

# Chapter2

## Chapter 2 R basics

In this book, we will be using the R software environment for all our analysis. You will learn R and data analysis techniques simultaneously. To follow along you will therefore need access to R. We also recommend the use of an integrated development environment (IDE), such as RStudio, to save your work. Note that it is common for a course or workshop to offer access to an R environment and an IDE through your web browser, as done by RStudio cloud<sup>12</sup>. If you have access to such a resource, you don't need to install R and RStudio. However, if you intend on becoming an advanced data analyst, we highly recommend installing these tools on your computer<sup>13</sup>. Both R and RStudio are free and available online.

### 2.1 Case study: US Gun Murders

Imagine you live in Europe and are offered a job in a US company with many locations across all states. It is a great job, but news with headlines such as US Gun Homicide Rate Higher Than Other Developed Countries<sup>14</sup> have you worried. Charts like this may concern you even more:

Or even worse, this version from everytown.org:

But then you remember that the US is a large and diverse country with 50 very different states as well as the District of Columbia (DC).

California, for example, has a larger population than Canada, and 20 US states have populations larger than that of Norway. In some respects, the variability across states in the US is akin to the variability across countries in Europe. Furthermore, although not included in the charts above, the murder rates in Lithuania, Ukraine, and Russia are higher than 4 per 100,000. So perhaps the news reports that worried you are too superficial. You have options of where to live and want to determine the safety of each particular state. We will gain some insights by examining data related to gun homicides in the US during 2010 using R.

Before we get started with our example, we need to cover logistics as well as some of the very basic building blocks that are required to gain more advanced R skills. Be aware that the usefulness of some of these building blocks may not be immediately obvious, but later in the book you will appreciate having mastered these skills.

### 2.2 The very basics

Before we get started with the motivating dataset, we need to cover the very basics of R.

#### 2.2.1 Objects

Suppose a high school student asks us for help solving several quadratic equations of the form  $ax^2+bx+c=0$ . The quadratic formula gives us the solutions:

$$\frac{-b-\sqrt{b^2-4ac}}{2a} \quad \text{and} \quad \frac{-b+\sqrt{b^2-4ac}}{2a}$$

which of course change depending on the values of a, b, and c. One advantage of programming languages is that we can define variables and write expressions with these variables, similar to how we do so in math,

but obtain a numeric solution. We will write out general code for the quadratic equation below, but if we are asked to solve  $x^2 + x - 1 = 0$

```
a <- 1
b <- 1
c <- -1
```

which stores the values for later use. We use `<-` to assign values to the variables.

We can also assign values using `=` instead of `<-`, but we recommend against using `=` to avoid confusion.

Copy and paste the code above into your console to define the three variables. Note that R does not print anything when we make this assignment. This means the objects were defined successfully. Had you made a mistake, you would have received an error message.

To see the value stored in a variable, we simply ask R to evaluate `a` and it shows the stored value:

```
a

## [1] 1
```

A more explicit way to ask R to show us the value stored in `a` is using `print` like this:

```
print(a)

## [1] 1
```

We use the term *object* to describe stuff that is stored in R. Variables are examples, but objects can also be more complicated entities such as functions, which are described later.

### 2.2.2 The workspace

As we define objects in the console, we are actually changing the workspace. You can see all the variables saved in your workspace by typing:

```
ls()

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

In RStudio, the Environment tab shows the values:

We should see `a`, `b`, and `c`. If you try to recover the value of a variable that is not in your workspace, you receive an error. For example, if you type `x` you will receive the following message: Error: object 'x' not found.

Now since these values are saved in variables, to obtain a solution to our equation, we use the quadratic formula:

```
(-b + sqrt(b^2 - 4*a*c) ) / ( 2*a )

## [1] 0.618034
```

```
(-b - sqrt(b^2 - 4*a*c) ) / ( 2*a )
```

```
## [1] -1.618034
```

### 2.2.3 Functions

Once you define variables, the data analysis process can usually be described as a series of functions applied to the data. R includes several predefined functions and most of the analysis pipelines we construct make extensive use of these.

We already used the `install.packages`, `library`, and `ls` functions. We also used the function `sqrt` to solve the quadratic equation above. There are many more prebuilt functions and even more can be added through packages. These functions do not appear in the workspace because you did not define them, but they are available for immediate use.

In general, we need to use parentheses to evaluate a function. If you type `ls`, the function is not evaluated and instead R shows you the code that defines the function. If you type `ls()` the function is evaluated and, as seen above, we see objects in the workspace.

Unlike `ls`, most functions require one or more arguments. Below is an example of how we assign an object to the argument of the function `log`. Remember that we earlier defined `a` to be 1:

