

Chapter 1

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. 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. 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 Countries¹⁴ 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 $a x^2 + b x + c = 0$. The quadratic formula gives us the solutions:

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

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 $x^2 + x - 1 = 0$, then we define:

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: `print(a)`

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

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:

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:

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:

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:

For most functions, we can also use this shorthand:

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:

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

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

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`:

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

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:

You can see the arithmetic operators by typing:

or

and the relational operators by typing:

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:

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:

R will show you Mauna Loa atmospheric CO2 concentration data.

Other prebuilt objects are mathematical quantities, such as the constant

π and

∞ :

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:

For more advice, we highly recommend studying Hadley Wickham's style guide¹⁵.

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

$3 \times 2 + 2 \times -1$, we can copy and paste the code above and then redefine the variables and recompute the solution:

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:

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?
2. Now use the same formula to compute the sum of the integers from 1 through 1,000.
3. Look at the result of typing the following code into R:

Based on the result, what do you think the functions `seq` and `sum` do? You can use `help`.

`sum` creates a list of numbers and `seq` adds them up. `seq` creates a list of numbers and `sum` adds them up. `seq` creates a random list and `sum` computes the sum of 1 through 1,000. `sum` always returns the same number. 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.

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.

```
log(10^x) log10(x^10) log(exp(x)) exp(log(x, base = 2))
```

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:

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

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

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`:

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:

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:

[Hide](#)

```
names(murders)
```

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:

[Hide](#)

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

This particular vector is numeric since population sizes are numbers:

[Hide](#)

```
class(pop)
```

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:

[Hide](#)

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

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.

[Hide](#)

```
z <- 3 == 2  
z  
class(z)
```

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:

[Hide](#)

```
?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)`

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

[Hide](#)

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

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:

[Hide](#)

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

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.

[Hide](#)

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

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: The function c is described in Section 2.6.

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

[Hide](#)

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

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

[Hide](#)

```
record$student_id
```

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

[Hide](#)

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

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.

[Hide](#)

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

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:

[Hide](#)

```
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.

[Hide](#)

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

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

[Hide](#)

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

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

[Hide](#)

```
mat[2, ]
```

Notice that this returns a vector, not a matrix.

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

[Hide](#)

```
mat[, 3]
```

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.

[Hide](#)

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

You can subset both rows and columns:

[Hide](#)

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

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

[Hide](#)

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

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

[Hide](#)

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

murders[2:3, ]
```

#####2.5 Exercises

1. Load the US murders dataset. 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?

The 51 states. The murder rates for all 50 states and DC. 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. str shows no relevant information.

[Hide](#)

```
library(dslabs)
data(murders)
str(murders)
#c
```

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

[Hide](#)

```
names(murders)
```

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

[Hide](#)

```
a <- murders$abb
class(a)
```

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.

[Hide](#)

```
b <- murders[["abb"]]
z <- a==b
z
```

5. We saw that the region column stores a factor. You can corroborate this by typing: With one line of code, use the function levels and length to determine the number of regions defined by this dataset.

[Hide](#)

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

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.

[Hide](#)

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

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:

[Hide](#)

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

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

[Hide](#)

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

In R you can also use single quotes:

[Hide](#)

```
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:

[Hide](#)

```
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:

[Hide](#)

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

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

[Hide](#)

```
class(codes)
```

but with names:

[Hide](#)

```
names(codes)
```

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

[Hide](#)

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

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:

Hide

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

2.6.3 Sequences Another useful function for creating vectors generates sequences:

Hide

```
seq(1, 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:

Hide

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

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

Hide

```
1:10
```

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

Hide

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

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

Hide

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

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:

Hide

```
codes[2]
```

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

Hide

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

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

Hide

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

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

[Hide](#)

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

2.7 Coercion

In general, coercion is an attempt by R to be flexible with data types. When an entry does not match the expected, some of the prebuilt R functions try to guess what was meant before throwing an error. This can also lead to confusion. Failing to understand coercion can drive programmers crazy when attempting to code in R since it behaves quite differently from most other languages in this regard. Let's learn about it with some examples.

