

# Basic Statistics

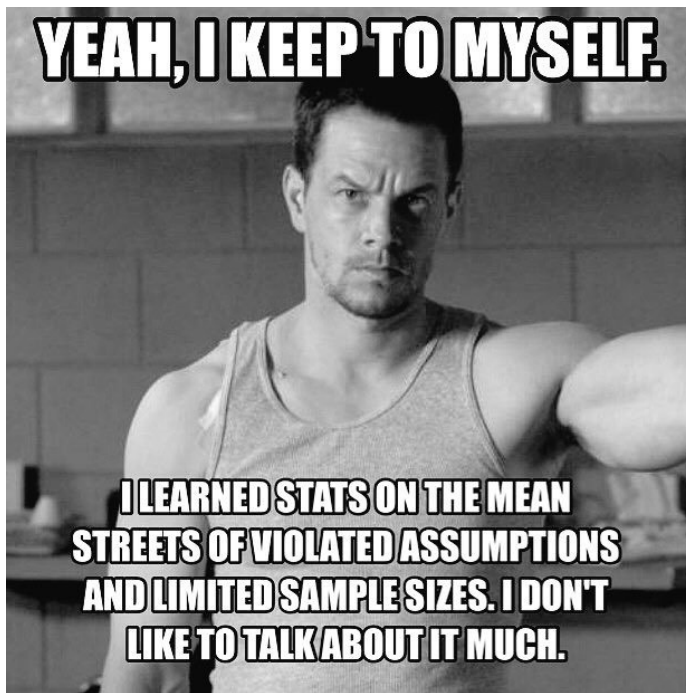
A primer in basic statistics for BCB (Hons) 2018

*AJ Smit and Robert Schlegel*

*2018-04-26*

# Contents

## Preface



This is a workshop about the practice of the basic statistics used by biologists, and not about the theory and mathematical underpinnings of the methods used. Each of the Chapters will cover a basic kind of statistical approach, and the main classes of data it applies to. Since much insight and understanding can be gained from visualising our data, we will also explore the main types of graphical summaries that best accompany the statistical methodologies. It is our intention to demonstrate how we go about analysing our data.

# Prerequisites

A prerequisite for this course is a basic proficiency in using R (?). The necessary experience will have been gained from completing the Intro R Workshop: Data Manipulation, Analysis, and Graphing<sup>1</sup> Workshop that was part of your BCB Core Honours module (i.e. Biostatistics). You will also need a laptop with R and RStudio installed as per the instructions provided in that workshop. If you do not have a personal laptop, most computers in the 5th floor lab will be correctly set up for this purpose.

---

<sup>1</sup>[https://robwschlegel.github.io/Intro\\_R\\_Workshop/](https://robwschlegel.github.io/Intro_R_Workshop/)

# 1

## Introduction

*“A scientist worthy of a lab coat should be able to make original discoveries while wearing a clown suit, or give a lecture in a high squeaky voice from inhaling helium. It is written nowhere in the math of probability theory that one may have no fun.”*

— Eliezer Yudkowsky

*“Prediction is very difficult, especially about the future.”*

— Niels Bohr

### 1.1 Venue, date and time

Basic Statistics is the second half of the BSc (Hons) Biostats core module, and will run from 12 April to 26 April 2018. This workshop will take place on Tuesdays from **13:00–17:00**, Thursdays from **10:40–17:00**, and Fridays from **08:30–17:00**. There will be an assignment due about six weeks after the end of this module, and it will provide the other half of the marks for the Biostats module. More on the assignment later.

### 1.2 Course outline

1. Introduction (this chapter)
2. Types of data
3. Descriptive statistics: Measures of location and dispersion
4. Representing data graphically
5. Distributions
6. One-sample and two-sample tests
7. Multi-sample ( $>2$ ) tests
8. Linear regression
9. Correlation
10. Confidence intervals

11. Transforming data
12. Generalised linear model (GLM)
13. Chi square tests

The course content can broadly be classified into two parts: *Descriptive Statistics* and *Inferential Statistics*.

Descriptive statistics and their associated statistical and graphical data summaries will be covered in Chapters 3 and 4. In Chapter 5 we will introduce the concepts of data distributions, knowledge of which is required to select the most appropriate inferential statistical methods.

Chapters 6-15 are about inferential statistics. Inferential tests allow us to evaluate hypotheses within a framework of probabilistic theory, which helps us infer the nature of a “population” based on a smaller representative set of samples. In particular, we can infer whether the property under scrutiny (arrived at by means of a designed experiment or a directed sampling programme) occurred as a result of deterministic influences, or whether it is as a result of chance.

### 1.3 About this Workshop

The aim of this workshop is to guide you through the outline given above. The workshop focuses broadly (and unequally) on three groups of concepts:

- Data and distributions
- Descriptive statistics and graphics
- Inferential statistics

Data and distributions are unsurprisingly about the data itself. Here we will talk about the various kinds of data that we will encounter as biologists. In the second part we will describe the data using a combination of numerical and graphical summaries. Third, all of this culminates in trying to infer from a small subset (a sample) of subjects if the characteristics under scrutiny also hold true for the entire population. We may also ask questions about probabilities, i.e. measuring the likelihood that an event will occur, or that an experiment has an outcome that is different from a situation where the influential factor(s) has no effect, or that some observation or outcome is non-random.

### 1.4 This is biology: why more R coding?

Please refer to the Intro R Workshop: Data Manipulation, Analysis and Graphing<sup>1</sup> for why we feel strongly that you use R (?) for the analyses that we will perform here. All of the reasons provided there are valid here too, but one reason perhaps more so than others—R and RStudio promote the principles of *reproducible research*, and in fact make it very easy to implement. We will focus on some of these principles throughout the workshop, and the assignments will require that you submit a fully functional working script, complete with all the notes, memos, examples, data, executable code, and output that will result from completing the course material.

What other options are there for analysing the kinds of data that we will encounter in biological research? Software packages like the ones you may be familiar with, such as Statistica and SPSS, are often used to perform many of these analyses. They are rather limited with regards to the full scope of modern statistical methods in use by biologists today, but many people still use these kinds of software as they provide the basic kinds analyses that still form the staple of the

---

<sup>1</sup>[https://robwschlegel.github.io/Intro\\_R\\_Workshop/](https://robwschlegel.github.io/Intro_R_Workshop/)

biological and medical sciences. For the many reasons provided above, we prefer to use R as the *engine* within which to do our biological data analysis. R is used by academic statisticians the world over, and it is therefore an excellent choice for our purpose here.

## 1.5 Installing R and RStudio

We assume that you already have R installed on your computer, as all of you will have already completed the the Intro R Workshop. If you need a refresher, please refer to Intro R Workshop: Data Manipulation, Analysis and Graphing<sup>2</sup> for the installation instructions.

## 1.6 Resources

- New users should introduce themselves to the R ecosystem<sup>3</sup>
- A fancy interactive website that covers a wide range of basic statistics<sup>4</sup>
- An easy to follow walkthrough for a statistical analysis<sup>5</sup>
- Learn more about tidy statistical inference<sup>6</sup>
- A thorough journey through the philosophy of data visualisation<sup>7</sup>
- Google<sup>8</sup>
- Stack Overflow<sup>9</sup>

## 1.7 Style and code conventions

Early on, develop the habit of unambiguous and consistent style and formatting when writing your code, or anything else for that matter. Pay attention to detail and be pedantic. This will benefit your scientific writing in general. Although many R commands rely on precisely formatted statements (code blocks), style can nevertheless to *some extent* have a personal flavour to it. The key is *consistency*. In this book we use certain conventions to improve readability. We also use a consistent set of conventions to refer to code, and in particular to typed commands and package names.

- Package names are shown in a bold font over a grey box, e.g. **tidyr**.
- Functions are shown in normal font followed by parentheses and also over a grey box, e.g. `plot()`, or `summary()`.
- Other R objects, such as data, function arguments or variable names are again in normal font over a grey box, but without parentheses, e.g. `x` and `apples`.
- Sometimes we might directly specify the package that contains the function by using two colons, e.g. `dplyr::filter()`.
- Commands entered onto the R command line (console) and the output that is returned will be shown in a code block, which is a light grey background with code font. The commands entered start at the beginning of a line and the output it produces is preceded by `R>`, like so:

<sup>2</sup>[https://robwschlegel.github.io/Intro\\_R\\_Workshop/](https://robwschlegel.github.io/Intro_R_Workshop/)

<sup>3</sup>[fg2re.sellorm.com](http://fg2re.sellorm.com)

<sup>4</sup><http://students.brown.edu/seeing-theory/>

<sup>5</sup><http://www.sthda.com/english/wiki/unpaired-two-samples-t-test-in-r>

<sup>6</sup><https://infer.netlify.com>

<sup>7</sup><http://www.serialmentor.com/dataviz/>

<sup>8</sup>[www.google.com](http://www.google.com)

<sup>9</sup>[www.stackoverflow.com](http://www.stackoverflow.com)

```
set.seed(666)
rnorm(n = 10, mean = 0, sd = 13)
```

```
R> [1]  9.7930436  26.1866107 -4.6167480  26.3661820 -28.8193679
R> [6]  9.8591503 -16.9804084 -10.4327544 -23.2991308 -0.5464219
```

Consult these resources for more about R code style :

- Google's R style guide<sup>10</sup>
- The tidyverse style guide<sup>11</sup>
- Hadley Wickham's advanced R style guide<sup>12</sup>

We may also insert maths expressions within the text, like this  $f(k) = \binom{n}{k} p^k (1-p)^{n-k}$  or on their own, like this:

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

## 1.8 Assessment and teaching philosophy

Grades will be based on the aggregate performance across two group projects; the first group project was completed after the Intro R Workshop. The project for this workshop will represent 35% of the total grade for BioStatistics. The remaining 15% will come from daily participation. This will be assessed by the R scripts produced in class and follows these five criteria:

- The script has been uploaded to GitHub
- The script covers the content of the day
- The code runs without errors
- Proper style conventions have been observed
- Liberally commented

### BONUS POINTS

- Additional analysis not performed in class
- Additional figure not created in class

In cases where students are borderline between lower and higher grades, a high level of participation in the class discussions and class in general will win the day for the higher grade.

The daily scripts are essential to understanding the material. Although they comprise only 15% of the final grade, performance on the projects is usually correlated with effort on the daily assignments.

Whereas plagiarism will not be tolerated, students ARE encouraged to work together to learn from one another and solve problems in a collaborative and collegial way.

## 1.9 About this document

This document, which is available as an HTML file that's viewable on a web browser of your choice (anything will do, but we discourage using Internet Explorer) and as a PDF (accessible from the link at the top of any of the website's pages) that may be printed, was prepared by the software tools available to R via RStudio. We use the package called `bookdown` that may

<sup>10</sup><https://google.github.io/styleguide/Rguide.xml>

<sup>11</sup><http://style.tidyverse.org>

<sup>12</sup><http://adv-r.had.co.nz/Style.html>



be accessed and read about here<sup>13</sup> to produce this documentation. The entire source code to reproduce this book is available from my GitHub repo<sup>14</sup>.

You will notice that this repository uses GitHub<sup>15</sup>, and you are advised to set up your own repository for R scripts and all your data. We will touch on GitHub and the principles of reproducible research later, and GitHub forms a core ingredient of such a workflow.

---

<sup>13</sup><https://bookdown.org/yihui/bookdown/>

<sup>14</sup>[https://github.com/ajsmitt/Basic\\_stats](https://github.com/ajsmitt/Basic_stats)

<sup>15</sup><https://github.com>

# 2

## Types of data

*“The plural of anecdote is not data.”*

— Roger Brinner

In this chapter we will, firstly, look at the different kinds of biological and environmental data that are typically encountered by most biologists. The data seen here are not an exhaustive list of all the various types of data out there, but it should represent the bulk of our needs.

After we have become familiar with the different kinds of data, we will look at summaries of these data, which is generally required as the starting point for our analyses. After summarising the data in tables and so forth, we may want to produce graphical summaries to see broad patterns and trends; visual data representations, which complement the tabulated data, will be covered in a later chapter (Chapter 4). Both of these approaches form the basis of “exploratory data analysis.”

### 2.1 Data classes

In biology we will encounter many kinds of data, and depending on which kind, the type of statistical analysis will be decided.

#### 2.1.1 Numerical data

Numerical data are quantitative in nature. They represent things that can be objectively counted, measured or calculated.

##### 2.1.1.1 *Nominal (discrete) data*

Integer data (discrete numbers or whole numbers), such as counts. For example, family A may have 3 children and family B may have 1 child, neither may have 2.3 children. Integer data usually answer the question, “how many?” In R integer data are called `int` or `<int>`.

### 2.1.1.2 Continuous data

These usually represent measured “things,” such as something’s heat content (temperature, measured in degrees Celsius) or distance (measured in metres or similar), etc. They can be rational numbers including integers and fractions, but typically they have an infinite number of “steps” that depends on rounding (they can even be rounded to whole integers) or considerations such as measurement precision and accuracy. Often, continuous data have upper and lower bounds that depend on the characteristics of the phenomenon being studied or the measurement being taken. In R, continuous data are denoted `num` or `<dbl>`.

The kinds of summaries that lend themselves to continuous data are:

- Frequency distributions
- Relative frequency distributions
- Cumulative frequency distributions
- Bar graphs
- Box plots
- Scatter plots

### 2.1.1.3 Dates

Dates are a special class of continuous data, and there are many different representations of the date classes. This is a complex group of data, and we will not cover much of it in this course.

## 2.1.2 Qualitative data

Qualitative data may be well-defined categories or they may be subjective, and generally include descriptive words for classes (e.g. mineral, animal, plant) or rankings (e.g. good, better, best).

### 2.1.2.1 Categorical data

Because there are categories, the number of members belonging to each of the categories can be counted. For example, there are three red flowers, 66 purple flowers, and 13 yellow flowers. The categories cannot be ranked relative to each other; in the example just provided, for instance, no value judgement can be assigned to the different colours. It is not better to be red than it is to be purple. There are just fewer red flowers than purple ones. Contrast this to another kind of categorical data called “ordinal data” (see next). This class of data in an R dataframe (or in a “tibble”) is indicated by `Factor` or `<fctr>`.

