# R programming

March 22, 2016

# Contents

# 1   Control flow

In the normal course of events, the statements in an R program are executed sequentially from the top of the program to the bottom. But there are times that you'll want to execute some statements repetitively, while only executing other statements if certain conditions are met.

For the syntax examples throughout this section, keep the following in mind:

- statement is a single R statement or a compound statement (a group of R statements enclosed in curly braces { } and separated by semicolons).

- cond is an expression that resolves to true or false.

- expr is a statement that evaluates to a number or character string.

- seq is a sequence of numbers or character strings.

Control structures in R allow you to control the flow of execution of the program, depending on runtime conditions. Common structures are

- if, else: testing a condition

- for: execute a loop a fixed number of times

- while: execute a loop while a condition is true

- repeat: execute an infinite loop

- break: break the execution of a loop

- next: skip an interation of a loop

- return: exit a function

Most control structures are not used in interactive sessions, but rather when writing functions or longer expresisons.

## 1.1   Repetition and looping

Looping constructs repetitively execute a statement or series of statements until a condition isn't true. These include the for and while structures.

### FOR

The for loop executes a statement repetitively until a variable's value is no longer contained in the sequence seq. The syntax is

`for (var in seq) statement`

A.

```
x <- c("a", "b", "c", "d")
for(i in 1:4) {
print(x[i])
}

## [1] "a"
## [1] "b"
## [1] "c"
## [1] "d"
```

```r
for(i in seq_along(x)) {
print(x[i])
}

## [1] "a"
## [1] "b"
## [1] "c"
## [1] "d"

for(letter in x) {
print(letter)
}

## [1] "a"
## [1] "b"
## [1] "c"
## [1] "d"

for(i in 1:4) print(x[i])

## [1] "a"
## [1] "b"
## [1] "c"
## [1] "d"
```

B.

```r
iTotal <-  0
for(i in 1:100)
{
    iTotal <- iTotal + i
}
cat("Sum of 1-100:",iTotal,"\n",sep="")

## Sum of 1-100:5050
```

C.

```r
szSymbols <- c("MSFT","GOOG","AAPL","INTL","ORCL","SYMC")
for(SymbolName in szSymbols)
{
    cat(SymbolName,"\n",sep="")
}

## MSFT
## GOOG
## AAPL
## INTL
## ORCL
## SYMC
```

D.

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

##      [,1] [,2] [,3]
## [1,]    1    3    5
## [2,]    2    4    6

for(i in seq_len(nrow(x))) {
        for(j in seq_len(ncol(x))) {
                print(x[i, j])
        }
}

## [1] 1
## [1] 3
## [1] 5
## [1] 2
## [1] 4
## [1] 6
```

## WHILE

A while loop executes a statement repetitively until the condition is no longer true. The syntax is

`while (cond) statement`

A.

```
count <- 0
while(count < 10) {
        print(count)
        count <- count + 1
}

## [1] 0
## [1] 1
## [1] 2
## [1] 3
## [1] 4
## [1] 5
## [1] 6
## [1] 7
## [1] 8
## [1] 9
```

B.

```
i <- 1
iTotal <- 0
while(i <= 100)
{
        iTotal <- iTotal + i
```

```
      i <- i + 1
}
cat("Total:",iTotal,"\n",sep="")

## Total:5050
```

C.

```
z <- 5
while(z >= 3 && z <= 10) {
        print(z)
        coin <- rbinom(1, 1, 0.5)
        if(coin == 1) { ## random walk
                z <- z + 1
        } else {
                z <- z - 1
        }
}

## [1] 5
## [1] 6
## [1] 7
## [1] 8
## [1] 7
## [1] 8
## [1] 9
## [1] 10
```

## REPEAT

```
i <- 1
iTotal <- 0
repeat
{
iTotal <- iTotal + i
i <- i + 1
if(i <= 100) next else break
}
cat("Total:",iTotal,"\n",sep="")

## Total:5050
```

Looping in R can be inefficient and time consuming when you're processing the rows or columns of large datasets. Whenever possible, it's better to use R's builtin numerical and character functions in conjunction with the apply family of functions.

## 1.2   Conditional execution

In conditional execution, a statement or statements are only executed if a specified condition is met. These constructs include if-else, ifelse, and switch.

## IF-ELSE

The if-else control structure executes a statement if a given condition is true. Optionally, a different statement is executed if the condition is false. The syntax is

```
#1
if (cond) statement
#2
if (cond) statement1 else statement2
#3
if(<condition1>) {
## do something
} else if(<condition2>) {
## do something different
}else {
## do something different
}
```

```
a <- 3
if( a == 1)
{
    print("a == 1")
}else
{
    print("a != 1")
}

## [1] "a != 1"
```

```
grade <- as.factor(c("grade1","grade2"))
if (!is.factor(grade))
{
    grade <- as.factor(grade)
}else
{
    print("Grade already is a factor")
}

## [1] "Grade already is a factor"
```

```
a <- 4
if( a == 1)
{
print("a == 1")
}else if( a == 2)
{
print("a == 2")
}else
{
print("Not 1 & 2")
}
```

```
## [1] "Not 1 & 2"
```

### IFELSE

The ifelse construct is a compact and vectorized version of the if-else construct . The syntax is

```
ifelse(cond, statement1, statement2)
```

```
x <- matrix(1:6, 2, 3)
ifelse(x >= 0, sqrt(x), NA)

##           [,1]     [,2]     [,3]
## [1,] 1.000000 1.732051 2.236068
## [2,] 1.414214 2.000000 2.449490
```

Use ifelse when you want to take a binary action or when you want to input and output vectors from the construct.

### SWITCH

switch chooses statements based on the value of an expression.

