

DR. NATHANIEL D. PHILLIPS

# YARRRR! THE PIRATE'S GUIDE TO R

Copyright © 2015 Dr. Nathaniel D. Phillips

PUBLISHED BY

<http://www.nathanieldphillips.com>

This document may not be used for any commercial purposes. All rights are reserved by Nathaniel Phillips.

*First printing, May 2015*

# Contents

<i>Introduction</i>	9
<i>1: Installing R, RStudio, and the pirate dataset</i>	13
<i>Getting help</i>	16
<i>Installing and loading packages</i>	16
<i>Downloading and loading the pirate dataset</i>	17
<i>The R Reference Card</i>	18
<i>2: Coding Basics</i>	21
<i>Defining objects with the &lt;- assignment</i>	22
<i>Data object types in R</i>	23
<i>Generating numeric vectors</i>	26
<i>3: Sampling data and Descriptive Statistics</i>	31
<i>Sampling data from probability distributions</i>	31
<i>Descriptive statistics</i>	36
<i>A worked example: A quick test of the law of large numbers</i>	39
<i>4: Indexing and comparing vectors</i>	41
<i>Indexing vectors with brackets</i>	41
<i>Creating logical vectors</i>	42
<i>Indexing data with logical vectors</i>	44

<i>Additional helpful vector functions</i>	46
<i>Set Functions</i>	47
<i>Using indexing to remove specific values of a vector</i>	47
<i>Taking the sum and mean of logical vectors to get counts and percentages</i>	48
<i>A worked example - Chicken Weights</i>	50
 <b>5: Matrices and Data Frames</b>	 53
<i>Creating matrices and dataframes</i>	53
<i>Data sets pre-loaded in R</i>	57
<i>Viewing matrices and dataframes</i>	57
<i>Loading data into R with read.table()</i>	59
 <b>6: Basic Dataframe Manipulation</b>	 63
<i>Getting information about matrices and dataframes</i>	63
<i>Indexing dataframes with brackets [rows, columns]</i>	64
<i>Adding new columns to a dataframe</i>	66
<i>Centering and standardizing (z-score) data</i>	67
<i>Subsetting dataframes with logical indexing and subset()</i>	69
<i>Combining indexing and descriptive statistics</i>	71
<i>A worked example: Credit default</i>	71
 <b>7: Plotting Basics</b>	 75
<i>High-level plotting functions</i>	75
<i>Symbol types: pch</i>	77
<i>Other high-Level plotting commands</i>	78
<i>Low-level plotting functions</i>	82
<i>Additional low-level plotting functions</i>	89
<i>Saving plots to a file</i>	90
<i>A worked example: Creating a plot with automated numeric labels</i>	91
<i>Additional Tips</i>	93

<b>8: Customizing Plots</b>	<b>95</b>
Colors in R	95
Plot margins	100
Arranging multiple plots with <i>par(mfrow)</i> and <i>layout</i>	102
Using alternative fonts in pdfs with the <i>extrafont</i> package	104
Additional Tips	107
<b>9: Advanced dataframe manipulation</b>	<b>109</b>
Recoding values in a vector	109
Splitting numerical data into groups using <i>cut()</i>	112
Grouped aggregation	114
Aggregation with <i>dplyr</i>	117
Merging two dataframes	121
<b>10: 1 and 2-sample Null-Hypothesis tests</b>	<b>125</b>
Warning about null-hypothesis tests with "frequentist" statistics	125
T-test	126
Correlation test	131
Chi-square test	133
<b>11: Regression and ANOVA</b>	<b>135</b>
The Linear Model	135
Generalized Linear Model (GLM)	137
ANOVA	139
<b>12: Writing your own functions (Coming Soon!)</b>	<b>141</b>
The basic structure of a function	141
Tips and tricks for complex functions	141
Storing and loading your functions to and from a function file	141

<b>13: Loops and Simulations (Coming Soon!)</b>	<b>143</b>
<i>When and when not to use loops</i>	143
<i>The list object</i>	143
<i>Simple loops over one index</i>	143
<i>Loops over multiple indices</i>	143
<i>Printing and saving temporary results</i>	143
<i>Parallel computing with snowfall()</i>	143
<b>14: Bayesian Inference (Coming Soon!)</b>	<b>145</b>
<i>What are Bayesian statistics?</i>	145
<i>Bayesian one and two sample tests</i>	145
<i>Bayesian general linear model</i>	145
<b>15: Model fitting (Coming Soon!)</b>	<b>147</b>
<i>What is a model?</i>	147
<i>What is a loss function?</i>	147
<i>Minimizing loss functions with optimization routines</i>	147
<i>A worked example: Prospect Theory</i>	147
<b>16: Writing and sharing your work (Coming Soon!)</b>	<b>149</b>
<i>RMarkdown</i>	149
<i>Shiny</i>	149
<i>Sweave (R and Latex)</i>	149
<b>Appendix</b>	<b>151</b>
<b>Index</b>	<b>155</b>

*This book is dedicated to my former statistics  
instructors Dr. Thomas Moore and Dr.*

*Wei Lin who taught me everything I know  
about statistics, and my PhD colleagues*

*Dr. Dirk Wulff and Dr. Stefan Herzog  
who taught me everything I know about  
R.*





# *Introduction*

## *Who am I?*

I am a pirate on the Bodensee in Konstanz Germany. When I started pirate training, I discovered R and have been hooked ever since. I'm now on a mission to convince everyone I can to make the switch from SPSS (or Excel, Matlab, JMP...) to R.

## *This book is in progress..*

If you haven't figured it out already, this book is very much a work in progress. I'm constantly experimenting with the material and the layout. If you have any recommendations for changes or spot any errors, please write me at [YaRrr.Book@gmail.com](mailto:YaRrr.Book@gmail.com) or tweet me @YaRrrBook

Email me with comments, recommendations or typos at:  
[YaRrr.Book@gmail.com](mailto:YaRrr.Book@gmail.com) or tweet me  
at @YaRrrBook

## *Who is this book for?*

Anyone who wants to learn R can benefit from this book. I will assume that you have taken an introductory course in statistics, but have no substantial programming experience. While the techniques in this book apply to most data analysis problems, because my background is in experimental psychology I will cater the course to solving analysis problems commonly faced in psychological research.

## *What this book is*

This book is meant to introduce you to the basic analytical tools in R, from basic coding and analyses, to data wrangling, plotting, and statistical inference.

## *What this book is **not***

This book does not cover any one topic in extensive detail. If you are interested in conducting analyses or creating plots not covered in the book, I'm sure you'll find the answer with a quick Google search!

## Why is R so great?

As you've already gotten this book, you probably already have some idea why R is so great. However, in order to help prevent you from giving up the first time you run into a programming wall, let me give you a few more reasons:

1. R is 100% free and as a result, has a huge support community. Unlike SPSS, Matlab, Excel and JMP, R is, and always will be completely free. This doesn't just help your wallet - it means that a huge community of R programmers will constantly develop and distribute new R functionality and packages at a speed that leaves all those other packages in the dust! Unlike Fight Club, the first rule of R is "Do talk about R!" The size of the R programming community is staggering. If you ever have a question about how to implement something in R, a quick Poog<sup>1</sup> search will lead you to your answer virtually every single time.
2. R is incredibly versatile. You can use R to do everything from calculating simple summary statistics, to performing complex simulations to creating gorgeous plots like the chord diagram in Figure 1. If you can imagine an analytical task, you can almost certainly implement it in R.
3. Using RStudio, You can easily and seamlessly combine R code, analyses, plots, and written text into elegant documents all in one place using Sweave (R and Latex) or RMarkdown. In fact, I wrote this entire book (the text, formatting, plots, code...yes, everything) in RStudio using Sweave. With RStudio and Sweave, instead of trying to manage two or three programs, say Excel, Word and (sigh) SPSS, where you find yourself spending half your time copying, pasting and formatting data, images and text, you can do everything in one place so nothing gets misread, mistyped, or forgotten.
4. Analyses conducted in R are transparent, easily shareable, and reproducible. If you ask an SPSS user how they conducted a specific analyses, they will either A) Not remember, B) Try (nervously) to construct an analysis procedure on the spot that makes sense - which may or may not correspond to what they actually did months or years ago, or C) Ask you what you are doing in their kitchen<sup>2</sup>. I used to primarily use SPSS, so I speak from experience on this. If you ask an R user (who uses good programming techniques!) how they conducted an analysis, they should always be able to show you the exact code they used. Of course, this doesn't mean that they used the appropriate analysis or interpreted it correctly, but with all the original code, any problems should be

<sup>1</sup> I am in the process of creating Poog<sup>le</sup> - Google for Pirates. Kickstarter page coming soon...

```
require("circlize")

## Loading required package: circlize

mat = matrix(sample(1:100, 18, replace = TRUE), 3, 6)
rownames(mat) = letters[1:3]
colnames(mat) = LETTERS[1:6]
chordDiagram(mat)
```

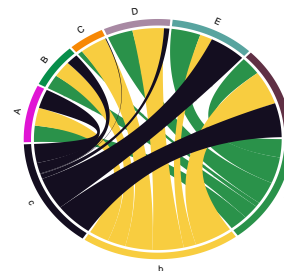


Figure 1: This is a chordDiagram plot that comes with the R package circlize.

<sup>2</sup> Get used to the bad jokes people. Lots more where that came from.

completely transparent!

5. And most importantly of all, R is the programming language of choice for pirates, (who prefer the "YaRrr!" pronunciation)

### *Code Chunks*

In this book, R code is (almost) always presented in a separate gray box like this one:

```
a <- 1 + 2 + 3 + 4 + 5
a
## [1] 15
```

This is called a *code chunk*. You should always be able to directly copy and paste code chunks directly into R. If you copy a chunk and it does not work for you, it is most likely because the code refers to a package, function, or object that I defined in a previous chunk. If so, read back and look for a previous chunk that contains the missing definition. As you'll soon learn, lines that begin with # are either comments or output from prior code that R will ignore.

As you'll notice, I'll include code chunks before all plots in the book. In early chapters, the code might not make sense just yet. However, I elected to always include plotting code so you have the option of re-creating (and tweaking) any plot in the book.

### *Additional Tips*

Because this is a beginner's book, I try to avoid bombarding you with too many details and tips when I first introduce a function. For this reason, at the end of every chapter, I present several tips and tricks that I think most R users should know once they get comfortable with the basics. I highly encourage you to read these additional tips as I expect you'll find at least some of them very useful if not invaluable.



# 1: Installing R, RStudio, and the pirate dataset

Now that I've convinced you to use R, let's get started! First, you'll need to install the base R software.

1. Download and install the base R software (around 50mb) + Windows <<http://cran.r-project.org/bin/windows/base/>> + Mac <<http://cran.r-project.org/bin/macosx/>>

See Figure 2 Here's how the base R software looks (on Mac). As you can see, it's very much a bare-bones software - just how we want it! No extra gimmicks or flashy bloatware needed!

While you can do pretty much everything you want within base R, you'll find that most people these days do their R programming in an application called RStudio. RStudio is a graphical user interface (GUI)-like interface for R that makes programming in R a bit easier. To download and install RStudio (around 40mb), go to <<http://www.rstudio.com/products/rstudio/download/>>

Once you've installed RStudio, you'll never need to open the base R application. Let's go ahead and boot up RStudio and see how she looks!

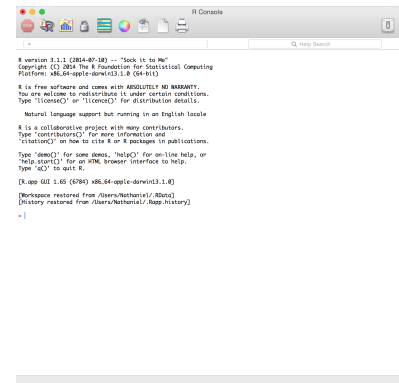


Figure 2: Here is how the standard R application looks. Not too exciting - just how we like it!

## *The four RStudio windows*

When you open RStudio, you'll see the following four windows (also called panes):

Note your windows might be in a different order. You can change the order of the windows under RStudio preferences.

## *Source - Your notepad for code*

The source pane is where you create and edit R Scripts - which are just text files with the ".R" extension. When you open RStudio, it will automatically start a new Untitled script. You will write 99% of your R code in a script in the source panel. However, your R code will not be evaluated until you 'send' the code to the Console.

You can send your code from the source to the Console by highlighting the code you wish to evaluate and clicking on the "Run" button on the top right of the Source. Alternatively, you can use the

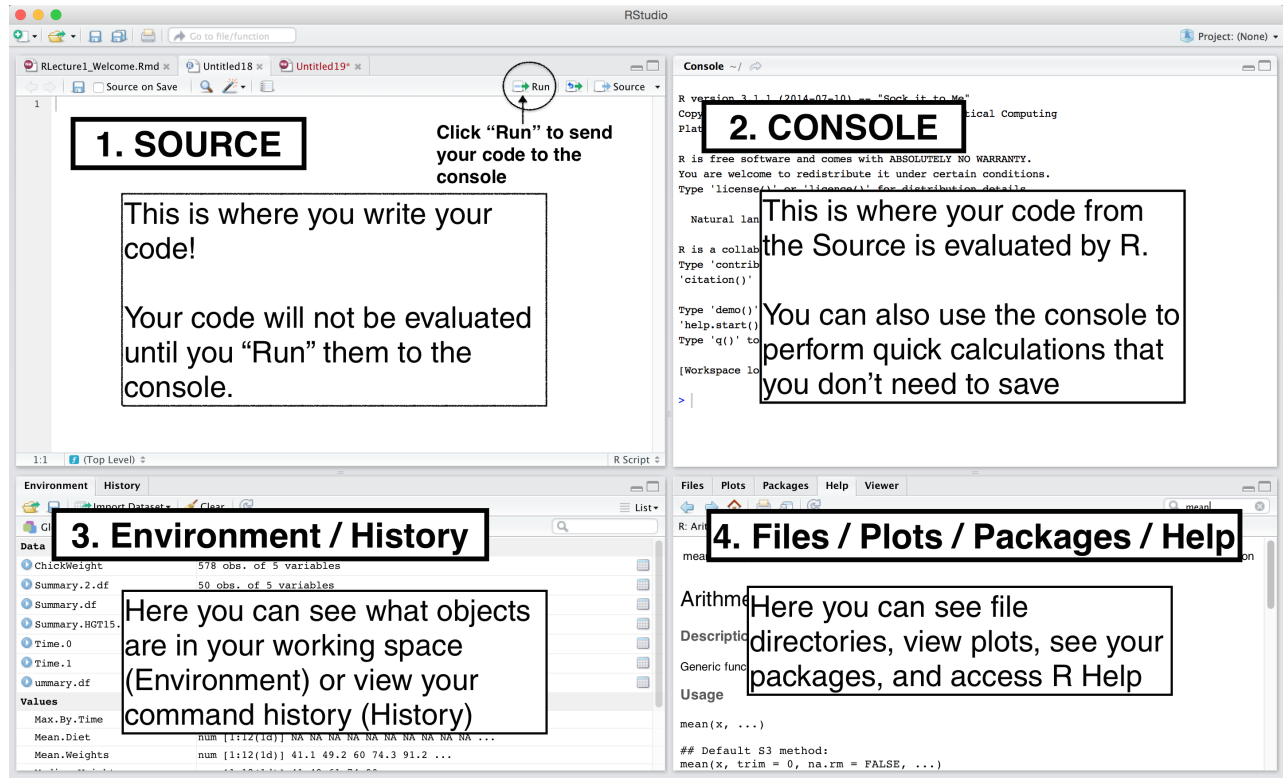


Figure 3: The four panes of RStudio.

hot-key "Command + Return" on Mac, or "Control + Enter" on PC to send code to the console.

### Console: The calculator

The console is where R actually executes (calculates) code. You can type code directly into the console and get an immediate response. For example, if you type `1+1` into the console, you'll see that R immediately gives an output of 2

```
1+1
## [1] 2
```

However, most of the time, you won't be typing directly into the console. Instead, you'll be writing code in the source and then "Running" it to the console. The reason for this is straightforward: If you type code into the console, it won't be saved (though you can look back on your command History). And if you make a mistake in typing code into the console, you'd have to re-type everything all over again. Instead, it's better to write all your code in the Source. When you are ready to execute some code, you can then send "Run"

Tip: Try to write most of your code in a document in the Source. Only type directly into the Console to de-bug or do quick analyses.

it to the console.

### *Environment / History*

The Environment tab of this panel shows you the names of all the data objects (like vectors, matrices, and dataframes) that you've defined in your current R session. The tab also has a few clickable actions like importing a new dataset. However, I almost never look at this menu.

The History tab of this panel simply shows you a history of your R commands. I never look at this. In fact, I didn't realize it was even there until I started writing this tutorial.

### *Files / Plots / Packages / Help*

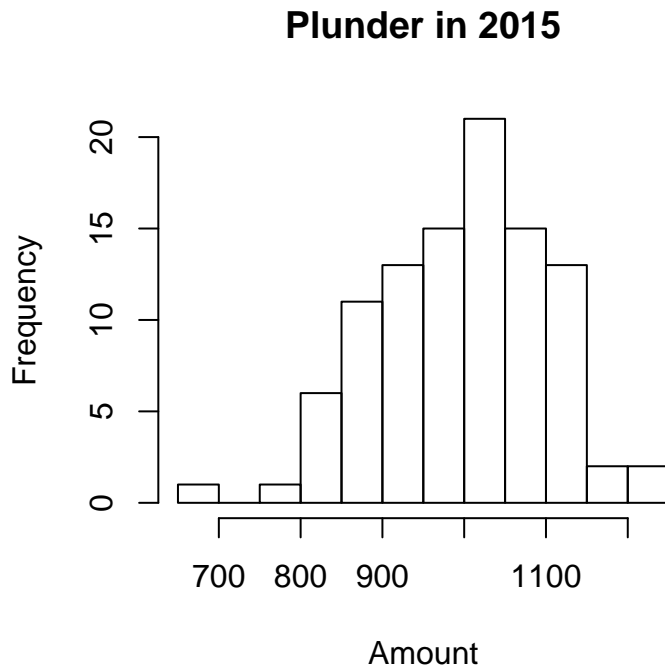
This panel shows you file directories, plots, your current packages, and help menus.

1. Files - Gives you access to the file directory on your harddrive. One nice feature of the "Files" panel is that you can use it to set your working directory - once you navigate to a folder you want to read and save files to, click "More" and then "Set As Working Directory."
2. Plots - Shows your plots. There are buttons for opening the plot in a separate window and exporting the plot as a pdf or jpeg (though you can also do this with code using the `pdf()` or `jpeg()` functions.)
3. Packages - Shows a list of all the R packages installed on your harddrive and indicates whether or not they are currently loaded. Packages that are loaded in the current session are checked while those that are installed but not yet loaded are unchecked.
4. Help - Help menu for R functions. You can either type the name of a function in the search window, or use the code `?function.name` to search for a function with the name `function.name`

Most - if not all - of the time when you perform actions using your mouse by pointing and clicking in RStudio, RStudio will perform the function by sending the appropriate R Code to the console. You can then copy and paste this code into your documents to automate the process later.

To see how plots are displayed try the following command which should display a histogram of 100 values randomly drawn from a standard normal distribution.

```
hist(x = rnorm(n = 100, mean = 1000, sd = 100),
     main = "Plunder in 2015",
     xlab = "Amount"
)
```



### Getting help

To get help and see documentation for a function, type `?fun`, where `fun` is the name of the function. For example, to get additional information on the histogram function, run the following code:

```
?hist
```

Tip: If you ever need to learn more about an R function: type `?functionname`, where `functionname` is the name of the function.

### Installing and loading packages

When you download and install R for the first time, you are installing the Base R software. Base R will contain most if not all the functions you need. However, one of the great things about R is that people are constantly writing and sharing new functions that you can use. When people share a new function, they usually do so in the form of an *R package* which contains anything from functions, to help menus, to vignettes (examples), to data. To install a new R package, you need to run the code `install.packages("package")`, where "package" is the name of the package. After you've installed the package, you need to *load* it into R by running the code `load("package")`. This will load the package into your current R session and allow you to use its contents.

Once you've installed a package on your computer, you never need to install it again. However, you do need to load the package every time you start a new R session.



For example, let's say you want to create a wordcloud - a graph that plots text in different sizes. You can certainly program this yourself in R, but thankfully someone has created a package called `wordcloud` with a function that will do this for you. Let's install the package, load it, and then use the `wordcloud` function:

```
install.packages("wordcloud")

##
## The downloaded binary packages are in
## /var/folders/yh/4fyk9h754h7fnmpgyc6p5whw0000gn/T//RtmpJ3VXTu/downloaded_packages

library("wordcloud") # Install the package

## Loading required package: RColorBrewer

par(mar = rep(0, 4))
wordcloud(words = c("sword", "YaRrr!", "eyepatch",
                    "parrot", "plunder", "treasure",
                    "chest", "scurvy"),
          freq = sample(50:1000, 8),
          colors = gray(runif(8, 0, 1)))
```



### *Downloading and loading the pirate dataset*

For much of this book, we will be referring to the pirate dataset. This dataset contains the results of a survey completed by 1,000 pirates that have been on my crew. To download this dataset, execute the following code:

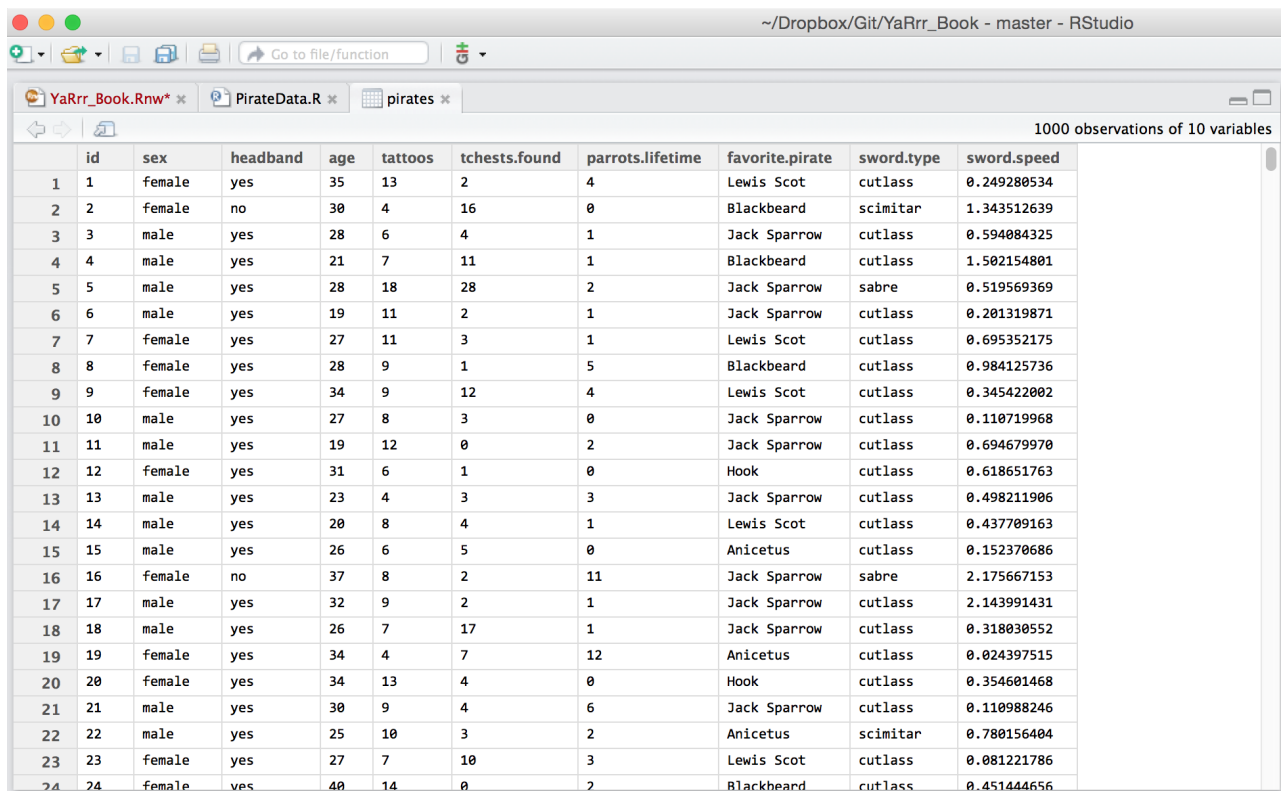
```
pirates <- read.table(file = "http://nathanielphillips.com/wp-content/uploads/2015/05/pirate_survey.txt",
                      header = T,
                      sep = "\t",
```

```
stringsAsFactors = F
)
```

If you'd like to see how the dataset looks, you can execute the `View()` function:

```
View(pirates)
```

When you run this command, you should see the first several rows and columns of the dataset (like this:)



	id	sex	headband	age	tattoos	tchests.found	parrots.lifetime	favorite.pirate	sword.type	sword.speed
1	1	female	yes	35	13	2	4	Lewis Scot	cutlass	0.249280534
2	2	female	no	30	4	16	0	Blackbeard	scimitar	1.343512639
3	3	male	yes	28	6	4	1	Jack Sparrow	cutlass	0.594084325
4	4	male	yes	21	7	11	1	Blackbeard	cutlass	1.502154801
5	5	male	yes	28	18	28	2	Jack Sparrow	sabre	0.519569369
6	6	male	yes	19	11	2	1	Jack Sparrow	cutlass	0.201319871
7	7	female	yes	27	11	3	1	Lewis Scot	cutlass	0.695352175
8	8	female	yes	28	9	1	5	Blackbeard	cutlass	0.984125736
9	9	female	yes	34	9	12	4	Lewis Scot	cutlass	0.345422002
10	10	male	yes	27	8	3	0	Jack Sparrow	cutlass	0.110719968
11	11	male	yes	19	12	0	2	Jack Sparrow	cutlass	0.694679970
12	12	female	yes	31	6	1	0	Hook	cutlass	0.618651763
13	13	male	yes	23	4	3	3	Jack Sparrow	cutlass	0.498211906
14	14	male	yes	20	8	4	1	Lewis Scot	cutlass	0.437709163
15	15	male	yes	26	6	5	0	Anicetus	cutlass	0.152370686
16	16	female	no	37	8	2	11	Jack Sparrow	sabre	2.175667153
17	17	male	yes	32	9	2	1	Jack Sparrow	cutlass	2.143991431
18	18	male	yes	26	7	17	1	Jack Sparrow	cutlass	0.318030552
19	19	female	yes	34	4	7	12	Anicetus	cutlass	0.024397515
20	20	female	yes	34	13	4	0	Hook	cutlass	0.354601468
21	21	male	yes	30	9	4	6	Jack Sparrow	cutlass	0.110988246
22	22	male	yes	25	10	3	2	Anicetus	scimitar	0.780156404
23	23	female	yes	27	7	10	3	Lewis Scot	cutlass	0.081221786
24	24	female	yes	40	14	0	2	Blackbeard	cutlass	0.451444656

Figure 4: The pirates dataset.

## The R Reference Card

Over the course of this book, you will be learning *lots* of new functions. Wouldn't it be nice if someone created a Cheatsheet / Notecard of many common R functions? Yes it would, and thankfully Tom Short has done this in his creation of the R Reference Card. You can download a copy at <https://dl.dropboxusercontent.com/u/7618380/RReferenceCard.pdf>. I highly encourage you to print this out and start highlighting functions as you learn them!

*Finished!*

That's it for this lecture! All you did was install the most powerful statistical package on the planet used by top universities and companies like Google. No big deal.



## 2: Coding Basics

### *Chapter Goals*

1. Accept that learning R will take time (and promise you'll never go back to SPSS!)
2. Know how to use comments and spaces in R code.
3. Be able to define and manipulate scalars and vectors
4. Generate vectors using `c()`, `:`, `rep()`, and `seq()`

### *Before we get started, a word of warning...*

So by now you've installed R and you're ready to get started. But first, let me give you a brief word of warning: Especially if this is your first experience programming, you are going to experience a *lot* of headaches when you get started. You will run into error after error and pound your fists against the table screaming: "WHY ISN'T MY CODE WORKING?!?!? There must be something wrong with this stupid software!!!" You will spend hours trying to find a bug in your code, only to find that - frustratingly enough, you had had an extra space or missed a comma somewhere. You'll then wonder why you ever decided to learn R when (::sigh::) SPSS was so "nice and easy."

This is perfectly normal! Don't get discouraged and DON'T GO BACK TO SPSS! Trust me, as you gain more programming experience, you'll experience fewer and fewer bugs (though they'll never go away completely). Once you get over the initial barriers, you'll find yourself conducting analyses much, much faster than you ever did before.