```
log(8)
```

```
## [1] 2.079442
```

```
log(a)
```

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

You can find out what the function expects and what it does by reviewing the very useful manuals included in R. You can get help by using the help function like this:

```
help(log)
```

```
## starting httpd help server ... done
```

For most functions, we can also use this shorthand:

```
?log
```

The help page will show you what arguments the function is expecting. For example, `log` needs `x` and `base` to run. However, some arguments are required and others are optional. You can determine which arguments are optional by noting in the help document that a default value is assigned with `=`. Defining these is optional. For example, the base of the function `log` defaults to `base = exp(1)` making `log` the natural log by default.

If you want a quick look at the arguments without opening the help system, you can type:

```
args(log)
```

```
## function (x, base = exp(1))  
## NULL
```

You can change the default values by simply assigning another object:

```
log(8, base=2)
```

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

Note that we have not been specifying the argument x as such:

```
log(x=8, base=2)
```

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

The above code works, but we can save ourselves some typing: if no argument name is used, R assumes you are entering arguments in the order shown in the help file or by args. So by not using the names, it assumes the arguments are x followed by base:

```
log(8,2)
```

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

If using the arguments' names, then we can include them in whatever order we want:

```
log(base=2, x=8)
```

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

To specify arguments, we must use =, and cannot use <-.

There are some exceptions to the rule that functions need the parentheses to be evaluated. Among these, the most commonly used are the arithmetic and relational operators. For example:

```
2^3
```

```
## [1] 8
```

You can see the arithmetic operators by typing:

```
help("+")
```

or

```
?"+"
```

and the relational operators by typing:

```
help(">")
```

or

```
?">"
```

## 2.2.4 Other prebuilt objects

There are several datasets that are included for users to practice and test out functions. You can see all the available datasets by typing:

```
data()
```

This shows you the object name for these datasets. These datasets are objects that can be used by simply typing the name. For example, if you type:

```
co2
```

##		Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct
##	1959	315.42	316.31	316.50	317.56	318.13	318.00	316.39	314.65	313.68	313.18
##	1960	316.27	316.81	317.42	318.87	319.87	319.43	318.01	315.74	314.00	313.68
##	1961	316.73	317.54	318.38	319.31	320.42	319.61	318.42	316.63	314.83	315.16
##	1962	317.78	318.40	319.53	320.42	320.85	320.45	319.45	317.25	316.11	315.27
##	1963	318.58	318.92	319.70	321.22	322.08	321.31	319.58	317.61	316.05	315.83
##	1964	319.41	320.07	320.74	321.40	322.06	321.73	320.27	318.54	316.54	316.71
##	1965	319.27	320.28	320.73	321.97	322.00	321.71	321.05	318.71	317.66	317.14
##	1966	320.46	321.43	322.23	323.54	323.91	323.59	322.24	320.20	318.48	317.94
##	1967	322.17	322.34	322.88	324.25	324.83	323.93	322.38	320.76	319.10	319.24
##	1968	322.40	322.99	323.73	324.86	325.40	325.20	323.98	321.95	320.18	320.09
##	1969	323.83	324.26	325.47	326.50	327.21	326.54	325.72	323.50	322.22	321.62
##	1970	324.89	325.82	326.77	327.97	327.91	327.50	326.18	324.53	322.93	322.90
##	1971	326.01	326.51	327.01	327.62	328.76	328.40	327.20	325.27	323.20	323.40
##	1972	326.60	327.47	327.58	329.56	329.90	328.92	327.88	326.16	324.68	325.04
##	1973	328.37	329.40	330.14	331.33	332.31	331.90	330.70	329.15	327.35	327.02
##	1974	329.18	330.55	331.32	332.48	332.92	332.08	331.01	329.23	327.27	327.21
##	1975	330.23	331.25	331.87	333.14	333.80	333.43	331.73	329.90	328.40	328.17
##	1976	331.58	332.39	333.33	334.41	334.71	334.17	332.89	330.77	329.14	328.78
##	1977	332.75	333.24	334.53	335.90	336.57	336.10	334.76	332.59	331.42	330.98
##	1978	334.80	335.22	336.47	337.59	337.84	337.72	336.37	334.51	332.60	332.38
##	1979	336.05	336.59	337.79	338.71	339.30	339.12	337.56	335.92	333.75	333.70
##	1980	337.84	338.19	339.91	340.60	341.29	341.00	339.39	337.43	335.72	335.84
##	1981	339.06	340.30	341.21	342.33	342.74	342.08	340.32	338.26	336.52	336.68
##	1982	340.57	341.44	342.53	343.39	343.96	343.18	341.88	339.65	337.81	337.69
##	1983	341.20	342.35	342.93	344.77	345.58	345.14	343.81	342.21	339.69	339.82
##	1984	343.52	344.33	345.11	346.88	347.25	346.62	345.22	343.11	340.90	341.18
##	1985	344.79	345.82	347.25	348.17	348.74	348.07	346.38	344.51	342.92	342.62
##	1986	346.11	346.78	347.68	349.37	350.03	349.37	347.76	345.73	344.68	343.99

