

Introduction to Probability and Statistics Using R

Third Edition

G. Jay Kerns

2017-02-08

Contents

Chapter 1

An Introduction to Probability and Statistics

This chapter has proved to be the hardest to write, by far. The trouble is that there is so much to say – and so many people have already said it so much better than I could. When I get something I like I will release it here.

In the meantime, there is a lot of information already available to a person with an Internet connection. I recommend to start at Wikipedia, which is not a flawless resource but it has the main ideas with links to reputable sources.

In my lectures I usually tell stories about Fisher, Galton, Gauss, Laplace, Quetelet, and the Chevalier de Mere.

1.1 Probability

The common folklore is that probability has been around for millennia but did not gain the attention of mathematicians until approximately 1654 when the Chevalier de Mere had a question regarding the fair division of a game's payoff to the two players, supposing the game had to end prematurely.

1.2 Statistics

Statistics concerns data; their collection, analysis, and interpretation. In this book we distinguish between two types of statistics: descriptive and inferential.

Descriptive statistics concerns the summarization of data. We have a data set and we would like to describe the data set in multiple ways. Usually this entails calculating numbers from the data, called descriptive measures, such as percentages, sums, averages, and so forth.

Inferential statistics does more. There is an inference associated with the data set, a conclusion drawn about the population from which the data originated.

There are two schools of thought regarding statistics: frequentist and bayesian. The difference between the schools is related to how the two groups interpret the underlying probability (see Section ??). The frequentist school gained a lot of ground among statisticians due in large part to the work of Fisher, Neyman, and Pearson in the early twentieth century. That dominance lasted until inexpensive computing power became widely available; nowadays the bayesian school is garnering more attention and at an increasing rate.

This book is devoted mostly to the frequentist viewpoint because that is how I was trained, with the conspicuous exception of Sections ?? and ??.

Chapter 2

Data Description

In this chapter we introduce the different types of data that a statistician is likely to encounter, and in each subsection we give some examples of how to display the data of that particular type. Once we see how to display data distributions, we next introduce the basic properties of data distributions. We qualitatively explore several data sets. Once that we have intuitive properties of data sets, we next discuss how we may numerically measure and describe those properties with descriptive statistics.

What do I want them to know?

- different data types, such as quantitative versus qualitative, nominal versus ordinal, and discrete versus continuous
- basic graphical displays for assorted data types, and some of their (dis)advantages
- fundamental properties of data distributions, including center, spread, shape, and crazy observations
- methods to describe data (visually/numerically) with respect to the properties, and how the methods differ depending on the data type
- all of the above in the context of grouped data, and in particular, the concept of a factor

2.1 Types of Data

Loosely speaking, a datum is any piece of collected information, and a data set is a collection of data related to each other in some way. We will categorize data into five types and describe each in turn:

- **Quantitative:** data associated with a measurement of some quantity on an observational unit,

- **Qualitative:** data associated with some quality or property of an observational unit,
- **Logical:** data which represent true or false and play an important role later,
- **Missing:** data which should be there but are not, and
- **Other types:** everything else under the sun.

In each subsection we look at some examples of the type in question and introduce methods to display them.

2.1.1 Quantitative data

Quantitative data are any data that measure or are associated with a measurement of the quantity of something. They invariably assume numerical values. Quantitative data can be further subdivided into two categories.

- *Discrete data* take values in a finite or countably infinite set of numbers, that is, all possible values could (at least in principle) be written down in an ordered list. Examples include: counts, number of arrivals, or number of successes. They are often represented by integers, say, 0, 1, 2, *etc.*
- *Continuous data* take values in an interval of numbers. These are also known as scale data, interval data, or measurement data. Examples include: height, weight, length, time, *etc.* Continuous data are often characterized by fractions or decimals: 3.82, 7.0001, $4\frac{5}{8}$, *etc.*

Note that the distinction between discrete and continuous data is not always clear-cut. Sometimes it is convenient to treat data as if they were continuous, even though strictly speaking they are not continuous. See the examples.

Example 2.1 (Annual Precipitation in US Cities). The vector `precip` contains average amount of rainfall (in inches) for each of 70 cities in the United States and Puerto Rico. Let us take a look at the data:

```
str(precip)
```

```
Named num [1:70] 67 54.7 7 48.5 14 17.2 20.7 13 43.4 40.2 ...
```

```
- attr(*, "names")= chr [1:70] "Mobile" "Juneau" "Phoenix" "Little Rock" ..
```

```
precip[1:4]
```

Mobile	Juneau	Phoenix	Little Rock
67.0	54.7	7.0	48.5

The output shows that `precip` is a numeric vector which has been *named*, that is, each value has a name associated with it (which can be set with the `names` function). These are quantitative continuous data.

Example 2.2 (Lengths of Major North American Rivers). The U.S. Geological Survey recorded the lengths (in miles) of several rivers in North America. They are stored in the vector `rivers` in the `datasets` package [`datasets`] (which ships with base R). See `?rivers`. Let us take a look at the data with the `str` function.

```
str(rivers)
```

```
num [1:141] 735 320 325 392 524 ...
```

The output says that `rivers` is a numeric vector of length 141, and the first few values are 735, 320, 325, *etc.* These data are definitely quantitative and it appears that the measurements have been rounded to the nearest mile. Thus, strictly speaking, these are discrete data. But we will find it convenient later to take data like these to be continuous for some of our statistical procedures.

Example 2.3 (Yearly Numbers of Important Discoveries). The vector `discoveries` contains numbers of “great” inventions/discoveries in each year from 1860 to 1959, as reported by the 1975 World Almanac. Let us take a look at the data:

```
str(discoveries)
```

```
Time-Series [1:100] from 1860 to 1959: 5 3 0 2 0 3 2 3 6 1 ...
```

The output is telling us that `discoveries` is a *time series* (see Section ?? for more) of length 100. The entries are integers, and since they represent counts this is a good example of discrete quantitative data. We will take a closer look in the following sections.

2.1.2 Displaying Quantitative Data

One of the first things to do when confronted by quantitative data (or any data, for that matter) is to make some sort of visual display to gain some insight into the data’s structure. There are almost as many display types from which to choose as there are data sets to plot. We describe some of the more popular alternatives.

Strip charts also known as Dot plots

These can be used for discrete or continuous data, and usually look best when the data set is not too large. Along the horizontal axis is a numerical scale above which the data values are plotted. We can do it in R with a call to the `stripchart` function. There are three available methods.

- `overplot`: plots ties covering each other. This method is good to display only the distinct values assumed by the data set.
- `jitter`: adds some noise to the data in the *y* direction in which case the data values are not covered up by ties.

- **stack**: plots repeated values stacked on top of one another. This method is best used for discrete data with a lot of ties; if there are no repeats then this method is identical to overplot.

See Figure ??, which was produced by the following code.

```
stripchart(precip, xlab="rainfall")
stripchart(rivers, method="jitter", xlab="length")
stripchart(discoveries, method="stack", xlab="number")
```

The leftmost graph is a strip chart of the `precip` data. The graph shows tightly clustered values in the middle with some others falling balanced on either side, with perhaps slightly more falling to the left. Later we will call this a symmetric distribution, see Section ?. The middle graph is of the `rivers` data, a vector of length 141. There are several repeated values in the rivers data, and if we were to use the `overplot` method we would lose some of them in the display. This plot shows a what we will later call a right-skewed shape with perhaps some extreme values on the far right of the display. The third graph strip charts `discoveries` data which are literally a textbook example of a right skewed distribution.

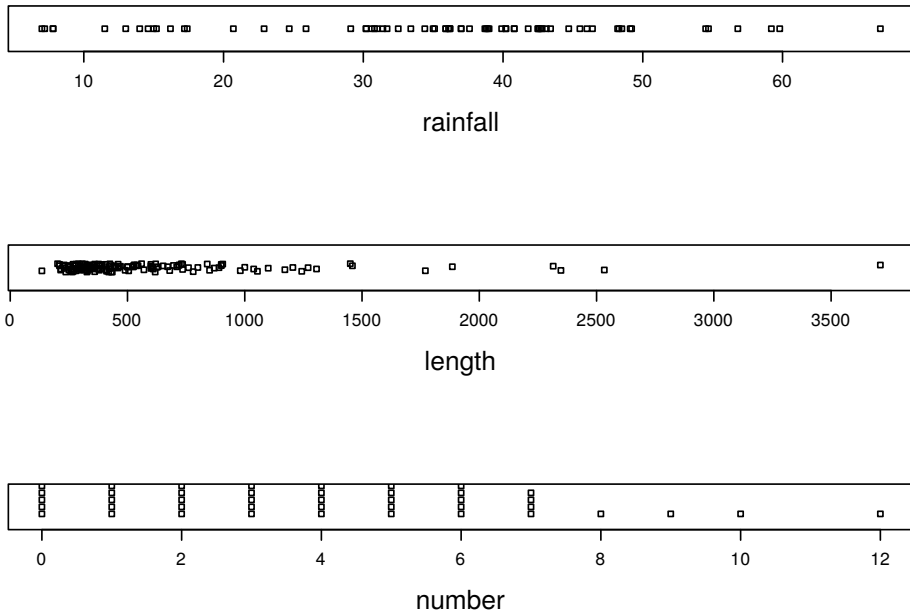


Figure 2.1: Three stripcharts of three data sets. The first graph uses the `overplot` method, the second the `jitter` method, and the third the `stack` method.

The `DOTplot` function in the `UsingR` package [[UsingR](#)] is another alternative.

Histogram

These are typically used for continuous data. A histogram is constructed by first deciding on a set of classes, or bins, which partition the real line into a set of boxes into which the data values fall. Then vertical bars are drawn over the bins with height proportional to the number of observations that fell into the bin.

These are one of the most common summary displays, and they are often misidentified as “Bar Graphs” (see below.) The scale on the y axis can be frequency, percentage, or density (relative frequency). The term histogram was coined by Karl Pearson in 1891, see [Miller].

Example 2.4 (Annual Precipitation in US Cities). We are going to take another look at the `precip` data that we investigated earlier. The strip chart in Figure ?? suggested a loosely balanced distribution; let us now look to see what a histogram says.

There are many ways to plot histograms in R, and one of the easiest is with the `hist` function. The following code produces the plots in Figure ??.

```
hist(precip, main = "")  
hist(precip, freq = FALSE, main = "")
```

Notice the argument `main = ""` which suppresses the main title from being displayed – it would have said “Histogram of `precip`” otherwise. The plot on the left is a frequency histogram (the default), and the plot on the right is a relative frequency histogram (`freq = FALSE`).

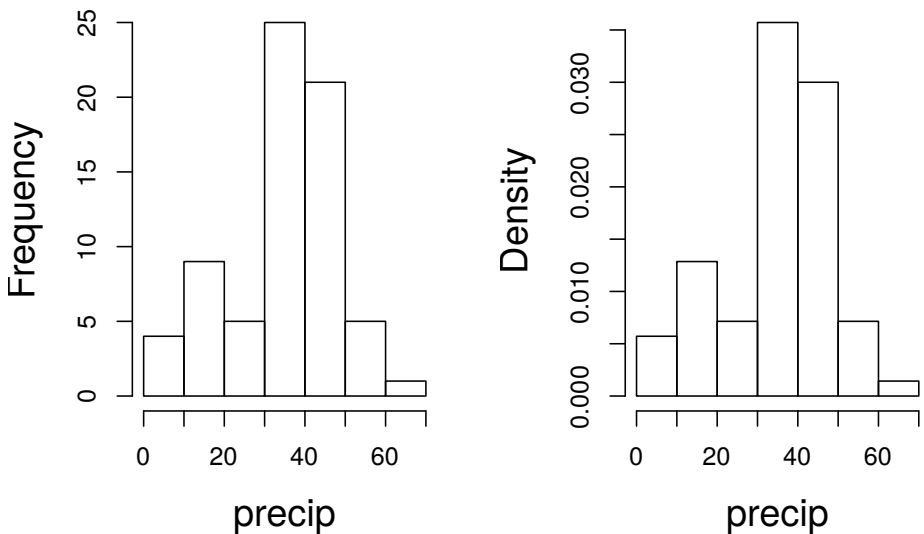


Figure 2.2: (Relative) frequency histograms of the `precip` data.

Please mind the biggest weakness of histograms: the graph obtained strongly depends

on the bins chosen. Choose another set of bins, and you will get a different histogram. Moreover, there are not any definitive criteria by which bins should be defined; the best choice for a given data set is the one which illuminates the data set's underlying structure (if any). Luckily for us there are algorithms to automatically choose bins that are likely to display well, and more often than not the default bins do a good job. This is not always the case, however, and a responsible statistician will investigate many bin choices to test the stability of the display.

Recall that the strip chart in Figure ?? suggested a relatively balanced shape to the precip data distribution. Watch what happens when we change the bins slightly (with the breaks argument to hist). See Figure ?? which was produced by the following code.

```
hist(precip, breaks = 10)
hist(precip, breaks = 25)
hist(precip, breaks = 50)
```

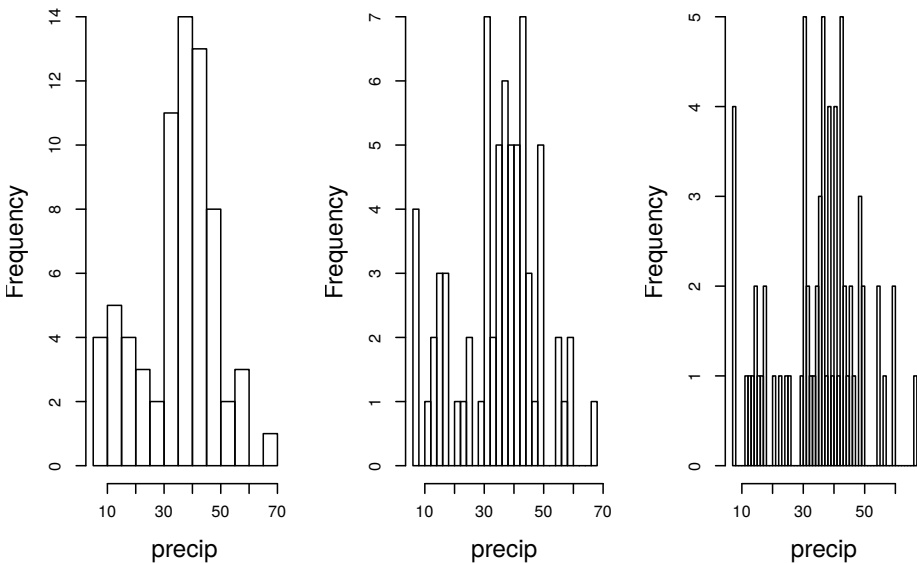


Figure 2.3: More histograms of the precip data.

The leftmost graph (with `breaks = 10`) shows that the distribution is not balanced at all. There are two humps: a big one in the middle and a smaller one to the left. Graphs like this often indicate some underlying group structure to the data; we could now investigate whether the cities for which rainfall was measured were similar in some way, with respect to geographic region, for example.

The rightmost graph in Figure ?? shows what happens when the number of bins is too large: the histogram is too grainy and hides the rounded appearance of the earlier histograms. If we were to continue increasing the number of bins we would eventually

get all observed bins to have exactly one element, which is nothing more than a glorified strip chart.

Stem-and-leaf displays (more to be said in Section ??)

Stem-and-leaf displays (also known as stemplots) have two basic parts: *stems* and *leaves*. The final digit of the data values is taken to be a *leaf*, and the leading digit(s) is (are) taken to be *stems*. We draw a vertical line, and to the left of the line we list the stems. To the right of the line, we list the leaves beside their corresponding stem. There will typically be several leaves for each stem, in which case the leaves accumulate to the right. It is sometimes necessary to round the data values, especially for larger data sets.

Example 2.5 (Driver Deaths in the United Kingdom). `UKDriverDeaths` is a time series that contains the total car drivers killed or seriously injured in Great Britain monthly from Jan 1969 to Dec 1984. See `?UKDriverDeaths`. Compulsory seat belt use was introduced on January 31, 1983. We construct a stem and leaf diagram in R with the `stem.leaf` function from the `aplpack` package [`aplpack`].

```
stem.leaf(UKDriverDeaths, depth = FALSE)
```

```
1 | 2: represents 120
  leaf unit: 10
      n: 192
10 | 57
11 | 136678
12 | 123889
13 | 0255666888899
14 | 00001222344444555556667788889
15 | 000011111222222344445555566677779
16 | 01222333444444555555678888889
17 | 11233344566667799
18 | 00011235568
19 | 01234455667799
20 | 0000113557788899
21 | 145599
22 | 013467
23 | 9
24 | 7
HI: 2654
```

The display shows a more or less balanced mound-shaped distribution, with one or maybe two humps, a big one and a smaller one just to its right. Note that the data have been rounded to the tens place so that each datum gets only one leaf to the right of the dividing line.

Notice that the depths have been suppressed. To learn more about this option and many others, see Section ???. Unlike a histogram, the original data values may be recovered from the stem-and-leaf display – modulo the rounding – that is, starting from the top and working down we can read off the data values 1050, 1070, 1110, 1130, and so forth.

Index plots

Done with the `plot` function. These are good for plotting data which are ordered, for example, when the data are measured over time. That is, the first observation was measured at time 1, the second at time 2, *etc.* It is a two dimensional plot, in which the index (or time) is the x variable and the measured value is the y variable. There are several plotting methods for index plots, and we mention two of them:

- **spikes:** draws a vertical line from the x -axis to the observation height.
- **points:** plots a simple point at the observation height.

Example 2.6 (Level of Lake Huron 1875-1972). Brockwell and Davis [Brockwell1991] give the annual measurements of the level (in feet) of Lake Huron from 1875–1972. The data are stored in the time series `LakeHuron`. See `?LakeHuron`. Figure ?? was produced with the following code:

```
plot(LakeHuron)
plot(LakeHuron, type = "p")
plot(LakeHuron, type = "h")
```

The plots show an overall decreasing trend to the observations, and there appears to be some seasonal variation that increases over time.

Density estimates

The default method uses a Gaussian kernel density estimate.

```
# The Old Faithful geyser data
d <- density(faithful$eruptions, bw = "sj")
d
plot(d)
hist(precip, freq = FALSE)
lines(density(precip))
```

2.1.3 Qualitative Data, Categorical Data, and Factors

Qualitative data are simply any type of data that are not numerical, or do not represent numerical quantities. Examples of qualitative variables include a subject's name, gender,

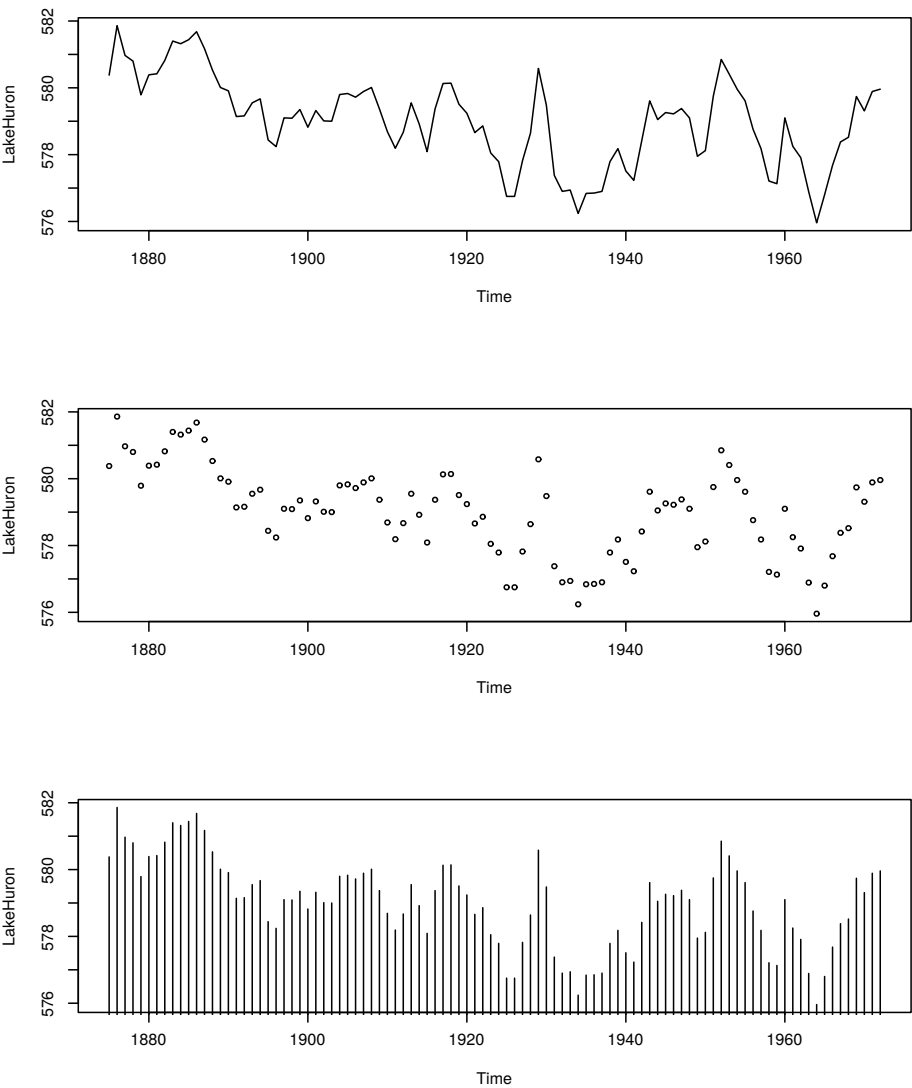


Figure 2.4: Index plots of the LakeHuron data.

race/ethnicity, political party, socioeconomic status, class rank, driver's license number, and social security number (SSN).

Please bear in mind that some data *look* to be quantitative but are *not*, because they do not represent numerical quantities and do not obey mathematical rules. For example, a person's shoe size is typically written with numbers: 8, or 9, or 12, or $12\frac{1}{2}$. Shoe size is not quantitative, however, because if we take a size 8 and combine with a size 9 we do not get a size 17.

Some qualitative data serve merely to *identify* the observation (such a subject's name, driver's license number, or SSN). This type of data does not usually play much of a role in statistics. But other qualitative variables serve to *subdivide* the data set into categories; we call these *factors*. In the above examples, gender, race, political party, and socioeconomic status would be considered factors (shoe size would be another one). The possible values of a factor are called its *levels*. For instance, the factor of gender would have two levels, namely, male and female. Socioeconomic status typically has three levels: high, middle, and low.

Factors may be of two types: *nominal* and *ordinal*. Nominal factors have levels that correspond to names of the categories, with no implied ordering. Examples of nominal factors would be hair color, gender, race, or political party. There is no natural ordering to "Democrat" and "Republican"; the categories are just names associated with different groups of people.

In contrast, ordinal factors have some sort of ordered structure to the underlying factor levels. For instance, socioeconomic status would be an ordinal categorical variable because the levels correspond to ranks associated with income, education, and occupation. Another example of ordinal categorical data would be class rank.

Factors have special status in R. They are represented internally by numbers, but even when they are written numerically their values do not convey any numeric meaning or obey any mathematical rules (that is, Stage III cancer is not Stage I cancer + Stage II cancer).

Example 2.7. The `state.abb` vector gives the two letter postal abbreviations for all 50 states.

```
str(state.abb)
```

```
chr [1:50] "AL" "AK" "AZ" "AR" "CA" "CO" ...
```

These would be ID data. The `state.name` vector lists all of the complete names and those data would also be ID.

Example 2.8 (U.S. State Facts and Features). The U.S. Department of Commerce of the U.S. Census Bureau releases all sorts of information in the *Statistical Abstract of the United States*, and the `state.region` data lists each of the 50 states and the region to which it belongs, be it Northeast, South, North Central, or West. See `?state.region`.

```
str(state.region)
```

```
Factor w/ 4 levels "Northeast","South",...: 2 4 4 2 4 4 1 2 2 2 ...
state.region[1:5]
```

```
[1] South West West South West
Levels: Northeast South North Central West
```

The `str` output shows that `state.region` is already stored internally as a factor and it lists a couple of the factor levels. To see all of the levels we printed the first five entries of the vector in the second line.

2.1.4 Displaying Qualitative Data

Tables

One of the best ways to summarize qualitative data is with a table of the data values. We may count frequencies with the `table` function or list proportions with the `prop.table` function (whose input is a frequency table). In the R Commander you can do it with **Statistics > Frequency Distribution...** Alternatively, to look at tables for all factors in the Active data set you can do **Statistics > Summaries > Active Dataset**.

```
Tbl <- table(state.division)
Tbl
```

```
state.division
      New England      Middle Atlantic      South Atlantic
              6              3              8
East South Central West South Central East North Central
              4              4              5
West North Central      Mountain      Pacific
              7              8              5
```

```
Tbl/sum(Tbl)      # relative frequencies
```

```
state.division
      New England      Middle Atlantic      South Atlantic
      0.12          0.06          0.16
East South Central West South Central East North Central
      0.08          0.08          0.10
West North Central      Mountain      Pacific
      0.14          0.16          0.10
```

```
prop.table(Tbl)      # same thing
```

state.division					
New England		Middle Atlantic		South Atlantic	
0.12		0.06		0.16	
East South Central		West South Central		East North Central	
0.08		0.08		0.10	
West North Central		Mountain		Pacific	
0.14		0.16		0.10	

Bar Graphs

A bar graph is the analogue of a histogram for categorical data. A bar is displayed for each level of a factor, with the heights of the bars proportional to the frequencies of observations falling in the respective categories. A disadvantage of bar graphs is that the levels are ordered alphabetically (by default), which may sometimes obscure patterns in the display.

Example 2.9 (U.S. State Facts and Features). The `state.region` data lists each of the 50 states and the region to which it belongs, be it Northeast, South, North Central, or West. See `?state.region`. It is already stored internally as a factor. We make a bar graph with the `barplot` function:

```
barplot(table(state.region), cex.names = 1.20)
barplot(prop.table(table(state.region)), cex.names = 1.20)
```

See Figure ?? . The display on the left is a frequency bar graph because the `y` axis shows counts, while the display on the left is a relative frequency bar graph. The only difference between the two is the scale. Looking at the graph we see that the majority of the fifty states are in the South, followed by West, North Central, and finally Northeast. Over 30% of the states are in the South.

Notice the `cex.names` argument that we used, above. It expands the names on the `x` axis by 20% which makes them easier to read. See `?par` for a detailed list of additional plot parameters.

Pareto Diagrams

A pareto diagram is a lot like a bar graph except the bars are rearranged such that they decrease in height going from left to right. The rearrangement is handy because it can visually reveal structure (if any) in how fast the bars decrease – this is much more difficult when the bars are jumbled.

Example 2.10 (U.S. State Facts and Features). The `state.division` data record the division (New England, Middle Atlantic, South Atlantic, East South Central, West South Central, East North Central, West North Central, Mountain, and Pacific) of the fifty states. We can make a pareto diagram with either the `RcmdrPlugin.IPSUR`

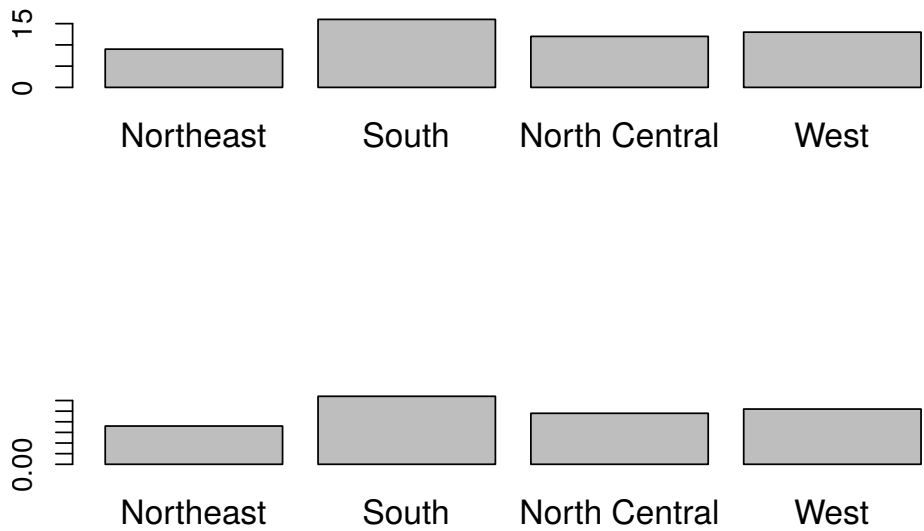


Figure 2.5: The top graph is a frequency barplot made with `table` and the bottom is a relative frequency barplot made with `prop.table`.

package `[RcmdrPlugin.IPSUR]` or with the `pareto.chart` function from the `qcc` package `[qcc]`. See Figure ???. The code follows.

```
pareto.chart(table(state.division), ylab="Frequency")
```

Pareto chart analysis for `table(state.division)`

	Frequency	Cum.Freq.	Percentage
South Atlantic	8	8	16
Mountain	8	16	16
West North Central	7	23	14
New England	6	29	12
East North Central	5	34	10
Pacific	5	39	10
East South Central	4	43	8
West South Central	4	47	8
Middle Atlantic	3	50	6

Pareto chart analysis for `table(state.division)`

	Cum.Percent.
South Atlantic	16
Mountain	32
West North Central	46
New England	58
East North Central	68
Pacific	78

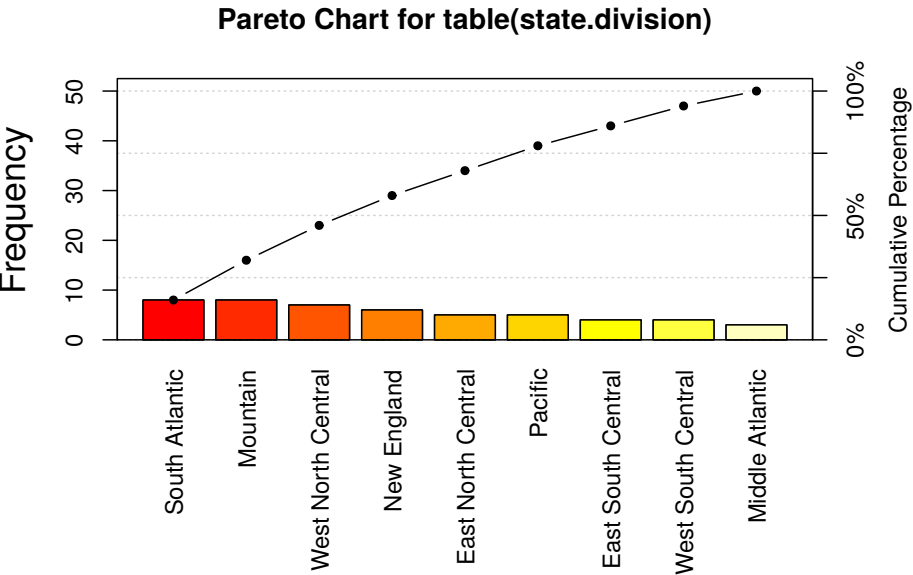


Figure 2.6: Pareto chart of the state.division data.