We said that vectors must be all of the same type. So if we try to combine, say, numbers and characters, you might expect an error:

[Hide](#)

```
x <- c(1, "canada", 3)
```

But we don't get one, not even a warning! What happened? Look at x and its class:

[Hide](#)

```
x  
class(x)
```

R coerced the data into characters. It guessed that because you put a character string in the vector, you meant the 1 and 3 to actually be character strings "1" and "3." The fact that not even a warning is issued is an example of how coercion can cause many unnoticed errors in R.

R also offers functions to change from one type to another. For example, you can turn numbers into characters with:

[Hide](#)

```
x <- 1:5  
y <- as.character(x)  
y
```

You can turn it back with `as.numeric`:

[Hide](#)

```
as.numeric(y)
```

This function is actually quite useful since datasets that include numbers as character strings are common.

2.7.1 Not availables (NA) When a function tries to coerce one type to another and encounters an impossible case, it usually gives us a warning and turns the entry into a special value called an NA for "not available." For example:

[Hide](#)

```
x <- c("1", "b", "3")  
as.numeric(x)
```

R does not have any guesses for what number you want when you type b, so it does not try.

As a data scientist you will encounter the NAs often as they are generally used for missing data, a common problem in real-world datasets.

2.8 Exercises

1. Use the function `c` to create a vector with the average high temperatures in January for Beijing, Lagos, Paris, Rio de Janeiro, San Juan, and Toronto, which are 35, 88, 42, 84, 81, and 30 degrees Fahrenheit. Call the object `temp`.

[Hide](#)

```
temp <- c('Beijing' = 35, 'Lagos' = 88, 'Paris' = 42, 'Rio de Janeiro' = 84, 'San Juan' = 81, 'Toronto' = 30)
```

2. Now create a vector with the city names and call the object `city`.

[Hide](#)

```
city <- c('Beijing', 'Lagos', 'Paris', 'Rio de Janeiro', 'San Juan', 'Toronto')
```

3. Use the `names` function and the objects defined in the previous exercises to associate the temperature data with its corresponding city.

[Hide](#)

```
names(temp) <- city  
temp
```

4. Use the `[` and `:` operators to access the temperature of the first three cities on the list.

[Hide](#)

```
temp[1:3]
```

5. Use the `[` operator to access the temperature of Paris and San Juan.

[Hide](#)

```
temp["Paris"]  
temp["San Juan"]
```

6. Use the `:` operator to create a sequence of numbers
12, 13, 14, ..., 73.

[Hide](#)

```
seq(12, 73)
```

7. Create a vector containing all the positive odd numbers smaller than 100.

[Hide](#)

```
c(seq(1, 99, 2))
```

8. Create a vector of numbers that starts at 6, does not pass 55, and adds numbers in increments of $4/7$: 6 , $6 + 4/7$, $6 + 8/7$, and so on. How many numbers does the list have? Hint: use `seq` and `length`.

Hide

```
length(seq(6,))
```

9. What is the class of the following object `a <- seq(1, 10, 0.5)`?

Hide

```
a <- seq(1, 10, 0.5)
class(a)
```

10. What is the class of the following object `a <- seq(1, 10)`?

Hide

```
a <- seq(1, 10)
class(a)
```

11. The class of `class(a<-1)` is numeric, not integer. R defaults to numeric and to force an integer, you need to add the letter L. Confirm that the class of 1L is integer.

Hide

```
class(a<-1)
class(a<-1L)
class(1L)
```

12. Define the following vector: and coerce it to get integers.

Hide

```
x <- c("1", "3", "5")
X <- as.integer
x
```

2.9 Sorting

Now that we have mastered some basic R knowledge, let's try to gain some insights into the safety of different states in the context of gun murders. 2.9.1 sort Say we want to rank the states from least to most gun murders. The function `sort` sorts a vector in increasing order. We can therefore see the largest number of gun murders by typing:

Hide

```
library(dslabs)
data(murders)
sort(murders$total)
```

However, this does not give us information about which states have which murder totals. For example, we don't know which state had 1257. 2.9.2 order The function `order` is closer to what we want. It takes a vector as input and returns the vector of indexes that sorts the input vector. This may sound confusing so let's look at a simple example. We can create a vector and sort it:

Hide

```
x <- c(31, 4, 15, 92, 65)
sort(x)
```

Rather than sort the input vector, the function `order` returns the index that sorts input vector:

[Hide](#)

```
index <- order(x)
x[index]
```

This is the same output as that returned by `sort(x)`. If we look at this index, we see why it works:

[Hide](#)

```
x
order(x)
```

The second entry of `x` is the smallest, so `order(x)` starts with 2. The next smallest is the third entry, so the second entry is 3 and so on.