The kinds of summaries that lend themselves to categorical data are:

- Frequency distributions
- Relative frequency distributions
- Bar graphs
- Pie graphs (!!!)
- Category statistics

### 2.1.2.2 Ordinal data

This is a type of categorical data where the classes are ordered (a synonym is “ranked”), typically from low to high (or *vice versa*), but where the magnitude between the ordered classes cannot be precisely measured or quantified. In other words, the difference between them is

somewhat subjective (i.e. it is qualitative rather than quantitative). These data are on an ordinal scale. The data may be entered as descriptive character strings (i.e. as words), or they may have been translated to an ordered vector of integers; for example, “1” for terrible, “2” for so-so, “3” for average, “4” for good and “5” for brilliant. Irrespective of how the data are present in the dataframe, computationally (for some calculations) they are treated as an ordered sequence of integers, but they are simultaneously treated as categories (say, where the number of responses that report “so-so” can be counted). Ordinal data usually answer questions such as, “how many categories can the phenomenon be divided into, and how does each category rank with respect to the others?” Columns containing this kind of data are named `ord.factor` or `<ord>`.

### 2.1.3 Binary data

Right or wrong? True or false? Accept or reject? Black or white? Positive or negative? Good or bad? You get the idea... In other words, these are observations or responses that can take only one of two mutually exclusive outcomes. In R these are treated as “Logical” data that take the values of `TRUE` or `FALSE` (note the case). In R, and computing generally, logical data are often denoted with 1 for `TRUE` and 0 for `FALSE`. This class of data is indicated by `logi` or `<lgl>`.

### 2.1.4 Character values

As the name implies, these are not numbers. Rather, they are human words that have found their way into R for one reason or another. In biology we most commonly encounter character values when we have a list of things, such as sites or species. These values will often be used as categorical or ordinal data.

### 2.1.5 Missing values

Unfortunately, one of the most reliable aspects of any biological dataset is that it will contain some missing data. But how can something contain missing data? One could be forgiven for assuming that if the data are missing, then they obviously aren’t contained in the dataset. To better understand this concept we must think back to the principles of tidy data. Every observation must be in a row, and every column in that row must contain a value. The combination of multiple observations then makes up our matrix of data. Because data are therefore presented in a two-dimensional format, any missing values from an observation will need to have an empty place-holder to ensure the consistency of the matrix. These are what we are referring to when we speak of “missing values”. In R these appear as a `NA` in a dataframe and are slightly lighter than the other values. These data are indicated in the Environment as `NA` and if a column contains only missing values it will be denoted as `<NA>`.

### 2.1.6 Complex numbers

*“And if you gaze long enough into an abyss, the abyss will gaze back into you.”*

— Friedrich Nietzsche

In an attempt to allow the shreds of our sanity to remain stitched together we will end here with data types. But be warned, ye who enter, that below countless rocks, and around a legion of corners, lay in wait a myriad of complex data types. We will encounter many of these at the end of this course when we encounter modeling, but by then we will have learned a few techniques that will prepare us for the encounter.

## 2.2 Viewing our data

There are many ways of finding broad views of our data in R. The first few functions that we will look at were designed to simply scrutinise the contents of the tibbles, which is the “tidyverse” name for the general “container” that holds our data in the software’s environment (i.e. in a block of the computer’s memory dedicated to the R software). Whatever data are in R’s environment will be seen in the “Environment” tab in the top right of RStudio’s four panes.

### 2.2.1 From the Environment pane

The first way to see what’s in the tibble is not really a function at all, but a convenient (and lazy) way of quickly seeing a few basic things about our data. Let us look at the `ChickWeight` data. Load it like so (you’ll remember from the Intro R Workshop):

```
# loads the tidyverse functions; it contains the 'as_tibble()' function
library(tidyverse)
# the 'ChickWeight' data are built into R;
# here we assign it as a tibble to an object named 'chicks'
chicks <- as_tibble(ChickWeight)
```

In the Environment pane, the object named `chicks` will now appear under the panel named Data. To the left of it is a small white arrow in a blue circular background. By default the arrow points to the right. Clicking on it causes it to point down, which denotes that the data contained within the tibble have become expanded. The names of the columns (more correctly called “variables”) can now be seen. There you can see the variables `weight`, `Time`, `Chick` and `Diet`. The class of data they represent can be seen too: there’s continuous data of class `num`, a variable of `ord.factor`, and a categorical variable of class `Factor`. Beneath these there’s a lot of attributes that denote some meta-data, which you may safely ignore for now.

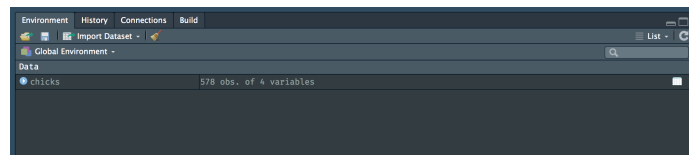


Figure 2.1: What is in the Chicks data?

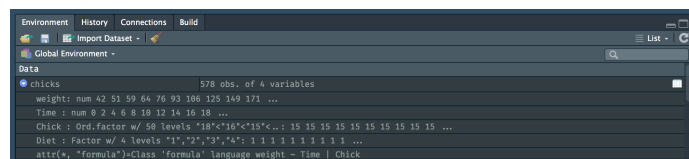


Figure 2.2: This is what is in the Chicks data.

### 2.2.2 `head()` and `tail()`

The `head()` and `tail()` functions simply display top and bottom portions of the tibble, and you may add the `n` argument and an integer to request that only a certain number of rows are returned; by default the top or bottom six rows are displayed.

There are various bits of additional information printed out. The display will change somewhat if there are many more variables than that which can comfortably fit within the width of the

output window (typically the Console). The same kinds of information as was returned with the Environment pane expansion arrow are displayed, but the data class is now accompanied by an angle bracket (i.e. `<...>`) notation. For example, `num` in the Environment pane and `<dbl>` as per the `head()` or `tail()` methods are exactly the same: both denote continuous (or “double precision”) data.

```
head(chicks)
```

```
R> # A tibble: 6 x 4
R>   weight Time Chick Diet
R>   <dbl> <dbl> <ord> <fct>
R> 1    42.    0. 1     1
R> 2    51.    2. 1     1
R> 3    59.    4. 1     1
R> 4    64.    6. 1     1
R> 5    76.    8. 1     1
R> 6    93.   10. 1     1
```

```
tail(chicks, n = 2)
```

```
R> # A tibble: 2 x 4
R>   weight Time Chick Diet
R>   <dbl> <dbl> <ord> <fct>
R> 1   264.   20. 50     4
R> 2   264.   21. 50     4
```

As an alternative to `head()`, you may also simply type the name of the object (here `chicks`) in the Console (or write it in the Source Editor if it is necessary to retain the function for future use) and the top portion of the tibble will be displayed, again trimmed to account for the width of the display.

### 2.2.3 `colnames()`

This function simply returns a listing of the variable (column) names.

```
colnames(chicks)
```

```
R> [1] "weight" "Time" "Chick" "Diet"
```

There is an equivalent function called `rownames()` that may be used to show the names of rows in your tibble, if these are present. Row names are generally discouraged, and we will refrain from using them here.

### 2.2.4 `summary()`

The next way to see the contents of the tibble is to apply the `summary()` function. Here we see something else. Some descriptive statistics that describe properties of the full set of data are now visible. These summary statistics condense each of the variables into numbers that describe some properties of the data within each column. You will already know the concepts of the “minimum,” “median,” “mean,” and “maximum.” These are displayed here.

```
summary(chicks)
```

```
R>   weight      Time      Chick      Diet
R> Min.   : 35.0   Min.   : 0.00   13      : 12   1:220
```

```
R> 1st Qu.: 63.0   1st Qu.: 4.00   9      : 12   2:120
R> Median :103.0   Median :10.00  20     : 12   3:120
R> Mean   :121.8   Mean    :10.72  10     : 12   4:118
R> 3rd Qu.:163.8   3rd Qu.:16.00  17     : 12
R> Max.    :373.0   Max.     :21.00  19     : 12
R>                                     (Other):506
```

This will serve well as an introduction to the next chapter, which is about descriptive statistics. What are they, and how do we calculate them?

# 3

## Descriptive statistics: central tendency and dispersion

*“I think it is much more interesting to live with uncertainty than to live with answers that might be wrong.”*

— Richard Feynman

In this Chapter we will focus on basic descriptions of the data, and these initial forays are built around measures of the central tendency of the data (the mean, median, mode) and the dispersion and variability of the data (standard deviations and their ilk). The materials covered in this and the next two chapters concern a broad discussion that will aid us in understanding our data better prior to analysing it, and before we can draw inference from it. In this work flow it emerges that descriptive statistics generally precede inferential statistics.

Let us now turn to some of the most commonly used descriptive statistics, and learn about how to calculate them.

### 3.1 Samples and populations

This is a simple toy example. In real life, however, our data will be available in a tibble (initially perhaps captured in MS Excel before importing it as a .csv file into R, where the tibble is created). To see how this can be done more realistically using actual data, let us turn to the `ChickenWeight` data, which, as before, we place in the object `chicks`. Recall the pipe operator (`%>%`, pronounced “then”) that we introduced in the Intro R Workshop—we will use that here, throughout. Let us calculate the sample size:

To determine the sample size we can use the `length()` or `n()` functions; the latter is for use within `dplyr`’s `summarise()` method, and it is applied without writing anything inside of the `()`, like this:



```
# first load the tidyverse packages that has the pipe operator, %>%
library(tidyverse)
chicks <- as_tibble(ChickWeight)

# how many weights are available across all Diets and Times?
chicks %>%
  summarise(length = n())

R> # A tibble: 1 x 1
R>   length
R>   <int>
R> 1     578

# the same as
length(chicks$weight)

R> [1] 578
```

Note that this gives us the number of all of the chicks in the experiment, irrespective of which diet they were given. It will make more sense to know how many chicks were assigned to each of the experimental diets they were raised on.

**Task:** Figure out the number of chicks in each of the feed categories.

## 3.2 Measures of central tendency

Table 3.1: Measures of central tendency.

Statistic	Function	Package
Mean	mean()	base
Median	median()	base
Skewness	skewness()	e1071
Kurtosis	kurtosis()	e1071

The measures of central tendency are also sometimes called “location” statistics. We have already seen summaries of the mean and the median when we called to `summary()` function on the chicks data in Chapter 2. Here we shall show you how they can be calculated using some built-in R functions.

### 3.2.1 The mean

The sample *mean* is the arithmetic average of the data, and it is calculated by summing all the data and dividing it by the sample size,  $n$ . The mean,  $\bar{x}$ , is calculated thus:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

where  $x_1 + x_2 + \cdots + x_n$  are the observations and  $n$  is the number of observations in the sample.

In R one can quickly apply the `mean()` function to some data. Let us create a vector of arbitrary numbers using the “combine” function, `c()`, and then apply the function for the mean:

```
# combine a series of numbers into a vector;
# hint: use this function in the exercises that we will require from you
# later on...
dat1 <- c(23, 45, 23, 66, 13)
mean(dat1)
```

```
R> [1] 34
```

Below, we use another tidyverse package, **dplyr** and its `summarise()` function, whose purpose it is to *summarise* the entire column into one summary statistic, in this case the mean:

```
chicks %>%
  summarise(mean_wt = mean(weight))
```

```
R> # A tibble: 1 x 1
R>   mean_wt
R>   <dbl>
R> 1    122.
```

We can achieve the same using the more traditional syntax, which in some instances may be slightly less verbose, but less user-friendly, especially when multiple summary statistics are required (we shall later on how we can summarise a vector of data into multiple statistics). The traditional syntax is:

```
# the '$' operator is used to denote a variable inside of the tibble
mean(chicks$weight)
```

```
R> [1] 121.8183
```

**Task:** How would you manually calculate the mean mass for the chicks? Do it now!

Notice above how the two approaches display the result differently: in the first instance, using `summarise()`, the answer is rounded to zero decimal places; in the second, it is displayed (here) at full precision. The precision of the answer that you require depends on the context of your study, so make sure that you use the appropriate number of significant digits. Using the `summarise()` approach again, here is how you can adjust the number of decimal places of the answer:

```
# the value printed in the HTML/PDF versions is incorrect;
# check in the console for correct output
chicks %>%
  summarise(mean_wt = round(mean(weight), 1))
```

```
R> # A tibble: 1 x 1
R>   mean_wt
R>   <dbl>
R> 1    122.
```

**Task:** What happens when there are missing values (NA)? Consult the help file for the `mean()` function, discuss amongst yourselves, and then provide a demonstration to the class of how you would handle missing values. Hint: use the `c()` function to capture a series of data that you can then use to demonstrate your understanding.

At this point it might be useful to point out that the mean (or any function for that matter, even one that does not yet exist) can be programatically calculated. Let us demonstrate the principle by reproducing the mean function from the constituent parts:

```
chicks %>%
  summarise(mean_wt = sum(weight) / n())
```

```
R> # A tibble: 1 x 1
R>   mean_wt
R>   <dbl>
R> 1    122.
```

The mean is quite sensitive to the presence of outliers or extreme values in the data, and it is advised that its use be reserved for normally distributed data from which the extremes/outliers have been removed. When extreme values are indeed part of our data and not simply “noise,” then we have to resort to a different measure of central tendency: the median.

\*Task:\* In statistics, what do we mean with “noise”?

### 3.2.2 The median

The *median* can be calculated by “hand” (if you have a small enough amount of data) by arranging all the numbers in sequence from low to high, and then finding the middle value. If there are five numbers, say 5, 2, 6, 13, 1, then you would arrange them from low to high, i.e. 1, 2, 5, 6, 13. The middle number is 5. This is the median. But there is no middle if we have an even number of values. What now? Take this example sequence of six integers (they may also be floating point numbers), which has already been ordered for your pleasure: 1, 2, 5, 6, 9, 13. Find the middle two numbers (i.e. 5, 6) and take the mean. It is 5.5. That is the median.

Let us find the median for the weights of the chickens in the ChickWeight data:

```
chicks %>%
  summarise(med_wt = median(weight))
```

```
R> # A tibble: 1 x 1
R>   med_wt
R>   <dbl>
R> 1    103.
```

The median is therefore the value that separates the lower half of the sample data from the upper half. In normally distributed continuous data the median is equal to the mean. Comparable concepts to the median are the *1st* and *3rd quartiles*, which, respectively, separate the first quarter of the data from the last quarter—see later. The advantage of the median over the mean is that it is unaffected (i.e. not skewed) by extreme values or outliers, and it gives an idea of the typical value of the sample. The median is also used to provide a robust description of non-parametric data (see Chapter 4 for a discussion on normal data and other data distributions).

### 3.2.3 Skewness

Skewness is a measure of symmetry, and it is best understood by understanding the location of the median relative to the mean. A negative skewness indicates that the mean of the data is less than their median—the data distribution is left-skewed. A positive skewness results from data that have a mean that is larger than their median; these data have a right-skewed distribution.

```
library(e1071)
skewness(faithful$eruptions)
```

```
R> [1] -0.4135498
```

Task: Is the distribution left or right skewed?