```
n <- 1
switch(n,
    print("option1"),
    print("option2"),
    print("option3")
)

## [1] "option1"
```

### Next

Next is used to skip an iteration of a loop

```
for(i in 1:100) {
if(i <= 20) {
## Skip the first 20 iterations
next
}
## Do something here
}
```

## 2   User-written functions

One of R's greatest strengths is the user's ability to add functions. In fact, many of the functions in R are functions of existing functions. The structure of a function looks like this:

```
myfunction <- function(arg1, arg2, ... ){ statements return(object) }
```

Functions in R are "first class objects", which means that they can be treated much like any other R object. Importantly,

- Functions can be passed as arguments to other functions

- Functions can be nested, so that you can define a function inside of another function

- The return value of a function is the last expression in the function body to be evaluated.

**Function Arguments**

Functions have named arguments which potentially have default values.

- The formal arguments are the arguments included in the function definition

- The formals function returns a list of all the formal arguments of a function

- Not every function call in R makes use of all the formal arguments

- Function arguments can be missing or might have default values

Example A:

```
"%g%" <- function(x,y)
{
print(x+y)
print(x-y)
print(x*y)
print(x/y)
}
3%g%5

## [1] 8
## [1] -2
## [1] 15
## [1] 0.6
```

Example B:

```
columnmean <- function(y){
        nc <- ncol(y)
        means <- numeric(nc)
        for(i in 1:nc){
                means[i] <- mean(y[,i])
        }
        means
}
columnmean(airquality)

## [1]       NA       NA  9.957516 77.882353  6.993464 15.803922

columnmean <- function(y,removeNA=TRUE){
        nc <- ncol(y)
        means <- numeric(nc)
        for(i in 1:nc){
```

```
            means[i] <- mean(y[,i], na.rm = removeNA)
        }
        means
}
columnmean(airquality)

## [1]  42.129310 185.931507   9.957516  77.882353   6.993464  15.803922

columnmean(airquality,FALSE)

## [1]        NA        NA  9.957516 77.882353  6.993464 15.803922
```

## Argument Matching

Function arguments can also be partially matched, which is useful for interactive work. The order of operations when given an argument is

1. Check for exact match for a named argument
2. Check for a partial match
3. Check for a positional match

```
mydata <- rnorm(100)
sd(mydata)

## [1] 0.9467659

sd(x = mydata)

## [1] 0.9467659

sd(x = mydata, na.rm = FALSE)

## [1] 0.9467659

sd(na.rm = FALSE, x = mydata)

## [1] 0.9467659

sd(na.rm = FALSE, mydata)

## [1] 0.9467659
```

Even though it's legal, I don't recommend messing around with the order of the arguments too much, since it can lead to some confusion.

You can mix positional matching with matching by name. When an argument is matched by name, it is "taken out" of the argument list and the remaining unnamed arguments are matched in the order that they are listed in the function definition.

```
args(lm)

## function (formula, data, subset, weights, na.action, method = "qr",
##     model = TRUE, x = FALSE, y = FALSE, qr = TRUE, singular.ok = TRUE,
##     contrasts = NULL, offset, ...)
## NULL
```

The following two calls are equivalent.

```
lm(data = mydata, y ~ x, model = FALSE, 1:100)
lm(y ~ x, mydata, 1:100, model = FALSE)
```

Most of the time, named arguments are useful on the command line when you have a long argument list and you want to use the defaults for everything except for an argument near the end of the list.

## Return

Return signals that a function should exit and return a given value

```
x <- c(1,9,2,8,3,7)
y <- c(9,2,8,3,7,2)
pmax(x,y)

## [1] 9 9 8 8 7 7

parmax <- function (a,b){
c <- pmax(a,b)
median(c)}
parmax(x,y)

## [1] 8
```

If you want to return two or more variables from a function you should use return with a list containing the variables to be returned.

```
parboth <- function (a,b) {
c <- pmax(a,b)
d <- pmin(a,b)
answer <- list(median(c),median(d))
names(answer)[[1]] <- "median of the parallel maxima"
names(answer)[[2]] <- "median of the parallel minima"
return(answer) }
parboth(x,y)

## $`median of the parallel maxima`
## [1] 8
##
## $`median of the parallel minima`
## [1] 2
```

## The "..." argument

The ... argument indicate a variable number of arguments that are usually passed on to other functions. ... is often used when extending another function and you don't want to copy the entire argument list of the original function.

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

The ... argument is also necessary when the number of arguments passed to the function cannot be known in advance.

```
args(paste)

## function (..., sep = " ", collapse = NULL)
## NULL

args(cat)

## function (..., file = "", sep = " ", fill = FALSE, labels = NULL,
##      append = FALSE)
## NULL
```

One catch with ... is that any arguments that appear after ... on the argument list must be named explicitly and cannot be partially matched.