```

## 1987 347.84 348.29 349.23 350.80 351.66 351.07 349.33 347.92 346.27 346.18
## 1988 350.25 351.54 352.05 353.41 354.04 353.62 352.22 350.27 348.55 348.72
## 1989 352.60 352.92 353.53 355.26 355.52 354.97 353.75 351.52 349.64 349.83
## 1990 353.50 354.55 355.23 356.04 357.00 356.07 354.67 352.76 350.82 351.04
## 1991 354.59 355.63 357.03 358.48 359.22 358.12 356.06 353.92 352.05 352.11
## 1992 355.88 356.63 357.72 359.07 359.58 359.17 356.94 354.92 352.94 353.23
## 1993 356.63 357.10 358.32 359.41 360.23 359.55 357.53 355.48 353.67 353.95
## 1994 358.34 358.89 359.95 361.25 361.67 360.94 359.55 357.49 355.84 356.00
## 1995 359.98 361.03 361.66 363.48 363.82 363.30 361.94 359.50 358.11 357.80
## 1996 362.09 363.29 364.06 364.76 365.45 365.01 363.70 361.54 359.51 359.65
## 1997 363.23 364.06 364.61 366.40 366.84 365.68 364.52 362.57 360.24 360.83
##      Nov      Dec
## 1959 314.66 315.43
## 1960 314.84 316.03
## 1961 315.94 316.85
## 1962 316.53 317.53
## 1963 316.91 318.20
## 1964 317.53 318.55
## 1965 318.70 319.25
## 1966 319.63 320.87
## 1967 320.56 321.80
## 1968 321.16 322.74
## 1969 322.69 323.95
## 1970 323.85 324.96
## 1971 324.63 325.85
## 1972 326.34 327.39
## 1973 327.99 328.48
## 1974 328.29 329.41
## 1975 329.32 330.59
## 1976 330.14 331.52
## 1977 332.24 333.68
## 1978 333.75 334.78
## 1979 335.12 336.56
## 1980 336.93 338.04
## 1981 338.19 339.44
## 1982 339.09 340.32
## 1983 340.98 342.82
## 1984 342.80 344.04
## 1985 344.06 345.38
## 1986 345.48 346.72
## 1987 347.64 348.78
## 1988 349.91 351.18
## 1989 351.14 352.37
## 1990 352.69 354.07
## 1991 353.64 354.89
## 1992 354.09 355.33
## 1993 355.30 356.78
## 1994 357.59 359.05
## 1995 359.61 360.74
## 1996 360.80 362.38
## 1997 362.49 364.34

```

R will show you Mauna Loa atmospheric CO2 concentration data.

Other prebuilt objects are mathematical quantities, such as the constant

and  $\infty$ :

```
pi
```

```
## [1] 3.141593
```

```
Inf+1
```

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

### 2.2.5 Variable names

We have used the letters a, b, and c as variable names, but variable names can be almost anything. Some basic rules in R are that variable names have to start with a letter, can't contain spaces, and should not be variables that are predefined in R. For example, don't name one of your variables `install.packages` by typing something like `install.packages <- 2`.

A nice convention to follow is to use meaningful words that describe what is stored, use only lower case, and use underscores as a substitute for spaces. For the quadratic equations, we could use something like this:

```
solution_1 <- (-b + sqrt(b^2 - 4*a*c)) / (2*a)
solution_2 <- (-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

For more advice, we highly recommend studying Hadley Wickham's style guide<sup>15</sup>.