Fun fact: SPSS stands for "Shitty Piece of Shitty Shit". True story.

### *The basics of R programming*

Ok, let's write some code! Again, we will write all our code in a script file in the Source pane of RStudio. When we want to execute it, we'll send it to the Console.

*R as a calculator*

At its heart, R is just a fancy calculator. Let's do some basic algebra, type the following command into the source, then highlight the text and click "Run" to execute it in the console:

```
1+1 # The result should be 2
## [1] 2
```

As you can see, R returns the (thankfully correct) value of 2. You'll notice that the console also returns the text [1]. This is just telling you you the index of the value next to it. Don't worry about this for now, it will make more sense later.

Additionally, you'll notice that I included a comment in the code using the # sign. R will ignore everything on a line after the # sign. So why do we use comments? Mainly to explain to others, including your future self, what you are trying to do with your code.

Let's try some more simple calculations

```
2 * 3 - 1 # R ignores spaces
## [1] 5

2 * (3 - 1) # R observes order of operations
## [1] 4
```

As you can see, R ignores spaces in between arguments in code. I recommend using spaces to make your code easier to look at. Personally, I include spaces between arithmetic operators (like + and -) and after commas (which we'll get to later).

*Defining objects with the <- assignment*

So far so good, you can use R as a simple calculator. Now, let's do our first *object assignment*. Object assignment is our way of storing information, such as a number or a statistical test, into something we can easily refer to later. Let's start by creating the object "mateys" and assigning the outcome of  $5 * 10$  to it:

```
# The symbol(s) "<-" mean "assign"
mateys <- 5 * 10 # Assign the value of 5*10 to a new object called mateys
```

Now, anytime we want to refer to the content of the object mateys, we can just type it. When you assign a value to an object, R won't automatically print it. If you want to see the value, you need to call the object by just executing its name:

Tip: To execute code from the source to the console, highlight it and use the hot-keys "Command-Return" on Mac or "Control-Enter" on PC.

Do your future self a favor and use comments to explain what you're doing with your code. Also, maybe go for a run once in a while.

Good object names strike a balance between being easy to type (i.e.; short names) and interpret. If you have several datasets, it's probably not a good idea to name them a, b, c because you'll forget which is which. However, using long names like March2015Group10OnlyFemales and March2015Group10OnlyMales will give you carpal tunnel syndrome.

```
mateys # What is the value of mateys?
## [1] 50
```

A few notes about defining objects: you can't start the name of an object with a number, and you can't have spaces or other 'weird' characters in the name. Here are some examples of *invalid* object names:

```
me mateys <- 50 # Can't have spaces
5.mateys <- 50 # Can't start a name with a number
mayeys! <- 5 0# Can't have an "!" in the object name
```

R is case-sensitive. If you define an object with uppercase letters, you must keep referring to it with uppercase letters!

```
Plunder <- 1
plunder <- 100
Plunder
## [1] 1
plunder
## [1] 100
```

Avoid using too many capital letters in object names because they require you to hold the shift key. This may sound silly, but you'd be surprised how much easier it is to type mydata than MyData 100 times.

Once you've defined an object, you can use it in other commands:

```
a <- 100
b <- 2
c <- a + b
c
## [1] 102
```

If you want to change an object, you can just reassign it. You can even refer to the same object when reassigning it:

```
a <- 2
a <- a + a
a
## [1] 4
```

You can use = instead of <- for object assignment but I recommend you stick with <- because the direction of the assignment is clear.

## Data object types in R

R stores everything as an object, and there are different types of objects. The first two objects we'll learn about are scalars and vectors.

Later on, we'll talk about more complicated objects like matrices, dataframes, hypothesis tests, etc.

### Two simple data objects: scalars and vectors

Two of the most common data objects in R are **scalars** and **vectors**. Let's discuss each in turn,

#### scalars

A **scalar** is just a single value. A scalar can either be *numeric* or *string*. A numeric scalar is a number, while a string scalar is a letter. We denote string scalars by using quotation marks. Here are some examples:

```
a <- 1
b <- 3 * 40
ship <- "Black Pearl"
```

It is important to note that once you've defined an object, you refer to it without quotation marks, even if the object is a character. For example, to refer to the object `ship` that I defined above, you need to write `ship` without quotations marks

```
ship # Print the value of the object ship
```

```
## [1] "Black Pearl"
```

```
"ship" # R thinks this is a new string called "ship", not the object called ship
```

```
## [1] "ship"
```

```
# scalar v vector v matrix
par(mar = rep(1, 4))
plot(1, xlim = c(0, 4), ylim = c(-.5, 5),
     xlab = "", ylab = "",
     xaxt = "n", yaxt = "n",
     bty = "n", type = "n")

# scalar
rect(rep(0, 1), rep(0, 1), rep(1, 1), rep(1, 1))
text(.5, -.5, "scalar")

# Vector
rect(rep(2, 5), 0:4, rep(3, 5), 1:5)
text(2.5, -.5, "Vector")
```

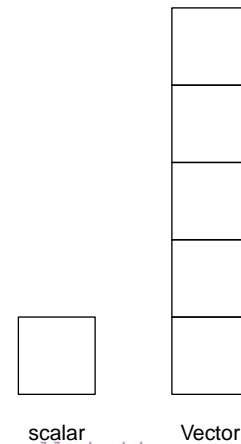


Figure 5: Visual depiction of a scalar and vector. Deep shit. Wait until we get to matrices - you're going to lose it.

#### Vectors

A **vector** is a combination of several scalars. For example, the numbers from one to ten could be a vector of length 10. Like scalars, you can also create character vectors that contain character scalars.

There are many ways to create vectors in R, here are the most common:

`c(a, b, c, ...)`

`a, b, c, ...`

One or more objects to be combined into a vector



The simplest way to create a vector is with the `c()` function. The `c` here stands for concatenate, which means "bring them together". When using `c()`, place a comma in between the objects (scalars or vectors) you want to combine:

The following code will create a vector of the integers from 1 to 5:

```
v <- c(1, 2, 3, 4, 5)
v
## [1] 1 2 3 4 5
```

`c(x, y, z)`: Create a vector with the `c()` command by separating elements with commas

You can also create vectors by combining objects you have already defined. Let's create a vector of the numbers from 1 to 10 by first generating a vector `a` from 1 to 5, and a vector `b` from 6 to 10 then combine them into a single vector `c`:

```
a <- c(1, 2, 3, 4, 5)
b <- c(6, 7, 8, 9, 10)
c <- c(a, b)
c
## [1] 1 2 3 4 5 6 7 8 9 10
```

You can also create string vectors containing only string values:

```
a <- c("this", "is", "a", "string", "vector")
a
## [1] "this" "is" "a" "string" "vector"
```

A vector can only contain one type of scalar: either numeric or character. If you try to create a vector with numeric and character scalars, then R will convert all of the numeric scalars to characters:

```
movie <- "Pirates of the Carribean"
revenue <- 634954111
c(movie, revenue) # Result is a string vector
## [1] "Pirates of the Carribean" "634954111"
```

Once you've created a vector, you can easily determine its length by using the `length()` function:

`length()`

```
length(c(1, 2, 3))
```

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

### *Generating numeric vectors*

While the `c()` operator is the most straightforward way to create a vector, it's also one of the most tedious. Let's say you wanted to create a vector of all integers from 1 to 100. You definitely don't want to have to type all the numbers into a `c()` operator. Instead, R has many simple built-in functions for generating numeric vectors. Let's start with three of them:

a:b	
a	The start of the sequence
b	The end of the sequence

The `a:b` function creates a vector of numbers from the starting point `a` to the ending point `b` in steps of 1:

```
1 : 10 # Integers from 1 to 10
```

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

```
10 : 1 # Integers from 10 to 1
```

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

```
20.1:30.1 # From 20.1 to 30.1
```

```
## [1] 20.1 21.1 22.1 23.1 24.1 25.1 26.1 27.1 28.1 29.1 30.1
```

## seq(from, to, by)

---

from

The start of the sequence

to

The end of the sequence

by

The step-size of the sequence

length.out

The desired length of the final sequence (only use if you don't specify by)

The `seq()` function allows you to create a sequence from a starting number to an ending number, in steps you specify. The function has three arguments, which are inputs to the function which changes how it works.

```
seq(from = 1, to = 10, by = 1)
## [1] 1 2 3 4 5 6 7 8 9 10
seq(from = 0, to = 100, by = 10)
## [1] 0 10 20 30 40 50 60 70 80 90 100
```

**seq(from, to, by)** - Creates a sequence between two numbers in steps that you specify.

from: The starting value

to: The ending value

by: The step size between begin and end

The `rep()` function allows you to repeat a number (or vector) a specified number of times. There are three arguments to the `rep` function:

## rep(x, times, each)

---

x

A scalar or vector of values to repeat

times

The number of times to repeat the sequence

each

The number of times to repeat each value within the sequence

```
rep(1:5, times = 2) # Repeat integers 1 to 5 two times
## [1] 1 2 3 4 5 1 2 3 4 5

rep(1:5, each = 2) # Repeat each integer from 1 to 5 two times
## [1] 1 1 2 2 3 3 4 4 5 5

rep(1:5, each = 2, times = 2) # Do both!
## [1] 1 1 2 2 3 3 4 4 5 5 1 1 2 2 3 3 4 4 5 5
```

**rep(x, times, each)** - Repeats the numbers in x in a manner you specify  
 times: The number of times the vector should be repeated  
 each: The number of times you want to repeat each element in the vector.

### *Arithmetic operations on scalars and vectors*

You can do basic arithmetic operations like +, -, \* and / on scalars and vectors. If you do an operation on a vector with a scalar, R will apply the scalar to each element in the vector:

```
a <- 1:5
a * 10
## [1] 10 20 30 40 50

a - 1
## [1] 0 1 2 3 4

a ^ 2
## [1] 1 4 9 16 25
```

If you do an operation on two vectors, R will try to apply the operation between the vectors by each item:

```
1:10 + 21:30
## [1] 22 24 26 28 30 32 34 36 38 40

(1:5) * (1:5)
## [1] 1 4 9 16 25

seq(10, 100, 10) + 1:10
## [1] 11 22 33 44 55 66 77 88 99 110
```

### *Additional Tips*

1. To get more tips on how good coding techniques, check out the R style guide at <http://adv-r.had.co.nz/Style.html>
2. For great blog articles on R, check out <http://www.r-bloggers.com/>

3. If you need to enter a lot of numeric data into R by hand you might want to use the `scan()` function. This function allows you to easily enter data using 10-key typing on a number pad. To do this, run the code `scan()` and then enter the data number by number. When you are finished, R will then print the appropriate code to store the data into a vector.
4. You can run several lines of code in one line by separating the code with the `;` key. For example, the following two chunks of code are the same:

```
a <- 1  
b <- 14  
c <- 67
```

```
a <- 1 ; b <- 14 ; c <- 67
```

However, I recommend you use the `;` key sparingly. If you get in the habit of trying to cram several lines of code in one line, your code will get cluttered and difficult to understand.



### 3: Sampling data and Descriptive Statistics

#### Chapter Goals

1. Learn functions for generating data from probability distributions: `rnorm()`, `runif()`, `sample()`
2. Learn functions for basic descriptive statistics: `mean()`, `median()`, `sd()`, `var()`, `min()`, `max()`

#### Sampling data from probability distributions

By now you know how to generate sequences of numbers with the functions `:`, `seq()`, and `rep()`. However, these functions don't generate very interesting data. Instead, we can use R to generate randomly sampled data from specified *probability distributions*. A probability distribution is simply an equation that indicates how likely certain numerical values are to be drawn. When you draw a *sample* of size N from a distribution, you are selecting N numerical values drawn according to that distribution's likelihood function.

For example, imagine you need to hire a new group of pirates for your crew. You have the option of hiring people from one of two different pirate training colleges that produce pirates of varying quality. One college "Pirate Training Unlimited" might tend to produce pirates that are generally ok - never great but never terrible. While another college "Unlimited Pirate Training" might produce pirates with a wide variety of quality, from very low to very high. In Figure 6 I plotted 5 example pirates from each college, where each pirate is shown as a ball with a number written on it. As you can see, pirates from PTU all tend to be clustered between 40 and 60 (not terrible but not great), while pirates from UPT are all over the map, from 0 to 100. We can use probability distributions (in this case, the uniform distribution) to

In the next section we'll go over some of the most commonly used sampling distributions: the Normal and Uniform distributions.

```
# Create blank plot
plot(1, xlim = c(0, 100), ylim = c(0, 100),
     xlab = "Pirate Quality", ylab = "", type = "n",
     main = "Two different Pirate colleges", yaxt = "n"
)

# Set colors
require("RColorBrewer")
col.vec <- brewer.pal(10, name = "Set3")[4:6]

# Draw Samples
samples.1 <- runif(n = 5, 40, 60)
samples.2 <- runif(n = 5, 0, 100)

text(50, 90, "Pirate Training Unlimited", font = 3)

for(i in 1:length(samples.1)) {
  points(samples.1[i], 75, pch = 21, bg = col.vec[1], cex = 3)
  text(samples.1[i], 75, round(samples.1[i], 0))
}

segments(40, 65, 60, 65, col = col.vec[1], lty = 1, lwd = 2)
text(50, 40, "Unlimited Pirate Training", font = 3)

for(i in 1:length(samples.2)) {
  points(samples.2[i], 25, pch = 21, bg = col.vec[2], cex = 3)
  text(samples.2[i], 25, round(samples.2[i], 0))
}

segments(10, 15, 90, 15, col = col.vec[2], lty = 1, lwd = 2)
```

Two different Pirate colleges

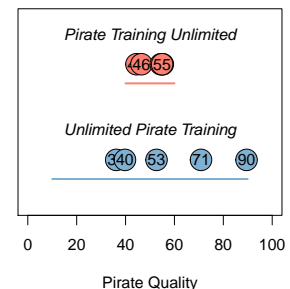


Figure 6: Sampling 5 potential pirates from two different pirate colleges. Pirate Training Unlimited (PTU) consistently produces average pirates (with scores between 40 and 60), while Unlimited Pirate Training (UPT), produces a wide range of pirates from 0 to 100.

### The Normal (Gaussian) distribution

Let's start with the most famous distribution in statistics: the Normal (or if you want to sound pretentious, the Gaussian) distribution. The Normal distribution is bell-shaped, and has two parameters: a mean and a standard deviation. See the margin figure 7 for plots of three different Normal distributions with different means and standard deviations.

To generate samples from a normal distribution in R, we use the function `rnorm()` this function has three arguments:

<code>rnorm()</code>	
<code>n</code>	The number of observations
<code>mean</code>	The mean of the Normal distribution from which samples are drawn (not the sample mean!!)
<code>sd</code>	The standard deviation of the Normal distribution from which samples are drawn

For example, let's draw 5 samples ( $n = 5$ ) from a normal distribution with mean 0 ( $\text{mean} = 0$ ) and standard deviation 1 ( $\text{sd} = 1$ )

```
rnorm(n = 5, mean = 0, sd = 1)
## [1] 0.375799 -1.367125 -1.289852 -1.430844 -2.457182
```

This code returns a vector of 5 values, where each value is a new random sample drawn from a Normal distribution with  $\text{mean} = 0$  and  $\text{standard deviation} = 1$ .

Because the sampling is done randomly, you'll get different values each time you run the `rnorm()` (or any other random sampling) function. To see this, let's create two different sets of samples from a normal distribution with mean 10 and standard deviation 5 and see how they compare:

```
a <- rnorm(5, mean = 10, sd = 5)
b <- rnorm(5, mean = 10, sd = 5)
a # print a
```

```
require("RColorBrewer")

# Create blank plot
plot(1, xlim = c(-5, 5), ylim = c(0, 1),
     xlab = "x", ylab = "dnorm(x)", type = "n",
     main = "Three Normal Distributions")

# Set up design matrix for loop
design.matrix <- data.frame("mean" = c(0, -2, 1),
                           "sd" = c(1, .5, 2))

# Set colors
col.vec <- brewer.pal(10, name = "Set3")[4:6]

# Start loop over distributions
for (i in 1:nrow(design.matrix)) {
  mean.i <- design.matrix$mean[i]
  sd.i <- design.matrix$sd[i]

  fun <- function(x) {
    dnorm(x, mean = mean.i, sd = sd.i)}

  curve(expr = fun,
        from = -5, to = 5,
        xlab = "x", lwd = 3,
        add = T, col = col.vec[i])

  samples <- rnorm(n = 10, mean = mean.i, sd = sd.i)

  segments(x0 = samples, y0 = rep(0, 10),
          x1 = samples, y1 = fun(samples),
          col = col.vec[i], lwd = 1, lty = 2)
}

legend.fun <- function(i) {
  paste("mean = ", design.matrix$mean[i],
        ", sd = ", design.matrix$sd[i], sep = "")}

legend("topright",
      legend = c(legend.fun(1),
                 legend.fun(2),
                 legend.fun(3)),
      lwd = rep(3, 3),
      col = col.vec[1:3])
```

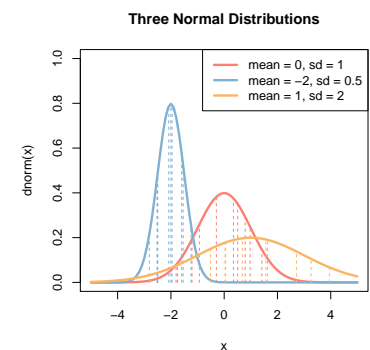


Figure 7: Three different normal distributions with different means and standard deviations.



```
## [1] 8.752808 9.329998 11.479484 16.710223 3.270555

b # print b

## [1] 7.159241 20.320119 13.405458 5.809321 5.997274
```

As you can see, even though I used the exact same code to generate the vectors *a* and *b*, the numbers in each sample are different. That's because the samples are each drawn randomly and independently from the normal distribution. To visualize the sampling process, run the code in the margin Figure 7 on your machine several times. You should see the sampling lines dance around the distribution.

### The Uniform distribution

Next, let's move on to the *uniform* distribution. The uniform distribution gives equal probability to all values between the minimum and maximum values.

To generate samples from a uniform distribution, we use the function `runif()`, the function has 3 arguments:

<code>runif()</code>	
<code>n</code>	The number of observations (i.e.; samples)
<code>min</code>	The lower bound of the Uniform distribution from which samples are drawn
<code>max</code>	The upper bound of the Uniform distribution from which samples are drawn

Let's draw 5 samples from two uniform distributions, one with bounds at 0 and 1, and one with bounds at -100 and 100:

```
runif(5, min = 0, max = 1) # 5 samples from U[0, 1]

## [1] 0.001454292 0.919849717 0.580778579 0.693033657 0.306065331

runif(5, min = -100, max = 100) # 5 samples from U[-100, 100]

## [1] 35.446931 53.111058 60.491567 -3.057751 13.683292
```

```
# uniform distribution
curve(dunif,
      from = 0, to = 1,
      xlim = c(-.5, 1.5),
      xlab = "x",
      lwd = 2,
      main = "Uniform\nmin = 0, max = 1")
```

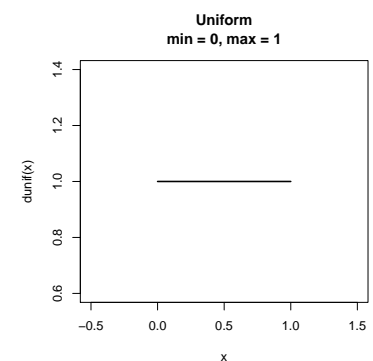


Figure 8: The Uniform distribution - known colloquially as the Anthony Davis distribution.

*Sampling from a set of values: sample()*

The next function we'll use is **sample()**. The sample function works a bit differently from `runif()` and `rnorm()` because it allows to you to define which values you want to sample and the probability associated with each value. For example, if you want to simulate the flip of a fair coin, you can tell the sample function to draw the value "Heads" with probability .50, and the value "Tails" with probability .50.

### sample()

---

**x**

A vector of outcomes you want to sample from. For example, to simulate coin flips, you'd enter `x = c("Heads", "Tails")`

**size**

The number of samples you want to draw.

**replace**

Should sampling be done with replacement? If T, then each individual sample will be replaced in the data vector. If F, then the same outcome will never be drawn more than once. Think about replacement like drawing different balls from a bag. Sampling with replacement (`replace = T`) means that each time you draw a ball, you return the ball back into the bag before drawing another ball. Sampling without replacement (`replace = F`) means that after you draw a ball, you remove that ball from the bag before drawing again.

**prob**

A vector of probabilities of the same length as `x` indicating how likely each outcome in "x" is. For example, to sample equally from two outcomes, you'd enter `prob = c(.5, .5)`. The first value corresponds to the first value of `x` and the second corresponds to the second value (etc.). The vector of probabilities you give as an argument should add up to one. However, if they don't, R will just rescale them so that they will sum to 1.

As a simple example, let's simulate 10 flips of a fair coin, where the probability of getting either a Head or Tail is .50:

```
sample(x = c("Heads", "Tails"), # The values you want to sample from
       size = 10, # The number of samples
       prob = c(.5, .5), # The probability of each value
       replace = T # Sampling with replacement
)

## [1] "Heads" "Heads" "Tails" "Heads" "Tails" "Tails" "Tails" "Tails"
## [9] "Tails" "Heads"
```

As you can see, our function returned a vector of 10 values corresponding to our sample size of 10. Keep in mind that, just like using `rnorm()` and `runif()`, the `sample()` function can give you different outcomes every time you run it.

### *Drawing coins from a treasure chest*

Now, let's sample drawing coins from a treasure chest. Let's say the chest has 100 coins: 20 gold, 30 silver, and 50 bronze. Let's draw 10 random coins from this chest. Because we remove coins when we draw them, we'll set `replace = F`.

```
# Create chest with the 100 coins

chest <- c(rep("gold", 20),
           rep("silver", 30),
           rep("bronze", 50)
          )

# Draw 10 coins from the chest without replacement

sample(x = chest,
       size = 10,
       prob = rep(1 / 100, times = 100),
       replace = F
)

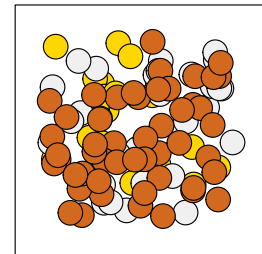
## [1] "bronze" "silver" "bronze" "bronze" "silver" "bronze" "bronze"
## [8] "bronze" "bronze" "silver"
```

The output of the `sample()` function above is a vector of 10 strings indicating the type of coin we drew on each sample. The order of these strings matter: the first one is the first coin we drew, and the last one is the 10th coin we drew. And like any random sampling function, this code will likely give you different results everytime you run it! See how long it takes you to get 10 gold coins...

```
par(mar = c(3, 3, 3, 3))
plot(1, xlim = c(0, 1), ylim = c(0, 1),
     xlab = "", ylab = "", xaxt = "n",
     yaxt = "n", type = "n",
     main = "Chest of 20 Gold, 30 Silver, and 50 Bronze Coins")

points(runif(100, .1, .9),
       runif(100, .1, .9),
       pch = 21, cex = 3,
       bg = c(rep("gold", 20),
              rep("gray94", 30),
              rep("chocolate", 50))
       )
```

**Chest of 20 Gold, 30 Silver,  
and 50 Bronze Coins**



### *Simulating Pinder Outcomes*

Let's simulate some Pinder outcomes. For those who don't know, Pinder is an app that allows Pirates to view profiles of potential dates. For each potential date, you can see their picture and either "like" them by swiping right, or "dislike" them by swiping left. If a pirate that you "liked" also "likes" you, then you've had a successful match and will be able to start chatting. let's say you "swipe right" on 20 Pinder profiles and the probability you get a match is 20%. We can simulate this using the sample function

```
sample(x = c("Match!!!", "No Match"),
       size = 20,
       replace = T, # Replace each sample back to the set
       prob = c(.2, .8) # Probability of Match! is .2, and No Match :( is .8)
)

## [1] "No Match" "Match!!!" "No Match" "Match!!!" "No Match" "No Match"
## [7] "No Match" "No Match" "No Match" "No Match" "No Match" "No Match"
## [13] "No Match" "Match!!!" "Match!!!" "No Match" "Match!!!" "No Match"
## [19] "No Match" "Match!!!"
```

The output of this function is a simulated response from 10 pirates that you liked.

### *Descriptive statistics*

Ok, now that we can generate some data, let's learn the basic descriptive statistics functions. We'll focus on the most common ones for numerical analyses.



Figure 9: A typical Pinder profile.

## Common Descriptive Statistics

---

`mean(x)`

The arithmetic mean of a vector `x`

`median(x)`

The median of a vector `x`. 50% of the data should be less than `median(x)` and 50% should be greater than `median(x)`.

`sd(x), var(x)`

The standard deviation and variance of a vector `x`.

`min(x), max(x)`

The minimum and maximum values of a vector `x`

`quantile(x, p)`

The `p`th sample quantile of a vector `x`. For example, `quantile(x, .2)` will tell you the value at which 20% of cases are less than `x`. The function `quantile(x, .5)` is identical to `median(x)`

`summary(x)`

Shows you several descriptive statistics of a vector `x`, including `min(x)`, `max(x)`, `median(x)`, `mean(x)`

Each of these functions takes a vector as an argument, and returns a scalar as a result. Let's calculate some descriptive statistics from some pirate related data. I'll create a vector called `data` that contains the number of tattoos from 10 random pirates

```
tattoos <- c(4, 50, 2, 39, 4, 20, 4, 8, 10, 100)
```

To calculate the mean of the data, we simply write:

```
mean(tattoos)
```

```
## [1] 24.1
```

The other descriptive statistics functions are just as easy: Let's test the `median()`, `sd()`, `min()`, and `max()` functions:

```
median(tattoos)
```

```
## [1] 9
```

```
sd(tattoos)
```

```
## [1] 31.32074
min(tattoos)
## [1] 2
max(tattoos)
## [1] 100
```

One important point about the descriptive statistics functions is that most (if not all) of them as a default will freak out if there is a missing (NA) value in the data. For example, the following code will return NA as a result because there is an NA value in the data vector:

```
mean(c(1, 5, NA, 2))
## [1] NA
```

Include the argument `na.rm = T` to ignore missing (NA) values when calculating a descriptive statistic.

To tell a descriptive statistic function to ignore missing (NA) values, include the argument `na.rm = T` in the function:

```
mean(c(1, 5, NA, 2), na.rm = T)
## [1] 2.666667
```

Now, the function will ignore NA and calculate the mean of the non-missing values. While this may seem trivial now (why did we include an NA value in the vector if we wanted to ignore it?!), it will become very important when we apply the function to large existing datasets that may contain missing values.

If you want to get many summary statistics from a vector, you can use the **summary()** function which gives you several key statistics:

```
summary(tattoos)
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      2.00   4.00   9.00  24.10  34.25  100.00
```

### *Other helpful vector functions*

Here are some other functions that you will find useful when managing numeric vectors:

### Other helpful numeric functions

`round(x, digits)`

Round values in a vector (or scalar) `x` to a certain number of digits.

`ceiling(x)`, `floor(x)`

Round a number to the next largest integer with `ceiling(x)` or down to the next lowest integer with `floor(x)`.

`x %% y`

Modular arithmetic (i.e.;  $x \bmod y$ ). You can interpret `x %% y` as "What is the remainder after dividing `x` by `y`?" For example, `10 %% 3` equals 1 because 3 times 3 is 9 (which leaves a remainder of 1).

### *A worked example: A quick test of the law of large numbers*

According to the law of large numbers, the larger our sample size, the closer our sample mean should be to the population mean. In other words, the more data (samples) you have, the more accurate your estimate should be. Let's test this by drawing either a small ( $N = 10$ ) or a large ( $N = 1,000,000$ ) number of observations from a Normal distribution with mean = 100 and sd = 20:

```
small <- rnorm(10, mean = 100, sd = 20) # 10 observations
large <- rnorm(1e6, mean = 100, sd = 20) # One million observations
```

Tip: You can easily write large powers of 10 by using the notation `1eN`, where `N` is the power of 10. For example: `1e6` is the same as 1,000,000

If our test worked, then the difference for the small sample should be larger than the large sample. Let's test this by calculating the mean of each sample and see how close they are to the true population mean of 100:

```
mean(small) # What is the mean of the small sample?
## [1] 97.75108

mean(large) # What is the mean of the large sample?
## [1] 100.0157

mean(small) - 100 # How far is the mean of Small from 100?
## [1] -2.248922

mean(large) - 100 # How far is the mean of Large from 100?
## [1] 0.01570442
```

Looks like the law of large numbers holds up!