```
paste("a", "b", sep = ":")

## [1] "a:b"

paste("a", "b", se = ":")

## [1] "a b :"
```

## Programming tips

- Know exactly what you are trying to achieve.

- Keep it simple.

- Clever is good, but clear is better.

- Test each line as you go along, to make sure it does what you want it to do.

- Put plenty of comments in the code, using # for documentation.

- Use variable names and function names that are self-explanatory.

- Do not use attach in programs. Use with, or refer to variables within named dataframes.

- Try different ways of doing the same thing, and select the fastest method.

- Use indents (tabs) to improve clarity of loops and if statements.

- Build up the program from small, independently tested functions.

- Stop tinkering once it works effectively.

## A Diversion on Binding Values to Symbol

How does R know which value to assign to which symbol? When I type

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

How does R know what value to assign to the symbol lm? Why doesn't it give it the value of lm that is in the stats package? When R tries to bind a value to a symbol, it searches through a series of environments to find the appropriate value. When you are working on the command line and need to retrieve the value of an R object, the order is roughly The search list can be found by using the search function.

1. Search the global environment for a symbol name matching the one requested.

2. Search the namespaces of each of the packages on the search list.

```
search()
```

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

## Exercise

Question: We have two sample:
  A 79.98, 80.04, 80.02, 80.04, 80.03, 80.03, 80.04, 79.97, 80.05, 80.03, 80.02, 80.00, 80.02
  B 80.02, 79.94, 79.98, 79.97, 79.97, 80.03, 79.95, 79.97
Find T statistics.
  Hint: T statistics are

$$T = \frac{(\bar{X} - \bar{Y})}{S\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

  where

$$S = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

Answer:

```
twosam <- function(y1, y2){
   n1 <- length(y1); n2 <- length(y2)
   yb1 <- mean(y1); yb2 <- mean(y2)
   s1 <- var(y1); s2 <- var(y2)
   s <- ((n1-1)*s1 + (n2-1)*s2)/(n1+n2-2)
   tst <- (yb1 - yb2)/sqrt(s*(1/n1 + 1/n2))
   tst
}
A <- c(79.98, 80.04, 80.02, 80.04, 80.03, 80.03, 80.04,
       79.97, 80.05, 80.03, 80.02, 80.00, 80.02)
B <- c(80.02, 79.94, 79.98, 79.97, 79.97, 80.03, 79.95,
       79.97)
twosam(A,B)
```

# 3   *apply functions

## 3.1   apply

For sums and means of matrix dimensions, we have some shortcuts.  apply is used to a 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

```
str(apply)

## function (X, MARGIN, FUN, ...)
```

- X is an array

- MARGIN is an integer vector indicating which margins should be "retained".

- FUN is a function to be applied

- ... is for other arguments to be passed to FUN

```
x <- matrix(1:24,nrow=4)
x

##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,]    1    5    9   13   17   21
## [2,]    2    6   10   14   18   22
## [3,]    3    7   11   15   19   23
## [4,]    4    8   12   16   20   24

rowSums = apply(x, 1, sum)
rowSums

## [1] 66 72 78 84

rowMeans = apply(x, 1, mean)
rowMeans

## [1] 11 12 13 14

colSums = apply(x, 2, sum)
colSums

## [1] 10 26 42 58 74 90

colMeans = apply(x, 2, mean)
colMeans

## [1]  2.5  6.5 10.5 14.5 18.5 22.5
```

The apply function is used for applying functions to the rows or columns of matrices or dataframes. Quantiles of the rows of a matrix.

```
x <- matrix(rnorm(200), 20, 10)
#?quantile
apply(x, 1, quantile, probs = c(0.25, 0.75))

##            [,1]       [,2]       [,3]       [,4]       [,5]       [,6]
## 25% -0.0497953 -0.8022563 -0.9165523 -0.9827469 -0.5212556 -0.5021466
## 75%  0.7826299  0.3231492  0.6085441  0.2463501  0.1655246  0.5439360
##            [,7]        [,8]      [,9]      [,10]      [,11]      [,12]
## 25% -1.0290199 -0.97685650 -0.618691 -0.3365819 -0.6432999 -0.7652058
## 75%  0.2339169  0.05407604  1.053643  0.2937521  0.5984581  0.3417121
##           [,13]      [,14]      [,15]      [,16]      [,17]      [,18]
## 25% -0.1956111 -0.7860212 -1.0213136 -0.2800948 -0.7231161 -0.7557842
## 75%  1.1500474  0.6629060  0.3117813  0.6599279  0.9879331 -0.5087939
##           [,19]      [,20]
## 25% -0.7533123 -0.1003522
## 75%  0.5754819  0.7471505
```

Average matrix in an array

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

##            [,1]         [,2]
## [1,] -0.3606650 -0.049059415
## [2,] -0.2131403 -0.002567112

rowMeans(a, dims = 2)

##            [,1]         [,2]
## [1,] -0.3606650 -0.049059415
## [2,] -0.2131403 -0.002567112
```

## 3.2   mapply

mapply is a multivariate apply of sorts which applies a function in parallel over a set of arguments.

```
str(mapply)

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

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

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

- FUN is a function to apply

- ... contains arguments to apply over

- MoreArgs is a list of other arguments to FUN.

- SIMPLIFY indicates whether the result should be simplified.

