

Syntax Stuff (from coursera)

mf

November 3, 2017

Source of this material (and further reading)

Coursera class “R programming” from Johns Hopkins University, by Roger Peng, Jeff Leek, and Brian Caffo.

<https://www.coursera.org/learn/r-programming/>

All their materials are on github at

https://github.com/DataScienceSpecialization/courses/tree/master/02_RProgramming

Objects

R has five basic or “atomic” classes of objects:

- ▶ character
- ▶ numeric (real numbers)
- ▶ integer
- ▶ complex
- ▶ logical (True/False)

The most basic object is a vector

- ▶ A vector can only contain objects of the same class
- ▶ BUT: The one exception is a *list*, which is represented as a vector but can contain objects of different classes (indeed, that’s usually why we use them) —

Attributes

R objects can have attributes - names, dimnames

- ▶ dimensions (e.g. matrices, arrays)
- ▶ class
- ▶ length
- ▶ other user-defined attributes/metadata

Attributes of an object can be accessed using the `attributes()` function.

Creating Vectors

The `c()` function can be used to create vectors of objects.

```
> x <- c(0.5, 0.6)      ## numeric
> x <- c(TRUE, FALSE)   ## logical
> x <- c(T, F)          ## logical
> x <- c("a", "b", "c") ## character
> x <- 9:29              ## integer
> x <- c(1+0i, 2+4i)     ## complex
```

If you mix objects, R will force them to all be the same class (usually strings)

Matrices

Matrices are vectors with a *dimension* attribute. The dimension attribute is itself an integer vector of length 2 (nrow, ncol)

```
m <- matrix(nrow = 2, ncol = 3)
```

```
m
```

```
##      [,1] [,2] [,3]
## [1,]   NA   NA   NA
## [2,]   NA   NA   NA
```

```
dim(m)
```

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

```
attributes(m)
```

```
## $dim
```

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

cbind-ing and rbind-ing

Matrices can be created by *column-binding* or *row-binding* with `cbind()` and `rbind()`.

```
> x <- 1:3
> y <- 10:12
> cbind(x, y)
      x  y
[1,] 1 10
[2,] 2 11
[3,] 3 12
> rbind(x, y)
  [,1] [,2] [,3]
x    1    2    3
y   10   11   12
```

Lists

Lists are a special type of vector that can contain elements of different classes. Lists are a very important data type in R and you should get to know them well.

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

[[2]]
[1] "a"

[[3]]
[1] TRUE

[[4]]
[1] 1+4i
```


Missing Values

Missing values are denoted by NA or NaN for undefined mathematical operations.

- ▶ `is.na()` is used to test objects if they are NA
- ▶ `is.nan()` is used to test for NaN (Not A Number (like divided by 0 etc.))
- ▶ NA values have a class also, so there are integer NA, character NA, etc.
- ▶ A NaN value is also NA but the converse is not true

Missing Values

```
> x <- c(1, 2, NA, 10, 3)
> is.na(x)
[1] FALSE FALSE  TRUE FALSE FALSE
> is.nan(x)
[1] FALSE FALSE FALSE FALSE FALSE
```

Data Frames

Data frames are used to store tabular data

- ▶ They are represented as a special type of list where every element of the list has to have the same length
- ▶ Each element of the list can be thought of as a column and the length of each element of the list is the number of rows
- ▶ Unlike matrices, data frames can store different classes of objects in each column (just like lists); matrices must have every element be the same class
- ▶ Data frames also have a special attribute called `row.names`
- ▶ Data frames are usually created by calling `read.table()` or `read.csv()`
- ▶ Can be converted to a matrix by calling `data.matrix()`

Data Frames

```
> x <- data.frame(foo = 1:4, bar = c(T, T, F, F))
> x
  foo  bar
1   1 TRUE
2   2 TRUE
3   3 FALSE
4   4 FALSE
> nrow(x)
[1] 4
> ncol(x)
[1] 2
```

Reading Data

There are a few principal functions reading data into R.

- `read.table`, `read.csv`, for reading tabular data
- `readLines`, for reading lines of a text file
- `source`, for reading in R code files (inverse of `dump`)
- `dget`, for reading in R code files (inverse of `dput`)
- `load`, for reading in saved workspaces
- `unserialize`, for reading single R objects in binary form

NOTE: help on any function can be looked up with

?name_of_function

Writing Data

There are analogous functions for writing data to files

- write.table
- writeLines
- dump
- dput
- save
- serialize

Reading Data Files with `read.table`

The `read.table` function is one of the most commonly used functions for reading data. It has a few important arguments:

- `file`, the name of a file, or a connection
- `header`, logical indicating if the file has a header line
- `sep`, a string indicating how the columns are separated
- `colClasses`, a character vector indicating the class of each column in the dataset
- `nrows`, the number of rows in the dataset
- `comment.char`, a character string indicating the comment character
- `skip`, the number of lines to skip from the beginning
- `stringsAsFactors`, should character variables be coded as factors?

Interfaces to the Outside World

Data are read in using *connection* interfaces. Connections can be made to files (most common) or to other more exotic things.

- ▶ `file`, opens a connection to a file
- ▶ `gzfile`, opens a connection to a file compressed with gzip
- ▶ `url`, opens a connection to a webpage

In general, connections are powerful tools that let you navigate files or other external objects. In practice, we often don't need to deal with the connection interface directly.