### 2.2.6 Saving your workspace

Values remain in the workspace until you end your session or erase them with the function `rm`. But workspaces also can be saved for later use. In fact, when you quit R, the program asks you if you want to save your workspace. If you do save it, the next time you start R, the program will restore the workspace.

We actually recommend against saving the workspace this way because, as you start working on different projects, it will become harder to keep track of what is saved. Instead, we recommend you assign the workspace a specific name. You can do this by using the function `save` or `save.image`. To load, use the function `load`. When saving a workspace, we recommend the suffix `rda` or `RData`. In RStudio, you can also do this by navigating to the Session tab and choosing Save Workspace as. You can later load it using the Load Workspace options in the same tab. You can read the help pages on `save`, `save.image`, and `load` to learn more.

### 2.2.7 Motivating scripts

To solve another equation such as  $3x^2 + 2x - 1$ , we can copy and paste the code above and then redefine the variables and recompute the solution:

```
a <- 3
b <- 2
c <- -1
(-b + sqrt(b^2 - 4*a*c)) / (2*a)
```

```
## [1] 0.3333333
```

```
(-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

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

By creating and saving a script with the code above, we would not need to retype everything each time and, instead, simply change the variable names. Try writing the script above into an editor and notice how easy it is to change the variables and receive an answer.

### 2.2.8 Commenting your code

If a line of R code starts with the symbol `#`, it is not evaluated. We can use this to write reminders of why we wrote particular code. For example, in the script above we could add:

```
## Code to compute solution to quadratic equation of the form ax^2 + bx + c  
## define the variables  
a <- 3  
b <- 2  
c <- -1  
  
## now compute the solution  
(-b + sqrt(b^2 - 4*a*c)) / (2*a)
```

```
## [1] 0.3333333
```

```
(-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

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

## 2.3 Exercises

1. What is the sum of the first 100 positive integers? The formula for the sum of integers 1 through  $n$  is  $n(n+1)/2$ . Define  $n = 100$  and then use R to compute the sum of 1 through 100 using the formula. What is the sum?

```
n <- 100  
n*(n+1)/2
```

```
## [1] 5050
```

2. Now use the same formula to compute the sum of the integers from 1 through 1,000.

```
n <- 1000  
n*(n+1)/2
```

```
## [1] 500500
```

3. Look at the result of typing the following code into R:



```
n <- 1000
x <- seq(1, n)
sum(x)
```

```
## [1] 500500
```

Based on the result, what do you think the functions `seq` and `sum` do? You can use `help`. a. `sum` creates a list of numbers and `seq` adds them up. b. `seq` creates a list of numbers and `sum` adds them up. c. `seq` creates a random list and `sum` computes the sum of 1 through 1,000. d. `sum` always returns the same number.

-> my answer: “b”

4. In math and programming, we say that we evaluate a function when we replace the argument with a given number. So if we type `sqrt(4)`, we evaluate the `sqrt` function. In R, you can evaluate a function inside another function. The evaluations happen from the inside out. Use one line of code to compute the log, in base 10, of the square root of 100.

```
log(sqrt(100), 10)
```

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

5. Which of the following will always return the numeric value stored in `x`? You can try out examples and use the help system if you want.

-> my prediction: “a”

(attempts)

```
x <- 123
log(10^x)
```

```
## [1] 283.218
```

```
log10(x^10)
```

```
## [1] 20.89905
```

```
log(exp(x))
```

```
## [1] 123
```

```
exp(log(x, base = 2))
```

```
## [1] 1035.37
```

-> answer: “c”

## 2.4 Data types

Variables in R can be of different types. For example, we need to distinguish numbers from character strings and tables from simple lists of numbers. The function `class` helps us determine what type of object we have:

```
a <- 2
class(a)
```

```
## [1] "numeric"
```

To work efficiently in R, it is important to learn the different types of variables and what we can do with these.