```
noise <- function(n, mean, sd) rnorm(n, mean, sd)
noise(5,1,2)

## [1] -2.413031  2.775439  5.719322  3.968982  5.350133

noise(1:5,1:5,2)

## [1] 2.973930 5.866200 4.458258 3.386506 3.949374

mapply(noise, 1:5, 1:5, 2)

## [[1]]
## [1] -3.10921
##
## [[2]]
## [1] 1.155555 1.049458
##
## [[3]]
## [1] -1.128978  1.569946  5.601192
##
## [[4]]
## [1] 5.600113 4.867628 5.636744 2.620906
##
## [[5]]
## [1] 4.973062 5.285809 4.213198 6.402439 3.736101

list(noise(1, 1, 2), noise(2, 2, 2), noise(3, 3, 2), noise(4, 4, 2), noise(5, 5, 2))
```

```
## [[1]]
## [1] -2.495152
##
## [[2]]
## [1]  1.88575619 -0.08583168
##
## [[3]]
## [1] 2.0571313 2.8158007 0.9964285
##
## [[4]]
## [1] 2.516064 4.168613 5.295537 4.692418
##
## [[5]]
## [1] 6.666417 5.040478 6.835452 3.194508 3.342273
```

## 3.3   lapply() and sapply()

Each of the *apply functions will SPLIT up some data into smaller pieces, APPLY a function to each piece, then COMBINE the results. lapply takes three arguments: (1) a list X; (2) a function (or the name of a function) FUN; (3) other arguments via its ... argument. If X is not a list, it will be coerced to a list using as.list. lapply always returns a list, regardless of the class of the input.

```
#1
x <- list(a = 1:5, b = rnorm(10))
lapply(x, mean)

## $a
## [1] 3
##
## $b
## [1] 0.2325889

x <- list(a = 1:4, b = rnorm(10), c = rnorm(20, 1), d = rnorm(100, 5))
lapply(x, mean)

## $a
## [1] 2.5
##
## $b
## [1] 0.7871178
##
## $c
## [1] 0.7734139
##
## $d
## [1] 4.961637

#2
x <- 1:4
lapply(x, runif)
```

```
## [[1]]
## [1] 0.5654771
##
## [[2]]
## [1] 0.6108811 0.8745072
##
## [[3]]
## [1] 0.4228842 0.9224485 0.6929324
##
## [[4]]
## [1] 0.6767389 0.3495941 0.2747887 0.4976514

lapply(x, runif, min = 0, max = 10)

## [[1]]
## [1] 4.890143
##
## [[2]]
## [1] 5.951393 8.216313
##
## [[3]]
## [1] 6.273454 3.727089 2.918946
##
## [[4]]
## [1] 4.590575 9.170664 8.662025 8.199267

#3
x <- list(a = matrix(1:4, 2, 2), b = matrix(1:6, 3, 2))
x

## $a
##      [,1] [,2]
## [1,]    1    3
## [2,]    2    4
##
## $b
##      [,1] [,2]
## [1,]    1    4
## [2,]    2    5
## [3,]    3    6

lapply(x, function(elt) elt[,1])

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

This dataset contains details of various nations and their flags. More information may be found here: http://archive.ics.uci.edu/ml/datasets/Flags.

```
flags<-read.csv(file="C:\\Users\\XXXHHF\\Documents\\R\\workfile\\flag.csv",header=FALSE)
names(flags)<-c("name","landmass","zone","area","population","language",
"religion","bars","stripes","colours","red","green","blue","gold","white",
"black","orange","mainhue","circles","crosses","saltires","quarter",
"sunstars","crescent","triangle","icon","animate","text","topleft","botright")
```

```
load("C:/Users/XXXHHF/Documents/LYX/8. R programming/Data/flags.RData")
head(flags)
```

```
##              name landmass zone area population language religion bars
## 1     Afghanistan        5    1  648         16       10        2    0
## 2         Albania        3    1   29          3        6        6    0
## 3         Algeria        4    1 2388         20        8        2    2
## 4 American-Samoa        6    3    0          0        1        1    0
## 5         Andorra        3    1    0          0        6        0    3
## 6          Angola        4    2 1247          7       10        5    0
##    stripes colours red green blue gold white black orange mainhue circles
## 1        3       5   1     1    0    1     1     1      0   green       0
## 2        0       3   1     0    0    1     0     1      0     red       0
## 3        0       3   1     1    0    0     1     0      0   green       0
## 4        0       5   1     0    1    1     1     0      1    blue       0
## 5        0       3   1     0    1    1     0     0      0    gold       0
## 6        2       3   1     0    0    1     0     1      0     red       0
##    crosses saltires quarter sunstars crescent triangle icon animate text
## 1        0        0       0        1        0        0    1       0    0
## 2        0        0       0        1        0        0    0       1    0
## 3        0        0       0        1        1        0    0       0    0
## 4        0        0       0        0        0        1    1       1    0
## 5        0        0       0        0        0        0    0       0    0
## 6        0        0       0        1        0        0    1       0    0
##    topleft botright
## 1    black    green
## 2      red      red
## 3    green    white
## 4     blue      red
## 5     blue      red
## 6      red    black
```

```
dim(flags)
```

```
## [1] 194  30
```

```
class(flags)
```

```
## [1] "data.frame"
```