How does this help us order the states by murders? First, remember that the entries of vectors you access with `$` follow the same order as the rows in the table. For example, these two vectors containing state names and abbreviations, respectively, are matched by their order:

[Hide](#)

```
murders$state[1:6]
murders$abb[1:6]
```

This means we can order the state names by their total murders. We first obtain the index that orders the vectors according to murder totals and then index the state names vector:

[Hide](#)

```
ind <- order(murders$total)
murders$abb[ind]
```

According to the above, California had the most murders.

2.9.3 `max` and `which.max` If we are only interested in the entry with the largest value, we can use `max` for the value:

[Hide](#)

```
max(murders$total)
#> [1] 1257
```

and `which.max` for the index of the largest value:

[Hide](#)

```
i_max <- which.max(murders$total)
murders$state[i_max]
```

For the minimum, we can use `min` and `which.min` in the same way.

Does this mean California is the most dangerous state? In an upcoming section, we argue that we should be considering rates instead of totals. Before doing that, we introduce one last order-related function: `rank`.

2.9.4 `rank` Although not as frequently used as `order` and `sort`, the function `rank` is also related to order and can be useful. For any given vector it returns a vector with the rank of the first entry, second entry, etc., of the input vector. Here is a simple example:

Hide

```
x <- c(31, 4, 15, 92, 65)
rank(x)
```

To summarize, let's look at the results of the three functions we have introduced:

original sort order rank 31 4 2 3 4 15 3 1 15 31 1 2 92 65 5 5 65 92 4 4

2.9.5 Beware of recycling Another common source of unnoticed errors in R is the use of recycling. We saw that vectors are added elementwise. So if the vectors don't match in length, it is natural to assume that we should get an error. But we don't. Notice what happens:

Hide

```
x <- c(1, 2, 3)
y <- c(10, 20, 30, 40, 50, 60, 70)
x+y
```

We do get a warning, but no error. For the output, R has recycled the numbers in x. Notice the last digit of numbers in the output.

2.10 Exercises For these exercises we will use the US murders dataset. Make sure you load it prior to starting.

Hide

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

1. Use the \$ operator to access the population size data and store it as the object pop. Then use the sort function to redefine pop so that it is sorted. Finally, use the [operator to report the smallest population size.

Hide

```
pop <- murders$population
sort(pop)
min(murders[["population"]])
```

2. Now instead of the smallest population size, find the index of the entry with the smallest population size. Hint: use order instead of sort.

Hide

```
x <- c(murders$population)
index <- order(x)
min(x[index])
```

3. We can actually perform the same operation as in the previous exercise using the function which.min. Write one line of code that does this.

Hide

```
i_min <- which.min(murders$population)
murders$population[i_min]
```

4. Now we know how small the smallest state is and we know which row represents it. Which state is it? Define a variable states to be the state names from the murders data frame. Report the name of the

state with the smallest population.

[Hide](#)

```
i_min <- which.min(murders$population)
murders$state[i_min]
```

5. You can create a data frame using the `data.frame` function. Here is a quick example:

[Hide](#)

```
temp <- c(35, 88, 42, 84, 81, 30)
city <- c("Beijing", "Lagos", "Paris", "Rio de Janeiro",
         "San Juan", "Toronto")
city_temps <- data.frame(name = city, temperature = temp)
```

Use the `rank` function to determine the population rank of each state from smallest population size to biggest. Save these ranks in an object called `ranks`, then create a data frame with the state name and its rank. Call the data frame `my_df`.

