

Introduction to Probability and Statistics

2022/23

Stephen Connor

Table of contents

| | |
|---|----------|
| Introduction | 3 |
| Computer labs | 4 |
| Computer Lab 1 | 5 |
| Introduction to R and RStudio | 5 |
| The Data: Dr. Arbuthnot's Baptism Records | 7 |
| Some Exploration | 8 |
| A newer data set | 12 |
| Simulating a probability experiment | 13 |
| On your own computer | 14 |

Introduction

Important

You are reading the work-in-progress lecture notes for Introduction to Probability and Statistics. This section is far from finished, and I don't recommend reading it.

Warning

Don't use any numbered environments before we get to the first numbered section: mismatch between html and pdf numbering!

This is a Quarto book.

To learn more about Quarto books visit <https://quarto.org/docs/books>.

Computer labs

Computer Lab 1

Introduction to R and RStudio ¹

The goal of this lab is to introduce you to R and RStudio, which you'll be using throughout the course, both to learn the statistical concepts discussed in the lectures and also to analyze real data and come to informed conclusions. To straighten out which is which:

- R is the name of the programming language itself;
- RStudio is a convenient interface.

The R language is the standard statistical tool used by most statisticians at universities. Feeling comfortable using it is not only important for this module and any further statistics modules you may take at the Department of Mathematics of the University of York, it can also be an important factor for your future career (see the article "[R skills attract the highest salaries](#)"). Even though R is specially designed for statistics, it is in the [list of the top ten most important programming languages](#) compiled by the IEEE spectrum magazine.

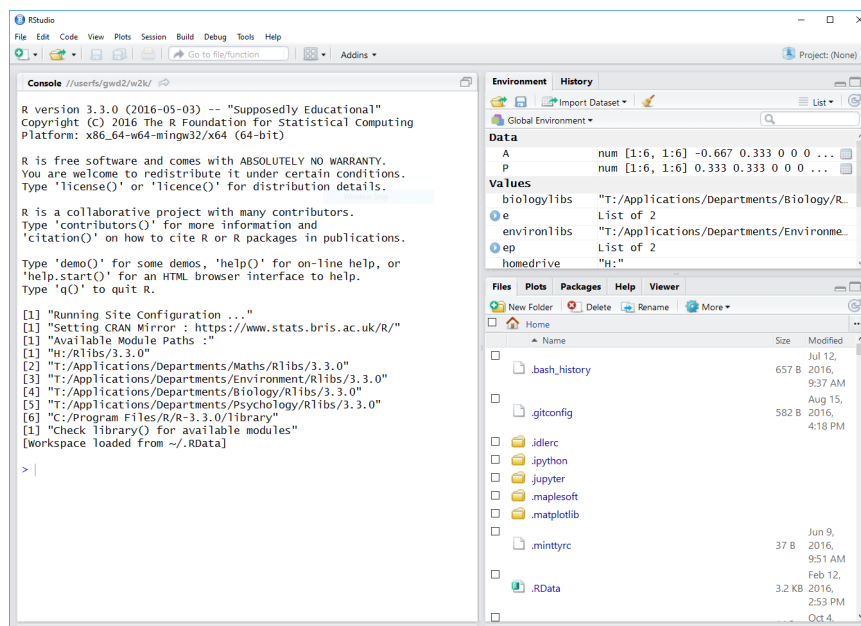
As the labs progress, you are encouraged to explore beyond what the labs dictate; a willingness to experiment will make you a much better programmer. Before we get to that stage, however, you need to build some basic fluency in R. Today we begin with the fundamental building blocks of R and RStudio: the interface, reading in data, and basic commands.

The first step is to open RStudio. If you are on a campus PC, RStudio is already installed and you can open it from the Windows Start menu. Just start typing 'RStudio' into the search box on the start menu and then click on RStudio when it shows up. This will open up a window similar to that depicted below.

If you get a popup asking you whether you want to upgrade to a newer version of RStudio, simply click the "Ignore update" button.

The way I would like you to work with this lab is to have this pdf file open on the left half of your screen and RStudio on the right half. On a Windows PC you can move a window to the left or right half of the screen by holding down the Windows key and pressing the left or right arrow key.