*Additional Tips*



## 4: Indexing and comparing vectors

Chapter Goals:

1. Use brackets [] and logical vectors to index vectors
2. Combine indexing with descriptive statistics
3. Learn indexing functions which(), sort()
4. Vector discrete summary functions table() and unique()
5. Set functions: intersect(), union(), setdiff(),

### *Indexing vectors with brackets*

When we have a vector, we will frequently want to access specific values of a vector. These might be values in a specific location in the vector (i.e.; the fifth element) or based on some criteria (i.e.; all values greater than 0). We can accomplish this using indexing.

#### **Indexing with brackets [ ]**

To get the *i*th value of a vector called *vec*, use the bracket notation *vec[i]*

**vector[index]**

There are two main ways that you can use indexing to access subsets of data in a vector: numerical and logical indexing.

### *Numerical Indexing*

With numerical indexing, you enter the integers corresponding to the values in the vector you want to access in the form *data[num.index]*, where *data* is the data vector, and *num.index* is a vector of index values. For example, to get the first value in a vector, you'd write *data[1]*. To get the first, second, and third value, you can either type *data[c(1, 2, 3)]* or *data[1:3]*.

Let's do a few more examples. We'll use the *tattoos* vector again and use indexing to extract specific values:

```
tattoos <- c(0, 50, 2, 39, 9, 20, 17, 8, 10, 100)
tattoos[1] # First element of tattoos
```

```
## [1] 0

tattoos[1:5] # 1-5 elements of tattoos

## [1] 0 50 2 39 9
```

If you have defined an object that is a vector of integers, you can then index a variable using that vector. For example, let's define an object called `index` and use this object to index our data vector:

```
get.these.values <- 6:10
tattoos[get.these.values] # Indexing with a named object

## [1] 20 17 8 10 100
```

You can also get random values from a vector by indexing a vector with the `sample()` function. Let's get 3 random values from the `tattoo` vector in 2 steps. First, we'll create 3 random indexing values using `sample()`. Second, we'll index the `tattoo` object with the indexing values we generated in the first step.

```
rand.values <- sample(x = 1:length(tattoos), # Step 1: Determine indexing values
                     size = 3,
                     replace = F)

tattoos[rand.values] # Step 2: Index tattoo with rand.values

## [1] 20 10 50
```

The result of our indexing is 3 randomly selected values from the `tattoos` vector. Of course, we also could have done the same thing in one step by just entering `tattoos` as an argument to `sample()` like this:

```
sample(x = tattoos, size = 3, replace = F)

## [1] 50 39 0
```

As you gain more experience with R, you'll realise that there are many ways to program the same result. The choice of which code you use comes down to a delicate balance of readability (How easily can your future self, and other people, understand what the code is doing?), simplicity (How many lines of code are necessary?), and processing speed (How quickly will R complete the task?).

## Creating logical vectors

Another way to index data vectors is with logical vectors. A logical vector is a vector that only contains TRUE and FALSE values. If you index a vector with a logical vector (of the same length), you will only receive the values for which the index is TRUE.

You can create a logical vector by using the comparison operators in Figure 10.

Let's start by creating single scalar logical values so you can see how they work. If you apply a comparison operator to a scalar, R will return a single logical value of TRUE or FALSE. Let's see if 3 truly equals 3 and if 3 is really not greater than 5.

```
3 == 3

## [1] TRUE

3 > 5

## [1] FALSE
```

```
par(mar = rep(.1, 4))
plot(1, xlim = c(0, 1.1), ylim = c(0, 9),
     xlab = "", ylab = "", xaxt = "n", yaxt = "n",
     type = "n")

text(rep(0, 8), 8:1,
     labels = c("==", "!=", "<", "<=",
               ">", ">=", "|", "!"),
     adj = 0, cex = 3)

text(rep(.2, 8), 8:1,
     labels = c("equal", "not equal", "less than",
               "less than or equal", "greater than",
               "greater than or equal", "or", "not"),
     adj = 0, cex = 3)
```

<code>==</code>	equal
<code>!=</code>	not equal
<code>&lt;</code>	less than
<code>&lt;=</code>	less than or equal
<code>&gt;</code>	greater than
<code>&gt;=</code>	greater than or equal
<code> </code>	or
<code>!</code>	not

Figure 10: Comparison operators in R

```
# Create blank plot with no margins
par(mar = rep(0, 4))
plot(1, xlim = c(0, 1), ylim = c(0, 13),
     bty = "n", xlab = "", ylab = "", main = "",
     type = "n", xaxt = "n", yaxt = "n")

# Add Main title
text(.5, 12.5, "log.vec <- data.vec > 0", cex = 2)

# Data vector
text(.3, 11.1, "data.vec", font = 2, cex = 1.6)
data.vec <- c(2, 7, -1, 5, -9, -2, 3, 0, 2, -2)
text(rep(.3, 10), 10:1, data.vec, cex = 1.6)
rect(.25, .5, .35, 10.5)
segments(rep(.25, 9), seq(1.5, 9.5, 1),
         rep(.35, 9), seq(1.5, 9.5, 1), lty = 2)
```

The negation operator `!` meaning NOT. To use it, place the statement you are testing in parentheses, and place the `!` operator before it:

```
pirate <- "david"
pirate == "jack"

## [1] FALSE

!(pirate == "jack")

## [1] TRUE

!(2 == 4)

## [1] TRUE
```

In addition to using single comparison operators, you can combine multiple logical comparisons using the OR `|` and AND `&` commands. The OR command will return TRUE if any of the values in the set is TRUE, while the AND command will only return TRUE if all of the values in the set are TRUE.

```
(1 < 3) # Is 1 less than 3?

## [1] TRUE

(4 < 2) # Is 4 less than 2?

## [1] FALSE

(1 < 3) & (4 < 2) # Is 1 less than 3 and is 4 less than 2?

## [1] FALSE

(1 < 3) | (4 < 2) # Is 1 less than 3 OR is 4 less than 2?

## [1] TRUE
```

If you apply a comparison operator between a scalar and a vector, R will return a logical vector of length equal to the length of the vector. For example, let's compare a vector of integers from 1 to 10 to a scalar value of three and look at the output:

```
1:10 == 3 # Are the values equal to 3?

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

Let's look at the outputs above: for each value of the object `vec`, R performs the comparison `== 3`. Because only the third element of the vector is equal to 3, R returns the value FALSE for all values except the third one.

You can also compare two vectors of equal length and obtain a single logical vector as a result. For example, let's say we have two data vectors (`data.1` and `data.2`) and we want a logical vector telling us which values of the two data vectors are equal. We can do this by just executing `data.1 == data.2`

```
data.1 <- c(1, 4, 2, 3, 3)
data.2 <- c(1, 2, 4, 3, 3)
data.1 == data.2

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

`x %in% y`

One very important function for creating logical indices is `%in%`. This function looks a bit different from other functions because it doesn't follow the typical format of `function(x, y)`. Instead, you place the function `%in%` between its arguments. When you execute `x %in% y`, R will evaluate, for each element in the vector `x`, if it is in the vector `y`. For example, let's create several vectors `x` and `y` and use the `%in%` function to test whether or not the elements of `x` are in `y`:

```
1 %in% c(1, 2, 3, 4, 5)

## [1] TRUE
```

In this example, R returns a single value of `TRUE` because it found the value of `1` in the second vector. However, you can also apply the `%in%` function to a vector `x` that is longer than `1`. When you do this, the `%in%` function will return a vector equal to the length of `x`. Now, let's try an example where we test whether each of several values are in a second set:

```
c(1, 2, 3, 77, 88, 99) %in% c(1, 2, 3, 4, 5)

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

In this example R checked, for each of the values in the first vector if it was in the second vector (`c(1, 2, 3, 4, 5)`). Because only the first three values (`1, 2` and `3`) were in the second vector, R returns a vector with `3` `TRUE` values and `3` `FALSE` values.

The `%in%` function is very handy for seeing which values in a vector are valid according to a criteria you specify. For example, imagine you conducted a survey where you asked 10 different pirates how many siblings they had and received the following responses:

```
siblings <- c(3, 2, 0, -5, 0, -20, 2, 3, 1, -200)
```

Of course, the only valid answers to this question should be `0, 1, 2, ...` up to a maximum of say `20`; but some of these values appear to be invalid (that is, negative). Let's use the `%in%` function to see which values in the survey are valid. We'll create a vector called `valid.responses` that represents all possible valid answers to the question (we'll limit the number of siblings to `20`). We'll then use `%in%` to create a logical vector indicating which responses were valid.

```
siblings <- c(3, 2, 0, -5, 0, -20, 2, 3, 1, -200)
valid.responses <- seq(0, 20, 1)
siblings %in% valid.responses

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

Because the fourth, sixth, and tenth values were not valid (they were negative), the final logical vector gives us `FALSE` values at those index values, and `TRUE` values for all others.

## Indexing data with logical vectors

Once we have a logical vector, we can use that vector as an indexing vector. That is, you can use it to select values of a vector that satisfy some criteria you specify. To do this, you create a logical vector containing `TRUE` and `FALSE` values. If you then index a data vector (with the same length as the logical vector), R will return the values of the data vector for all `TRUE` values of the logical vector. See Figure to see visually how this works.

For example, let's say that we have the following set of data

```
# Create blank plot with no margins
par(mar = rep(0, 4))
plot(1, xlim = c(0, 1), ylim = c(0, 13),
     bty = "n", xlab = "", ylab = "", main = "",
     type = "n", xaxt = "n", yaxt = "n")

# Add Main title
text(.5, 12.5, "output.vec <- data.vec[log.vec]", cex = 2)

# Data vector
text(.2, 11.1, "data.vec", font = 2, cex = 1.6)
data.vec <- c(2, 7, -1, 5, -9, -2, 3, 0, 2, -2)
text(rep(.2, 10), 10:1, data.vec, cex = 1.6)
rect(.15, .5, .25, 10.5)
segments(rep(.15, 9), seq(1.5, 9.5, 1),
         rep(.25, 9), seq(1.5, 9.5, 1), lty = 2)
text(rep(.12, 10), 10:1, 1:10, cex = .8)

# Comparisons
text(rep(.32, 10), 1:10, "> 0", col = gray(.5))

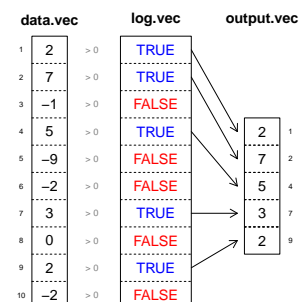
# Logical vector
text(.5, 11.1, "log.vec", font = 2, cex = 1.6)
index.text <- rep("FALSE", 10)
index.text[data.vec > 0] <- "TRUE"
col.vec <- rep("red", 10)
col.vec[data.vec > 0] <- "blue"
text(rep(.5, 10), 10:1,
     index.text,
     col = col.vec, cex = 1.6)

rect(.4, .5, .6, 10.5)
segments(rep(.4, 9), seq(1.5, 9.5, 1),
         rep(.6, 9), seq(1.5, 9.5, 1), lty = 2)

# Output vector
text(.8, 11.1, "output.vec", font = 2, cex = 1.6)
output.text <- data.vec[data.vec > 0]
text(rep(.8, 5), 7:3, output.text, cex = 1.6)
rect(.75, 2.5, .85, 7.5)
segments(rep(.75, 9), seq(3.5, 6.5, 1),
         rep(.85, 9), seq(3.5, 6.5, 1), lty = 2)
text(rep(.88, 5), 7:3, which(data.vec > 0), cex = .8)

# Arrows connecting log.vec to output.vec
arrows(rep(.6, 5),
       11 - which(data.vec > 0),
       rep(.73, 5),
       7:3, lwd = .5, length = .15)
)
```

```
output.vec <- data.vec[log.vec]
```



```
tattoos <- c(4, 50, 2, 39, 4, 20, 4, 8, 10, 100)
```

Now, let's say that we want to access just the data points that are less than 10. We'll start by first creating a logical indexing vector that tells us whether each value is less than 1. Then, we'll index the original data vector using this logical vector:

```
log.vec <- tattoos < 10 # Step 1: Create logical vector
tattoos[log.vec] # Step 2: Index the original data by the logical vector

## [1] 4 2 4 4 8
```

Logical vectors aren't just good for indexing, you can also use them to figure out which values in a vector satisfy some criteria. To do this, use the function `which()`

## which(log.vec)

If you apply the function `which()` to a logical vector, R will tell you which values of the index are TRUE. For example, let's create a logical vector and then see which index values are TRUE

```
log.vec <- c(TRUE, TRUE, FALSE, TRUE, FALSE)
which(log.vec)

## [1] 1 2 4
```

By using the `which()` function, we know that the first, second, and fourth elements of the logical vector are TRUE.

Let's take the example of comparing the treasure chest finding ability of 10 pirates. In each of two years - 2014 and 2015 - I measured how many chests 10 pirates found over the entire year. I recorded these values in two vectors, where the first value of each vector corresponds to the first pirate, and the last value corresponds to the last pirate:

```
pirate.names <- c("Andrew", "Heidi", "Madisen", "Becki", "Jack Dyanamite")
chests.2014 <- c(0, 10, 1, 2, 5)
chests.2015 <- c(0, 6, 3, 0, 20)
```

Ok, so let's see which pirates improved their chest findind ability. I'll start by finding the index values where the number of chests found increased between the two years

```
improve.log <- chests.2015 > chests.2014 # create logical vector
improve.log # print values

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

If I want to know the index values of the pirates who improved, I can use the `which()` function. The `which` function will tell me the index of each TRUE value in a logical vector:

```
which(improve.log)

## [1] 3 5
```

This vector tells us that the 3rd and 5th pirates found more chests in 2015 than 2014. Now I can use this index value to figure out the names of those pirates:

```
pirate.names[which(improve.log)]

## [1] "Madisen" "Jack Dyanamite"
```

Because you can index vectors with logical vectors, I could get the same results by just indexing `pirate.names` with `improve.log`.

```
pirate.names[improve.log]
## [1] "Madisen"      "Jack Dyanamite"
```

For this example, the `which()` command was unnecessary, but it's important to understand the logic of both methods.

### *Additional helpful vector functions*

Here are some other functions you might find useful when dealing with vectors:

#### Other Helpful Vector Functions

```
length(x)
  The length of a vector
sort(x)
  Sort a vector x. Add the argument decreasing = T to sort in decreasing order.
rev(x)
  Reverse the order of a vector x
unique(x)
  Determine all unique values in a vector x
table(x)
  Determine the number of counts for all unique values in a vector x
```

Once you have a vector of data, you may want to sort it in order to see, for example, the largest and smallest values. You can do this using the `sort()` function. Let's look back on my summer joke data and sort the results:

```
tattoos <- c(4, 50, 2, 39, 4, 20, 4, 8, 10, 100)
sort(tattoos, decreasing = T) # Sort decreasing
## [1] 100 50 39 20 10 8 4 4 4 2

sort(tattoos, decreasing = F) # Sort increasing
## [1] 2 4 4 4 8 10 20 39 50 100
```

You'll notice that the `sort` function has an argument `decreasing` which you can set to `TRUE` or `FALSE`.

The function `unique(x)` will tell you all the unique values in the vector, but won't tell you anything about how often each value occurs.

```
unique(c(1, 1, 2, 2, 2, 4, 500))
## [1] 1 2 4 500

unique(c("a", "A", "A", "A", "b", "b", "b", "c"))
## [1] "a" "A" "b" "c"
```

**unique(x):** Gives you all unique values in a vector, ignoring the number of times each value occurs.

The function `table()` does the same thing as `unique()`, but goes a step further in telling you how often each of the unique values occurs:

```
table(c(1, 1, 1, 2, 2, 5, 5, 700, 700, 1000))

##
##      1      2      5    700   1000
##      3      2      2      2      1

table(c("a", "A", "A", "A", "b", "b", "b", "c"))

##
## a A b c
## 1 3 3 1
```

**table(x):** Gives you all unique values in a vector and tells you how often each value occurs.

## Set Functions

R contains many functions that allow you to compare two sets (vectors) of data. See [margin Figure](#) for a visual depiction. Here are the most common ones:

### Set Functions

**union(x, y)**  
Tells you all unique values included in *either* the vector x or y.

**intersect(x, y)**  
Tells you all values common in *both* the vectors x and y.

**setdiff(x)**  
Tells you which values are in the vector x but *not* in the vector y. Keep in mind that `setdiff(x, y)` is *not* the same as `setdiff(y, x)`!

**setequal(x)**  
Returns TRUE if the two vectors x and y are identical (ignoring order) and FALSE if they are not identical.

```
require("plotrix")

## Loading required package: plotrix

require("RColorBrewer")

Transparent <- function(orig.col = "red", trans.val = 1, maxColorValue = 255) {
  if(length(orig.col) == 1) {orig.col <- col2rgb(orig.col)}
  if(!(length(orig.col) %in% c(1, 3))) {return(paste("length of original color must be 1 or 3"))}
  final.col <- rgb(orig.col[1], orig.col[2], orig.col[3], alpha = trans.val * 255)
  return(final.col)
}

color.vec <- brewer.pal(12, "Set3")
par(mar = rep(0, 4))
plot(1, xlim = c(0, 1), ylim = c(0, 1),
     bty = "n", xlab = "", ylab = "", main = "",
     type = "n", xaxt = "n", yaxt = "n")

draw.circle(x = .35, y = .5, radius = .35, col = Transparent(color.vec[4], .3), lty = 1)
draw.circle(x = .65, y = .5, radius = .35, col = Transparent(color.vec[5], .3), lty = 1)

text(.35, .1, "Set X", cex = 1.5)
text(.65, .1, "Set Y", cex = 1.5)

text(.5, .5, "intersect(x, y)")
text(.15, .5, "setdiff(x, y)")
text(.85, .5, "setdiff(y, x)")
text(.5, .9, "union(x, y)")
```

## Using indexing to remove specific values of a vector

Sometimes you might want to remove values of a vector before performing some analyses. This might be because some of the values are invalid or just not values that you want to include in your analyses. For example, let's say you asked 7 people how happy they were on a scale of 1 to 5 and received the following responses:

```
happy <- c(1, 4, 2, 999, 2, 3, -2)
```

As you can see, we have some invalid values (999 and -2) in this vector. We can use logical indexing to create a new vector called `happy.valid` that only contains values 1 through 5.

```
valid.log <- happy %in% c(1, 2, 3, 4, 5)
happy.valid <- happy[valid.log]
happy.valid

## [1] 1 4 2 2 3
```

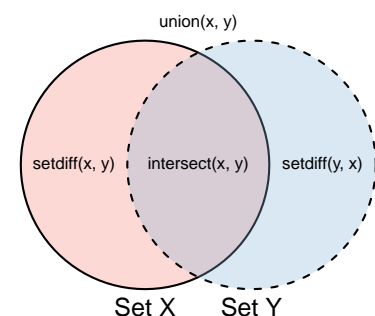


Figure 13: Common set functions in R.

As you can see, the new vector `happy.valid` only contains values from the original vector that are integers from 1 to 5.

R has special functions for testing whether or not values in a dataset are either missing (or infinite). Here are some you can use:

### Logical testing functions

`is.integer(x)`

Tests if values in a vector are integers

`is.na(x)`, `is.null(x)`

Tests if values in a vector are NA or NULL

`is.finite(x)`

Tests if a value is a finite numerical value. If a value is NA, NULL, Inf, or -Inf, `is.finite()` will return FALSE.

`duplicated(x)`

Returns FALSE at the first location of each unique value in `x`, and TRUE for all future locations of unique values. For example, `duplicated(c(1, 2, 1, 2, 3))` returns (FALSE, FALSE, TRUE, TRUE, FALSE). If you want to remove duplicated values from a vector, just run `x <- x[!duplicated(x)]`

You can use these functions to generate logical indices for indexing. For example, let's say you had a vector of data with several missing values. To create a new vector of data that does not contain the original NA values, we can index the original data vector with `is.finite(data)`:

```
data <- c(5, 2, NA, 3, NA, 10, NA)
data.finite <- data[is.finite(data)]
data.finite

## [1] 5 2 3 10
```

### *Taking the sum and mean of logical vectors to get counts and percentages*

Many (if not all) R functions that take numeric data as inputs will interpret TRUE values as 1 and FALSE values as 0. This allows us to easily answer questions like "How many values in a data vector are greater than 0?" or "What percentage of values are equal to 5?" by applying the `sum()` or `mean()` function to a logical vector.

Let's use this logic to see how many of the integers from 1 to 100 are greater than 0, 50, and 100:

```
sum(1:100 > 0) # How many values in 1:100 are greater than 0?

## [1] 100

sum(1:100 > 50) # How many values in 1:100 are greater than 50?

## [1] 50

sum(1:100 > 100) # How many values in 1:100 are greater than 100?

## [1] 0
```



These results should make sense: every value from 1:100 is greater than 0, 50 are greater than 50, and non are greater than 100. Now, let's do the same thing but calculate percentages instead of counts using `mean()` instead of `sum()`:

```
mean(1:100 > 0) # How many values in 1:100 are greater than 0?
## [1] 1
mean(1:100 > 50) # How many values in 1:100 are greater than 50?
## [1] 0.5
mean(1:100 > 100) # How many values in 1:100 are greater than 100?
## [1] 0
```

So far so good, now let's try this on our tattoo data:

```
tattoos <- c(4, 50, 2, 39, 4, 20, 4, 8, 10, 100)
```

Let's see how many of these 10 pirates have more than 10 tattoos. We'll do this in two steps; First, we'll create a logical vector indicating which values are greater than 10. Second, we'll take the sum of this logical vector. This will tell us how many TRUE values there are in the logical vector:

```
log.vec <- tattoos > 10 # Step 1: Which values are > 10?
sum(log.vec) # Step 2: How many TRUE values are there?
## [1] 4
```

Looks like 4 pirates have more than 10 tattoos. Now, let's test what percent of pirates have 5 tattoos or less. We'll do this by first creating the logical vector, and then calculating the `mean()` of this vector. We can do this because the mean of a vector of 0s and 1s is identical to the percentage of 1s:

```
log.vec <- tattoos <= 5 # Step 1: Which values are <= 5?
mean(log.vec) # Step 2: What percent of values are TRUE?
## [1] 0.4
```

Looks like 40% of pirates have 5 tattoos or less.

### Additional Tips

- If you have a vector of values and you want to know which values are duplicates of previous values, you can use the `uplicated` function. This function will go through the vector from beginning to end and tag the first unique instance of a value as TRUE and all repeated instances of a value as FALSE:

```
vec <- c("a", "b", "a", "a", "c")
uplicated(c("a", "b", "a", "a", "c"))
## [1] FALSE FALSE TRUE TRUE FALSE
```

If you want to remove duplicated values from a vector, you can just index the vector by `!duplicated`:

```
vec[!duplicated(vec)]
## [1] "a" "b" "c"
```

However, you can do the same thing with `unique()`!

To see what percentage of values are TRUE in a logical vector, just take the mean of the vector. For example, the command `mean(c(-1, -2, 1, 1) > 0)` will return 0.50, telling you that half of the values are positive.

## A worked example - Chicken Weights

A farmer is testing the effectiveness of three different diets on the weight gain of chickens. When they are born, 50 chicks are randomly assigned to one of 4 diets. Over several time periods, the farmer weighs each chicken. These data are contained in the dataset `ChickWeight`. Because the data are stored in a dataframe, which we haven't learned yet, we'll convert the four columns in the dataset to vectors as follows:

```
weights <- ChickWeight$weight
time <- ChickWeight$Time
chick <- as.numeric(paste(ChickWeight$Chick))
diet <- as.numeric(ChickWeight$Diet)
```

Let's answer 5 questions with these vectors:

1. What are the first 10 elements of the `weights` vector and the last 10 elements of the `weights` vector?

```
weights[1:10]
## [1] 42 51 59 64 76 93 106 125 149 171
weights[(length(weights) - 9):length(weights)]
## [1] 67 84 105 122 155 175 205 234 264 264
```

To answer the second question, I used the `length()` function to index index `weights` to go from 9 elements *before* the end of the vector, to the end of the vector.

2. Which chicks were given diets 1 and 2?

```
unique(chick[diet == 1])
## [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
unique(chick[diet == 2])
## [1] 21 22 23 24 25 26 27 28 29 30
```

3. What was the mean weight across all time periods separately for diets 3 and 4?

```
mean(weights[diet == 3])
## [1] 142.95
mean(weights[diet == 4])
## [1] 135.2627
```

First, I indexed the `weights` vector with a logical vector created from from `diet`. I then calculated the mean of this indexed vector.

4. What was the standard deviation of weights for diets 1 and 2 at time < 10?

```
sd(weights[diet <= 2 & time < 10])
## [1] 17.2321
```

5. What was the median weight for chicks 10, 20, and 30 for time periods greater than 10?

```
median(weights[chick %in% c(10, 20, 30) & time > 10])
## [1] 115
```

6. Which chicks did not make it until the final time period?

```
# Step 1: Create a vector of all chicks (all.chicks)
all.chicks <- sort(unique(chick))

# Step 2: Create a vector of all chicks that survive until the end (surviving.chicks)
surviving.chicks <- sort(unique(chick[time == max(time)]))

# Step 3: For each chick, see if it is present in the vector of surviving chicks
survived.log <- all.chicks %in% surviving.chicks
```

This one is a bit tricky. First, I need a vector of all chicks in the study (`all.chicks`). Next, I need a vector of all chicks that survived to the last time point (`surviving.chicks`). Third, I need to test, for each chick, whether they are present in the vector of surviving chicks (`survived.log`). Finally, I index the vector of all chicks where the logical index is `FALSE` (because we want chicks that did not survive).

```
# Step 4: Index the vector of all chicks by the logical vector  
all.chicks[survived.log == FALSE]  
## [1]  8 15 16 18 44
```



## 5: Matrices and Data Frames

### Chapter Goals

1. Learn about the matrix and dataframe data objects
2. Create matrices with `matrix()`, `cbind()`, and `data.frame()`
3. Index matrices/dataframes with brackets `[]`, and `$`
4. Use matrix/dataframe functions `dim()`, `nrow()`, `ncol()`, `head()`, and `tail()`
5. Import datasets

### Creating matrices and dataframes

By now, you should be comfortable with scalars and vectors. Next, we'll cover the next two most common data objects in R, **matrices** and **dataframes**

Matrices and dataframes are both two dimensional objects that contain rows and columns. Really, they're just like spreadsheets in Excel. Each matrix or dataframe contains a certain number of rows (call that number *m*) and columns (*n*). You can think of a matrix as a combination of *n* vectors, where each vector has a length of *m*. See Figure 14 to see the difference.

You can use several functions in R to create matrices and dataframes. In the next sections we'll cover the most common ones.

### `cbind()` and `rbind()`

`cbind()` and `rbind()` both create matrices by combining several vectors together into a single matrix. `cbind()` combines vectors as columns in the matrix, while `rbind()` combines them as rows.

```
# scalar v vector v matrix
par(mar = rep(1, 4))
plot(1, xlim = c(0, 10), ylim = c(-.5, 5),
     xlab = "", ylab = "",
     xaxt = "n", yaxt = "n",
     bty = "n", type = "n")

# scalar
rect(rep(0, 1), rep(0, 1), rep(1, 1), rep(1, 1))
text(.5, -.5, "scalar")

# Vector
rect(rep(2, 5), 0:4, rep(3, 5), 1:5)
text(2.5, -.5, "Vector")

# Matrix
rect(rep(4:8, each = 5),
     rep(0:4, times = 5),
     rep(5:9, each = 5),
     rep(1:5, times = 5))
text(6.5, -.5, "Matrix / Data Frame")
)
```

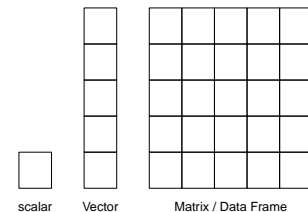


Figure 14: scalar, Vector, Matrix...  
::drops mike::

## cbind(), rbind()

---

`x, y, ...`

One or more vectors to be combined into a matrix

Let's use these functions to create a matrix with the numbers 1 through 30. First, we'll create three vectors of length 10, then we'll combine them into one matrix.