[Hide](#)

```
population <- c(murders$population)
state <- c(murders$state)
state_pop <- data.frame(names = state, population)
ranks <- rank(state_pop)
my_df <- data.frame(ranks)
```

6. Repeat the previous exercise, but this time order `my_df` so that the states are ordered from least populous to most populous. Hint: create an object `ind` that stores the indexes needed to order the population values. Then use the bracket operator `[]` to re-order each column in the data frame.

[Hide](#)

```
ind <- order(my_df)
my_df[ind]
```

7. The `na_example` vector represents a series of counts. You can quickly examine the object using:

[Hide](#)

```
data("na_example")
str(na_example)
```

However, when we compute the average with the function `mean`, we obtain an NA:

[Hide](#)

```
mean(na_example)
```

The `is.na` function returns a logical vector that tells us which entries are NA. Assign this logical vector to an object called `ind` and determine how many NAs does `na_example` have.

[Hide](#)

```
ind <- is.na(na_example)
num <- table(ind)
num[["TRUE"]]
```

8. Now compute the average again, but only for the entries that are not NA. Hint: remember the ! operator.

Hide

```
!
```

2.11 Vector arithmetics California had the most murders, but does this mean it is the most dangerous state? What if it just has many more people than any other state? We can quickly confirm that California indeed has the largest population:

Hide

```
library(dslabs)
data("murders")
murders$state[which.max(murders$population)]
#> [1] "California"
```

with over 37 million inhabitants. It is therefore unfair to compare the totals if we are interested in learning how safe the state is. What we really should be computing is the murders per capita. The reports we describe in the motivating section used murders per 100,000 as the unit. To compute this quantity, the powerful vector arithmetic capabilities of R come in handy.

2.11.1 Rescaling a vector In R, arithmetic operations on vectors occur element-wise. For a quick example, suppose we have height in inches:

Hide

```
inches <- c(69, 62, 66, 70, 70, 73, 67, 73, 67, 70)
```

and want to convert to centimeters. Notice what happens when we multiply inches by 2.54:

Hide

```
inches * 2.54
```

In the line above, we multiplied each element by 2.54. Similarly, if for each entry we want to compute how many inches taller or shorter than 69 inches, the average height for males, we can subtract it from every entry like this:

Hide

```
inches - 69
```

2.11.2 Two vectors If we have two vectors of the same length, and we sum them in R, they will be added entry by entry as follows:

$$\begin{pmatrix} | & | & | & | \end{pmatrix} \begin{pmatrix} a & b & c & d \end{pmatrix} + \begin{pmatrix} | & | & | & | \end{pmatrix} \begin{pmatrix} e & f & g & h \end{pmatrix} = \begin{pmatrix} | & | & | & | \end{pmatrix} \begin{pmatrix} a + e & b + f & c + g & d + h \end{pmatrix}$$

The same holds for other mathematical operations, such as -, * and /.

This implies that to compute the murder rates we can simply type:

Hide

```
murder_rate <- murders$total / murders$population * 100000
```

Once we do this, we notice that California is no longer near the top of the list. In fact, we can use what we have learned to order the states by murder rate:

Hide

```
murders$abb[order(murder_rate)]
```

2.12 Exercises 1. Previously we created this data frame:

Hide

```
temp <- c(35, 88, 42, 84, 81, 30)
city <- c("Beijing", "Lagos", "Paris", "Rio de Janeiro",
         "San Juan", "Toronto")
city_temps <- data.frame(name = city, temperature = temp)
```

Remake the data frame using the code above, but add a line that converts the temperature from Fahrenheit to Celsius. The conversion is

$$C = 5/9 \times (F - 32) .$$

Hide

```
tempF <- c(35, 88, 42, 84, 81, 30)
tempC <- 5/9*(tempF-32)
city <- c("Beijing", "Lagos", "Paris", "Rio de Janeiro",
         "San Juan", "Toronto")
city_temps <- data.frame(name = city, temperature = tempC)
```

2. What is the following sum

1

- $1/2^2$
- $1/3^2$
- ... $1/100^2$? Hint: thanks to Euler, we know it should be close to $\pi^2/6$.