### 3.2.4 Kurtosis

Kurtosis describes the tail shape of the data's distribution. A normal distribution has zero kurtosis and thus the standard tail shape (mesokurtic). Negative kurtosis indicates data with a thin-tailed (platykurtic) distribution. Positive kurtosis indicates a fat-tailed distribution (leptokurtic).

```
# library(e1071)
kurtosis(faithful$eruptions)
```

```
R> [1] -1.511605
```

## 3.3 Measures of variation and spread

Since the mean or median does not tell us everything there is to know about data, we will also have to determine some statistics that inform us about the variation (or spread or dispersal or inertia) around the central/mean value.

Table 3.2: Measures of variation and spread.

Statistic	Function
Variance	<code>var()</code>
Standard deviation	<code>sd()</code>
Minimum	<code>min()</code>
Maximum	<code>max()</code>
Range	<code>range()</code>
Quantile	<code>quantile()</code>
Covariance	<code>cov()</code>
Correlation	<code>cor()</code>

### 3.3.1 The variance and standard deviation

The variance and standard deviation are examples of interval estimates. The *sample variance*,  $S^2$ , may be calculated according to the following formula:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

This reads: “the sum of the squared differences from the mean, divided by the sample size minus 1.”

To get the *standard deviation*,  $S$ , we take the square root of the variance, i.e.  $S = \sqrt{S^2}$ .

No need to plug these equations into MS Excel. Let us quickly calculate  $S$  in R. Again, we use the `chicks` data:

```
chicks %>%
  summarise(sd_wt = sd(weight))
```

```
R> # A tibble: 1 x 1
R>   sd_wt
R>   <dbl>
R> 1  71.1
```

The interpretation of the concepts of mean and median are fairly straight forward and intuitive. Not so for the measures of variance. What does  $S$  represent? Firstly, the unit of measurement of  $S$  is the same as that of  $\bar{x}$  (but the variance doesn't share this characteristic). If temperature is measured in  $^{\circ}\text{C}$ , then  $S$  also takes a unit of  $^{\circ}\text{C}$ . Since  $S$  measures the dispersion *around* the mean, we write it as  $\bar{x} \pm S$  (note that often the mean and standard deviation are written with the letters *mu* and *sigma*, respectively; i.e.  $\mu \pm \sigma$ ). The smaller  $S$  the closer the sample data are to  $\bar{x}$ , and the larger the value is the further away they will spread out from  $\bar{x}$ . So, it tells us about the proportion of observations above and below  $\bar{x}$ . But what proportion? We invoke the 68-95-99.7 rule: ~68% of the population (as represented by a random sample of  $n$  observations taken from the population) falls within  $1S$  of  $\bar{x}$  (i.e. ~34% below  $\bar{x}$  and ~34% above  $\bar{x}$ ); ~95% of the population falls within  $2S$ ; and ~99.7% falls within  $3S$ .

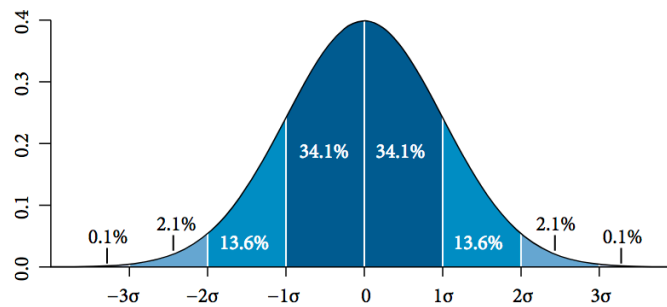


Figure 3.1: Caption here.

Like the mean,  $S$  is affected by extreme values and outliers, so before we attach  $S$  as a summary statistic to describe some data, we need to ensure that the data are in fact normally distributed. We will talk about how to do this in Chapter 6, where we will go over the numerous ways to check the assumption of normality. When the data are found to be non-normal, we need to find appropriate ways to express the spread of the data. Enter the quartiles.

### 3.3.2 Quartiles

A more forgiving approach (forgiving of the extremes, often called “robust”) is to divide the distribution of ordered data into quarters, and find the points below which 25% (0.25, the first quartile), 50% (0.50, the median) and 75% (0.75, the third quartile) of the data are distributed. These are called *quartiles* (for “quarter;” not to be confused with *quantile*, which is a more general form of the function that can be used to divide the distribution into any arbitrary proportion from 0 to 1). In R we use the `quantile()` function to provide the quartiles; we demonstrate two approaches:

```
quantile(chicks$weight)
```

```
R>      0%      25%      50%      75%     100%
R>  35.00  63.00 103.00 163.75 373.00
```

```
chicks %>%
  summarise(min_wt = min(weight),
            qrt1_wt = quantile(weight, p = 0.25),
            med_wt = median(weight),
            qrt3_wt = median(weight, p = 0.75),
            max_wt = max(weight))
```

```
R> # A tibble: 1 x 5
R>   min_wt qrt1_wt med_wt qrt3_wt max_wt
R>   <dbl>  <dbl>  <dbl>  <dbl>  <dbl>
R> 1    35.    63.   103.   163.  373.
```

```
# note median(weight) is the same as quantile(weight, p = 0.5)
# in the summarise() implementation, above
```

*Task:* What is different about the `quantile()` function that caused us to specify the calculation in the way in which we have done so above? You will have to consult the help file, read it, understand it, think about it, and experiment with the ideas. Take 15 minutes to figure it out and report back to the class.

### 3.3.3 The minimum, maximum and range

A description of the extent of the data can also be provided by the functions `min()`, `max()` and `range()`.

These statistics apply to data of any distribution, and not only to normal data. This is often the first place you want to start when looking at the data for the first time. We've seen above how to use `min()` and `max()`, so below we will quickly look at how to use `range()` in both the base R and tidy methods:

```
range(chicks$weight)
```

```
R> [1] 35 373
```

```
chicks %>%
  summarise(lower_wt = range(weight)[1],
            upper_wt = range(weight)[2])
```

```
R> # A tibble: 1 x 2
R>   lower_wt upper_wt
R>   <dbl>    <dbl>
R> 1     35.    373.
```

Note that `range()` actually gives us the minimum and maximum values, and not the difference between them. To find the range value properly we must be a bit more clever:

```
range(chicks$weight)[2] - range(chicks$weight)[1]
```

```
R> [1] 338
```

```
chicks %>%
  summarise(range_wt = range(weight)[2] - range(weight)[1])
```

```
R> # A tibble: 1 x 1
R>   range_wt
R>   <dbl>
R> 1     338.
```

### 3.3.4 Covariance

### 3.3.5 Correlation

The correlation coefficient of two matched (paired) variables is equal to their covariance divided by the product of their individual standard deviations. It is a normalised measurement of how linearly related the two variables are.

Graphical displays of correlations are provided by scatter plots as can be seen in Section X.

## 3.4 Missing values

As mentioned in Chapter 2, missing data are pervasive in the biological sciences. Happily for us, R is designed to handle these data easily. It is important to note here explicitly that all of the basic functions in R will by default *NOT* ignore missing data. This has been done so as to prevent the user from accidentally forgetting about the missing data and potentially making errors in later stages in an analysis. Therefore, we must explicitly tell R when we want it to omit missing values from a calculation. Let's create a small vector of data to demonstrate this:

```
dat1 <- c(NA, 12, 76, 34, 23)

# Without telling R to omit missing data
mean(dat1)
```

```
R> [1] NA
```

```
# Omitting the missing data
mean(dat1, na.rm = TRUE)
```

```
R> [1] 36.25
```

Note that this argument, `na.rm = TRUE` may be used in all of the functions we have seen thus far in this chapter.

## 3.5 Descriptive statistics by group

Above we have revised the basic kinds of summary statistics, and how to calculate them. This is nice. But it can be more useful. The real reason why we might want to see the descriptive statistics is to facilitate comparisons between groups. In the chicks data we calculated the mean (etc.) for all the chickens, over all the diet groups to which they had been assigned (there are four factors, i.e. Diets 1 to 4), and over the entire duration of the experiment (the experiment lasted 21 days). It would be more useful to see what the weights are of the chickens in each of the four groups at the end of the experiment — we can compare means ( $\pm$  SD) and medians ( $\pm$  interquartile ranges, etc.), for instance. You'll notice now how the measures of central tendency is being combined with the measures of variability/range. Further, we can augment this statistical summary with many kinds of graphical summaries, which will be far more revealing

of differences (if any) amongst groups. We will revise how to produce the group statistics and show a range of graphical displays.

### 3.5.1 Groupwise summary statistics

At this point you need to refer to Chapter 10<sup>1</sup> (Tidy data) and Chapter 11<sup>2</sup> (Tidier data) in the **Intro R Workshop** to remind yourself about in what format the data need to be before we can efficiently work with it. A hint: one observation in a row, and one variable per column. From this point, it is trivial to do the various data descriptions, visualisations, and analyses. Thankfully, the `chicks` data are already in this format.

So, what are the summary statistics for the chickens for each diet group at day 21?

```
grp_stat <- chicks %>%
  filter(Time == 21) %>%
  group_by(Diet, Time) %>%
  summarise(mean_wt = round(mean(weight, na.rm = TRUE), 2),
            med_wt = median(weight, na.rm = TRUE),
            sd_wt = round(sd(weight, na.rm = TRUE), 2),
            sum_wt = sum(weight),
            min_wt = min(weight),
            qrt1_wt = quantile(weight, p = 0.25),
            med_wt = median(weight),
            qrt3_wt = median(weight, p = 0.75),
            max_wt = max(weight),
            n_wt = n())

grp_stat

R> # A tibble: 4 x 11
R> # Groups:   Diet [?]
R>   Diet Time mean_wt med_wt sd_wt sum_wt min_wt qrt1_wt qrt3_wt max_wt
R>   <fct> <dbl>   <dbl>   <dbl> <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
R> 1 1      21.    178.    166.   58.7  2844.    96.    138.    166.    305.
R> 2 2      21.    215.    212.   78.1  2147.    74.    169.    212.    331.
R> 3 3      21.    270.    281.   71.6  2703.   147.    229.    281.    373.
R> 4 4      21.    239.    237.   43.4  2147.   196.    204.    237.    322.
R> # ... with 1 more variable: n_wt <int>
```

### 3.5.2 Displays of group summaries

There are several kinds of graphical displays for your data. We will show some which are able to display the spread of the raw data, the mean or median, as well as the appropriate accompanying indicators of variation around the mean or median.

```
library(ggpubr) # needed for arranging multi-panel plots
plt1 <- chicks %>%
  filter(Time == 21) %>%
  ggplot(aes(x = Diet, y = weight)) +
  geom_point(data = grp_stat, aes(x = Diet, y = mean_wt),
            col = "black", fill = "red", shape = 23, size = 3) +
```

<sup>1</sup>[https://robwschlegel.github.io/Intro\\_R\\_Workshop/tidy.html](https://robwschlegel.github.io/Intro_R_Workshop/tidy.html)

<sup>2</sup>[https://robwschlegel.github.io/Intro\\_R\\_Workshop/tidier.html](https://robwschlegel.github.io/Intro_R_Workshop/tidier.html)



```

geom_jitter(width = 0.05) + # geom_point() if jitter not required
labs(y = "Chicken mass (g)") +
theme_pubr()

plt2 <- ggplot(data = grp_stat, aes(x = Diet, y = mean_wt)) +
  geom_bar(position = position_dodge(), stat = "identity",
           col = NA, fill = "salmon") +
  geom_errorbar(aes(ymin = mean_wt - sd_wt, ymax = mean_wt + sd_wt),
               width = .2) +
  labs(y = "Chicken mass (g)") +
  theme_pubr()
# position_dodge() places bars side-by-side
# stat = "identity" prevents the default count from being plotted

# a description of the components of a boxplot is provided in the help file
# geom_boxplot()
plt3 <- chicks %>%
  filter(Time == 21) %>%
  ggplot(aes(x = Diet, y = weight)) +
  geom_boxplot(fill = "salmon") +
  geom_jitter(width = 0.05, fill = "white", col = "blue", shape = 21) +
  labs(y = "Chicken mass (g)") +
  theme_pubr()

plt4 <- chicks %>%
  filter(Time %in% c(10, 21)) %>%
  ggplot(aes(x = Diet, y = weight, fill = as.factor(Time))) +
  geom_boxplot() +
  geom_jitter(shape = 21, width = 0.1) +
  labs(y = "Chicken mass (g)", fill = "Time") +
  theme_pubr()

ggarrange(plt1, plt2, plt3, plt4, ncol = 2, nrow = 2, labels = "AUTO")

```

## 3.6 Exercises

### 3.6.1 Exercise 1

Notice how the data summary for chicken weights contained within `wt_summary` is very similar to the summary returned for `weight` when we apply `summary(chicks)`. Please use the `summarise()` approach and construct a data summary with exactly the same summary statistics for `weight` as that which `summary()` returns.

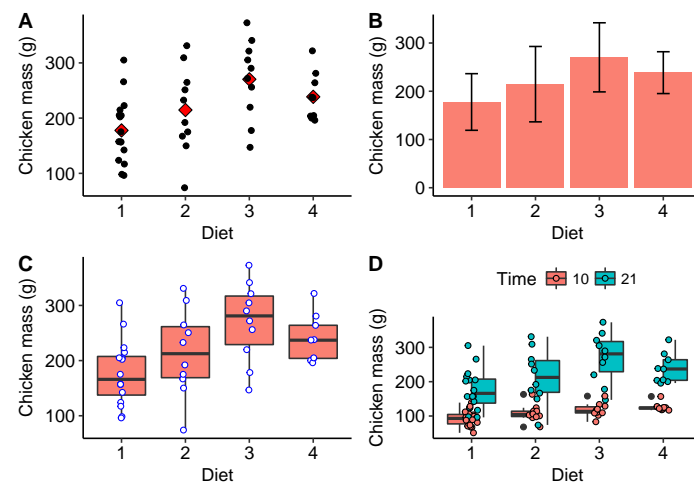


Figure 3.2: A) Scatterplot of the mean and raw chicken mass values. B) Bar graph of the chicken mass values, showing whiskers indicating  $1 \pm SD$ . C) Box and whisker plot of the chicken mass data.

# 4

## Graphical data displays

```
library(tidyverse)
library(ggpubr)
library(RColorBrewer)
library(ggthemes)
```

Here we shall provide examples of many kinds of graphical data summaries. We use **ggplot2** for these figures and will refrain from using the base R graphs. We shall provide examples of various themes so you can see what is available to use for your own plots. We also include various modifications of the default plots so you can get an idea of how to modify some of the plot characteristics, and we always make sure that our graphs are publication ready. The best way to learn is to work by example. Deeper understanding will emerge only from working through all of these example plots, and making your own changes here and there to see how your own modification will affect your graphs' appearance. Then find your own data and plot them. As always, liberally make use of the built-in help facility (the more you do, the easier it becomes, like riding a bicycle). Also, don't be shy to use Google.