This tells us that there are 194 rows, or observations, and 30 columns, or variables. Each observation is a country and each variable describes some characteristic of that country or its flag. The lapply() function takes a list as input, applies a function to each element of the list, then returns a list of the same length as the original one.

```
cls_list <- lapply(flags, class)
#cls_list
class(cls_list)

## [1] "list"

as.character(cls_list)

##  [1] "factor"  "integer" "integer" "integer" "integer" "integer" "integer"
##  [8] "integer" "integer" "integer" "integer" "integer" "integer" "integer"
## [15] "integer" "integer" "integer" "factor"  "integer" "integer" "integer"
## [22] "integer" "integer" "integer" "integer" "integer" "integer" "integer"
## [29] "factor"  "factor"
```

As expected, we got a list of length 30 – one element for each variable/column. Sapply() allows you to automate this process by calling lapply() behind the scenes, but then attempting to simplify (hence the 's' in 'sapply') the result for you. Use sapply() the same way you used lapply() to get the class of each column of the flags dataset and store the result in cls_vect. If you need help, type ?sapply to bring up the documentation.

```
cls_vect <- sapply(flags, class)
cls_vect

##       name    landmass       zone        area  population    language
##   "factor"   "integer"  "integer"   "integer"   "integer"   "integer"
##    religion        bars     stripes     colours         red       green
##   "integer"   "integer"  "integer"   "integer"   "integer"   "integer"
##        blue        gold       white       black      orange     mainhue
##   "integer"   "integer"  "integer"   "integer"   "integer"    "factor"
##     circles     crosses    saltires     quarter    sunstars    crescent
##   "integer"   "integer"  "integer"   "integer"   "integer"   "integer"
##    triangle        icon     animate        text     topleft    botright
##   "integer"   "integer"  "integer"   "integer"    "factor"    "factor"

class(cls_vect)

## [1] "character"
```

Therefore, if we want to know the total number of countries (in our dataset) with, for example, the color orange on their flag, we can just add up all of the 1s and 0s in the 'orange' column.

```
sum(flags$orange)

## [1] 26
```

Now we want to repeat this operation for each of the colors recorded in the dataset.

```
flag_colors <- flags[, 11:17]
head(flag_colors)

##   red green blue gold white black orange
## 1   1     1    0    1     1     1      0
## 2   1     0    0    1     0     1      0
```

```
## 3   1    1    0    0    1    0    0
## 4   1    0    1    1    1    0    1
## 5   1    0    1    1    0    0    0
## 6   1    0    0    1    0    1    0
```

```
lapply(flag_colors, sum)
```

```
## $red
## [1] 153
##
## $green
## [1] 91
##
## $blue
## [1] 99
##
## $gold
## [1] 91
##
## $white
## [1] 146
##
## $black
## [1] 52
##
## $orange
## [1] 26
```

```
sapply(flag_colors, sum)
```

```
##    red  green   blue   gold  white  black orange
##    153     91     99     91    146     52     26
```

```
sapply(flag_colors, mean)
```

```
##       red     green      blue      gold     white     black    orange
## 0.7886598 0.4690722 0.5103093 0.4690722 0.7525773 0.2680412 0.1340206
```

This tells us that of the 194 flags in our dataset, 153 contain the color red, 91 contain green, 99 contain blue, and so on. In the examples we've looked at so far, sapply() has been able to simplify the result to vector. That's because each element of the list returned by lapply() was a vector of length one. Recall that sapply() instead returns a matrix when each element of the list returned by lapply() is a vector of the same length ($> 1$).

```
flag_shapes <- flags[, 19:23]
lshape <- lapply(flag_shapes, range)
mshape <- sapply(flag_shapes, range)
class(lshape)
```

```
## [1] "list"
```

```
class(mshape)
```

```
## [1] "matrix"
```

```
dim(mshape)
```

```
## [1] 2 5
```

When given a vector, the unique() function returns a vector with all duplicate elements removed. We want to know the unique values for each variable in the flags dataset.

```
unique(c(3, 4, 5, 5, 5, 6, 6))
```

```
## [1] 3 4 5 6
```

```
unique_vals <- lapply(flags, unique)
```

Since unique_vals is a list, you can use what you've learned to determine the length of each element of unique_vals (i.e. the number of unique values for each variable).

```
sapply(unique_vals, length)
```

```
##       name   landmass       zone       area population   language
##        194          6          4        136         48         10
##   religion       bars    stripes    colours        red      green
##          8          5         12          8          2          2
##       blue       gold      white      black     orange    mainhue
##          2          2          2          2          2          8
##    circles     crosses   saltires    quarter   sunstars   crescent
##          4          3          2          3         14          2
##   triangle       icon    animate       text    topleft   botright
##          2          2          2          2          7          8
```

## 3.4   vapply() and tapply()

Whereas sapply() tries to 'guess' the correct format of the result, vapply() allows you to specify it explicitly. If the result doesn't match the format you specify, vapply() will throw an error, causing the operation to stop.

```
vapply(flags, unique, numeric(1))
```

which says that you expect each element of the result to be a numeric vector of length 1.

```
cflags<-vapply(flags, class, character(1))
class(cflags)
```

```
## [1] "character"
```

The 'character(1)' argument tells R that we expect the class function to return a character vector of length 1 when applied to EACH column of the flags dataset.