Hide

```
x <- (1/n)^2
n <- 1:100
sum(x)
```

3. Compute the per 100,000 murder rate for each state and store it in the object murder_rate. Then compute the average murder rate for the US using the function mean. What is the average?

Hide

```
murder_rate <- murders$total / murders$population * 100000
sum(murder_rate)/51
```

2.13 Indexing R provides a powerful and convenient way of indexing vectors. We can, for example, subset a vector based on properties of another vector. In this section, we continue working with our US murders example, which we can load like this:

Hide

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

2.13.1 Subsetting with logicals We have now calculated the murder rate using:

Hide

```
murder_rate <- murders$total / murders$population * 100000
```

Imagine you are moving from Italy where, according to an ABC news report, the murder rate is only 0.71 per 100,000. You would prefer to move to a state with a similar murder rate. Another powerful feature of R is that we can use logicals to index vectors. If we compare a vector to a single number, it actually performs the test for each entry. The following is an example related to the question above:

Hide

```
ind <- murder_rate < 0.71
```

If we instead want to know if a value is less or equal, we can use:

Hide

```
ind <- murder_rate <= 0.71
```

Note that we get back a logical vector with TRUE for each entry smaller than or equal to 0.71. To see which states these are, we can leverage the fact that vectors can be indexed with logicals.

Hide

```
murders$state[ind]
```

In order to count how many are TRUE, the function `sum` returns the sum of the entries of a vector and logical vectors get coerced to numeric with TRUE coded as 1 and FALSE as 0. Thus we can count the states using:

Hide

```
sum(ind)
```

2.13.2 Logical operators Suppose we like the mountains and we want to move to a safe state in the western region of the country. We want the murder rate to be at most 1. In this case, we want two different things to be true. Here we can use the logical operator `and`, which in R is represented with `&`. This operation results in TRUE only when both logicals are TRUE. To see this, consider this example:

Hide

```
TRUE & TRUE  
TRUE & FALSE  
FALSE & FALSE
```

For our example, we can form two logicals:

Hide

```
west <- murders$region == "West"  
safe <- murder_rate <= 1
```

and we can use the `&` to get a vector of logicals that tells us which states satisfy both conditions:

Hide

```
ind <- safe & west  
murders$state[ind]
```

2.13.3 which Suppose we want to look up California's murder rate. For this type of operation, it is convenient to convert vectors of logicals into indexes instead of keeping long vectors of logicals. The function which tells us which entries of a logical vector are TRUE. So we can type:

Hide

```
ind <- which(murders$state == "California")
murder_rate[ind]
```

2.13.4 match If instead of just one state we want to find out the murder rates for several states, say New York, Florida, and Texas, we can use the function match. This function tells us which indexes of a second vector match each of the entries of a first vector:

Hide

```
ind <- match(c("New York", "Florida", "Texas"), murders$state)
ind
```

Now we can look at the murder rates:

Hide

```
murder_rate[ind]
```

2.13.5 %in% If rather than an index we want a logical that tells us whether or not each element of a first vector is in a second, we can use the function %in%. Let's imagine you are not sure if Boston, Dakota, and Washington are states. You can find out like this:

Hide

```
c("Boston", "Dakota", "Washington") %in% murders$state
```

Note that we will be using %in% often throughout the book.

Advanced: There is a connection between match and %in% through which. To see this, notice that the following two lines produce the same index (although in different order):

Hide

```
match(c("New York", "Florida", "Texas"), murders$state)
which(murders$state%in%c("New York", "Florida", "Texas"))
```

2.14 Exercises Start by loading the library and data.

Hide

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

1. Compute the per 100,000 murder rate for each state and store it in an object called murder_rate. Then use logical operators to create a logical vector named low that tells us which entries of murder_rate are lower than 1.

Hide

```
murder_rate <- murders$total / murders$population * 100000
ind <- murder_rate < 1
low <- c(murders$state[ind])
low
```

2. Now use the results from the previous exercise and the function which to determine the indices of murder_rate associated with values lower than 1.

