

Module 1

Overview and History of R

R is a dialect of the S language.

What is S?

S is a language that was developed by John Chambers and others at Bell Labs. S was initiated in 1976 as an internal statistical analysis environment - originally implemented as Fortran libraries. Early versions of the language did not contain functions for statistical modeling. In 1988 the system was rewritten in C and began to resemble the system we see today. In 1993, Bell Labs gave StatSci an exclusive license to develop and sell the S language. In 2004 Insightful purchased the S language from Lucent and is the current owner. In 2006 Alcatel purchased Lucent technologies. Insightful sells its implementation of S under the name S-PLUS. In 2008, Insightful is acquired by TIBCO for \$25 million. The makers of S wanted to create an interactive environment for users to work with data and perform basic analyses. Later the user could use the language to develop their own tools.

A short history of R

R was created in 1991 by Ross Ihaka and Robert Gentleman. In 1993, R was announced to the public for the first time. In 1995, Martin Machler convinced Ross and Robert to use the GNU General Public License; thus making it a free software. In 1996, a mailing list was created for developers using R. This led to the formation of the R core group in 1997; this group controls the source code of R

Features of R

- Syntax is very similar to S, making migration easier.
- Semantics are superficially similar to S, but in reality are quite different.
- Can run on any standard computing platform or operating system.
- There is active development with frequent bugfix releases.
- The software itself is quite lean; the functionalities are organized in modular packages
- It has much better graphic and visualization capabilities when compared to most other stat packages.
- While it is useful for working interactively (Markdown and Console), it is also a powerful programming language to develop new tools
- It has a vibrant and active user community.
- It is free!!

What is meant by “free software”?

Free software, libre software, libreware sometimes known as freedom-respecting software is computer software distributed under terms that allow users to run the software for any purpose as well as to study, change, and distribute it and any adapted versions. Free software is a matter of liberty, not price; all users are legally free to do what they want with their copies of free software (including profiting from them) regardless of how much is paid to obtain the program.

- The freedom to run the program, for any purpose.
- The freedom to study how the program works; adapt it to your needs. Access to the source code is a precondition for this.
- The freedom to redistribute copies so you can help your neighbour.
- The freedom to improve the program, and release your improvements to the public so that the whole community benefits.

Drawbacks of R

- Essentially based on a 40 year old technology.
- Minimal built-in support for dynamic or 3D graphics.
- New functionality is based on consumer demand and user contributions.
- Objects must be stored in physical memory; however there are methods to deal with this.
- Not ideal for all situations.

Design of the R system

The R system divided into two conceptual parts: base R and everything else. R functionality is divided into a number of packages.

- The “base” R system contains the **base** package required to run R, and contains the most fundamental actions
- The other packages contained in the base system include utils, stats, datasets, graphics, grDevices, grid, methods, tools, parallel, compiler, splines, tcltk, stats4.
- There are also recommended packages like boot, class, cluster, codetools, foreign, KernSmooth, lattice, mgcv, nime, rpart, survival, MASS, spatial, nnet, Matrix.

There are about 4000 packages on CRAN that have been developed by users and programmers around the world. There are also many packages associated with the Bioconductor project. Additionally, a lot of people also make packages available on their personal websites; there is no reliable way to keep track of how many packages exist in this way.

R console input and evaluation

Expressions can be typed into the R prompt for evaluation. The `<-` operator is used to assign a value to a variable

```
x <- 1
x
```

```
## [1] 1
```

```
msg <- "Hello"
```

Entering the name of a variable in the console will display its value. One can also use the `print()` function.

```
print(x)
```

```
## [1] 1
```

When a complete expression is entered at the prompt, it is evaluated and the result of the expression is returned. It may sometimes be auto-printed (e.g. when entering a variable, we get its value).

```
msg
```

```
## [1] "Hello"
```

[1] indicates that msg has one element and that it is the string Hello.

The colon is used to create a sequence of integers.

```
x <- 6:25  
x
```

```
## [1] 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
```

The result has [1] before printing the elements.

Data types in R

Objects and Attributes

R has five basic categories of objects: Character, numeric, integer, complex, logical. The most basic object in R is a vector. A vector can only contain objects of the same class. The only exception is a list; which can be a collection of objects belonging to different classes.

Numbers in R are generally treated as numeric objects i.e. double precision real numbers. If you need an integer, specify the L suffix. 1 will return a numeric type while 1L will return an integer type. R also lets you use a special number called Inf. It can be used in regular calculations. 1/0 will return Inf while 1/Inf will return 0. NaN stands for undefined values. 0/0 will return NaN.