¹This tutorial is adapted from OpenIntro and is released under a Creative Commons Attribution-ShareAlike 3.0 Unported (<http://creativecommons.org/licenses/by-sa/3.0/>). This lab was adapted for OpenIntro by Andrew Bray and Mine Çetinkaya-Rundel from a lab written by Mark Hansen of UCLA Statistics; it was extended for the University of York by Gustav Delius, and then by Stephen Connor.



The panel in the upper right of the RStudio window contains your *Environment* as well as a *History* of the commands that you've previously entered. The lower right panel has several tabs, including *Plots* where any plots that you generate will show up.

The panel on the left is where the action happens. It's called the *Console*. Every time you launch RStudio, it will have text at the top of the console giving lots of information that you can mostly ignore, including the version of R that you're running. Below that information is the *prompt*. As its name suggests, this prompt is really a request, a request for a command. Initially, interacting with R is all about typing commands and interpreting the output. These commands and their syntax have evolved over decades (literally) and now provide what many users feel is a fairly natural way to access data and organize, describe, and invoke statistical computations.

To get you started, enter the following command at the R prompt (i.e. right after `>` on the console). You can either type it in manually or copy and paste it from this document. (Just clicking on the link in this document is going to open the file in your browser, which is not what you want. Instead you want to load it into R, therefore you need to enter the command on the R console, not just click on it here.)

```
source("http://www.openintro.org/stat/data/arbutnot.R")
```

This command instructs R to access the OpenIntro website and fetch some data: the Arbuthnot baptism counts for boys and girls. You should see that the environment area in the upper right hand corner of the RStudio window now lists a data set called `arbutnot` that has 82 observations on 3 variables.

As you interact with R, you will create a series of objects. Sometimes you load them as we have

done here, and sometimes you create them yourself as the byproduct of a computation or some analysis you have performed.

Note that because it is accessing data on the web, the above command (and the entire assignment) will work in a computer lab, in the library, or in your dorm room; anywhere you have access to the Internet.

The Data: Dr. Arbuthnot's Baptism Records

The Arbuthnot data set was compiled by Dr. John Arbuthnot, an 18th century physician, writer, and mathematician. He was interested in the ratio of newborn boys to newborn girls, so he gathered the baptism records for children born in London for every year from 1629 to 1710. We can take a look at the data by typing its name into the console and hitting Enter.

```
arbuthnot
```

Or we can view just the first six lines by typing

```
head(arbuthnot)
#>   year boys girls
#> 1 1629 5218 4683
#> 2 1630 4858 4457
#> 3 1631 4422 4102
#> 4 1632 4994 4590
#> 5 1633 5158 4839
#> 6 1634 5035 4820
```

What you should see are four columns of numbers, each row representing a different year: the first entry in each row is simply the row number (an index we can use to access the data from individual years if we want), the second is the year, and the third and fourth are the numbers of boys and girls baptized that year, respectively. Use the scroll bar on the right side of the console window to examine the complete data set.

Note that the row numbers in the first column are not part of Arbuthnot's data. R adds them as part of its printout to help you make visual comparisons. You can think of them as the index that you see on the left side of a spreadsheet. In fact, the comparison to a spreadsheet will generally be helpful. R has stored Arbuthnot's data in a kind of spreadsheet or table called a *data frame*.

You can see the dimensions of this data frame by typing:

```
dim(arbuthnot)
#> [1] 82  3
```

This indicates that there are 82 rows and 3 columns (we'll get to what the [1] means in a bit), just as it says next to the object in your Environment tab. You can see the names of these columns (or variables) by typing:

```
names(arbuthnot)
#> [1] "year" "boys" "girls"
```

You should see that the data frame contains the columns year, boys, and girls. By this point, you might have noticed that many of the commands in R look a lot like functions; that is, invoking R commands means supplying a function with some number of arguments. The dim and names commands, for example, each took a single argument, the name of a data frame.

One advantage of RStudio is that it comes with a built-in data viewer. Click on the name arbuthnot in the upper right window that lists the objects in your environment. This will bring up an alternative display of the Arbuthnot counts in the upper left panel of the RStudio window.

Some Exploration

Let's start to examine the data a little more closely. We can access the data in a single column of a data frame separately using a command like

```
arbuthnot$boys
```

This command will only show the number of boys baptized each year.