### 4.1 Qualitative data

Qualitative data that describe group representivity to various categories are best presented as frequency distribution histograms (I interchangeably use histograms, frequency histograms, and frequency distribution histograms). Histograms apply to categorical data. Although it can be presented numerically in tabular form, one more frequently creates a bar or pie graph of the number of occurrences in a collection of non-overlapping classes or categories. Both the data and graphical displays will be demonstrated here.

The first case of a frequency distribution histogram is one that shows the raw counts per each of the categories that are represented in the data. The count within each of the categories (represented by a bar graph called a histogram) sums to the sample size,  $n$ . In the second case, we may want to report those data as proportions. Here we show the frequency proportion in a

collection of non-overlapping categories. For example, we have a sample size of 12 ( $n = 12$ ). In this sample, two are coloured blue, six red, and five purple. The relative proportions are  $2/12 = 0.1666667$  blue,  $6/12 = 0.5$  red, and  $5/12 = 0.4166667$  purple. The important thing to note here is that the relative proportions sum to 1, i.e.  $0.1666667 + 0.5 + 0.4166667 = 1$ . These data may be presented as a table or as a graph.

Let us demonstrate the numerical and graphical summaries using the built-in iris data:

```
# the numerical summary produced by a piped series of functions;
# create a summary of the data (i.e. number of replicates per species)
# used for (A), (B) and (C), below
iris.cnt <- iris %>%
  count(Species) %>% # automagically creates a column, n, with the counts
  mutate(prop = n / sum(n)) # creates the relative proportion of each species
iris.cnt
```

```
R> # A tibble: 3 x 3
R>   Species      n prop
R>   <fct>      <int> <dbl>
R> 1 setosa      50 0.333
R> 2 versicolor  50 0.333
R> 3 virginica   50 0.333
```

```
# a stacked bar graph with the cumulative sum of observations
plt1 <- ggplot(data = iris.cnt, aes(x = "", y = n, fill = Species)) +
  geom_bar(width = 1, stat = "identity") +
  labs(title = "Stacked bar graph", subtitle = "cumulative sum",
       x = NULL, y = "Count") +
  theme_pubclean() + scale_color_few() +
  scale_fill_few()
```

```
# a stacked bar graph with the relative proportions of observations
plt2 <- ggplot(data = iris.cnt, aes(x = "", y = prop, fill = Species)) +
  geom_bar(width = 1, stat = "identity") +
  scale_y_continuous(breaks = c(0.00, 0.33, 0.66, 1.00)) +
  labs(title = "Stacked bar graph", subtitle = "relative proportions",
       x = NULL, y = "Proportion") +
  theme_pubclean() + scale_color_few() +
  scale_fill_few()
```

```
# a basic pie chart
plt3 <- plt1 + coord_polar("y", start = 0) +
  labs(title = "Friends don't let...", subtitle = "...friends make pie charts",
       x = NULL, y = NULL) +
  scale_fill_brewer(palette = "Blues") +
  theme_minimal()
```

```
# if you seriously want a pie chart, rather use the base R function, `pie()`
```

```
# here now a bar graph...
```

```
# the default mapping of `geom_bar` is `stat = count`, which is a
# bar for each fo the categories (`Species`), with `count` along y
plt4 <- ggplot(data = iris, aes(x = Species, fill = Species)) +
```

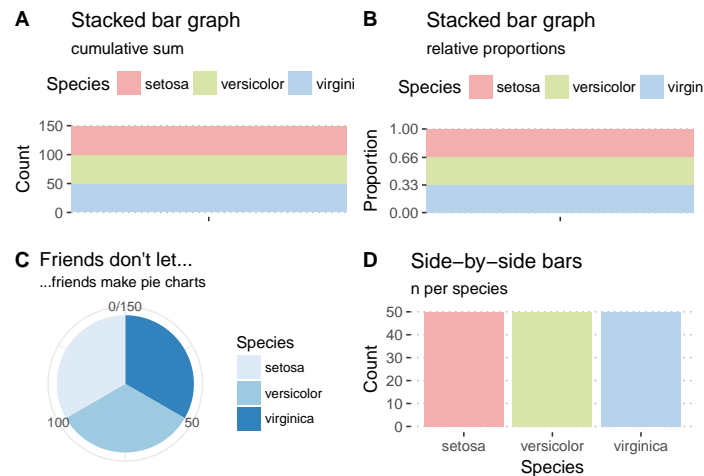


Figure 4.1: Examples of histograms for the built-in Iris data. A) A default frequency histogram showing the count of samples for each of the three species. B) A relative frequency histogram of the same data; here, the sum of counts of samples available for each of the three species is 1. C) A boring pie chart. D) A frequency histogram of raw data counts shown as a series of side-by-side bars.

```
geom_bar(show.legend = FALSE) +
  labs(title = "Side-by-side bars", subtitle = "n per species", y = "Count") +
  theme_pubclean() + scale_color_few() +
  scale_fill_few()

ggarrange(plt1, plt2, plt3, plt4, nrow = 2, ncol = 2, labels = "AUTO")
```

## 4.2 Continuous data

### 4.2.1 Frequency distributions (histograms)

As with discrete data, we have a choice of absolute (Fig. 4.2A) and relative (Fig. 4.2 B-C) frequency histograms. There's also the empirical cumulative distribution function (ECDF) (Fig. 4.2 D) that uses relative proportions, but in this instance it is the relative proportion that each individual observation has towards the sample. Since the purpose of frequency histograms is to count the number of times something takes place or occurs within a category, what do we do when we are faced with continuous data where no categories are available? We can create our own categories, called *bins*. See the Old Faithful data, for example. The eruptions last between 1.6 and 5.1 minutes. So, we create intervals of time spanning these times, and within each count the number of times an event lasts as long as denoted by the intervals. Here we might choose intervals of 1-2 minutes, 2-3 minutes, 3-4 minutes, 4-5 minutes, and 5-6 minutes. The `ggplot2` `geom_histogram()` function automatically creates the bins, but we may specify our own. It is best to explain these principles by example (see Figure 4.2 A-D).

```
# a normal frequency histogram, with count along y
hist1 <- ggplot(data = faithful, aes(x = eruptions)) +
  geom_histogram(colour = "black", fill = "salmon", alpha = 0.6) +
  labs(title = "Old Faithful data",
```

```

    subtitle = "A vanilla frequency histogram",
    x = "Eruption duration (min)",
    y = "Count") + theme_pubclean()

# when the binwidth is 1, the density histogram *is* the relative
# frequency histogram
hist2 <- ggplot(data = faithful, aes(x = eruptions)) +
  geom_histogram(aes(y = ..density..),
    position = 'identity', binwidth = 1,
    colour = "black", fill = "salmon", alpha = 0.6) +
  labs(title = "Old Faithful data",
    subtitle = "Relative frequency histogram",
    x = "Eruption duration (min)",
    y = "Count") + theme_pubclean()

# if binwidth is something other than 1, the relative frequency in
# a histogram is ..density.. * binwidth
hist3 <- ggplot(data = faithful, aes(x = eruptions)) +
  geom_histogram(aes(y = 0.5 * ..density..),
    position = 'identity', binwidth = 0.5,
    colour = "black", fill = "salmon", alpha = 0.6) +
  labs(title = "Old Faithful data",
    subtitle = "Relative frequency histogram",
    x = "Eruption duration (min)",
    y = "Relative contribution") + theme_pubclean()

# ECDF
hist4 <- ggplot(data = faithful, aes(x = eruptions)) +
  stat_ecdf() +
  labs(title = "Old Faithful data",
    subtitle = "ECDF",
    x = "Eruption duration (min)",
    y = "Relative contribution") + theme_pubclean()

ggarrange(hist1, hist2, hist3, hist4, ncol = 2, nrow = 2, labels = "AUTO")

```

What if we have continuous data belonging with multiple categories? The *iris* data provide a nice set of measurements that we may use to demonstrate a grouped frequency histogram. These data are size measurements (cm) of the variables sepal length and width and petal length and width, respectively, for 50 flowers from each of three species of *Iris*. The species are *Iris setosa*, *I. versicolor*, and *I. virginica*.

```

# first we make long data
iris.long <- iris %>%
  gather(key = "variable", value = "size", -Species)

ggplot(data = iris.long, aes(x = size)) +
  geom_histogram(position = "dodge", # omitting this creates a stacked histogram
    colour = NA, bins = 20,
    aes(fill = Species)) +
  facet_wrap(~variable) +

```

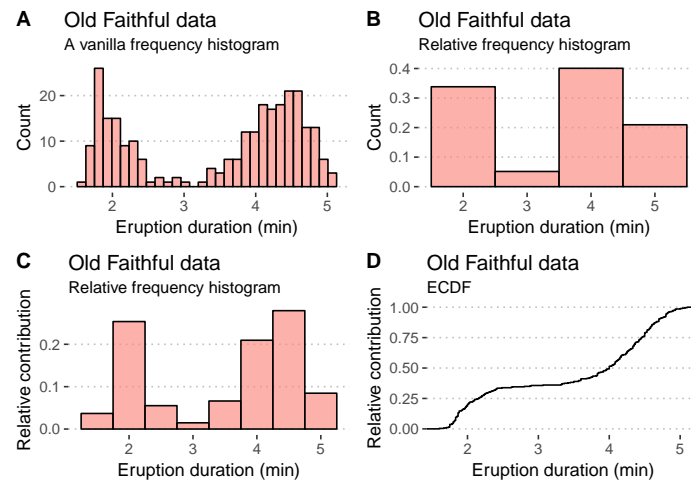


Figure 4.2: Example histograms for the Old Faithful data. A) A default frequency histogram with the count of eruption times falling within the specified bins. B) A relative frequency histogram with bins adjusted to a width of 1 minute intervals; here, the sum of counts within each of the four bins is 1. C) Another relative frequency histogram, but with the bins adjusted to each be 0.5 minute increments; again the sum of counts represented by each bin is equal to 1.

```
labs(title = "Iris data",
      subtitle = "Grouped frequency histogram",
      x = "Size (mm)",
      y = "Count") +
theme_pubclean()
```

## 4.2.2 Box plots

Box plots are sometimes called box-and-whisker plots. These graphs are a graphical representation of the data based on its quartiles as well as its smallest and largest values. The keen eye can glance the “shape” of the data distribution; they provide an alternative view to that given by the frequency distribution. A variation of the basic box-and-whisker plot is to superimpose a jittered scatter plot of the raw data on each bar.

From the `geom_boxplot` documentation, which says it best (type `?geom_boxplot`):

“The lower and upper hinges correspond to the first and third quartiles (the 25th and 75th percentiles).”

“The upper whisker extends from the hinge to the largest value no further than  $1.5 * \text{IQR}$  from the hinge (where IQR is the inter-quartile range, or distance between the first and third quartiles). The lower whisker extends from the hinge to the smallest value at most  $1.5 * \text{IQR}$  of the hinge. Data beyond the end of the whiskers are called ‘outlying’ points and are plotted individually.”

“In a notched box plot, the notches extend  $1.58 * \text{IQR} / \sqrt{n}$ . This gives a roughly 95% confidence interval for comparing medians.”

Here be examples:

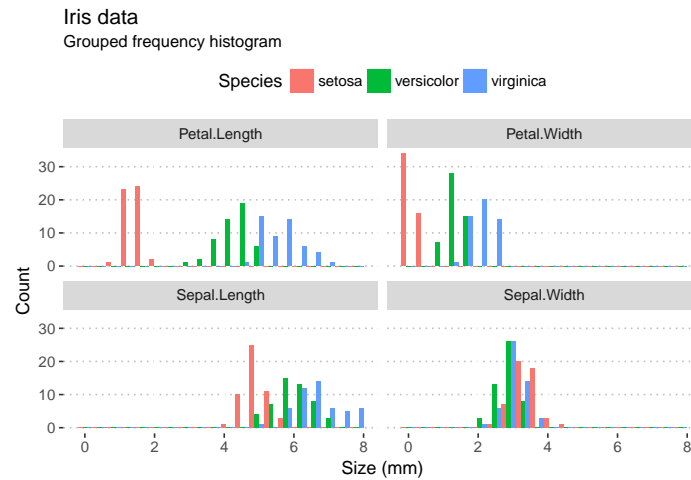


Figure 4.3: Panelled grouped histograms for the four Iris variables.

```
plt1 <- ggplot(data = iris, aes(x = Species, y = Sepal.Length, fill = Species)) +
  geom_boxplot(show.legend = FALSE, notch = FALSE) + theme_pubclean() +
  labs(y = "Sepal length (mm)") +
  theme(axis.text.x = element_text(face = "italic"))

plt2 <- ggplot(data = iris.long, aes(x = Species, y = size)) +
  geom_boxplot(fill = "red", alpha = 0.4, notch = TRUE) +
  geom_jitter(width = 0.1, shape = 21, colour = "blue", fill = NA, alpha = 0.2) +
  facet_wrap(~variable, nrow = 1) +
  labs(y = "Size (mm)") + theme_pubclean() +
  theme(axis.text.x = element_text(face = "italic")) +
  theme(axis.ticks.length=unit(-0.25, "cm"), axis.ticks.margin=unit(0.5, "cm"))

ggarrange(plt1, plt2, nrow = 2, ncol = 1, labels = "AUTO")
```

Box-and-whisker plots have traditionally been used to display data that are not normally distributed, but I like to use them for any old data, even normal data. I prefer these over the old-fashioned bar graphs (as seen later in this section).

The **ggpubr** package provides many convenience functions for the drawing of publication quality graphs, many of which include summaries of pairwise comparisons (e.g. in t-tests and ANOVAs). Please see [here](http://www.sthda.com/english/articles/24-ggpubr-publication-ready-plots/)<sup>1</sup> and [here](http://www.sthda.com/english/rpkgs/ggpubr/)<sup>2</sup>.