East South Central	86
West South Central	94
Middle Atlantic	100

Dot Charts

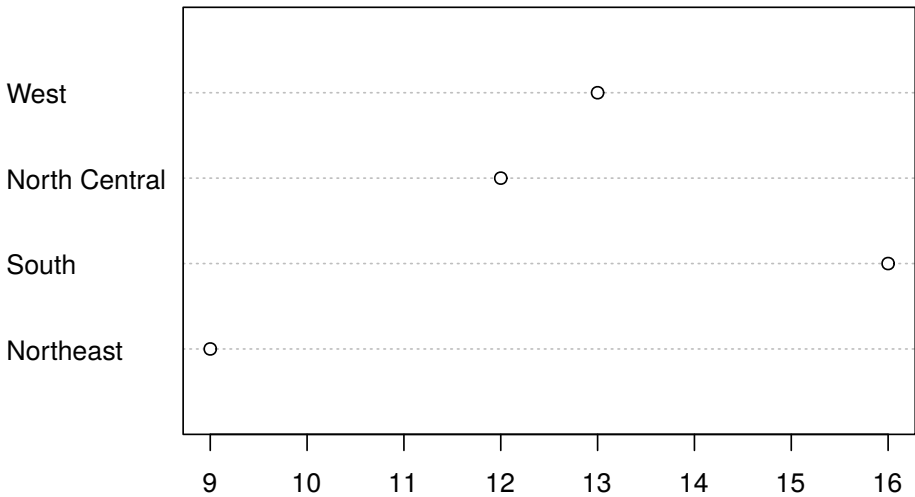
These are a lot like a bar graph that has been turned on its side with the bars replaced by dots on horizontal lines. They do not convey any more (or less) information than the associated bar graph, but the strength lies in the economy of the display. Dot charts are so compact that it is easy to graph very complicated multi-variable interactions together in one graph. See Section ?? . We will give an example here using the same data as above for comparison. The graph was produced by the following code.

```
x <- table(state.region)
dotchart(as.vector(x), labels = names(x), cex.lab = cexlab)
```

See Figure ?? . Compare it to Figure ?? .

Pie Graphs

These can be done with R and the R Commander, but they fallen out of favor in recent years because researchers have determined that while the human eye is good at judging linear measures, it is notoriously bad at judging relative areas (such as those displayed

Figure 2.7: Dot chart of the `state.region` data.

by a pie graph). Pie charts are consequently a very bad way of displaying information unless the number of categories is two or three. A bar chart or dot chart is a preferable way of displaying qualitative data. See `?pie` for more information.

We are not going to do any examples of a pie graph and discourage their use elsewhere.

2.1.5 Logical Data

There is another type of information recognized by R which does not fall into the above categories. The value is either `TRUE` or `FALSE` (note that equivalently you can use `1 = TRUE`, `0 = FALSE`). Here is an example of a logical vector:

```
x <- 5:9
y <- (x < 7.3)
y
```

```
[1] TRUE TRUE TRUE FALSE FALSE
```

Many functions in R have options that the user may or may not want to activate in the function call. For example, the `stem.leaf` function has the `depths` argument which is `TRUE` by default. We saw in Section ?? how to turn the option off, simply enter `stem.leaf(x, depths = FALSE)` and they will not be shown on the display.

We can swap `TRUE` with `FALSE` with the exclamation point `!`.

```
!y
```

```
[1] FALSE FALSE FALSE TRUE TRUE
```

2.1.6 Missing Data

Missing data are a persistent and prevalent problem in many statistical analyses, especially those associated with the social sciences. R reserves the special symbol `NA` to representing missing data.

Ordinary arithmetic with `NA` values give `NA`'s (addition, subtraction, *etc.*) and applying a function to a vector that has an `NA` in it will usually give an `NA`.

```
x <- c(3, 7, NA, 4, 7)
y <- c(5, NA, 1, 2, 2)
x + y
```

```
[1] 8 NA NA 6 9
```

Some functions have a `na.rm` argument which when `TRUE` will ignore missing data as if they were not there (such as `mean`, `var`, `sd`, `IQR`, `mad`, ...).

```
sum(x)
```

```
[1] NA
```

```
sum(x, na.rm = TRUE)
```

```
[1] 21
```

Other functions do not have a `na.rm` argument and will return `NA` or an error if the argument has `NA`s. In those cases we can find the locations of any `NA`s with the `is.na` function and remove those cases with the `[]` operator.

```
is.na(x)
```

```
[1] FALSE FALSE TRUE FALSE FALSE
```

```
z <- x[!is.na(x)]
sum(z)
```

```
[1] 21
```

The analogue of `is.na` for rectangular data sets (or data frames) is the `complete.cases` function. See Appendix ??.

2.1.7 Other Data Types

2.2 Features of Data Distributions

Given that the data have been appropriately displayed, the next step is to try to identify salient features represented in the graph. The acronym to remember is *C*-enter, *U*-nusual features, *S*-pread, and *S*-hape. (CUSS).

2.2.1 Center

One of the most basic features of a data set is its center. Loosely speaking, the center of a data set is associated with a number that represents a middle or general tendency of the data. Of course, there are usually several values that would serve as a center, and our later tasks will be focused on choosing an appropriate one for the data at hand. Judging from the histogram that we saw in Figure ??, a measure of center would be about 35.

2.2.2 Spread

The spread of a data set is associated with its variability; data sets with a large spread tend to cover a large interval of values, while data sets with small spread tend to cluster tightly around a central value.

2.2.3 Shape

When we speak of the *shape* of a data set, we are usually referring to the shape exhibited by an associated graphical display, such as a histogram. The shape can tell us a lot about any underlying structure to the data, and can help us decide which statistical procedure we should use to analyze them.

Symmetry and Skewness

A distribution is said to be *right-skewed* (or *positively skewed*) if the right tail seems to be stretched from the center. A *left-skewed* (or *negatively skewed*) distribution is stretched to the left side. A symmetric distribution has a graph that is balanced about its center, in the sense that half of the graph may be reflected about a central line of symmetry to match the other half.

We have already encountered skewed distributions: both the discoveries data in Figure ?? and the precip data in Figure ?? appear right-skewed. The UKDriverDeaths data in Example ?? is relatively symmetric (but note the one extreme value 2654 identified at the bottom of the stem-and-leaf display).

Kurtosis

Another component to the shape of a distribution is how “peaked” it is. Some distributions tend to have a flat shape with thin tails. These are called *platykurtic*, and an example of a platykurtic distribution is the uniform distribution; see Section ?. On the other end of the spectrum are distributions with a steep peak, or spike, accompanied by heavy tails; these are called *leptokurtic*. Examples of leptokurtic distributions are

2.2.5 Extreme Observations and other Unusual Features

Extreme observations fall far from the rest of the data. Such observations are troublesome to many statistical procedures; they cause exaggerated estimates and instability. It is important to identify extreme observations and examine the source of the data more closely. There are many possible reasons underlying an extreme observation:

- *Maybe the value is a typographical error.* Especially with large data sets becoming more prevalent, many of which being recorded by hand, mistakes are a common problem. After closer scrutiny, these can often be fixed.
- *Maybe the observation was not meant for the study,* because it does not belong to the population of interest. For example, in medical research some subjects may have relevant complications in their genealogical history that would rule out their participation in the experiment. Or when a manufacturing company investigates the properties of one of its devices, perhaps a particular product is malfunctioning and is not representative of the majority of the items.
- *Maybe it indicates a deeper trend or phenomenon.* Many of the most influential scientific discoveries were made when the investigator noticed an unexpected result, a value that was not predicted by the classical theory. Albert Einstein, Louis Pasteur, and others built their careers on exactly this circumstance.

2.3 Descriptive Statistics

One of my favorite professors would repeatedly harp, “You cannot do statistics without data.”

What do I want them to know?

- The fundamental data types we encounter most often, how to classify given data into a likely type, and that sometimes the distinction is blurry.

2.3.1 Frequencies and Relative Frequencies

These are used for categorical data. The idea is that there are a number of different categories, and we would like to get some idea about how the categories are represented in the population.

2.3.2 Measures of Center

The *sample mean* is denoted \bar{x} (read “x-bar”) and is simply the arithmetic average of the observations:

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i. \quad (2.1)$$

- Good: natural, easy to compute, has nice mathematical properties
- Bad: sensitive to extreme values

It is appropriate for use with data sets that are not highly skewed without extreme observations.

The *sample median* is another popular measure of center and is denoted \tilde{x} . To calculate its value, first sort the data into an increasing sequence of numbers. If the data set has an odd number of observations then \tilde{x} is the value of the middle observation, which lies in position $(n + 1)/2$; otherwise, there are two middle observations and \tilde{x} is the average of those middle values.

- Good: resistant to extreme values, easy to describe
- Bad: not as mathematically tractable, need to sort the data to calculate

One desirable property of the sample median is that it is *resistant* to extreme observations, in the sense that the value of \tilde{x} depends only on those data values in the middle, and is quite unaffected by the actual values of the outer observations in the ordered list. The same cannot be said for the sample mean. Any significant changes in the magnitude of an observation x_k results in a corresponding change in the value of the mean. Consequently, the sample mean is said to be *sensitive* to extreme observations.

The *trimmed mean* is a measure designed to address the sensitivity of the sample mean to extreme observations. The idea is to “trim” a fraction (less than 1/2) of the observations off each end of the ordered list, and then calculate the sample mean of what remains. We will denote it by $\bar{x}_{t=0.05}$.

- Good: resistant to extreme values, shares nice statistical properties
- Bad: need to sort the data

How to do it with R

- You can calculate frequencies or relative frequencies with the `table` function, and relative frequencies with `prop.table(table())`.
- You can calculate the sample mean of a data vector `x` with the command `mean(x)`.
- You can calculate the sample median of `x` with the command `median(x)`.
- You can calculate the trimmed mean with the `trim` argument; `mean(x, trim = 0.05)`.

2.3.3 Order Statistics and the Sample Quantiles

A common first step in an analysis of a data set is to sort the values. Given a data set x_1, x_2, \dots, x_n , we may sort the values to obtain an increasing sequence

$$x_{(1)} \leq x_{(2)} \leq x_{(3)} \leq \cdots \leq x_{(n)} \quad (2.2)$$

and the resulting values are called the *order statistics*. The k^{th} entry in the list, $x_{(k)}$, is the k^{th} order statistic, and approximately $100(k/n)\%$ of the observations fall below $x_{(k)}$. The order statistics give an indication of the shape of the data distribution, in the sense that a person can look at the order statistics and have an idea about where the data are concentrated, and where they are sparse.

The *sample quantiles* are related to the order statistics. Unfortunately, there is not a universally accepted definition of them. Indeed, R is equipped to calculate quantiles using nine distinct definitions! We will describe the default method (`type = 7`), but the interested reader can see the details for the other methods with `?quantile`.

Suppose the data set has n observations. Find the sample quantile of order p ($0 < p < 1$), denoted \tilde{q}_p , as follows:

- **First step:** sort the data to obtain the order statistics $x_{(1)}, x_{(2)}, \dots, x_{(n)}$.
- **Second step:** calculate $(n - 1)p + 1$ and write it in the form $k.d$, where k is an integer and d is a decimal.
- **Third step:** The sample quantile \tilde{q}_p is

$$\tilde{q}_p = x_{(k)} + d(x_{(k+1)} - x_{(k)}). \quad (2.3)$$

The interpretation of \tilde{q}_p is that approximately $100p\%$ of the data fall below the value \tilde{q}_p .

Keep in mind that there is not a unique definition of percentiles, quartiles, *etc.* Open a different book, and you'll find a different definition. The difference is small and seldom plays a role except in small data sets with repeated values. In fact, most people do not even notice in common use.

Clearly, the most popular sample quantile is $\tilde{q}_{0.50}$, also known as the sample median, \tilde{x} . The closest runners-up are the *first quartile* $\tilde{q}_{0.25}$ and the *third quartile* $\tilde{q}_{0.75}$ (the *second quartile* is the median).

How to do it with R

At the command prompt We can find the order statistics of a data set stored in a vector `x` with the command `sort(x)`.

We can calculate the sample quantiles of any order p where $0 < p < 1$ for a data set stored in a data vector `x` with the `quantile` function, for instance, the command `quantile(x, probs = c(0, 0.25, 0.37))` will return the smallest observation, the first quartile, $\tilde{q}_{0.25}$, and the 37th sample quantile, $\tilde{q}_{0.37}$. For \tilde{q}_p simply change the values in the `probs` argument to the value p .

With the R Commander we can find the order statistics of a variable in the Active data set by doing Data ► Manage variables in Active data set... ► Compute new variable... In the Expression to compute dialog simply type `sort(varname)`, where `varname` is the variable that it is desired to sort.

In Rcmdr, we can calculate the sample quantiles for a particular variable with the sequence Statistics ► Summaries ► Numerical Summaries... We can automatically calculate the quartiles for all variables in the Active data set with the sequence Statistics ► Summaries ► Active Dataset.

2.3.4 Measures of Spread

Sample Variance and Standard Deviation

The *sample variance* is denoted s^2 and is calculated with the formula

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (2.4)$$

The *sample standard deviation* is $s = \sqrt{s^2}$. Intuitively, the sample variance is approximately the average squared distance of the observations from the sample mean. The sample standard deviation is used to scale the estimate back to the measurement units of the original data.

- Good: tractable, has nice mathematical/statistical properties
- Bad: sensitive to extreme values

We will spend a lot of time with the variance and standard deviation in the coming chapters. In the meantime, the following two rules give some meaning to the standard deviation, in that there are bounds on how much of the data can fall past a certain distance from the mean.

Fact: (Chebychev's Rule). The proportion of observations within k standard deviations of the mean is at least $1 - 1/k^2$, i.e., at least 75%, 89%, and 94% of the data are within 2, 3, and 4 standard deviations of the mean, respectively.

Note that Chebychev's Rule does not say anything about when $k = 1$, because $1 - 1/1^2 = 0$, which states that at least 0% of the observations are within one standard deviation of the mean (which is not saying much).

Chebychev's Rule applies to any data distribution, any list of numbers, no matter where it came from or what the histogram looks like. The price for such generality is that the bounds are not very tight; if we know more about how the data are shaped then we can say more about how much of the data can fall a given distance from the mean.

Fact: (Empirical Rule). If data follow a bell-shaped curve, then approximately 68%, 95%, and 99.7% of the data are within 1, 2, and 3 standard deviations of the mean, respectively.

Interquartile Range

Just as the sample mean is sensitive to extreme values, so the associated measure of spread is similarly sensitive to extremes. Further, the problem is exacerbated by the fact that the extreme distances are squared. We know that the sample quartiles are resistant to extremes, and a measure of spread associated with them is the *interquartile range* (*IQR*) defined by $IQR = q_{0.75} - q_{0.25}$.

- Good: stable, resistant to outliers, robust to nonnormality, easy to explain
- Bad: not as tractable, need to sort the data, only involves the middle 50% of the data.

Median Absolute Deviation

A measure even more robust than the *IQR* is the *median absolute deviation* (*MAD*). To calculate it we first get the median \tilde{x} , next the *absolute deviations* $|x_1 - \tilde{x}|, |x_2 - \tilde{x}|, \dots, |x_n - \tilde{x}|$, and the *MAD* is proportional to the median of those deviations:

$$MAD \propto \text{median}(|x_1 - \tilde{x}|, |x_2 - \tilde{x}|, \dots, |x_n - \tilde{x}|). \quad (2.5)$$

That is, the $MAD = c \cdot \text{median}(|x_1 - \tilde{x}|, |x_2 - \tilde{x}|, \dots, |x_n - \tilde{x}|)$, where c is a constant chosen so that the *MAD* has nice properties. The value of c in R is by default $c = 1.4286$. This value is chosen to ensure that the estimator of σ is correct, on the average, under suitable sampling assumptions (see Section ??).

- Good: stable, very robust, even more so than the *IQR*.
- Bad: not tractable, not well known and less easy to explain.

Comparing Apples to Apples

We have seen three different measures of spread which, for a given data set, will give three different answers. Which one should we use? It depends on the data set. If the data are well behaved, with an approximate bell-shaped distribution, then the sample mean and sample standard deviation are natural choices with nice mathematical properties. However, if the data have an unusual or skewed shape with several extreme values, perhaps the more resistant choices among the *IQR* or *MAD* would be more appropriate.

However, once we are looking at the three numbers it is important to understand that the estimators are not all measuring the same quantity, on the average. In particular, it can be shown that when the data follow an approximately bell-shaped distribution, then

on the average, the sample standard deviation s and the *MAD* will be the approximately the same value, namely, σ , but the *IQR* will be on the average 1.349 times larger than s and the *MAD*. See ?? for more details.

How to do it with R

At the command prompt we may compute the sample range with `range(x)` and the sample variance with `var(x)`, where x is a numeric vector. The sample standard deviation is `sqrt(var(x))` or just `sd(x)`. The *IQR* is `IQR(x)` and the median absolute deviation is `mad(x)`.

With the R Commander we can calculate the sample standard deviation with the **Statistics** \triangleright **Summaries** \triangleright **Numerical Summaries...** combination. R Commander does not calculate the *IQR* or *MAD* in any of the menu selections, by default.

2.3.5 Measures of Shape

Sample Skewness

The *sample skewness*, denoted by g_1 , is defined by the formula

$$g_1 = \frac{1}{n} \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{s^3}. \quad (2.6)$$

The sample skewness can be any value $-\infty < g_1 < \infty$. The sign of g_1 indicates the direction of skewness of the distribution. Samples that have $g_1 > 0$ indicate right-skewed distributions (or positively skewed), and samples with $g_1 < 0$ indicate left-skewed distributions (or negatively skewed). Values of g_1 near zero indicate a symmetric distribution. These are not hard and fast rules, however. The value of g_1 is subject to sampling variability and thus only provides a suggestion to the skewness of the underlying distribution.

We still need to know how big is “big”, that is, how do we judge whether an observed value of g_1 is far enough away from zero for the data set to be considered skewed to the right or left? A good rule of thumb is that data sets with skewness larger than $2\sqrt{6/n}$ in magnitude are substantially skewed, in the direction of the sign of g_1 . See Tabachnick & Fidell [Tabachnick2006] for details.

Sample Excess Kurtosis

The *sample excess kurtosis*, denoted by g_2 , is given by the formula

$$g_2 = \frac{1}{n} \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{s^4} - 3. \quad (2.7)$$

The sample excess kurtosis takes values $-2 \leq g_2 < \infty$. The subtraction of 3 may seem mysterious but it is done so that mound shaped samples have values of g_2 near zero. Samples with $g_2 > 0$ are called *leptokurtic*, and samples with $g_2 < 0$ are called *platykurtic*. Samples with $g_2 \approx 0$ are called *mesokurtic*.

As a rule of thumb, if $|g_2| > 4\sqrt{6/n}$ then the sample excess kurtosis is substantially different from zero in the direction of the sign of g_2 . See Tabachnick & Fidell [Tabachnick2006] for details.

Notice that both the sample skewness and the sample kurtosis are invariant with respect to location and scale, that is, the values of g_1 and g_2 do not depend on the measurement units of the data.

How to do it with R

The `e1071` package [e1071] has the `skewness` function for the sample skewness and the `kurtosis` function for the sample excess kurtosis. Both functions have a `na.rm` argument which is `FALSE` by default.

Example 2.11. We said earlier that the `discoveries` data looked positively skewed; let's see what the statistics say:

```
skewness(discoveries)
```

```
[1] 1.2076
```

```
2*sqrt(6/length(discoveries))
```

```
[1] 0.4898979
```

The data are definitely skewed to the right. Let us check the sample excess kurtosis of the `UKDriverDeaths` data:

```
kurtosis(UKDriverDeaths)
```

```
[1] 0.07133848
```

```
4*sqrt(6/length(UKDriverDeaths))
```

```
[1] 0.7071068
```

so that the `UKDriverDeaths` data appear to be mesokurtic, or at least not substantially leptokurtic.

2.4 Exploratory Data Analysis

This field was founded (mostly) by John Tukey (1915-2000). Its tools are useful when not much is known regarding the underlying causes associated with the data set, and are often used for checking assumptions. For example, suppose we perform an experiment and collect some data. . . now what? We look at the data using exploratory visual tools.

2.4.1 More About Stem-and-leaf Displays

There are many bells and whistles associated with stemplots, and the `stem.leaf` function can do many of them.

- **Trim Outliers:** Some data sets have observations that fall far from the bulk of the other data (in a sense made more precise in Section ??). These extreme observations often obscure the underlying structure to the data and are best left out of the data display. The `trim.outliers` argument (which is `TRUE` by default) will separate the extreme observations from the others and graph the stemplot without them; they are listed at the bottom (respectively, top) of the stemplot with the label `HI` (respectively `LO`).
- **Split Stems:** The standard stemplot has only one line per stem, which means that all observations with first digit 3 are plotted on the same line, regardless of the value of the second digit. But this gives some stemplots a “skyscraper” appearance, with too many observations stacked onto the same stem. We can often fix the display by increasing the number of lines available for a given stem. For example, we could make two lines per stem, say, `3*` and `3..` Observations with second digit 0 through 4 would go on the upper line, while observations with second digit 5 through 9 would go on the lower line. (We could do a similar thing with five lines per stem, or even ten lines per stem.) The end result is a more spread out stemplot which often looks better. A good example of this was shown on page ??.
- **Depths:** these are used to give insight into the balance of the observations as they accumulate toward the median. In a column beside the standard stemplot, the frequency of the stem containing the sample median is shown in parentheses. Next, frequencies are accumulated from the outside inward, including the outliers. Distributions that are more symmetric will have better balanced depths on either side of the sample median.

How to do it with R

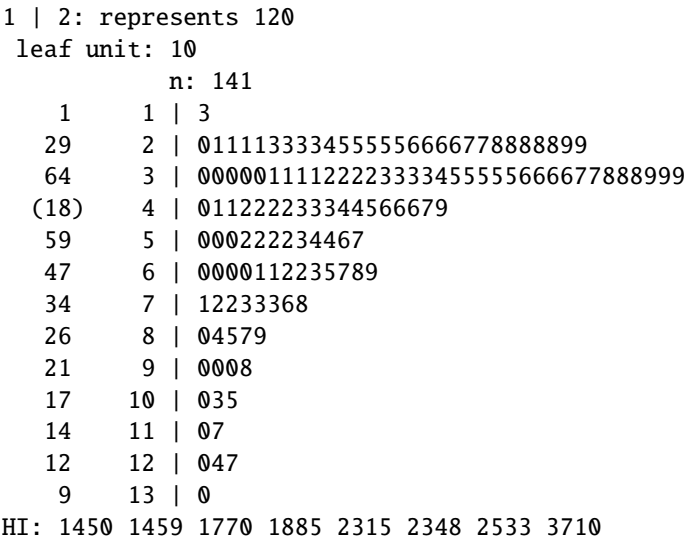
The basic command is `stem(x)` or a more sophisticated version written by Peter Wolf called `stem.leaf(x)` in the R Commander. We will describe `stem.leaf` since that is the one used by R Commander.

WARNING: Sometimes when making a stem-and-leaf display the result will not be what you expected. There are several reasons for this:

- Stemplots by default will trim extreme observations (defined in Section ??) from the display. This in some cases will result in stemplots that are not as wide as expected.
- The leafs digit is chosen automatically by `stem.leaf` according to an algorithm that the computer believes will represent the data well. Depending on the choice of the digit, `stem.leaf` may drop digits from the data or round the values in unexpected ways.

Let us take a look at the `rivers` data set.

```
stem.leaf(rivers)
```



The stem-and-leaf display shows a right-skewed shape to the `rivers` data distribution. Notice that the last digit of each of the data values were dropped from the display. Notice also that there were eight extreme observations identified by the computer, and their exact values are listed at the bottom of the stemplot. Look at the scale on the left of the stemplot and try to imagine how ridiculous the graph would have looked had we tried to include enough stems to include these other eight observations; the stemplot would have stretched over several pages. Notice finally that we can use the depths to approximate the sample median for these data. The median lies in the row identified by (18), which means that the median is the average of the ninth and tenth observation on that row. Those two values correspond to 43 and 43, so a good guess for the median would be 430. (For the record, the sample median is $\tilde{x} = 425$. Recall that stemplots round the data to the nearest stem-leaf pair.)

Next let us see what the `precip` data look like.

```
stem.leaf(precip)
```

```
1 | 2: represents 12
leaf unit: 1
      n: 70
LO: 7 7.2 7.8 7.8
      8   1* | 1344
      13  1. | 55677
      16  2* | 024
      18  2. | 59
      28  3* | 0000111234
(15)  3. | 555566677788899
      27  4* | 0000122222334
      14  4. | 56688899
      6   5* | 44
      4   5. | 699
HI: 67
```

Here is an example of split stems, with two lines per stem. The final digit of each datum has been dropped for the display. The data appear to be left skewed with four extreme values to the left and one extreme value to the right. The sample median is approximately 37 (it turns out to be 36.6).

2.4.2 Hinges and the Five Number Summary

Given a data set x_1, x_2, \dots, x_n , the hinges are found by the following method:

- Find the order statistics $x_{(1)}, x_{(2)}, \dots, x_{(n)}$.
- The lower hinge h_L is in position $L = \lfloor (n + 3)/2 \rfloor / 2$, where the symbol $\lfloor x \rfloor$ denotes the largest integer less than or equal to x . If the position L is not an integer, then the hinge h_L is the average of the adjacent order statistics.
- The upper hinge h_U is in position $n + 1 - L$.

Given the hinges, the five number summary (5NS) is

$$5NS = (x_{(1)}, h_L, \tilde{x}, h_U, x_{(n)}). \tag{2.8}$$

An advantage of the 5NS is that it reduces a potentially large data set to a shorter list of only five numbers, and further, these numbers give insight regarding the shape of the data distribution similar to the sample quantiles in Section ??.

How to do it with R

If the data are stored in a vector \mathbf{x} , then you can compute the $5NS$ with the `fivenum` function.

2.4.3 Boxplots

A boxplot is essentially a graphical representation of the $5NS$. It can be a handy alternative to a stripchart when the sample size is large.

A boxplot is constructed by drawing a box alongside the data axis with sides located at the upper and lower hinges. A line is drawn parallel to the sides to denote the sample median. Lastly, whiskers are extended from the sides of the box to the maximum and minimum data values (more precisely, to the most extreme values that are not potential outliers, defined below).

Boxplots are good for quick visual summaries of data sets, and the relative positions of the values in the $5NS$ are good at indicating the underlying shape of the data distribution, although perhaps not as effectively as a histogram. Perhaps the greatest advantage of a boxplot is that it can help to objectively identify extreme observations in the data set as described in the next section.

Boxplots are also good because one can visually assess multiple features of the data set simultaneously:

- **Center:** can be estimated by the sample median, \tilde{x} .
- **Spread:** can be judged by the width of the box, $h_U - h_L$. We know that this will be close to the IQR , which can be compared to s and the MAD , perhaps after rescaling if appropriate.
- **Shape:** is indicated by the relative lengths of the whiskers, and the position of the median inside the box. Boxes with unbalanced whiskers indicate skewness in the direction of the long whisker. Skewed distributions often have the median tending in the opposite direction of skewness. Kurtosis can be assessed using the box and whiskers. A wide box with short whiskers will tend to be platykurtic, while a skinny box with wide whiskers indicates leptokurtic distributions.
- **Extreme observations:** are identified with open circles (see below).

2.4.4 Outliers

A *potential outlier* is any observation that falls beyond 1.5 times the width of the box on either side, that is, any observation less than $h_L - 1.5(h_U - h_L)$ or greater than $h_U + 1.5(h_U - h_L)$. A *suspected outlier* is any observation that falls beyond 3 times the width of the box on either side. In R, both potential and suspected outliers (if present) are denoted by open circles; there is no distinction between the two.

When potential outliers are present, the whiskers of the boxplot are then shortened to extend to the most extreme observation that is not a potential outlier. If an outlier is displayed in a boxplot, the index of the observation may be identified in a subsequent plot in Rcmdr by clicking the **Identify outliers with mouse** option in the **Boxplot** dialog.

What do we do about outliers? They merit further investigation. The primary goal is to determine why the observation is outlying, if possible. If the observation is a typographical error, then it should be corrected before continuing. If the observation is from a subject that does not belong to the population of interest, then perhaps the datum should be removed. Otherwise, perhaps the value is hinting at some hidden structure to the data.

How to do it with R

The quickest way to visually identify outliers is with a boxplot, described above. Another way is with the `boxplot.stats` function.

Example 2.12 (Lengths of Major North American Rivers). We will look for potential outliers in the `rivers` data.

```
boxplot.stats(rivers)$out
```

```
[1] 1459 1450 1243 2348 3710 2315 2533 1306 1270 1885 1770
```

We may change the `coef` argument to 3 (it is 1.5 by default) to identify suspected outliers.

```
boxplot.stats(rivers, coef = 3)$out
```

```
[1] 2348 3710 2315 2533 1885
```

2.4.5 Standardizing variables

It is sometimes useful to compare data sets with each other on a scale that is independent of the measurement units. Given a set of observed data x_1, x_2, \dots, x_n we get z scores, denoted z_1, z_2, \dots, z_n , by means of the following formula

$$z_i = \frac{x_i - \bar{x}}{s}, \quad i = 1, 2, \dots, n.$$

How to do it with R

The `scale` function will rescale a numeric vector (or data frame) by subtracting the sample mean from each value (column) and/or by dividing each observation by the sample standard deviation.