Exercise

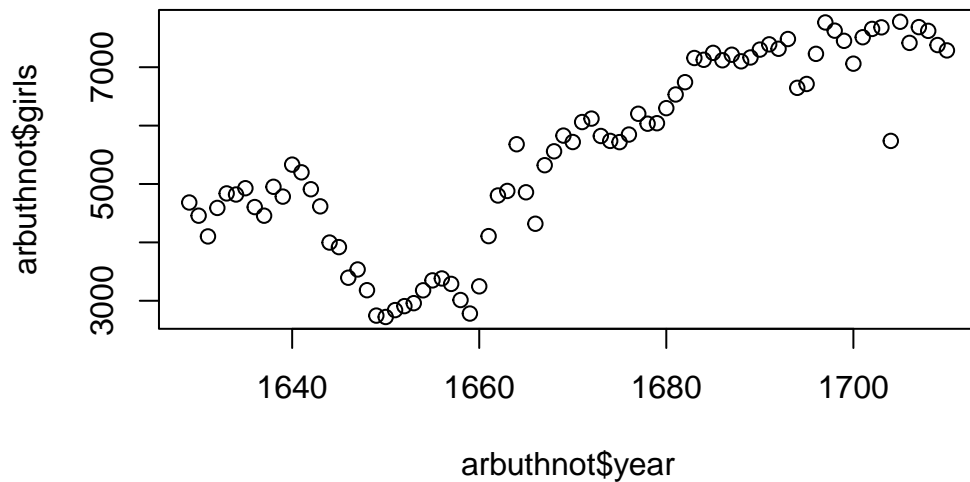
What command would you use to extract just the counts of girls baptized each year? Try it! Also, to register that you are still paying attention, please enter your command into the first question on the quiz accompanying this lab.

Notice that the way R has printed these data is different. When we looked at the complete data frame, we saw 82 rows, one on each line of the display. These data are no longer structured in a table with other variables, so they are displayed one right after another.

Objects that print out in this way are called *vectors*; they represent a set of numbers. R has added numbers in [brackets] along the left side of the printout to indicate locations within the vector. For example, 5218 follows [1], indicating that 5218 is the first entry in the vector. And if [43] starts a line, then that would mean the first number on that line would represent the 43rd entry in the vector.

R has some powerful functions for making graphics. We can create a simple plot of the number of girls baptized per year with the command

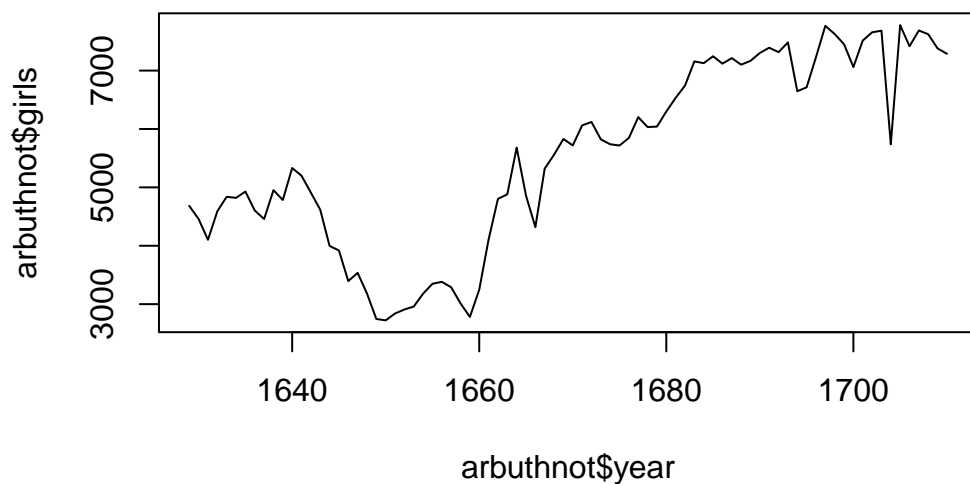

```
plot(x = arbuthnot$year, y = arbuthnot$girls)
```



By default, R creates a scatterplot with each x,y pair indicated by an open circle. The plot itself should appear under the *Plots* tab of the lower right panel of RStudio.

Notice that the command above again looks like a function, this time with two arguments separated by a comma. The first argument in the plot function specifies the variable for the x-axis and the second for the y-axis. If we wanted to connect the data points with lines, we could add a third argument, the letter *l* for line.

```
plot(x = arbuthnot$year, y = arbuthnot$girls, type = "l")
```



You might wonder how you are supposed to know that it was possible to add that third argu-

ment. Thankfully, R documents all of its functions extensively. To read what a function does and learn the arguments that are available to you, just type in a question mark followed by the name of the function that you're interested in. Try the following.

```
?plot
```

Can you figure out how to produce a plot that shows both the points and the lines connecting them?