### 2.4.1 Data frames

Up to now, the variables we have defined are just one number. This is not very useful for storing data. The most common way of storing a dataset in R is in a data frame. Conceptually, we can think of a data frame as a table with rows representing observations and the different variables reported for each observation defining the columns. Data frames are particularly useful for datasets because we can combine different data types into one object.

A large proportion of data analysis challenges start with data stored in a data frame. For example, we stored the data for our motivating example in a data frame. You can access this dataset by loading the `dslabs` library and loading the `murders` dataset using the `data` function:

```
library(dslabs)
data(murders)
```

To see that this is in fact a data frame, we type:

```
class(murders)
```

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

### 2.4.2 Examining an object

The function `str` is useful for finding out more about the structure of an object:

```
str(murders)
```

```
## 'data.frame':   51 obs. of  5 variables:
## $ state      : chr  "Alabama" "Alaska" "Arizona" "Arkansas" ...
## $ abb       : chr  "AL" "AK" "AZ" "AR" ...
## $ region    : Factor w/ 4 levels "Northeast","South",...: 2 4 4 2 4 4 1 2 2 2 ...
## $ population: num  4779736 710231 6392017 2915918 37253956 ...
## $ total     : num   135  19  232  93 1257 ...
```

This tells us much more about the object. We see that the table has 51 rows (50 states plus DC) and five variables. We can show the first six lines using the function `head`:

```
head(murders)
```

```
##      state abb region population total
## 1  Alabama AL  South   4779736   135
## 2  Alaska  AK   West    710231    19
## 3  Arizona AZ   West   6392017   232
## 4  Arkansas AR  South   2915918    93
## 5 California CA  West  37253956  1257
## 6  Colorado CO   West   5029196    65
```

In this dataset, each state is considered an observation and five variables are reported for each state.

Before we go any further in answering our original question about different states, let's learn more about the components of this object.

### 2.4.3 The accessor: \$

For our analysis, we will need to access the different variables represented by columns included in this data frame. To do this, we use the accessor operator `$` in the following way:

```
murders$population
```

```
## [1] 4779736 710231 6392017 2915918 37253956 5029196 3574097 897934
## [9] 601723 19687653 9920000 1360301 1567582 12830632 6483802 3046355
## [17] 2853118 4339367 4533372 1328361 5773552 6547629 9883640 5303925
## [25] 2967297 5988927 989415 1826341 2700551 1316470 8791894 2059179
## [33] 19378102 9535483 672591 11536504 3751351 3831074 12702379 1052567
## [41] 4625364 814180 6346105 25145561 2763885 625741 8001024 6724540
## [49] 1852994 5686986 563626
```

But how did we know to use `population`? Previously, by applying the function `str` to the object `murders`, we revealed the names for each of the five variables stored in this table. We can quickly access the variable names using:

```
names(murders)
```

```
## [1] "state"      "abb"        "region"     "population" "total"
```

It is important to know that the order of the entries in `murders$population` preserves the order of the rows in our data table. This will later permit us to manipulate one variable based on the results of another. For example, we will be able to order the state names by the number of murders.

Tip: R comes with a very nice auto-complete functionality that saves us the trouble of typing out all the names. Try typing `murders$p` then hitting the tab key on your keyboard. This functionality and many other useful auto-complete features are available when working in RStudio.

### 2.4.4 Vectors: numerics, characters, and logical

The object `murders$population` is not one number but several. We call these types of objects vectors. A single number is technically a vector of length 1, but in general we use the term vectors to refer to objects with several entries. The function `length` tells you how many entries are in the vector:

```
pop <- murders$population
length(pop)
```

```
## [1] 51
```

This particular vector is numeric since population sizes are numbers:

```
class(pop)
```

```
## [1] "numeric"
```

In a numeric vector, every entry must be a number.

To store character strings, vectors can also be of class character. For example, the state names are characters:

```
class(murders$state)
```

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

As with numeric vectors, all entries in a character vector need to be a character.