You might think of vapply() as being 'safer' than sapply(), since it requires you to specify the format of the output in advance, instead of just allowing R to 'guess' what you wanted. In addition, vapply() may perform faster than sapply() for large datasets. However, when doing data analysis interactively (at the prompt), sapply() saves you some typing and will often be good enough.

As a data analyst, you'll often wish to split your data up into groups based on the value of some variable, then apply a function to the members of each group. The next function we'll look at, tapply(), does exactly that.

```
str(tapply)

## function (X, INDEX, FUN = NULL, ..., simplify = TRUE)
```

- X is a vector

- INDEX is a factor or a list of factors (or else they are coerced to factors)

- FUN is a function to be applied

- ... contains other arguments to be passed FUN

- simplify, should we simplify the result?

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

##  [1] -1.35786813 -0.14913200 -0.61935278  0.02564552 -0.94273179
##  [6]  0.29520945  0.03335189 -0.44569929 -0.19884647 -0.07787020
## [11]  0.07375128  0.06218643  0.65970996  0.76737767  0.59409123
## [16]  0.25572701  0.79466919  0.47844719  0.53884598  0.60606506
## [21]  0.68949798  1.27428513  1.44591688  1.62678667  1.58256284
## [26]  2.25619603  1.37249891  1.17428430  1.35961172 -0.30548013

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
## -0.3437294  0.4830871  1.2476160

tapply(x, f, range)

## $`1`
## [1] -1.3578681  0.2952095
##
## $`2`
## [1] 0.06218643 0.79466919
##
## $`3`
## [1] -0.3054801  2.2561960

tapply(x, f, mean, simplify = FALSE)

## $`1`
## [1] -0.3437294
##
## $`2`
## [1] 0.4830871
##
## $`3`
## [1] 1.247616
```

The 'landmass' variable in our dataset takes on integer values between 1 and 6, each of which represents a different part of the world. The 'animate' variable in our dataset takes the value 1 if a country's flag contains an animate image (e.g. an eagle, a tree, a human hand) and 0 otherwise. If you take the arithmetic mean of a bunch of 0s and 1s, you get the proportion of 1s.

```
table(flags$landmass)

##
##  1  2  3  4  5  6
## 31 17 35 52 39 20

table(flags$animate)

##
##   0   1
## 155  39

tapply(flags$animate, flags$landmass, mean)

##         1         2         3         4         5         6
## 0.4193548 0.1764706 0.1142857 0.1346154 0.1538462 0.3000000
```

The above command apply the mean function to the 'animate' variable separately for each of the six landmass groups, thus giving us the proportion of flags containing an animate image WITHIN each landmass group.

Similarly, we can look at a summary of population values (in round millions) for countries with and without the color red on their flag.

```
tapply(flags$population, flags$red, summary)

## $`0`
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    0.00    0.00    3.00   27.63    9.00  684.00
##
## $`1`
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##     0.0     0.0     4.0    22.1    15.0  1008.0
```

Question: What is the median population (in millions) for countries *without* the color red on their flag?

## 3.5 split

split takes a vector or other objects and splits it into groups determined by a factor or list of factors.

```
str(split)

## function (x, f, drop = FALSE, ...)
```

- x is a vector (or list) or data frame
- f is a factor (or coerced to one) or a list of factors
- drop indicates whether empty factors levels should be dropped

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

## $`1`
##  [1] -0.642074795 -1.080305314  0.427325494 -1.873209657 -1.741598699
##  [6]  0.708129654 -1.110219225  0.007001781 -1.098660795  1.914250311
##
## $`2`
##  [1] 0.80514072 0.87066157 0.13616392 0.00710275 0.74310475 0.89354998
##  [7] 0.78449901 0.25709270 0.02149347 0.80960750
##
## $`3`
##  [1]  1.0741530  0.9109033  0.4381801  1.0105033 -0.3620661  1.9046239
##  [7]  1.9507956  2.6108858  1.1765143  0.4832202

lapply(split(x, f), mean)

## $`1`
## [1] -0.4489361
##
## $`2`
## [1] 0.5328416
##
## $`3`
## [1] 1.119771
```

## An Example

```
library(datasets)
head(airquality)

##   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

s <- split(airquality, airquality$Month)
lapply(s, function(x) colMeans(x[, c("Ozone", "Solar.R", "Wind")]))

## $`5`
##   Ozone  Solar.R     Wind
##      NA       NA 11.62258
##
## $`6`
##   Ozone   Solar.R      Wind
##      NA 190.16667  10.26667
##
```

```
## $`7`
##      Ozone    Solar.R       Wind
##         NA 216.483871   8.941935
##
## $`8`
##    Ozone  Solar.R     Wind
##       NA       NA 8.793548
##
## $`9`
##    Ozone  Solar.R     Wind
##       NA 167.4333  10.1800

sapply(s, function(x) colMeans(x[, c("Ozone", "Solar.R", "Wind")]))

##                  5         6         7         8         9
## Ozone          NA        NA        NA        NA        NA
## Solar.R        NA 190.16667 216.483871        NA 167.4333
## Wind    11.62258  10.26667   8.941935 8.793548  10.1800

sapply(s, function(x) colMeans(x[, c("Ozone", "Solar.R", "Wind")], na.rm = TRUE))

##                  5         6          7          8         9
## Ozone    23.61538  29.44444  59.115385  59.961538  31.44828
## Solar.R 181.29630 190.16667 216.483871 171.857143 167.43333
## Wind     11.62258  10.26667   8.941935   8.793548  10.18000
```