```
con <- file("foo.txt", "r")  
data <- read.csv(con)  
close(con)
```

is the same as

```
data <- read.csv("foo.txt")
```


Reading Lines of a Text File

readLines can be useful for reading in lines of webpages

```
## This might take time
con <- url("http://www.ifpri.org", "r")
x <- readLines(con)
head(x,10)
close(con)
```

Subsetting

There are a number of operators that can be used to extract subsets of R objects.

- ▶ `[]` always returns an object of the same class as the original; can be used to select more than one element
- ▶ `[[` is used to extract elements of a list or a data frame; it can only be used to extract a single element and the class of the returned object will not necessarily be a list or data frame
- ▶ `$` is used to extract elements of a list or data frame by name; semantics are similar to that of `[[`.

Subsetting 2

```
> x <- c("a", "b", "c", "c", "d", "a")
```

Try all of the following:

```
> x[1]
```

```
> x[2]
```

```
> x[1:4]
```

```
> x[x > "a"]
```

and some more:

```
> u <- x > "a"
```

```
> u
```

```
> x[u]
```

Subsetting a Matrix

(This is starting to be useful for data manipulation)

Matrices can be subsetting in the usual way with (i,j) type indices.

```
> x <- matrix(1:6, 2, 3)
```

```
> x[1, 2]
```

```
> x[2, 1]
```

Missing index means the whole row/column:

```
> x[1, ]
```

```
> x[, 2]
```

Subsetting Lists

```
> 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
> x["bar"]
$bar
[1] 0.6
```

Removing NA Values

A common task is to remove missing values (NAs).

```
> x <- c(1, 2, NA, 4, NA, 5)
> bad <- is.na(x)
> x[!bad]
[1] 1 2 4 5
```

```
> x <- c(1, 2, NA, 4, NA, 5)
> y <- c("a", "b", NA, "d", NA, "f")
> good <- complete.cases(x, y)
> good
[1] TRUE TRUE FALSE TRUE FALSE TRUE
> x[good]
[1] 1 2 4 5
> y[good]
[1] "a" "b" "d" "f"
```

Removing NA Values

```
> airquality[1:6, ]
  Ozone Solar.R Wind Temp Month Day
1    41     190  7.4   67     5   1
2    36     118  8.0   72     5   2
3    12     149 12.6   74     5   3
4    18     313 11.5   62     5   4
5    NA      NA 14.3   56     5   5
6    28      NA 14.9   66     5   6
```

```
> good <- complete.cases(airquality)
> airquality[good, ][1:6, ]
  Ozone Solar.R Wind Temp Month Day
1    41     190  7.4   67     5   1
2    36     118  8.0   72     5   2
3    12     149 12.6   74     5   3
4    18     313 11.5   62     5   4
7    23     299  8.6   65     5   7
```

How R does looping

`apply` is used to evaluate a function (often an anonymous one) over the margins of an array.

- ▶ It is most often used to apply a function to the rows or columns of a matrix
- ▶ It can be used with general arrays, e.g. taking the average of an array of matrices
- ▶ It is not really faster than writing a loop, but it works in one line!

`apply(data, dimension, function)`

For sums and means of matrix dimensions, we have some shortcuts.

- `rowSums = apply(x, 1, sum)`
- `rowMeans = apply(x, 1, mean)`
- `colSums = apply(x, 2, sum)`
- `colMeans = apply(x, 2, mean)`

Other “apply-type” functions:

- ▶ `lapply`: Loop over a list and evaluate a function on each element
- ▶ `sapply`: Same as `lapply` but try to simplify the result
- ▶ `apply`: Apply a function over the margins of an array
- ▶ `tapply`: Apply a function over subsets of a vector
- ▶ `mapply`: Multivariate version of `lapply`

For, while, etc.

Those also exist, but the apply functions are what is new for us stata people.

```
for (year in c(2010,2011,2012,2013,2014,2015)){  
  print(paste("The year is", year))  
}
```

```
## [1] "The year is 2010"  
## [1] "The year is 2011"  
## [1] "The year is 2012"  
## [1] "The year is 2013"  
## [1] "The year is 2014"  
## [1] "The year is 2015"
```

Further training

These slides were cherry-picked material from a 4-week course.

Take it if you want more practice:

Coursera class “R programming” from Johns Hopkins University, by Roger Peng, Jeff Leek, and Brian Caffo.

<https://www.coursera.org/learn/r-programming/>

Or just read their slides

https://github.com/DataScienceSpecialization/courses/tree/master/02_RProgramming