2.5 Multivariate Data and Data Frames

We have had experience with vectors of data, which are long lists of numbers. Typically, each entry in the vector is a single measurement on a subject or experimental unit in the study. We saw in Section ?? how to form vectors with the `c` function or the `scan` function.

However, statistical studies often involve experiments where there are two (or more) measurements associated with each subject. We display the measured information in a rectangular array in which each row corresponds to a subject, and the columns contain the measurements for each respective variable. For instance, if one were to measure the height and weight and hair color of each of 11 persons in a research study, the information could be represented with a rectangular array. There would be 11 rows. Each row would have the person's height in the first column and hair color in the second column.

The corresponding objects in R are called *data frames*, and they can be constructed with the `data.frame` function. Each row is an observation, and each column is a variable.

Example 2.13. Suppose we have two vectors `x` and `y` and we want to make a data frame out of them.

```
x <- 5:8
y <- letters[3:6]
A <- data.frame(v1 = x, v2 = y)
```

Notice that `x` and `y` are the same length. This is *necessary*. Also notice that `x` is a numeric vector and `y` is a character vector. We may choose numeric and character vectors (or even factors) for the columns of the data frame, but each column must be of exactly one type. That is, we can have a column for `height` and a column for `gender`, but we will get an error if we try to mix function `height` (numeric) and `gender` (character or factor) information in the same column.

Indexing of data frames is similar to indexing of vectors. To get the entry in row i and column j do `A[i, j]`. We can get entire rows and columns by omitting the other index.

```
A[3, ]
```

```
  v1 v2
3  7  e
```

```
A[, 1]
```

```
[1] 5 6 7 8
```

```
A[, 2]
```

```
[1] c d e f
Levels: c d e f
```

There are several things happening above. Notice that `A[3,]` gave a data frame (with the same entries as the third row of `A`) yet `A[, 1]` is a numeric vector. `A[, 2]` is a factor vector because the default setting for `data.frame` is `stringsAsFactors = TRUE`.

Data frames have a `names` attribute and the names may be extracted with the `names` function. Once we have the names we may extract given columns by way of the dollar sign.

```
names(A)
```

```
[1] "v1" "v2"
```

```
A['v1']
```

```
  v1
1  5
2  6
3  7
4  8
```

The above is identical to `A[, 1]`.

2.5.1 Bivariate Data

- Stacked bar charts
- odds ratio and relative risk
- Introduce the sample correlation coefficient.

The *sample Pearson product-moment correlation coefficient*:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- independent of scale
- $-1 < r < 1$
- measures *strength* and *direction* of linear association
- Two-Way Tables. Done with `table`, or in the R Commander by following Statistics ▸ Contingency Tables ▸ Two-way Tables. You can also enter and analyze a two-way table.
- `table`
- `prop.table`
- `addmargins`
- `rowPercents (Rcmdr)`

- colPercents (Rcmdr)
- totPercents(Rcmdr)
- `A <- xtabs(~ gender + race, data = RcmdrTestDrive)`
- `xtabs(Freq ~ Class + Sex, data = Titanic)` from built in table
- `barplot(A, legend.text=TRUE)`
- `barplot(t(A), legend.text=TRUE)`
- `barplot(A, legend.text=TRUE, beside = TRUE)`
- `spineplot(gender ~ race, data = RcmdrTestDrive)`
- Spine plot: plots categorical versus categorical
- Scatterplot: look for linear association and correlation.
- `carb ~ optden, data = Formaldehyde (boring)`
- `conc ~ rate, data = Puromycin`
- `xyplot(accel ~ dist, data = attenu)` nonlinear association
- `xyplot(eruptions ~ waiting, data = faithful)` (linear, two groups)
- `xyplot(Petal.Width ~ Petal.Length, data = iris)`
- `xyplot(pressure ~ temperature, data = pressure)` (exponential growth)
- `xyplot(weight ~ height, data = women)` (strong positive linear)

2.5.2 Multivariate Data

Multivariate Data Display

- Multi-Way Tables. You can do this with `table`, or in R Commander by following `Statistics > Contingency Tables > Multi-way Tables`.
- Scatterplot matrix. used for displaying pairwise scatterplots simultaneously. Again, look for linear association and correlation.
- 3D Scatterplot. See Figure ??
- `plot(state.region, state.division)`
- `barplot(table(state.division, state.region), legend.text=TRUE)`

```
require(graphics)
mosaicplot(HairEyeColor)
x <- apply(HairEyeColor, c(1, 2), sum)
x
mosaicplot(x, main = "Relation between hair and eye color")
y <- apply(HairEyeColor, c(1, 3), sum)
y
```

```
mosaicplot(y, main = "Relation between hair color and sex")
z <- apply(HairEyeColor, c(2, 3), sum)
z
mosaicplot(z, main = "Relation between eye color and sex")
```

2.6 Comparing Populations

Sometimes we have data from two or more groups (or populations) and we would like to compare them and draw conclusions. Some issues that we would like to address:

- Comparing centers and spreads: variation within versus between groups
- Comparing clusters and gaps
- Comparing outliers and unusual features
- Comparing shapes.

2.6.1 Numerically

I am thinking here about the `Statistics > Numerical Summaries > Summarize by groups` option or the `Statistics > Summaries > Table of Statistics` option.

2.6.2 Graphically

- Boxplots
- Variable width: the width of the drawn boxplots are proportional to $\sqrt{n_i}$, where n_i is the size of the i^{th} group. Why? Because many statistics have variability proportional to the reciprocal of the square root of the sample size.
- Notches: extend to $1.58 \cdot (h_U - h_L) / \sqrt{n}$. The idea is to give roughly a 95% confidence interval for the difference in two medians. See Chapter ??.
- Stripcharts
- `stripchart(weight ~ feed, method = "stack", data = chickwts)`
- Bar Graphs
- `barplot(xtabs(Freq ~ Admit + Gender, data = UCBAAdmissions))` stacked bar chart
- `barplot(xtabs(Freq ~ Admit, data = UCBAAdmissions))`
- `barplot(xtabs(Freq ~ Gender + Admit, data = UCBAAdmissions, legend = TRUE, beside = TRUE))` oops, discrimination.
- `barplot(xtabs(Freq ~ Admit+Dept, data = UCBAAdmissions), legend = TRUE, beside = TRUE)` different departments have different standards
- `barplot(xtabs(Freq ~ Gender+Dept, data = UCBAAdmissions), legend = TRUE, beside = TRUE)` men mostly applied to easy departments, women mostly applied to difficult departments

- `barplot(xtabs(Freq ~ Gender+Dept, data = UCBAAdmissions), legend = TRUE, beside = TRUE)`
- `barchart(Admit ~ Freq, data = C)`
- `barchart(Admit ~ Freq|Gender, data = C)`
- `barchart(Admit ~ Freq | Dept, groups = Gender, data = C)`
- `barchart(Admit ~ Freq | Dept, groups = Gender, data = C, auto.key = TRUE)`
- Histograms
- `~ breaks | wool{*}tension, data = warpbreaks`
- `~ weight | feed, data = chickwts`
- `~ weight | group, data = PlantGrowth`
- `~ count | spray, data = InsectSprays`
- `~ len | dose, data = ToothGrowth`
- `~ decrease | treatment, data = OrchardSprays` (or `rowpos` or `colpos`)
- Scatterplots

```
xyplot(Petal.Width ~ Petal.Length, data = iris, group = Species)
```

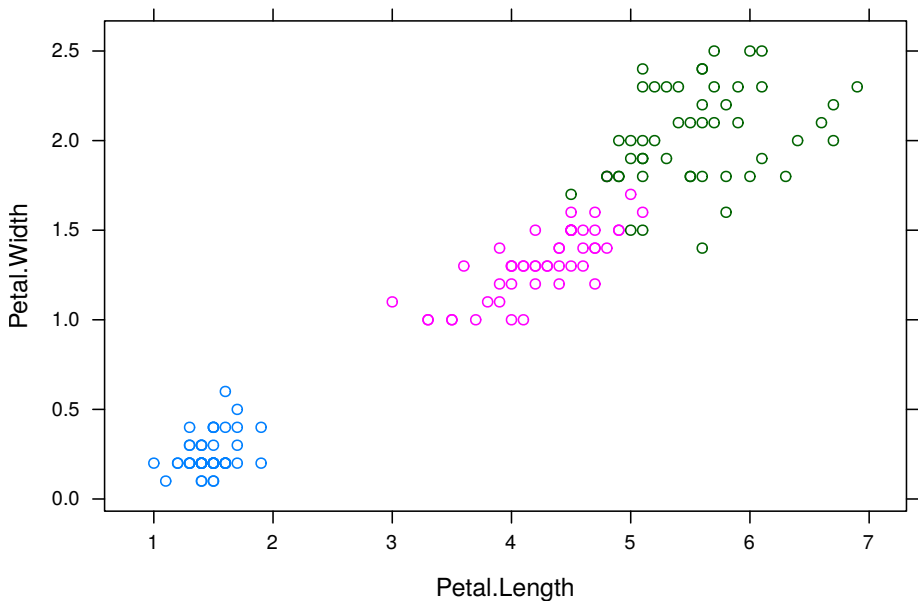


Figure 2.8: Scatterplot of Petal width versus length in the `iris` data.

- Scatterplot matrices
- `sploM(~ cbind(GNP.deflator, GNP, Unemployed, Armed.Forces, Population, Year, Employed), data = longley)`
- `sploM(~ cbind(pop15, pop75, dpi), data = LifeCycleSavings)`
- `sploM(~ cbind(Murder, Assault, Rape), data = USArrests)`

- `spiom(~ cbind(CONT, INTG, DMNR), data = USJudgeRatings)`
- `spiom(~ cbind(area,peri,shape,perm), data = rock)`
- `spiom(~ cbind(Air.Flow, Water.Temp, Acid.Conc., stack.loss), data = stackloss)`
- `spiom(~ cbind(Fertility,Agriculture,Examination,Education,Catholic,Infant.Mortality), data = swiss)`
- `spiom(~ cbind(Fertility,Agriculture,Examination), data = swiss)` (positive and negative)
- Dot charts
- `dotchart(USPersonalExpenditure)`
- `dotchart(t(USPersonalExpenditure))`
- `dotchart(WorldPhones)` (transpose is no good)
- `freeny.x` is no good, neither is `volcano`
- `dotchart(UCBAdmissions{[,1]})`
- `dotplot(Survived ~ Freq | Class, groups = Sex, data = B)`
- `dotplot(Admit ~ Freq | Dept, groups = Gender, data = C)`
- Mosaic plot
- `mosaic(~ Survived + Class + Age + Sex, data = Titanic)` (or just `mosaic(Titanic)`)
- `mosaic(~ Admit + Dept + Gender, data = UCBAdmissions)`
- Spine plots
- `spineplot(xtabs(Freq ~ Admit + Gender, data = UCBAdmissions))`
- rescaled barplot
- Quantile-quantile plots: There are two ways to do this. One way is to compare two independent samples (of the same size). `qqplot(x,y)`. Another way is to compare the sample quantiles of one variable to the theoretical quantiles of another distribution.

Given two samples x_1, x_2, \dots, x_n , and y_1, y_2, \dots, y_n , we may find the order statistics $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ and $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n)}$. Next, plot the n points $(x_{(1)}, y_{(1)})$, $(x_{(2)}, y_{(2)})$, \dots , $(x_{(n)}, y_{(n)})$.

It is clear that if $x_{(k)} = y_{(k)}$ for all $k = 1, 2, \dots, n$, then we will have a straight line. It is also clear that in the real world, a straight line is NEVER observed, and instead we have a scatterplot that hopefully had a general linear trend. What do the rules tell us?

- If the y -intercept of the line is greater (less) than zero, then the center of the Y data is greater (less) than the center of the X data.

- If the slope of the line is greater (less) than one, then the spread of the Y data is greater (less) than the spread of the X data.

2.6.3 Lattice Graphics

The following types of plots are useful when there is one variable of interest and there is a factor in the data set by which the variable is categorized.

It is sometimes nice to set `lattice.options(default.theme = "col.whitebg")`

Side by side boxplots

```
bwplot(~weight | feed, data = chickwts)
```

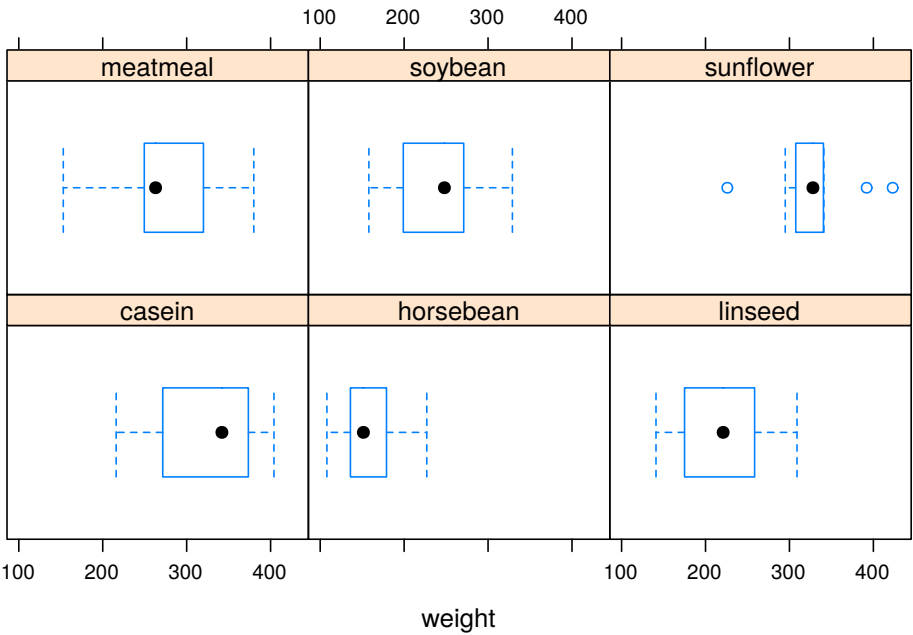


Figure 2.9: Boxplots of weight by feed type in the chickwts data.

Histograms

```
histogram(~age | education, data = infert)
```

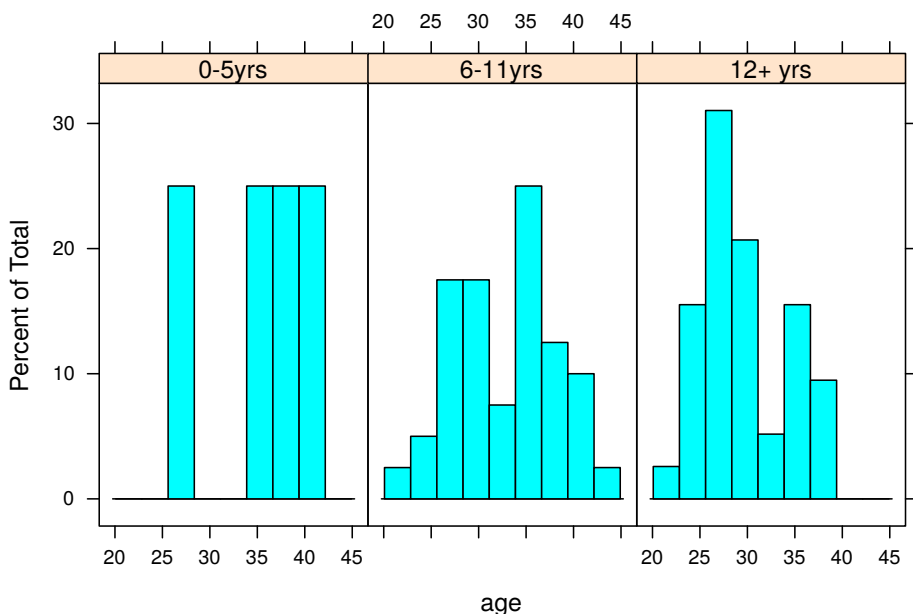


Figure 2.10: Histograms of age by education level from the infert data.

Scatterplots

```
xyplot(Petal.Length ~ Petal.Width | Species, data = iris)
```

Coplots

```
coplot(conc ~ uptake | Type * Treatment, data = CO2)
```

NULL

2.7 Exercises

Open R and issue the following commands at the command line to get started. Note that you need to have the `RcmdrPlugin.IPSUR` package [`RcmdrPlugin.IPSUR`] installed, and for some exercises you need the `e1071` package [`e1071`].

```
library("RcmdrPlugin.IPSUR")
data(RcmdrTestDrive)
```

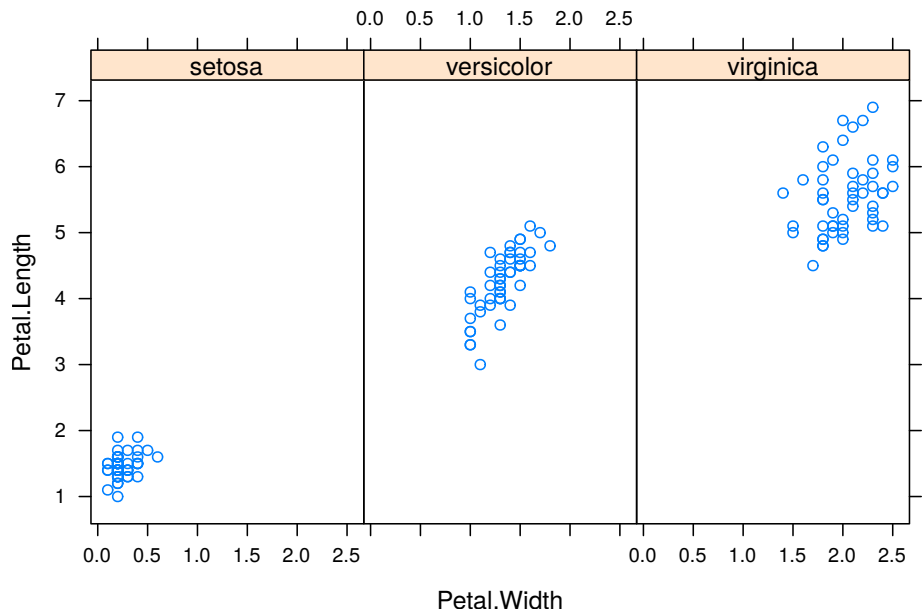


Figure 2.11: An xyplot of Petal.Length versus Petal.Width by Species in the iris data.

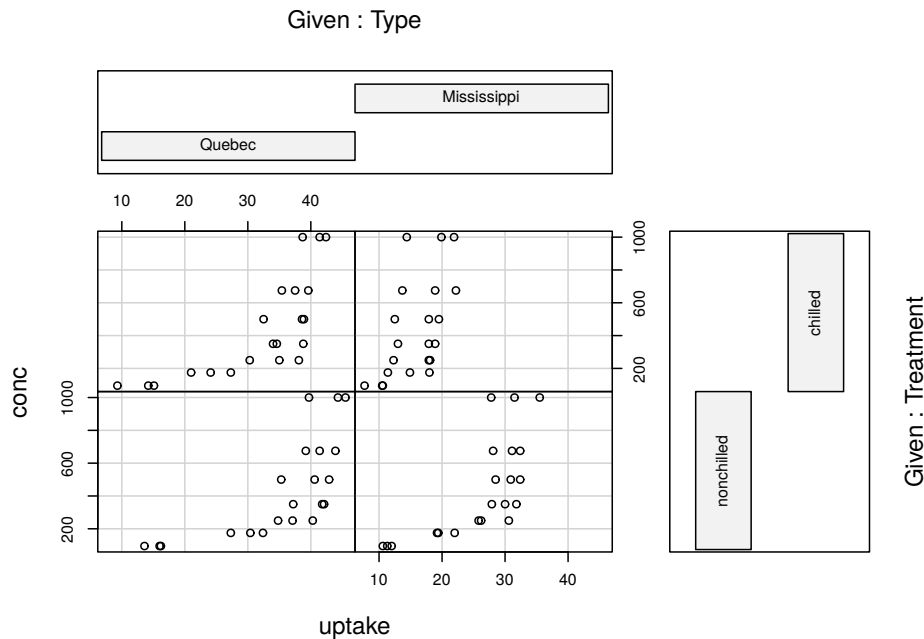


Figure 2.12: A coplot of conc versus uptake by Type and Treatment.

```
attach(RcmdrTestDrive)
names(RcmdrTestDrive)
```

To load the data in the R Commander (Rcmdr), click the Data Set button, and select RcmdrTestDrive as the active data set. To learn more about the data set and where it comes from, type `?RcmdrTestDrive` at the command line.

Exercise 2.1. Perform a summary of all variables in RcmdrTestDrive. You can do this with the command `summary(RcmdrTestDrive)`.

Alternatively, you can do this in the Rcmdr with the sequence Statistics ► Summaries ► Active Data Set. Report the values of the summary statistics for each variable.

Exercise 2.2. Make a table of the race variable. Do this with Statistics ► Summaries ► Frequency Distributions - IPSUR... 1. Which ethnicity has the highest frequency? 1. Which ethnicity has the lowest frequency? 1. Include a bar graph of race. Do this with Graphs ► IPSUR - Bar Graph...

Exercise 2.3. Calculate the average salary by the factor gender. Do this with Statistics ► Summaries ► Table of Statistics...

1. Which gender has the highest mean salary?
2. Report the highest mean salary.
3. Compare the spreads for the genders by calculating the standard deviation of salary by gender. Which gender has the biggest standard deviation?
4. Make boxplots of salary by gender with the following method:
 - i) On the Rcmdr, click Graphs ► IPSUR - Boxplot...
 - ii) In the Variable box, select salary.
 - iii) Click the Plot by groups... box and select gender. Click OK.
 - iv) Click OK to graph the boxplot.

How does the boxplot compare to your answers to (1) and (3)?

Exercise 2.4. For this problem we will study the variable reduction.

1. Find the order statistics and store them in a vector `x`. *Hint:* `x <- sort(reduction)`
2. Find $x_{(137)}$, the 137th order statistic.
3. Find the IQR.
4. Find the Five Number Summary (5NS).
5. Use the 5NS to calculate what the width of a boxplot of reduction would be.
6. Compare your answers (3) and (5). Are they the same? If not, are they close?
7. Make a boxplot of reduction, and include the boxplot in your report. You can do this with the `boxplot` function, or in Rcmdr with Graphs ► IPSUR - Boxplot...

8. Are there any potential/suspected outliers? If so, list their values. *Hint*: use your answer to (a).
9. Using the rules discussed in the text, classify answers to (8), if any, as *potential* or *suspected* outliers.

Exercise 2.5. In this problem we will compare the variables `before` and `after`. Don't forget `library("e1071")`.

1. Examine the two measures of center for both variables. Judging from these measures, which variable has a higher center?
2. Which measure of center is more appropriate for `before`? (You may want to look at a boxplot.) Which measure of center is more appropriate for `after`?
3. Based on your answer to (2), choose an appropriate measure of spread for each variable, calculate it, and report its value. Which variable has the biggest spread? (Note that you need to make sure that your measures are on the same scale.)
4. Calculate and report the skewness and kurtosis for `before`. Based on these values, how would you describe the shape of `before`?
5. Calculate and report the skewness and kurtosis for `after`. Based on these values, how would you describe the shape of `after`?
6. Plot histograms of `before` and `after` and compare them to your answers to (4) and (5).

Exercise 2.6. Describe the following data sets just as if you were communicating with an alien, but one who has had a statistics class. Mention the salient features (data type, important properties, anything special). Support your answers with the appropriate visual displays and descriptive statistics.

1. Conversion rates of Euro currencies stored in `euro`.
2. State abbreviations stored in `state.abb`.
3. Areas of the world's landmasses stored in `islands`.
4. Areas of the 50 United States stored in `state.area`.
5. Region of the 50 United States stored in `state.region`.

Chapter 3

Probability

In this chapter we define the basic terminology associated with probability and derive some of its properties. We discuss three interpretations of probability. We discuss conditional probability and independent events, along with Bayes' Theorem. We finish the chapter with an introduction to random variables, which paves the way for the next two chapters.

In this book we distinguish between two types of experiments: *deterministic* and *random*. A *deterministic* experiment is one whose outcome may be predicted with certainty beforehand, such as combining Hydrogen and Oxygen, or adding two numbers such as $2 + 3$. A *random* experiment is one whose outcome is determined by chance. We posit that the outcome of a random experiment may not be predicted with certainty beforehand, even in principle. Examples of random experiments include tossing a coin, rolling a die, and throwing a dart on a board, how many red lights you encounter on the drive home, how many ants traverse a certain patch of sidewalk over a short period, *etc.*

What do I want them to know?

- that there are multiple interpretations of probability, and the methods used depend somewhat on the philosophy chosen
- nuts and bolts of basic probability jargon: sample spaces, events, probability functions, *etc.*
- how to count
- conditional probability and its relationship with independence
- Bayes' Rule and how it relates to the subjective view of probability
- what we mean by 'random variables', and where they come from

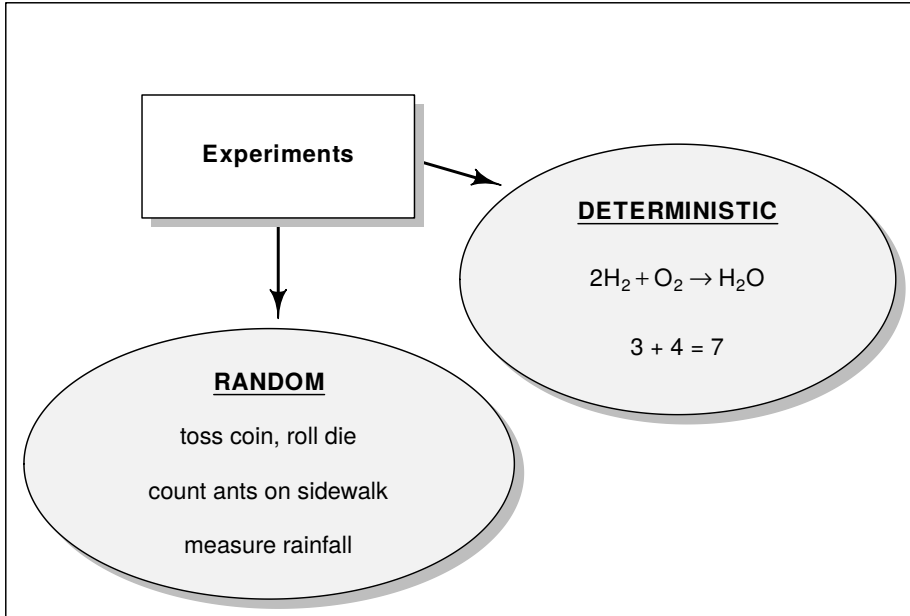


Figure 3.1: Two types of experiments.

3.1 Sample Spaces

For a random experiment E , the set of all possible outcomes of E is called the *sample space* and is denoted by the letter S . For a coin-toss experiment, S would be the results “Head” and “Tail”, which we may represent by $S = \{H, T\}$. Formally, the performance of a random experiment is the unpredictable selection of an outcome in S .

3.1.1 How to do it with R

Most of the probability work in this book is done with the `prob` package [`prob`]. A sample space is (usually) represented by a *data frame*, that is, a rectangular collection of variables (see Section ??). Each row of the data frame corresponds to an outcome of the experiment. The data frame choice is convenient both for its simplicity and its compatibility with the R Commander. Data frames alone are, however, not sufficient to describe some of the more interesting probabilistic applications we will study later; to handle those we will need to consider a more general *list* data structure. See Section ?? for details.

Example 3.1. Consider the random experiment of dropping a Styrofoam cup onto the floor from a height of four feet. The cup hits the ground and eventually comes to rest. It

could land upside down, right side up, or it could land on its side. We represent these possible outcomes of the random experiment by the following.

```
S <- data.frame(lands = c("down", "up", "side"))
S
```

```
  lands
1 down
2  up
3 side
```

The sample space S contains the column `lands` which stores the outcomes `down`, `up`, and `side`.

Some sample spaces are so common that convenience wrappers were written to set them up with minimal effort. The underlying machinery that does the work includes the `expand.grid` function in the base package [**base**], `combn` in the `combinat` package [**combinat**], and `permsn` in the `prob` package [**prob**]¹.

Consider the random experiment of tossing a coin. The outcomes are H and T . We can set up the sample space quickly with the `tosscoin` function:

```
tosscoin(1)
```

```
  toss1
1     H
2     T
```

The number 1 tells `tosscoin` that we only want to toss the coin once. We could toss it three times:

```
tosscoin(3)
```

```
  toss1 toss2 toss3
1     H     H     H
2     T     H     H
3     H     T     H
4     T     T     H
5     H     H     T
6     T     H     T
7     H     T     T
8     T     T     T
```

Alternatively we could roll a fair die:

```
rolldie(1)
```

```
X1
```

¹The seasoned R user can get the job done without the convenience wrappers. I encourage the beginner to use them to get started, but I also recommend that introductory students wean themselves as soon as possible. The wrappers were designed for ease and intuitive use, not for speed or efficiency.

```
1 1
2 2
3 3
4 4
5 5
6 6
```

The `rolldie` function defaults to a 6-sided die, but we can specify others with the `nsides` argument. The command `rolldie(3, nsides = 4)` would be used to roll a 4-sided die three times.

Perhaps we would like to draw one card from a standard set of playing cards (it is a long data frame):

```
head(cards())
```

```
rank suit
1    2 Club
2    3 Club
3    4 Club
4    5 Club
5    6 Club
6    7 Club
```

The `cards` function that we just used has optional arguments `jokers` (if you would like Jokers to be in the deck) and `makespace` which we will discuss later. There is also a `roulette` function which returns the sample space associated with one spin on a roulette wheel. There are EU and USA versions available. Interested readers may contribute any other game or sample spaces that may be of general interest.

3.1.2 Sampling from Urns

This is perhaps the most fundamental type of random experiment. We have an urn that contains a bunch of distinguishable objects (balls) inside. We shake up the urn, reach inside, grab a ball, and take a look. That's all.

