

R Tutorial: Creating Data Objects in R

Stat 2080: Statistical Modeling

Joe Roith, Fall 2018

1. Creating Data in R:

R is a *functional* language as most of the operations in R are carried out through the use of *functions*: either internal R functions or user-constructed functions. The function `c()` is an example of a simple R function. It is used to create *vector* functions as follows:

```
y <- c(5.0, 7.1, 13.0)
```

```
y
```

```
## [1] 5.0 7.1 13.0
```

```
regions <- c("Northeast", "Midwest", "Southeast", "Northwest", "Southwest")
```

```
regions
```

```
## [1] "Northeast" "Midwest" "Southeast" "Northwest" "Southwest"
```

When the number of values gets larger, the `scan()` function is more useful. You may input data values separated by blanks at the keyboard using `scan()`. In this form, `scan()` prompts the user to enter data with a varying number of values to a line until an entire blank line is entered to signal the end of data entry:

```
z <- scan()
```

```
1: 3.7 6.5 1.9 8.4 10
```

```
6: 6.2 15 11.4
```

```
10:
```

```
Read 9 items
```

```
z
```

```
## [1] 3.7 6.5 1.9 8.4 10.0 6.2 15.0 11.4
```

In the following example, character strings with embedded spaces are to be entered as data values, and hence a *separator* different from a blank space is required.

```
actors <- scan(what = "", sep = "#")
```

```
1: Dwayne Johnson#Tom Hanks#Denzel Washington
```

```
4: Meryl Streep#Cate Blanchett
```

```
6: Octavia Spencer#Benedict Cumberbatch
```

```
8:
```

```
actors
```

```
## [1] "Dwayne Johnson" "Tom Hanks" "Denzel Washington"
```

```
## [4] "Meryl Streep" "Cate Blanchett" "Octavia Spencer"
```

```
## [7] "Benedict Cumberbatch"
```

Note that the value `"",` specified as the value of the argument `what =,` indicates that these values are to be read as character strings. Alternatively, the *type* may be specified as `what = "character".` By default, the *type* of data values is assumed to be numeric (stored as a double). Finally, `sep = "#"` specifies that the alternative *separator* to be used is the `#` symbol.

2. Reading Data from Text Files:

Typically we use the `Import Dataset` option in the Environment window. Then select the type of file. But the information below is how that process is automated. It is still good to know.

Most text data files are in “ASCII Format”, and are usually produced by a text editor such as Notepad or Emacs.

Since the text file **blood.txt** is in the current working directory, it can be accessed within R simply as shown in the example below. Otherwise, it is necessary to enter the complete pathname of the location of the text file within the double quotes. Additionally, when entering the *complete pathname* you will need to change the folder separators from a single backslash, `\` to either a double backslash, `\\` or a forward slash, `/`. `scan()` reads the data from the text file row-wise and creates a vector object.

```
blood <- scan("U:\\Stat 2080\\blood.txt")
```

This data set consists of values for two diets, 1 and 2, randomly assigned to 10 mice, and blood coagulation times measured for each animal. Two separate vector objects can be created to represent them using the `scan()` function in the following way:

```
blood1 <- scan("blood.txt", what = list(diet = 0, time = 0))
```

```
## $diet
## [1] 1 1 1 1 2 2 2 2 2 2
##
## $time
## [1] 62 60 63 59 63 67 71 64 65 66
```

The `what = list(diet = 0, time = 0)` argument to `scan()` specifies that two vectors named `diet` and `time` will be created as components of a *list* object. `diet = 0, time = 0` indicates that values for these vectors are to be read as numeric values. It is important that the data be alternating between `diet` and `time` in the original file for this to work. `blood1` is an example of a list object that will be discussed in detail later. The two vectors named `diet` and `time` are components of the list `blood1` and may be referenced as `blood1$diet` and `blood1$time`, thus:

```
blood1$diet
```

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

```
blood1$time
```

```
## [1] 62 60 63 59 63 67 71 64 65 66
```

Typically, `scan()` is used to feed data into another function to create an object that has a desired structure. For example, the data in the text named **insulin.txt** available in the current working directory can be used to create a *matrix* object. In the following example, the `scan()` function is used directly as an argument to the `matrix()` function, to create the matrix object named *insulin*.

```
insulin <- matrix(scan("insulin.txt"), ncol = 3, byrow = T)
insulin
```

```
##      [,1] [,2] [,3]
## [1,] 2.02 1.49 0.33
## [2,] 3.83 2.67 1.67
## [3,] 6.67 4.62 4.67
## [4,] 5.38 4.18 2.45
## [5,] 5.49 2.78 2.29
## [6,] 3.50 2.56 1.95
## [7,] 5.90 4.46 0.49
## [8,] 4.89 3.79 1.81
```

A *data frame* is an R object very similar in appearance to a matrix object except that columns in a data frame can be of different types: *numeric* or *character*. This allows standard rectangular data sets to be represented as a data frame. Data frames are a required form of input data for many statistical functions in R used for data analysis and modeling. The `read.table()` function allows us to convert data available in an external text file directly into a data frame. In the following example, we use the text file **insulin.txt** to create a data frame in R.

```
insulin1 <- read.table("insulin.txt")
insulin1
```

```
##      V1    V2    V3
## 1 2.02  1.49  0.33
## 2 3.83  2.67  1.67
## 3 6.67  4.62  4.67
## 4 5.38  4.18  2.45
## 5 5.49  2.78  2.29
## 6 3.50  2.56  1.95
## 7 5.90  4.46  0.49
## 8 4.89  3.79  1.81
```