### 4.2.3 Pairwise Scatter plots

This graph shows the relationship between two (matched) continuous variables. The statistical strength of the relationship can be indicated by a correlation (no causal relationship implied as is the case here) or a regression (when a causal link of  $x$  on  $y$  is demonstrated).

```
plt1 <- ggplot(data = iris, aes(x = Petal.Length, y = Petal.Width, colour = Species)) +
  geom_point() +
  labs(x = "Petal length (mm)", y = "Petal width (mm)") +
```

<sup>1</sup><http://www.sthda.com/english/articles/24-ggpubr-publication-ready-plots/>

<sup>2</sup><http://www.sthda.com/english/rpkgs/ggpubr/>



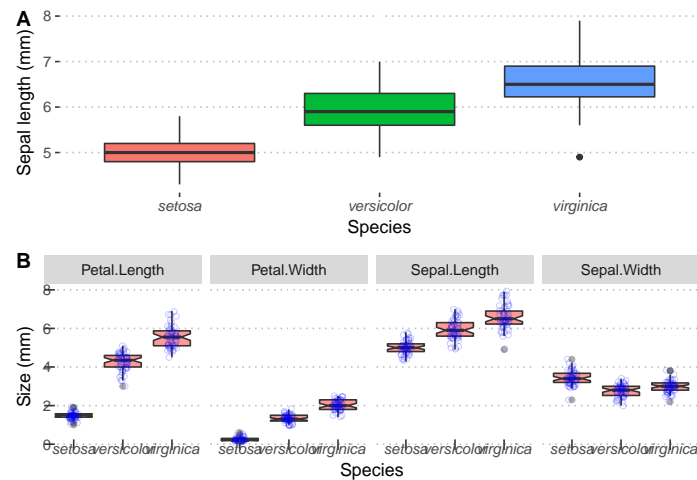


Figure 4.4: Examples of box plots made for the Iris data. A) A default box plot for one of the variables only. B) A panelled collection of box plots, one for each of the four variables, with a scatterplot to indicate the spread of the actual replicates.

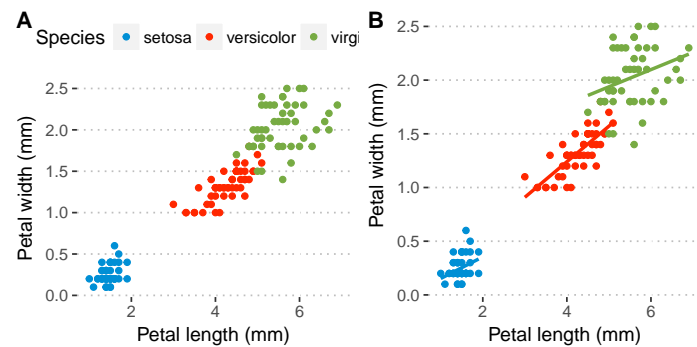
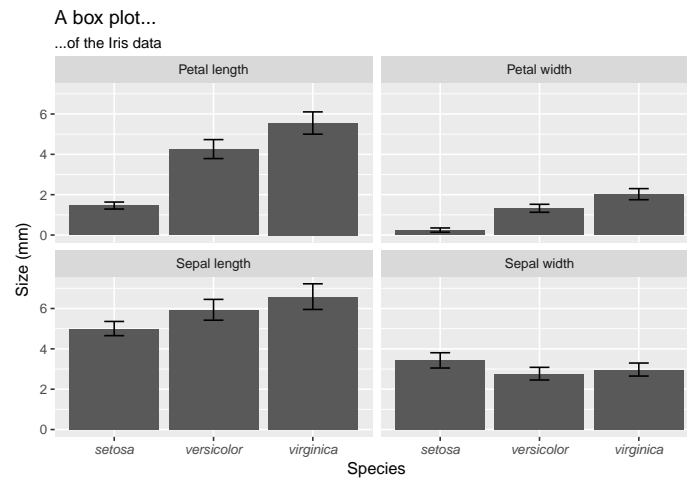


Figure 4.5: Examples of scatterplots made for the Iris data. A) A default scatter plot showing the relationship between petal length and width. B) The same as (A) but with a correlation line added.

```
theme(legend.position = c(0.18, 0.85)) +
scale_color_fivethirtyeight() +
scale_fill_fivethirtyeight() +
theme_pubclean()

plt2 <- ggplot(data = iris, aes(x = Petal.Length, y = Petal.Width, colour = Species)) +
  geom_point(show.legend = FALSE) +
  geom_smooth(method = "lm", se = FALSE, show.legend = FALSE) +
  scale_color_fivethirtyeight() +
  scale_fill_fivethirtyeight() +
  labs(x = "Petal length (mm)", y = "Petal width (mm)") +
  theme_pubclean()

ggarrange(plt1, plt2, ncol = 2, nrow = 1, labels = "AUTO")
```

Figure 4.6: Box plots of the mean  $\pm$  SD of the four Iris variables.

#### 4.2.4 Bar graphs

Bar graphs display the mean plus/minus some measure of variation around the mean—typically the standard error or the standard deviation. The mean $\pm$ SE and mean $\pm$ SD are typically used for normally-distributed data. Here I provide an example bar graph for one of the Iris data set's variables:

```
# first make nice labels for the facets because the default ones
# in the dataframe are not so nice; use the `labeller()` function
# to receive the new variable names defined here
facet.names <- c(Petal.Length = "Petal length",
                 Petal.Width = "Petal width",
                 Sepal.Length = "Sepal length",
                 Sepal.Width = "Sepal width")

# start with the `iris.long` long data that were produced above
# we create summaries of mean and SD and squirt it directly
# into the ggplot functions
iris.long %>%
  group_by(Species, variable) %>%
  summarise(mean.size = mean(size),
            sd.size = sd(size)) %>%
  ggplot(aes(x = Species, y = mean.size)) +
  geom_bar(stat = "identity") +
  geom_errorbar(aes(ymin = mean.size - sd.size, ymax = mean.size + sd.size), width = 0.2) +
  facet_wrap(~variable, labeller = labeller(variable = facet.names)) +
  labs(y = "Size (mm)", title = "A box plot...", subtitle = "...of the Iris data") +
  theme(axis.text.x = element_text(face = "italic"))
```

#### 4.2.5 Density graphs

Often when we are displaying a distribution of data we are interested in the “shape” of the data more than the actual count of values in a specific category, as shown by a standard histogram. When one wishes to more organically visualise the frequency of values in a sample set

a density graphs is used. These may also be thought of as smooth histograms. These work well with histograms and rug plots, as we may see in the figure below. It is important to note with density plots that they show the relative density of the distribution along the Y axis, and *not* the counts of the data. This can of course be changed, as seen below, but is not the default setting. Sometimes it can be informative to see how different the count and density distributions appear.

```
# a normal density graph
dens1 <- ggplot(data = faithful, aes(x = eruptions)) +
  geom_density(colour = "black", fill = "salmon", alpha = 0.6) +
  labs(title = "Old Faithful data",
        subtitle = "A vanilla density plot",
        x = "Eruption duration (min)",
        y = "Density") + theme_pubr()

# a density and rug plot combo
dens2 <- ggplot(data = faithful, aes(x = eruptions)) +
  geom_density(colour = "black", fill = "salmon", alpha = 0.6) +
  geom_rug(colour = "red") +
  labs(title = "Old Faithful data",
        subtitle = "A density and rug plot",
        x = "Eruption duration (min)",
        y = "Density") + theme_pubr()

# a relative frequency histogram overlayed with a density plot
dens3 <- ggplot(data = faithful, aes(x = eruptions)) +
  geom_histogram(aes(y = ..density..),
                 position = 'identity', binwidth = 1,
                 colour = "black", fill = "turquoise", alpha = 0.6) +
  geom_density(colour = "black", fill = "salmon", alpha = 0.6) +
  labs(title = "Old Faithful data",
        subtitle = "Relative frequency with density",
        x = "Eruption duration (min)",
        y = "Density") + theme_pubr()

# a normal frequency histogram with density overlayed
# note that the density curve must be adjusted by
# the number of data points times the bin width
dens4 <- ggplot(data = faithful, aes(x = eruptions)) +
  geom_histogram(aes(y = ..count..),
                 binwidth = 0.2, colour = "black", fill = "turquoise", alpha = 0.6) +
  geom_density(aes(y = ..density.. * nrow(datasets::faithful) * 0.2), position = "identity",
               colour = "black", fill = "salmon", alpha = 0.6) +
  labs(title = "Old Faithful data",
        subtitle = "Frequency with density",
        x = "Eruption duration (min)",
        y = "Count") + theme_pubr()

ggarrange(dens1, dens2, dens3, dens4, ncol = 2, nrow = 2, labels = "AUTO")
```

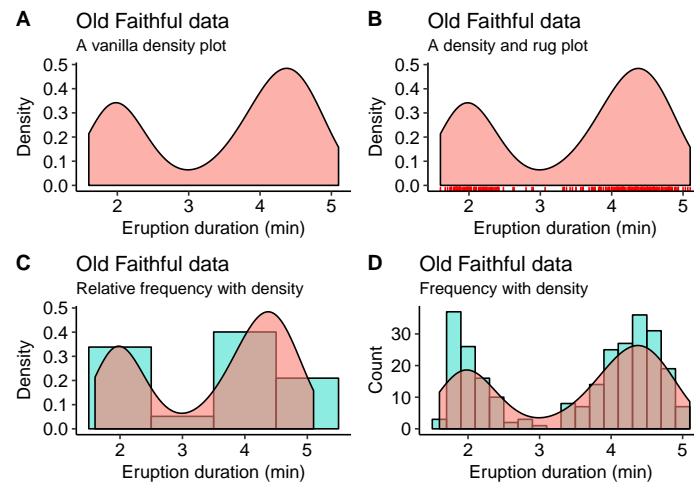


Figure 4.7: A bevy of density graphs option based on the iris data. A) A lone density graph. B) A density graph accompanied by a rug plot. C) A histogram with a density graph overlay. D) A ridge plot.

#### 4.2.6 Violin plots

The density graph is not limited to its use with histograms. We may combine this concept with box plots, too. These are known as violin plots and are very useful when we want to show the distribution of multiple categories of the same variable alongside one another. Violin plots may show the same information as box plots but take things one step further by allowing the shape of the boxplot to also show the distribution of the data within the sample set. We will use the iris data below to highlight the different types of violin plots one may use.

```
# A basic violin plot
vio1 <- ggplot(data = iris, aes(x = Species, y = Sepal.Length, fill = Species)) +
  geom_violin() +
  theme_pubclean() + theme(legend.position = "none") +
  labs(title = "Iris data",
        subtitle = "Basic violin plot", y = "Sepal length (mm)") +
  theme(axis.text.x = element_text(face = "italic"))

# A violin plot showing the quartiles as lines
vio2 <- ggplot(data = iris, aes(x = Species, y = Sepal.Length, fill = Species)) +
  geom_violin(show.legend = FALSE, draw_quantiles = c(0.25, 0.5, 0.75)) +
  theme_pubclean() + theme(legend.position = "none") +
  labs(title = "Iris data",
        subtitle = "Violin plot with quartiles", y = "Sepal length (mm)") +
  theme(axis.text.x = element_text(face = "italic"))

# Box plots nested within violin plots
vio3 <- ggplot(data = iris, aes(x = Species, y = Sepal.Length, colour = Species)) +
  geom_violin(fill = "grey70") +
  geom_boxplot(width = 0.1, colour = "grey30", fill = "white") +
  theme_pubclean() + theme(legend.position = "none") +
  labs(title = "Iris data",
        subtitle = "Box plots nested within violin plots", y = "Sepal length (mm)") +
```

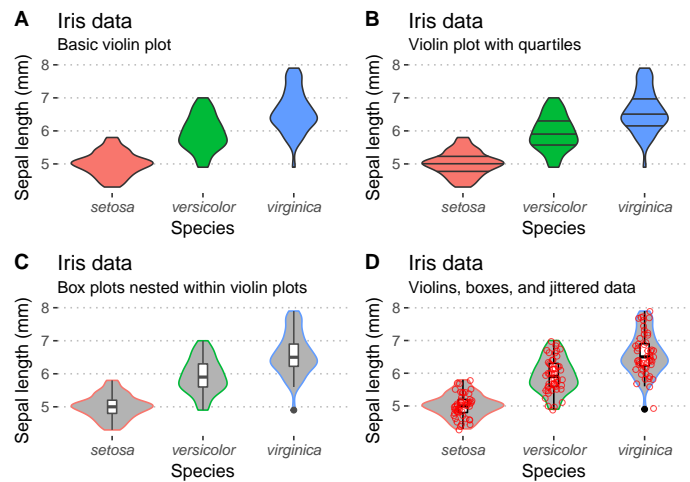


Figure 4.8: Variations of violin plots.

```

theme(axis.text.x = element_text(face = "italic"))

# Boxes in violins with the raw data jittered about
vio4 <- ggplot(data = iris, aes(x = Species, y = Sepal.Length, colour = Species)) +
  geom_violin(fill = "grey70") +
  geom_boxplot(width = 0.1, colour = "black", fill = "white") +
  geom_jitter(shape = 1, width = 0.1, colour = "red", alpha = 0.7, fill = NA) +
  theme_pubclean() + theme(legend.position = "none") +
  labs(title = "Iris data",
        subtitle = "Violins, boxes, and jittered data", y = "Sepal length (mm)") +
  theme(axis.text.x = element_text(face = "italic"))

ggarrange(vio1, vio2, vio3, vio4, ncol = 2, nrow = 2, labels = "AUTO")

```

## 4.3 Exercises

### 4.3.1 Exercise 1

Choose a dataset, either one of the many built into R or one of your own, and create four distinctly different figures. Use `ggarrange()` to stitch them together in a meaningful way.

# 5

## Distributions

Therefore, we must next learn about the different types of data distributions we are likely to encounter in the wild.

### 5.1 Discrete distributions

A discrete random variable has a finite or countable number of possible values. As the name suggests, it models integer data. Below we provide options to generate and visualise data belonging to several classes of discrete distributions. Later (Chapter 11) we will learn how to transform these data prior to performing the appropriate statistical analysis.