```
x <- 1:10
y <- 11:20
z <- 21:30

matrix.1 <- rbind(x, y, z)
matrix.1

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## x      1   2   3   4   5   6   7   8   9   10
## y     11  12  13  14  15  16  17  18  19  20
## z     21  22  23  24  25  26  27  28  29  30

matrix.2 <- cbind(x, y, z)
matrix.2

##           x  y  z
## [1,]   1 11 21
## [2,]   2 12 22
## [3,]   3 13 23
## [4,]   4 14 24
## [5,]   5 15 25
## [6,]   6 16 26
## [7,]   7 17 27
## [8,]   8 18 28
## [9,]   9 19 29
## [10,]  10 20 30
```

As you can see, the `rbind()` function combined the vectors as rows in the final matrix, while the `cbind()` function combined them as columns.

If you want to create a matrix from a single vector of data, you can do this using the `matrix()` function.

## matrix()

**data**

A vector of data

**nrow**

The number of rows in the final matrix

**ncol**

The number of columns in the final matrix

**byrow**

A logical value indicating whether to fill the matrix by row or column

Let's use the `matrix()` function to re-create a matrix containing the values from 1 to 30.

```
matrix.1 <- matrix(data = 1:30,
                   nrow = 10,
                   ncol = 3)

matrix.1

##      [,1] [,2] [,3]
## [1,]   1   11  21
## [2,]   2   12  22
## [3,]   3   13  23
## [4,]   4   14  24
## [5,]   5   15  25
## [6,]   6   16  26
## [7,]   7   17  27
## [8,]   8   18  28
## [9,]   9   19  29
## [10,]  10  20  30

matrix.2 <- matrix(data = 1:30,
                   nrow = 3,
                   ncol = 10)

matrix.2

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]   1   4   7  10  13  16  19  22  25  28
## [2,]   2   5   8  11  14  17  20  23  26  29
## [3,]   3   6   9  12  15  18  21  24  27  30
```

Keep in mind that matrices can either contain numbers or characters. If you try to create a matrix with both numbers and characters, it will turn all the numbers into characters:

```
cbind(1:5, c("a", "b", "c", "d", "e"))

##      [,1] [,2]
## [1,] "1"  "a"
## [2,] "2"  "b"
## [3,] "3"  "c"
## [4,] "4"  "d"
## [5,] "5"  "e"
```

*Dataframe: An  $m \times n$  object containing numbers, strings and factors*

A dataframe looks a lot like a matrix at first: it is also rectangular and has  $m$  rows and  $n$  columns. However, unlike matrices, dataframes can contain *both* string vectors and numeric vectors within the same object. For this reason, most large datasets in R, for example, a survey including numeric data and text data, will be stored as dataframes.

## data.frame()

To create a dataframe, you can use the `data.frame()` function. Let's create a dataframe of fictional survey data. I'll create 5 entries for Males and 5 entries for Females. I'll then generate 10 heights from a normal distribution with mean 150 and standard deviation 10.

```
survey <- data.frame("gender" = rep(c("Female", "Male"), each = 10),
                     "height" = rnorm(20, mean = 150, sd = 10),
                     stringsAsFactors = F # don't convert strings to factors
                     )
survey # Print the dataframe

##   gender  height
## 1 Female 157.3150
## 2 Female 155.3883
## 3 Female 144.3909
## 4 Female 145.0709
## 5 Female 155.6452
## 6 Female 149.5139
## 7 Female 154.8151
## 8 Female 137.8721
## 9 Female 146.3488
## 10 Female 135.1628
## 11 Male 163.7332
## 12 Male 149.0106
## 13 Male 138.2191
## 14 Male 140.5986
## 15 Male 153.2361
## 16 Male 151.7072
## 17 Male 151.0853
## 18 Male 160.3671
## 19 Male 143.5117
## 20 Male 143.4082
```

You'll notice I included the argument `stringsAsFactors = F`, this tells R to NOT convert the strings (the Gender column) to a factor

A dataframe is just a more flexible matrix that allows you to combine both character and numeric vectors into the same data object. Because dataframes are more flexible than matrices, Most datafiles you use will be stored as dataframes.



datatype. For now, don't worry about what factors are. Just know that you don't want to use them just yet!

### *Data sets pre-loaded in R*

Until now, we've used the functions `matrix()` and `dataframe()` to manually create our own datasets within R. However, for demonstration purposes, it's frequently easier to use existing datasets. Thankfully, R has us covered: R has several datasets that come pre-installed in a package called `datasets`. While you probably won't make any major scientific discoveries with these datasets, they allow all R users to test and compare code on the same sets of data. Here are a few datasets that we will be using in future examples:

- `ChickWeight`: Weight versus age of chicks on four different diets
- `InsectSprays`: Effectiveness of six different types of insect sprays
- `ToothGrowth`: The effects of different levels of vitamin C on the tooth growth of guinea pigs.

Since these datasets are preloaded in R, you can always access them by name. We'll use them in the following examples.

To see a complete list of all the datasets included in the `datasets` package, run the code: `library(help = "datasets")`

### *Viewing matrices and dataframes*

When you start working with a new dataset loaded as a matrix or dataframe, you'll usually want to get a quick visual look at it to make sure it looks ok. There are two functions that I use to do this:

#### `head(x)`

The function `head(x)` will show you the first few rows of a matrix / dataframe. Personally, I am constantly using this function to make sure that I didn't screw up a dataset when I'm working on it. Let's look at the first few rows of the dataframe `ChickWeight`, which contains data on the growth of chickens on several different diets.

```
head(ChickWeight)
```

```
##   weight Time Chick Diet
## 1     42    0     1    1
## 2     51    2     1    1
## 3     59    4     1    1
## 4     64    6     1    1
## 5     76    8     1    1
## 6     93   10     1    1
```

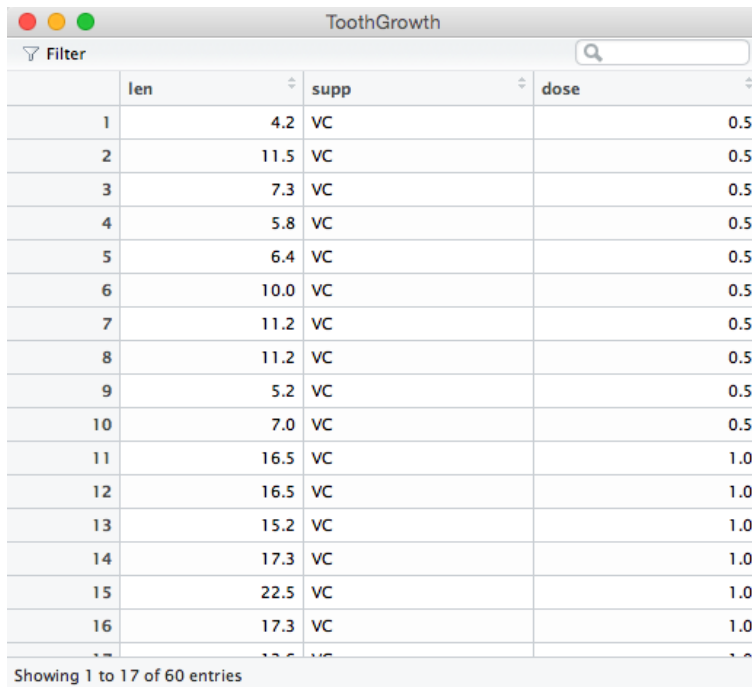
The `head()` function only shows you the first few rows of a dataframe, but usually this is enough to get a visual sense of the names of the dataframe, the number of columns, and the type of data in each column. But what if you want to see all values? You could print the entire dataframe into the console, but the console isn't a very friendly environment to view data. Instead, you can use the `View()` function, which will print the entire dataframe into a spreadsheet-like window:

## View(x)

Let's use the `View()` function to look at the entire `ToothGrowth` dataframe:

```
View(ToothGrowth)
```

When you run this code, you should see a separate window open (see Figure 15). You can use this window to scroll through the data, sort it via column values (by clicking on the column name), and even apply filters using the filter button on the top left of the screen. However, keep in mind that anything you do in the `View()` window will *not* change the actual dataframe in any way. You cannot add or remove data using the window, and any sorting or filtering you apply won't be replicated in the actual data.



	len	supp	dose
1	4.2	VC	0.5
2	11.5	VC	0.5
3	7.3	VC	0.5
4	5.8	VC	0.5
5	6.4	VC	0.5
6	10.0	VC	0.5
7	11.2	VC	0.5
8	11.2	VC	0.5
9	5.2	VC	0.5
10	7.0	VC	0.5
11	16.5	VC	1.0
12	16.5	VC	1.0
13	15.2	VC	1.0
14	17.3	VC	1.0
15	22.5	VC	1.0
16	17.3	VC	1.0
17	12.5	VC	1.0

Showing 1 to 17 of 60 entries

Figure 15: Screenshot of the window from `View(ToothGrowth)`. You can use this window to visually sort and filter the data to get an idea of how it looks, but you can't add or remove data and nothing you do will actually change the dataframe.

### Loading data into R with `read.table()`

So far we've used either randomly generated data, or datasets pre-loaded in R. But how do you get an existing dataset into R? For the most part, getting datasets into R isn't that tricky - but only if your data is already in a 'nice' format. By 'nice,' I mean a text file with tab (or comma) separated columns. If your data is in another format (like Excel or Shitty Piece of Shitty Shit), I strongly recommend first exporting the data to a tab-delimited text file, and only then loading the data into R. That said, if for some reason you absolutely have to load a non-text file into R, look at the *Additional Tips* sections for instructions.

Once you have a text file, you can load it into R using the `read.table()` function. To use the `read.table()` function, you need to know where the text file is located on your computer. To do this, find the file on your harddrive then right-click it and view its properties. You should be able to see its file-path there. For example, the file path of a text file called `mydata` on my desktop is `"Users/Nathaniel/Desktop/mydata.txt"`.

Here are the main arguments to `read.table()` (to see all of them, run `?read.table`)

Import data into R as comma or tab-delimited text files whenever possible. If you need to load data in another format (e.g.; Excel), save it as a text file from the original program first.

#### `read.table()`

##### `file`

The document's file path (make sure to enter as a string with quotation marks!) OR an html link to a file.

##### `header`

A logical value indicating whether the data has a header row or not.

##### `ncol`

The number of columns in the final matrix

##### `sep`

A string indicating how the columns are separated. For comma separated files, use `" , "`, for tab-delimited files, use `"\t"`

##### `stringsAsFactors`

A logical value indicating whether or not to convert strings to factors. I always set this to `FALSE` (because I don't like using factors)

To test this function, let's read in the datafile called `Flights.txt`. This dataset contains data on all flights leaving the Houston airport in 2011. You can access this data in one of two ways: First, you can download this file from: <http://nathanielphillips.com/wp-content/uploads/2015/04/Flights.txt>) and a note of its directory on your computer (on my computer, the path is `/Users/Nathaniel/Dropbox/Public/Flights.txt`). You can then load the data into R by using `read.table`:

```
Flights <- read.table(file = "/Users/Nathaniel/Dropbox/Public/Flights.txt",
                      header = T,
                      sep = "\t", # tab-delimited
                      stringsAsFactors = F
                      )
```

If you receive an error, it's probably because you entered the file path incorrectly. One trick to get the file path easily is by using RStudio's **Import Dataset** menu (see *Additional Tips*). If you got the directory location correct, and the file exists, then you should not receive any error warning after executing `read.table()`.

Alternatively, you can load the dataset directly into R by entering the HTML link as the file argument to `read.table`

```
Flights <- read.table(file = "http://nathanielphillips.com/wp-content/uploads/2015/04/Flights.txt",
                      header = T,
                      sep = "\t", # tab-delimited
                      stringsAsFactors = F
                      )
```

The data is now stored as a dataframe and you can now access it via the object name you assigned it to (in my case, I called it `Flights`). To make sure it loaded correctly, try seeing the first few rows with `head()`

```
head(Flights)
```

##		date	hour	minute	dep	arr	dep_delay	arr_delay	carrier
## 1	2011-01-01	12:00:00	14	0	1400	1500	0	-10	AA
## 2	2011-01-02	12:00:00	14	1	1401	1501	1	-9	AA
## 3	2011-01-03	12:00:00	13	52	1352	1502	-8	-8	AA
## 4	2011-01-04	12:00:00	14	3	1403	1513	3	3	AA
## 5	2011-01-05	12:00:00	14	5	1405	1507	5	-3	AA
## 6	2011-01-06	12:00:00	13	59	1359	1503	-1	-7	AA

##	flight	dest	plane	cancelled	time	dist
## 1	428	DFW	N576AA	0	40	224
## 2	428	DFW	N557AA	0	45	224
## 3	428	DFW	N541AA	0	48	224
## 4	428	DFW	N403AA	0	39	224
## 5	428	DFW	N492AA	0	44	224
## 6	428	DFW	N262AA	0	45	224

### *Additional tips*

- If you're like me, and you hate figuring out (and typing) the directory of a file, you can use RStudio's menu to help you. If you

click on the Environment window and click the button Import Dataset, you'll activate a menu that will allow you to select the file using your computer's finder. You'll then be greeted with a graphical interface for setting the import parameters. When you are finished, RStudio not only import the dataset, but it will paste the R code needed to import the data into the console. You can then copy the code (which includes the file path) and paste it into your R document so the next time you use the document you can just run the code to import the data.

- There are many functions other than `read.table()` for importing data. For example, the functions `read.csv` and `read.delim` are specific for importing comma-separated and tab-separated text files. In practice, these functions do the same thing as `read.table`, but they don't require you to specify a `sep` argument. Personally, I always use `read.table()` because it always works and I don't like trying to remember unnecessary functions.
- If you absolutely have to read a non-text file into R, check out the package called `foreign`. This package has functions for importing Stata, SAS and Shitty Piece of Shitty Shit files directly into R. To read Excel files, try the package `xlsx`



## 6: Basic Dataframe Manipulation

### Chapter Goals

1. Getting basic information about dataframes: `dim()`, `nrow()`, `ncol()`, `summary()`
2. Indexing dataframes with brackets `[],` and `$`
3. Subsetting dataframes with logical indexing and `subset()`
4. Recoding values in a dataframe with indexing

In this chapter we'll cover how to do some basic analyses on dataframes. We'll focus on dataframes, and not on matrices, because most datasets you use will be stored as dataframes. However, if you do find yourself working with matrices, many of the techniques you'll learn in this chapter will also apply to them.

### *Getting information about matrices and dataframes*

When you are working with dataframes, you will frequently want to know its general attributes, such as the number of rows and columns it has. Here are some common functions to get basic information about a dataframe:

- `dim(x)`: Number of rows and columns in a dataframe `x` (returns a vector of length 2)
- `nrow(x)` `ncol(x)`: How many rows or columns are there in a dataframe `x` (each function returns a scalar)
- `summary(x)`: Summary of information about each column in a dataframe `x`

```
dim(survey) # How many rows and columns?
```

```
## [1] 20  2
```

```
nrow(survey) # How many rows?
```

```
## [1] 20

ncol(survey) # How many columns?

## [1] 2

summary(survey) # Summary information on each column

##      gender      height
## Length:20      Min.   :135.2
## Class :character 1st Qu.:143.5
## Mode  :character Median :149.3
##                Mean   :148.8
##                3rd Qu.:155.0
##                Max.   :163.7
```

While you might not see the benefits of these functions now, they will become invaluable later if you conduct simulations on datasets.

Next let's start with the basics of indexing dataframes using brackets.

### *Indexing dataframes with brackets [rows, columns]*

Just like vectors, you can access specific data in dataframes using brackets. But now, instead of just using one indexing vector, we use two indexing vectors: one for the rows and one for the columns. To do this, use the notation `data[rows, columns]`, where rows and columns are scalars or vectors of the row and column numbers you want to get.

Let's try this on the `ChickWeight` dataframe that comes with R.

```
ChickWeight[1:5, 1] # Give me rows 1 through 5 in column 1

## [1] 42 51 59 64 76

ChickWeight[2:6, 2:3] # Give me rows 2 through 6 in columns 2 and 3

##      Time Chick
## 2      2      1
## 3      4      1
## 4      6      1
## 5      8      1
## 6     10      1

ChickWeight[seq(from = 1, to = nrow(ChickWeight), by = 10), 3] # Give me every 10th row in the 3rd column

## [1] 1 1 2 3 4 5 6 6 7 8 9 10 11 11 12 13 14 15 16 17 19 20 21
## [24] 21 22 23 24 25 26 26 27 28 29 30 31 31 32 33 34 35 36 36 37 38 39 40
## [47] 41 41 42 43 44 45 46 47 47 48 49 50
## 50 Levels: 18 < 16 < 15 < 13 < 9 < 20 < 10 < 8 < 17 < 19 < 4 < ... < 48
```

If you want an entire row or column, you can simply leave one of the indices blank. For example, if I want the entire first row of `ChickWeight` and all of the columns, I can simply leave the column index blank:



Figure 16: I tried ordering a vegan alternative at a traditional German restaurant and this is what I got.

Here you can see the benefits of using `nrow()` - I used it to make sure I gave valid index values to `ChickWeight`



```
ChickWeight[1,]

##   weight Time Chick Diet
## 1     42   0     1     1
```

You can use the same logic to get an entire column of a dataframe by leaving the index for rows blank. If you leave both index values blank, you'll get the entire dataframe back.

### *Accessing dataframe columns by column name and \$*

One of the nice things about dataframes is that each column will have a name. You can then use this name to access specific columns without having to index columns by numbers. To access the names of a dataframe, use the function `names()`. This will return a string vector with the names of the dataframe.

Let's use `names()` to get the names of some of the dataframes stored in R.

```
names(ChickWeight)

## [1] "weight" "Time"   "Chick"  "Diet"

names(InsectSprays)

## [1] "count" "spray"

names(ToothGrowth)

## [1] "len"   "supp"  "dose"
```

To access a specific column in a dataframe by name, you use the the `$` operator:

### **dataframe\$colname**

where `dataframe` is the name of the dataframe, and `colname` is the name of the column you are interested in. When you apply the `$` operator to a dataframe, it will return a vector. Let's access some of the vectors in the dataframes `ChickWeight`, and `ToothGrowth`

```
ChickWeight$weight[1:20] # Just the first 20 values

## [1] 42 51 59 64 76 93 106 125 149 171 199 205 40 49 58 72 84
## [18] 103 122 138

ToothGrowth$len

## [1] 4.2 11.5 7.3 5.8 6.4 10.0 11.2 11.2 5.2 7.0 16.5 16.5 15.2 17.3
## [15] 22.5 17.3 13.6 14.5 18.8 15.5 23.6 18.5 33.9 25.5 26.4 32.5 26.7 21.5
## [29] 23.3 29.5 15.2 21.5 17.6 9.7 14.5 10.0 8.2 9.4 16.5 9.7 19.7 23.3
## [43] 23.6 26.4 20.0 25.2 25.8 21.2 14.5 27.3 25.5 26.4 22.4 24.5 24.8 30.9
## [57] 26.4 27.3 29.4 23.0
```

Because the `$` operator returns a vector, you can easily calculate descriptive statistics on columns of a dataframe using `$`:

```
mean(ChickWeight$weight)

## [1] 121.8183

median(ToothGrowth$len)

## [1] 19.25
```

### *Adding new columns to a dataframe*

You can easily add columns to a dataframe using the `$` and assignment `<-` operators. To do this, just use the `dataframe$colname` notation and assign a new vector to it. Let's test this by adding a new column to `ChickWeight` called `height` which indicates the height of each chick at each time point. Now, because I don't have the actual height data, I'll just sample some random data from a normal distribution with mean 5 and standard deviation 1.

```
ChickWeight$height <- rnorm(n = nrow(ChickWeight), mean = 5, sd = 1)
```

Let's look at the first few rows of `ChickWeight` to make sure it worked:

```
head(ChickWeight)

##   weight Time Chick Diet   height
## 1    42    0     1    1 5.617418
## 2    51    2     1    1 4.719336
## 3    59    4     1    1 5.475568
## 4    64    6     1    1 4.449876
## 5    76    8     1    1 4.742738
## 6    93   10     1    1 5.641287
```

As you can see, the `ChickWeight` dataframe now has a column named `height` with the random data we generated.

Of course, because columns of dataframes are just vectors, you can add columns that are some function of existing columns. For example, let's add a column to `ChickWeight` called `density` which we'll define as each chick's height divided by its weight:

```
ChickWeight$density <- ChickWeight$height / ChickWeight$weight
head(ChickWeight)

##   weight Time Chick Diet   height   density
## 1    42    0     1    1 5.617418 0.13374805
## 2    51    2     1    1 4.719336 0.09253599
## 3    59    4     1    1 5.475568 0.09280624
## 4    64    6     1    1 4.449876 0.06952931
## 5    76    8     1    1 4.742738 0.06240445
## 6    93   10     1    1 5.641287 0.06065900
```

When you are conducting analyses on dataframes, it's important that you always repeat the name of the dataframe when accessing its columns. If you don't, R will assume the column name is a totally different object. For example, the following code *won't work* because R thinks that height and weight are totally separate objects from ChickWeight

```
ChickWeight$density <- height / weight # BAD CODE!
```

However, there is a function `with()` that can help prevent you from having to repeat the name of a dataframe over and over again.

### `with(x, ...)`

The function `with()` allows you to specify a dataframe (or any other object in R) once, and R will assume you're referring to that object in an expression.

For example, let's repeat the `ChickWeight$density` calculation using `with()`. We'll set the name of the dataframe as the first argument, then do our regular calculations on the column names.

```
ChickWeight$density <- with(ChickWeight, height / weight)
```

`with(x, ...)`: Simplifies your code for dataframe manipulation by allowing you to just enter the name of the dataframe once.

As you can imagine, if you're performing a long set of calculations on many columns of a dataframe, the `with()` function can save you lots of typing!

### *Centering and standardizing (z-score) data*

Centering and standardizing are two common methods of transforming data. Centering data simply means transforming the data so that the mean is 0, while standardizing data means centering the data and dividing all data points by the standard deviation of the data. Here's how to do each:

#### Centering

Centering data is quite easy. All you need to do is calculate the mean of a vector, then subtract that mean from all data in the vector.

Generally, if we have a dataframe called `df`, and we want to center a column called `x`, we'd run the following code:

```
df$x.centered <- with(df, x - mean(x))
```

Let's use this method to center the weight data from `ChickWeight` - we'll call the new column `weight.c`

```
ChickWeight$weight.c <- with(ChickWeight, weight - mean(weight))
```

To see if this worked, let's compare the mean of `weight` and `weight.c`

```
mean(ChickWeight$weight)
## [1] 121.8183
mean(ChickWeight$weight.c)
## [1] -1.219277e-15
```

I know what you're thinking..."But wait!!! The mean of `ChickWeight$weight.c` isn't exactly 0!!!" Don't worry,  $-1.2919 \times 10^{-15}$  is  $-0.00000000000000012919$ . For all intents and purposes, that's equal to 0 - the reason it's not *exactly* 0 is due to peculiarities about how computers represent numbers. Don't ask me why, I'm just a pirate.

#### Standardizing

Standardizing data is almost as easy as centering. The only difference is that, in addition to subtracting the mean from the data, we need to divide the data by its standard deviation. If you have a dataframe `df` and you want to standardize a column `x` into a new column called `x.z`, we use the following code:

```
# Create a standardized version of column x in a dataframe df
df$x.z <- with(df, (x - mean(x)) / sd(x))
```

Let's use this method to standardize the weight data from `ChickWeight` - we'll call the new column `weight.z`

```
ChickWeight$weight.z <- with(ChickWeight, (weight - mean(weight)) / sd(weight))
```

To see if this worked, let's compare the mean of `weight`, `weight.c`, and `weight.z`. The mean of `weight.z` should be 0 and its standard deviation should be 1:

```
c(mean(ChickWeight$weight), sd(ChickWeight$weight))
## [1] 121.81834 71.07196
c(mean(ChickWeight$weight.c), sd(ChickWeight$weight.c))
## [1] -1.219277e-15 7.107196e+01
c(mean(ChickWeight$weight.z), sd(ChickWeight$weight.z))
## [1] -1.194348e-17 1.000000e+00
```

### Subsetting dataframes with logical indexing and subset()

Frequently you will want to access specific rows of a dataframe based on some criteria - this is called subsetting. For example, we may want to look just at the data from females in our survey data. To do this, we can use one of two methods: indexing with logical vectors, or the `subset()` function.

Indexing dataframes with logical vectors is very similar to indexing data vectors. First, we create a logical vector. Next, we index the dataframe using that logical vector. Let's use indexing to access just the data from Chick 1 in `ChickWeight`:

```
chick1.log <- ChickWeight$Chick == 1 # Step 1: Create a logical vector
chick1.data <- ChickWeight[chick1.log,] # Step 2: Index dataframe by logical vector
chick1.data # Print the result
```

	weight	Time	Chick	Diet	height	density	weight.c	weight.z
## 1	42	0	1	1	5.617418	0.13374805	-79.818339	-1.12306372
## 2	51	2	1	1	4.719336	0.09253599	-70.818339	-0.99643150
## 3	59	4	1	1	5.475568	0.09280624	-62.818339	-0.88386952
## 4	64	6	1	1	4.449876	0.06952931	-57.818339	-0.81351829
## 5	76	8	1	1	4.742738	0.06240445	-45.818339	-0.64467533
## 6	93	10	1	1	5.641287	0.06065900	-28.818339	-0.40548114
## 7	106	12	1	1	3.082565	0.02908081	-15.818339	-0.22256793
## 8	125	14	1	1	3.784443	0.03027554	3.181661	0.04476675
## 9	149	16	1	1	6.087122	0.04085317	27.181661	0.38245267
## 10	171	18	1	1	5.127636	0.02998618	49.181661	0.69199810
## 11	199	20	1	1	6.022833	0.03026549	77.181661	1.08596500
## 12	205	21	1	1	4.743539	0.02313921	83.181661	1.17038648

If you'd like, you can also combine the two steps in one line. For example, the following code gives the same result as the previous:

```
chick1.data <- ChickWeight[ChickWeight$Chick == 1, ] # Two steps in 1
```

Now, let's try indexing the `ChickWeight` data using a slightly more complicated index. For example, let's access just the data for the chicks where `Time` is less than 10

```
ChickWeight.lt10 <- ChickWeight[ChickWeight$Time < 10,]
head(ChickWeight.lt10)
```

	weight	Time	Chick	Diet	height	density	weight.c	weight.z
## 1	42	0	1	1	5.617418	0.13374805	-79.81834	-1.1230637
## 2	51	2	1	1	4.719336	0.09253599	-70.81834	-0.9964315
## 3	59	4	1	1	5.475568	0.09280624	-62.81834	-0.8838695
## 4	64	6	1	1	4.449876	0.06952931	-57.81834	-0.8135183
## 5	76	8	1	1	4.742738	0.06240445	-45.81834	-0.6446753
## 13	40	0	2	1	5.495245	0.13738111	-81.81834	-1.1512042

Indexing with brackets is the standard way to slice and dice dataframes. However, if you are working on data that is all in the same dataframe, it can get a bit tiresome to have to constantly repeat the name of the dataframe. For example, let's say we wanted

to get data from `ChickWeight` where `Time < 10` and `Diet == 1`. The following code works, but it's a bit tedious:

```
mydata <- ChickWeight[ChickWeight$Time < 10 & ChickWeight$Diet == 1,]
```

A way to get around having to repeat the name of the dataframe over and over is to use the `subset()` function.

## subset()

`x`

The data (usually a dataframe)

`subset`

A logical vector indicating which rows you want to select

`select`

An optional vector of the columns you want to select

For example, let's get the `ChickWeight` data for `Diet == 1` and `Time > 15`

```
data <- subset(x = ChickWeight,
               subset = (Diet == 1 & Time > 15)
               )
head(data)
```

##	weight	Time	Chick	Diet	height	density	weight.c	weight.z
## 9	149	16	1	1	6.087122	0.04085317	27.18166	0.3824527
## 10	171	18	1	1	5.127636	0.02998618	49.18166	0.6919981
## 11	199	20	1	1	6.022833	0.03026549	77.18166	1.0859650
## 12	205	21	1	1	4.743539	0.02313921	83.18166	1.1703865
## 21	162	16	2	1	4.671985	0.02883942	40.18166	0.5653659
## 22	187	18	2	1	4.280018	0.02288780	65.18166	0.9171220

In the example above, I didn't specify an input to the `select` argument because I wanted all columns. However, if you just want certain columns, you can just name the columns you want. For example, let's say I just want the `weight` and `Time` columns from the previous analysis. To do this, I'll just add the column names as inputs to the `select` argument:

```
data <- subset(x = ChickWeight,
               subset = (Diet == 1 & Time > 15),
               select = c(weight, Time)
               )
head(data)
```

##	weight	Time
----	--------	------