Another important type of vectors are logical vectors. These must be either TRUE or FALSE.

```
z <- 3 == 2
z
```

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

```
class(z)
```

```
## [1] "logical"
```

Here the == is a relational operator asking if 3 is equal to 2. In R, if you just use one =, you actually assign a variable, but if you use two == you test for equality.

You can see the other relational operators by typing:

```
?Comparison
```

In future sections, you will see how useful relational operators can be.

We discuss more important features of vectors after the next set of exercises.

Advanced: Mathematically, the values in pop are integers and there is an integer class in R. However, by default, numbers are assigned class numeric even when they are round integers. For example, class(1) returns numeric. You can turn them into class integer with the as.integer() function or by adding an L like this: 1L. Note the class by typing: class(1L)

```
class(1)
```

```
## [1] "numeric"
```

```
as.integer()
```

```
## integer(0)
```

```
class(1)
```

```
## [1] "numeric"
```

```
class(1L)
```

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

### 2.4.5 Factors

In the murders dataset, we might expect the region to also be a character vector. However, it is not:

```
class(murders$region)
```

```
## [1] "factor"
```

It is a factor. Factors are useful for storing categorical data. We can see that there are only 4 regions by using the levels function:

```
levels(murders$region)
```

```
## [1] "Northeast"      "South"           "North Central"  "West"
```

In the background, R stores these levels as integers and keeps a map to keep track of the labels. This is more memory efficient than storing all the characters.

Note that the levels have an order that is different from the order of appearance in the factor object. The default in R is for the levels to follow alphabetical order. However, often we want the levels to follow a different order. You can specify an order through the levels argument when creating the factor with the factor function. For example, in the murders dataset regions are ordered from east to west. The function reorder lets us change the order of the levels of a factor variable based on a summary computed on a numeric vector. We will demonstrate this with a simple example, and will see more advanced ones in the Data Visualization part of the book.

Suppose we want the levels of the region by the total number of murders rather than alphabetical order. If there are values associated with each level, we can use the reorder and specify a data summary to determine the order. The following code takes the sum of the total murders in each region, and reorders the factor following these sums.

```
region <- murders$region
value <- murders$total
region <- reorder(region, value, FUN = sum)
levels(region)
```

```
## [1] "Northeast"      "North Central"  "West"           "South"
```

The new order is in agreement with the fact that the Northeast has the least murders and the South has the most.

Warning: Factors can be a source of confusion since sometimes they behave like characters and sometimes they do not. As a result, confusing factors and characters are a common source of bugs.

### 2.4.6 Lists

Data frames are a special case of lists. Lists are useful because you can store any combination of different types. You can create a list using the list function like this:

```
record <- list(name = "John Doe",
              student_id = 1234,
              grades = c(95, 82, 91, 97, 93),
              final_grade = "A")
```

The function c is described in Section 2.6.

This list includes a character, a number, a vector with five numbers, and another character.

```
record
```

```
## $name
## [1] "John Doe"
##
## $student_id
## [1] 1234
##
## $grades
## [1] 95 82 91 97 93
##
## $final_grade
## [1] "A"
```

```
class(record)
```

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

As with data frames, you can extract the components of a list with the accessor \$.

```
record$student_id
```

```
## [1] 1234
```

We can also use double square brackets ([[]]) like this:

```
record[["student_id"]]
```

```
## [1] 1234
```

You should get used to the fact that in R, there are often several ways to do the same thing, such as accessing entries.

You might also encounter lists without variable names.

```
record2 <- list("John Doe", 1234)
record2
```

```
## [[1]]
## [1] "John Doe"
##
## [[2]]
## [1] 1234
```

If a list does not have names, you cannot extract the elements with \$, but you can still use the brackets method and instead of providing the variable name, you provide the list index, like this:

```
record2[[1]]
```

```
## [1] "John Doe"
```

We won't be using lists until later, but you might encounter one in your own exploration of R. For this reason, we show you some basics here.

### 2.4.7 Matrices

Matrices are another type of object that are common in R. Matrices are similar to data frames in that they are two-dimensional: they have rows and columns. However, like numeric, character and logical vectors, entries in matrices have to be all the same type. For this reason data frames are much more useful for storing data, since we can have characters, factors, and numbers in them.