But there are all sorts of variations on this theme. Maybe we would like to grab more than one ball – say, two balls. What are all of the possible outcomes of the experiment now? It depends on how we sample. We could select a ball, take a look, put it back, and sample again. Another way would be to select a ball, take a look – but do not put it back – and sample again (equivalently, just reach in and grab two balls). There are certainly more possible outcomes of the experiment in the former case than in the latter. In the first (second) case we say that sampling is done *with* (*without*) *replacement*.

There is more. Suppose we do not actually keep track of which ball came first. All we observe are the two balls, and we have no idea about the order in which they were selected. We call this *unordered sampling* (in contrast to *ordered*) because the order of

the selections does not matter with respect to what we observe. We might as well have selected the balls and put them in a bag before looking.

Note that this one general class of random experiments contains as a special case all of the common elementary random experiments. Tossing a coin twice is equivalent to selecting two balls labeled H and T from an urn, with replacement. The die-roll experiment is equivalent to selecting a ball from an urn with six elements, labeled 1 through 6.

How to do it with R

The `prob` package [`prob`] accomplishes sampling from urns with the `urnsamples` function, which has arguments `x`, `size`, `replace`, and `ordered`. The argument `x` represents the urn from which sampling is to be done. The `size` argument tells how large the sample will be. The `ordered` and `replace` arguments are logical and specify how sampling will be performed. We will discuss each in turn.

Example 3.2. Let our urn simply contain three balls, labeled 1, 2, and 3, respectively. We are going to take a sample of size 2 from the urn.

Ordered, With Replacement

If sampling is with replacement, then we can get any outcome 1, 2, or 3 on any draw. Further, by “ordered” we mean that we shall keep track of the order of the draws that we observe. We can accomplish this in R with

```
urnsamples(1:3, size = 2, replace = TRUE, ordered = TRUE)
```

	X1	X2
1	1	1
2	2	1
3	3	1
4	1	2
5	2	2
6	3	2
7	1	3
8	2	3
9	3	3

Notice that rows 2 and 4 are identical, save for the order in which the numbers are shown. Further, note that every possible pair of the numbers 1 through 3 are listed. This experiment is equivalent to rolling a 3-sided die twice, which we could have accomplished with `rolldie(2, nsides = 3)`.

Ordered, Without Replacement

Here sampling is without replacement, so we may not observe the same number twice in any row. Order is still important, however, so we expect to see the outcomes 1, 2 and 2, 1 somewhere in our data frame.

```
urnsamples(1:3, size = 2, replace = FALSE, ordered = TRUE)
```

	X1	X2
1	1	2
2	2	1
3	1	3
4	3	1
5	2	3
6	3	2

This is just as we expected. Notice that there are less rows in this answer due to the more restrictive sampling procedure. If the numbers 1, 2, and 3 represented “Fred”, “Mary”, and “Sue”, respectively, then this experiment would be equivalent to selecting two people of the three to serve as president and vice-president of a company, respectively, and the sample space shown above lists all possible ways that this could be done.

Unordered, Without Replacement

Again, we may not observe the same outcome twice, but in this case, we will only retain those outcomes which (when jumbled) would not duplicate earlier ones.

```
urnsamples(1:3, size = 2, replace = FALSE, ordered = FALSE)
```

	X1	X2
1	1	2
2	1	3
3	2	3

This experiment is equivalent to reaching in the urn, picking a pair, and looking to see what they are. This is the default setting of `urnsamples`, so we would have received the same output by simply typing `urnsamples(1:3, 2)`.

Unordered, With Replacement

The last possibility is perhaps the most interesting. We replace the balls after every draw, but we do not remember the order in which the draws came.

```
urnsamples(1:3, size = 2, replace = TRUE, ordered = FALSE)
```

	X1	X2
1	1	1

```

2  1  2
3  1  3
4  2  2
5  2  3
6  3  3

```

We may interpret this experiment in a number of alternative ways. One way is to consider this as simply putting two 3-sided dice in a cup, shaking the cup, and looking inside – as in a game of *Liar's Dice*, for instance. Each row of the sample space is a potential pair we could observe. Another way is to view each outcome as a separate method to distribute two identical golf balls into three boxes labeled 1, 2, and 3. Regardless of the interpretation, `urnsamples` lists every possible way that the experiment can conclude.

Note that the urn does not need to contain numbers; we could have just as easily taken our urn to be `x = c("Red", "Blue", "Green")`. But, there is an *important* point to mention before proceeding. Astute readers will notice that in our example, the balls in the urn were *distinguishable* in the sense that each had a unique label to distinguish it from the others in the urn. A natural question would be, “What happens if your urn has indistinguishable elements, for example, what if `x = c("Red", "Red", "Blue")`?” The answer is that `urnsamples` behaves as if each ball in the urn is distinguishable, regardless of its actual contents. We may thus imagine that while there are two red balls in the urn, the balls are such that we can tell them apart (in principle) by looking closely enough at the imperfections on their surface.

In this way, when the `x` argument of `urnsamples` has repeated elements, the resulting sample space may appear to be `ordered = TRUE` even when, in fact, the call to the function was `urnsamples(..., ordered = FALSE)`. Similar remarks apply for the `replace` argument.

3.2 Events

An *event* A is merely a collection of outcomes, or in other words, a subset of the sample space². After the performance of a random experiment E we say that the event A *occurred* if the experiment's outcome belongs to A . We say that a bunch of events A_1, A_2, A_3, \dots are *mutually exclusive* or *disjoint* if $A_i \cap A_j = \emptyset$ for any distinct pair $A_i \neq A_j$. For instance, in the coin-toss experiment the events $A = \{\text{Heads}\}$ and $B = \{\text{Tails}\}$ would be mutually exclusive. Now would be a good time to review the algebra of sets in Appendix ??.

²This naive definition works for finite or countably infinite sample spaces, but is inadequate for sample spaces in general. In this book, we will not address the subtleties that arise, but will refer the interested reader to any text on advanced probability or measure theory.

3.2.1 How to do it with R

Given a data frame sample/probability space S , we may extract rows using the `[]` operator:

```
S <- tosscoin(2, makespace = TRUE)
S[1:3, ]
```

	toss1	toss2	probs
1	H	H	0.25
2	T	H	0.25
3	H	T	0.25

```
S[c(2,4), ]
```

	toss1	toss2	probs
2	T	H	0.25
4	T	T	0.25

and so forth. We may also extract rows that satisfy a logical expression using the `subset` function, for instance

```
S <- cards()
```

```
subset(S, suit == "Heart")
```

	rank	suit
27	2	Heart
28	3	Heart
29	4	Heart
30	5	Heart
31	6	Heart
32	7	Heart
33	8	Heart
34	9	Heart
35	10	Heart
36	J	Heart
37	Q	Heart
38	K	Heart
39	A	Heart

```
subset(S, rank %in% 7:9)
```

	rank	suit
6	7	Club
7	8	Club
8	9	Club
19	7	Diamond
20	8	Diamond

```

21    9 Diamond
32    7  Heart
33    8  Heart
34    9  Heart
45    7  Spade
46    8  Spade
47    9  Spade

```

We could continue indefinitely. Also note that mathematical expressions are allowed:

```
subset(rolldie(3), X1+X2+X3 > 16)
```

```

      X1 X2 X3
180    6  6  5
210    6  5  6
215    5  6  6
216    6  6  6

```

3.2.2 Functions for Finding Subsets

It does not take long before the subsets of interest become complicated to specify. Yet the main idea remains: we have a particular logical condition to apply to each row. If the row satisfies the condition, then it should be in the subset. It should not be in the subset otherwise. The ease with which the condition may be coded depends of course on the question being asked. Here are a few functions to get started.

The %in% function

The function %in% helps to learn whether each value of one vector lies somewhere inside another vector.

```

x <- 1:10
y <- 8:12
y %in% x

```

```
[1] TRUE TRUE TRUE FALSE FALSE
```

Notice that the returned value is a vector of length 5 which tests whether each element of y is in x, in turn.

The isin function

It is more common to want to know whether the *whole* vector y is in x. We can do this with the isin function.

```
isin(x,y)
```

```
[1] FALSE
```

Of course, one may ask why we did not try something like `all(y %in% x)`, which would give a single result, `TRUE`. The reason is that the answers are different in the case that `y` has repeated values. Compare:

```
x <- 1:10
y <- c(3,3,7)
```

```
all(y %in% x)
```

```
[1] TRUE
```

```
isin(x,y)
```

```
[1] FALSE
```

The reason for the above is of course that `x` contains the value 3, but `x` does not have *two* 3's. The difference is important when rolling multiple dice, playing cards, *etc.* Note that there is an optional argument `ordered` which tests whether the elements of `y` appear in `x` in the order in which they appear in `y`. The consequences are

```
isin(x, c(3,4,5), ordered = TRUE)
```

```
[1] TRUE
```

```
isin(x, c(3,5,4), ordered = TRUE)
```

```
[1] FALSE
```

The connection to probability is that have a data frame sample space and we would like to find a subset of that space. A `data.frame` method was written for `isin` that simply applies the function to each row of the data frame. We can see the method in action with the following:

```
S <- rolldie(4)
subset(S, isin(S, c(2,2,6), ordered = TRUE))
```

	X1	X2	X3	X4
188	2	2	6	1
404	2	2	6	2
620	2	2	6	3
836	2	2	6	4
1052	2	2	6	5
1088	2	2	1	6
1118	2	1	2	6
1123	1	2	2	6
1124	2	2	2	6
1125	3	2	2	6

1126	4	2	2	6
1127	5	2	2	6
1128	6	2	2	6
1130	2	3	2	6
1136	2	4	2	6
1142	2	5	2	6
1148	2	6	2	6
1160	2	2	3	6
1196	2	2	4	6
1232	2	2	5	6
1268	2	2	6	6

There are a few other functions written to find useful subsets, namely, `countrep` and `isrep`. Essentially these were written to test for (or count) a specific number of designated values in outcomes. See the documentation for details.

3.2.3 Set Union, Intersection, and Difference

Given subsets A and B , it is often useful to manipulate them in an algebraic fashion. To this end, we have three set operations at our disposal: union, intersection, and difference. Below is a table that summarizes the pertinent information about these operations.

Table 3.1: Basic set operations. The first column lists the name, the second shows the typical notation, the third describes set membership, and the fourth shows how to accomplish it with R.

Name	Denoted	Defined by elements	Code
Union	$A \cup B$	in A or B or both	<code>union(A,B)</code>
Intersection	$A \cap B$	in both A and B	<code>intersect(A,B)</code>
Difference	$A \setminus B$	in A but not in B	<code>setdiff(A,B)</code>

Some examples follow.

```
S <- cards()
A <- subset(S, suit == "Heart")
B <- subset(S, rank %in% 7:9)
```

We can now do some set algebra:

```
union(A,B)
```

	rank	suit
6	7	Club
7	8	Club
8	9	Club

```

19    7 Diamond
20    8 Diamond
21    9 Diamond
27    2 Heart
28    3 Heart
29    4 Heart
30    5 Heart
31    6 Heart
32    7 Heart
33    8 Heart
34    9 Heart
35   10 Heart
36    J Heart
37    Q Heart
38    K Heart
39    A Heart
45    7 Spade
46    8 Spade
47    9 Spade

```

```
intersect(A,B)
```

```

      rank  suit
32      7 Heart
33      8 Heart
34      9 Heart

```

```
setdiff(A,B)
```

```

      rank  suit
27      2 Heart
28      3 Heart
29      4 Heart
30      5 Heart
31      6 Heart
35     10 Heart
36      J Heart
37      Q Heart
38      K Heart
39      A Heart

```

```
setdiff(B,A)
```

```

      rank  suit
6       7 Club
7       8 Club
8       9 Club

```

```

19      7 Diamond
20      8 Diamond
21      9 Diamond
45      7 Spade
46      8 Spade
47      9 Spade

```

Notice that `setdiff` is not symmetric. Further, note that we can calculate the *complement* of a set A , denoted A^c and defined to be the elements of S that are not in A simply with `setdiff(S,A)`. There have been methods written for `intersect`, `setdiff`, `subset`, and `union` in the case that the input objects are of class `ps`. See Section ??.

Note 3.1. When the `prob` package [`@prob`] loads you will notice a message: The following object(s) are masked from `package:base`: `intersect`, `setdiff`. The reason for this message is that there already exist methods for the functions `intersect`, `setdiff`, `subset`, and `union` in the base package which ships with R. However, these methods were designed for when the arguments are vectors of the same mode. Since we are manipulating sample spaces which are data frames and lists, it was necessary to write methods to handle those cases as well. When the `prob` package is loaded, R recognizes that there are multiple versions of the same function in the search path and acts to shield the new definitions from the existing ones. But there is no cause for alarm, thankfully, because the `prob` functions have been carefully defined to match the usual base package definition in the case that the arguments are vectors.

3.3 Model Assignment

Let us take a look at the coin-toss experiment more closely. What do we mean when we say “the probability of Heads” or write $\mathbb{P}(\text{Heads})$? Given a coin and an itchy thumb, how do we go about finding what $\mathbb{P}(\text{Heads})$ should be?

3.3.1 The Measure Theory Approach

This approach states that the way to handle $\mathbb{P}(\text{Heads})$ is to define a mathematical function, called a *probability measure*, on the sample space. Probability measures satisfy certain axioms (to be introduced later) and have special mathematical properties, so not just any mathematical function will do. But in any given physical circumstance there are typically all sorts of probability measures from which to choose, and it is left to the experimenter to make a reasonable choice – one usually based on considerations of objectivity. For the tossing coin example, a valid probability measure assigns probability p to the event $\{\text{Heads}\}$, where p is some number $0 \leq p \leq 1$. An experimenter that wishes to incorporate the symmetry of the coin would choose $p = 1/2$ to balance the likelihood of $\{\text{Heads}\}$ and $\{\text{Tails}\}$.

Once the probability measure is chosen (or determined), there is not much left to do. All assignments of probability are made by the probability function, and the experimenter needs only to plug the event {Heads} into the probability function to find $\mathbb{P}(\text{Heads})$. In this way, the probability of an event is simply a calculated value, nothing more, nothing less. Of course this is not the whole story; there are many theorems and consequences associated with this approach that will keep us occupied for the remainder of this book. The approach is called *measure theory* because the measure (probability) of a set (event) is associated with how big it is (how likely it is to occur).

The measure theory approach is well suited for situations where there is symmetry to the experiment, such as flipping a balanced coin or spinning an arrow around a circle with well-defined pie slices. It is also handy because of its mathematical simplicity, elegance, and flexibility. There are literally volumes of information that one can prove about probability measures, and the cold rules of mathematics allow us to analyze intricate probabilistic problems with vigor.

The large degree of flexibility is also a disadvantage, however. When symmetry fails it is not always obvious what an “objective” choice of probability measure should be; for instance, what probability should we assign to {Heads} if we spin the coin rather than flip it? (It is not $1/2$.) Furthermore, the mathematical rules are restrictive when we wish to incorporate subjective knowledge into the model, knowledge which changes over time and depends on the experimenter, such as personal knowledge about the properties of the specific coin being flipped, or of the person doing the flipping.

The mathematician who revolutionized this way to do probability theory was Andrey Kolmogorov, who published a landmark monograph in 1933. See <http://www-history.mcs.st-andrews.ac.uk/Mathematicians/Kolmogorov.html> for more information.

3.3.2 Relative Frequency Approach

This approach states that the way to determine $\mathbb{P}(\text{Heads})$ is to flip the coin repeatedly, in exactly the same way each time. Keep a tally of the number of flips and the number of Heads observed. Then a good approximation to $\mathbb{P}(\text{Heads})$ will be

$$\mathbb{P}(\text{Heads}) \approx \frac{\text{number of observed Heads}}{\text{total number of flips}}. \quad (3.1)$$

The mathematical underpinning of this approach is the celebrated *Law of Large Numbers* which may be loosely described as follows. Let E be a random experiment in which the event A either does or does not occur. Perform the experiment repeatedly, in an identical manner, in such a way that the successive experiments do not influence each other. After each experiment, keep a running tally of whether or not the event A occurred. Let S_n count the number of times that A occurred in the n experiments. Then the law of large numbers says that

$$\frac{S_n}{n} \rightarrow \mathbb{P}(A) \text{ as } n \rightarrow \infty. \quad (3.2)$$

As the reasoning goes, to learn about the probability of an event A we need only repeat the random experiment to get a reasonable estimate of the probability's value, and if we are not satisfied with our estimate then we may simply repeat the experiment more times all the while confident that with more and more experiments our estimate will stabilize to the true value.

The frequentist approach is good because it is relatively light on assumptions and does not worry about symmetry or claims of objectivity like the measure-theoretic approach does. It is perfect for the spinning coin experiment. One drawback to the method is that one can never know the exact value of a probability, only a long-run approximation. It also does not work well with experiments that can not be repeated indefinitely, say, the probability that it will rain today, the chances that you get will get an A in your Statistics class, or the probability that the world is destroyed by nuclear war.

This approach was espoused by Richard von Mises in the early twentieth century, and some of his main ideas were incorporated into the measure theory approach. See <http://www-history.mcs.st-andrews.ac.uk/Biographies/Mises.html> for more.

3.3.3 The Subjective Approach

The subjective approach interprets probability as the experimenter's *degree of belief* that the event will occur. The estimate of the probability of an event is based on the totality of the individual's knowledge at the time. As new information becomes available, the estimate is modified accordingly to best reflect his/her current knowledge. The method by which the probabilities are updated is commonly done with Bayes' Rule, discussed in Section ??.

So for the coin toss example, a person may have $\mathbb{P}(\text{Heads}) = 1/2$ in the absence of additional information. But perhaps the observer knows additional information about the coin or the thrower that would shift the probability in a certain direction. For instance, parlor magicians may be trained to be quite skilled at tossing coins, and some are so skilled that they may toss a fair coin and get nothing but Heads, indefinitely. I have *seen* this. It was similarly claimed in *Bringing Down the House* [Mezrich2003] that MIT students were accomplished enough with cards to be able to cut a deck to the same location, every single time. In such cases, one clearly should use the additional information to assign $\mathbb{P}(\text{Heads})$ away from the symmetry value of $1/2$.

This approach works well in situations that cannot be repeated indefinitely, for example, to assign your probability that you will get an A in this class, the chances of a devastating nuclear war, or the likelihood that a cure for the common cold will be discovered.

The roots of subjective probability reach back a long time. See http://en.wikipedia.org/wiki/Subjective_probability for a short discussion and links to references about the

subjective approach.

3.3.4 Equally Likely Model (ELM)

We have seen several approaches to the assignment of a probability model to a given random experiment and they are very different in their underlying interpretation. But they all cross paths when it comes to the equally likely model which assigns equal probability to all elementary outcomes of the experiment.

The ELM appears in the measure theory approach when the experiment boasts symmetry of some kind. If symmetry guarantees that all outcomes have equal “size”, and if outcomes with equal “size” should get the same probability, then the ELM is a logical objective choice for the experimenter. Consider the balanced 6-sided die, the fair coin, or the dart board with equal-sized wedges.

The ELM appears in the subjective approach when the experimenter resorts to indifference or ignorance with respect to his/her knowledge of the outcome of the experiment. If the experimenter has no prior knowledge to suggest that (s)he prefer Heads over Tails, then it is reasonable for the him/her to assign equal subjective probability to both possible outcomes.

The ELM appears in the relative frequency approach as a fascinating fact of Nature: when we flip balanced coins over and over again, we observe that the proportion of times that the coin comes up Heads tends to $1/2$. Of course if we assume that the measure theory applies then we can prove that the sample proportion must tend to $1/2$ as expected, but that is putting the cart before the horse, in a manner of speaking.

The ELM is only available when there are finitely many elements in the sample space.

How to do it with R

In the **prob** package [**prob**], a probability space is an object of outcomes **S** and a vector of probabilities (called **probs**) with entries that correspond to each outcome in **S**. When **S** is a data frame, we may simply add a column called **probs** to **S** and we will be finished; the probability space will simply be a data frame which we may call **S**. In the case that **S** is a list, we may combine the outcomes and **probs** into a larger list, **space**; it will have two components: **outcomes** and **probs**. The only requirements we need are for the entries of **probs** to be nonnegative and **sum(probs)** to be one.

To accomplish this in R, we may use the **probspace** function. The general syntax is **probspace(x, probs)**, where **x** is a sample space of outcomes and **probs** is a vector (of the same length as the number of outcomes in **x**). The specific choice of **probs** depends on the context of the problem, and some examples follow to demonstrate some of the more common choices.

Example 3.3. The Equally Likely Model asserts that every outcome of the sample space has the same probability, thus, if a sample space has n outcomes, then `probs` would be a vector of length n with identical entries $1/n$. The quickest way to generate `probs` is with the `rep` function. We will start with the experiment of rolling a die, so that $n = 6$. We will construct the sample space, generate the `probs` vector, and put them together with `probspace`.

```
outcomes <- rolldie(1)
p <- rep(1/6, times = 6)
probspace(outcomes, probs = p)
```

	X1	probs
1	1	0.1666667
2	2	0.1666667
3	3	0.1666667
4	4	0.1666667
5	5	0.1666667
6	6	0.1666667

The `probspace` function is designed to save us some time in the most common situations. For example, due to the especial simplicity of the sample space in this case, we could have achieved the same result with only (note the name change for the first column)

```
probspace(1:6, probs = p)
```

	x	probs
1	1	0.1666667
2	2	0.1666667
3	3	0.1666667
4	4	0.1666667
5	5	0.1666667
6	6	0.1666667

Further, since the equally likely model plays such a fundamental role in the study of probability the `probspace` function will assume that the equally model is desired if no `probs` are specified. Thus, we get the same answer with only

```
probspace(1:6)
```

	x	probs
1	1	0.1666667
2	2	0.1666667
3	3	0.1666667
4	4	0.1666667
5	5	0.1666667
6	6	0.1666667

And finally, since rolling dice is such a common experiment in probability classes,

the `rolldie` function has an additional logical argument `makespace` that will add a column of equally likely `probs` to the generated sample space:

```
rolldie(1, makespace = TRUE)
```

```

X1      probs
1  1 0.1666667
2  2 0.1666667
3  3 0.1666667
4  4 0.1666667
5  5 0.1666667
6  6 0.1666667

```

or just `rolldie(1, TRUE)`. Many of the other sample space functions (`tosscoin`, `cards`, `roulette`, *etc.*) have similar `makespace` arguments. Check the documentation for details.

One sample space function that does NOT have a `makespace` option is the `urnsamples` function. This was intentional. The reason is that under the varied sampling assumptions the outcomes in the respective sample spaces are NOT, in general, equally likely. It is important for the user to carefully consider the experiment to decide whether or not the outcomes are equally likely and then use `probspace` to assign the model.

Example 3.4 (An unbalanced coin). While the `makespace` argument to `tosscoin` is useful to represent the tossing of a *fair* coin, it is not always appropriate. For example, suppose our coin is not perfectly balanced, for instance, maybe the *H* side is somewhat heavier such that the chances of a *H* appearing in a single toss is 0.70 instead of 0.5. We may set up the probability space with

```
probspace(tosscoin(1), probs = c(0.70, 0.30))
```

```

toss1 probs
1      H   0.7
2      T   0.3

```

The same procedure can be used to represent an unbalanced die, roulette wheel, *etc.*

3.3.5 Words of Warning

It should be mentioned that while the splendour of R is uncontested, it, like everything else, has limits both with respect to the sample/probability spaces it can manage and with respect to the finite accuracy of the representation of most numbers (see the R FAQ 7.31). When playing around with probability, one may be tempted to set up a probability space for tossing 100 coins or rolling 50 dice in an attempt to answer some scintillating question. (Bear in mind: rolling a die just 9 times has a sample space with over *10 million* outcomes.)

Alas, even if there were enough RAM to barely hold the sample space (and there were enough time to wait for it to be generated), the infinitesimal probabilities that are associated with *so many* outcomes make it difficult for the underlying machinery to handle reliably. In some cases, special algorithms need to be called just to give something that holds asymptotically. User beware.

3.4 Properties of Probability

3.4.1 Probability Functions

A *probability function* is a rule that associates with each event A of the sample space a single number $\mathbb{P}(A) = p$, called the *probability of A* . Any probability function \mathbb{P} satisfies the following three Kolmogorov Axioms:

Axiom 3.2. $\mathbb{P}(A) \geq 0$ for any event $A \subset S$.

Axiom 3.3. $\mathbb{P}(S) = 1$.

Axiom 3.4. If the events A_1, A_2, A_3, \dots are disjoint then

$$\mathbb{P}\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n \mathbb{P}(A_i) \text{ for every } n, \quad (3.3)$$

and furthermore,

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mathbb{P}(A_i). \quad (3.4)$$

The intuition behind the axioms goes like this: first, the probability of an event should never be negative. Second, since the sample space contains all possible outcomes, its probability should be one, or 100%. The last axiom may look intimidating but it simply means that in a sequence of disjoint events (in other words, sets that do not overlap), the total probability (measure) should equal the sum of its parts. For example, the chance of rolling a 1 or a 2 on a die should be the chance of rolling a 1 plus the chance of rolling a 2.

3.4.2 Properties

For any events A and B ,

1. $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$. **Proof:** Since $A \cup A^c = S$ and $A \cap A^c = \emptyset$, we have

$$1 = \mathbb{P}(S) = \mathbb{P}(A \cup A^c) = \mathbb{P}(A) + \mathbb{P}(A^c).$$

2. $\mathbb{P}(\emptyset) = 0$. **Proof:** Note that $\emptyset = S^c$, and use Property 1.
 3. If $A \subset B$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$. **Proof:** Write $B = A \cup (B \cap A^c)$, and notice that $A \cap (B \cap A^c) = \emptyset$; thus

$$\mathbb{P}(B) = \mathbb{P}(A \cup (B \cap A^c)) = \mathbb{P}(A) + \mathbb{P}(B \cap A^c) \geq \mathbb{P}(A),$$

since $\mathbb{P}(B \cap A^c) \geq 0$.

4. $0 \leq \mathbb{P}(A) \leq 1$. **Proof:** The left inequality is immediate from Axiom ??, and the second inequality follows from Property 3 since $A \subset S$.
 5. **The General Addition Rule.**

$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B). \quad (3.5)$$

More generally, for events $A_1, A_2, A_3, \dots, A_n$,

$$\mathbb{P}\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n \mathbb{P}(A_i) - \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{P}(A_i \cap A_j) + \dots + (-1)^{n-1} \mathbb{P}\left(\bigcap_{i=1}^n A_i\right) \quad (3.6)$$

6. **The Theorem of Total Probability.** Let B_1, B_2, \dots, B_n be mutually exclusive and exhaustive. Then

$$\mathbb{P}(A) = \mathbb{P}(A \cap B_1) + \mathbb{P}(A \cap B_2) + \dots + \mathbb{P}(A \cap B_n). \quad (3.7)$$

3.4.3 Assigning Probabilities

A model of particular interest is the *equally likely model*. The idea is to divide the sample space S into a finite collection of elementary events $\{a_1, a_2, \dots, a_N\}$ that are equally likely in the sense that each a_i has equal chances of occurring. The probability function associated with this model must satisfy $\mathbb{P}(S) = 1$, by Axiom 2. On the other hand, it must also satisfy

$$\mathbb{P}(S) = \mathbb{P}(\{a_1, a_2, \dots, a_N\}) = \mathbb{P}(a_1 \cup a_2 \cup \dots \cup a_N) = \sum_{i=1}^N \mathbb{P}(a_i),$$

by Axiom 3. Since $\mathbb{P}(a_i)$ is the same for all i , each one necessarily equals $1/N$.

For an event $A \subset S$, we write A as a collection of elementary outcomes: if $A = \{a_{i_1}, a_{i_2}, \dots, a_{i_k}\}$ then A has k elements and

$$\begin{aligned}
\mathbb{P}(A) &= \mathbb{P}(a_{i_1}) + \mathbb{P}(a_{i_2}) + \cdots + \mathbb{P}(a_{i_k}), \\
&= \frac{1}{N} + \frac{1}{N} + \cdots + \frac{1}{N}, \\
&= \frac{k}{N} = \frac{\#(A)}{\#(S)}.
\end{aligned}$$

In other words, under the equally likely model, the probability of an event A is determined by the number of elementary events that A contains.

Example 3.5. Consider the random experiment E of tossing a coin. Then the sample space is $S = \{H, T\}$, and under the equally likely model, these two outcomes have $\mathbb{P}(H) = \mathbb{P}(T) = 1/2$. This model is taken when it is reasonable to assume that the coin is fair.

Example 3.6. Suppose the experiment E consists of tossing a fair coin twice. The sample space may be represented by $S = \{HH, HT, TH, TT\}$. Given that the coin is fair and that the coin is tossed in an independent and identical manner, it is reasonable to apply the equally likely model.

What is $\mathbb{P}(\text{at least 1 Head})$? Looking at the sample space we see the elements HH , HT , and TH have at least one Head; thus, $\mathbb{P}(\text{at least 1 Head}) = 3/4$.

What is $\mathbb{P}(\text{no Heads})$? Notice that the event $\{\text{no Heads}\} = \{\text{at least one Head}\}^c$, which by Property ?? means $\mathbb{P}(\text{no Heads}) = 1 - \mathbb{P}(\text{at least one head}) = 1 - 3/4 = 1/4$. It is obvious in this simple example that the only outcome with no Heads is TT , however, this complementation trick can be handy in more complicated problems.

Example 3.7. Imagine a three child family, each child being either Boy (B) or Girl (G). An example sequence of siblings would be BGB . The sample space may be written

$$S = \{BBB, BGB, GBB, GGB, BBG, BGG, GBG, GGG, \}.$$

Note that for many reasons (for instance, it turns out that girls are slightly more likely to be born than boys), this sample space is *not* equally likely. For the sake of argument, however, we will assume that the elementary outcomes each have probability $1/8$.

What is $\mathbb{P}(\text{exactly 2 Boys})$? Inspecting the sample space reveals three outcomes with exactly two boys: $\{BBG, BGB, GBB\}$. Therefore $\mathbb{P}(\text{exactly 2 Boys}) = 3/8$.