```
## 9      149    16
## 10     171    18
## 11     199    20
## 12     205    21
## 21     162    16
## 22     187    18
```

### *Combining indexing and descriptive statistics*

Once you know how to index a dataframe to get the data vectors you want, you can then easily calculate descriptive statistics based on specific criteria. For example, let's calculate the mean weight of the chicks on `Diet == 1`. To show you that there are many ways to do this, I'll write the code in three different ways:

```
# What is the mean weight of chicks on the first diet?

mean(ChickWeight$weight[ChickWeight$Diet == 1]) # Using logical indexing

## [1] 102.6455

with(ChickWeight, mean(weight[Diet == 1])) # Logical indexing and with()

## [1] 102.6455

mean(subset(x = ChickWeight, subset = Diet == 1)$weight) # Using subset()

## [1] 102.6455
```

As you can see, there are many ways to do the same thing in R. Ultimately, the choice of which specific code and functions you use is up to you.

### *A worked example: Credit default*

For this example, we'll work with a dataset called `credit`. This dataset contains information about German loan borrowers (IVs) and whether or not the borrower defaulted on their loan (DV). To load the dataset, either download and load the data from the link <http://goo.gl/a7umut>, or simply run the following code:

```
credit <- read.table("http://nathanielphillips.com/wp-content/uploads/2015/05/credit.csv",
                     sep = ",", header = T, stringsAsFactors = F)
```

Here is a screenshot of the dataset:

```
View(credit)
```

The dataset has 17 total columns - to see their names, execute `names(credit)`

	checking_balance	months_loan_duration	credit_history	purpose	amount	savings_balance	employment_duration	percent_of_income	years_at_residence	age	other_credit	housing	existing_loans_count	job	dependents	phone	default
1	< 0 DM	6	critical	furniture/appliances	1189	unknown	> 7 years	4	4	57	none	own	1	skilled	1	yes	no
2	1 - 200 DM	48	good	furniture/appliances	5951	< 100 DM	1 - 4 years	2	2	22	none	own	1	skilled	1	no	yes
3	unknown	12	critical	education	2096	< 100 DM	4 - 7 years	2	3	49	none	own	1	unskilled	2	no	no
4	< 0 DM	42	good	furniture/appliances	7882	< 100 DM	4 - 7 years	2	4	40	none	other	1	skilled	2	no	no
5	< 0 DM	24	poor	car	4870	< 100 DM	1 - 4 years	3	4	53	none	other	2	skilled	2	no	yes
6	unknown	36	good	education	9055	unknown	1 - 4 years	2	4	35	none	other	1	unskilled	2	yes	no
7	unknown	24	good	furniture/appliances	2857	100 - 1000 DM	> 7 years	3	4	50	none	own	1	skilled	1	no	no
8	1 - 200 DM	36	good	car	6948	< 100 DM	1 - 4 years	2	2	35	none	rent	1	management	1	yes	no
9	unknown	12	good	furniture/appliances	3059	< 100 DM	4 - 7 years	2	4	61	none	own	1	unskilled	1	no	no
10	1 - 200 DM	36	critical	car	5234	< 100 DM	unemployed	4	2	26	none	own	2	management	1	no	yes
11	1 - 200 DM	12	good	car	1295	< 100 DM	< 1 year	3	1	25	none	rent	1	skilled	1	no	yes
12	< 0 DM	48	good	business	4308	< 100 DM	< 1 year	3	4	24	none	rent	1	skilled	1	no	yes
13	1 - 200 DM	12	good	furniture/appliances	1567	< 100 DM	1 - 4 years	1	1	22	none	own	1	skilled	1	yes	no
14	< 0 DM	24	critical	car	1199	< 100 DM	> 7 years	4	4	60	none	own	2	unskilled	1	no	yes
15	< 0 DM	15	good	car	1403	< 100 DM	1 - 4 years	2	4	28	none	rent	1	skilled	1	no	no
16	< 0 DM	24	good	furniture/appliances	1282	100 - 500 DM	1 - 4 years	4	2	32	none	own	1	unskilled	1	no	yes
17	unknown	24	critical	furniture/appliances	2454	unknown	> 7 years	4	4	50	none	own	2	skilled	1	no	no
18	< 0 DM	30	perfect	business	8072	unknown	< 1 year	2	3	25	bank	own	3	skilled	1	no	no
19	1 - 200 DM	24	good	car	13279	< 100 DM	> 7 years	4	2	44	none	other	1	management	1	yes	yes
20	unknown	24	good	furniture/appliances	3450	100 - 1000 DM	> 7 years	3	2	31	none	own	1	skilled	2	yes	no
21	unknown	9	critical	car	2134	< 100 DM	1 - 4 years	4	4	48	none	own	3	skilled	1	yes	no
22	< 0 DM	6	good	furniture/appliances	2647	100 - 1000 DM	1 - 4 years	2	3	44	none	rent	1	skilled	2	no	no
23	< 0 DM	18	critical	car	2241	< 100 DM	< 1 year	1	3	48	none	rent	2	unskilled	2	no	no
24	1 - 200 DM	12	critical	car	1864	100 - 500 DM	< 1 year	3	4	44	none	own	1	skilled	1	no	no
25	unknown	10	critical	furniture/appliances	2069	unknown	1 - 4 years	2	1	26	none	rent	2	skilled	1	no	no
26	< 0 DM	6	good	furniture/appliances	1374	< 100 DM	1 - 4 years	3	2	36	bank	own	1	unskilled	1	yes	no
27	unknown	6	perfect	furniture/appliances	426	< 100 DM	> 7 years	4	4	39	none	own	1	unskilled	1	no	no
28	> 200 DM	12	very good	furniture/appliances	409	< 100 DM	1 - 4 years	3	3	42	none	rent	2	skilled	1	no	no
29	1 - 200 DM	7	good	furniture/appliances	2415	< 100 DM	1 - 4 years	3	2	34	none	own	1	skilled	1	no	no
30	< 0 DM	60	poor	business	6806	< 100 DM	> 7 years	3	4	65	none	own	2	skilled	1	yes	yes
31	1 - 200 DM	18	good	business	1913	< 100 DM	< 1 year	3	3	38	bank	own	1	skilled	1	yes	no
32	< 0 DM	24	good	furniture/appliances	4020	< 100 DM	1 - 4 years	2	2	27	other	own	1	skilled	1	no	no
33	1 - 200 DM	18	good	car	5860	100 - 1000 DM	1 - 4 years	2	2	30	none	own	2	skilled	1	yes	no
34	unknown	12	critical	business	1264	unknown	> 7 years	4	4	57	none	rent	1	unskilled	1	no	no
35	> 200 DM	12	good	furniture/appliances	1474	< 100 DM	< 1 year	4	1	33	bank	own	1	management	1	yes	no
36	1 - 200 DM	45	critical	furniture/appliances	4748	< 100 DM	< 1 year	4	2	20	none	own	2	unskilled	1	no	yes
37	unknown	48	critical	education	6110	< 100 DM	1 - 4 years	1	3	31	bank	other	1	skilled	1	yes	yes
38	> 200 DM	18	good	furniture/appliances	2100	< 100 DM	1 - 4 years	4	2	37	other	own	1	skilled	1	no	yes
39	> 200 DM	10	good	furniture/appliances	1221	< 100 DM	1 - 4 years	2	2	37	none	own	1	skilled	1	yes	no
40	1 - 200 DM	9	good	furniture/appliances	458	< 100 DM	1 - 4 years	4	3	24	none	own	1	skilled	1	no	no
41	unknown	30	poor	furniture/appliances	2331	300 - 1000 DM	> 7 years	4	2	30	bank	own	1	management	1	no	no

Figure 17: Screenshot of the credit dataset.

```
names(credit)

## [1] "checking_balance" "months_loan_duration" "credit_history"
## [4] "purpose" "amount" "savings_balance"
## [7] "employment_duration" "percent_of_income" "years_at_residence"
## [10] "age" "other_credit" "housing"
## [13] "existing_loans_count" "job" "dependents"
## [16] "phone" "default"
```

Let's answer 5 questions with this dataset:

1. Was there a relationship between the size of the loan and whether or not it defaulted?

```
amount.default <- credit$amount[credit$default == "yes"]
amount.nodefault <- credit$amount[credit$default == "no"]

summary(amount.default)

##      Min. 1st Qu.  Median      Mean 3rd Qu.     Max.
##      433   1352   2574     3938   5142   18420

summary(amount.nodefault)

##      Min. 1st Qu.  Median      Mean 3rd Qu.     Max.
##      250   1376   2244     2985   3635   15860
```

```
amount.default <- credit$amount[credit$default == "yes"]
amount.nodefault <- credit$amount[credit$default == "no"]

require(beanplot)

## Loading required package: beanplot

beanplot(amount ~ default, data = credit,
  col = c("white", gray(.8), gray(.8), "black"),
  names = c("No", "Yes"),
  main = "Loan size by default",
  xlab = "Did the loan default?",
  ylab = "Loan size (log-transformed)",
  what = c(1, 1, 1, 0)
)

## log="y" selected
```

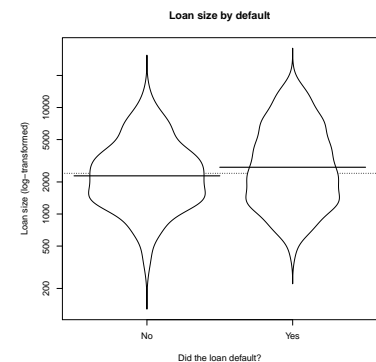


Figure 18: Distributions of loan sizes separated for loans that defaulted and those that did not.

We calculated the median loan size separately for people whose loans defaulted and those whose loans did not default. The loan amounts of loans that defaulted (median of 2574) tended to be a bit larger than those that did not (median of 2244). However, looking at the full amount distributions (see Figure 18), it is unclear if the difference is really very meaningful.



## 6: BASIC I

### 2. Was the age of the borrower related to the loan amount?

```
summary(credit$amount[credit$age <= median(credit$age)])

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      276   1364   2310   3193   3924   18420

summary(credit$amount[credit$age > median(credit$age)])

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      250   1370   2324   3354   4042   15940

cor(credit$age, credit$amount)

## [1] 0.03271642
```

To answer this, we separately calculated the median loan amount for borrowers below and above the median age (of 33). Borrowers below the median age had a median loan amount of 2310, while borrowers above the median age had a median loan amount of 2324. This suggests that age was unrelated to loan amount. Examining the scatterplot in Figure 19, we do not find strong evidence for an effect of borrower age on loan amount.

### 3. Was there a relationship between whether or not someone had a phone and whether or not their loan defaulted?

```
with(credit, table(phone, default))

##      default
## phone  no yes
##   no  409 187
##   yes 291 113

with(credit[credit$phone == "yes",], mean(default == "yes"))

## [1] 0.279703

with(credit[credit$phone == "no",], mean(default == "yes"))

## [1] 0.3137584
```

We separately the proportion of people who defaulted on their loans separately between those who own a phone and those who do not. We found that people without a phone were slightly more likely to default (31.38%) than people with a phone (27.97%) (a mosaic.plot of the data is presented in margin Figure ).

```
# Main Plot
plot(credit$age, credit$amount,
     pch = 16, col = gray(.5, alpha = .2),
     main = "Loan amount by borrower age",
     xlab = "Borrower Age", ylab = "Loan Amount (DM)"
)

par(xpd=NA)
segments(25, 20000, 35, 20000, col = "red", lwd = 2)
text(50, 20000, "5 year average")

# Create factor from age
age.cut <- cut(credit$age, breaks = seq(20, 70, 5))

# Determine mean loan by age factor
amount.cut <- tapply(credit$amount, age.cut, mean)

# Add mean lines
lines(seq(22.5, 67.5, 5), amount.cut,
      col = "red", lwd = 2, type = "b", pch = 16)
```

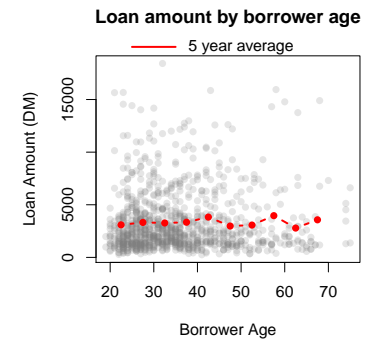


Figure 19: Distributions of loan sizes as a function of the age of the borrower.

### Additional Tips

- If you want to change the names of columns in a dataframe, you can do this by reassigning elements of the `names()` function. For example, let's change the names of the first two columns of our dataframe survey to "Sex" and "Height.cm"

```
names(survey)

## [1] "gender" "height"

names(survey)[1:2] <- c("Sex", "Height.cm")
names(survey)

## [1] "Sex"      "Height.cm"
```

```
require(RColorBrewer)
with(credit, mosaicplot(table(phone, default),
  main = "Phone Ownership and Loan Default",
  xlab = "Own Phone?",
  ylab = "Loan Default?", color = brewer.pal(12, "Set3")[5:4]
))

defper.withphone <- mean(credit$default[credit$phone == "yes"] == "yes")
defper.nophone <- mean(credit$default[credit$phone == "no"] == "yes")

text(mean(credit$phone == "no") / 2,
  defper.nophone / 2,
  paste(100 * round(defper.nophone, 2), "%", sep = ""))

text(1 - mean(credit$phone == "yes") / 2,
  defper.withphone / 2,
  paste(100 * round(defper.withphone, 2), "%", sep = ""))
```

Phone Ownership and Loan Default

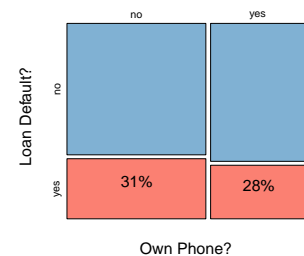


Figure 20: Mosaic plot of the contingency between phone ownership and loan defaults in the `credit` dataset. People who own a phone are slightly less likely to default on their loans than people who do not own a phone.

## 7: Plotting Basics

### Chapter Goals

1. High-level plotting commands: `plot()`, `hist()`, `boxplot`, `barplot()`
2. Main plotting parameters: `main`, `xlab`, `ylab`, `xlim`, `ylim`
3. Low-level plotting functions: `abline()`, `points()`, `text()`, `legend()`
4. Saving plots with `pdf()` and `jpg()`

Sammy Davis Jr. was one of the greatest performers of all time. If you don't know him already, Sammy was an American entertainer who lived from 1925 to 1990. The range of his talents was just incredible. He could sing, dance, act, and play multiple instruments with ease. So how is R like Sammy Davis Jr.? Like Sammy Davis Jr., R is incredibly good at doing many different things. R does data analysis like Sammy dances, and creates plot like Sammy sings. If Sammy and R did just one of these things, they'd be great. The fact that they can do both is pretty amazing.

Plotting in R works like putting paint on a canvas. You start by creating a canvas and plotting the main elements using a *high-level* plotting command. In these high-level plotting commands, you specify things like the x and y coordinates of the plot, the plot titles, and the main data in the plot. Next, you use *low-level* plotting commands to sequentially add as many additional individual elements as you'd like, from lines to arrows to text. Once you are done, you can export the plot as a jpg or pdf file. In the next section, we'll cover the most common high-level plotting functions

### High-level plotting functions

The most common high-level plotting function is `plot(x, y)`. While its name sounds like it can make any kind of plot, the `plot()` command creates a scatterplot from two vectors x and y:



Figure 21: The great Sammy Davis Jr. Do yourself a favor and spend an evening watching videos of him performing on YouTube. Image used entirely without permission.

`plot(x, y)`: Create a scatterplot from two vectors x and y.  
  `main`: Title of plot  
  `xlab`, `ylab`: axes labels  
  `xlim`, `ylim`: Limits of axes  
  `xaxt`, `yaxt`: Set to "n" to remove the axes  
  `cex`: Size of the plotting points  
  `pch`: Type of plotting points (see ?points)

## plot()

---

**x, y**

Two vectors of data on the x and y-axes

**main**

The title of the plot

**xlab, ylab**

Labels for the x and y-axes.

**xlim, ylim**

A vector of length two containing the minimum and maximum values of the x and y-axes. For example: `xlim = c(0, 100)`, `ylim = c(50, 60)` will set the x limits to [0, 100] and the y limits to [50, 60].

**col**

The color of the plotting points. For example `col = "red"` will create red plotting points.

**pch**

An integer indicating the type of plotting symbols (see `?points` and section below), or a string specifying symbols as text.

**cex**

The size of the symbols (from 0 to Inf). The default size is 1.

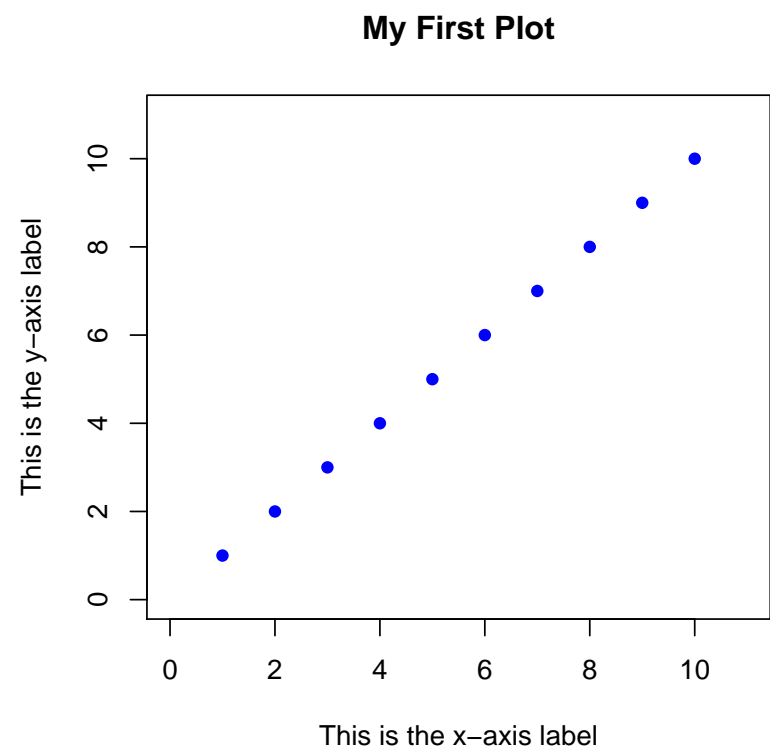
**type**

The type of plot. Use "p" for points (the default), "l" for lines, "b" for points and lines, and "\n" for no plotting

The `plot()` function, like many plotting functions, has several optional arguments that allow you to change aspects of the plot. There are so many ways to customize the look of a plot that the number of optional arguments can be overwhelming at first. Let's start by looking at an example of a simple scatterplot showing ten data points: [1, 1], [2, 2] ... [10, 10].

```
plot(x = 1:10,
     y = 1:10,
     main = "My First Plot",
     xlab = "This is the x-axis label",
     ylab = "This is the y-axis label",
```

```
xlim = c(0, 11), # Min and max values for x-axis
ylim = c(0, 11), # Min and max values for y-axis
col = "blue", # Color of the points
pch = 16, # Type of symbol (Filled circle)
cex = 1, # Size of the symbols,
type = "p" # Plot
)
```



Aside from the x and y arguments, all of the arguments are optional. If you don't specify a specific argument, then R will use a default value, or try to come up with a value that makes sense. For example, if you don't specify the `xlim` and `ylim` arguments, R will set the limits so that all the points fit inside the plot.

*Symbol types: pch*

When you create a plot with `plot(x, y)`, you can specify the type of symbol with the `pch` argument. You can specify the symbol type in one of two ways: with an integer, or with a string. If you use a string (like "p"), R will use that text as the plotting symbol. If you use an integer value, you'll get the symbol that correspond to that number. See Figure 22 in the margin.

Symbols differ in their border shape and how the filling is done.

```
par(mar = rep(0, 4))

plot(x = rep(1:5, each = 5),
     y = rep(5:1, times = 5),
     pch = 1:25,
     xlab = "", ylab = "", xaxt = "n", yaxt = "n",
     xlim = c(.5, 5.5),
     ylim = c(0, 6),
     bty = "n", bg = "gray", cex = 1.4
)

text(x = rep(1:5, each = 5) - .35,
     y = rep(5:1, times = 5),
     labels = 1:25, cex = 1.2
)
```

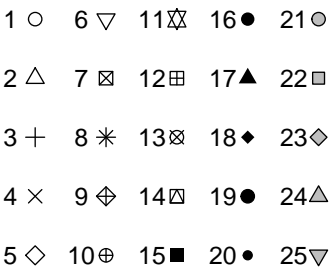


Figure 22: The plot symbols available in R

Symbols 1 through 14 only have borders and are always empty, while symbols 15 through 20 don't have a border and are always filled. Symbols 21 through 25 have both a border and a filling.

To specify the border color for all symbols, use the `col` argument. For symbols 21 through 25, you can additionally set the color of the fill using the `bg` ("background") argument.

### Other high-Level plotting commands

While `plot()` is the most widely used high-level plotting command, there are several (perhaps even hundreds) of additional ones. I'll briefly highlight a few additional ones that you may wish to use

#### Histograms `hist()`

The function `hist()` is a high-level plotting command that creates (wait for it...) a histogram. Here are the main arguments for `hist()`:

#### `hist()`

`x`

A vector of data

`breaks`

One of several values that defines how bins are created. The most common argument is a single number giving the number of bins you want in the histogram. See `?hist` for additional ways to specify this.

`col`

The color of the filling of the bars. (e.g.; `col = "red"`)

`border`

The color of the border of the bars. (e.g.; `border = "green"`)

`probability`

A logical value indicating whether to plot the results as probabilities (the default is `FALSE`)

`main`, `xlab`, `ylab`, `xlim`, `ylim` ...

Other standard plotting arguments

```
hist(x = ChickWeight$weight,
     main = "Chick Weights",
     xlab = "Weights (all time points)"
)
```

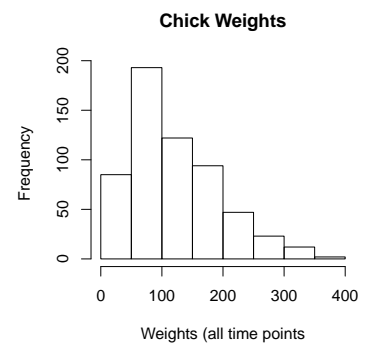


Figure 23: Histogram of weight data from `ChickWeight`

You can see an example of a histogram in the margin Figure 23.

### Boxplots `boxplot()`

Boxplots aren't used so often anymore (for reasons that I'll show you shortly), but I think it's good to know how to make them, even if it's just for historical purposes. To create a boxplot, use the `boxplot()` function:

#### `boxplot()`

##### formula, data

A formula in the form `formula = dv ~ iv` indicating the dependent variable and independent variable, and a dataframe containing the variables in the formula. For example `formula = height ~ sex`

##### subset

An optional logical vector indicating a subset of the data to plot. For example, the command `subset = gender == "male" & weight < 120`, will only plot data for males with weight less than 120.

##### border, col

The color of the borders (`border`) and filling (`col`) of the boxes.

##### names

A string vector indicating the names of the boxes. E.g.; `names = c("males", "females")`

##### horizontal

A logical value indicating whether to plot the boxes horizontally.

```
boxplot(x = ChickWeight$weight,
        names = "All Data", ylab = "Weight",
        main = "Plot 1: All Weights")
```

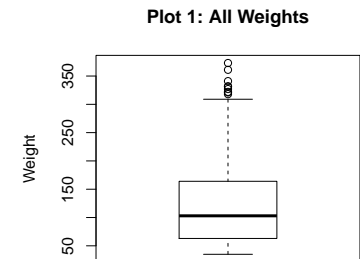


Figure 24: Plotting data from a single vector

```
boxplot(weight ~ Diet, Formula: DV is weight and IV is Diet
        data = ChickWeight, dataframe xlab = "Diet", ylab = "Weight",
        main = "Plot 2: Weight separated by Diet") @
```

Figure 25: Plotting data as a function of levels of an independent variable using the `y ~ x` formula notation.

When you use `boxplot()`, you can either specify a single vector of data to plot, or you can use a formula to indicate a dependent and independent variable. If you do this, R will add separate boxes for all values of the independent variable.

Let's go through two examples of boxplots in Figure ?? . In the first plot, I just entered a single vector of data: `ChickWeight$weight` representing all weight data in the dataframe. In the second plot, I plotted separate boxes for the different levels of `Diet` using the

formula notation `weight ~ Diet`. If you're wondering how R knows that I'm referring to the `ChickWeight` dataframe when using the formula notation, the answer is that I had to specify the name of the dataframe `ChickWeight` as an additional data argument. This argument tells R that the objects in the formula are names in the `ChickWeight` dataframe.

*Beanplots: `beanplot()`*

The last high-level plotting I want to show you is the `beanplot()` function. This function creates a beanplot, which (like boxplots and histograms), shows you a distribution of sample data. However, as you can see in Figure 26, they look much, much cooler than a boxplot or histogram. What's really great about beanplots is that they show you a combination of three elements: raw data, smoothed distribution lines, and group averages. This means that you can quickly detect outliers, multiple modes, or missing data much better than you can with boxplots.

To use the `beanplot()` function, you first need to download the `beanplot` package:

```
install.packages("beanplot")
```

Here are some of the main arguments for `beanplot()`. Check out the help menu (`?beanplot`) to see several additional arguments

```
require("beanplot")

bean.cols <- lapply(brewer.pal(4, "Set3"),
  function(x) {return(c(x, "black", "black", "black"))})

beanplot(weight ~ Diet,
  data = ChickWeight,
  main = "Beans",
  xlab = "Diet",
  ylab = "Weight",
  col = bean.cols,
  lwd = 1,
  what = c(1, 1, 1, 1), log = ""
)
```

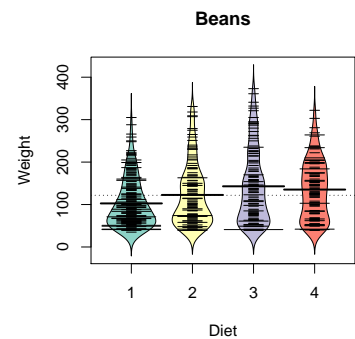


Figure 26: Beanplots from the `beanplot` package. Beanplots are great at simultaneously showing raw data, smoothed distributions, and group averages.



## beanplot()

---

### formula, data

A formula in the form `formula = dv ~ iv` indicating the dependent variable and independent variable, and a dataframe containing the variables in the formula. For example `formula = height ~ sex, data = survey`

### subset

An optional logical vector indicating a subset of the data to plot. For example, the command `subset = gender == "male" & weight < 120`, will only plot data for males with weight less than 120.

### what

A vector of four boolean (0 or 1) values indicating what to plot in the following order: the total average line, the beans, the bean average, and the beanlines. For example, to plot everything, use `what = c(1, 1, 1, 1)`. To plot just the beans, use `what = c(0, 0, 1, 0)`

### color

The colors in the plot. A vector of up to four colors can be used representing the areas of the beans, the lines inside the beans, the lines outside the beans, and the average line per bean. If you want to make each bean a different color, you have to specify a list of separate color vectors, one for each bean. Look at my code in Figure 26 for a way to do this using `lapply`.

### names

A vector of names of the beans.

### overallline

A method for determining the overall line (either "mean" or "median")

The code for creating beanplots is very similar to the code for boxplots. The following code creates the plot in Figure 26.