Objects in R can have attributes associated with them. Attributes are: names or dimnames, dimensions (for arrays and matrices), class, length and any other attribute that is user-defined. Attributes of an object can be accessed using the attributes() function.

Vectors and Lists

The c() function is used to create vectors of objects. It concatenates the elements together.

```
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
```

We can also use the vector() function to create a vector of a certain type and length. In the example below we see a numeric vector of size 10. It will be initialized with the default value of its type. For numeric the default is 0.

```
x <- vector("numeric", length = 10)
x
```

```
## [1] 0 0 0 0 0 0 0 0 0 0
```

Apart from the `:` operator, the `seq` function can also be used to create sequences of numbers.

```
seq(0,10, by=0.5)
```

```
## [1] 0.0 0.5 1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0 5.5 6.0 6.5 7.0
## [16] 7.5 8.0 8.5 9.0 9.5 10.0
```

The `by` argument is used to specify the value by which we should increment / decrement from one element to the next. So the code above gives a vector that starts at 0 and goes up to 10, increasing by 0.5 each time. Sometimes, we may only want a certain number of elements, regardless of the increment. The code below gives a vector that starts from 5 (first element is 5) and ends at 10 with length 30 (thirtieth element is 10). It calculates the appropriate value to increment by in order to satisfy this condition.

```
seq(5,10, length=30)
```

```
## [1] 5.000000 5.172414 5.344828 5.517241 5.689655 5.862069 6.034483
## [8] 6.206897 6.379310 6.551724 6.724138 6.896552 7.068966 7.241379
## [15] 7.413793 7.586207 7.758621 7.931034 8.103448 8.275862 8.448276
## [22] 8.620690 8.793103 8.965517 9.137931 9.310345 9.482759 9.655172
## [29] 9.827586 10.000000
```

```
my_seq <- 10:15 #This is a vector having 6 elements
```

```
seq_along(my_seq)
```

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

`seq_along` is a handy function when you want to create a vector starting from 1 upto a certain length. In the above code, we can see it creates a vector of the exact same size as the argument that it gets passed.

When you mix different types of objects in a single vector, the general rule is that R will try to coerce the vector to get each element into the same class. This is called coercion.

```
y <- c(1.7, "a")
y
```

```
## [1] "1.7" "a"
```

```
z <- c(TRUE, 3)
z
```

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

```
x <- c("a", FALSE)
x
```

```
## [1] "a"      "FALSE"
```

We can see in the first example that the resultant vector is a vector of strings. This is because R can convert 1.7 into character but not convert “a” into a number. In the second example, the logical value TRUE is converted to the number 1; FALSE would have been converted to zero. In the third example, the logical value FALSE is converted to a string.

Logical Vectors

```
num_vect <- c(0.5, 55, -10, 6)
tf <- num_vect < 1
```

This created a vector of numbers and then a second vector with the result of the logical expression `<1`. So the result we get is a logical vector, with each element being either TRUE or FALSE. The value depends on the result of the expression.

```
num_vect >= 6
```

```
## [1] FALSE TRUE FALSE TRUE
```

Character vectors

```
my_char <- c("My", "name", "is", "Stewart")
my_char
```

```
## [1] "My"      "name"    "is"      "Stewart"
```

```
paste(my_char, collapse = " ")
```

```
## [1] "My name is Stewart"
```

The collapse argument to the `paste()` function tells R that when we join the elements of the `my_char` vector, they should be separated with single spaces.

Classes and coercion

```
x <- 0:6
x
```

```
## [1] 0 1 2 3 4 5 6
```

```
class(x)
```

```
## [1] "integer"
```

```
as.logical(x)
```

```
## [1] FALSE TRUE TRUE TRUE TRUE TRUE TRUE
```

```
as.character(x)
```

```
## [1] "0" "1" "2" "3" "4" "5" "6"
```

Explicit coercion can be done using the `as.*` functions in R.

```
x <- c("A", "B", "C")
as.numeric(x)
```

```
## Warning: NAs introduced by coercion
```

```
## [1] NA NA NA
```

```
as.logical(x)
```

```
## [1] NA NA NA
```

If a conversion is attempted that is not possible, we get NaN and corresponding warnings.