## Splitting on More than One Level

```
x <- rnorm(10)
f1 <- gl(2, 5)
f2 <- gl(5, 2)
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] 0.899 -1.215
##  $ 2.1: num(0)
##  $ 1.2: num [1:2] -0.243 1.566
##  $ 2.2: num(0)
##  $ 1.3: num 0.0322
##  $ 2.3: num 1.56
##  $ 1.4: num(0)
##  $ 2.4: num [1:2] -0.67711 -0.00528
##  $ 1.5: num(0)
##  $ 2.5: num [1:2] -0.229 0.432

str(split(x, list(f1, f2), drop = TRUE))
```

```
## List of 6
##  $ 1.1: num [1:2] 0.899 -1.215
##  $ 1.2: num [1:2] -0.243 1.566
##  $ 1.3: num 0.0322
##  $ 2.3: num 1.56
##  $ 2.4: num [1:2] -0.67711 -0.00528
##  $ 2.5: num [1:2] -0.229 0.432
```

## 3.6 Summary

Writing for, while loops is useful when programming but not particularly easy when working interactively on the command line. There are some functions which implement looping to make life easier.

- 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

## 3.7 str function

```
str(str)

## function (object, ...)

str(lm)

## function (formula, data, subset, weights, na.action, method = "qr",
##     model = TRUE, x = FALSE, y = FALSE, qr = TRUE, singular.ok = TRUE,
##     contrasts = NULL, offset, ...)

x <- rnorm(100, 2, 4)
summary(x)

##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -6.5890 -0.2183  2.1220  2.5870  5.2230 11.7100

str(x)

##  num [1:100] -2.24 -3.79 -0.43 4.71 -2.47 ...

f <- gl(40, 10)
str(f)

##  Factor w/ 40 levels "1","2","3","4",..: 1 1 1 1 1 1 1 1 1 1 ...

summary(f)
```

```
##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
## 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10
## 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40
## 10 10 10 10 10 10 10 10 10 10 10 10 10 10 10
```

```
library(datasets)
head(airquality)
```

```
##   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
```

```
str(airquality)
```

```
## 'data.frame': 153 obs. of  6 variables:
##  $ Ozone  : int  41 36 12 18 NA 28 23 19 8 NA ...
##  $ Solar.R: int  190 118 149 313 NA NA 299 99 19 194 ...
##  $ Wind   : num  7.4 8 12.6 11.5 14.3 14.9 8.6 13.8 20.1 8.6 ...
##  $ Temp   : int  67 72 74 62 56 66 65 59 61 69 ...
##  $ Month  : int  5 5 5 5 5 5 5 5 5 5 ...
##  $ Day    : int  1 2 3 4 5 6 7 8 9 10 ...
```

```
m <- matrix(rnorm(100), 10, 10)
str(m)
```

```
##  num [1:10, 1:10] 0.745 0.895 2.319 -0.984 0.583 ...
```

```
m[,1]
```

```
##  [1]  0.7446878  0.8949779  2.3192975 -0.9844778  0.5830348 -0.5133931
##  [7]  1.4112811  1.1223591 -0.4994462  0.4478064
```

```
s <- split(airquality, airquality$Month)
str(s)
```

```
## List of 5
##  $ 5:'data.frame': 31 obs. of  6 variables:
##   ..$ Ozone  : int [1:31] 41 36 12 18 NA 28 23 19 8 NA ...
##   ..$ Solar.R: int [1:31] 190 118 149 313 NA NA 299 99 19 194 ...
##   ..$ Wind   : num [1:31] 7.4 8 12.6 11.5 14.3 14.9 8.6 13.8 20.1 8.6 ...
##   ..$ Temp   : int [1:31] 67 72 74 62 56 66 65 59 61 69 ...
##   ..$ Month  : int [1:31] 5 5 5 5 5 5 5 5 5 5 ...
##   ..$ Day    : int [1:31] 1 2 3 4 5 6 7 8 9 10 ...
##  $ 6:'data.frame': 30 obs. of  6 variables:
##   ..$ Ozone  : int [1:30] NA NA NA NA NA NA 29 NA 71 39 ...
##   ..$ Solar.R: int [1:30] 286 287 242 186 220 264 127 273 291 323 ...
##   ..$ Wind   : num [1:30] 8.6 9.7 16.1 9.2 8.6 14.3 9.7 6.9 13.8 11.5 ...
##   ..$ Temp   : int [1:30] 78 74 67 84 85 79 82 87 90 87 ...
##   ..$ Month  : int [1:30] 6 6 6 6 6 6 6 6 6 6 ...
##   ..$ Day    : int [1:30] 1 2 3 4 5 6 7 8 9 10 ...
```

```
##  $ 7:'data.frame': 31 obs. of  6 variables:
##   ..$ Ozone  : int [1:31] 135 49 32 NA 64 40 77 97 97 85 ...
##   ..$ Solar.R: int [1:31] 269 248 236 101 175 314 276 267 272 175 ...
##   ..$ Wind   : num [1:31] 4.1 9.2 9.2 10.9 4.6 10.9 5.1 6.3 5.7 7.4 ...
##   ..$ Temp   : int [1:31] 84 85 81 84 83 83 88 92 92 89 ...
##   ..$ Month  : int [1:31] 7 7 7 7 7 7 7 7 7 7 ...
##   ..$ Day    : int [1:31] 1 2 3 4 5 6 7 8 9 10 ...
##  $ 8:'data.frame': 31 obs. of  6 variables:
##   ..$ Ozone  : int [1:31] 39 9 16 78 35 66 122 89 110 NA ...
##   ..$ Solar.R: int [1:31] 83 24 77 NA NA NA 255 229 207 222 ...
##   ..$ Wind   : num [1:31] 6.9 13.8 7.4 6.9 7.4 4.6 4 10.3 8 8.6 ...
##   ..$ Temp   : int [1:31] 81 81 82 86 85 87 89 90 90 92 ...
##   ..$ Month  : int [1:31] 8 8 8 8 8 8 8 8 8 8 ...
##   ..$ Day    : int [1:31] 1 2 3 4 5 6 7 8 9 10 ...
##  $ 9:'data.frame': 30 obs. of  6 variables:
##   ..$ Ozone  : int [1:30] 96 78 73 91 47 32 20 23 21 24 ...
##   ..$ Solar.R: int [1:30] 167 197 183 189 95 92 252 220 230 259 ...
##   ..$ Wind   : num [1:30] 6.9 5.1 2.8 4.6 7.4 15.5 10.9 10.3 10.9 9.7 ...
##   ..$ Temp   : int [1:30] 91 92 93 93 87 84 80 78 75 73 ...
##   ..$ Month  : int [1:30] 9 9 9 9 9 9 9 9 9 9 ...
##   ..$ Day    : int [1:30] 1 2 3 4 5 6 7 8 9 10 ...
```