Notice that the help file replaces the plot in the lower right panel. You can toggle between plots and help files using the tabs at the top of that panel.

Exercise

Is there an apparent trend in the number of girls baptized over the years? Answer quiz question 2. Can you also guess, just by looking at the graph, when the English civil war started?

Now, suppose we want to plot the total number of baptisms. To compute this, we could use the fact that R is really just a big calculator. We can type in mathematical expressions like

```
5218 + 4683
```

to see the total number of baptisms in 1629. We could repeat this once for each year, but there is a faster way. If we add the vector for baptisms for boys and girls, R will compute all sums simultaneously.

```
arbuthnot$boys + arbuthnot$girls
```

What you will see are 82 numbers (in that packed display, because we aren't looking at a data frame here), each one representing the sum we're after. Take a look at a few of them and verify that they are right.

We can now make a plot of the total number of baptisms per year with the command

```
plot(arbuthnot$year, arbuthnot$boys + arbuthnot$girls, type = "l")
```

This time, note that we left out the names of the first two arguments. We can do this because the help file shows that the default for `plot` is for the first argument to be the x-variable and the second argument to be the y-variable.

Next we calculate the proportion of the baptized children that are boys. We can do this for the year 1629 with the command

```
5218 / (5218 + 4683)
```

but this may also be computed for all years simultaneously:

```
arbutnnot$boys / (arbutnnot$boys + arbutnnot$girls)
```

Note that with R, as with your calculator, you need to be conscious of the order of operations. Here, we want to divide the number of boys by the total number of newborns, so we have to use parentheses. Without them, R will first do the division, then the addition, giving you something that is not a proportion.

Exercise

Now, make a plot of the proportion of boys over time. The command for making the plot will be similar to the plot command you used earlier, just with a different expression for the y argument. Tip: If you use the up and down arrow keys, you can scroll through your previous commands, your so-called command *history*. You can also access it by clicking on the history tab in the upper right panel. This will save you a lot of typing in the future. Now answer quiz question 3.

In addition to simple mathematical operators like subtraction and division, you can ask R to make comparisons like greater than, >, less than, <, and equality, == (Note that it has to be a double equal sign, not a single equal sign). For example, we can ask if boys outnumber girls in each year with the expression

```
arbutnnot$boys > arbutnnot$girls
```

This command returns 82 values of either TRUE if that year had more boys than girls, or FALSE if that year did not (the answer may surprise you). This output shows a different kind of data than we have considered so far. In the arbutnnot data frame our values are numerical (the year, the number of boys and girls). Here, we've asked R to create *logical* data, data where the values are either TRUE or FALSE. In general, data analysis will involve many different kinds of data types, and one reason for using R is that it is able to represent and compute with many of them.

You can count the number of entries for which the condition is TRUE by just summing the entries in the vector

```
sum(arbutnnot$boys > arbutnnot$girls)
```

The reason this works is that R automatically converts TRUE to 1 and FALSE to 0 when asked to do a numerical calculation with these values.

Exercise

Above you have seen how to calculate the proportion of newborns that are boys. You have also learned how to count the number of entries in the data that satisfy a particular condition. Now combine those two to answer quiz question 4.

A newer data set

In the previous few pages, you recreated some of the displays and preliminary analysis of Arbuthnot's baptism data. To practice your new skills, you will now repeat these steps, but for present day birth records in the United States. Load up the present day data with the following command.

```
source("http://www.openintro.org/stat/data/present.R")
```

The data are stored in a data frame called `present`.

Exercise

1. What years are included in this data set? What are the dimensions of the data frame and what are the variable or column names?
2. How do these counts compare to Arbuthnot's? Are they on a similar scale?
3. Does Arbuthnot's observation about boys being born in greater proportion than girls hold up in the U.S.?
4. Make a plot that displays the boy-to-girl ratio for every year in the data set. What do you see?
5. What was the largest total number of births in a single year in the U.S. during the period covered by the dataset? You can refer to the help files or the [R reference card](#) to find helpful commands.

Exercise

Now answer questions 5 and 6 in the quiz.