Lists are a special type of vector that can contain elements of different classes

```
x <- list(1, "a", TRUE, 1+2i, FALSE)
x
```

```
## [[1]]
## [1] 1
##
## [[2]]
## [1] "a"
##
## [[3]]
## [1] TRUE
##
## [[4]]
## [1] 1+2i
##
## [[5]]
## [1] FALSE
```

Here you can see that the first element is a vector of integers, the second is a vector of characters, the third is a vector of logical TRUE and so on. Elements of a list will have double brackets around them.

```
x <- list("a", "b", 0.5, 1, FALSE, "d")
x
```

```
## [[1]]
## [1] "a"
##
## [[2]]
## [1] "b"
##
## [[3]]
## [1] 0.5
##
## [[4]]
## [1] 1
##
## [[5]]
## [1] FALSE
##
## [[6]]
## [1] "d"
```

Matrices

Matrices are vectors with a dimension attribute. The dimension attribute is 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
```

Matrices are constructed column-wise, so entries can be thought of starting in the upper-left corner and running down the columns.

```
m <- matrix(1:6, nrow=2, ncol=3)
m
```

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

Matrices can also 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
```

Factors

Factors are used to represent categorical data. Factors can be unordered or ordered. Gender (male, female) would be an unordered factor whereas something like a grade (A, B, C, D) is ordered because it has a rank. One can think of a factor as an integer vector where each integer has a label. Factors are treated specially by modeling functions like `lm()` and `glm()`. Using factors with labels is better than using integers because factors are self-describing; having a variable with values “male” and “female” is better than a variable with values 1 and 2.

A factor can be created using the `factor()` function in R. Its input should usually be a character vector

```
x <- factor(c("yes", "yes", "no", "yes", "no"))
x
```

```
## [1] yes yes no  yes no
## Levels: no yes
```

The function `table` will give a frequency count of how many elements of each level are present in the factor.

```
table(x)
```

```
## x
##  no yes
##   2  3
```

`unclass` reduces the factor to an integer vector. In the example below, yes is coded as 2 and no is coded as 1. A factor at its core is an integer vector with the levels attribute.

```
unclass(x)
```

```
## [1] 2 2 1 2 1
## attr("levels")
## [1] "no"  "yes"
```


The order of the levels of a factor can be set by using the `levels()` argument to a `factor()`. This can be important in linear modeling because the first level is used as the baseline level. If levels are not mentioned then R processes them alphabetically. So because N comes before Y, no would be the baseline level and yes will be the second level

```
x <- factor(c("yes", "yes", "no", "yes", "no", "yes", "no"),
            levels = c("yes", "no"))
x
```

```
## [1] yes yes no  yes no  yes no
## Levels: yes no
```

In the factor initialized above, yes is the baseline level and no is the second level. It can be verified using the `unclass` function.

```
unclass(x)
```

```
## [1] 1 1 2 1 2 1 2
## attr("levels")
## [1] "yes" "no"
```

Missing values

Missing values in R are denoted by either NA or NaN. NaN is used for undefined mathematical operations. And NA is pretty much used for everything else. There's a function in R called `is.na` and `is.nan` which is used to test objects to see if there are missing values in that object. NA values also have a class, so there can be integer NA, character NA, etc. A NaN value is also NA but the converse is not true

```
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
```

The example below illustrates how NaN is considered to be NA, but NA is not a NaN

```
Y <- c(1, 2, NA, NaN, 4)
is.na(Y)
```

```
## [1] FALSE FALSE  TRUE  TRUE FALSE
```

```
is.nan(Y)
```

```
## [1] FALSE FALSE FALSE  TRUE FALSE
```

```

y <- rnorm(1000)
z <- rep(NA,1000)
my_data <- sample(c(y,z), 100)
my_na <- is.na(my_data)
my_na

```

```

## [1] FALSE TRUE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE TRUE TRUE
## [13] FALSE FALSE FALSE TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [25] FALSE FALSE TRUE TRUE FALSE FALSE TRUE TRUE FALSE FALSE FALSE TRUE TRUE
## [37] FALSE FALSE TRUE FALSE FALSE FALSE TRUE FALSE FALSE FALSE TRUE FALSE
## [49] FALSE FALSE TRUE FALSE FALSE FALSE FALSE TRUE TRUE FALSE TRUE FALSE
## [61] FALSE FALSE FALSE TRUE FALSE TRUE TRUE TRUE FALSE FALSE TRUE FALSE
## [73] FALSE TRUE FALSE TRUE FALSE FALSE FALSE TRUE TRUE TRUE FALSE FALSE
## [85] TRUE TRUE TRUE FALSE TRUE TRUE FALSE FALSE TRUE FALSE FALSE FALSE
## [97] TRUE TRUE TRUE TRUE

```

Trying to use the == operator to check for NAs

```
my_data == NA
```

```

## [1] NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA
## [26] NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA
## [51] NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA
## [76] NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA NA

```

The reason you got a vector of all NAs is that NA is not really a value, but just a placeholder for a quantity that is not available. Therefore the logical expression is incomplete and R has no choice but to return a vector of the same length as my_data that contains all NAs.