Yet matrices have a major advantage over data frames: we can perform matrix algebra operations, a powerful type of mathematical technique. We do not describe these operations in this book, but much of what happens in the background when you perform a data analysis involves matrices. We cover matrices in more detail in Chapter 33.1 but describe them briefly here since some of the functions we will learn return matrices.

We can define a matrix using the matrix function. We need to specify the number of rows and columns.

```
mat <- matrix(1:12, 4, 3)
mat
```

```
##      [,1] [,2] [,3]
## [1,]    1    5    9
## [2,]    2    6   10
## [3,]    3    7   11
## [4,]    4    8   12
```

You can access specific entries in a matrix using square brackets ([]). If you want the second row, third column, you use:

```
mat[2,3]
```

```
## [1] 10
```

If you want the entire second row, you leave the column spot empty:

```
mat[2, ]
```

```
## [1] 2 6 10
```

Notice that this returns a vector, not a matrix.

Similarly, if you want the entire third column, you leave the row spot empty:

```
mat[,3]
```

```
## [1] 9 10 11 12
```

This is also a vector, not a matrix.

You can access more than one column or more than one row if you like. This will give you a new matrix.

```
mat[, 2:3]
```

```
##      [,1] [,2]
## [1,]    5    9
## [2,]    6   10
## [3,]    7   11
## [4,]    8   12
```

You can subset both rows and columns:

```
mat[1:2, 2:3]
```

```
##      [,1] [,2]
## [1,]    5    9
## [2,]    6   10
```

We can convert matrices into data frames using the function `as.data.frame`:

```
as.data.frame(mat)
```

```
##   V1 V2 V3
## 1  1  5  9
## 2  2  6 10
## 3  3  7 11
## 4  4  8 12
```

You can also use single square brackets (`()`) to access rows and columns of a data frame:

```
data("murders")
murders[25,1]
```

```
## [1] "Mississippi"
```



```
murders[2:3, ]
```

```
##      state abb region population total
## 2  Alaska  AK   West      710231     19
## 3 Arizona  AZ   West     6392017    232
```

## 2.5 Exercises

1. Load the US murders dataset.

```
library(dslabs)
data(murders)
```

Use the function `str` to examine the structure of the `murders` object. Which of the following best describes the variables represented in this data frame?

```
str(murders)
```

```
## 'data.frame':   51 obs. of  5 variables:
## $ state      : chr  "Alabama" "Alaska" "Arizona" "Arkansas" ...
## $ abb       : chr  "AL" "AK" "AZ" "AR" ...
## $ region    : Factor w/ 4 levels "Northeast","South",...: 2 4 4 2 4 4 1 2 2 2 ...
## $ population: num  4779736 710231 6392017 2915918 37253956 ...
## $ total     : num  135 19 232 93 1257 ...
```

- a. The 51 states.
- b. The murder rates for all 50 states and DC.
- c. The state name, the abbreviation of the state name, the state's region, and the state's population and total number of murders for 2010.
- d. `str` shows no relevant information.

-> answer: "c"

2. What are the column names used by the data frame for these five variables?

-> answer: state, abb, region, population, total

3. Use the accessor `$` to extract the state abbreviations and assign them to the object `a`. What is the class of this object?

```
a <- murders$abb
a
```

```
## [1] "AL" "AK" "AZ" "AR" "CA" "CO" "CT" "DE" "DC" "FL" "GA" "HI" "ID" "IL" "IN"
## [16] "IA" "KS" "KY" "LA" "ME" "MD" "MA" "MI" "MN" "MS" "MO" "MT" "NE" "NV" "NH"
## [31] "NJ" "NM" "NY" "NC" "ND" "OH" "OK" "OR" "PA" "RI" "SC" "SD" "TN" "TX" "UT"
## [46] "VT" "VA" "WA" "WV" "WI" "WY"
```

```
class(a)
```

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

4. Now use the square brackets to extract the state abbreviations and assign them to the object b. Use the identical function to determine if a and b are the same.

```
murders[, 2]
```

```
## [1] "AL" "AK" "AZ" "AR" "CA" "CO" "CT" "DE" "DC" "FL" "GA" "HI" "ID" "IL" "IN"  
## [16] "IA" "KS" "KY" "LA" "ME" "MD" "MA" "MI" "MN" "MS" "MO" "MT" "NE" "NV" "NH"  
## [31] "NJ" "NM" "NY" "NC" "ND" "OH" "OK" "OR" "PA" "RI" "SC" "SD" "TN" "TX" "UT"  
## [46] "VT" "VA" "WA" "WV" "WI" "WY"
```

```
b <- murders[, 2]  
a == b
```

```
## [1] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE  
## [16] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE  
## [31] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE  
## [46] TRUE TRUE TRUE TRUE TRUE TRUE
```

