.483
## $ 1.2: num [1:2] 0.192 0.191
## $ 1.3: num -0.402
## $ 2.3: num 1.55
## $ 2.4: num [1:2] -1.376 -0.408
## $ 2.5: num [1:2] 0.694 0.359

```

`mapply` is a multivariate version. If there are two lists and the elements of each list go into the function, `mapply` may apply. In this case,

```

str(mapply)

## function (FUN, ..., MoreArgs = NULL, SIMPLIFY = TRUE, USE.NAMES = TRUE)

```

Where `FUN` is the function to be applied, `...` contains the arguments for it to be applied over; `MoreArgs` are more arguments to pass to the function.

```
list(rep(1, 4), rep(2, 3), rep(4, 1))
```

```
## [[1]]  
## [1] 1 1 1 1  
##  
## [[2]]  
## [1] 2 2 2  
##  
## [[3]]  
## [1] 4
```

```
mapply(rep, 1:4, 4:1)
```

```
## [[1]]  
## [1] 1 1 1 1  
##  
## [[2]]  
## [1] 2 2 2  
##  
## [[3]]  
## [1] 3 3  
##  
## [[4]]  
## [1] 4
```

1.5 Logical expressions

if we have two logical expressions A and B, $A \mid B$ will ask whether at least one is true (union); $A \& B$ will ask whether both are TRUE (ie intersection). $\neg A$ is the negation of A.

1.6 Character vectors

Use the paste function to put words that are in a character vector together or to join the elements of multiple character vectors. For example,

```
my_name <- c("My name", "is", "Rob")  
paste(my_name, collapse = " ") # the collapse is the  
  
## [1] "My name is Rob"
```

```
# space between.
```

Alternatively,

```
paste("Hello", "world!", sep = " ") # now sep is the separate.  
## [1] "Hello world!"
```

1.7 Debugging

Something is wrong

- message: notification/diagnostic
- warning: indication that something is wrong, but not fatally.
- error: a fatal problem has occurred
- condition: these are all conditions that can be created in the programme

`invisible` is a function that prevents auto-printing. If this is called on the return object, this will prevent it being returned.

Debugging tools in R

- `traceback` prints out the function call stack
- `debug` flags a function for debug mode which allows step through.
- `browser` suspends the execution of the function and puts it in debug mode.
- `trace` allows you to insert debugging code
- `recover` allows modification of the error behaviour so that function call stack can be reviewed.

```
mean(f)  
  
## Warning: argument is not numeric or logical: returning NA  
## [1] NA  
  
traceback()  
  
## No traceback available
```


traceback will list that past error.

debug will open a workspace within a workspace. This environment is the function environment. By passing `n` it will move to the next line until the error is reached. The debugger can build like a nest on top of another debugger. You can call the debug function on functions within the function.

recover will provide options for going into the code to find out what went wrong.