To calculate how many NAs in my_data : The trick is to recognize that underneath the surface, R represents TRUE as the number 1 and FALSE as the number 0. Therefore, if we take the sum of a bunch of TRUEs and FALSEs, we get the total number of TRUEs.

```
sum(my_na)
```

```
## [1] 44
```

```
my_data
```

```

## [1] 1.36775080      NA 0.56753024 -0.63459978 -1.22551838 -1.07748228
## [7] -0.94049717      NA 0.91892093 2.00027955      NA      NA
## [13] -0.87560151 1.80934791 0.21922056      NA 0.39933813      NA
## [19]      NA      NA      NA      NA      NA      NA
## [25] 2.03084077 0.65677710      NA      NA 0.40880176 -0.10528897
## [31]      NA      NA 1.00169168 -1.04735081 -1.98823766      NA
## [37] -0.09202365 -0.80987710      NA 0.70342492 2.21191367 0.72041983
## [43]      NA 0.66561652 -0.09561535 0.45100658      NA 2.47137998
## [49] 0.02984437 -0.66548078      NA -1.24277085 0.94793404 -1.05243409
## [55] 0.56582590      NA      NA 0.56735099      NA -1.42606944
## [61] -0.28663065 -0.11314081 -0.48063282      NA -0.27103201      NA
## [67]      NA      NA -0.71749269 -1.19921654      NA 1.36190341

```

```
## [73] -0.73164730      NA -0.72872951      NA -0.16421026 -0.78168718
## [79] -0.42560950      NA      NA      NA -0.55760263  2.38480915
## [85]      NA      NA      NA  1.58402747      NA      NA
## [91] -0.02264847 -0.77597694      NA -0.55423602 -0.06158198 -0.67065656
## [97]      NA      NA      NA      NA
```

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 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. Data frames also have a special attribute called `row.names()`

```
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
```

Names

R objects can also have names. This can be used for writing readable code and self-describing objects.

```
x <- 1:3
names(x) <- c("foo", "bar", "norf")
x
```

```
##   foo   bar norf
##    1    2    3
```

Lists can also have names

```
x <- list(a=1, b=2, c=3)
x
```

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

```
## [1] 2
##
## $c
## [1] 3
```

Names for matrices can be set as follows. Provide a list with two vectors, one for the rows and one for columns.

```
m <- matrix(1:4, nrow=2, ncol=2)
dimnames(m) <- list(c("a", "b"), c("c", "d"))
m
```

```
##    c d
## a 1 3
## b 2 4
```

We can also convert a vector into a matrix by setting dimensions for it

```
my_vector <- 1:20
my_vector
```

```
## [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
```

```
dim(my_vector)
```

```
## NULL
```

```
dim(my_vector) <- c(4,5)
my_vector
```

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

```
class(my_vector)
```

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

Reading and Writing Data in R: The basics

The `read.table()` and `read.csv()` functions are used to read tabular data. They are the most commonly used functions for reading data into R. They can read any kind of files that are formatted as rows and columns and return a data frame in R. The function `readLines()` is for reading lines of a text file; this returns text as a character vector in R. The function `source()` is used for reading R code. The `dget()` function is also used for reading R code files, specifically for R objects that have been deparsed into text files. The `unserialize()` function is used to read R objects in binary form. The analogous functions for writing data are `write.table()`, `writeLines()`, `dump()`, `dput()`, `save()` and `serialize()`.

Reading Tabular data

The `read.table` function is the most commonly used for reading data into R. It has some 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`: integer, the number of lines to skip from the beginning
- `stringAsFactors`: logical, should the character variables be coded as factors? It defaults to `true`,

For small to moderately sized datasets, you can usually call `read.table()` without specifying any other arguments. R will automatically skip lines that begin with a `#`, figure out how many rows are there (and how much memory needs to be allocated), figure what type of variable is in each column of the table. However, telling R all of these things directly makes it run faster and more efficiently. This is why its important to specify these arguments when the dataset is large. `read.csv` is identical to `read.table` except that the default separator is a comma. For `read.table` the default separator is a single space.

Reading large tables

Textual data formats

Connections : Interfaces

Subsetting