[Hide](#)

```
ind <- which(murders$low)
```

3. Use the results from the previous exercise to report the names of the states with murder rates lower than 1.

[Hide](#)

```
low
```

4. Now extend the code from exercises 2 and 3 to report the states in the Northeast with murder rates lower than 1. Hint: use the previously defined logical vector low and the logical operator &.

[Hide](#)

```
Northeast <- murders$region == "Northeast"
low <- murder_rate < 1
ind <- Northeast & low
murders$state[ind]
```

5. In a previous exercise we computed the murder rate for each state and the average of these numbers. How many states are below the average?

[Hide](#)

```
murder_rate <- murders$total / murders$population * 100000
ind <- murder_rate < 2.78
sum(ind)
```

6. Use the match function to identify the states with abbreviations AK, MI, and IA. Hint: start by defining an index of the entries of murders\$abb that match the three abbreviations, then use the [operator to extract the states.

[Hide](#)

```
ind <- match(c("AK", "MI", "IA"), murders$abb)
murders$state[ind]
```

7. Use the %in% operator to create a logical vector that answers the question: which of the following are actual abbreviations: MA, ME, MI, MO, MU?

[Hide](#)

```
c("MA", "ME", "MI", "MO", "MU") %in% murders$abb
```

8. Extend the code you used in exercise 7 to report the one entry that is not an actual abbreviation. Hint: use the ! operator, which turns FALSE into TRUE and vice versa, then which to obtain an index.

2.15 Basic plots In Chapter 7 we describe an add-on package that provides a powerful approach to producing plots in R. We then have an entire part on Data Visualization in which we provide many examples. Here we briefly describe some of the functions that are available in a basic R installation.

2.15.1 plot The plot function can be used to make scatterplots. Here is a plot of total murders versus population.

[Hide](#)

```
x <- murders$population / 10^6
y <- murders$total
plot(x, y)
```

For a quick plot that avoids accessing variables twice, we can use the with function:

[Hide](#)

```
with(murders, plot(population, total))
```

The function with lets us use the murders column names in the plot function. It also works with any data frames and any function.

2.15.2 hist We will describe histograms as they relate to distributions in the Data Visualization part of the book. Here we will simply note that histograms are a powerful graphical summary of a list of numbers that gives you a general overview of the types of values you have. We can make a histogram of our murder rates by simply typing:

[Hide](#)

```
x <- with(murders, total / population * 100000)
hist(x)
```

We can see that there is a wide range of values with most of them between 2 and 3 and one very extreme case with a murder rate of more than 15:

[Hide](#)

```
murders$state[which.max(x)]
```

2.15.3 boxplot Boxplots will also be described in the Data Visualization part of the book. They provide a more terse summary than histograms, but they are easier to stack with other boxplots. For example, here we can use them to compare the different regions:

[Hide](#)

```
murders$rate <- with(murders, total / population * 100000)
boxplot(rate~region, data = murders)
```

We can see that the South has higher murder rates than the other three regions.

2.15.4 image The image function displays the values in a matrix using color. Here is a quick example:

[Hide](#)

```
x <- matrix(1:120, 12, 10)
image(x)
```

2.16 Exercises 1. We made a plot of total murders versus population and noted a strong relationship. Not surprisingly, states with larger populations had more murders.

[Hide](#)

```
library(dslabs)
data(murders)
population_in_millions <- murders$population/10^6
total_gun_murders <- murders$total
plot(population_in_millions, total_gun_murders)
```

Keep in mind that many states have populations below 5 million and are bunched up. We may gain further insights from making this plot in the log scale. Transform the variables using the log10 transformation and then plot them.

[Hide](#)

```
population_in_millions <- murders$population/10^6
total_gun_murders <- log10(murders$total)
plot(population_in_millions, total_gun_murders)
```

2. Create a histogram of the state populations.

[Hide](#)

```
x <- c(murders$population/10^6)
y <- c(murders$state)

hist(x)
```

3. Generate boxplots of the state populations by region.

[Hide](#)

```
Northeast <- length(levels(murders$region[Northeast]))
South <- length(levels(murders$region[south]))
North_Central <- length(levels(murders$region[North Central]))
West <- length(levels(murders$region[West]))
```