These data come from a report by the [Centers for Disease Control](#). Check it out if you would like to read more about an analysis of sex ratios at birth in the United States.

The above has given you a glimpse of how to work with data in R. Next we will get a glimpse of how R can simulate probability experiments.

Simulating a probability experiment

We will use R to simulate the probability experiment from Example 2.17 from the lecture. The experiment consisted of drawing one ball at random from a bag containing 4 red, 6 green and 3 blue balls.

We can use the `rep` function to create a vector with repeated entries. For example `rep("red", 4)`.

Note that in this section the worksheet also shows the output of the commands, in the lines starting with `#>`, so that you can compare with what you get when you run the commands yourself.

We can use the `c` function to *concatenate* several vectors. We can combine these two commands to create our bag; we will store this in a variable, that we choose to call `bag`, so that we can use it in what follows.

```
bag <- c(rep("red", 4), rep("green", 6), rep("blue", 3))
```

Note how this variable also shows up in the Environment tab in RStudio.

Now we can use the `sample` function to draw a single ball from the bag: `sample(bag, size = 1)`. The `size` argument specifies that we want to draw a single ball. Now when you run this command, you may get a different answer, just as when you draw an actual ball from an actual bag.

If you run the above argument repeatedly, you will get lots of different outcomes. We can easily repeat the experiment 100 times.

```
sample(bag, size = 100, replace = TRUE)
#> [1] "red" "green" "blue" "green" "green" "green" "blue" "green" "red"
#> [10] "blue" "green" "green" "green" "green" "green" "green" "green" "red"
#> [19] "blue" "green" "green" "red" "red" "blue" "green" "green" "red"
#> [28] "blue" "green" "red" "green" "green" "blue" "green" "red" "blue"
#> [37] "green" "blue" "green" "green" "blue" "green" "blue" "green" "red"
#> [46] "red" "green" "blue" "green" "blue" "green" "red" "green" "red"
#> [55] "blue" "blue" "red" "red" "green" "red" "blue" "green" "green"
#> [64] "green" "green" "red" "green" "green" "red" "green" "blue" "green"
#> [73] "red" "green" "green" "blue" "blue" "green" "red" "red" "blue"
#> [82] "blue" "green" "red" "green" "green" "red" "green" "green" "green"
#> [91] "red" "red" "green" "red" "red" "green" "green" "green" "blue"
#> [100] "green"
```

This was a lot faster than by hand taking a ball out of a bag, recording its colour, putting it back and then repeating that 100 times. In the above command the `replace = TRUE` argument meant that the ball was replaced before the next draw was made.

Now let us repeat the experiment 10000 times and calculate the frequency with which we get a red ball.

```
x <- sample(bag, size = 10000, replace = TRUE)
sum(x == "red") / 10000
#> [1] 0.3105
```

This gives us an estimate of the probability of drawing a red ball, for reasons that we will discuss theoretically later in the module. We can obviously do the same for the other colours:

```
sum(x == "green") / length(x)
#> [1] 0.4595
sum(x == "blue") / length(x)
#> [1] 0.23
```

Here we've also used the `length` command to calculate the number of entries in `x`; we know this is 10,000, of course, but using this command helps to prevent the possibility of making a typo and dividing by 1,000 or 100,000, etc.

Exercise

Now add five additional yellow balls to the bag you used so far. Then record the outcome of 100,000 repetitions of the experiment of drawing a ball from that bag. Calculate the proportion of those 100,000 draws that gave a yellow ball. Then answer the last quiz question.

‘ That was a short introduction to R and RStudio, but we will provide you with more functions and a more complete sense of the language as the course progresses. Feel free to browse around the websites for [R](#) and [RStudio](#) if you're interested in learning more.

To exit RStudio you can click the `x` in the upper right corner of the whole window. You will be prompted to save your workspace. If you click save, RStudio will save the history of your commands and all the objects in your workspace so that the next time you launch RStudio, you will see `arbuthnot` and you will have access to the commands you typed in your previous session.

On your own computer

If you would like to work on your own computer, you can download and install R from [here](#) and then download and install RStudio from [here](#). Both are free and open-source and available for Windows, Mac and Linux.

It is also possible to use R from within a browser window in case you find yourself on a computer that does not have R installed. The best option for that is probably [CoCalc](#), which offers

an online notebook interface not only to R but also to Sage (an open source alternative to Maple) and many other tools useful for mathematicians.