An object *type* can be checked with the `class()` function. An object can be changed typically with the `as.####` prefix. For example, a matrix can be converted to a data frame with `as.data.frame()` and vice versa with `as.matrix()`.

```
class(insulin)
```

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

```
class(insulin1)
```

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

```
insulin2 <- as.data.frame(insulin1)
```

```
class(insulin2)
```

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

In the absence of user-specified names for the columns (variables) and the rows (observations), by default R assigns the names V1, V2, etc. for the columns and the names 1, 2, etc. for the rows. The user can assign names by including them in the text file along with the data itself or by using the keyword parameters `col.names` and/or `row.names` when using the `read.table` command, as illustrated below:

```
insulin1 <- read.table("insulin.txt", col.names = c("Week1", "Week2", "Week3"))
insulin1
```

```
##   Week1 Week2 Week3
## 1  2.02  1.49  0.33
## 2  3.83  2.67  1.67
## 3  6.67  4.62  4.67
## 4  5.38  4.18  2.45
## 5  5.49  2.78  2.29
## 6  3.50  2.56  1.95
## 7  5.90  4.46  0.49
## 8  4.89  3.79  1.81
```

If the data were separated by another character, such as a comma, then a different command can be used, such as:

```
insulin2 <- read.table("insulin.txt", sep=",")
```

As observed before, this is an option if one or more of the variables were of character type that included spaces, because blanks cannot be used as separators in this case.

A number of built-in datasets are supplied with R. Sample datasets are also available with other packages that may be already installed in R. To see a list of the datasets available with the base system, use the command `data()`. To load one of the datasets listed into R, use the same function. For example, `data(cars)` or `data(DDT, package="MASS")`:

```
data(cars)
head(cars)
```

```
##   speed dist
## 1     4     2
## 2     4    10
## 3     7     4
## 4     7    22
## 5     8    16
## 6     9    10
```

```
data(DDT, package = "MASS")
DDT
```

```
## [1] 2.79 2.93 3.22 3.78 3.22 3.38 3.18 3.33 3.34 3.06 3.07 3.56 3.08 4.64
## [15] 3.34
```

`cars` is an R data frame object. Attaching the data frame using the `attach(cars)` command, allows the columns in the data frame to be accessed by simply giving their names. Thus,

```
attach(cars)
speed
```

```
## [1]  4  4  7  7  8  9 10 10 10 11 11 12 12 12 12 13 13 13 13 14 14 14 14
## [24] 15 15 15 16 16 17 17 17 18 18 18 18 19 19 19 20 20 20 20 20 22 23 24
## [47] 24 24 24 25
```

3. Generating Data in R:

The simplest operator for generating data vectors in R is the colon `:` called the sequence operator. For example,

```
1:10
```

```
## [1]  1  2  3  4  5  6  7  8  9 10
```

```
index <- 4:-5
index
```

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

This operator generates vectors of numbers incremented by either 1 or -1. You may use the `seq()` function to generate sequences incremented by specific constants:

```
seq(from = 0, to = 1, by = 0.1)
```

```
## [1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
```

```
seq(0, 1, 0.1)
```

```
## [1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
```

```
seq(0, 1, length = 10)
```

```
## [1] 0.0000000 0.1111111 0.2222222 0.3333333 0.4444444 0.5555556 0.6666667  
## [8] 0.7777778 0.8888889 1.0000000
```

The arguments to an R function may be specified as *named* arguments i.e., in the form `name = value` or just by specifying the values if they are provided in the same order as given in the function specifications. For `seq()` function, it is recommended that the third parameter always be named. The `rep()` function is another function for generating data in R. In the simplest use, this function repeats values in the first argument a number of times equal to the second argument. Some examples are:

```
rep(3, times = 5)
```

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

```
rep(1:5, times = 3)
```

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

```
rep(1:5, each = 2)
```

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

```
rep(1:5, rep(2, 5))
```

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

```
rep(1:5, each = 2, times = 3)
```

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

4. Generating Random Data in R:

Random samples from standard distributions are used in statistics for large scale simulations using Monte Carlo methods. They can be generated using the functions such as `runif()`, `rnorm()`, `rgamma()`, `rbinom()`, `rexp()`, and `rt()`. As you may observe, the function name is made by combining the letter “r” with a mnemonic name identifying the target distribution.

```
runif(n = 10)
```

```
## [1] 0.9939893 0.9605696 0.3030899 0.6363989 0.9889042 0.5918193 0.4216367  
## [8] 0.3698891 0.4014895 0.9572764
```

```
rnorm(n = 5, mean = 0, sd = 1)
```

```
## [1] 1.0026318 0.2485702 -0.2086163 0.4011047 -1.9472902
```

```
rnorm(5, 10, 2)
```

```
## [1] 12.273915 7.098129 12.190367 8.350579 12.120488
```

```
rt(n = 8, df = 10)
```

```
## [1] 0.02888670 0.02038437 -0.29292711 0.22074428 0.24923834 0.62040332  
## [7] -3.06765801 -0.60047063
```

```
rbinom(n = 20, size = 5, prob = 0.5)
```

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

```
stem(rnorm(100, 5, 2))
```

```
##
## The decimal point is at the |
##
## -0 | 8
## 0 | 1389
## 1 | 556
## 2 | 03579
## 3 | 00001112445666899
## 4 | 03466778888
## 5 | 0001222233344444455556667889
## 6 | 123333344455558899
## 7 | 11237778
## 8 | 0226
## 9 | 0
```