```
require("beanplot")
bean.cols <- lapply(brewer.pal(4, "Set3"),
  function(x) {return(c(x, "black", "black", "black"))})

beanplot(weight ~ Diet,
```

```
data = ChickWeight,
main = "Beans",
xlab = "Diet",
ylab = "Weight",
col = bean.cols ,
lwd = 1,
what = c(1, 1, 1, 1), log = ""
)
```

One major new argument is `what`, which dictates what exactly is plotted. You specify `what` using a vector of four Boolean (0 or 1) values. In the plot in Figure 26, I've set all values to 1 which means that the function will include all four plot elements. If you'd like to remove certain elements, like the individual lines or the average lines, you can remove them by replacing the respective 1s to 0s.

### *Low-level plotting functions*

Once you've created a plot with a high-level plotting function, you can add additional elements, like additional data points, reference lines, text, and legends using low-level plotting functions. There are many low-level plotting functions, I will focus on those that I frequently use.

### *Starting with a blank plot*

I like using low-level plotting functions so much that I frequently like to start with a (mostly) blank plotting space, and then add the main plot elements using low-level plotting functions. To start with a blank plot, use the `plot()` function combined with the arguments `type = "n"`, `xaxt = "n"`, `yaxt = "n"` and all labels set to `""`. See margin Figure 27 for an example

Once you've created a blank plot, you can proceed to add all the elements you'd like with low-level plotting commands. Let's start with `points()`, which adds points to an existing plot

```
# Create a blank plot
plot(x = 1, y = 1, xlab = "", ylab = "",
     xaxt = "n", yaxt = "n", type = "n",
     xlim = c(0, 100), ylim = c(0, 100), main = "Blank Plot")
```

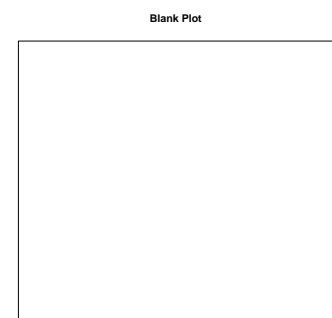


Figure 27: A blank plot. Useful to start with before adding elements with low-level plotting commands. Just make sure to set the axis limits to values that make sense for your future data.

*points()*

## points()

*x, y*

Two vectors corresponding to the x and y values of the points

*pch, col, bg*

Type of plotting symbols (*pch*), color of the plotting symbols (*col*), and the color of the filling of the plotting symbols (*bg*) for plotting symbols 21 through 25

For example, to add red circle points to a plot where *x.vals* are the x-values and *y.vals* are the y-values, you can run the code:

```
points(x = x.vals, # x-values
       y = y.vals, # y-values
       col = "red", # Symbol color
       pch = 16 # Symbol type (circles)
)
```

Because you can continue adding as many low-level plotting commands to a plot as you'd like, you can keep adding different types or colors of points by adding additional `points()` functions. However, keep in mind that because R plots each element on top of the previous one, early calls to `points()` might be covered by later calls. So add the points that you want in the foreground at the end!

In margin Figure 28, I use the `points` function to plot data from `ChickWeight`, where chicks on diet 1 are plotted in red, and chicks on diet 2 are plotted in skyblue.

Next, we'll look at `abline()` which adds straight lines to a plot:

```
# Get subsets of data
diet.1 <- subset(ChickWeight, Diet == 1)
diet.2 <- subset(ChickWeight, Diet == 2)

# Create a blank plot
plot(x = 1, y = 1, xlab = "Time", ylab = "Weight",
     type = "n", main = "Even More Chicken Weights",
     xlim = c(0, 23), ylim = c(0, 315))

# Add red points for diet 1
points(diet.1$Time, diet.1$weight, pch = 16, col = "red")

# Add skyblue points for diet 2
points(diet.2$Time, diet.2$weight, pch = 16, col = "skyblue")
```

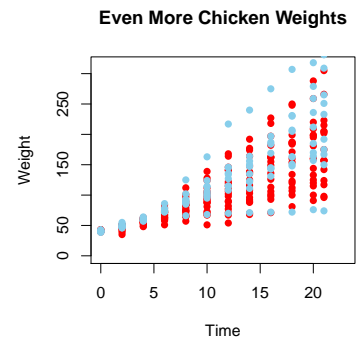


Figure 28: Adding additional points to an existing plot with `points()`

```
abline(v = 1:6, lty = 1:6, lwd = 2)

mtext(1:6,
      side = 3,
      at = 1:6, cex = 1.5, line = 1
    )

mtext("lty = ...", side = 3, at = 3.5, line = 4, cex = 2)
```

*abline()***abline()****a, b**

Numeric scalars or vectors indicating the slope (a) and intercept b of the line(s)

**h, v**

Numeric scalars or vectors indicating the y-value of horizontal lines (h) or x-values of vertical lines v. For example, `abline(h = 1)` will add a horizontal line at  $y = 1$ , while `abline(v = 10)` will add a vertical line at  $x = 1$

**lty, lwd**

Type (lty) and width (lwd) of line. See margin Figure to see line types.

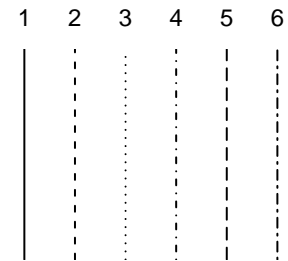
**lty = ...**

Figure 29: Line types generated from arguments to lty.

For example, to add a vertical line at an x-value of 0 or a horizontal line at a y-value at 100 you'd enter

```
abline(v = 0) # Add a vertical line at x = 0
abline(h = 100) # Add a horizontal line at y = 100
```

You can easily use `abline()` to add gridlines to plots by entering vectors in the `h` and `v` arguments. For example, to add gridlines to a plot at x-values and y-values from 0 to 10 in steps of 1, you'd enter

```
abline(v = 1:10) # Add vertical lines from 1 to 10
abline(h = 1:10) # Add horizontal lines from 1 to 10
```

In margin Figure 30 I add gridlines and a diagonal reference line to a plot before adding points.

Next, we'll move on to text, which adds text to a plot

*text()*

With `text()`, you can add text to a plot. You can use `text()` to highlight specific points of interest in the plot, or to add information (like a third variable) for every point in a plot. Here are the main arguments to `text()`

```
# Create a blank plot
plot(x = 1, y = 1, xlab = "Group", ylab = "Length",
     type = "n", main = "Gridlines with abline()",
     xlim = c(0, 10), ylim = c(0, 10))

# Add horizontal gridlines
abline(h = 1:10, lwd = 1, col = gray(.8))

# Add vertical gridlines
abline(v = 1:10, col = gray(.8))

# Add main diagonal reference line
abline(a = 0, b = 1, lwd = 2, lty = 2)

# Create data
x.data <- rnorm(100, mean = 5, sd = 2)
y.data <- x.data + rnorm(100, mean = 0, sd = 3)

# Add points
points(x = x.data,
      y = y.data, pch = 16,
      col = gray(.4, alpha = .5),
      cex = c(runif(90, 0, 2), runif(10, 3, 4))
    )
```

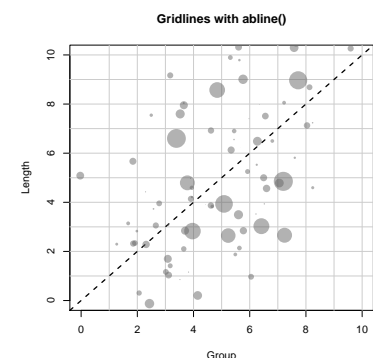


Figure 30: Adding gridlines to a plot with `abline()`.

## text()

**x, y**

Numeric scalars or vectors specifying the coordinates of the labels

**labels**

String vector of the text you're plotting. Use the `paste()` function to create multiple strings or combine strings with numeric objects.

**cex**

Numeric scalar or vector specifying the size of the labels

**adj**

A numerical value between 0 and 1 specifying the horizontal and/or vertical justification of text. Use 0 for left justification, .5 for centering, and 1 for right justification.

**pos**

Specifies the position of the text relative to the x-y coordinates. Values of 1, 2, 3 and 4 respectively indicate below, to the left, above, and to the right of the x-y coordinates.

**font**

The font face. 1 = plain, 2 = bold, 3 = italic, 4 = bold-italic.

For example, if you want to add the text "This is the center of the plot" to a plot at the coordinates (0, 0), you'd enter

```
text(x = 0, y = 0, labels = "This is the center of the plot")
```

Alternatively, let's say you have a scatterplot and wanted to add the x-values in text right above (`pos = 3`) each point, you could do this by using the code:

```
text(x = x.data, # X-values of data
     y = y.data, # X-values of data
     labels = x.data, # Add text of the x-values
     pos = 3 # Put the text right above the points
)
```

To see `text()` in action, look at margin Figure 31 where I put the x-values of some random data right above their points:

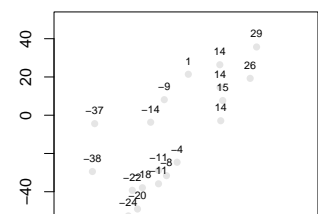
```
# Step 1: Generate Data
x.data <- rnorm(20, mean = 0, sd = 20)
y.data <- x.data + rnorm(20, mean = 0, sd = 20)

# Step 2: Create a blank plot
plot(x = 1, xlab = "", ylab = "",
     type = "n", main = "Adding text with text()",
     xlim = c(-50, 50), ylim = c(-50, 50))

# Step 3: Add points
points(x = x.data, y = y.data,
      pch = 16, col = gray(.5, alpha = .2))

# Step #4: Add x-coordinates in text above points
text(x = x.data,
     y = y.data,
     labels = round(x.data, 0),
     pos = 3, # put coordinates below the points
     cex = .7
)
```

Adding text with `text()`



When entering text in the `labels` argument, keep in mind that R will, by default, plot the entire text in one line. However, if you are adding a long text string (like a sentence), you may want to separate the text into separate lines. To do this, add the text `"\n"` where you want new lines to start. Look at Figure 32 for an example.

### Formatting text for plotting

A common way to use text in a plot, either in the main title of a plot or using the `text()` function, is to combine text with numerical data. For example, you may want to include the text "Mean = 3.14" in a plot to show that the mean of the data is 3.14. But how can we combine numerical data with text? In R, we can do this with the `paste()` function:

paste()	
...	
	One or more scalars or vectors (numeric or string) to be combined. For example <code>paste("The mean of x is ", mean(x), sep = "")</code> will create a string combining text and a statistic calculated from data.
sep	A character string that separates the arguments. Set to "" for no separation

The `paste` function will be helpful to you anytime you want to combine either multiple strings, or text and strings together. For example, let's say you want to write text in a plot that says The mean of these data are XXX, where XXX is replaced by the group mean. To do this, just include the main text and the object referring to the numerical mean as arguments to `paste()`:

```
data <- rnorm(200, mean = 20, sd = 10)
mean(data)

## [1] 20.56928

paste("The mean of the group is", mean(data)) # No rounding

## [1] "The mean of the group is 20.5692797809263"

paste("The mean of the group is", round(mean(data), 2)) # No rounding

## [1] "The mean of the group is 20.57"
```

To plot text on separate lines in a plot, put the tag `"\n"` between lines.

```
plot(1, type = "n", main = "The \\n tag",
     xlab = "", ylab = "")

# Text without \\n breaks
text(x = 1, y = 1.3, labels = "Text without \\n", font = 2)
text(x = 1, y = 1.2,
     labels = "Haikus are easy. But sometimes they don't make sense. Refrigerator")

# Text with \\n breaks
abline(h = 1, lty = 2)
text(x = 1, y = .92, labels = "Text with \\n", font = 2)
text(x = 1, y = .7,
     labels = "Haikus are easy\\nBut sometimes they don't make sense\\nRefrigerator")
```

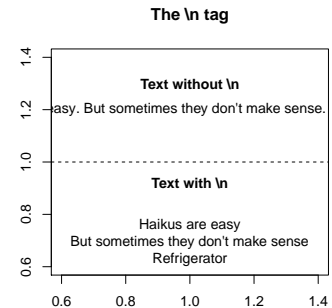


Figure 32: Using the `"\n"` tag to plot text on separate lines.

When you include descriptive statistics in a plot, you will almost always want to use the `round(x, digits)` function to reduce the number of digits in the statistic.

You can also use vectors as arguments to the `paste()` function. For example, let's say that you want to create a vector of labels for 5 groups, and you want each group to be labelled "Group X". We can easily do this with `paste()`

```
paste("Group", 1:5, sep = " ")

## [1] "Group 1" "Group 2" "Group 3" "Group 4" "Group 5"
```

### `curve()`

The `curve()` function allows you to add a line showing a specific function or equation to a plot

```
plot(1, xlim = c(-5, 5), ylim = c(-5, 5),
     type = "n", main = "Plotting function lines with curve()")
abline(h = 0)
abline(v = 0)

require("RColorBrewer")
col.vec <- brewer.pal(12, name = "Set3")[4:7]

curve(expr = x^2, from = -5, to = 5,
      add = T, lwd = 2, col = col.vec[1])
curve(expr = x^5, from = 0, to = 5,
      add = T, lwd = 2, col = col.vec[2])
curve(expr = sin, from = -5, to = 5,
      add = T, lwd = 2, col = col.vec[3])

my.fun <- function(x) {return(dnorm(x, mean = 2, sd = .2))}
curve(expr = my.fun, from = -5, to = 5,
      add = T, lwd = 2, col = col.vec[4])

legend("bottomright",
      legend = c("x^2", "x^5", "sin(x)", "dnorm(x, 2, .2)"),
      col = col.vec[1:4], lwd = 2,
      lty = 1, cex = .8, bty = "n")
```

### `curve()`

#### `expr`

The name of a function written as a function of  $x$  that returns a single vector. You can either use base functions in R like `expr = x^2`, `expr = x + 4 - 2`, or use your own custom functions such as `expr = my.fun`, where `my.fun` is previously defined (e.g.; `my.fun <- function(x) dnorm(x, mean = 10, sd = 3)`)

#### `from, to`

The starting (`from`) and ending (`to`) value of  $x$  to be plotted.

#### `add`

A logical value indicating whether or not to add the curve to an existing plot. If `add = FALSE`, then `curve()` will act like a high-level plotting function and create a new plot. If `add = TRUE`, then `curve()` will act like a low-level plotting function.

#### `lty, lwd, col`

Additional arguments such as `lty`, `col`, `lwd`, ...

Plotting function lines with `curve()`

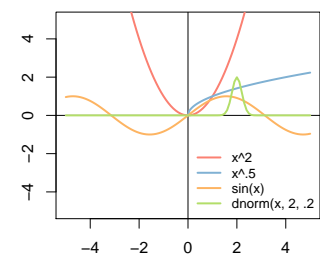


Figure 33: Using `curve()` to easily create lines of functions

For example, to add the function  $x^2$  to a plot from the  $x$ -values -10 to 10, you can run the code:

```
curve(expr = x^2, from = -10, to = 10)
```

If you want to add a custom function to a plot, you can define the function and then use that function name as the argument to `expr`. For example, to plot the normal distribution with a mean of 10 and standard deviation of 3, you can use this code:

```
my.fun <- function(x) {dnorm(x, mean = 10, sd = 3)}
curve(expr = my.fun, from = -10, to = 10)
```

In Figure 33, I use the `curve()` function to create curves of several mathematical formulas.

## `legend()`

The last low-level plotting function that we'll go over in detail is `legend()` which adds a legend to a plot. This function has the following arguments

### `legend()`

`x, y`

Coordinates of the legend - for example, `x = 0, y = 0` will put the text at the coordinates (0, 0). Alternatively, you can enter a string indicating where to put the legend (i.e.; "topright", "topleft"). For example, "bottomright" will always put the legend at the bottom right corner of the plot.

`labels`

A string vector specifying the text in the legend. For example, `legend = c("Males", "Females")` will create two groups with names Males and Females.

`pch, lty, lwd, col, pt.bg, ...`

Additional arguments specifying symbol types (`pch`), line types (`lty`), line widths (`lwd`), background color of symbol types 21 through 25 (`pt.bg`) and several other optional arguments. See `?legend` for a complete list

```
# Generate some random data
female.x <- rnorm(100)
female.y <- female.x + rnorm(100)
male.x <- rnorm(100)
male.y <- male.x + rnorm(100)

# Create plot with data from females
plot(female.x, female.y, pch = 16, col = 'blue',
     xlab = "x", ylab = "y", main = "Adding a legend with legend()")

# Add data from males
points(male.x, male.y, pch = 16, col = 'orange')

# Add legend
legend("bottomright",
     legend = c("Females", "Males"),
     col = c('blue', 'orange'),
     pch = c(16, 16),
     bg = "white")
```

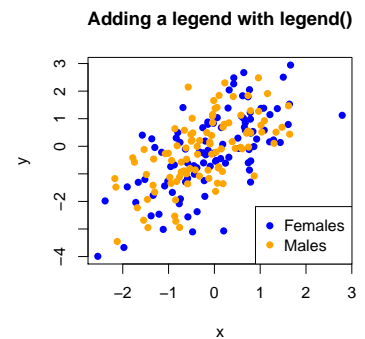


Figure 34: Creating a legend labeling the symbol types from different groups

For example, to add a legend to the bottom-right of an existing graph where data from females are plotted in blue circles and data from males are plotted in pink circles, you'd use the following code:

```
legend("bottomright", # Put legend in bottom right of graph
     legend = c("Females", "Males"), # Names of groups
     col = c("blue", "orange"), # Colors of symbols
     pch = c(16, 16) # Point types
)
```



In margin Figure I use this code to add a legend to plot containing data from males and females.

### *Additional low-level plotting functions*

There are many more low-level plotting functions that can add additional elements to your plots. Here are some I use. To see examples of how to use each one, check out their associated help menus.

#### Additional low-level plotting functions

##### `rect()`

Add rectangles to a plot at coordinates specified by `xleft`, `ybottom`, `xright`, `ybottom`. For example, to add a rectangle with corners at (0, 0) and c(10, 10), specify `xleft = 0`, `ybottom = 0`, `xright = 10`, `yttop = 10`. Additional arguments like `col`, `border` change the color of the rectangle.

##### `polygon()`

Add a polygon to a plot at coordinates specified by vectors `x` and `y`. Additional arguments such as `col`, `border` change the color of the inside and border of the polygon

##### `segments()`, `arrows()`

Add segments (lines with fixed endings), or arrows to a plot.

##### `symbols(add = T)`

Add symbols (circles, squares, rectangles, stars, thermometers) to a plot. The dimensions of each symbol are specified with specific input types. See `?symbols` for details. Specify `add = T` to add to an existing plot or `add = F` to create a new plot.

##### `axis()`

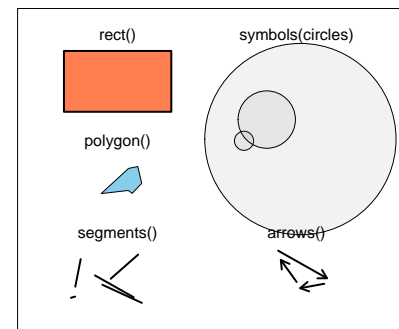
Add an additional axis to a plot (or add fully customizable `x` and `y` axes). Usually you only use this if you set `xaxt = "n"`, `yaxt = "n"` in the original high-level plotting function.

##### `mtext()`

Add text to the margins of a plot. Look at the help menu for `mtext()` to see parameters for this function.

```
par(mar = c(0, 0, 3, 0))
plot(1, xlim = c(1, 100), ylim = c(1, 100),
     type = "n", xaxt = "n", yaxt = "n",
     ylab = "", xlab = "", main = "Adding simple figures to a plot")
text(25, 95, labels = "rect()")
rect(xleft = 10, ybottom = 70,
     xright = 40, ytop = 90, lwd = 2, col = "coral")
text(25, 60, labels = "polygon()")
polygon(x = runif(6, 15, 35),
        y = runif(6, 40, 55),
        col = "skyblue")
# polygon(x = c(15, 35, 25, 15),
#         y = c(40, 40, 55, 40),
#         col = "skyblue")
text(25, 30, labels = "segments()")
segments(x0 = runif(5, 10, 40),
         y0 = runif(5, 5, 25),
         x1 = runif(5, 10, 40),
         y1 = runif(5, 5, 25), lwd = 2)
text(75, 95, labels = "symbols(circles)")
symbols(x = runif(3, 60, 90),
        y = runif(3, 60, 70),
        circles = c(1, .1, .3),
        add = T, bg = gray(.5, .1))
text(75, 30, labels = "arrows()")
arrows(x0 = runif(3, 60, 90),
       y0 = runif(3, 10, 25),
       x1 = runif(3, 60, 90),
       y1 = runif(3, 10, 25),
       length = .1, lwd = 2)
```

Adding simple figures to a plot



*Saving plots to a file*

Once you've created a plot in R, you may wish to save it to a file so you can use it in another document. To do this, you'll use either the `pdf()` or `jpeg()` functions. These functions will save your plot to either a .pdf or jpeg file.

### pdf() and jpeg()

---

**file**

The name and file destination of the final plot entered as a string. For example, to put a plot on my desktop, I'd write `file = "/Users/Nathaniel/Desktop/plot.pdf"` when creating a pdf, and `file = "/Users/Nathaniel/Desktop/plot.jpg"` when creating a jpeg.

**width, height**

The width and height of the final plot in inches.

**family()**

An optional name of the font family to use for the plot. For example, `family = "Helvetica"` will use the Helvetica font for all text (assuming you have Helvetica on your system). For more help on using different fonts, look at section "Using extra fonts in R" in Chapter XX

**dev.off()**

This is *not* an argument to `pdf()` and `jpeg()`. You just need to execute this code after creating the plot to finish creating the image file (see examples below).

To use these functions to save files, you need to follow 3 steps

1. Execute the `pdf()` or `jpeg()` functions with `file`, `width` and `height` arguments.
2. Execute all your plotting code.
3. Complete the file by executing the command `dev.off()`. This tells R that you're done creating the file.

Here's an example of the three steps.

```
# Step 1: Call the pdf command
pdf(file = "/Users/Nathaniel/Desktop/My Plot.pdf", # The directory you want to save the file in
    width = 4, # The width of the plot in inches
    height = 4 # The height of the plot in inches
)

# Step 2: Create the plot
plot(1:10, 1:10)
abline(v = 0) # Additional low-level plotting commands
text(x = 0, y = 1, labels = "Random text")

# Step 3: Run dev.off() to create the file!
dev.off()
```

You'll notice that after you close the plot with `dev.off()`, you'll see a message in the prompt like "null device".

Using the command `pdf()` will save the file as a pdf. If you use `jpeg()`, it will be saved as a jpeg.

### *A worked example: Creating a plot with automated numeric labels*

Let's use the `paste()` command to create a histogram with labels indicating the mean, median, min, and mean of the dataset. We'll do this in 5 steps

1. Generate the data and the histogram
2. Add text and reference line for the mean
3. Add text and reference line for the minimum
4. Add text and reference line for the maximum
5. Add a subtitle in full sentences with each summary statistic.

```
# Step 1: Generate data and main histogram
data <- rnorm(100, mean = 20, sd = 2)

title.text <- paste(
  "Note: There were ", length(data), " data points. The mean and median of the data were ",
  round(mean(data), 2), " and ", round(median(data), 2)), ".\n\nThe minimum and maximum values were ",
  round(min(data), 2), " and ", round(max(data), 2), " .", sep = ""

)

hist(data,
  xlim = c(10, 30),
  ylim = c(0, 40),
  main = title.text,
  cex.main = .7
)

# Step 2: Add mean text and line
text(mean(data), 38,
  labels = paste("Mean\n", round(mean(data), 2), sep = ""),
  adj = 0,
```

```

    pos = 4
  )
  abline(v = mean(data), lty = 2)

# Step 3: Add minimum text and line
text(min(data), 25,
     labels = paste("Min\n", round(min(data), 2), sep = ""),
     adj = 0,
     pos = 2
  )

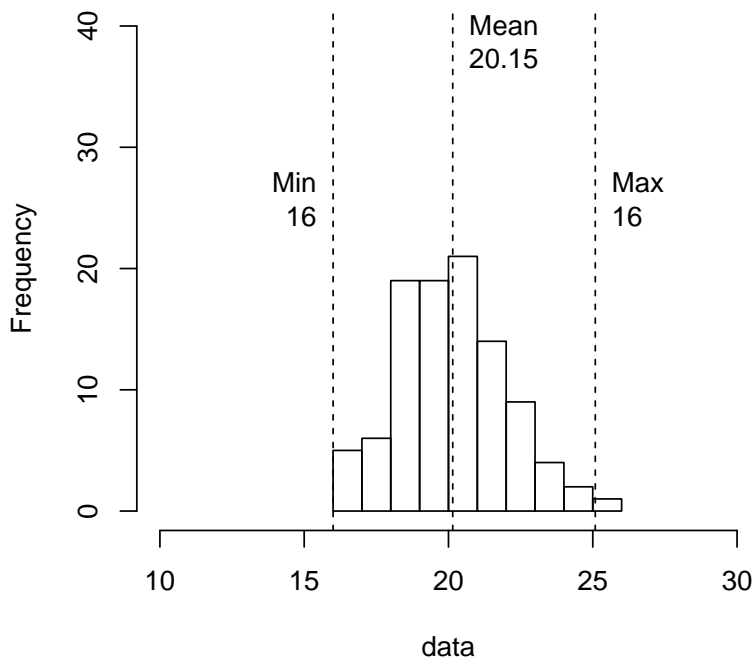
  abline(v = min(data), lty = 2)

# Step 4: Add maximum text and line
text(max(data), 25,
     labels = paste("Max\n", round(max(data), 2), sep = ""),
     adj = 0,
     pos = 4
  )

  abline(v = max(data), lty = 2)

```

**Note:** There were 100 data points. The mean and median of the data were 20.15 and 20. The minimum and maximum values were 16 and 25.09.



The benefit of using `paste()` over hard-coding the text (for example by typing `labels = "Mean = 20"`) is that the code will automatically change the value of the mean when the data changes. To see this in action, run the code above several times and see how the numbers automatically update. In later chapters, when use loops to

create multiple graphs over different sets of data, this will become extremely helpful!

### *Additional Tips*

- Many high-level plotting functions can be used like low-level plotting functions if you add an additional argument like `new = F` or `add = T`. Look at the help menu for specific high-level functions to see which arguments allow you to add a high-level plot to an existing plot.



## 8: Customizing Plots

1. Specifying and creating colors
2. Specifying plot margins with `par(mar)`
3. Putting several plots together with `par(mfrow)` and `layout`
4. Using different fonts in plots

### *Colors in R*

There are many ways to specify colors in R. If you want to specify a color directly, you can do that in one of the following ways:

#### *Specifying colors as a string*

The easiest way to specify a color is to just write its name as a string. For example, you can write "blue", "lightgreen", "red", among many other colors. To see a list of all the named colors, look at the vector `colors()` which contains all 657 named colors in R. Here is a random sample of 10 of them (to see all the colors, look at the color graph in the Appendix)

```
colors()[sample(1:length(colors()), 10)]  
  
## [1] "coral2"      "gold1"      "grey72"     "darkseagreen"  
## [5] "dodgerblue"  "gray79"     "violetred"  "orange4"  
## [9] "darkorange2" "red1"
```

#### *Shades of gray with `gray()`*

If you're a lonely, sexually repressed, 50+ year old housewife, then you might want to stick with shades of gray. If so, the function `gray(x)` is your answer. `gray()` is a function that takes a number (or vector of numbers) between 0 and 1 as an argument, and returns a shade of gray (or many shades of gray with a vector input). A value of 1 is equivalent to "white" while 0 is equivalent to "black". This function is very helpful if you want to create shades of gray

depending on the value of a numeric vector. For example, if you had survey data and plotted income on the x-axis and happiness on the y-axis of a scatterplot, you could determine the darkness of each point as a function of a third quantitative variable (such as number of children or amount of travel time to work). I plotted an example of this in Figure 35.

### RGB values: `rgb()`

Every color can be defined by its RGB ("Red", "Green", "Blue") value. This value specifies the combination of shades of Red, Green and Blue that create that color. Traditionally, each color shade is defined on a scale from 0 to 255. For example, the RGB value [255, 0, 0] is pure Red, while [0, 255, 0] is pure Green.

To create a color from RGB values, use the function `rgb()`

```
inc <- rnorm(n = 200, mean = 50, sd = 10)
hap <- inc + rnorm(n = 200, mean = 0, sd = 15)
drive <- inc + rnorm(n = 200, mean = 0, sd = 5)

plot(x = inc, y = hap, pch = 16,
     col = gray((drive - min(drive)) / max(drive - min(drive))), alpha = .4),
     cex = 1.5,
     xlab = "income", ylab = "happiness"
)
```

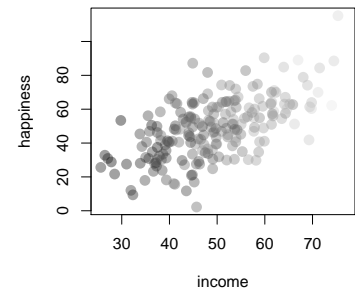


Figure 35: Using the `gray()` function to easily create shades of gray in plotting symbols based on numerical data.

### `rgb()`

red, green blue

Numeric arguments indicating the strength of red, green, and blue hues

maxColorValue

A number indicating the maximum possible hue value. The default is 1 - however, most people use maxColorValue = 255

alpha

The opacity of the color(s) inputted as a number between 0 and maxColorValue

When you use the function `rgb()`, the function will return a string as output. The string will look like nonsense to you, but that's just how R names colors:

```
rgb(red = 0, green = 255, blue = 0, maxColorValue = 255) # pure red
## [1] "#00FF00"

rgb(red = 0, green = 255, blue = 0, alpha = 100, maxColorValue = 255) # transparent red
## [1] "#00FF0064"

rgb(red = 100, green = 100, blue = 100, maxColorValue = 255) # even mixture
## [1] "#646464"
```