What is $\mathbb{P}(\text{at most 2 Boys})$? One way to solve the problem would be to count the outcomes that have 2 or less Boys, but a quicker way would be to recognize that the only way that the event $\{\text{at most 2 Boys}\}$ does *not* occur is the event $\{\text{all Boys}\}$.

Thus

$$\mathbb{P}(\text{at most 2 Boys}) = 1 - \mathbb{P}(BBB) = 1 - 1/8 = 7/8.$$

Example 3.8. Consider the experiment of rolling a six-sided die, and let the outcome be the face showing up when the die comes to rest. Then $S = \{1, 2, 3, 4, 5, 6\}$. It is usually reasonable to suppose that the die is fair, so that the six outcomes are equally likely.

Example 3.9. Consider a standard deck of 52 cards. These are usually labeled with the four *suits*: Clubs, Diamonds, Hearts, and Spades, and the 13 *ranks*: 2, 3, 4, ..., 10, Jack (J), Queen (Q), King (K), and Ace (A). Depending on the game played, the Ace may be ranked below 2 or above King.

Let the random experiment E consist of drawing exactly one card from a well-shuffled deck, and let the outcome be the face of the card. Define the events $A = \{\text{draw an Ace}\}$ and $B = \{\text{draw a Club}\}$. Bear in mind: we are only drawing one card.

Immediately we have $\mathbb{P}(A) = 4/52$ since there are four Aces in the deck; similarly, there are 13 Clubs which implies $\mathbb{P}(B) = 13/52$.

What is $\mathbb{P}(A \cap B)$? We realize that there is only one card of the 52 which is an Ace and a Club at the same time, namely, the Ace of Clubs. Therefore $\mathbb{P}(A \cap B) = 1/52$.

To find $\mathbb{P}(A \cup B)$ we may use the above with the General Addition Rule to get

$$\begin{aligned}\mathbb{P}(A \cup B) &= \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B), \\ &= 4/52 + 13/52 - 1/52, \\ &= 16/52.\end{aligned}$$

Example 3.10. Staying with the deck of cards, let another random experiment be the selection of a five card stud poker hand, where “five card stud” means that we draw exactly five cards from the deck without replacement, no more, and no less. It turns out that the sample space S is so large and complicated that we will be obliged to settle for the trivial description $S = \{\text{all possible 5 card hands}\}$ for the time being. We will have a more precise description later.

What is $\mathbb{P}(\text{Royal Flush})$, or in other words, $\mathbb{P}(A, K, Q, J, 10 \text{ all in the same suit})$?

It should be clear that there are only four possible royal flushes. Thus, if we could only count the number of outcomes in S then we could simply divide four by that number and we would have our answer under the equally likely model. This is the subject of Section ??.

How to do it with R

Probabilities are calculated in the `prob` package [`prob`] with the `Prob` function.

Consider the experiment of drawing a card from a standard deck of playing cards. Let's denote the probability space associated with the experiment as S , and let the subsets A and B be defined by the following:

```
S <- cards(makespace = TRUE)
A <- subset(S, suit == "Heart")
B <- subset(S, rank %in% 7:9)
```

Now it is easy to calculate

```
Prob(A)
```

```
[1] 0.25
```

Note that we can get the same answer with

```
Prob(S, suit == "Heart")
```

```
[1] 0.25
```

We also find $\text{Prob}(B) = 0.23$ (listed here approximately, but $12/52$ actually). Internally, the `Prob` function operates by summing the `probs` column of its argument. It will find subsets on-the-fly if desired.

We have as yet glossed over the details. More specifically, `Prob` has three arguments: x , which is a probability space (or a subset of one), `event`, which is a logical expression used to define a subset, and `given`, which is described in Section ??.

WARNING. The `event` argument is used to define a subset of x , that is, the only outcomes used in the probability calculation will be those that are elements of x and satisfy `event` simultaneously. In other words, `Prob(x, event)` calculates `Prob(intersect(x, subset(x, event)))`.

Consequently, x should be the entire probability space in the case that `event` is non-null.

3.5 Counting Methods

The equally-likely model is a convenient and popular way to analyze random experiments. And when the equally likely model applies, finding the probability of an event A amounts to nothing more than counting the number of outcomes that A contains (together with the number of events in S). Hence, to be a master of probability one must be skilled at counting outcomes in events of all kinds.

Proposition 3.1 (The Multiplication Principle). Suppose that an experiment is composed of two successive steps. Further suppose that the first step may be performed in n_1 distinct ways while the second step may be performed in n_2 distinct ways. Then the experiment may be performed in $n_1 n_2$ distinct ways. More generally, if the experiment

is composed of k successive steps which may be performed in n_1, n_2, \dots, n_k distinct ways, respectively, then the experiment may be performed in $n_1 n_2 \cdots n_k$ distinct ways.

Example 3.11. We would like to order a pizza. It will be sure to have cheese (and marinara sauce), but we may elect to add one or more of the following five (5) available toppings:

pepperoni, sausage, anchovies, olives, and green peppers.

How many distinct pizzas are possible?

There are many ways to approach the problem, but the quickest avenue employs the Multiplication Principle directly. We will separate the action of ordering the pizza into a series of stages. At the first stage, we will decide whether or not to include pepperoni on the pizza (two possibilities). At the next stage, we will decide whether or not to include sausage on the pizza (again, two possibilities). We will continue in this fashion until at last we will decide whether or not to include green peppers on the pizza.

At each stage we will have had two options, or ways, to select a pizza to be made. The Multiplication Principle says that we should multiply the 2's to find the total number of possible pizzas: $2 \cdot 2 \cdot 2 \cdot 2 \cdot 2 = 2^5 = 32$.

Example 3.12. We would like to buy a desktop computer to study statistics. We go to a website to build our computer our way. Given a line of products we have many options to customize our computer. In particular, there are 2 choices for a processor, 3 different operating systems, 4 levels of memory, 4 hard drives of differing sizes, and 10 choices for a monitor. How many possible types of computer must the company be prepared to build? **Answer:** $2 \cdot 3 \cdot 4 \cdot 4 \cdot 10 = 960$

3.5.1 Ordered Samples

Imagine a bag with n distinguishable balls inside. Now shake up the bag and select k balls at random. How many possible sequences might we observe?

Proposition 3.2. The number of ways in which one may select an ordered sample of k subjects from a population that has n distinguishable members is

- n^k if sampling is done with replacement,
- $n(n-1)(n-2) \cdots (n-k+1)$ if sampling is done without replacement.

Recall from calculus the notation for *factorials*:

$$\begin{aligned}
1! &= 1, \\
2! &= 2 \cdot 1 = 2, \\
3! &= 3 \cdot 2 \cdot 1 = 6, \\
&\vdots \\
n! &= n(n-1)(n-2) \cdots 3 \cdot 2 \cdot 1.
\end{aligned}$$

Fact: The number of permutations of n elements is $n!$.

Example 3.13. Take a coin and flip it 7 times. How many sequences of Heads and Tails are possible? **Answer:** $2^7 = 128$.

Example 3.14. In a class of 20 students, we randomly select a class president, a class vice-president, and a treasurer. How many ways can this be done? **Answer:** $20 \cdot 19 \cdot 18 = 6840$.

Example 3.15. We rent five movies to watch over the span of two nights. We wish to watch 3 movies on the first night. How many distinct sequences of 3 movies could we possibly watch? **Answer:** $5 \cdot 4 \cdot 3 = 60$.

3.5.2 Unordered Samples

Proposition 3.3. The number of ways in which one may select an unordered sample of k subjects from a population that has n distinguishable members is

- $(n-1+k)!/[(n-1)!k!]$ if sampling is done with replacement,
- $n!/ [k!(n-k)!]$ if sampling is done without replacement.

The quantity $n!/ [k!(n-k)!]$ is called a *binomial coefficient* and plays a special role in mathematics; it is denoted

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (3.8)$$

and is read “ n choose k ”.

Example 3.16. You rent five movies to watch over the span of two nights, but only wish to watch 3 movies the first night. Your friend, Fred, wishes to borrow some movies to watch at his house on the first night. You owe Fred a favor, and allow him to select 2 movies from the set of 5. How many choices does Fred have? **Answer:** $\binom{5}{2} = 10$.

Example 3.17. Place 3 six-sided dice into a cup. Next, shake the cup well and pour out the dice. How many distinct rolls are possible? **Answer:** $(6-1+3)!/[(6-1)!3!] = \binom{8}{5} = 56$.

How to do it with R

The factorial $n!$ is computed with the command `factorial(n)` and the binomial coefficient $\binom{n}{k}$ with the command `choose(n,k)`.

The sample spaces we have computed so far have been relatively small, and we can visually study them without much trouble. However, it is *very* easy to generate sample spaces that are prohibitively large. And while R is wonderful and powerful and does almost everything except wash windows, even R has limits of which we should be mindful.

But we often do not need to actually generate the sample space; it suffices to count the number of outcomes. The `nsamp` function will calculate the number of rows in a sample space made by `urnsamples` without actually devoting the memory resources necessary to generate the space. The arguments are `n`, the number of (distinguishable) objects in the urn, `k`, the sample size, and `replace`, `ordered`, as above.

(#tab-sampling-k-from-n)

Table 3.2: Sampling k from n objects with `urnsamples`.

	ordered = TRUE	ordered = FALSE
replace = TRUE	n^k	$(n - 1 + k)! / [(n - 1)!k!]$
replace = FALSE	$n! / (n - k)!$	$\binom{n}{k}$

Example 3.18. We will compute the number of outcomes for each of the four `urnsamples` examples that we saw in Example ?? . Recall that we took a sample of size two from an urn with three distinguishable elements.

```
nsamp(n=3, k=2, replace = TRUE, ordered = TRUE)
```

[1] 9

```
nsamp(n=3, k=2, replace = FALSE, ordered = TRUE)
```

[1] 6

```
nsamp(n=3, k=2, replace = FALSE, ordered = FALSE)
```

[1] 3

```
nsamp(n=3, k=2, replace = TRUE, ordered = FALSE)
```

[1] 6

Compare these answers with the length of the data frames generated above.

The Multiplication Principle

A benefit of `nsamp` is that it is *vectorized* so that entering vectors instead of numbers for `n`, `k`, `replace`, and `ordered` results in a vector of corresponding answers. This becomes particularly convenient for combinatorics problems.

Example 3.19. There are 11 artists who each submit a portfolio containing 7 paintings for competition in an art exhibition. Unfortunately, the gallery director only has space in the winners' section to accommodate 12 paintings in a row equally spread over three consecutive walls. The director decides to give the first, second, and third place winners each a wall to display the work of their choice. The walls boast 31 separate lighting options apiece. How many displays are possible?

Answer: The judges will pick 3 (ranked) winners out of 11 (with `rep = FALSE`, `ord = TRUE`). Each artist will select 4 of his/her paintings from 7 for display in a row (`rep = FALSE`, `ord = TRUE`), and lastly, each of the 3 walls has 31 lighting possibilities (`rep = TRUE`, `ord = TRUE`). These three numbers can be calculated quickly with

```
n <- c(11, 7, 31)
k <- c(3, 4, 3)
r <- c(FALSE, FALSE, TRUE)

x <- nsamp(n, k, rep = r, ord = TRUE)
```

(Notice that `ordered` is always `TRUE`; `nsamp` will recycle `ordered` and `replace` to the appropriate length.) By the Multiplication Principle, the number of ways to complete the experiment is the product of the entries of `x`:

```
prod(x)
```

```
[1] 24774195600
```

Compare this with the some other ways to compute the same thing:

```
(11*10*9)*(7*6*5*4)*31^3
```

```
[1] 24774195600
```

or alternatively

```
prod(9:11)*prod(4:7)*31^3
```

```
[1] 24774195600
```

or even

```
prod(factorial(c(11, 7))/factorial(c(8, 3)))*31^3
```

```
[1] 24774195600
```

As one can guess, in many of the standard counting problems there aren't substantial savings in the amount of typing; it is about the same using `nsamp` versus `factorial`

and choose. But the virtue of `nsamp` lies in its collecting the relevant counting formulas in a one-stop shop. Ultimately, it is up to the user to choose the method that works best for him/herself.

Example 3.20 (The Birthday Problem). Suppose that there are n people together in a room. Each person announces the date of his/her birthday in turn. The question is: what is the probability of at least one match? If we let the event A represent

$$A = \{\text{there is at least one match}\},$$

then we are looking for $\mathbb{P}(A)$, but as we soon will see, it will be more convenient for us to calculate $\mathbb{P}(A^c)$.

For starters we will ignore leap years and assume that there are only 365 days in a year. Second, we will assume that births are equally distributed over the course of a year (which is not true due to all sorts of complications such as hospital delivery schedules). See http://en.wikipedia.org/wiki/Birthday_problem for more.

Let us next think about the sample space. There are 365 possibilities for the first person's birthday, 365 possibilities for the second, and so forth. The total number of possible birthday sequences is therefore $\#(S) = 365^n$.

Now we will use the complementation trick we saw in Example ???. We realize that the only situation in which A does *not* occur is if there are *no* matches among all people in the room, that is, only when everybody's birthday is different, so

$$\mathbb{P}(A) = 1 - \mathbb{P}(A^c) = 1 - \frac{\#(A^c)}{\#(S)},$$

since the outcomes are equally likely. Let us then suppose that there are no matches. The first person has one of 365 possible birthdays. The second person must not match the first, thus, the second person has only 364 available birthdays from which to choose. Similarly, the third person has only 363 possible birthdays, and so forth, until we reach the n^{th} person, who has only $365 - n + 1$ remaining possible days for a birthday. By the Multiplication Principle, we have $\#(A^c) = 365 \cdot 364 \cdots (365 - n + 1)$, and

$$\mathbb{P}(A) = 1 - \frac{365 \cdot 364 \cdots (365 - n + 1)}{365^n} = 1 - \frac{364}{365} \cdot \frac{363}{365} \cdots \frac{(365 - n + 1)}{365}. \quad (3.9)$$

As a surprising consequence, consider this: how many people does it take to be in the room so that the probability of at least one match is at least 0.50? Clearly, if there is only $n = 1$ person in the room then the probability of a match is zero, and when there are $n = 366$ people in the room there is a 100% chance of a match (recall that we are ignoring leap years). So how many people does it take so that there is an equal chance of a match and no match?

When I have asked this question to students, the usual response is “somewhere around $n = 180$ people” in the room. The reasoning seems to be that in order to get a 50%

chance of a match, there should be 50% of the available days to be occupied. The number of students in a typical classroom is 25, so as a companion question I ask students to estimate the probability of a match when there are $n = 25$ students in the room. Common estimates are a 1%, or 0.5%, or even 0.1% chance of a match. After they have given their estimates, we go around the room and each student announces their birthday. More often than not, we observe a match in the class, to the students' disbelief.

Students are usually surprised to hear that, using the formula above, one needs only $n = 23$ students to have a greater than 50% chance of at least one match. Figure ?? shows a graph of the birthday probabilities:

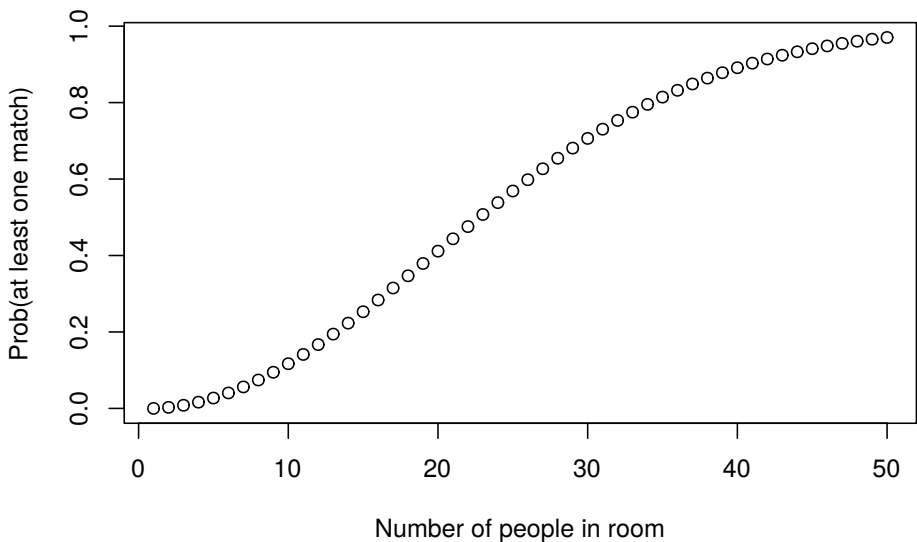


Figure 3.2: The birthday problem. The horizontal line is at $p = 0.50$ and the vertical line is at $n = 23$.

How to do it with R

We can make the plot in Figure ?? with the following sequence of commands.

```
library(RcmdrPlugin.IPSUR)
g <- Vectorize(pbirthday.ipsur)
plot(1:50, g(1:50), xlab = "Number of people in room",
     ylab = "Prob(at least one match)" )
abline(h = 0.5)
abline(v = 23, lty = 2)
remove(g)
```

There is a `Birthday` problem item in the Probability menu of `RcmdrPlugin.IPSUR`. In the base R version, one can compute approximate probabilities for the more general case of probabilities other than $1/2$, for differing total number of days in the year, and even for more than two matches.

3.6 Conditional Probability

Consider a full deck of 52 standard playing cards. Now select two cards from the deck, in succession. Let $A = \{\text{first card drawn is an Ace}\}$ and $B = \{\text{second card drawn is an Ace}\}$. Since there are four Aces in the deck, it is natural to assign $\mathbb{P}(A) = 4/52$. Suppose we look at the first card. What now is the probability of B ? Of course, the answer depends on the value of the first card. If the first card is an Ace, then the probability that the second also is an Ace should be $3/51$, but if the first card is not an Ace, then the probability that the second is an Ace should be $4/51$. As notation for these two situations we write

$$\mathbb{P}(B|A) = 3/51, \quad \mathbb{P}(B|A^c) = 4/51.$$

Definition 3.1. The conditional probability of B given A , denoted $\mathbb{P}(B|A)$, is defined by

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(A)}, \quad \text{if } \mathbb{P}(A) > 0. \quad (3.10)$$

We will not be discussing a conditional probability of B given A when $\mathbb{P}(A) = 0$, even though this theory exists, is well developed, and forms the foundation for the study of stochastic processes³.

Toss a coin twice. The sample space is given by $S = \{HH, HT, TH, TT\}$. Let $A = \{\text{a head occurs}\}$ and $B = \{\text{a head and tail occur}\}$. It should be clear that $\mathbb{P}(A) = 3/4$, $\mathbb{P}(B) = 2/4$, and $\mathbb{P}(A \cap B) = 2/4$. What now are the probabilities $\mathbb{P}(A|B)$ and $\mathbb{P}(B|A)$?

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{2/4}{2/4} = 1,$$

in other words, once we know that a Head and Tail occur, we may be certain that a Head occurs. Next

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(A)} = \frac{2/4}{3/4} = \frac{2}{3},$$

³Conditional probability in this case is defined by means of *conditional expectation*, a topic that is well beyond the scope of this text. The interested reader should consult an advanced text on probability theory, such as Billingsley, Resnick, or Ash Dooleans-Dade.

which means that given the information that a Head has occurred, we no longer need to account for the outcome TT , and the remaining three outcomes are equally likely with exactly two outcomes lying in the set B .

Example 3.21. Toss a six-sided die twice. The sample space consists of all ordered pairs (i, j) of the numbers $1, 2, \dots, 6$, that is, $S = \{(1, 1), (1, 2), \dots, (6, 6)\}$. We know from Section ?? that $\#(S) = 6^2 = 36$. Let $A = \{\text{outcomes match}\}$ and $B = \{\text{sum of outcomes at least } 8\}$. The sample space may be represented by a matrix:

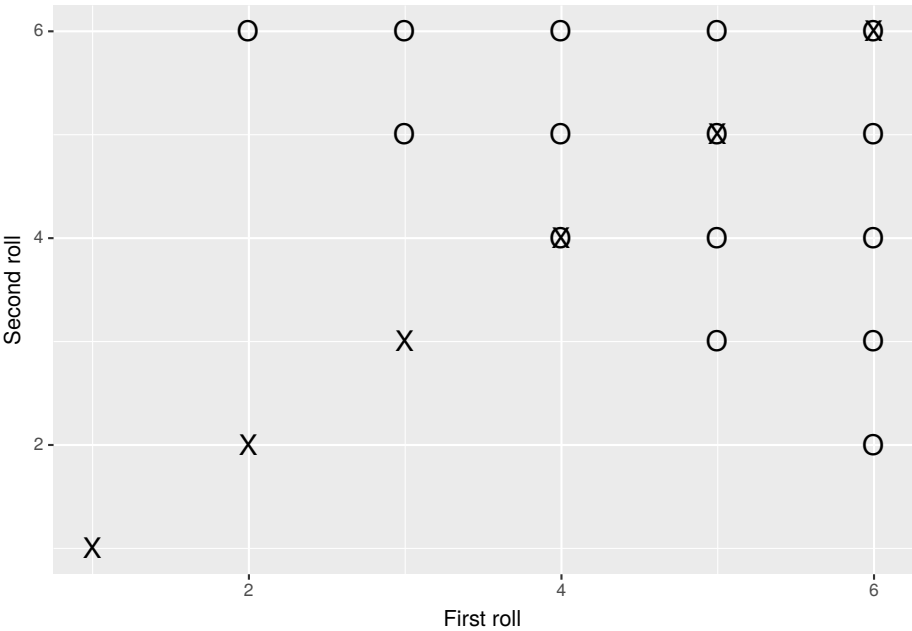


Figure 3.3: Rolling two dice. The outcomes in A are marked with X , the outcomes in B are marked with O .

The outcomes lying in the event A are marked with the symbol “ X ”, the outcomes falling in B are marked with “ O ”, and the outcomes in $A \cap B$ are those where the letters overlap. Now it is clear that $\mathbb{P}(A) = 6/36$, $\mathbb{P}(B) = 15/36$, and $\mathbb{P}(A \cap B) = 3/36$. Finally,

$$\mathbb{P}(A|B) = \frac{3/36}{15/36} = \frac{1}{5}, \quad \mathbb{P}(B|A) = \frac{3/36}{6/36} = \frac{1}{2}.$$

Again, we see that given the knowledge that B occurred (the 15 outcomes in the upper right triangle), there are 3 of the 15 that fall into the set A , thus the probability is $3/15$. Similarly, given that A occurred (we are on the diagonal), there are 3 out of 6 outcomes that also fall in B , thus, the probability of B given A is $1/2$.

3.6.1 How to do it with R

Continuing with Example ??, the first thing to do is set up the probability space with the `rolldie` function.

```
S <- rolldie(2, makespace = TRUE) # assumes ELM
head(S)                          # first few rows
```

```
  X1 X2      probs
1  1  1 0.02777778
2  2  1 0.02777778
3  3  1 0.02777778
4  4  1 0.02777778
5  5  1 0.02777778
6  6  1 0.02777778
```

Next we define the events

```
A <- subset(S, X1 == X2)
B <- subset(S, X1 + X2 >= 8)
```

And now we are ready to calculate probabilities. To do conditional probability, we use the `given` argument of the `Prob` function:

```
Prob(A, given = B)
```

```
[1] 0.2
```

```
Prob(B, given = A)
```

```
[1] 0.5
```

Note that we do not actually need to define the events A and B separately as long as we reference the original probability space S as the first argument of the `Prob` calculation:

```
Prob(S, X1==X2, given = (X1 + X2 >= 8) )
```

```
[1] 0.2
```

```
Prob(S, X1+X2 >= 8, given = (X1==X2) )
```

```
[1] 0.5
```

3.6.2 Properties and Rules

The following theorem establishes that conditional probabilities behave just like regular probabilities when the conditioned event is fixed.

Theorem 3.1. For any fixed event A with $\mathbb{P}(A) > 0$,

1. $\mathbb{P}(B|A) \geq 0$, for all events $B \subset S$,
2. $\mathbb{P}(S|A) = 1$, and
3. If B_1, B_2, B_3, \dots are disjoint events, then

$$\mathbb{P}\left(\bigcup_{k=1}^{\infty} B_k \mid A\right) = \sum_{k=1}^{\infty} \mathbb{P}(B_k|A). \quad (3.11)$$

In other words, $\mathbb{P}(\cdot|A)$ is a legitimate probability function. With this fact in mind, the following properties are immediate:

Proposition 3.4. For any events A, B , and C with $\mathbb{P}(A) > 0$,

1. $\mathbb{P}(B^c|A) = 1 - \mathbb{P}(B|A)$.
2. If $B \subset C$ then $\mathbb{P}(B|A) \leq \mathbb{P}(C|A)$.
3. $\mathbb{P}[(B \cup C)|A] = \mathbb{P}(B|A) + \mathbb{P}(C|A) - \mathbb{P}[(B \cap C)|A]$.
4. **The Multiplication Rule.** For any two events A and B ,

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B|A). \quad (3.12)$$

And more generally, for events $A_1, A_2, A_3, \dots, A_n$,

$$\mathbb{P}(A_1 \cap A_2 \cap \dots \cap A_n) = \mathbb{P}(A_1)\mathbb{P}(A_2|A_1) \cdots \mathbb{P}(A_n|A_1 \cap A_2 \cap \dots \cap A_{n-1}). \quad (3.13)$$

The Multiplication Rule is very important because it allows us to find probabilities in random experiments that have a sequential structure, as the next example shows.

Example 3.22. At the beginning of the section we drew two cards from a standard playing deck. Now we may answer our original question, what is $\mathbb{P}(\text{both Aces})$?

$$\mathbb{P}(\text{both Aces}) = \mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B|A) = \frac{4}{52} \cdot \frac{3}{51} \approx 0.00452.$$

How to do it with R

Continuing Example ??, we set up the probability space by way of a three step process. First we employ the `cards` function to get a data frame `L` with two columns: `rank` and `suit`. Both columns are stored internally as factors with 13 and 4 levels, respectively.

Next we sample two cards randomly from the `L` data frame by way of the `urnsamples` function. It returns a list `M` which contains all possible pairs of rows from `L` (there are `choose(52, 2)` of them). The sample space for this experiment is exactly the list `M`.

At long last we associate a probability model with the sample space. This is right down the `probspace` function's alley. It assumes the equally likely model by default. We call this result `N` which is an object of class `ps` – short for “probability space”.

But do not be intimidated. The object `N` is nothing more than a list with two elements: `outcomes` and `probs`. The `outcomes` element is itself just another list, with

`choose(52, 2)` entries, each one a data frame with two rows which correspond to the pair of cards chosen. The `probs` element is just a vector with `choose(52, 2)` entries all the same: $1/\text{choose}(52, 2)$.

Putting all of this together we do

```
L <- cards()
M <- urnsamples(L, size = 2)
N <- probspace(M)
N[[1]][[1]]; N$probs[1]
```

```
      rank suit
1      2 Club
2      3 Club

[1] 0.0007541478
```

Now that we have the probability space `N` we are ready to do some probability. We use the `Prob` function, just like before. The only trick is to specify the event of interest correctly, and recall that we were interested in $\mathbb{P}(\text{both Aces})$. But if the cards are both Aces then the rank of both cards should be A, which sounds like a job for the `all` function:

```
Prob(N, all(rank == "A"))

[1] 0.004524887
```

Note that this value matches what we found in Example ??, above. We could calculate all sorts of probabilities at this point; we are limited only by the complexity of the event's computer representation.

Example 3.23. Consider an urn with 10 balls inside, 7 of which are red and 3 of which are green. Select 3 balls successively from the urn. Let $A = \{1^{\text{st}} \text{ ball is red}\}$, $B = \{2^{\text{nd}} \text{ ball is red}\}$, and $C = \{3^{\text{rd}} \text{ ball is red}\}$. Then

$$\mathbb{P}(\text{all 3 balls are red}) = \mathbb{P}(A \cap B \cap C) = \frac{7}{10} \cdot \frac{6}{9} \cdot \frac{5}{8} \approx 0.2917.$$

How to do it with R

Example ?? is similar to Example ??, but it is even easier. We need to set up an urn (vector `L`) to hold the balls, we sample from `L` to get the sample space (data frame `M`), and we associate a probability vector (column `probs`) with the outcomes (rows of `M`) of the sample space. The final result is a probability space (an ordinary data frame `N`).

It is easier for us this time because our urn is a vector instead of a `cards()` data frame. Before there were two dimensions of information associated with the outcomes (rank and suit) but presently we have only one dimension (color).


```
L <- rep(c("red", "green"), times = c(7, 3))
M <- urnsamples(L, size = 3, replace = FALSE, ordered = TRUE)
N <- probspace(M)
```

Now let us think about how to set up the event {all 3 balls are red}. Rows of N that satisfy this condition have $X1=="red"$ & $X2=="red"$ & $X3=="red"$, but there must be an easier way. Indeed, there is. The `isrep` function (short for “is repeated”) in the `prob` package was written for this purpose. The command `isrep(N, "red", 3)` will test each row of N to see whether the value `red` appears 3 times. The result is exactly what we need to define an event with the `Prob` function. Observe

```
Prob(N, isrep(N, "red", 3))
```

```
[1] 0.2916667
```

Note that this answer matches what we found in Example ???. Now let us try some other probability questions. What is the probability of getting two red’s?

```
Prob(N, isrep(N, "red", 2))
```

```
[1] 0.525
```

Note that the exact value is $21/40$; we will learn a quick way to compute this in Section ??. What is the probability of observing red, then green, then red?

```
Prob(N, isin(N, c("red", "green", "red"), ordered = TRUE))
```

```
[1] 0.175
```

Note that the exact value is $7/40$ (do it with the Multiplication Rule). What is the probability of observing red, green, and red, in no particular order?

```
Prob(N, isin(N, c("red", "green", "red")))
```

```
[1] 0.525
```

We already knew this. It is the probability of observing two red’s, above.