#### 5.1.1 Bernoulli distribution

A Bernoulli random variable,  $x$ , takes the value 1 with probability  $p$  and the value 0 with probability  $q = 1 - p$ . It is used to represent data resulting from a *single* experiment with binary (yes or no; black or white; positive or negative; success or failure; dead or alive;) outcomes, such as a coin toss—there are only two options, heads or tails. Nothing else. Here,  $p$  represents the probability of the one outcome and  $q$  the probability of the other outcome. The distribution of the possible outcomes,  $x$ , is given by:

$$f(x; p) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$$

#### 5.1.2 Binomial distribution

A binomial random variable,  $x$ , is the sum of  $n$  independent Bernoulli random variables with parameter  $p$ . This data distribution results from repeating identical experiments that produce a binary outcome with probability  $p$  a specified number of times, and choosing  $n$  samples at random. As such, it represents a collection of Bernoulli trials.

$$f(x; n, p) = \binom{n}{x} p^x (1-p)^{n-x}$$

### 5.1.3 Negative binomial distribution

A negative binomial random variable,  $x$ , counts the number of successes in a sequence of independent Bernoulli trials with probability  $p$  before  $r$  failures occur. This distribution could for example be used to predict the number of heads that result from a series of coin tosses before three tails are observed:

$$f(x; n, r, p) = \binom{x+r-1}{x} p^x (1-p)^r$$

where  $x$  is the number of successes,  $r$  is the number of failures, and  $p$  is the probability of success

### 5.1.4 Geometric distribution

A geometric random variable,  $x$ , represents the number of trials that are required to observe a single success. Each trial is independent and has success probability  $p$ . As an example, the geometric distribution is useful to model the number of times a die must be tossed in order for a six to be observed. It is given by:

$$f(x; p) = (1-p)^x p$$

### 5.1.5 Poisson distribution

A Poisson random variable,  $x$ , tallies the number of events occurring in a fixed interval of time or space, given that these events occur with an average rate  $\lambda$ . Poisson distributions can be used to model events such as meteor showers and or number of people entering a shopping mall. This equation describes the Poisson distribution:

$$f(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$

## 5.2 Continuous distributions

### 5.2.1 Normal distribution

Another name for this kind of distribution is a Gaussian distribution. A random sample with a Gaussian distribution is normally distributed. These values are identically distributed and independent—we say they are independent and identically distributed random variables (i.i.d.), and they have an expected mean given by  $\mu$  (or  $\hat{x}$  in Chapter 3.2.1) and a finite variance given by  $\sigma^2$  (or  $S^2$  in Chapter 3.3.1); if the number of samples drawn from a population is sufficiently large, the estimated mean and SD will be indistinguishable from the population (as per the central limit theorem).

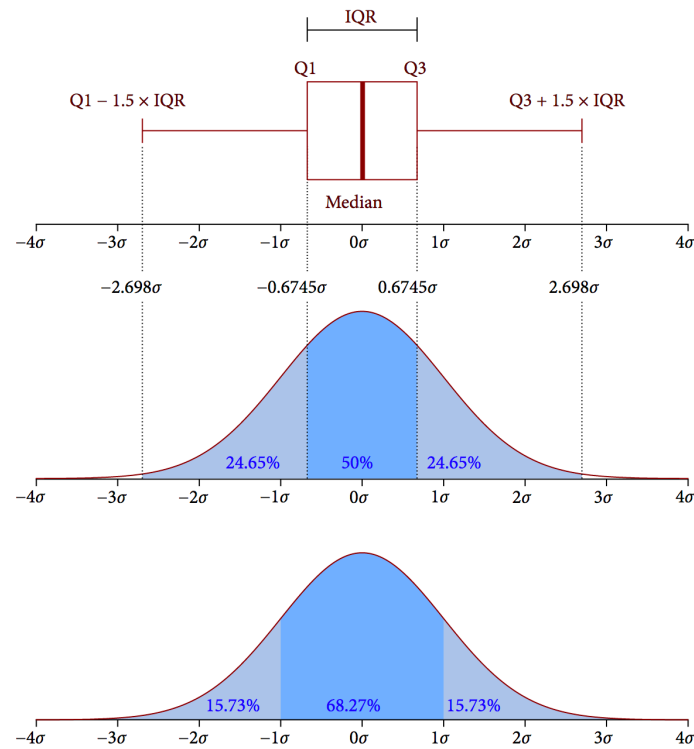


Figure 5.1: Boxplot and probability density function of a normal distribution. Credit: Wikipedia.

### 5.2.2 Uniform distribution

The continuous uniform distribution is sometime called a rectangular distribution. Simply, it states that all measurements of the same magnitude included with this distribution are equally probable. This is basically random numbers.

### 5.2.3 Student T distribution

This is a continuous probability distribution that arises when estimating the mean of a normally distributed population in situations where the sample size is small and population standard deviation is unknown. It is used in the statistical significance testing between the means of different sets of samples, and not much so in the modelling of natural phenomena.

### 5.2.4 Chi-squared distribution

Mostly used in hypothesis testing, but not to encapsulate the distribution of data drawn to represent natural phenomena.

### 5.2.5 Exponential distribution

This is a probability distribution that describes the time between events in a Poisson point process, i.e., a process in which events occur continuously and independently at a constant average rate.

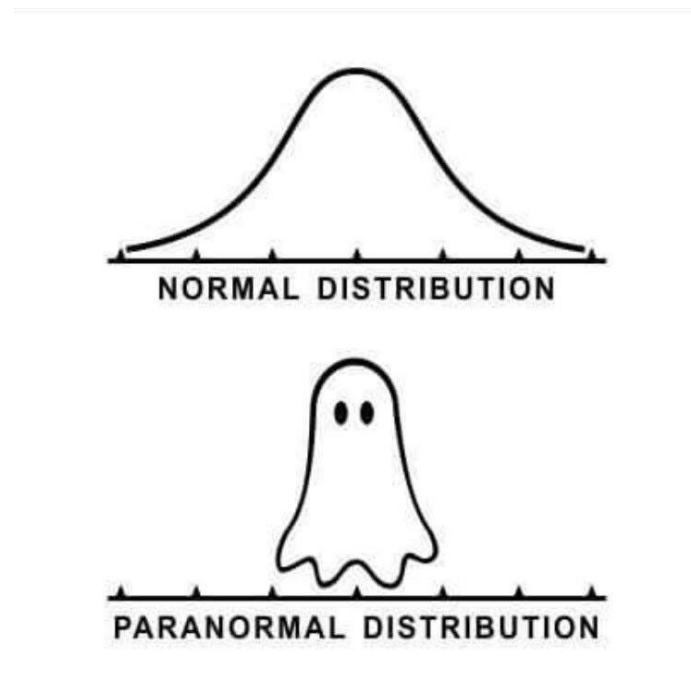


### 5.2.6 F distribution

### 5.2.7 Gamma distribution

### 5.2.8 Beta distribution

### 5.2.9 Paranormal distributions



## 5.3 Finding one's data distribution

data belonging to a sample will never exactly follow a specific distribution, even when the test for normality says it does—there will always be a small probability that they are non-normal and is in fact better described by some other distribution. In other words, data are only *compatible* with a certain distribution, and one can never answer the question “Does my data follow the distribution *xy* exactly?” as simply as providing a yes/no answer. So what now? How does one find one's data distribution? We can use the *Cullen and Frey graph* function that lives in the **fitdistrplus** package. This graph tells us whether the skewness and kurtosis of our data are consistent with that of a particular distribution. We will demonstrate by generating various data distributions and testing them using the Cullen and Frey graph.

```
library(fitdistrplus)
```

```
## Loading required package: MASS
```

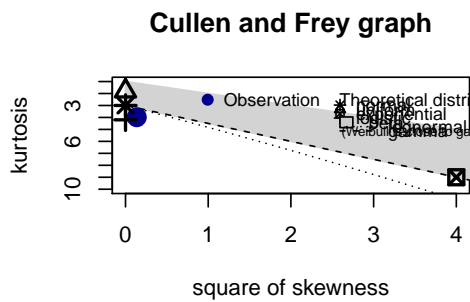
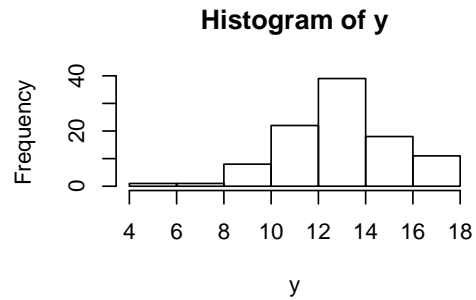
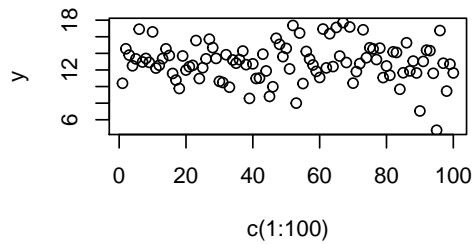
```
## Loading required package: survival
```

```
library(logspline)
```

```
# Generate log-normal data
```

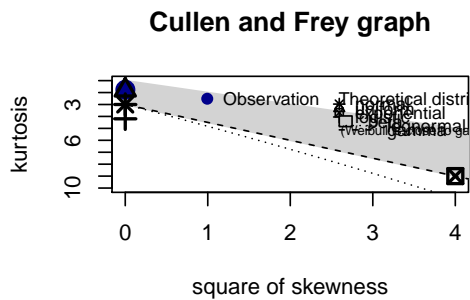
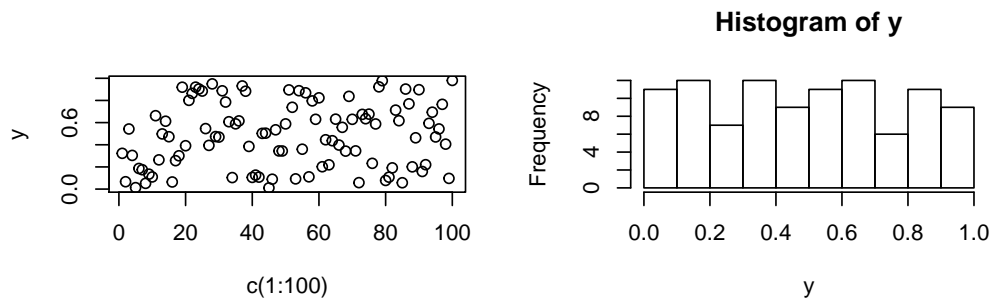


```
## min: 4.737749 max: 17.71292
## median: 12.90348
## mean: 12.91636
## estimated sd: 2.320356
## estimated skewness: -0.3688259
## estimated kurtosis: 3.981041
```



```
# uniformly distributed data
y <- runif(100)
par(mfrow = c(2, 2))
plot(x = c(1:100), y = y)
hist(y)
descdist(y, discrete = FALSE)
```

```
## summary statistics
## -----
## min: 0.01084804 max: 0.9804608
## median: 0.4858252
## mean: 0.4852429
## estimated sd: 0.2916798
## estimated skewness: 0.03884635
## estimated kurtosis: 1.761942
```



```
# uniformly distributed data
```

```
y <- rexp(100, 0.7)
```

```
par(mfrow = c(2, 2))
```

```
plot(x = c(1:100), y = y)
```

```
hist(y)
```

```
descdist(y, discrete = FALSE)
```

```
## summary statistics
```

```
## -----
```

```
## min: 0.02697365 max: 6.15109
```

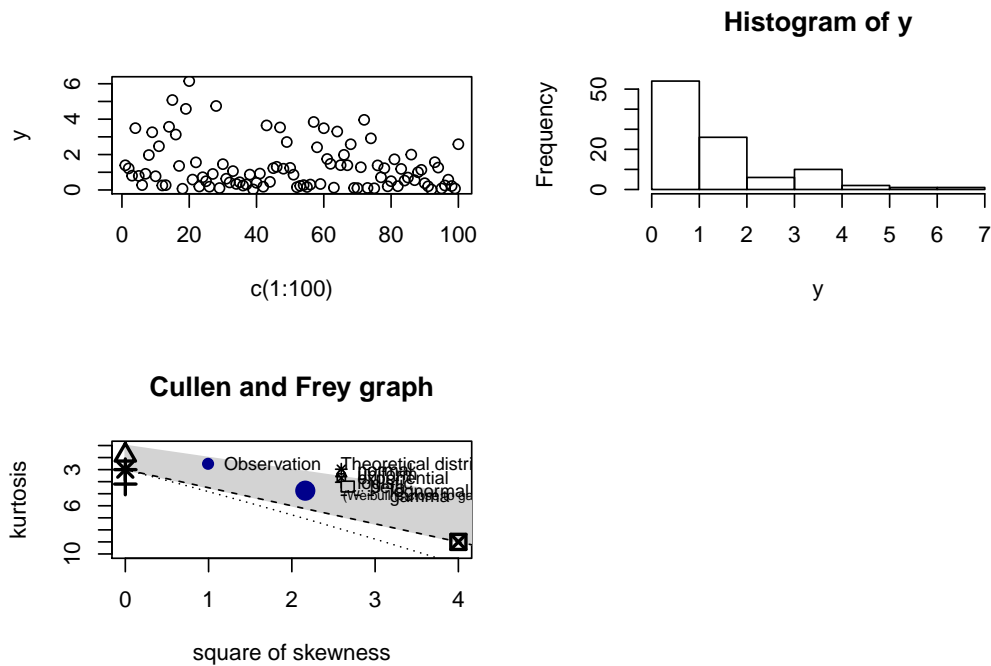
```
## median: 0.8649283
```

```
## mean: 1.290057
```

```
## estimated sd: 1.318786
```

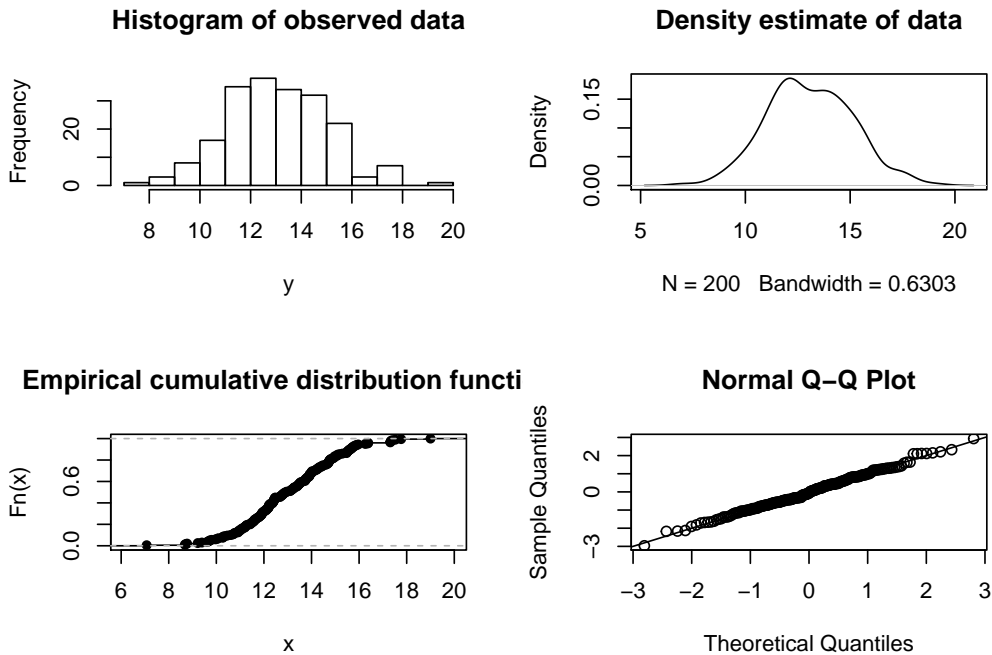
```
## estimated skewness: 1.470216
```

```
## estimated kurtosis: 4.747477
```



There is also a whole bunch of other approaches to use to try and identify the data distribution. Let us start with the gold standard first: normal data. We will demonstrate some visualisation approaches. The one that you already know is a basic histogram; it tells us something about the distribution's skewness, the tails, the mode(s) of the data, outliers, etc. Histograms can be compared to shapes associated with idealistic (simulated) distributions, as we will do here.

```
y <- rnorm(n = 200, m = 13, sd = 2)
par(mfrow = c(2, 2))
# using some basic base graphics as ggplot2 is overkill;
# we can get a histogram using hist() statement
hist(y, main = "Histogram of observed data")
plot(density(y), main = "Density estimate of data")
plot(ecdf(y), main = "Empirical cumulative distribution function")
# standardise the data
z.norm <- (y - mean(y)) / sd(y)
# make a qqplot
qqnorm(z.norm)
# add a 45-degree reference line
abline(0, 1)
```



Above we have also added a diagonal line to the qqplot. If the sampled data come from the population with the chosen distribution, the points should fall approximately along this reference line. The greater the departure from this reference line, the greater the evidence for the conclusion that the data set have come from a population with a different distribution.

```
# curve(dnorm(100, m = 10, sd = 2), from = 0, to = 20, main = "Normal distribution")
# curve(dgamma(100, scale = 1.5, shape = 2), from = 0, to = 15, main = "Gamma distribution")
# curve(dweibull(100, scale = 2.5, shape = 1.5), from = 0, to = 15, main = "Weibull distribution")
```

## 5.4 Exercises

### 5.4.1 Exercise 1

Choose two different datasets and plot them as histograms with density curves overlaid. Label them with the distribution they appear to be and stitch them together with `ggarrange()`.