```
my.colors <- brewer.pal(4, "Set1") # 4 colors from Set1
my.colors

## [1] "#E41A1C" "#377EB8" "#4DAF4A" "#984EA3"
```

I know the results look like gibberish, but trust me, R will interpret them as the colors in the palette. Once you store the output of the `brewer.pal()` function as a vector (something like `my.colors`), you can then use this vector as an argument for the colors in your plot.

### Numerically defined color gradients with `colorRamp2`

My favorite way to generate colors that represent numerical data is with the function `colorRamp2` in the `circlize` package (the same package that creates that really cool `chordDiagram` from Chapter 1). The `colorRamp2` function allows you to easily generate shades of colors based on numerical data.

The best way to explain how `colorRamp2` works is by giving you an example. Let's say that you want to plot data showing the relationship between the number of drinks someone has on average per week and the resulting risk of some adverse health effect. Further, let's say you want to color the points as a function of the number of packs of cigarettes per week that person smokes, where a value of 0 packs is colored Blue, 10 packs is Orange, and 30 packs is Red. Moreover, you want the values in between these *break points* of 0, 10 and 30 to be a mix of the colors. For example, the value of 5 (half way between 0 and 10) should be an equal mix of Blue and Orange.

`colorRamp2` allows you to do exactly this. The function has three arguments:

- `breaks`: A vector indicating the break points
- `colors`: A vector of colors corresponding to each value in breaks
- `transparency`: A value between 0 and 1 indicating the transparency (1 means fully transparent)

When you run the function, the function will actually *return* another function that you can then use to generate colors. Once you store the resulting function as an object (something like `my.color.fun`). You can then apply this new function on numerical data (in our example, the number of cigarettes someone smokes) to obtain the correct color for each data point.

For example, let's create the color ramp function for our smoking data points. I'll use `colorRamp2` to create a function that I'll call

```
require("RColorBrewer")
require("circlize")

# Create Data
drinks <- sample(1:30, size = 100, replace = T)
smokes <- sample(1:30, size = 100, replace = T)
risk <- 1 / (1 + exp(-drinks / 20 + rnorm(100, mean = 0, sd = 1)))

# Create color function from colorRamp2
smoking.colors <- colorRamp2(breaks = c(0, 15, 30),
  colors = c("blue", "orange", "red"),
  transparency = .3
)

# Set up plot layout
layout(mat = matrix(c(1, 2), nrow = 2, ncol = 1),
  heights = c(2.5, 5), widths = 4)

# Top Plot
par(mar = c(4, 4, 2, 1))
plot(1, xlim = c(-.5, 31.5), ylim = c(0, .3),
  type = "n", xlab = "Cigarette Packs",
  yaxt = "n", ylab = "", bty = "n",
  main = "colorRamp2 Example")

segments(x0 = c(0, 15, 30),
  y0 = rep(0, 3),
  x1 = c(0, 15, 30),
  y1 = rep(.1, 3),
  lty = 2)

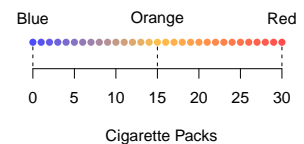
points(x = 0:30,
  y = rep(.1, 31), pch = 16,
  col = smoking.colors(0:30))

text(x = c(0, 15, 30), y = rep(.2, 3),
  labels = c("Blue", "Orange", "Red"))

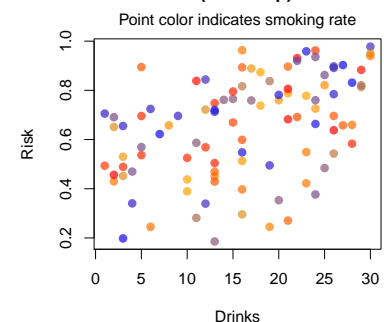
# Bottom Plot
par(mar = c(4, 4, 5, 1))
plot(x = drinks, y = risk, col = smoking.colors(smokes),
  pch = 16, cex = 1.2, main = "Plot of (Made-up) Data",
  xlab = "Drinks", ylab = "Risk")

mtext(text = "Point color indicates smoking rate", line = .5, side = 3)
```

colorRamp2 Example



Plot of (Made-up) Data



`smoking.colors` which takes a number as an argument, and returns the corresponding color:

```
smoking.colors <- colorRamp2(breaks = c(0, 15, 30),
                             colors = c("blue", "orange", "red"),
                             transparency = .3
                             )

smoking.colors(0) # Equivalent to blue

## [1] "#0000FFB3"

smoking.colors(20) # Mix of orange and red

## [1] "#FF6E00B3"
```

To see this function in action, check out the the margin Figure for an example, and check out the help menu `?colorRamp2` for more information and examples.

### *Stealing any color from your screen with a kuler*

One of my favorite tricks for getting great colors in R is to use a *color kuler*. A color kuler is a tool that allows you to determine the exact RGB values for a color on a screen. For example, let's say that you wanted to use the exact colors used in the Google logo. To do this, you need to use an app that allows you to pick colors off your computer screen. On a Mac, you can use the program called "Digital Color Meter." If you then move your mouse over the color you want, the software will tell you the exact RGB values of that color. In the image below, you can see me figuring out that the RGB value of the G in Google is R: 19, G: 72, B: 206. Using this method, I figured out the four colors of Google! Check out the margin Figure for the grand result.

```
google.colors <- c(
  rgb(19, 72, 206, maxColorValue = 255),
  rgb(206, 45, 35, maxColorValue = 255),
  rgb(253, 172, 10, maxColorValue = 255),
  rgb(18, 140, 70, maxColorValue = 255))

par(mar = rep(0, 4))

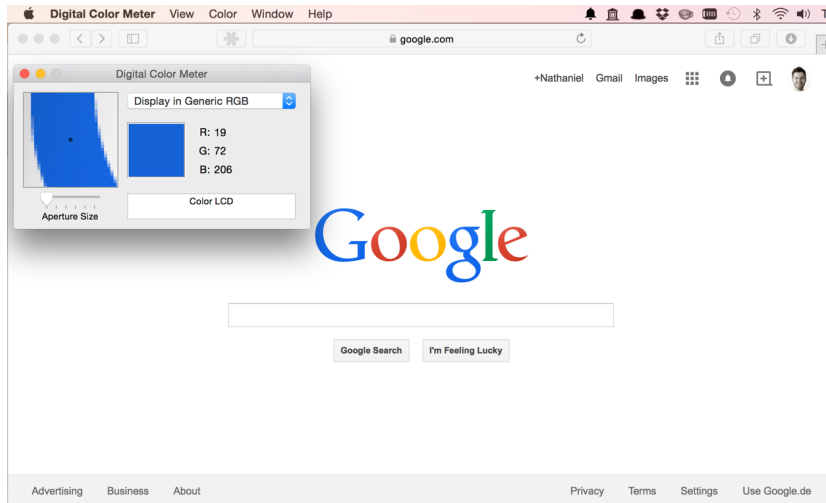
plot(1, xlim = c(0, 7), ylim = c(0, 1),
     xlab = "", ylab = "", xaxt = "n", yaxt = "n",
     type = "n", bty = "n"
     )

points(1:6, rep(.5, 6),
       pch = c(15, 16, 16, 17, 18, 15),
       col = google.colors[c(1, 2, 3, 1, 4, 2)],
       cex = 2.5)

text(3.5, .7, "Look familiar?", cex = 1.5)
```

Look familiar?





## Plot margins

All plots in R have margins surrounding them that separate the main plotting space from the area where the axes, labels and additional text lie.. To visualize how R creates plot margins, look at margin Figure .

You can adjust the size of the margins by specifying a margin parameter using the syntax `par(mar = c(a, b, c, d))` before you execute your first high-level plotting function, where a, b, c and d are the size of the margins on the bottom, left, top, and right of the plot. Let's see how this works by creating two plots with different margins:

In the plot on the left, I'll set the margins to 3 on all sides. In the plot on the right, I'll set the margins to 6 on all sides.

```
par(mar = rep(8, 4))

x.vals <- rnorm(500)
y.vals <- x.vals + rnorm(500, sd = .5)

plot(x.vals, y.vals, xlim = c(-2, 2), ylim = c(-2, 2),
     main = "", xlab = "", ylab = "", xaxt = "n",
     yaxt = "n", bty = "n", pch = 16, col = gray(.5, alpha = .2))

axis(1, at = seq(-2, 2, .5), col.axis = gray(.8), col = gray(.8))
axis(2, at = seq(-2, 2, .5), col.axis = gray(.8), col = gray(.8))

par(new = T)
par(mar = rep(0, 4))
plot(1, xlim = c(0, 1), ylim = c(0, 1), type = "n",
     main = "", bty = "n", xlab = "", ylab = "", xaxt = "n", yaxt = "n")

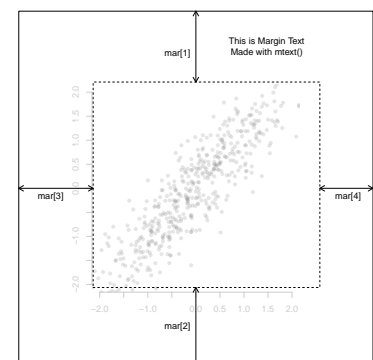
rect(0, 0, 1, 1)

rect(.21, .22, .85, .8, lty = 2)

arrows(c(.5, .5, 0, .85),
       c(.8, .22, .5, .5),
       c(.5, .5, .21, 1),
       c(1, 0, .5, .5),
       code = 3, length = .1
       )

text(c(.5, .5, .09, .93),
     c(.88, .11, .5, .5),
     labels = c("mar[1]", "mar[2]", "mar[3]", "mar[4]"),
     pos = c(2, 2, 1, 1)
     )

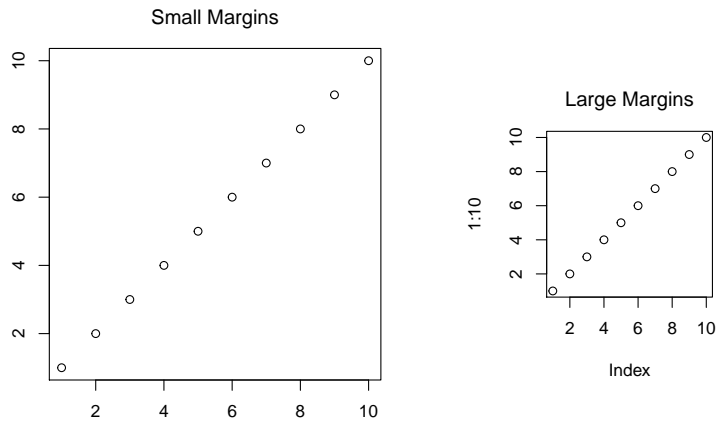
text(.7, .9, "This is Margin Text\nMade with mtext()")
```



```
par(mfrow = c(1, 2)) # Put plots next to each other

# First Plot
par(mar = rep(2, 4)) # Set the margin on all sides to 2
plot(1:10)
mtext("Small Margins", side = 3, line = 1, cex = 1.2)

# Second Plot
par(mar = rep(6, 4)) # Set the margin on all sides to 6
plot(1:10)
mtext("Large Margins", side = 3, line = 1, cex = 1.2)
```



You'll notice that the margins are so small in the first plot that you can't even see the axis labels, while in the second plot there is plenty (probably too much) white space around the plotting region.

In addition to using the `mar` parameter, you can also specify margin sizes with the `mai` parameter. This acts just like `mar` except that the values for `mai` set the margin size in inches.

The default value for `mar` is `c(5.1, 4.1, 4.1, 2.1)`

## Arranging multiple plots with `par(mfrow)` and `layout`

R makes it easy to arrange multiple plots in the same plotting space. The most common ways to do this is with the `par(mfrow)` parameter, and the `layout()` function. Let's go over each in turn:

### Simple plot layouts with `par(mfrow)` and `par(mfcol)`

The `mfrow` and `mfcol` parameters allow you to create a matrix of plots in one plotting space. Both parameters take a vector of length two as an argument, corresponding to the number of rows and columns in the resulting plotting matrix. For example, the following code sets up a 3 x 3 plotting matrix.

```
par(mfrow = c(3, 3)) # Create a 3 x 3 plotting matrix
```

When you execute this code, you won't see anything happen. However, when you execute your first high-level plotting command, you'll see that the plot will show up in the space reserved for the first plot (the top left). When you execute a second high-level plotting command, R will place that plot in the second place in the plotting matrix - either the top middle (if using `par(mfrow)`) or the left middle (if using `par(mfcol)`). As you continue to add high-level plots, R will continue to fill the plotting matrix.

So what's the difference between `par(mfrow)` and `par(mfcol)`? The only difference is that while `par(mfrow)` puts sequential plots into the plotting matrix by row, `par(mfcol)` will fill them by column.

When you are finished using a plotting matrix, be sure to reset the plotting parameter back to its default state:

```
par(mfrow = c(1, 1))
```

If you don't reset the `mfrow` parameter, R will continue creating new plotting matrices.

### Complex plot layouts with `layout()`

While `par(mfrow)` allows you to create matrices of plots, it does not allow you to create plots of different sizes. In order to arrange plots in different sized plotting spaces, you need to use the `layout()` function. Unlike `par(mfrow)`, `layout` is not a plotting parameter, rather it is a function all on its own. Let's go through the main arguments of `layout()`:

```
layout(mat, widths, heights)
```

```
par(mfrow = c(3, 3))
par(mar = rep(2.5, 4))

for(i in 1:9) { # Loop across plots

  # Generate data
  x <- rnorm(100)
  y <- x + rnorm(100)

  # Plot data
  plot(x, y, xlim = c(-2, 2), ylim = c(-2, 2),
       col.main = "gray",
       pch = 16, col = gray(.0, alpha = .1),
       xaxt = "n", yaxt = "n"
  )

  # Add a regression line for fun
  abline(lm(y ~ x), col = "gray", lty = 2)

  # Add gray axes
  axis(1, col.axis = "gray",
       col.lab = gray(.1), col = "gray")
  axis(2, col.axis = "gray",
       col.lab = gray(.1), col = "gray")

  # Add large index text
  text(0, 0, i, cex = 7)

  # Create box around border
  box(which = "figure", lty = 2)

}
```

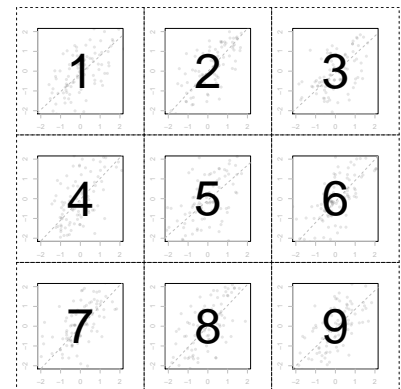


Figure 37: A matrix of plotting regions created by `par(mfrow = c(3, 3))`

- **mat:** A matrix indicating the location of the next N figures in the global plotting space. Each value in the matrix must be 0 or a positive integer. R will plot the first plot in the entries of the matrix with 1, the second plot in the entries with 2,...
- **widths:** A vector of values for the widths of the columns of the plotting space.
- **heights:** A vector of values for the heights of the rows of the plotting space.

The `layout()` function can be a bit confusing at first, so I think it's best to start with an example. Let's say you want to place histograms next to a scatterplot: Let's do this using `layout`

We'll begin by creating the *layout matrix*, this matrix will tell R in which order to create the plots:

```
layout.matrix <- matrix(c(0, 2, 3, 1), nrow = 2, ncol = 2)
layout.matrix
##      [,1] [,2]
## [1,]    0    3
## [2,]    2    1
```

Looking at the values of `layout.matrix`, you can see that we've told R to put the first plot in the bottom right, the second plot on the bottom left, and the third plot in the top right. Because we put a 0 in the first element, R knows that we don't plan to put anything in the top left area.

Now, because our layout matrix has two rows and two columns, we need to set the widths and heights of the two columns. We do this using a numeric vector of length 2. I'll set the heights of the two rows to 1 and 2 respectively, and the widths of the columns to 1 and 2 respectively. Now, when I run the code `layout.show(3)`, R will show us the plotting region we set up (see margin Figure 38)

Now we're ready to put the plots together

```
layout.matrix <- matrix(c(2, 1, 0, 3), nrow = 2, ncol = 2)

layout(mat = layout.matrix,
       heights = c(1, 2), # Heights of the two rows
       widths = c(2, 1) # Widths of the two columns
)

x.vals <- rnorm(100, mean = 100, sd = 10)
y.vals <- x.vals + rnorm(100, mean = 0, sd = 10)

# Plot 1: Scatterplot
par(mar = c(5, 4, 0, 0))
plot(x.vals, y.vals)
```

```
layout.matrix <- matrix(c(2, 1, 0, 3), nrow = 2, ncol = 2)
layout(mat = layout.matrix,
       heights = c(1, 2), # Heights of the two rows
       widths = c(2, 2) # Widths of the two columns
)
layout.show(3)
```

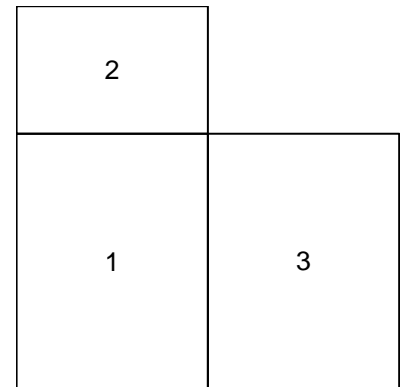


Figure 38: A plotting layout created by setting a layout matrix and specific heights and widths.

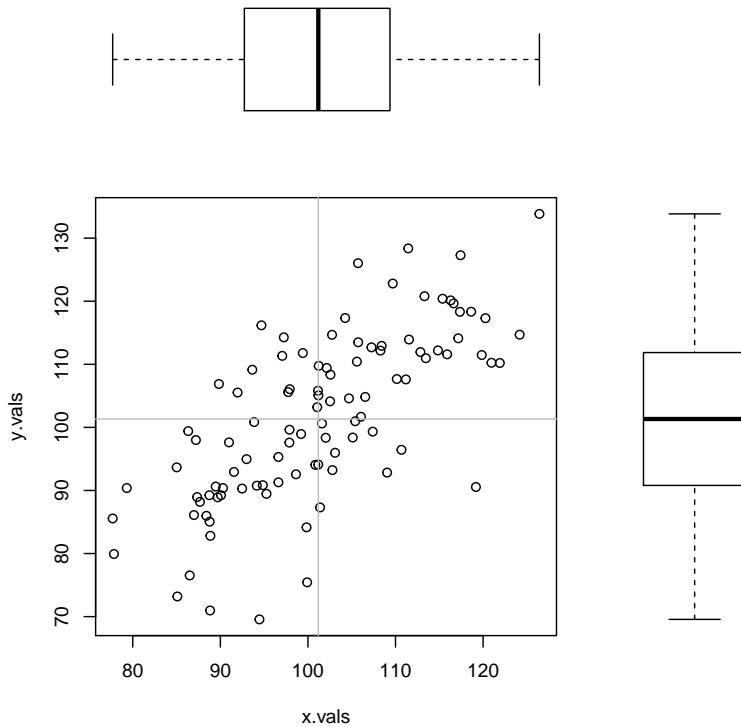
```

abline(h = median(y.vals), lty = 1, col = "gray")
abline(v = median(x.vals), lty = 1, col = "gray")

# Plot 2: X boxplot
par(mar = c(0, 4, 0, 0))
boxplot(x.vals, xaxt = "n",
        yaxt = "n", bty = "n", yaxt = "n",
        col = "white", frame = F, horizontal = T)

# Plot 3: Y boxplot
par(mar = c(5, 0, 0, 0))
boxplot(y.vals, xaxt = "n",
        yaxt = "n", bty = "n", yaxt = "n",
        col = "white", frame = F)

```



### *Using alternative fonts in pdfs with the extrafont package*

If you don't like the default font that R uses in creating plots, you can use the `extrafont` package to use additional fonts in plots saved as .pdf files. However, let me warn you that it's not as easy as just selecting a font from a drop-down menu (like in Word). To use other fonts, follow these four steps:

First, install and load the `extrafont` package



```
install.packages("extrafont")
library("extrafont")
```

### ## Registering fonts with R

Second, import the fonts on your computer into R by running the `font_import()` function. When you execute this, you'll receive a warning telling you that importing the fonts may take a few minutes. Type "y" and watch R do its magic. Don't worry if you see some warnings or if it takes a few minutes, once you've run `font_import()` once you *won't* need to run it again on your machine.

```
font_import()
```

Third, load your fonts into your current R session using the function `loadfonts()`. Unfortunately, you DO need to run this in each R session.

```
loadfonts()
```

Now you're ready to go. To see which fonts are available to use in your plots, use the `fonts()` function. When you execute this function, you'll see a table with all the fonts you can use. Let's do this on my system (I'll just print the first 50 values here)

```
fonts()[1:50] # Show me the first 50 fonts on my system

## [1] ".Keyboard"          "Andale Mono"
## [3] "Apple Braille"       "AppleMyungJo"
## [5] "Arial Black"         "Arial"
## [7] "Arial Narrow"        "Arial Rounded MT Bold"
## [9] "Arial Unicode MS"    "Brush Script MT"
## [11] "Comic Sans MS"       "Courier New"
## [13] "DIN Alternate"       "DIN Condensed"
## [15] "Georgia"            "Impact"
## [17] "Khmer Sangam MN"     "Lao Sangam MN"
## [19] "Microsoft Sans Serif" "Myanmar Sangam MN"
## [21] "Tahoma"             "Times New Roman"
## [23] "Trabuchet MS"        "Verdana"
## [25] "Webdings"            "Wingdings"
## [27] "Wingdings 2"         "Wingdings 3"
## [29] "Helvetica LT Std steevo harvie" "Helvetica World"
## [31] "HelvFB"              "HelvFE"
## [33] "Helvetica-Black-SemiBold" "Helvetica"
## [35] "Helvetica-Condensed-Black-Se" "Helvetica-Condensed-Light-Li"
## [37] "Helvetica-Condensed-Light-Light" "Helvetica-Condensed-Thin"
## [39] "Helvetica-Conth"      "Helvetica-Light-Light-Italic"
## [41] "Helvetica-Normal"     "HelveticaExt0-No"
## [43] "HelveticaExt0 2"      "HelveticaExt0 3"
## [45] "HelveticaExt0 4"      "HelveticaInserat-Roman-SemiB"
## [47] "HelveticaInserat-Roman-SemiBold" "HelveticaNeueLT Com 55 Roman"
## [49] "HelveticaNeueLT Com 57 Cn" "HelveticaNeueLT Com 53 Ex"
```

You will likely have more or less fonts than I have on my system - if you want more fonts, you'll need to download them. Now that we have a list of fonts we can use, we can finally create a plot using the new font. To do this, we need to add two special arguments when creating our plot:

1. In the `pdf()` function, add the argument `family = "fontname"`, where `fontname` is the name of the font you want to use.

2. After executing `dev.off()` to finish the plot, execute the command `embed_fonts("filename.pdf")`. This command will embed the font in the pdf file.

Let's follow these steps to create a plot using the Helvetica Light font. Again, I found this font on my computer by running `fonts()`. If you don't have this font on your computer, then it won't work for you. Instead, replace the argument "HelvLight" with a different font on your system:

```
pdf("/Users/Nathaniel/Dropbox/Git/YaRrr_Book/media/helveticalight.pdf",
    width = 4, height = 4,
    family = "HelvLight" # Specify the font in the plot
)

hist(x = rnorm(100), # some random data
     col = "skyblue",
     main = "Helvetica Light font")

dev.off()

## pdf
## 2

embed_fonts("/Users/Nathaniel/Dropbox/Git/YaRrr_Book/media/helveticalight.pdf") # Embed the fonts in the p
```

Look at Figure X to see the plot that this code created

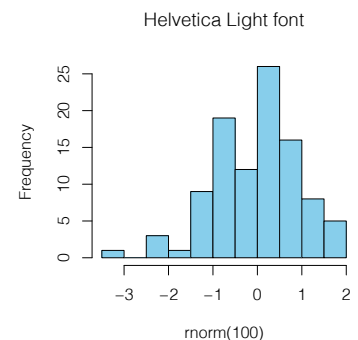
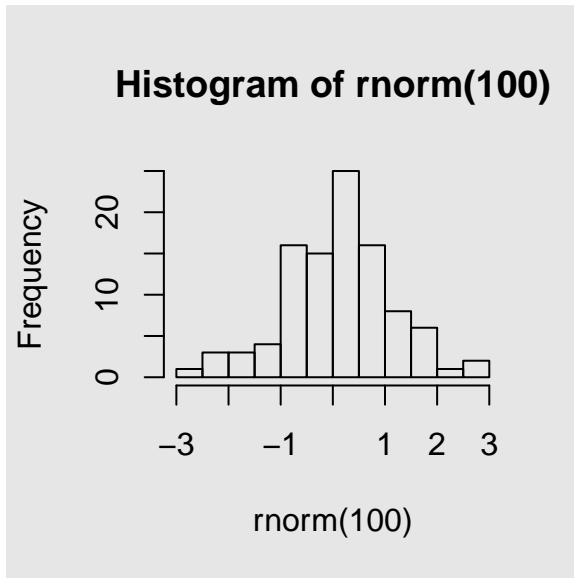


Figure 39: Plot created with Helvetica Light font (see the main text for plotting code).

### Additional Tips

- To change the background color of a plot, add the command `par(bg = mycolor)` (where `my.color` is the color you want to use) prior to creating the plot. For example, the following code will put a light gray background behind a histogram:

```
par(bg = gray(.9))
hist(x = rnorm(100))
```



See Figure 40 for a nicer example.

- Sometimes you'll mess so much with plotting parameters that you may want to set things back to their default value. To see the default values for all the plotting parameters, execute the code `par()` to print the default parameter values for all plotting parameters to the console.

```
pdf("/Users/Nathaniel/Dropbox/Git/YaRrr_Book/media/parrothelvetica.pdf",
    width = 8, height = 6, family = "HelvLight")

parrot.data <- data.frame(
  "parrots" = 0:6,
  "female" = c(200, 150, 100, 175, 55, 25, 10),
  "male" = c(150, 125, 180, 242, 10, 62, 5)
)

n.data <- nrow(parrot.data)

par(bg = rgb(61, 55, 72, maxColorValue = 255),
    mar = c(8, 6, 6, 3)
)

plot(1, xlab = "", ylab = "", xaxt = "n",
     yaxt = "n", main = "", bty = "n", type = "n",
     ylim = c(0, 250), xlim = c(.5, n.data + .5)
)

abline(h = seq(0, 250, 50), lty = 3, col = gray(.95), lwd = 1)

mtext(text = seq(50, 250, 50),
      side = 2, at = seq(50, 250, 50),
      las = 1, line = 1, col = gray(.95))

mtext(text = paste(0:(n.data - 1), " Parrots"),
      side = 1, at = 1:n.data, las = 1,
      line = 1, col = gray(.95))

female.col <- gray(1, alpha = .7)
male.col <- rgb(226, 89, 92, maxColorValue = 255, alpha = 220)

rect.width <- .35
rect.space <- .04

rect(1:n.data - rect.width - rect.space / 2,
     rep(0, n.data),
     1:n.data - rect.space / 2,
     parrot.data$female,
     col = female.col, border = NA
)

rect(1:n.data + rect.space / 2,
     rep(0, n.data),
     1:n.data + rect.width + rect.space / 2,
     parrot.data$male,
     col = male.col, border = NA
)

legend(n.data - 1, 250, c("Male Pirates", "Female Pirates"),
       col = c(female.col, male.col), pch = rep(15, 2),
       bty = "n", pt.cex = 1.5, text.col = "white"
)

mtext("Number of parrots owned by pirates", side = 3,
      at = n.data + .5, adj = 1, cex = 1.2, col = "white")

mtext("Source: Drunken survey on 22 May 2015", side = 1,
      at = 0, adj = 0, line = 3, font = 3, col = "white")

dev.off()

## pdf
## 2

embed_fonts("/Users/Nathaniel/Dropbox/Git/YaRrr_Book/media/parrothelvetica.pdf")
```

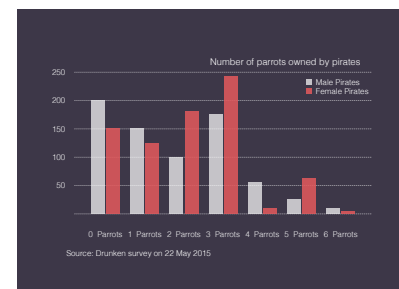


Figure 40: Use `par(bg = my.color)` before creating a plot to add a colored background. The

design of this plot was inspired by

<http://www.vox.com/2015/5/20/8625785/expensive-wine>



## 9: Advanced dataframe manipulation

### Chapter Goals

1. Grouped aggregation with `aggregate()` and `dplyr`
2. Merging datasets with `merge()`
3. Common data management tasks: rescaling, recoding

In this chapter, we'll cover how to do advanced dataframe manipulation. You'll learn how to recode values of a dataframe and quickly and easily calculate summary data from a dataframe,

### *Recoding values in a vector*

Let's say that you conducted an online survey where you asked 100 people some basic demographic information and then asked them how happy they are with their life on a scale from 1 to 8. When you look at the dataset, you realize that many people gave invalid answers to the questions. For example, some people said they had a negative income and had a negative number of siblings. How can you fix these data? In this chapter, we'll cover the basics of data-recoding

There are three basic steps to recoding values in a vector:

#### *3 Steps to recoding values*

1. Copy the original column (`original`) to a new column with a new name (`new`) using assignment.
2. Create a logical vector indicating which values of `new` should be converted.
3. Index `new` by the logical vector and assign the updated value to it.

Let's start with a dummy example where we take the integers from 1 to 10, then convert all values greater than 5 to 999

#### **3 Steps to recoding values**

1. Copy the original column (`original`) to a new column with a new name (`new`) using assignment.
2. Create a logical vector indicating which values of `new` should be converted.
3. Index `new` by the logical vector and assign the updated value to it.