5. We saw that the region column stores a factor. You can corroborate this by typing:

```
class(murders$region)
```

```
## [1] "factor"
```

with one line of code, use the function levels and length to determine the number of regions defined by this dataset.

```
length(levels(murders$region))
```

```
## [1] 4
```

6. The function table takes a vector and returns the frequency of each element. You can quickly see how many states are in each region by applying this function. Use this function in one line of code to create a table of states per region.

```
table(murders$region)
```

```
##  
## Northeast South North Central West  
##          9      17          12      13
```

## 2.6 Vectors

In R, the most basic objects available to store data are vectors. As we have seen, complex datasets can usually be broken down into components that are vectors. For example, in a data frame, each column is a vector. Here we learn more about this important class.

### 2.6.1 Creating vectors

We can create vectors using the function `c`, which stands for concatenate. We use `c` to concatenate entries in the following way:

```
codes <- c(380, 124, 818)
codes
```

```
## [1] 380 124 818
```

We can also create character vectors. We use the quotes to denote that the entries are characters rather than variable names.

```
country <- c("italy", "canada", "egypt")
```

In R you can also use single quotes:

```
country <- c('italy', 'canada', 'egypt')
```

But be careful not to confuse the single quote `'` with the back quote ```.

By now you should know that if you type:

```
country <- c(italy, canada, egypt)
```

you receive an error because the variables `italy`, `canada`, and `egypt` are not defined. If we do not use the quotes, R looks for variables with those names and returns an error.

### 2.6.2 Names

Sometimes it is useful to name the entries of a vector. For example, when defining a vector of country codes, we can use the names to connect the two:

```
codes <- c(italy = 380, canada = 124, egypt = 818)
codes
```

```
##   italy canada  egypt
##    380    124    818
```

The object `codes` continues to be a numeric vector:

```
class(codes)
```

```
## [1] "numeric"
```

but with names:

```
names(codes)
```

```
## [1] "italy" "canada" "egypt"
```

If the use of strings without quotes looks confusing, know that you can use the quotes as well:

```
codes <- c("italy" = 380, "canada" = 124, "egypt" = 818)
codes
```

```
## italy canada egypt
##    380    124    818
```

There is no difference between this function call and the previous one. This is one of the many ways in which R is quirky compared to other languages.

We can also assign names using the `names` functions:

```
codes <- c(380, 124, 818)
country <- c("italy", "canada", "egypt")
names(codes) <- country
codes
```

```
## italy canada egypt
##    380    124    818
```

### 2.6.3 Sequences

Another useful function for creating vectors generates sequences:

```
seq(1,10)
```

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

The first argument defines the start, and the second defines the end which is included. The default is to go up in increments of 1, but a third argument lets us tell it how much to jump by:

```
seq(1,10,2)
```

```
## [1] 1 3 5 7 9
```

If we want consecutive integers, we can use the following shorthand:

```
1:10
```

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

When we use these functions, R produces integers, not numerics, because they are typically used to index something:

```
class(1:10)
```

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

However, if we create a sequence including non-integers, the class changes:

```
class(seq(1,10,0.5))
```

```
## [1] "numeric"
```

#### 2.6.4 Subsetting

We use square brackets to access specific elements of a vector. For the vector codes we defined above, we can access the second element using:

```
codes[2]
```

```
## canada  
##      124
```

You can get more than one entry by using a multi-entry vector as an index:

```
codes[c(1,3)]
```

```
## italy egypt  
##    380   818
```

The sequences defined above are particularly useful if we want to access, say, the first two elements:

```
codes[1:2]
```

```
## italy canada  
##    380   124
```

If the elements have names, we can also access the entries using these names. Below are two examples.

```
codes["canada"]
```

```
## canada  
##      124
```

```
codes[c("egypt", "italy")]
```

```
## egypt italy  
##    818   380
```