Example 3.24. Consider two urns, the first with 5 red balls and 3 green balls, and the second with 2 red balls and 6 green balls. Your friend randomly selects one ball from the first urn and transfers it to the second urn, without disclosing the color of the ball. You select one ball from the second urn. What is the probability that the selected ball is red?

Let $A = \{\text{transferred ball is red}\}$ and $B = \{\text{selected ball is red}\}$. Write

$$\begin{aligned} B &= S \cap B \\ &= (A \cup A^c) \cap B \\ &= (A \cap B) \cup (A^c \cap B) \end{aligned}$$

and notice that $A \cap B$ and $A^c \cap B$ are disjoint. Therefore

$$\begin{aligned}\mathbb{P}(B) &= \mathbb{P}(A \cap B) + \mathbb{P}(A^c \cap B) \\ &= \mathbb{P}(A)\mathbb{P}(B|A) + \mathbb{P}(A^c)\mathbb{P}(B|A^c) \\ &= \frac{5}{8} \cdot \frac{3}{9} + \frac{3}{8} \cdot \frac{2}{9} \\ &= \frac{21}{72}\end{aligned}$$

(which is $7/24$ in lowest terms).

Example 3.25. We saw the `RcmdrTestDrive` data set in Chapter ?? in which a two-way table of the smoking status versus the gender was

```
library(RcmdrPlugin.IPSUR)
data(RcmdrTestDrive)
.Table <- xtabs(~ smoking + gender, data = RcmdrTestDrive)
addmargins(.Table) # Table with marginal distributions
```

	gender		
smoking	Female	Male	Sum
Nonsmoker	61	75	136
Smoker	9	23	32
Sum	70	98	168

If one person were selected at random from the data set, then we see from the two-way table that $\mathbb{P}(\text{Female}) = 70/168$ and $\mathbb{P}(\text{Smoker}) = 32/168$. Now suppose that one of the subjects quits smoking, but we do not know the person's gender. If we now select one nonsmoker at random, what would be $\mathbb{P}(\text{Female})$? This example is just like the last example, but with different labels. Let $A = \{\text{the quitter is a female}\}$ and $B = \{\text{selected nonsmoker is a female}\}$. Write

$$\begin{aligned}B &= S \cap B \\ &= (A \cup A^c) \cap B \\ &= (A \cap B) \cup (A^c \cap B)\end{aligned}$$

and notice that $A \cap B$ and $A^c \cap B$ are disjoint. Therefore

$$\begin{aligned}
\mathbb{P}(B) &= \mathbb{P}(A \cap B) + \mathbb{P}(A^c \cap B), \\
&= \mathbb{P}(A)\mathbb{P}(B|A) + \mathbb{P}(A^c)\mathbb{P}(B|A^c), \\
&= \frac{9}{32} \cdot \frac{62}{137} + \frac{23}{32} \cdot \frac{76}{137}, \\
&= \frac{2306}{4384},
\end{aligned}$$

(which is 1153/2192 in lowest terms).

Using the same reasoning, we can return to the example from the beginning of the section and show that

$$\mathbb{P}(\{\text{second card is an Ace}\}) = 4/52.$$

3.7 Independent Events

Toss a coin twice. The sample space is $S = \{HH, HT, TH, TT\}$. We know that $\mathbb{P}(1^{\text{st}} \text{ toss is } H) = 2/4$, $\mathbb{P}(2^{\text{nd}} \text{ toss is } H) = 2/4$, and $\mathbb{P}(\text{both } H) = 1/4$. Then

$$\begin{aligned}
\mathbb{P}(2^{\text{nd}} \text{ toss is } H \mid 1^{\text{st}} \text{ toss is } H) &= \frac{\mathbb{P}(\text{both } H)}{\mathbb{P}(1^{\text{st}} \text{ toss is } H)}, \\
&= \frac{1/4}{2/4}, \\
&= \mathbb{P}(2^{\text{nd}} \text{ toss is } H).
\end{aligned}$$

Intuitively, this means that the information that the first toss is H has no bearing on the probability that the second toss is H . The coin does not remember the result of the first toss.

Definition 3.2. Events A and B are said to be *independent* if

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B). \quad (3.14)$$

Otherwise, the events are said to be *dependent*.

The connection with the above example stems from the following. We know from Section ?? that when $\mathbb{P}(B) > 0$ we may write

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}. \quad (3.15)$$

In the case that A and B are independent, the numerator of the fraction factors so that $\mathbb{P}(B)$ cancels with the result:

$$\mathbb{P}(A|B) = \mathbb{P}(A) \text{ when } A, B \text{ are independent.} \quad (3.16)$$

The interpretation in the case of independence is that the information that the event B occurred does not influence the probability of the event A occurring. Similarly, $\mathbb{P}(B|A) = \mathbb{P}(B)$, and so the occurrence of the event A likewise does not affect the probability of event B . It may seem more natural to define A and B to be independent when $\mathbb{P}(A|B) = \mathbb{P}(A)$; however, the conditional probability $\mathbb{P}(A|B)$ is only defined when $\mathbb{P}(B) > 0$. Our definition is not limited by this restriction. It can be shown that when $\mathbb{P}(A), \mathbb{P}(B) > 0$ the two notions of independence are equivalent.

Proposition 3.5. If the events A and B are independent then

- A and B^c are independent,
- A^c and B are independent,
- A^c and B^c are independent.

Proof. Suppose that A and B are independent. We will show the second one; the others are similar. We need to show that

$$\mathbb{P}(A^c \cap B) = \mathbb{P}(A^c)\mathbb{P}(B).$$

To this end, note that the Multiplication Rule, Equation (??) implies

$$\begin{aligned} \mathbb{P}(A^c \cap B) &= \mathbb{P}(B)\mathbb{P}(A^c|B), \\ &= \mathbb{P}(B)[1 - \mathbb{P}(A|B)], \\ &= \mathbb{P}(B)\mathbb{P}(A^c). \end{aligned}$$

□

Definition 3.3. The events A , B , and C are *mutually independent* if the following four conditions are met:

$$\begin{aligned} \mathbb{P}(A \cap B) &= \mathbb{P}(A)\mathbb{P}(B), \\ \mathbb{P}(A \cap C) &= \mathbb{P}(A)\mathbb{P}(C), \\ \mathbb{P}(B \cap C) &= \mathbb{P}(B)\mathbb{P}(C), \end{aligned}$$

and

$$\mathbb{P}(A \cap B \cap C) = \mathbb{P}(A)\mathbb{P}(B)\mathbb{P}(C).$$

If only the first three conditions hold then A , B , and C are said to be independent *pairwise*. Note that pairwise independence is not the same as mutual independence when the number of events is larger than two.

We can now deduce the pattern for n events, $n > 3$. The events will be mutually independent only if they satisfy the product equality pairwise, then in groups of three, in groups of four, and so forth, up to all n events at once. For n events, there will be $2^n - n - 1$ equations that must be satisfied (see Exercise ??). Although these requirements for a set of events to be mutually independent may seem stringent, the good news is that for most of the situations considered in this book the conditions will all be met (or at least we will suppose that they are).

Example 3.26. Toss ten coins. What is the probability of observing at least one Head?

Answer: Let $A_i = \{\text{the } i^{\text{th}} \text{ coin shows } H\}$, $i = 1, 2, \dots, 10$. Supposing that we toss the coins in such a way that they do not interfere with each other, this is one of the situations where all of the A_i may be considered mutually independent due to the nature of the tossing. Of course, the only way that there will not be at least one Head showing is if all tosses are Tails. Therefore,

$$\begin{aligned}\mathbb{P}(\text{at least one } H) &= 1 - \mathbb{P}(\text{all } T), \\ &= 1 - \mathbb{P}(A_1^c \cap A_2^c \cap \dots \cap A_{10}^c), \\ &= 1 - \mathbb{P}(A_1^c)\mathbb{P}(A_2^c) \dots \mathbb{P}(A_{10}^c), \\ &= 1 - \left(\frac{1}{2}\right)^{10},\end{aligned}$$

which is approximately 0.9990234.

3.7.1 How to do it with R

Example 3.27. Toss ten coins. What is the probability of observing at least one Head?

```
S <- tosscoin(10, makespace = TRUE)
A <- subset(S, isrep(S, vals = "T", nrep = 10))
1 - Prob(A)
```

```
[1] 0.9990234
```

Compare this answer to what we got in Example ??.

3.7.2 Independent, Repeated Experiments

Generalizing from above it is common to repeat a certain experiment multiple times under identical conditions and in an independent manner. We have seen many examples

of this already: tossing a coin repeatedly, rolling a die or dice, *etc.*

The `iidspace` function was designed specifically for this situation. It has three arguments: `x`, which is a vector of outcomes, `ntrials`, which is an integer telling how many times to repeat the experiment, and `probs` to specify the probabilities of the outcomes of `x` in a single trial.

Example 3.28 (An unbalanced coin). (Continued, see Example ??). It was easy enough to set up the probability space for one unbalanced toss, however, the situation becomes more complicated when there are many tosses involved. Clearly, the outcome *HHH* should not have the same probability as *TTT*, which should again not have the same probability as *HTH*. At the same time, there is symmetry in the experiment in that the coin does not remember the face it shows from toss to toss, and it is easy enough to toss the coin in a similar way repeatedly.

We may represent tossing our unbalanced coin three times with the following:

```
iidspace(c("H","T"), ntrials = 3, probs = c(0.7, 0.3))
```

	X1	X2	X3	probs
1	H	H	H	0.343
2	T	H	H	0.147
3	H	T	H	0.147
4	T	T	H	0.063
5	H	H	T	0.147
6	T	H	T	0.063
7	H	T	T	0.063
8	T	T	T	0.027

As expected, the outcome *HHH* has the largest probability, while *TTT* has the smallest. (Since the trials are independent, $\mathbb{P}(HHH) = 0.7^3$ and $\mathbb{P}(TTT) = 0.3^3$, *etc.*) Note that the result of the function call is a probability space, not a sample space (which we could construct already with the `tosscoin` or `urnsamples` functions). The same procedure could be used to model an unbalanced die or any other experiment that may be represented with a vector of possible outcomes.

Note that `iidspace` will assume `x` has equally likely outcomes if no `probs` argument is specified. Also note that the argument `x` is a *vector*, not a data frame. Something like `:iidspace(tosscoin(1),...)` would give an error.

3.8 Bayes' Rule

We mentioned the subjective view of probability in Section ?? . In this section we introduce a rule that allows us to update our probabilities when new information becomes available.

Theorem 3.2 (Bayes' Rule). Let B_1, B_2, \dots, B_n be mutually exclusive and exhaustive and let A be an event with $\mathbb{P}(A) > 0$. Then

$$\mathbb{P}(B_k|A) = \frac{\mathbb{P}(B_k)\mathbb{P}(A|B_k)}{\sum_{i=1}^n \mathbb{P}(B_i)\mathbb{P}(A|B_i)}, \quad k = 1, 2, \dots, n. \quad (3.17)$$

Proof. The proof follows from looking at $\mathbb{P}(B_k \cap A)$ in two different ways. For simplicity, suppose that $\mathbb{P}(B_k) > 0$ for all k . Then

$$\mathbb{P}(A)\mathbb{P}(B_k|A) = \mathbb{P}(B_k \cap A) = \mathbb{P}(B_k)\mathbb{P}(A|B_k).$$

Since $\mathbb{P}(A) > 0$ we may divide through to obtain

$$\mathbb{P}(B_k|A) = \frac{\mathbb{P}(B_k)\mathbb{P}(A|B_k)}{\mathbb{P}(A)}.$$

Now remembering that $\{B_k\}$ is a partition, the Theorem of Total Probability (Equation (??)) gives the denominator of the last expression to be

$$\mathbb{P}(A) = \sum_{k=1}^n \mathbb{P}(B_k \cap A) = \sum_{k=1}^n \mathbb{P}(B_k)\mathbb{P}(A|B_k).$$

□

What does it mean? Usually in applications we are given (or know) *a priori* probabilities $\mathbb{P}(B_k)$. We go out and collect some data, which we represent by the event A . We want to know: how do we *update* $\mathbb{P}(B_k)$ to $\mathbb{P}(B_k|A)$? The answer: Bayes' Rule.

Example 3.29 (Misfiling Assistants). In this problem, there are three assistants working at a company: Moe, Larry, and Curly. Their primary job duty is to file paperwork in the filing cabinet when papers become available. The three assistants have different work schedules:

Table 3.3: Misfiling assistants: workload.

	Moe	Larry	Curly
Workload	60%	30%	10%

That is, Moe works 60% of the time, Larry works 30% of the time, and Curly does the remaining 10%, and they file documents at approximately the same speed. Suppose a person were to select one of the documents from the cabinet at random. Let M be the event

$$M = \{\text{Moe filed the document}\}$$

and let L and C be the events that Larry and Curly, respectively, filed the document.

What are these events’ respective probabilities? In the absence of additional information, reasonable prior probabilities would just be

Table 3.4: Misfiling assistants: prior.

	Moe	Larry	Curly
Prior Probability	0.6	0.3	0.1

Now, the boss comes in one day, opens up the file cabinet, and selects a file at random. The boss discovers that the file has been misplaced. The boss is so angry at the mistake that (s)he threatens to fire the one who erred. The question is: who misplaced the file?

The boss decides to use probability to decide, and walks straight to the workload schedule. (S)he reasons that, since the three employees work at the same speed, the probability that a randomly selected file would have been filed by each one would be proportional to his workload. The boss notifies *Moe* that he has until the end of the day to empty his desk.

But Moe argues in his defense that the boss has ignored additional information. Moe’s likelihood of having misfiled a document is smaller than Larry’s and Curly’s, since he is a diligent worker who pays close attention to his work. Moe admits that he works longer than the others, but he doesn’t make as many mistakes as they do. Thus, Moe recommends that – before making a decision – the boss should update the probability (initially based on workload alone) to incorporate the likelihood of having observed a misfiled document.

And, as it turns out, the boss has information about Moe, Larry, and Curly’s filing accuracy in the past (due to historical performance evaluations). The performance information may be represented by the following table:

Table 3.5: Misfiling assistants: misfile rate.

	Moe	Larry	Curly
Misfile Rate	0.003	0.007	0.010

In other words, on the average, Moe misfiles 0.3% of the documents he is supposed to file. Notice that Moe was correct: he is the most accurate filer, followed by Larry, and lastly Curly. If the boss were to make a decision based only on the worker’s overall accuracy, then *Curly* should get the axe. But Curly hears this and interjects that he only works a short period during the day, and consequently makes mistakes only very rarely; there is only the tiniest chance that he misfiled this particular document.

The boss would like to use this updated information to update the probabilities for the three assistants, that is, (s)he wants to use the additional likelihood that the document

was misfiled to update his/her beliefs about the likely culprit. Let A be the event that a document is misfiled. What the boss would like to know are the three probabilities

$$\mathbb{P}(M|A), \mathbb{P}(L|A), \text{ and } \mathbb{P}(C|A).$$

We will show the calculation for $\mathbb{P}(M|A)$, the other two cases being similar. We use Bayes' Rule in the form

$$\mathbb{P}(M|A) = \frac{\mathbb{P}(M \cap A)}{\mathbb{P}(A)}.$$

Let's try to find $\mathbb{P}(M \cap A)$, which is just $\mathbb{P}(M) \cdot \mathbb{P}(A|M)$ by the Multiplication Rule. We already know $\mathbb{P}(M) = 0.6$ and $\mathbb{P}(A|M)$ is nothing more than Moe's misfile rate, given above to be $\mathbb{P}(A|M) = 0.003$. Thus, we compute

$$\mathbb{P}(M \cap A) = (0.6)(0.003) = 0.0018.$$

Using the same procedure we may calculate

$$\mathbb{P}(L \cap A) = 0.0021 \text{ and } \mathbb{P}(C \cap A) = 0.0010.$$

Now let's find the denominator, $\mathbb{P}(A)$. The key here is the notion that if a file is misplaced, then either Moe or Larry or Curly must have filed it; there is no one else around to do the misfiling. Further, these possibilities are mutually exclusive. We may use the Theorem of Total Probability (??) to write

$$\mathbb{P}(A) = \mathbb{P}(A \cap M) + \mathbb{P}(A \cap L) + \mathbb{P}(A \cap C).$$

Luckily, we have computed these above. Thus

$$\mathbb{P}(A) = 0.0018 + 0.0021 + 0.0010 = 0.0049.$$

Therefore, Bayes' Rule yields

$$\mathbb{P}(M|A) = \frac{0.0018}{0.0049} \approx 0.37.$$

This last quantity is called the posterior probability that Moe misfiled the document, since it incorporates the observed data that a randomly selected file was misplaced (which is governed by the misfile rate). We can use the same argument to calculate

Table 3.6: Misfiling assistants: posterior.

	Moe	Larry	Curly
Posterior Probability	≈ 0.37	≈ 0.43	≈ 0.20

The conclusion: **Larry** gets the axe. What is happening is an intricate interplay between

the time on the job and the misfile rate. It is not obvious who the winner (or in this case, loser) will be, and the statistician needs to consult Bayes' Rule to determine the best course of action.

Example 3.30. Suppose the boss gets a change of heart and does not fire anybody. But the next day (s)he randomly selects another file and again finds it to be misplaced. To decide whom to fire now, the boss would use the same procedure, with one small change. (S)he would not use the prior probabilities 60%, 30%, and 10%; those are old news. Instead, she would replace the prior probabilities with the posterior probabilities just calculated. After the math she will have new posterior probabilities, updated even more from the day before.

In this way, probabilities found by Bayes' rule are always on the cutting edge, always updated with respect to the best information available at the time.

3.8.1 How to do it with R

There are not any special functions for Bayes' Rule in the `prob` package [`prob`], but problems like the ones above are easy enough to do by hand.

Example 3.31 (Misfiling Assistants). Continued from Example ?? . We store the prior probabilities and the likelihoods in vectors and go to town.

```
prior <- c(0.6, 0.3, 0.1)
like <- c(0.003, 0.007, 0.010)
post <- prior # like
post / sum(post)
```

```
[1] 0.6 0.3 0.1
```

Compare these answers with what we got in Example ?? . We would replace `prior` with `post` in a future calculation. We could raise `like` to a power to see how the posterior is affected by future document mistakes. (Do you see why? Think back to Section ??.)

Example 3.32. Let us incorporate the posterior probability (`post`) information from the last example and suppose that the assistants misfile seven more documents. Using Bayes' Rule, what would the new posterior probabilities be?

```
newprior <- post
post <- newprior * like^7
post / sum(post)
```

```
[1] 0.001051126 0.197907584 0.801041290
```

We see that the individual with the highest probability of having misfiled all eight documents given the observed data is no longer Larry, but Curly.

There are two important points. First, we did not divide `post` by the sum of its entries until the very last step; we do not need to calculate it, and it will save us computing time to postpone normalization until absolutely necessary, namely, until we finally want to interpret them as probabilities.

Second, the reader might be wondering what the boss would get if (s)he skipped the intermediate step of calculating the posterior after only one misfiled document. What if she started from the *original* prior, then observed eight misfiled documents, and calculated the posterior? What would she get? It must be the same answer, of course.

```
fastpost <- prior * like^8
fastpost / sum(fastpost)
```

```
[1] 0.0003355044 0.1473949328 0.8522695627
```

Compare this to what we got in Example ??.

3.9 Random Variables

We already know about experiments, sample spaces, and events. In this section, we are interested in a *number* that is associated with the experiment. We conduct a random experiment E and after learning the outcome ω in S we calculate a number X . That is, to each outcome ω in the sample space we associate a number $X(\omega) = x$.

Definition 3.4. A *random variable* X is a function $X : S \rightarrow \mathbb{R}$ that associates to each outcome $\omega \in S$ exactly one number $X(\omega) = x$.

We usually denote random variables by uppercase letters such as X , Y , and Z , and we denote their observed values by lowercase letters x , y , and z . Just as S is the set of all possible outcomes of E , we call the set of all possible values of X the *support* of X and denote it by S_X .

Example 3.33. Let E be the experiment of flipping a coin twice. We have seen that the sample space is $S = \{HH, HT, TH, TT\}$. Now define the random variable X = the number of heads. That is, for example, $X(HH) = 2$, while $X(HT) = 1$. We may make a table of the possibilities:

{#tab-flip-coin-twice}

Table 3.7: Flipping a coin twice.

$\omega \in S$	HH	HT	TH	TT
$X(\omega) = x$	2	1	1	0

Taking a look at the second row of the table, we see that the support of X – the set of all

numbers that X assumes – would be $S_X = \{0, 1, 2\}$.

Example 3.34. Let E be the experiment of flipping a coin repeatedly until observing a Head. The sample space would be $S = \{H, TH, TTH, TTTH, \dots\}$. Now define the random variable Y = the number of Tails before the first head. Then the support of Y would be $S_Y = \{0, 1, 2, \dots\}$.

Example 3.35. Let E be the experiment of tossing a coin in the air, and define the random variable Z = the time (in seconds) until the coin hits the ground. In this case, the sample space is inconvenient to describe. Yet the support of Z would be $(0, \infty)$. Of course, it is reasonable to suppose that the coin will return to Earth in a short amount of time; in practice, the set $(0, \infty)$ is admittedly too large. However, we will find that in many circumstances it is mathematically convenient to study the extended set rather than a restricted one.

There are important differences between the supports of X , Y , and Z . The support of X is a finite collection of elements that can be inspected all at once. And while the support of Y cannot be exhaustively written down, its elements can nevertheless be listed in a naturally ordered sequence. Random variables with supports similar to those of X and Y are called *discrete random variables*. We study these in Chapter ??.

In contrast, the support of Z is a continuous interval, containing all rational and irrational positive real numbers. For this reason⁴, random variables with supports like Z are called *continuous random variables*, to be studied in Chapter ??.

3.9.1 How to do it with R

The primary vessel for this task is the `addrv` function. There are two ways to use it, and we will describe both.

Supply a Defining Formula

The first method is based on the `transform` function. See `?transform`. The idea is to write a formula defining the random variable inside the function, and it will be added as a column to the data frame. As an example, let us roll a 4-sided die three times, and let us define the random variable $U = X_1 - X_2 + X_3$.

```
S <- rolldie(3, nsides = 4, makespace = TRUE)
S <- addrv(S, U = X1-X2+X3)
```

Now let's take a look at the values of U . In the interest of space, we will only reproduce the first few rows of S (there are $4^3 = 64$ rows in total).

⁴This isn't really the reason, but it serves as an effective litmus test at the introductory level. See Billingsley [Billingsley1995] or Resnick [Resnick1999].

```
head(S)
```

	X1	X2	X3	U	probs
1	1	1	1	1	0.015625
2	2	1	1	2	0.015625
3	3	1	1	3	0.015625
4	4	1	1	4	0.015625
5	1	2	1	0	0.015625
6	2	2	1	1	0.015625

We see from the U column it is operating just like it should. We can now answer questions like

```
Prob(S, U > 6)
```

```
[1] 0.015625
```

Supply a Function

Sometimes we have a function laying around that we would like to apply to some of the outcome variables, but it is unfortunately tedious to write out the formula defining what the new variable would be. The `addrv` function has an argument `FUN` specifically for this case. Its value should be a legitimate function from R, such as `sum`, `mean`, `median`, and so forth. Or, you can define your own function. Continuing the previous example, let's define $V = \max(X1, X2, X3)$ and $W = X1 + X2 + X3$.

```
S <- addrv(S, FUN = max, invars = c("X1", "X2", "X3"), name = "V")
S <- addrv(S, FUN = sum, invars = c("X1", "X2", "X3"), name = "W")
head(S)
```

	X1	X2	X3	U	V	W	probs
1	1	1	1	1	1	3	0.015625
2	2	1	1	2	2	4	0.015625
3	3	1	1	3	3	5	0.015625
4	4	1	1	4	4	6	0.015625
5	1	2	1	0	2	4	0.015625
6	2	2	1	1	2	5	0.015625

Notice that `addrv` has an `invars` argument to specify exactly to which columns one would like to apply the function `FUN`. If no input variables are specified, then `addrv` will apply `FUN` to all non-probs columns. Further, `addrv` has an optional argument `name` to give the new variable; this can be useful when adding several random variables to a probability space (as above). If not specified, the default name is `X`.

3.9.2 Marginal Distributions

As we can see above, often after adding a random variable V to a probability space one will find that V has values that are repeated, so that it becomes difficult to understand what the ultimate behavior of V actually is. We can use the `marginal` function to aggregate the rows of the sample space by values of V , all the while accumulating the probability associated with V 's distinct values. Continuing our example from above, suppose we would like to focus entirely on the values and probabilities of $V = \max(X1, X2, X3)$.

```
marginal(S, vars = "V")
```

	V	probs
1	1	0.015625
2	2	0.109375
3	3	0.296875
4	4	0.578125

We could save the probability space of V in a data frame and study it further, if we wish. As a final remark, we can calculate the marginal distributions of multiple variables desired using the `vars` argument. For example, suppose we would like to examine the joint distribution of V and W .

```
marginal(S, vars = c("V", "W"))
```

	V	W	probs
1	1	3	0.015625
2	2	4	0.046875
3	2	5	0.046875
4	3	5	0.046875
5	2	6	0.015625
6	3	6	0.093750
7	4	6	0.046875
8	3	7	0.093750
9	4	7	0.093750
10	3	8	0.046875
11	4	8	0.140625
12	3	9	0.015625
13	4	9	0.140625
14	4	10	0.093750
15	4	11	0.046875
16	4	12	0.015625

Note that the default value of `vars` is the names of all columns except `probs`. This can be useful if there are duplicated rows in the probability space.

3.10 Exercises

Exercise 3.1. Prove the assertion given in the text: the number of conditions that the events A_1, A_2, \dots, A_n must satisfy in order to be mutually independent is $2^n - n - 1$. (*Hint:* think about Pascal's triangle.)

Chapter 4

Multiple Linear Regression

We know a lot about simple linear regression models, and a next step is to study multiple regression models that have more than one independent (explanatory) variable. In the discussion that follows we will assume that we have p explanatory variables, where $p > 1$.

The language is phrased in matrix terms – for two reasons. First, it is quicker to write and (arguably) more pleasant to read. Second, the matrix approach will be required for later study of the subject; the reader might as well be introduced to it now.

Most of the results are stated without proof or with only a cursory justification. Those yearning for more should consult an advanced text in linear regression for details, such as *Applied Linear Regression Models* [Neter1996] or *Linear Models: Least Squares and Alternatives* [Rao1999].

What do I want them to know?

- the basic MLR model, and how it relates to the SLR
- how to estimate the parameters and use those estimates to make predictions
- basic strategies to determine whether or not the model is doing a good job
- a few thoughts about selected applications of the MLR, such as polynomial, interaction, and dummy variable models
- some of the uses of residuals to diagnose problems
- hints about what will be coming later

4.1 The Multiple Linear Regression Model

The first thing to do is get some better notation. We will write

$$\mathbf{Y}_{n \times 1} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \text{and} \quad \mathbf{X}_{n \times (p+1)} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{p1} \\ 1 & x_{12} & x_{22} & \cdots & x_{p2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{pn} \end{bmatrix}. \quad (4.1)$$

The vector \mathbf{Y} is called the *response vector* and the matrix \mathbf{X} is called the *model matrix*. As in Chapter ??, the most general assumption that relates \mathbf{Y} to \mathbf{X} is

$$\mathbf{Y} = \mu(\mathbf{X}) + \epsilon, \quad (4.2)$$

where μ is some function (the *signal*) and ϵ is the *noise* (everything else). We usually impose some structure on μ and ϵ . In particular, the standard multiple linear regression model assumes

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon, \quad (4.3)$$

where the parameter vector β looks like

$$\beta_{(p+1) \times 1} = [\beta_0 \quad \beta_1 \quad \cdots \quad \beta_p]^T, \quad (4.4)$$

and the random vector $\epsilon_{n \times 1} = [\epsilon_1 \quad \epsilon_2 \quad \cdots \quad \epsilon_n]^T$ is assumed to be distributed

$$\epsilon \sim \text{mvnrm}(\text{mean} = \mathbf{0}_{n \times 1}, \text{sigma} = \sigma^2 \mathbf{I}_{n \times n}). \quad (4.5)$$

The assumption on ϵ is equivalent to the assumption that $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are IID $\text{norm}(\text{mean} = 0, \text{sd} = \sigma)$. It is a linear model because the quantity $\mu(\mathbf{X}) = \mathbf{X}\beta$ is linear in the parameters $\beta_0, \beta_1, \dots, \beta_p$. It may be helpful to see the model in expanded form; the above matrix formulation is equivalent to the more lengthy

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_p x_{pi} + \epsilon_i, \quad i = 1, 2, \dots, n. \quad (4.6)$$

Example 4.1 (Girth, Height, and Volume for Black Cherry trees). Measurements were made of the girth, height, and volume of timber in 31 felled black cherry trees. Note that girth is the diameter of the tree (in inches) measured at 4 ft 6 in above the ground. The variables are

1. **Girth:** tree diameter in inches (denoted x_1)
2. **Height:** tree height in feet (x_2).
3. **Volume:** volume of the tree in cubic feet. (y)

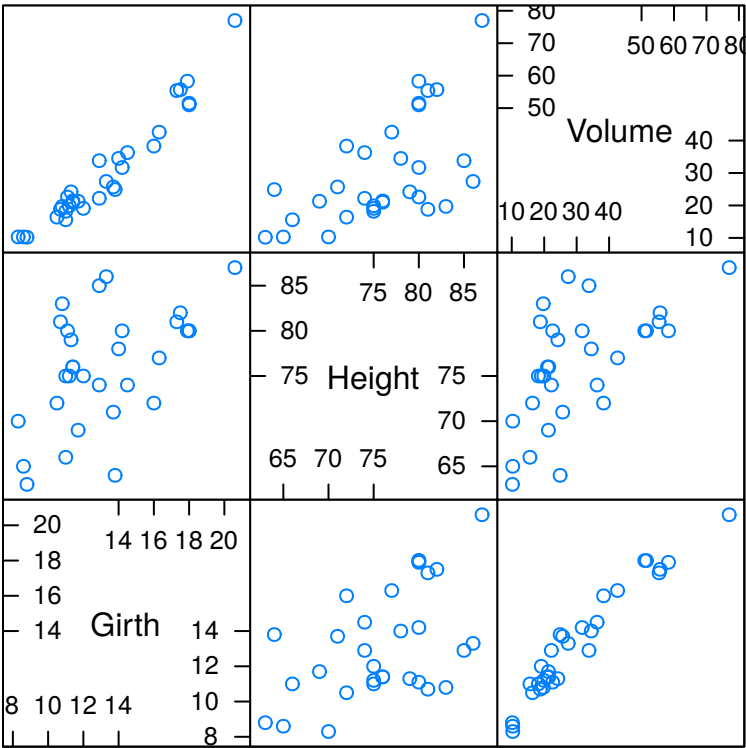
The data are in the `datasets` package [`datasets`] and are already on the search path; they can be viewed with

```
head(trees)
```

	Girth	Height	Volume
1	8.3	70	10.3
2	8.6	65	10.3
3	8.8	63	10.2
4	10.5	72	16.4
5	10.7	81	18.8
6	10.8	83	19.7

Let us take a look at a visual display of the data. For multiple variables, instead of a simple scatterplot we use a scatterplot matrix which is made with the `splo`m function in the `lattice` package [`lattice`] as shown below. The plot is shown in Figure ??.

```
splom(trees)
```



Scatter Plot Matrix

Figure 4.1: A scatterplot matrix of the `trees` data.