# 6

## Inferences about one or two populations

```
library(tidyverse)
library(plotly)
```

At the heart of many basic scientific inquiries is the simple question “Is A different from B?” The scientific notation for this question is:

- $H_0$ : Group A is not different from group B.
- $H_1$ : Group A is different from group B.

More formally, one would say:

1.  $H_0 : \bar{A} = \bar{B}$  vs. the alternative hypothesis that  $H_a : \bar{A} \neq \bar{B}$ .
2.  $H_0 : \bar{A} \leq \bar{B}$  vs. the alternative hypothesis that  $H_a : \bar{A} > \bar{B}$ .
3.  $H_0 : \bar{A} \geq \bar{B}$  vs. the alternative hypothesis that  $H_a : \bar{A} < \bar{B}$ .

*NOTE:* Hypothesis 1 is a two-sided  $t$ -test and hypotheses 2 and 3 are one-sided tests.

Biologists typically define the probability of one in twenty (0.05) as the cutoff level to reject the null hypothesis.

To answer this fundamental question one often uses a  $t$ -test. There are several variations of  $t$ -tests, depending on the nature of our samples and the type of question being asked:

- **One-sample  $t$ -tests:** only one sample set of data that we wish to compare against a known population mean:
  - one-sided one-sample  $t$ -tests
  - two-sided one-sample  $t$ -tests
- **Two-sample  $t$ -tests:** the means of two groups are compared against each other:
  - independent sample  $t$ -tests
    - one-sided two-sample  $t$ -tests

- two-sided two-sample  $t$ -tests
- paired sample  $t$ -tests
  - one-sided
  - two-sided

Before we cover each of these, we need to understand some of the assumptions behind  $t$ -tests. We shall cover that next.

## 6.1 Assumptions

Irrespective of the kind of  $t$ -test, we have to make a couple of important assumptions that are *not* guaranteed to be true. In fact, these assumptions are often violated because real data, especially biological data, are messy. In order to use a  $t$ -test to determine if a significant difference between two sample sets of data exists we must first establish that:

- the dependent variable must be continuous (i.e. it is measured at the interval or ratio level),
- the observations in the groups being compared are independent of each other,
- the data are **normally distributed**, and
- that the data are **homoscedastic**, and in particular, that there are no outliers.

### 6.1.1 Normality

Remember from Chapter 5 what a normal distribution is/looks like? Let's have a peek below to remind ourselves:

```
# Random normal data
set.seed(666)
r_dat <- data.frame(dat = c(rnorm(n = 1000, mean = 10, sd = 3),
                           rnorm(n = 1000, mean = 8, sd = 2)),
                   sample = c(rep("A", 1000), rep("B", 1000)))

# Create histogram
h <- ggplot(data = r_dat, aes(x = dat, fill = sample)) +
  geom_histogram(position = "dodge", binwidth = 1, alpha = 0.8) +
  geom_density(aes(y = 1*..count.., fill = sample), colour = NA, alpha = 0.4) +
  labs(x = "value")
h
```

Whereas histograms may be a pretty way to check the normality of our data, there is actually a statistical test for this, which is preferable to a visual inspection alone. But remember that you should *always* visualise your data before performing any statistics on them. To check the normality of the data we use the Shapiro-Wilk test. This test produces a  $W$  value and a  $p$ -value. We are only interested in the  $p$ -value as this is how we are to determine the normality of the data. The Shapiro-Wilk test tests the null hypothesis that a sample  $x_1, \dots, x_n$  comes from a normally distributed population. i.e. that the normality *does not* differ significantly from normal. If the  $p$ -value is **above** 0.05 we may assume the data to be normally distributed. In order to demonstrate what the output of `shapiro.test()` looks like we will run it on all of the random data we generated.

```
shapiro.test(r_dat$dat)
```

```
R>
```



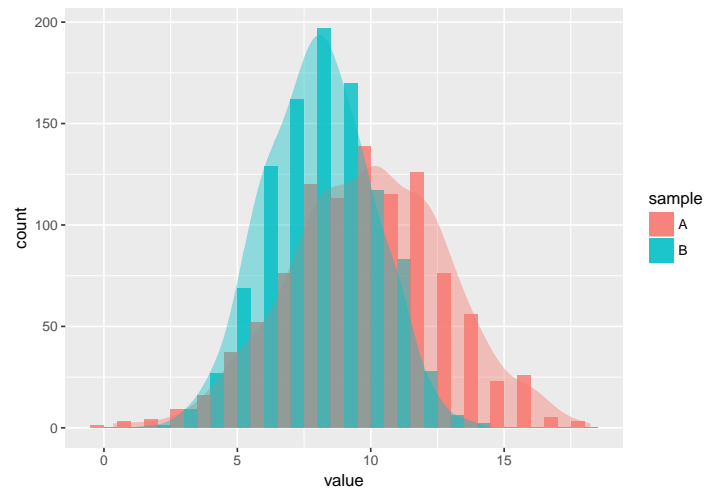


Figure 6.1: Interactive histogram showing two randomly generated normal distributions.

```
R> Shapiro-Wilk normality test
R>
R> data:  r_dat$dat
R> W = 0.9942, p-value = 4.649e-07
```

Note that this shows that the data are *not* normally distributed. This is because we have incorrectly run this function simultaneously on two different samples of data. To perform this test correctly, and in the tidy way, we need to select only the second piece of information from the `shapiro.test()` output and ensure that it is presented as a numeric value:

```
# we use the square bracket notation to select only the p-value;
# had we used `[1]` we'd have gotten W
r_dat %>%
  group_by(sample) %>%
  summarise(norm_dat = as.numeric(shapiro.test(dat)[2]))
```

```
R> # A tibble: 2 x 2
R>   sample norm_dat
R>   <fct>    <dbl>
R> 1 A      0.375
R> 2 B      0.461
```

Now we see that our two sample sets are indeed normally distributed.

What if we find that the data are not normally distributed? Although there are many options, the easiest is to perform a Wilcoxon Rank Sum test, which is the non-parametric equivalent to a *t*-test (see Section X). Another option is to transform the data (Chapter 11).

### 6.1.2 Homoscedasticity

Besides requiring that our data are normally distributed, we must also ensure that they are homoscedastic. This word means that the scedasticity (variance) of things are homogeneous (similar). In practical terms this means that the variance of the samples we are comparing should not be more than two to four times greater than one another. In R, we use the function `var()` to check the variance in a sample:

```
r_dat %>%
  group_by(sample) %>%
  summarise(sample_var = var(dat))
```

```
R> # A tibble: 2 x 2
R>   sample sample_var
R>   <fct>      <dbl>
R> 1 A         8.72
R> 2 B         3.97
```

Above we see that the variance of our two samples are homoscedastic because the variance of one is not more than two to four times greater than the other.

What if our data are not equal in their variances? This is easier to fix as the solution is built right into the  $t$ -test function; all we need to do is to perform Welch Two Sample  $t$ -test (the default) in the `t.test()` function that we shall use below. If the variances are equal, we simply perform a normal Student's  $t$ -test by setting the argument `var.equal = TRUE` in the function call (see below).

### 6.1.3 Two for one

Because these two assumptions of normality and homoscedasticity are performed in tandem with one another, it is helpful to have a function that checks for both simultaneously. Below we see how just such a function would be written:

```
two_assum <- function(x) {
  x_var <- var(x)
  x_norm <- as.numeric(shapiro.test(x)[2])
  result <- c(x_var, x_norm)
  return(result)
}
```

To use our new function in a tidy way we use the following code:

```
r_dat %>%
  group_by(sample) %>%
  summarise(sample_var = two_assum(dat)[1],
            sample_norm = two_assum(dat)[2])
```

```
R> # A tibble: 2 x 3
R>   sample sample_var sample_norm
R>   <fct>      <dbl>      <dbl>
R> 1 A         8.72         0.375
R> 2 B         3.97         0.461
```

Do these data meet our assumptions? How do we know this?

Once we have tested our assumptions we may perform a  $t$ -test to ascertain whether or not our samples are significantly different from one another. The base R function for this is `t.test()`; however, by utilising the **ggpubr** package we gain access to `compare_means()`, which allows us to perform any sort of test that compares sample sets of data and outputs the results as a dataframe. We will see throughout this and the following chapters why this is so useful.

```
library(ggpubr)
```

## 6.2 One-sample $t$ -tests

Generally when we use a  $t$ -test it will be a two-sample  $t$ -test (see below). Occasionally, however, we may have only one sample set of data that we wish to compare against a known population mean, which is generally denoted as  $\mu$ , or `mu` in the function call to the  $t$ -test in R:

```
# create a single sample of random normal data
set.seed(666)
r_one <- data.frame(dat = rnorm(n = 20, mean = 20, sd = 5),
                    sample = "A")

# check normality
shapiro.test(r_one$dat)
```

```
R>
R> Shapiro-Wilk normality test
R>
R> data:  r_one$dat
R> W = 0.94911, p-value = 0.3538
```

```
# No variance to compare
# ...

# compare random data against a population mean of 20
t.test(r_one$dat, mu = 20)
```

```
R>
R> One Sample t-test
R>
R> data:  r_one$dat
R> t = 0.0048653, df = 19, p-value = 0.9962
R> alternative hypothesis: true mean is not equal to 20
R> 95 percent confidence interval:
R>  16.91306 23.10133
R> sample estimates:
R> mean of x
R>  20.00719
```

```
# compare random data against a population mean of 30
t.test(r_one$dat, mu = 30)
```

```
R>
R> One Sample t-test
R>
R> data:  r_one$dat
R> t = -6.7596, df = 19, p-value = 1.858e-06
R> alternative hypothesis: true mean is not equal to 30
R> 95 percent confidence interval:
R>  16.91306 23.10133
R> sample estimates:
R> mean of x
R>  20.00719
```

What do the results of these two different tests show? Let's visualise these data to get a better

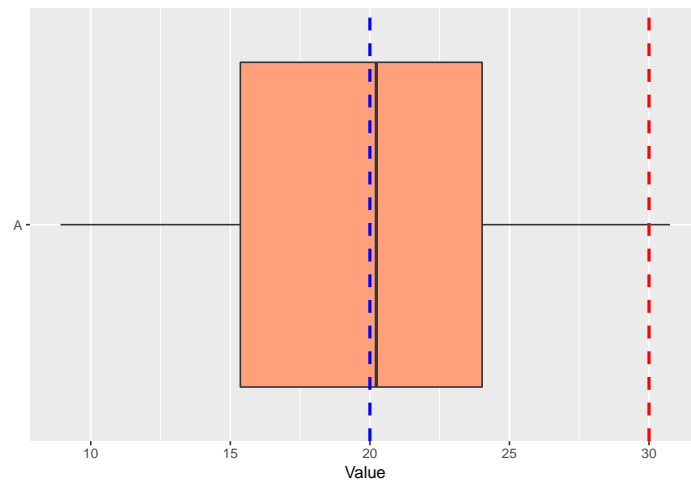


Figure 6.2: Boxplot of random normal data with. A hypothetical population mean of 20 is shown as a blue line, with the red line showing a mean of 30.

understanding.

```
ggplot(data = r_one, aes(y = dat, x = sample)) +
  geom_boxplot(fill = "lightsalmon") +
  # population mean (mu) = 20
  geom_hline(yintercept = 20, colour = "blue",
    size = 1, linetype = "dashed") +
  # population mean (mu) = 30
  geom_hline(yintercept = 30, colour = "red",
    size = 1, linetype = "dashed") +
  labs(y = "Value", x = NULL) +
  coord_flip()
```

The boxplot above shows the distribution of our random data against two potential population means. Does this help now to illustrate the results of our one-sample  $t$ -tests?

### 6.2.1 One-sided one-sample $t$ -tests

As we may remember from Chapter 5, a distribution has two tails. When we are testing for significance we are generally looking for a result that lays in the far end of either of these tails. Occasionally, however, we may want to know if the result lays specifically in one of the tails. Explicitly the leading or trailing tail. For example, is the mean value of our sample population significantly *greater than* the value  $\mu$ ? Or, is the mean value of our sample population significantly *less than* the value  $\mu$ ? To specify this in R we must add an argument as seen below:

```
# check against the trailing tail
t.test(r_one$dat, mu = 30, alternative = "less")
```

```
R>
R> One Sample t-test
R>
R> data:  r_one$dat
R> t = -6.7596, df = 19, p-value = 9.292e-07
R> alternative hypothesis: true mean is less than 30
```

```

R> 95 percent confidence interval:
R>      -Inf 22.56339
R> sample estimates:
R> mean of x
R>  20.00719

# check against the leading tail
t.test(r_one$dat, mu = 30, alternative = "greater")

R>
R> One Sample t-test
R>
R> data:  r_one$dat
R> t = -6.7596, df = 19, p-value = 1
R> alternative hypothesis: true mean is greater than 30
R> 95 percent confidence interval:
R>  17.451      Inf
R> sample estimates:
R> mean of x
R>  20.00719

```

Are these the results we would have expected? Why does the second test not return a significant result?