# 4   Debugging

## 4.1   Something's Wrong!

Something's Wrong! Indications that something's not right:

- message: A generic notification/diagnostic message produced by the message function; execution of the function continues

- warning: An indication that something is wrong but not necessarily fatal; execution of the function continues; generated by the warning function

- error: An indication that a fatal problem has occurred; execution stops; produced by the stop function

- condition: A generic concept for indicating that something unexpected can occur; programmers can create their own conditions

```
log(-1)
```

**The Essence of Debugging: The Principle of Confirmation**

Fixing a buggy program is a process of confirming, one by one, that the many things you believe to be true about the code actually are true. When you find that one of your assumptions is not true, you have found a clue to the location (if not the exact nature) of a bug.

**Start Small**

At least at the beginning of the debugging process, stick to small, simple test cases. Working with large data objects may make it harder to think about the problem. Of course, you should eventually test your code on large, complicated cases, but start small.

**Debug in a Modular, Top-Down Manner**

Most good software developers agree that code should be written in a modular manner. The functions should not be too lengthy and should call other functions if necessary. This makes the code easier to organize during the writing stage and easier for others to understand when it comes time for the code to be extended.

**Antibugging**

You may adopt some "antibugging" strategies as well. Suppose you have a section of code in which a variable $x$ should be positive. You could insert this line:

```
stopifnot(x > 0)
```

If there is a bug earlier in the code that renders x equal to, say, $-12$, the call to stopifnot() will bring things to a halt right there, with an error message.

**Why Use a Debugging Tool?**

For example, say you have the following code:

```
x <- y^2 + 3*g(z,2)
w <- 28
if (w+q > 0) u <- 1 else v <- 10
```

In the old days, programmers would perform the debugging confirmation process by temporarily inserting print statements into their code and rerunning the program to see what printed out.

```
x <- y^2 + 3*g(z,2)
cat("x =",x,"\n")
w <- 28
if (w+q > 0) {
u <- 1
print("the 'if' was done")
} else {
v <- 10
print("the 'else' was done")
}
```

We would rerun the program and inspect the feedback printed out. We would then remove the print statements and put in new ones to track down the next bug. This manual process is fine for one or two cycles, but it gets really tedious during a long debugging session. And worse, all that editing work distracts your attention, making it harder to concentrate on finding the bug.

## An example

```
printmessage <- function(x){
if (x > 0)
print("x is greater than zero")
else print("x is less than or equal to zero")
invisible(x)
}
```

```
printmessage(1)
printmessage(NA)
```

```
printmessage2 <- function(x) {
if(is.na(x))
print("x is a missing value!")
else if(x > 0) print("x is greater than zero")
else
print("x is less than or equal to zero")
invisible(x)
}
x <- log(-1)
printmessage2(x)
```

How do you know that something is wrong with your function?

- What was your input?

- How did you call the function?

- What were you expecting?

- Output, messages, other results?

- What did you get?

- How does what you get differ from what you were expecting?

- Were your expectations correct in the first place?

- Can you reproduce the problem (exactly)?

## 4.2   Debugging Tools in R

The primary tools for debugging functions in R are

- traceback: prints out the function call stack after an error occurs; does nothing if there's no error

- debug: flags a function for "debug" mode which allows you to step through execution of a function one line at a time

- browser: suspends the execution of a function wherever it is called and puts the function in debug mode

- trace: allows you to insert debugging code into a function a specific places

- recover: allows you to modify the error behavior so that you can browse the function call stack

## Debugging with RStudio

https://support.rstudio.com/hc/en-us/articles/205612627-Debugging-with-RStudio