The dependent (response) variable `Volume` is listed in the first row of the scatterplot matrix. Moving from left to right, we see an approximately linear relationship between `Volume` and the independent (explanatory) variables `Height` and `Girth`. A first guess at a model for these data might be

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon, \quad (4.7)$$

in which case the quantity $\mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ would represent the mean value of Y at the point (x_1, x_2) .

4.1.1 What does it mean?

The interpretation is simple. The intercept β_0 represents the mean `Volume` when all other independent variables are zero. The parameter β_i represents the change in mean `Volume` when there is a unit increase in x_i , while the other independent variable is held constant. For the `trees` data, β_1 represents the change in average `Volume` as `Girth` increases by one unit when the `Height` is held constant, and β_2 represents the change in average `Volume` as `Height` increases by one unit when the `Girth` is held constant.

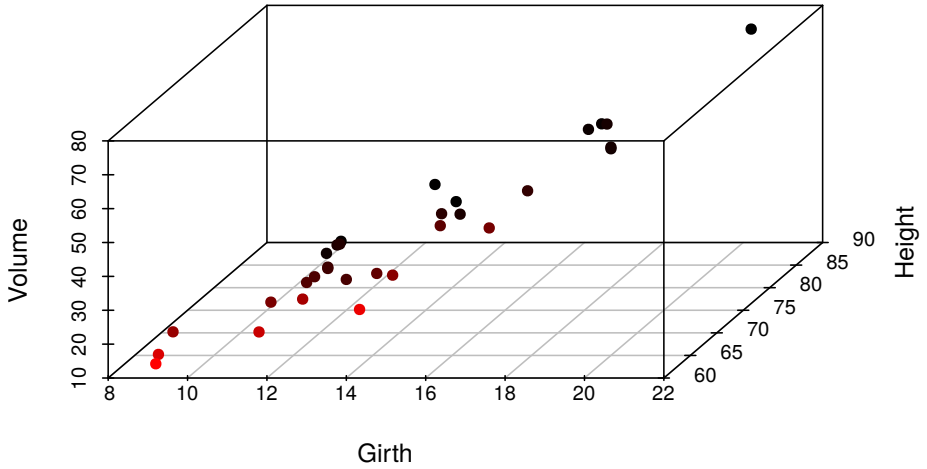
In simple linear regression, we had one independent variable and our linear regression surface was 1D, simply a line. In multiple regression there are many independent variables and so our linear regression surface will be many-D... in general, a hyperplane. But when there are only two explanatory variables the hyperplane is just an ordinary plane and we can look at it with a 3D scatterplot.

One way to do this is with the R Commander in the `Rcmdr` package [**Rcmdr**]. It has a 3D scatterplot option under the `Graphs` menu. It is especially great because the resulting graph is dynamic; it can be moved around with the mouse, zoomed, *etc.* But that particular display does not translate well to a printed book.

Another way to do it is with the `scatterplot3d` function in the `scatterplot3d` package [**scatterplot3d**]. The code follows, and the result is shown in Figure ??.

```
s3d <- with(trees, scatterplot3d(Girth, Height, Volume,
                                pch = 16, highlight.3d = TRUE, angle = 60))
```

Looking at the graph we see that the data points fall close to a plane in three dimensional space. (The plot looks remarkably good. In the author's experience it is rare to see points fit so well to the plane without some additional work.)

Figure 4.2: A 3D scatterplot with regression plane for the *trees* data.

4.2 Estimation and Prediction

4.2.1 Parameter estimates

We will proceed exactly like we did in Section ???. We know

$$\epsilon \sim \text{mvnorm}(\text{mean} = \mathbf{0}_{n \times 1}, \text{sigma} = \sigma^2 \mathbf{I}_{n \times n}), \quad (4.8)$$

which means that $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ has an $\text{mvnorm}(\text{mean} = \mathbf{X}\beta, \text{sigma} = \sigma^2 \mathbf{I}_{n \times n})$ distribution. Therefore, the likelihood function is

$$L(\beta, \sigma) = \frac{1}{2\pi^{n/2}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta)\right\}. \quad (4.9)$$

To *maximize* the likelihood in β , we need to *minimize* the quantity $g(\beta) = (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta)$. We do this by differentiating g with respect to β . (It may be a good idea to brush up on the material in Appendices ??? and ???.) First we will rewrite g :

$$g(\beta) = \mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{X}\beta - \beta^T \mathbf{X}^T \mathbf{Y} + \beta^T \mathbf{X}^T \mathbf{X}\beta, \quad (4.10)$$

which can be further simplified to $g(\beta) = \mathbf{Y}^T \mathbf{Y} - 2\beta^T \mathbf{X}^T \mathbf{Y} + \beta^T \mathbf{X}^T \mathbf{X}\beta$ since $\beta^T \mathbf{X}^T \mathbf{Y}$ is 1×1 and thus equal to its transpose. Now we differentiate to get

$$\frac{\partial g}{\partial \beta} = \mathbf{0} - 2\mathbf{X}^T \mathbf{Y} + 2\mathbf{X}^T \mathbf{X}\beta, \quad (4.11)$$

since $\mathbf{X}^T\mathbf{X}$ is symmetric. Setting the derivative equal to the zero vector yields the so called “normal equations”

$$\mathbf{X}^T\mathbf{X}\beta = \mathbf{X}^T\mathbf{Y}. \quad (4.12)$$

In the case that $\mathbf{X}^T\mathbf{X}$ is invertible¹, we may solve the equation for β to get the maximum likelihood estimator of β which we denote by \mathbf{b} :

$$\mathbf{b} = (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T\mathbf{Y}. \quad (4.13)$$

Remark 4.1. The formula in Equation (??) is convenient for mathematical study but is inconvenient for numerical computation. Researchers have devised much more efficient algorithms for the actual calculation of the parameter estimates, and we do not explore them here.

Remark 4.2. We have only found a critical value, and have not actually shown that the critical value is a minimum. We omit the details and refer the interested reader to [Rao1999].

How to do it with R

We do all of the above just as we would in simple linear regression. The powerhouse is the `lm` function. Everything else is based on it. We separate explanatory variables in the model formula by a plus sign.

```
trees.lm <- lm(Volume ~ Girth + Height, data = trees)
trees.lm
```

Call:

```
lm(formula = Volume ~ Girth + Height, data = trees)
```

Coefficients:

(Intercept)	Girth	Height
-57.9877	4.7082	0.3393

We see from the output that for the `trees` data our parameter estimates are

$$\mathbf{b} = \begin{bmatrix} -58.0 & 4.7 & 0.3 \end{bmatrix},$$

and consequently our estimate of the mean response is $\hat{\mu}$ given by

¹We can find solutions of the normal equations even when $\mathbf{X}^T\mathbf{X}$ is not of full rank, but the topic falls outside the scope of this book. The interested reader can consult an advanced text such as Rao [Rao1999].

$$\hat{\mu}(x_1, x_2) = b_0 + b_1 x_1 + b_2 x_2, \quad (4.14)$$

$$\approx -58.0 + 4.7x_1 + 0.3x_2. \quad (4.15)$$

We could see the entire model matrix \mathbf{X} with the `model.matrix` function, but in the interest of brevity we only show the first few rows.

```
head(model.matrix(trees.lm))
```

	(Intercept)	Girth	Height
1	1	8.3	70
2	1	8.6	65
3	1	8.8	63
4	1	10.5	72
5	1	10.7	81
6	1	10.8	83

4.2.2 Point Estimates of the Regression Surface

The parameter estimates \mathbf{b} make it easy to find the fitted values, $\hat{\mathbf{Y}}$. We write them individually as \hat{Y}_i , $i = 1, 2, \dots, n$, and recall that they are defined by

$$\hat{Y}_i = \hat{\mu}(x_{1i}, x_{2i}), \quad (4.16)$$

$$= b_0 + b_1 x_{1i} + b_2 x_{2i}, \quad i = 1, 2, \dots, n. \quad (4.17)$$

They are expressed more compactly by the matrix equation

$$\hat{\mathbf{Y}} = \mathbf{X}\mathbf{b}. \quad (4.18)$$

From Equation (??) we know that $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$, so we can rewrite

$$\hat{\mathbf{Y}} = \mathbf{X} \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \right], \quad (4.19)$$

$$= \mathbf{H}\mathbf{Y}, \quad (4.20)$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is appropriately named *the hat matrix* because it “puts the hat on \mathbf{Y} ”. The hat matrix is very important in later courses. Some facts about \mathbf{H} are

- \mathbf{H} is a symmetric square matrix, of dimension $n \times n$.
- The diagonal entries h_{ii} satisfy $0 \leq h_{ii} \leq 1$ (compare to Equation (??)).

- The trace is $\text{tr}(\mathbf{H}) = p$.
- \mathbf{H} is *idempotent* (also known as a *projection matrix*) which means that $\mathbf{H}^2 = \mathbf{H}$. The same is true of $\mathbf{I} - \mathbf{H}$.

Now let us write a column vector $\mathbf{x}_0 = (x_{10}, x_{20})^T$ to denote given values of the explanatory variables `Girth` == x_{10} and `Height` == x_{20} . These values may match those of the collected data, or they may be completely new values not observed in the original data set. We may use the parameter estimates to find $\hat{Y}(\mathbf{x}_0)$, which will give us 1. an estimate of $\mu(\mathbf{x}_0)$, the mean value of a future observation at \mathbf{x}_0 , and 2. a prediction for $Y(\mathbf{x}_0)$, the actual value of a future observation at \mathbf{x}_0 .

We can represent $\hat{Y}(\mathbf{x}_0)$ by the matrix equation

$$\hat{Y}(\mathbf{x}_0) = \mathbf{x}_0^T \mathbf{b}, \quad (4.21)$$

which is just a fancy way to write

$$\hat{Y}(x_{10}, x_{20}) = b_0 + b_1 x_{10} + b_2 x_{20}. \quad (4.22)$$

Example 4.2. If we wanted to predict the average volume of black cherry trees that have `Girth` = 15 in and are `Height` = 77 ft tall then we would use the estimate

$$\begin{aligned} \hat{\mu}(15, 77) &= -58 + 4.7(15) + 0.3(77), \\ &\approx 35.6 \text{ ft}^3. \end{aligned}$$

We would use the same estimate $\hat{Y} = 35.6$ to predict the measured Volume of another black cherry tree – yet to be observed – that has `Girth` = 15 in and is `Height` = 77 ft tall.

How to do it with R

The fitted values are stored inside `trees.lm` and may be accessed with the `fitted` function. We only show the first five fitted values.

```
fitted(trees.lm)[1:5]
```

```
      1      2      3      4      5
4.837660 4.553852 4.816981 15.874115 19.869008
```

The syntax for general prediction does not change much from simple linear regression. The computations are done with the `predict` function as described below.

The only difference from SLR is in the way we tell R the values of the explanatory variables for which we want predictions. In SLR we had only one independent variable

but in MLR we have many (for the `trees` data we have two). We will store values for the independent variables in the data frame `new`, which has two columns (one for each independent variable) and three rows (we shall make predictions at three different locations).

```
new <- data.frame(Girth = c(9.1, 11.6, 12.5),
                  Height = c(69, 74, 87))
```

We can view the locations at which we will predict:

```
new
```

```
   Girth Height
1    9.1     69
2   11.6     74
3   12.5     87
```

We continue just like we would have done in SLR.

```
predict(trees.lm, newdata = new)
```

```
      1      2      3
8.264937 21.731594 30.379205
```

Example. Using the `trees` data,

1. Report a point estimate of the mean `Volume` of a tree of `Girth` 9.1 in and `Height` 69 ft. The fitted value for $x_1 = 9.1$ and $x_2 = 69$ is 8.3, so a point estimate would be 8.3 cubic feet.
2. Report a point prediction for and a 95% prediction interval for the `Volume` of a hypothetical tree of `Girth` 12.5 in and `Height` 87 ft. The fitted value for $x_1 = 12.5$ and $x_2 = 87$ is 30.4, so a point prediction for the `Volume` is 30.4 cubic feet.

4.2.3 Mean Square Error and Standard Error

The residuals are given by

$$\mathbf{E} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{H}\mathbf{Y} = (\mathbf{I} - \mathbf{H})\mathbf{Y}. \quad (4.23)$$

Now we can use Theorem ?? to see that the residuals are distributed

$$\mathbf{E} \sim \text{mvnorm}(\text{mean} = \mathbf{0}, \text{sigma} = \sigma^2(\mathbf{I} - \mathbf{H})), \quad (4.24)$$

since $(\mathbf{I} - \mathbf{H})\mathbf{X}\beta = \mathbf{X}\beta - \mathbf{X}\beta = \mathbf{0}$ and $(\mathbf{I} - \mathbf{H})(\sigma^2\mathbf{I})(\mathbf{I} - \mathbf{H})^T = \sigma^2(\mathbf{I} - \mathbf{H})^2 = \sigma^2(\mathbf{I} - \mathbf{H})$. The sum of squared errors SSE is just

$$SSE = \mathbf{E}^T \mathbf{E} = \mathbf{Y}^T (\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) \mathbf{Y} = \mathbf{Y}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y}. \quad (4.25)$$

Recall that in SLR we had two parameters (β_0 and β_1) in our regression model and we estimated σ^2 with $s^2 = SSE/(n - 2)$. In MLR, we have $p + 1$ parameters in our regression model and we might guess that to estimate σ^2 we would use the *mean square error* S^2 defined by

$$S^2 = \frac{SSE}{n - (p + 1)}. \quad (4.26)$$

That would be a good guess. The *residual standard error* is $S = \sqrt{S^2}$.

How to do it with R

The residuals are also stored with `trees.lm` and may be accessed with the `residuals` function. We only show the first five residuals.

```
residuals(trees.lm)[1:5]
```

```
      1      2      3      4      5
5.4623403 5.7461484 5.3830187 0.5258848 -1.0690084
```

The `summary` function output (shown later) lists the Residual Standard Error which is just $S = \sqrt{S^2}$. It is stored in the `sigma` component of the `summary` object.

```
treesumry <- summary(trees.lm)
treesumry$sigma
```

```
[1] 3.881832
```

For the `trees` data we find $s \approx 3.882$.

4.2.4 Interval Estimates of the Parameters

We showed in Section ?? that $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$, which is really just a big matrix – namely $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ – multiplied by \mathbf{Y} . It stands to reason that the sampling distribution of \mathbf{b} would be intimately related to the distribution of \mathbf{Y} , which we assumed to be

$$\mathbf{Y} \sim \text{mvn}(\text{mean} = \mathbf{X}\beta, \text{sigma} = \sigma^2 \mathbf{I}). \quad (4.27)$$

Now recall Theorem ?? that we said we were going to need eventually (the time is now). That proposition guarantees that

$$\mathbf{b} \sim \text{mvnrm}\left(\text{mean} = \beta, \text{sigma} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\right), \quad (4.28)$$

since

$$\mathbb{E}\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta) = \beta, \quad (4.29)$$

and

$$\text{Var}(\mathbf{b}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\sigma^2 \mathbf{I}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}, \quad (4.30)$$

the first equality following because the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$ is symmetric.

There is a lot that we can glean from Equation (??). First, it follows that the estimator \mathbf{b} is unbiased (see Section ??). Second, the variances of b_0, b_1, \dots, b_n are exactly the diagonal elements of $\sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$, which is completely known except for that pesky parameter σ^2 . Third, we can estimate the standard error of b_i (denoted S_{b_i}) with the mean square error S (defined in the previous section) multiplied by the corresponding diagonal element of $(\mathbf{X}^T \mathbf{X})^{-1}$. Finally, given estimates of the standard errors we may construct confidence intervals for β_i with an interval that looks like

$$b_i \pm t_{\alpha/2}(\text{df} = n - p - 1) S_{b_i}. \quad (4.31)$$

The degrees of freedom for the Student's t distribution² are the same as the denominator of S^2 .

How to do it with R

To get confidence intervals for the parameters we need only use `confint` :

```
confint(trees.lm)
```

	2.5 %	97.5 %
(Intercept)	-75.68226247	-40.2930554
Girth	4.16683899	5.2494820
Height	0.07264863	0.6058538

For example, using the calculations above we say that for the regression model `Volume ~ Girth + Height` we are 95% confident that the parameter β_1 lies somewhere in the interval 4.2 to 5.2.

²We are taking great leaps over the mathematical details. In particular, we have yet to show that s^2 has a chi-square distribution and we have not even come close to showing that b_i and s_{b_i} are independent. But these are entirely outside the scope of the present book and the reader may rest assured that the proofs await in later classes. See C.R. Rao for more.

4.2.5 Confidence and Prediction Intervals

We saw in Section ?? how to make point estimates of the mean value of additional observations and predict values of future observations, but how good are our estimates? We need confidence and prediction intervals to gauge their accuracy, and lucky for us the formulas look similar to the ones we saw in SLR.

In Equation (??) we wrote $\hat{Y}(\mathbf{x}_0) = \mathbf{x}_0^T \mathbf{b}$, and in Equation (??) we saw that

$$\mathbf{b} \sim \text{mvnorm}(\text{mean} = \beta, \text{sigma} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}). \quad (4.32)$$

The following is therefore immediate from Theorem ??:

$$\hat{Y}(\mathbf{x}_0) \sim \text{mvnorm}(\text{mean} = \mathbf{x}_0^T \beta, \text{sigma} = \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0). \quad (4.33)$$

It should be no surprise that confidence intervals for the mean value of a future observation at the location $\mathbf{x}_0 = [x_{10} \ x_{20} \ \dots \ x_{p0}]^T$ are given by

$$\hat{Y}(\mathbf{x}_0) \pm t_{\alpha/2}(\text{df} = n - p - 1) S \sqrt{\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}. \quad (4.34)$$

Intuitively, $\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0$ measures the distance of \mathbf{x}_0 from the center of the data. The degrees of freedom in the Student's t critical value are $n - (p + 1)$ because we need to estimate $p + 1$ parameters.

Prediction intervals for a new observation at \mathbf{x}_0 are given by

$$\hat{Y}(\mathbf{x}_0) \pm t_{\alpha/2}(\text{df} = n - p - 1) S \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}. \quad (4.35)$$

The prediction intervals are wider than the confidence intervals, just as in Section ??.

How to do it with R

The syntax is identical to that used in SLR, with the proviso that we need to specify values of the independent variables in the data frame `new` as we did in Section ?? (which we repeat here for illustration).

```
new <- data.frame(Girth = c(9.1, 11.6, 12.5),
                  Height = c(69, 74, 87))
```

Confidence intervals are given by

```
predict(trees.lm, newdata = new, interval = "confidence")
```

	fit	lwr	upr
1	8.264937	5.77240	10.75747
2	21.731594	20.11110	23.35208
3	30.379205	26.90964	33.84877

Prediction intervals are given by

```
predict(trees.lm, newdata = new, interval = "prediction")
```

	fit	lwr	upr
1	8.264937	-0.06814444	16.59802
2	21.731594	13.61657775	29.84661
3	30.379205	21.70364103	39.05477

As before, the interval type is decided by the `interval` argument and the default confidence level is 95% (which can be changed with the `level` argument).

Example. Using the `trees` data,

1. Report a 95% confidence interval for the mean `Volume` of a tree of `Girth` 9.1 in and `Height` 69 ft. The 95% CI is given by 5.8 to 10.8, so with 95% confidence the mean `Volume` lies somewhere between 5.8 cubic feet and 10.8 cubic feet.
2. Report a 95% prediction interval for the `Volume` of a hypothetical tree of `Girth` 12.5 in and `Height` 87 ft. The 95% prediction interval is given by 26.9 to 33.8, so with 95% confidence we may assert that the hypothetical `Volume` of a tree of `Girth` 12.5 in and `Height` 87 ft would lie somewhere between 26.9 cubic feet and 33.8 feet.

4.3 Model Utility and Inference

4.3.1 Multiple Coefficient of Determination

We saw in Section ?? that the error sum of squares SSE can be conveniently written in MLR as

$$SSE = \mathbf{Y}^T(\mathbf{I} - \mathbf{H})\mathbf{Y}. \quad (4.36)$$

It turns out that there are equally convenient formulas for the total sum of squares $SSTO$ and the regression sum of squares SSR . They are:

$$SSTO = \mathbf{Y}^T \left(\mathbf{I} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y} \quad (4.37)$$

and

$$SSR = \mathbf{Y}^T \left(\mathbf{H} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y}. \quad (4.38)$$

(The matrix \mathbf{J} is defined in Appendix ??.) Immediately from Equations (??), (??), and (??) we get the *Anova Equality*

$$SSTO = SSE + SSR. \quad (4.39)$$

(See Exercise ??.) We define the *multiple coefficient of determination* by the formula

$$R^2 = 1 - \frac{SSE}{SSTO}. \quad (4.40)$$

We interpret R^2 as the proportion of total variation that is explained by the multiple regression model. In MLR we must be careful, however, because the value of R^2 can be artificially inflated by the addition of explanatory variables to the model, regardless of whether or not the added variables are useful with respect to prediction of the response variable. In fact, it can be proved that the addition of a single explanatory variable to a regression model will increase the value of R^2 , *no matter how worthless* the explanatory variable is. We could model the height of the ocean tides, then add a variable for the length of cheetah tongues on the Serengeti plain, and our R^2 would inevitably increase.

This is a problem, because as the philosopher, Occam, once said: “causes should not be multiplied beyond necessity”. We address the problem by penalizing R^2 when parameters are added to the model. The result is an *adjusted* R^2 which we denote by \bar{R}^2 .

$$\bar{R}^2 = \left(R^2 - \frac{p}{n-1} \right) \left(\frac{n-1}{n-p-1} \right). \quad (4.41)$$

It is good practice for the statistician to weigh both R^2 and \bar{R}^2 during assessment of model utility. In many cases their values will be very close to each other. If their values differ substantially, or if one changes dramatically when an explanatory variable is added, then (s)he should take a closer look at the explanatory variables in the model.

How to do it with R

For the `trees` data, we can get R^2 and \bar{R}^2 from the `summary` output or access the values directly by name as shown (recall that we stored the `summary` object in `treesumry`).

```
treesumry$r.squared
```

```
[1] 0.94795
```

```
treesumry$adj.r.squared
```

```
[1] 0.9442322
```

High values of R^2 and \bar{R}^2 such as these indicate that the model fits very well, which agrees with what we saw in Figure ??.

4.3.2 Overall F-Test

Another way to assess the model's utility is to test the hypothesis

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0 \text{ versus } H_1 : \text{ at least one } \beta_i \neq 0.$$

The idea is that if all β_i 's were zero, then the explanatory variables X_1, \dots, X_p would be worthless predictors for the response variable Y . We can test the above hypothesis with the overall F statistic, which in MLR is defined by

$$F = \frac{SSR/p}{SSE/(n - p - 1)}. \tag{4.42}$$

When the regression assumptions hold and under H_0 , it can be shown that $F \sim f(\text{df1} = p, \text{df2} = n - p - 1)$. We reject H_0 when F is large, that is, when the explained variation is large relative to the unexplained variation.

How to do it with R

The overall F statistic and its associated p -value is listed at the bottom of the `summary` output, or we can access it directly by name; it is stored in the `fstatistic` component of the `summary` object.

```
treesumry$fstatistic
```

```
      value      numdf      dendif
254.9723    2.0000    28.0000
```

For the `trees` data, we see that ($F =$) 254.9723374 with a p -value $=< 2.2\text{e-}16$. Consequently we reject H_0 , that is, the data provide strong evidence that not all β_i 's are zero.

4.3.3 Student's t Tests

We know that

$$\mathbf{b} \sim \text{mvnorm}\left(\text{mean} = \beta, \text{sigma} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\right) \quad (4.43)$$

and we have seen how to test the hypothesis $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$, but let us now consider the test

$$H_0 : \beta_i = 0 \text{ versus } H_1 : \beta_i \neq 0, \quad (4.44)$$

where β_i is the coefficient for the i^{th} independent variable. We test the hypothesis by calculating a statistic, examining its null distribution, and rejecting H_0 if the p -value is small. If H_0 is rejected, then we conclude that there is a significant relationship between Y and x_i in the regression model $Y \sim (x_1, \dots, x_p)$. This last part of the sentence is very important because the significance of the variable x_i sometimes depends on the presence of other independent variables in the model³.

To test the hypothesis we go to find the sampling distribution of (b_i) , the estimator of the corresponding parameter β_i , when the null hypothesis is true. We saw in Section ?? that

$$T_i = \frac{b_i - \beta_i}{S_{b_i}} \quad (4.45)$$

has a Student's t distribution with $n - (p + 1)$ degrees of freedom. (Remember, we are estimating $p + 1$ parameters.) Consequently, under the null hypothesis $H_0 : \beta_i = 0$ the statistic $t_i = b_i / S_{b_i}$ has a $t(\text{df} = n - p - 1)$ distribution.

How to do it with R

The Student's t tests for significance of the individual explanatory variables are shown in the summary output.

```
treesumry
```

Call:

```
lm(formula = Volume ~ Girth + Height, data = trees)
```

Residuals:

Min	1Q	Median	3Q	Max
-----	----	--------	----	-----

³In other words, a variable might be highly significant one moment but then fail to be significant when another variable is added to the model. When this happens it often indicates a problem with the explanatory variables, such as *multicollinearity*. See Section ??.

-6.4065 -2.6493 -0.2876 2.2003 8.4847

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-57.9877	8.6382	-6.713	2.75e-07	***
Girth	4.7082	0.2643	17.816	< 2e-16	***
Height	0.3393	0.1302	2.607	0.0145	*

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.882 on 28 degrees of freedom

Multiple R-squared: 0.948, Adjusted R-squared: 0.9442

F-statistic: 255 on 2 and 28 DF, p-value: < 2.2e-16

We see from the *p-values* that there is a significant linear relationship between Volume and Girth and between Volume and Height in the regression model `Volume ~ Girth + Height`. Further, it appears that the Intercept is significant in the aforementioned model.

4.4 Polynomial Regression

4.4.1 Quadratic Regression Model

In each of the previous sections we assumed that μ was a linear function of the explanatory variables. For example, in SLR we assumed that $\mu(x) = \beta_0 + \beta_1 x$, and in our previous MLR examples we assumed $\mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$. In every case the scatterplots indicated that our assumption was reasonable. Sometimes, however, plots of the data suggest that the linear model is incomplete and should be modified.

```
plot(Volume ~ Girth, data = trees)
```

For example, let us examine a scatterplot of Volume versus Girth a little more closely. See Figure ???. There might be a slight curvature to the data; the volume curves ever so slightly upward as the girth increases. After looking at the plot we might try to capture the curvature with a mean response such as

$$\mu(x_1) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2. \quad (4.46)$$

The model associated with this choice of μ is

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \epsilon. \quad (4.47)$$

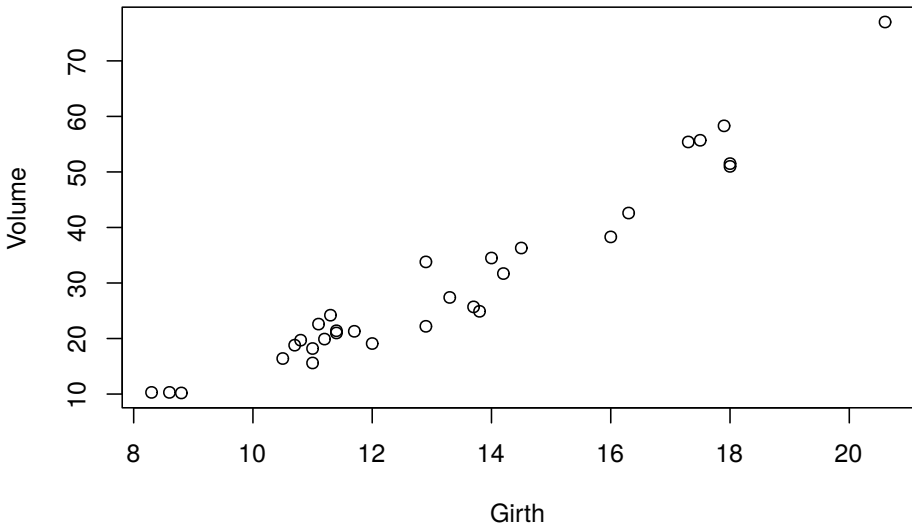


Figure 4.3: Scatterplot of Volume versus Girth for the trees data.

The regression assumptions are the same. Almost everything indeed is the same. In fact, it is still called a “linear regression model”, since the mean response μ is linear in the parameters β_0 , β_1 , and β_2 .

However, there is one important difference. When we introduce the squared variable in the model we inadvertently also introduce strong dependence between the terms which can cause significant numerical problems when it comes time to calculate the parameter estimates. Therefore, we should usually rescale the independent variable to have mean zero (and even variance one if we wish) *before* fitting the model. That is, we replace the x_i ’s with $x_i - \bar{x}$ (or $(x_i - \bar{x})/s$) before fitting the model⁴.