**TASK:** Create a visualisation to graphically demonstrate the outcome of this  $t$ -test.

## 6.2.2 Two-sided one-sample $t$ -tests

In R, the default setting for any comparison of means test is that it is two-sided so we do not need to state this explicitly. For the sake of thoroughness let's see how to do this below. Note that the results for the two following tests are identical:

```

# R assumes we want a to-sided test
t.test(r_one$dat, mu = 30)

R>
R> One Sample t-test
R>
R> data:  r_one$dat
R> t = -6.7596, df = 19, p-value = 1.858e-06
R> alternative hypothesis: true mean is not equal to 30
R> 95 percent confidence interval:
R>  16.91306 23.10133
R> sample estimates:
R> mean of x
R>  20.00719

# but we can be explicit as we choose
t.test(r_one$dat, mu = 30, alternative = "two.sided")

R>
R> One Sample t-test
R>
R> data:  r_one$dat
R> t = -6.7596, df = 19, p-value = 1.858e-06

```

```

R> alternative hypothesis: true mean is not equal to 30
R> 95 percent confidence interval:
R>  16.91306 23.10133
R> sample estimates:
R> mean of x
R>  20.00719

```

### 6.3 Two-sample $t$ -tests

A two-sample  $t$ -test is used when we have samples from two different populations that we would like to compare against one another. This is the most common use of a  $t$ -test. Two sample  $t$ -tests can accommodate samples with equal variances or samples with unequal variances (as determined by the test for heteroscedasticity that we performed earlier).

In the case of samples that share the same variance we perform a classical  $t$ -test (otherwise known as Student's  $t$ -test); the equation to calculate the  $t$ -statistic is this one:

$$t = \frac{\bar{A} - \bar{B}}{\sqrt{\frac{S^2}{n} + \frac{S^2}{m}}}$$

$\bar{A}$  and  $\bar{B}$  are the means for groups  $A$  and  $B$ , respectively;  $n$  and  $m$  are the sample sizes of the two sets of samples, respectively; and  $S^2$  is the pooled variance, which is calculated as:

$$S^2 = \frac{(n-1)S_A^2 + (m-1)S_B^2}{n+m-2}$$

The degrees of freedom,  $d.f.$ , in the equation for the shared variance is  $n_A + n_B - 2$ .

When variances are unequal we perform the Welch's  $t$ -test; Welch's  $t$ -statistics is calculated as per this equation:

$$t = \frac{\bar{A} - \bar{B}}{\sqrt{\frac{S_A^2}{n} + \frac{S_B^2}{m}}}$$

Above,  $S_A$  and  $S_B$  are the variances of groups  $A$  and  $B$ , respectively (see Section X). The degrees of freedom to use with Welch's  $t$ -test is obtained using the Welch–Satterthwaite equation as:

$$d.f. = \frac{\left(\frac{S_A^2}{n} + \frac{S_B^2}{m}\right)^2}{\left(\frac{S_A^4}{n-1} + \frac{S_B^4}{m-1}\right)}$$

What do we do with this  $t$ -statistic? In the olden days we had to calculate the  $t$ -statistics and the  $d.f.$  by hand. These two values, the  $d.f.$  and  $t$ -value had to be read off a table of pre-calculated  $t$ -values, probabilities and degrees of freedom as in here<sup>1</sup>. Luckily, the  $t$ -test function nowadays does this all automagically. But if you are feeling nostalgic over times that you have sadly never experienced, please calculate the  $t$ -statistic and the  $d.f.$  yourself and give the table a go.

<sup>1</sup><https://home.ubalt.edu/ntsbarsh/Business-stat/StatistialTables.pdf>

Back to the present day and the wonders of modern technology. Let's generate some new random normal data and test to see if the data belonging to the two groups differ significantly from one-another. First, we apply the  $t$ -test function as usual:

```
# random normal data
set.seed(666)
r_two <- data.frame(dat = c(rnorm(n = 20, mean = 4, sd = 1),
                           rnorm(n = 20, mean = 5, sd = 1)),
                   sample = c(rep("A", 20), rep("B", 20)))

# perform t-test
# note how we set the `var.equal` argument to TRUE because we know
# our data has the same SD (they are simulated as such!)
t.test(dat ~ sample, data = r_two, var.equal = TRUE)

R>
R> Two Sample t-test
R>
R> data:  dat by sample
R> t = -1.9544, df = 38, p-value = 0.05805
R> alternative hypothesis: true difference in means is not equal to 0
R> 95 percent confidence interval:
R>  -1.51699175  0.02670136
R> sample estimates:
R> mean in group A mean in group B
R>      4.001438      4.746584

# if the variances are not equal, simply set `var.equal` to false
# and a Welch's t-test will be performed
```

The first argument we see in `t.test()` is `dat ~ sample`. Usually in R when we see a `~` (tilde) we are creating what is known as a formula. A formula tells R how it should look for interactions between data and factors. For example `Y ~ X` reads:  $Y$  as a function of  $X$ . In our code above we see `dat ~ sample`. This means we are telling R that the  $t$ -test we want it to perform is when the `dat` column is a function of the `sample` column. In plain English we are dividing up the `dat` column into the two different samples we have, and then running a  $t$ -test on these samples. Another way of stating this is that the value of `dat` depends on the grouping it belong to (A or B). We will see this same formula notation cropping up later under ANOVAs, linear models, etc.

Now that we have seen the nitty gritty of how a  $t$ -test is meant to work, [click here<sup>2</sup>](#) to watch a visualisation that demonstrates how the relationships between two different sample sets (based on their mean and variance) influence the results.

**TASK:** Create a visualisation to graphically demonstrate the outcome of this  $t$ -test.

Now we do the same test using a convenient function that comes with the **ggpubr** package, called `compare_means()`, to perform the same  $t$ -test:

```
# perform t-test using `compare_means()`
# note that compare_means() takes the same arguments as t.test()
compare_means(dat ~ sample, data = r_two, method = "t.test", var.equal = TRUE)
```

```
R> # A tibble: 1 x 8
```

```
R> .y. group1 group2      p p.adj p.format p.signif method
R> <chr> <chr> <chr> <dbl> <dbl> <chr> <chr> <chr>
R> 1 dat  A      B      0.0580 0.0580 0.058 ns      T-test
```

Note above that in order to tell `compare_means()` to perform a  $t$ -test we feed it the argument `method = "t.test"`. The output is similar to that of the familiar `t.test()` function that we used earlier, but the output is more abbreviated and less useful. Typically, the output of the  $t$ -tests that we need to report in the results sections of our papers include the  $t$ -statistic, the  $P$ -value, and the degrees of freedom,  $d.f.$ , and these are absent from the `compare_means()` function's output.

### 6.3.1 One-sided two-sample $t$ -tests

Just as with the one-sample  $t$ -tests above, we may also specify which tail of the distribution we are interested in when we compare the means of our two samples. We do so by providing the same arguments as previously:

```
# is the mean of sample B smaller than that of sample A?
compare_means(dat ~ sample, data = r_two, method = "t.test", var.equal = TRUE, alternative = "less")

R> # A tibble: 1 x 8
R> .y. group1 group2      p p.adj p.format p.signif method
R> <chr> <chr> <chr> <dbl> <dbl> <chr> <chr> <chr>
R> 1 dat  A      B      0.971 0.971 0.97 ns      T-test

# is the mean of sample B greater than that of sample A?
compare_means(dat ~ sample, data = r_two, method = "t.test", var.equal = TRUE, alternative = "greater")

R> # A tibble: 1 x 8
R> .y. group1 group2      p p.adj p.format p.signif method
R> <chr> <chr> <chr> <dbl> <dbl> <chr> <chr> <chr>
R> 1 dat  A      B      0.0290 0.0290 0.029 *      T-test
```

What do these results show? Is this surprising?

### 6.3.2 Two-sided two-sample $t$ -tests

Again, as stated above, the default setting in R for comparisons of means is that the test is two-sided. If one wants to state this explicitly it may be done as previously. Note that the results are identical.

```
# default settings
compare_means(dat ~ sample, data = r_two, method = "t.test")

R> # A tibble: 1 x 8
R> .y. group1 group2      p p.adj p.format p.signif method
R> <chr> <chr> <chr> <dbl> <dbl> <chr> <chr> <chr>
R> 1 dat  A      B      0.0584 0.0584 0.058 ns      T-test

# explicitly state a two-sided test
compare_means(dat ~ sample, data = r_two, method = "t.test", alternative = "two.sided")

R> # A tibble: 1 x 8
R> .y. group1 group2      p p.adj p.format p.signif method
R> <chr> <chr> <chr> <dbl> <dbl> <chr> <chr> <chr>
R> 1 dat  A      B      0.0584 0.0584 0.058 ns      T-test
```



## 6.4 Paired $t$ -tests

Paired  $t$ -tests are done when comparing matched samples, and in other words, when our second assumption of  $t$ -tests is violated: the observations are independent of one another—in paired samples, clearly they are not independent. This test is also sometimes called a dependent sample  $t$ -test.

For example, we design a survey to determine if, in a group of 20 people, individuals' right arms differ in length from that of their left arms. For person A, we measure her right arm and her left arm. For person B we measure his right arm and his left arm. So we go all the way to person 20. A right arm belonging with one individual is always matched against a left arm in the same individual. The samples are paired so we use a paired  $t$ -test. Another example: we follow the growth of a sample of 20 piglets over three weeks to see if they weigh more after three weeks than they did at the start of the assessment period. We measure the first piglet, named Halal, at the start of the three week period and again after. We do the same for the second piglet, Kosher. And so it goes. Each piglet has a paired set of measurements, one before matched with one after. In both these examples the data in the two groups (left arm and right arm; or before and after) are not independent, so we need to account for this in the analysis. In practice, how do we perform such a  $t$ -test? Who can think of a dataset we've used in the past that we would use a paired  $t$ -test for?

```
compare_means(dat ~ sample, data = r_two, method = "t.test", paired = TRUE)
```

```
R> # A tibble: 1 x 8
R>   .y. group1 group2      p p.adj p.format p.signif method
R>   <chr> <chr> <chr>   <dbl> <dbl> <chr>   <chr>   <chr>
R> 1 dat   A      B      0.0391 0.0391 0.039   *      T-test
```

## 6.5 Comparison of two population proportions

All of the tests we covered above are designed to deal with continuous data, such as fish lengths or chlorophyll content. If we want to compare proportions (probabilities of success) of different samples against each other, or some known population mean, we need `prop.test()`. Let's create a dummy dataset to get a better idea of how this function works. Below we create some data showing the result of placing two different subjects, Jack and Jill, in separate sealed rooms for two hours (120 minutes). Once every minute a mosquito is let into the room before being extracted again. The columns `yes` and `no` show if the mosquito bit the subject during that one minute. Who says science can't be fun!

```
mosquito <- matrix(c(70, 85, 50, 35), ncol = 2)
colnames(mosquito) <- c("yes", "no")
rownames(mosquito) <- c("Jack", "Jill")
mosquito
```

```
R>      yes no
R> Jack  70 50
R> Jill  85 35
```

### 6.5.1 One-sample and two-sample tests

As with  $t$ -tests, proportion tests may also be based on one sample, or two. If we have only one sample we must specify the total number of trials as well as what the expected population

probability of success is. Because these are individual values, and not matrices, we will show what this would look like without using any objects but will rather give each argument within `prop.test()` a single exact value. In the arguments within `prop.test()`, `x` denotes the number of successes recorded, `n` shows the total number of individual trials performed, and `p` is the expected probability. It is easiest to consider this as though it were a series of 100 coin tosses.

```
# When the probability matches the population
```

```
prop.test(x = 45, n = 100, p = 0.5)
```

```
R>
R> 1-sample proportions test with continuity correction
R>
R> data: 45 out of 100, null probability 0.5
R> X-squared = 0.81, df = 1, p-value = 0.3681
R> alternative hypothesis: true p is not equal to 0.5
R> 95 percent confidence interval:
R> 0.3514281 0.5524574
R> sample estimates:
R> p
R> 0.45
```

```
# When it doesn't
```

```
prop.test(x = 33, n = 100, p = 0.5)
```

```
R>
R> 1-sample proportions test with continuity correction
R>
R> data: 33 out of 100, null probability 0.5
R> X-squared = 10.89, df = 1, p-value = 0.0009668
R> alternative hypothesis: true p is not equal to 0.5
R> 95 percent confidence interval:
R> 0.2411558 0.4320901
R> sample estimates:
R> p
R> 0.33
```

If we have two samples that we would like to compare against one another we enter them into the function as follows:

```
# NB: Note that the `mosquito` data are a matrix, NOT a data.frame
```

```
prop.test(mosquito)
```

```
R>
R> 2-sample test for equality of proportions with continuity
R> correction
R>
R> data: mosquito
R> X-squared = 3.5704, df = 1, p-value = 0.05882
R> alternative hypothesis: two.sided
R> 95 percent confidence interval:
R> -0.253309811 0.003309811
R> sample estimates:
R> prop 1 prop 2
R> 0.5833333 0.7083333
```

Do mosquito's bite Jack and Jill at different proportions?

## 6.5.2 One-sided and two-sided tests

AS with all other tests that compare values, proportion tests may be specified as either one or two-sided. Just to be clear, the default setting for `prop.test()`, like everything else, is a two-sided test. See code below to confirm that the results are identical with or without the added argument:

```
# Default
prop.test(mosquito)

R>
R> 2-sample test for equality of proportions with continuity
R> correction
R>
R> data: mosquito
R> X-squared = 3.5704, df = 1, p-value = 0.05882
R> alternative hypothesis: two.sided
R> 95 percent confidence interval:
R> -0.253309811 0.003309811
R> sample estimates:
R> prop 1 prop 2
R> 0.5833333 0.7083333
```

```
# Explicitly state two-sided test
prop.test(mosquito, alternative = "two.sided")

R>
R> 2-sample test for equality of proportions with continuity
R> correction
R>
R> data: mosquito
R> X-squared = 3.5704, df = 1, p-value = 0.05882
R> alternative hypothesis: two.sided
R> 95 percent confidence interval:
R> -0.253309811 0.003309811
R> sample estimates:
R> prop 1 prop 2
R> 0.5833333 0.7083333
```

Should we want to specify only one of the tails to be considered, we do so precisely the same as with t-tests. Below are examples of what this code would look like:

```
# Jack is bit less than Jill
prop.test(mosquito, alternative = "less")

R>
R> 2-sample test for equality of proportions with continuity
R> correction
R>
R> data: mosquito
R> X-squared = 3.5704, df = 1, p-value = 0.02941
R> alternative hypothesis: less
```