How to do it with R

There are multiple ways to fit a quadratic model to the variables Volume and Girth using R.

1. One way would be to square the values for Girth and save them in a vector `Girthsq`. Next, fit the linear model `Volume ~ Girth + Girthsq`.
2. A second way would be to use the `insulate` function in R, denoted by `I`: `Volume ~ Girth + I(Girth^2)` The second method is shorter than the first but the

⁴Rescaling the data gets the job done but a better way to avoid the multicollinearity introduced by the higher order terms is with *orthogonal polynomials*, whose coefficients are chosen just right so that the polynomials are not correlated with each other. This is beginning to linger outside the scope of this book, however, so we will content ourselves with a brief mention and then stick with the rescaling approach in the discussion that follows. A nice example of orthogonal polynomials in action can be run with `example(cars)`.

end result is the same. And once we calculate and store the fitted model (in, say, `treesquad.lm`) all of the previous comments regarding R apply.

3. A third and “right” way to do it is with orthogonal polynomials: `Volume ~ poly(Girth, degree = 2)` See `?poly` and `?cars` for more information. Note that we can recover the approach in 2 with `poly(Girth, degree = 2, raw = TRUE)`.

Example 4.3. We will fit the quadratic model to the `trees` data and display the results with `summary`, being careful to rescale the data before fitting the model. We may rescale the `Girth` variable to have zero mean and unit variance on-the-fly with the `scale` function.

```
treesquad.lm <- lm(Volume ~ scale(Girth) + I(scale(Girth)^2),
                  data = trees)
summary(treesquad.lm)
```

Call:

```
lm(formula = Volume ~ scale(Girth) + I(scale(Girth)^2), data = trees)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-5.4889	-2.4293	-0.3718	2.0764	7.6447

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	27.7452	0.8161	33.996	< 2e-16 ***
scale(Girth)	14.5995	0.6773	21.557	< 2e-16 ***
I(scale(Girth)^2)	2.5067	0.5729	4.376	0.000152 ***

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.335 on 28 degrees of freedom

Multiple R-squared: 0.9616, Adjusted R-squared: 0.9588

F-statistic: 350.5 on 2 and 28 DF, p-value: < 2.2e-16

We see that the F statistic indicates the overall model including `Girth` and `Girth^2` is significant. Further, there is strong evidence that both `Girth` and `Girth^2` are significantly related to `Volume`. We may examine a scatterplot together with the fitted quadratic function using the `lines` function, which adds a line to the plot tracing the estimated mean response.

```
plot(Volume ~ scale(Girth), data = trees)
lines(fitted(treesquad.lm) ~ scale(Girth), data = trees)
```

The plot is shown in Figure ?? . Pay attention to the scale on the x-axis: it is on the scale of the transformed `Girth` data and not on the original scale.

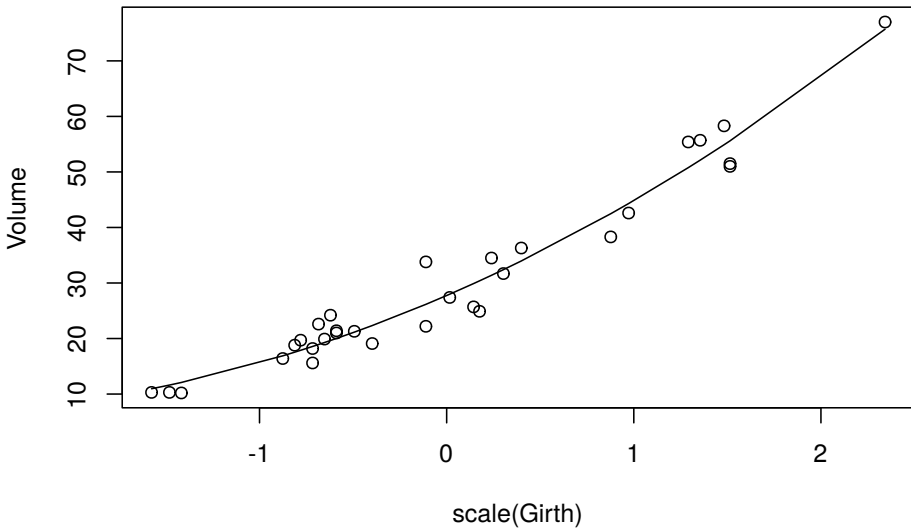


Figure 4.4: A quadratic model for the trees data.

Remark 4.3. When a model includes a quadratic term for an independent variable, it is customary to also include the linear term in the model. The principle is called *parsimony*. More generally, if the researcher decides to include x^m as a term in the model, then (s)he should also include all lower order terms x, x^2, \dots, x^{m-1} in the model.

We do estimation/prediction the same way that we did in Section ??, except we do not need a Height column in the dataframe =new= since the variable is not included in the quadratic model.

```
new <- data.frame(Girth = c(9.1, 11.6, 12.5))
predict(treesquad.lm, newdata = new, interval = "prediction")
```

	fit	lwr	upr
1	11.56982	4.347426	18.79221
2	20.30615	13.299050	27.31325
3	25.92290	18.972934	32.87286

The predictions and intervals are slightly different from what they were previously. Notice that it was not necessary to rescale the Girth prediction data before input to the predict function; the model did the rescaling for us automatically.

Remark 4.4. We have mentioned on several occasions that it is important to rescale the explanatory variables for polynomial regression. Watch what happens if we ignore this advice:

```
summary(lm(Volume ~ Girth + I(Girth^2), data = trees))
```

Call:

```
lm(formula = Volume ~ Girth + I(Girth^2), data = trees)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-5.4889	-2.4293	-0.3718	2.0764	7.6447

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	10.78627	11.22282	0.961	0.344728
Girth	-2.09214	1.64734	-1.270	0.214534
I(Girth^2)	0.25454	0.05817	4.376	0.000152 ***

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.335 on 28 degrees of freedom
 Multiple R-squared: 0.9616, Adjusted R-squared: 0.9588
 F-statistic: 350.5 on 2 and 28 DF, p-value: < 2.2e-16

Now nothing is significant in the model except Girth^2 . We could delete the Intercept and Girth from the model, but the model would no longer be *parsimonious*. A novice may see the output and be confused about how to proceed, while the seasoned statistician recognizes immediately that Girth and Girth^2 are highly correlated (see Section ??). The only remedy to this ailment is to rescale Girth, which we should have done in the first place.

In Example ?? of Section ??.

Note: The *trees* data do not have any qualitative explanatory variables, so we will construct one for illustrative purposes⁵. We will leave the Girth variable alone, but we will replace the variable Height by a new variable Tall which indicates whether or not the cherry tree is taller than a certain threshold (which for the sake of argument will be the sample median height of 76 ft). That is, Tall will be defined by

$$\text{Tall} = \begin{cases} \text{yes,} & \text{if Height} > 76, \\ \text{no,} & \text{if Height} \leq 76. \end{cases} \quad (4.48)$$

We can construct Tall very quickly in R with the cut function:

```
trees$Tall <- cut(trees$Height, breaks = c(-Inf, 76, Inf),
                 labels = c("no", "yes"))
trees$Tall[1:5]
```

⁵This procedure of replacing a continuous variable by a discrete/qualitative one is called *binning*, and is almost *never* the right thing to do. We are in a bind at this point, however, because we have invested this chapter in the *trees* data and I do not want to switch mid-discussion. I am currently searching for a data set with pre-existing qualitative variables that also conveys the same points present in the *trees* data, and when I find it I will update this chapter accordingly.

```
[1] no no no no yes
Levels: no yes
```

Note that `Tall` is automatically generated to be a factor with the labels in the correct order. See `?cut` for more.

Once we have `Tall`, we include it in the regression model just like we would any other variable. It is handled internally in a special way. Define a “dummy variable” `Tallyes` that takes values

$$\text{Tallyes} = \begin{cases} 1, & \text{if Tall} = \text{yes}, \\ 0, & \text{otherwise.} \end{cases} \quad (4.49)$$

That is, `Tallyes` is an *indicator variable* which indicates when a respective tree is tall. The model may now be written as

$$\text{Volume} = \beta_0 + \beta_1 \text{Girth} + \beta_2 \text{Tallyes} + \epsilon. \quad (4.50)$$

Let us take a look at what this definition does to the mean response. Trees with `Tall = yes` will have the mean response

$$\mu(\text{Girth}) = (\beta_0 + \beta_2) + \beta_1 \text{Girth}, \quad (4.51)$$

while trees with `Tall = no` will have the mean response

$$\mu(\text{Girth}) = \beta_0 + \beta_1 \text{Girth}. \quad (4.52)$$

In essence, we are fitting two regression lines: one for tall trees, and one for short trees. The regression lines have the same slope but they have different y intercepts (which are exactly $|\beta_2|$ far apart).

4.4.2 How to do it with R

The important thing is to double check that the qualitative variable in question is stored as a factor. The way to check is with the `class` command. For example,

```
class(trees$Tall)
```

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

If the qualitative variable is not yet stored as a factor then we may convert it to one with the `factor` command. See Section ?? . Other than this we perform MLR as we normally would.

```
treesdummy.lm <- lm(Volume ~ Girth + Tall, data = trees)
summary(treesdummy.lm)
```

Call:

```
lm(formula = Volume ~ Girth + Tall, data = trees)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-5.7788	-3.1710	0.4888	2.6737	10.0619

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-34.1652	3.2438	-10.53	3.02e-11 ***
Girth	4.6988	0.2652	17.72	< 2e-16 ***
Tallyes	4.3072	1.6380	2.63	0.0137 *

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.875 on 28 degrees of freedom

Multiple R-squared: 0.9481, Adjusted R-squared: 0.9444

F-statistic: 255.9 on 2 and 28 DF, p-value: < 2.2e-16

From the output we see that all parameter estimates are statistically significant and we conclude that the mean response differs for trees with `Tall = yes` and trees with `Tall = no`.

Remark 4.5. We were somewhat disingenuous when we defined the dummy variable `Tallyes` because, in truth, R defines `Tallyes` automatically without input from the user⁶. Indeed, the author fit the model beforehand and wrote the discussion afterward with the knowledge of what R would do so that the output the reader saw would match what (s)he had previously read. The way that R handles factors internally is part of a much larger topic concerning *contrasts*, which falls outside the scope of this book. The interested reader should see Neter et al [Neter1996] or Fox [Fox1997] for more.

Remark 4.6. In general, if an explanatory variable `foo` is qualitative with n levels `bar1`, `bar2`, ..., `bar n` then R will by default automatically define $n - 1$ indicator variables in the following way:

$$\text{foobar2} = \begin{cases} 1, & \text{if } \text{foo} = \text{"bar2"}, \\ 0, & \text{otherwise.} \end{cases}, \dots, \text{foobarn} = \begin{cases} 1, & \text{if } \text{foo} = \text{"bar n "}, \\ 0, & \text{otherwise.} \end{cases}$$

⁶That is, R by default handles contrasts according to its internal settings which may be customized by the user for fine control. Given that we will not investigate contrasts further in this book it does not serve the discussion to delve into those settings, either. The interested reader should check `?contrasts` for details.

The level `bar1` is represented by `foobar2 = ... = foobarn = 0`. We just need to make sure that `foo` is stored as a factor and R will take care of the rest.

4.4.3 Graphing the Regression Lines

We can see a plot of the two regression lines with the following mouthful of code.

```
treesTall <- split(trees, trees$Tall)
treesTall[["yes"]]$Fit <- predict(treesdummy.lm, treesTall[["yes"]])
treesTall[["no"]]$Fit <- predict(treesdummy.lm, treesTall[["no"]])
plot(Volume ~ Girth, data = trees)
points(Volume ~ Girth, data = treesTall[["yes"]], pch = 1)
points(Volume ~ Girth, data = treesTall[["no"]], pch = 2)
lines(Fit ~ Girth, data = treesTall[["yes"]])
lines(Fit ~ Girth, data = treesTall[["no"]])
```

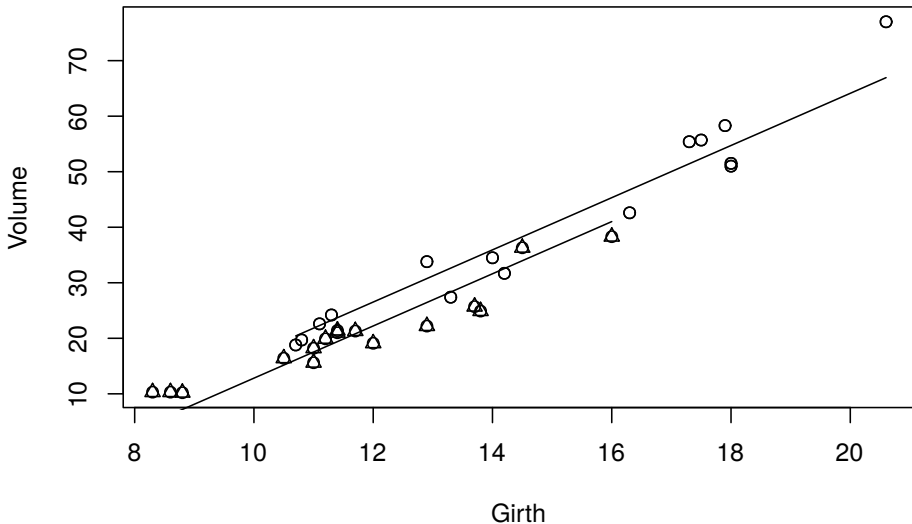


Figure 4.5: A dummy variable model for the `trees` data.

It may look intimidating but there is reason to the madness. First we `split` the `trees` data into two pieces, with groups determined by the `Tall` variable. Next we add the fitted values to each piece via `predict`. Then we set up a plot for the variables `Volume` versus `Girth`, but we do not plot anything yet (`type = n`) because we want to use different symbols for the two groups. Next we add `points` to the plot for the `Tall = yes` trees and use an open circle for a plot character (`pch = 1`), followed by `points` for the `Tall = no` trees with a triangle character (`pch = 2`). Finally, we add regression lines to the plot, one for each group.

There are other – shorter – ways to plot regression lines by groups, namely the

scatterplot function in the `car` package [`car`] and the `xyplot` function in the `lattice` package [`lattice`]. We elected to introduce the reader to the above approach since many advanced plots in R are done in a similar, consecutive fashion.

4.5 Partial F Statistic

We saw in Section ?? how to test $H_0 : \beta_0 = \beta_1 = \dots = \beta_p = 0$ with the overall F statistic and we saw in Section ?? how to test $H_0 : \beta_i = 0$ that a particular coefficient β_i is zero. Sometimes, however, we would like to test whether a certain part of the model is significant. Consider the regression model

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \beta_{j+1} x_{j+1} + \dots + \beta_p x_p + \epsilon, \quad (4.53)$$

where $j \geq 1$ and $p \geq 2$. Now we wish to test the hypothesis

$$H_0 : \beta_{j+1} = \beta_{j+2} = \dots = \beta_p = 0 \quad (4.54)$$

versus the alternative

$$H_1 : \text{at least one of } \beta_{j+1}, \beta_{j+2}, \dots, \beta_p \neq 0. \quad (4.55)$$

The interpretation of H_0 is that none of the variables x_{j+1}, \dots, x_p is significantly related to Y and the interpretation of H_1 is that at least one of x_{j+1}, \dots, x_p is significantly related to Y . In essence, for this hypothesis test there are two competing models under consideration:

$$\text{the full model: } y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon, \quad (4.56)$$

$$\text{the reduced model: } y = \beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \epsilon, \quad (4.57)$$

Of course, the full model will always explain the data *better* than the reduced model, but does the full model explain the data *significantly better* than the reduced model? This question is exactly what the partial F statistic is designed to answer.

We first calculate $SS E_f$, the unexplained variation in the full model, and $SS E_r$, the unexplained variation in the reduced model. We base our test on the difference $SS E_r - SS E_f$ which measures the reduction in unexplained variation attributable to the variables x_{j+1}, \dots, x_p . In the full model there are $p + 1$ parameters and in the reduced model there are $j + 1$ parameters, which gives a difference of $p - j$ parameters (hence degrees of freedom). The partial F statistic is

$$F = \frac{(SSE_r - SSE_f)/(p - j)}{SSE_f/(n - p - 1)}. \quad (4.58)$$

It can be shown when the regression assumptions hold under H_0 that the partial F statistic has an $f(df1 = p - j, df2 = n - p - 1)$ distribution. We calculate the p -value of the observed partial F statistic and reject H_0 if the p -value is small.

4.5.1 How to do it with R

The key ingredient above is that the two competing models are *nested* in the sense that the reduced model is entirely contained within the complete model. The way to test whether the improvement is significant is to compute `lm` objects both for the complete model and the reduced model then compare the answers with the `anova` function.

Example 4.4 (The Trees data). For the `trees` data, let us fit a polynomial regression model and for the sake of argument we will ignore our own good advice and fail to rescale the explanatory variables.

```
treesfull.lm <- lm(Volume ~ Girth + I(Girth^2) + Height +
                  I(Height^2), data = trees)
summary(treesfull.lm)
```

Call:

```
lm(formula = Volume ~ Girth + I(Girth^2) + Height + I(Height^2),
    data = trees)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-4.368	-1.670	-0.158	1.792	4.358

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.955101	63.013630	-0.015	0.988
Girth	-2.796569	1.468677	-1.904	0.068 .
I(Girth^2)	0.265446	0.051689	5.135	2.35e-05 ***
Height	0.119372	1.784588	0.067	0.947
I(Height^2)	0.001717	0.011905	0.144	0.886

Signif. codes:

```
0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 2.674 on 26 degrees of freedom

Multiple R-squared: 0.9771, Adjusted R-squared: 0.9735

F-statistic: 277 on 4 and 26 DF, p-value: < 2.2e-16

In this ill-formed model nothing is significant except `Girth` and `Girth^2`. Let us continue down this path and suppose that we would like to try a reduced model which contains nothing but `Girth` and `Girth^2` (not even an Intercept). Our two models are now

$$\begin{aligned}\text{the full model: } Y &= \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2 + \epsilon, \\ \text{the reduced model: } Y &= \beta_1 x_1 + \beta_2 x_1^2 + \epsilon,\end{aligned}$$

We fit the reduced model with `lm` and store the results:

```
treesreduced.lm <- lm(Volume ~ -1 + Girth + I(Girth^2), data = trees)
```

To delete the intercept from the model we used `-1` in the model formula. Next we compare the two models with the `anova` function. The convention is to list the models from smallest to largest.

```
anova(treesreduced.lm, treesfull.lm)
```

Analysis of Variance Table

```
Model 1: Volume ~ -1 + Girth + I(Girth^2)
Model 2: Volume ~ Girth + I(Girth^2) + Height + I(Height^2)
  Res.Df    RSS Df Sum of Sq    F    Pr(>F)
1      29 321.65
2      26 185.86  3    135.79 6.3319 0.002279 **
---
Signif. codes:
0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We see from the output that the complete model is highly significant compared to the model that does not incorporate `Height` or the Intercept. We wonder (with our tongue in our cheek) if the `Height^2` term in the full model is causing all of the trouble. We will fit an alternative reduced model that only deletes `Height^2`.

```
treesreduced2.lm <- lm(Volume ~ Girth + I(Girth^2) + Height,
                      data = trees)
anova(treesreduced2.lm, treesfull.lm)
```

Analysis of Variance Table

```
Model 1: Volume ~ Girth + I(Girth^2) + Height
Model 2: Volume ~ Girth + I(Girth^2) + Height + I(Height^2)
  Res.Df    RSS Df Sum of Sq    F    Pr(>F)
1      27 186.01
2      26 185.86  1    0.14865 0.0208 0.8865
```

In this case, the improvement to the reduced model that is attributable to Height^2 is not significant, so we can delete Height^2 from the model with a clear conscience. We notice that the p -value for this latest partial F test is 0.8865, which seems to be remarkably close to the p -value we saw for the univariate t test of Height^2 at the beginning of this example. In fact, the p -values are *exactly* the same. Perhaps now we gain some insight into the true meaning of the univariate tests.

4.6 Residual Analysis and Diagnostic Tools

We encountered many, many diagnostic measures for simple linear regression in Sections ?? and ?. All of these are valid in multiple linear regression, too, but there are some slight changes that we need to make for the multivariate case. We list these below, and apply them to the trees example.

- **Shapiro-Wilk, Breusch-Pagan, Durbin-Watson:** unchanged from SLR, but we are now equipped to talk about the Shapiro-Wilk test statistic for the residuals. It is defined by the formula

$$W = \frac{\mathbf{a}^T \mathbf{E}^*}{\mathbf{E}^T \mathbf{E}}, \quad (4.59)$$

where \mathbf{E}^* is the sorted residuals and $\mathbf{a}_{1 \times n}$ is defined by

$$\mathbf{a} = \frac{\mathbf{m}^T \mathbf{V}^{-1}}{\sqrt{\mathbf{m}^T \mathbf{V}^{-1} \mathbf{V}^{-1} \mathbf{m}}}, \quad (4.60)$$

where $\mathbf{m}_{n \times 1}$ and $\mathbf{V}_{n \times n}$ are the mean and covariance matrix, respectively, of the order statistics from an mvnrm (mean = $\mathbf{0}$, sigma = \mathbf{I}) distribution.

- **Leverages:** are defined to be the diagonal entries of the hat matrix \mathbf{H} (which is why we called them h_{ii} in Section ?). The sum of the leverages is $\text{tr}(\mathbf{H}) = p + 1$. One rule of thumb considers a leverage extreme if it is larger than double the mean leverage value, which is $2(p + 1)/n$, and another rule of thumb considers leverages bigger than 0.5 to indicate high leverage, while values between 0.3 and 0.5 indicate moderate leverage.
- **Standardized residuals:** unchanged. Considered extreme if $|R_i| > 2$.
- **Studentized residuals:** compared to a $t(\text{df} = n - p - 2)$ distribution.
- **DFBETAS:** The formula is generalized to

$$(\text{DFBETAS})_{j(i)} = \frac{b_j - b_{j(i)}}{S_{(i)} \sqrt{c_{jj}}}, \quad j = 0, \dots, p, \quad i = 1, \dots, n, \quad (4.61)$$

where c_{jj} is the j^{th} diagonal entry of $(\mathbf{X}^T \mathbf{X})^{-1}$. Values larger than one for small data sets or $2/\sqrt{n}$ for large data sets should be investigated.

- **DFFITS:** unchanged. Larger than one in absolute value is considered extreme.
- **Cook's D:** compared to an $f(\text{df1} = p + 1, \text{df2} = n - p - 1)$ distribution. Observations falling higher than the 50th percentile are extreme. Note that plugging the value $p = 1$ into the formulas will recover all of the ones we saw in Chapter ?.

4.7 Additional Topics

4.7.1 Nonlinear Regression

We spent the entire chapter talking about the `trees` data, and all of our models looked like `Volume ~ Girth + Height` or a variant of this model. But let us think again: we know from elementary school that the volume of a rectangle is $V = lwh$ and the volume of a cylinder (which is closer to what a black cherry tree looks like) is

$$V = \pi r^2 h \quad \text{or} \quad V = 4\pi dh, \quad (4.62)$$

where r and d represent the radius and diameter of the tree, respectively. With this in mind, it would seem that a more appropriate model for μ might be

$$\mu(x_1, x_2) = \beta_0 x_1^{\beta_1} x_2^{\beta_2}, \quad (4.63)$$

where β_1 and β_2 are parameters to adjust for the fact that a black cherry tree is not a perfect cylinder.

How can we fit this model? The model is not linear in the parameters any more, so our linear regression methods will not work... or will they? In the `trees` example we may take the logarithm of both sides of Equation (??) to get

$$\mu^*(x_1, x_2) = \ln[\mu(x_1, x_2)] = \ln\beta_0 + \beta_1 \ln x_1 + \beta_2 \ln x_2, \quad (4.64)$$

and this new model μ^* is linear in the parameters $\beta_0^* = \ln\beta_0$, $\beta_1^* = \beta_1$ and $\beta_2^* = \beta_2$. We can use what we have learned to fit a linear model `log(Volume) ~ log(Girth) + log(Height)`, and everything will proceed as before, with one exception: we will need to be mindful when it comes time to make predictions because the model will have been fit on the log scale, and we will need to transform our predictions back to the original scale (by exponentiating with `exp`) to make sense.

```
treesNonlin.lm <- lm(log(Volume) ~ log(Girth) +
                     log(Height), data = trees)
summary(treesNonlin.lm)
```

Call:

```
lm(formula = log(Volume) ~ log(Girth) + log(Height), data = trees)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-0.168561	-0.048488	0.002431	0.063637	0.129223

Coefficients:

```

              Estimate Std. Error t value Pr(>|t|)
(Intercept) -6.63162      0.79979  -8.292 5.06e-09 ***
log(Girth)   1.98265      0.07501  26.432 < 2e-16 ***
log(Height)  1.11712      0.20444   5.464 7.81e-06 ***

```

```
---
```

```
Signif. codes:
```

```
0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

Residual standard error: 0.08139 on 28 degrees of freedom
Multiple R-squared:  0.9777,    Adjusted R-squared:  0.9761
F-statistic: 613.2 on 2 and 28 DF,  p-value: < 2.2e-16

```

This is our best model yet (judging by R^2 and \bar{R}^2), all of the parameters are significant, it is simpler than the quadratic or interaction models, and it even makes theoretical sense. It rarely gets any better than that.

We may get confidence intervals for the parameters, but remember that it is usually better to transform back to the original scale for interpretation purposes:

```
exp(confint(treesNonlin.lm))
```

```

              2.5 %      97.5 %
(Intercept) 0.0002561078 0.006783093
log(Girth)   6.2276411645 8.468066317
log(Height)  2.0104387829 4.645475188

```

(Note that we did not update the row labels of the matrix to show that we exponentiated and so they are misleading as written.) We do predictions just as before. Remember to transform the response variable back to the original scale after prediction.

```

new <- data.frame(Girth = c(9.1, 11.6, 12.5),
                  Height = c(69, 74, 87))
exp(predict(treesNonlin.lm, newdata = new,
            interval = "confidence"))

```

```

      fit      lwr      upr
1 11.90117 11.25908 12.57989
2 20.82261 20.14652 21.52139
3 28.93317 27.03755 30.96169

```

The predictions and intervals are slightly different from those calculated earlier, but they are close. Note that we did not need to transform the `Girth` and `Height` arguments in the dataframe `new`. All transformations are done for us automatically.

4.7.2 Real Nonlinear Regression

We saw with the `trees` data that a nonlinear model might be more appropriate for the data based on theoretical considerations, and we were lucky because the functional form of μ allowed us to take logarithms to transform the nonlinear model to a linear one. The same trick will not work in other circumstances, however. We need techniques to fit general models of the form

$$\mathbf{Y} = \mu(\mathbf{X}) + \epsilon, \quad (4.65)$$

where μ is some crazy function that does not lend itself to linear transformations.

There are a host of methods to address problems like these which are studied in advanced regression classes. The interested reader should see Neter *et al* [Neter1996] or Tabachnick and Fidell [Tabachnick2006].

It turns out that John Fox has posted an Appendix to his book [Fox2002] which discusses some of the methods and issues associated with nonlinear regression; see <http://cran.r-project.org/doc/contrib/Fox-Companion/appendix.html> for more. Here is an example of how it works, based on a question from R-help.

```
# fake data
set.seed(1)
x <- seq(from = 0, to = 1000, length.out = 200)
y <- 1 + 2*(sin((2*pi*x/360) - 3))^2 + rnorm(200, sd = 2)
# plot(x, y)
acc.nls <- nls(y ~ a + b*(sin((2*pi*x/360) - c))^2,
               start = list(a = 0.9, b = 2.3, c = 2.9))
summary(acc.nls)
```

Formula: $y \sim a + b * (\sin((2 * \pi * x/360) - c))^2$

Parameters:

	Estimate	Std. Error	t value	Pr(> t)
a	0.95884	0.23097	4.151	4.92e-05 ***
b	2.22868	0.37114	6.005	9.07e-09 ***
c	3.04343	0.08434	36.084	< 2e-16 ***

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.865 on 197 degrees of freedom

Number of iterations to convergence: 3

Achieved convergence tolerance: 6.554e-08

```
#plot(x, fitted(acc.nls))
```

4.7.3 Multicollinearity

A multiple regression model exhibits *multicollinearity* when two or more of the explanatory variables are substantially correlated with each other. We can measure multicollinearity by having one of the explanatory play the role of “dependent variable” and regress it on the remaining explanatory variables. The the R^2 of the resulting model is near one, then we say that the model is multicollinear or shows multicollinearity.

Multicollinearity is a problem because it causes instability in the regression model. The instability is a consequence of redundancy in the explanatory variables: a high R^2 indicates a strong dependence between the selected independent variable and the others. The redundant information inflates the variance of the parameter estimates which can cause them to be statistically insignificant when they would have been significant otherwise. To wit, multicollinearity is usually measured by what are called *variance inflation factors*.

Once multicollinearity has been diagnosed there are several approaches to remediate it. Here are a couple of important ones.

- **Principal Components Analysis.** This approach casts out two or more of the original explanatory variables and replaces them with new variables, derived from the original ones, that are by design uncorrelated with one another. The redundancy is thus eliminated and we may proceed as usual with the new variables in hand. Principal Components Analysis is important for other reasons, too, not just for fixing multicollinearity problems.
- **Ridge Regression.** The idea of this approach is to replace the original parameter estimates with a different type of parameter estimate which is more stable under multicollinearity. The estimators are not found by ordinary least squares but rather a different optimization procedure which incorporates the variance inflation factor information.

We decided to omit a thorough discussion of multicollinearity because we are not equipped to handle the mathematical details. Perhaps the topic will receive more attention in a later edition.

- What to do when data are not normal
- Bootstrap (see Chapter ??).

4.7.4 Akaike’s Information Criterion

$$AIC = -2 \ln L + 2(p + 1)$$

4.8 Exercises

Exercise. Use Equations (??), (??), and (??) to prove the Anova Equality:

$$SSTO = SSE + SSR.$$