```
original <- 1:10

new.vec <- original # Step 1
log.vec <- new.vec > 5 # Step 2
new.vec[log.vec] <- 999 # Step 3

new.vec

## [1] 1 2 3 4 5 999 999 999 999 999
```

As you can see, our new vector `new.vec` is identical to the integers 1 to 10 up until the value of 6, where all values have been converted to 999.

### *Recoding values from the pirate survey*

Let's start by using the dataset `pirate_survey_werrors.txt` which is stored at [http://nathanieldphillips.com/wp-content/uploads/2015/05/pirate\\_survey\\_werrors.txt](http://nathanieldphillips.com/wp-content/uploads/2015/05/pirate_survey_werrors.txt). This is a tab-delimited file containing results from 1,000 pirates.

To load the data, execute the following commands:

```
pirates.errors <- read.table("http://nathanieldphillips.com/wp-content/uploads/2015/05/pirate_survey_werro
```

Unfortunately, some pirates decided to supply some invalid answers. We'll start with the `sex` variable. To see the values that are currently in this column, we'll run `table()`:

```
table(pirates.errors$sex)

##
## depends on who is offering          female
##              3              477
##              male          other/NA
##              474              44
##      sure I'll have some
##              2
```

Very funny guys. There should only be three valid responses to this question: female, male, and other/NA. Let's recode all the other values to NA. We'll do this in three steps:

1. Copy the values of `sex` into a new vector called `sex.r` standing for "sex recoded"
2. Create a logical vector called `log.vec` that indicates which values of `sex.r` are valid - that is, in the set of values `c("female", "male", "other/NA")`

3. Assign a value of NA to all values of `sex.r` for which `log.vec` is FALSE.

```
# Step 1: Copy sex to a new column called sex.r
pirates.errors$sex.r <- pirates.errors$sex

# Step 2: Create a logical vector indicating which values are valid
log.vec <- pirates.errors$sex.r %in% c("female", "male", "other/NA")

# Step 3: Recode all invalid values to NA
pirates.errors$sex.r[log.vec == FALSE] <- NA # Step 3
```

Let's make sure this worked by looking at the values of `sex.r`. Hopefully, we will only see valid entries now:

```
table(pirates.errors$sex.r, useNA = "ifany")

##
##  female      male other/NA      <NA>
##    477       474       44        5
```

As you can see, our new column `sex.r` only contains valid entries, so it looks like our recoding did what we wanted.

Let's do the same with `age` and `tattoos`. The column `age` should be an integer between 18 and 99, and the values of `tattoos` should be an integer between 0 and a maximum of (let's say) 100. Let's start by looking at the current values of each column using the `table()` function:

```
table(pirates.errors$age)

##
## -99    0   12   13   14   15   16   17   18   19   20   21
##   2    5    2    3    2    7   10   11   17   24   30   48
##  22   23   24   25   26   27   28   29   30   31   32   33
##  45   43   65   57   62   81   71   63   46   62   58   37
##  34   35   36   37   38   39   40   41   42   44   49   500
##  37   25   20   22   19    6    3    1    1    1    1    6
## 999 12345
##   4    3

table(pirates.errors$tattoos)

##
## -10    0    1    2    3    4    5    6    7    8    9   10
##   2   11    5   13   20   25   51   68   91  104  112  112
##  11   12   13   14   15   16   17   18   19   20 1e+06
## 119   80   58   45   37   22    8    4    4    1    5
```

Note that the `table()` function does not (by default), show NA values. To see how many NA values are in the dataset, you can include the argument `useNA = "ifany"` as in `table(pirates.errors$sex.r, useNA = "ifany")`

We can see some problems here: `age` has some negative values and a few very large values. Further, `tattoos` has a few negative values and a few values greater than 100. Let's convert all these troublesome values of `age` to NA.

```
# Step 1: Copy age to a new column age.r
pirates.errors$age.r <- pirates.errors$age

# Step 2: See which values of age.r are valid
age.valid <- pirates.errors$age.r >= 18 & pirates.errors$age.r <= 99 # Step 2

# Step 3: Recode non-valid values to NA
pirates.errors$age.r[age.valid == FALSE] <- NA # Step 3
```

Now let's do the same for tattoos:

```
# Step 1: Copy tattoos to a new column tattoos.r
pirates.errors$tattoos.r <- pirates.errors$tattoos

# Step 2: See which values of tattoos.r are valid
tattoos.valid <- pirates.errors$tattoos.r %in% 0:100 # Step 2

# Step 3: Recode non-valid values to NA
pirates.errors$tattoos.r[tattoos.valid == FALSE] <- NA # Step 3
```

Now, let's look at the frequency tables of these recoded values:

```
table(pirates.errors$age.r, useNA = "always")

##
##  18  19  20  21  22  23  24  25  26  27  28  29  30  31  32
##  17  24  30  48  45  43  65  57  62  81  71  63  46  62  58
##  33  34  35  36  37  38  39  40  41  42  44  49 <NA>
##  37  37  25  20  22  19   6   3   1   1   1   1  55

table(pirates.errors$tattoos.r, useNA = "always")

##
##   0   1   2   3   4   5   6   7   8   9  10  11  12  13  14
##  11   5  13  20  25  51  68  91 104 112 112 119  80  58  45
##  15  16  17  18  19  20 <NA>
##  37  22   8   4   4   1  10
```

In these examples, we converted invalid values to NA. Of course, there's no reason why we couldn't recode the values to other values. For example, we might recode values of tattoos greater than 50 to the maximum possible value of 50. To do this, we can just change the assignment in Step 3 to a value of 50 instead of NA.

### *Splitting numerical data into groups using cut()*

When we create some plots and analyses, we may want to group numerical data into bins of similar values. For example, in our pirate survey, we might want to group pirates into age decades, where all pirates in their 20s are in one group, all those in their 30s go into another group, etc. Once we have these bins, we can calculate aggregate statistics for each group.



R has a handy function for grouping numerical data called `cut()`

### `cut()`

**x**

A vector of numeric data

**breaks**

Either a numeric vector of two or more unique cut points or a single number (greater than or equal to 2) giving the number of intervals into which `x` is to be cut. For example, `breaks = 1:10` will put break points at all integers from 1 to 10, while `breaks = 5` will split the data into 5 equal sections.

**labels**

An optional string vector of labels for each grouping. By default, labels are constructed using "(a,b]" interval notation. If `labels = FALSE`, simple integer codes are returned instead of a factor.

**right**

A logical value indicating if the intervals should be closed on the right (and open on the left) or vice versa.

Let's try a simple example by converting the integers from 1 to 50 into bins of size 10:

```
cut(1:50, seq(0, 50, 10))

## [1] (0,10] (0,10] (0,10] (0,10] (0,10] (0,10] (0,10] (0,10] (0,10]
## [9] (0,10] (0,10] (10,20] (10,20] (10,20] (10,20] (10,20] (10,20] (10,20]
## [17] (10,20] (10,20] (10,20] (10,20] (10,20] (20,30] (20,30] (20,30] (20,30]
## [25] (20,30] (20,30] (20,30] (20,30] (20,30] (20,30] (30,40] (30,40]
## [33] (30,40] (30,40] (30,40] (30,40] (30,40] (30,40] (30,40] (30,40] (30,40]
## [41] (40,50] (40,50] (40,50] (40,50] (40,50] (40,50] (40,50] (40,50] (40,50]
## [49] (40,50] (40,50]
## Levels: (0,10] (10,20] (20,30] (30,40] (40,50]
```

As you can see, our result is a vector of factors, where the first ten elements are (0, 10], the next ten elements are (10, 20], and so on. In other words, the new vector treats all numbers from 1 to 10 as being the same, and all numbers from 11 to 20 as being the same.

Let's test the `cut()` function on the age data from pirates. We'll add a new column to the dataset called `age.binned`, which separates the age data into bins of size 10. This means that every pirate between the ages of 10 and 20 will be in the first bin, those between

the ages of 21 and 30 will be in the second bin, and so on. To do this, we'll enter `pirates$age` as the `x` argument, and `seq(10, 100, 10)` as the `breaks` argument:

```
pirates$age.cut <- cut(x = pirates$age, # The raw data
                      breaks = seq(10, 100, 10) # The break points of the cuts
                      )
```

To show you how this worked, let's look at the first few rows of the columns `age` and `age.cut`

```
head(pirates[c("age", "age.cut")])

##   age age.cut
## 1  23 (20,30]
## 2  28 (20,30]
## 3  35 (30,40]
## 4  34 (30,40]
## 5  25 (20,30]
## 6  36 (30,40]
```

As you can see, `age.cut` has correctly converted the original age variable to a factor.

From these data, we can now easily calculate how many pirates are in each age group using `table()`

```
table(pirates$age.cut)

##
##  (10,20] (20,30] (30,40] (40,50] (50,60] (60,70] (70,80] (80,90]
##      111     569     294         5         0         0         0         0
## (90,100]
##         0
```

## Grouped aggregation

As you can see, we have quite a bit of data in our `Flights` dataset. Many of the questions we might want to answer with this dataset have to do with comparisons between groups. For example, "What is the flight cancellation rate for each carrier?" or "What is the arrival delay for each destination?" In each of these questions, we want to know a descriptive statistic of a numeric variable (cancellation rate and arrival delay) as a function of a nominal independent variable (carrier and destination). By now, your R skills are good enough that you *could* answer these questions already. You could do this by using `subset()` or logical indexing to slice and dice the data set for each

level of the independent variable. However, it would be a pain to have to manually create new subsets or indexes for each level of the independent variable. Thankfully, R contains many functions that will help you do this in a snap.

The first function we'll cover is `aggregate()`. The function `aggregate()` takes three arguments, a formula in the form `y ~ x1 + x2` defining the dependent (Y) and one or more independent variables (`x1`, `x2`, ...), a function (FUN), and a dataframe (data). When you execute `aggregate(y ~ x1 + x2 + ..., data, FUN)`, R will apply the input function (FUN) to the dependent variable (Y) *separately* for each level(s) of the independent variable(s) (`x1`, `x2`, ...). Let's see how it works:

### `aggregate()`

#### formula

A formula in the form `y ~ x1 + x2 + ...` where `y` is the dependent variable, and `x1`, `x2`, ... are the index (independent) variables. For example, `salary ~ sex + age` will aggregate a salary column at every unique combination of age and sex

#### FUN

A function that you want to apply to `x` at every level of the independent variables. For example, `FUN = mean` will calculate the mean of the data

#### data

The dataset containing the variables in formula

...

Optional arguments passed on to FUN (like `na.rm = T` to ignore NA values in `x`)

Let's give `aggregate()` a whirl. We'll use the function to answer the question "What is the flight cancellation rate for each carrier?" For this question, we'll set the value of Y to `cancelled`, `x1` to `carrier`, and FUN to `mean`

```
aggregate(formula = cancelled ~ carrier, # DV is cancelled, IV is carrier
          FUN = mean, # Calculate the mean of the DV for each IV level
          na.rm = T, # Ignore NA values when calculating the mean
          data = Flights # IV and DV are located in the Flights dataframe
          )
```

```
## carrier cancelled
## 1 AA 0.018495684
## 2 AS 0.000000000
## 3 B6 0.025899281
## 4 CO 0.006782614
## 5 DL 0.015903067
## 6 EV 0.034482759
## 7 F9 0.007159905
## 8 FL 0.009817672
## 9 MQ 0.029044750
## 10 OO 0.013946828
## 11 UA 0.016409266
## 12 US 0.011268986
## 13 WN 0.015504047
## 14 XE 0.015495599
## 15 YV 0.012658228
```

As you can see, the `aggregate()` function has returned a dataframe with a column for the independent variable (`carrier`), and a column for the results of the function mean applied to each level of the independent variable. We can easily plot these data using the `barplot()` function, which plots a numeric variable as a function of a nominal variable (see margin Figure 41)

You can include any function as the argument to `FUN` as long as the function takes a single numeric argument and returns a single scalar. For example, let's use `aggregate()` to now get the median arrival and departure delays for each airline carrier:

```
# Calculate median departure delay by carrier
med.arrdelay <- aggregate(formula = arr_delay ~ carrier, # DV is arr_delay, IV is carrier
  FUN = median, # Calculate the median arr_delay
  data = Flights, # Columns are in the Flights dataframe
  na.rm = T) # Ignore NA values

# Calculate median departure delay by carrier
med.depdelay <- aggregate(formula = dep_delay ~ carrier,
  FUN = median,
  data = Flights,
  na.rm = T) # Ignore NA values
```

Let's combine two results from `aggregate()` into one plot. In margin Figure 42, I compare the median arrival and departure delays for each airline. Interestingly, we can see that the carrier `US` has both a very low median arrival and departure delay. In fact, they're both negative, suggesting that `US` flights tend to leave and arrive early. In contrast, the carrier `UA` has a median arrival and departure delay of 0. If you ask me, `UA` has the better numbers as their flights, on average, depart and arrive when they are supposed to!

While `aggregate()` is good for calculating summary statistics for a single dependent variable, it can't handle multiple dependent variables. For example, if you had a survey and you wanted to calculate summary statistics for multiple dependent variables (like

```
aggregated.data <- aggregate(formula = cancelled ~ carrier,
  FUN = mean, na.rm = T, data = Flights)

barplot(height = aggregated.data$cancelled,
  names.arg = aggregated.data$carrier)
```

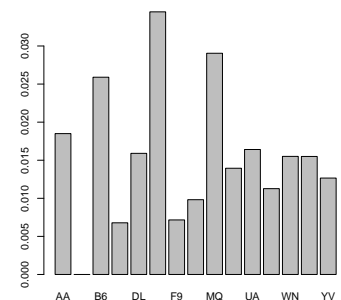


Figure 41: Barplot showing the mean flight cancellation rate for airlines in the `Flights.txt` dataset

```
require("RColorBrewer")

# Step 1: Calculate median arrival delays
med.arrdelay <- aggregate(formula = arr_delay ~ carrier,
  FUN = median,
  data = Flights,
  na.rm = T) # Ignore NA values

# Step 2: Calculate median departure delays
med.depdelay <- aggregate(formula = dep_delay ~ carrier,
  FUN = median,
  data = Flights,
  na.rm = T) # Ignore NA values

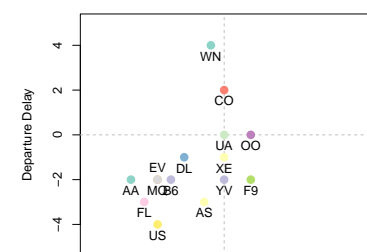
# Step 3: Set up plot and add reference line
plot(1, xlim = c(-10, 10), ylim = c(-5, 5),
  type = "n", main = "Aggregate Carrier Flight Delays",
  xlab = "Arrival Delay", ylab = "Departure Delay")

abline(h = 0, lty = 2, col = "gray")
abline(v = 0, lty = 2, col = "gray")

# Step 4: Add points
points(x = med.arrdelay$arr_delay,
  y = med.depdelay$dep_delay,
  col = brewer.pal(12, "Set3"),
  pch = 16, cex = 1.5)

# Step 5: Add carrier labels
text(x = med.arrdelay$arr_delay,
  y = med.depdelay$dep_delay,
  labels = med.depdelay$carrier,
  pos = c(rep(1, 5), 3, rep(1, 9)) # Insert 3 for carrier EV
)
```

Aggregate Carrier Flight Delays



income, weight, age, etc.), you'd need to execute an `aggregate()` command for each dependent variable, and then combine the results into a single dataframe. Thankfully, a recently released R package called `dplyr` makes this process very simple!

### *Aggregation with dplyr*

The `dplyr` package is a new package that allows you to do data 'wrangling' (manipulating datasets) quickly and easily. In this section, we'll go over a very brief overview of how the functions in `dplyr` work. However, this will be a very brief overview and I strongly recommend you look at the help menu for `dplyr` for additional descriptions and examples.

`dplyr` works by combining objects (dataframes and columns in dataframes), functions (mean, median, etc.), and *verbs* (special commands in `dplyr`). In between these commands is a new operator called the *pipe* which looks like this: `%>%`. The pipe simply tells R that you want to continue executing some functions or verbs on the object you are working on. You can think about this pipe as meaning 'and then...'

To aggregate data with `dplyr`, your code will look something like the following code. In this example, assume that the original (raw) dataframe is called `my.df`, the variable you want to collapse the data over is called `grouping.column`, and the columns you want to aggregate are called `col.a`, `col.b` and `col.c`

```
my.df %>% # Specify original dataframe
  group_by(grouping.column) %>% # Grouping variable
  summarise(
    a.mean = mean(col.a), # calculate mean of column col.a in my.df
    b.sd = sd(col.b),      # calculate sd of column col.b in my.df
    c.max = max(col.c)     # calculate max on column col.c in my.df, ...
  )
```

Here's how you should think about the code above:

Start with the dataframe `my.df`. Then, group `my.df` by the grouping variable `grouping.variable`. Then, calculate the following summary columns in the dataset: `a.mean` should be the mean of `col.a` in `my.df`, `b.sd` should be the standard deviation of `col.b` in `my.df`, and `c.max` should be the maximum value of `col.c` in `my.df`.

When you use `dplyr`, you write code that sounds like: "The original dataframe is XXX, now filter the dataframe to only include rows that satisfy the conditions YYY, now group the data at each level of the variable(s) ZZZ, now summarize the data and calculate summary functions XXX..."

Let's start with an example: Let's create a dataframe of aggregated data from the `Flights` dataset, where each row is an airline carrier, and each column is a different summary statistic of some data across all flights. Specifically, let's create 5 columns: `dep.delay.med`: The median departure delay of that airline, `dep.delay.sd`: The standard deviation of departure delays of that airline, `arr.delay.max`: The maximum departure delay of that airline, and `dist.mean`: The average flight distance of that airline.

To create this aggregated data frame, I will use the new function `group_by` and the verb `summarise`. I will assign the result to a new dataframe called `flights.agg`:

```
require(dplyr)
flights.agg <- Flights %>% # Define dataframe, THEN...
  group_by(carrier) %>% # Define the grouping variable, THEN...
  summarise( # Tell R you are going to calculate summaries
    dep.delay.med = median(dep_delay, na.rm = T), # Define first summary...
    dep.delay.sd = sd(dep_delay, na.rm = T), # Define second summary...
    arr.delay.max = max(arr_delay, na.rm = T),
    dist.mean = mean(dist, na.rm = T)
  ) # End

flights.agg # Print result!
```

```
## Source: local data frame [15 x 5]
##
##   carrier dep.delay.med dep.delay.sd arr.delay.max dist.mean
## 1      AA           -2    35.35627         978   483.8212
## 2      AS           -3    20.29107         183  1874.0000
## 3      B6           -2    42.86727         335  1428.0000
## 4      CO            2    25.89858         957  1098.0549
## 5      DL           -1    39.94030         701   723.2832
## 6      EV           -2    42.44875         469   775.6815
## 7      F9           -2    23.71563         277   882.7411
## 8      FL           -3    31.64888         500   685.4063
## 9      MQ           -2    43.72860         918   650.5310
## 10     OO            0    27.53757         380   819.7279
## 11     UA            0    45.64341         861  1177.8388
## 12     US           -4    22.80906         433   981.4677
## 13     WN            4    29.38784         499   606.6218
## 14     XE           -1    28.09030         634   589.0326
## 15     YV           -2    13.65217          72   938.6709
```

As you can see, our final object `flights.agg` is the aggregated dataframe we want which aggregates all the columns we wanted at the level of the airline carrier. Let's walk through the code

- First, we define the original dataframe that we are basing our summary statistics on. In this case, the original dataframe is `Flights`. We then include the pipe `%>%` to tell R we are still working.
- Second, we define the grouping variable using the `group_by` function. This tells R to group the results at the level of each carrier. We then use the `%>%` pipe.

When you use `dplyr`, the output will always be an object called a *local data frame*. A local dataframe is identical to a regular dataframe, except that it looks a bit nicer if you print the entire dataframe into the console. This means you don't have to use the `head()` function when looking at a local dataframe.

- Next, we call the `summarise` function, which tells R that the following functions will be summaries of the grouping variable. Because all the arguments to the `summarise` function are within the parentheses, we don't need to use a pipe.
- Finally, we define the summary columns in our final dataframe. For each column, we give it a name (e.g.; `dep.delay.med`), and then write the calculation as a function of the appropriate columns in the original dataframe. For example, to define `dep.delay.med`, we write `median(dep_delay, na.rm = T)`.

Hopefully you can see that this `dplyr` code is *much* simpler than the code we'd have to use if we wanted to create all these summary columns using `aggregate`.

### *The 5 verbs in dplyr*

In the example above, we used the `dplyr` verb `summarise`. However, `dplyr` has other verbs that are just as useful:

dplyr verbs	
<code>filter</code>	Select a subset of rows in a dataframe. For example, the following code will restrict data to males
<code>arrange</code>	Reorders rows of a dataframe.
<code>select</code>	Select specific columns of a dataframe.
<code>mutate</code>	Add a column to a dataframe.
<code>summarise</code>	Creates summary columns as a function of columns in the original dataframe. Note: Only use after specifying <code>group_by</code> variables.

Let's do an example where we combine multiple verbs into one chunk of code. We'll create a new dataframe called `flights.cancelled` that gives us data on the cancelled flights from each airline. However,

let's add some additional data filters this time. We'll filter the data to only include flights that were scheduled to depart before 12pm, and whose destination is either DFW (Dallas Fort Worth), MIA (Miami), or DSM (Des Moines).

```
require(dplyr)
flights.agg <- Flights %>% # Step 1
  filter(hour < 12 & dest %in% c("DFW", "MIA", "DSM")) %>% # Step 2
  group_by(carrier) %>% # Step 3
  summarise( # Step 4
    cancelled.p = mean(cancelled), # Step 5
    n = length(cancelled) # Step 6
  ) %>%
  arrange(n) # Step 7

flights.agg # Print the result!

## Source: local data frame [5 x 3]
##
##   carrier cancelled.p    n
## 1      00 0.0000000000    79
## 2      XE 0.0000000000   199
## 3      CO 0.0000000000   675
## 4      MQ 0.0008598452  1163
## 5      AA 0.0000000000  1598
```

As you can see, our result is a dataframe with 5 rows and 3 columns. The reason why we only have 5 rows is because only 5 carriers had flights to one of the three airports we specified. Let's walk through the code line by line:

1. First, we define the original dataframe as `Flights`, (`%>%` then...)
2. Next we filter the `Flights` dataframe by only including rows where the hour is less than 12, and the destination is in the set DFW, MIA and DSM. (`%>%` then...)
3. We group the data according to `carrier` (`%>%` then...)
4. We call the `summarise` verb, telling `dplyr` that the next commands are summary functions of `Flights`. These will be the columns in our new aggregated dataframe. (`%>%` then...)
5. The first column in our new aggregated dataset is called `cancelled.p` and is defined as the mean value of the `cancelled` column in `Flights`.
6. The second column is called `n` and is simply the length of the vector `cancelled`. Because this function is called for each level of `carrier`, this will give us the number of observations for each level of `carrier`.
7. Arrange the final dataframe by the new column `n`



## Merging two dataframes

Merging two dataframes together allows you to combine information from both dataframes into one. For example, a teacher might have a dataframe called `students` containing information about her class. She then might have another dataframe called `exam1scores` showing the scores each student received on an exam. To combine these data into one dataframe, you can use the `merge()` function. For those of you who are used to working with Excel, `merge()` works a lot like `vlookup` in Excel:

### `merge()`

`x, y`

2 dataframes to be merged

`by, by.x, by.y`

The names of the columns that will be used for merging. If the the merging columns have the same names in both dataframes, you can just use `by = c("col.1", "col.2"...)`. If the merging columns have different names in both dataframes, use `by.x` to name the columns in the `x` dataframe, and `by.y` to name the columns in the `y` dataframe. For example, if the merging column is called `STUDENT.NAME` in dataframe `x`, and `name` in dataframe `y`, you can enter `by.x = "STUDENT.NAME", by.y = "name"`

`all.x, all.y`

A logical value indicating whether or not to include non-matching rows of the dataframes in the final output. The default value is `all.y = FALSE`, such that any non-matching rows in `y` are not included in the final merged dataframe.

A generic use of `merge()`, looks like this:

```
new.df <- merge(x = df.1, # First dataframe
                y = df.2, # Second dataframe
                by = "column" # Common column name in both x and y
                )
```

where `df.1` is the first dataframe, `df.2` is the second dataframe, and `"column"` is the name of the column that is common to both dataframes.

For example, let's say that we wanted to add a column to our happiness dataframe `hsurvey` showing the continent that each person was from. In the dataframe, we only have country information, so how can we add continent info? We can do this by merging the `hsurvey` dataframe with a new dataframe called `continents`. By merging the dataframes, we will add all the information from `continents` to the `hsurvey` dataframe as additional rows. We'll also match the two dataframes with the column `country`.

```
continents <- data.frame(country = c("USA", "Canada", "Mexico", "India", "Portugal"),
                        continent = c("North America", "North America", "North America", "Asia", "Europe"),
                        gdp = c(53042, 51964.3, 10307.3, 1497.5, 21738.3),
                        stringsAsFactors = F
                        )

hsurvey <- merge(x = hsurvey, # The first dataframe
               y = continents, # The second dataframe
               by = "country" # The name of the matching column
               )

## Error in merge(x = hsurvey, y = continents, by = "country"): object
## 'hsurvey' not found
```

To see if this worked, let's look at the first few rows of `hsurvey`

```
head(hsurvey)

## Error in head(hsurvey): object 'hsurvey' not found
```

As you can see, our new merged dataframe has added the information from `continents` to the `hsurvey` dataframe by matching rows with the `country` column.

### *Easily recode values in a dataframe with `merge()`*

You can also use the `merge()` function to quickly recode values in a vector in a dataframe. For example, let's say some drunk pirate accidentally entered the wrong country names for each person in the column `country` in `hsurvey()`. Thankfully, when we got him sober enough he could tell us which (incorrect) entry goes to which correct entry. For example: the entry `USA` should be `Canada`, `Mexico` should be `India`, etc. We can quickly correct this in our original dataframe using `merge()`. We'll do this in two steps:

1. Create a lookup table that shows the original (incorrect) and the new (correct) values
2. Merge the the original data with the lookup table

Let's create the lookup table called `lookup`. This will be a dataframe with two columns: `country.wrong` - the original country values, and

`country.true` the corrected values. Each row in the dataframe will connect the original incorrect value with the new, correct value.

```
lookup <- data.frame("country.wrong" = c("USA", "Canada", "Mexico", "India", "Portugal"),
                     "country.true" = c("Canada", "Mexico", "India", "Portugal", "USA")
)
```

Now, let's merge `hsurvey` with `lookup`. To do this, we'll need to specify that the matching column in `hsurvey` is called `country`, and the matching column in `lookup` is called `country.wrong`. When we merge these two dataframes, R will add the `country.true` column to `hsurvey`

```
hsurvey <- merge(x = hsurvey,
                 y = lookup,
                 by.x = "country",
                 by.y = "country.wrong"
)
```

```
## Error in merge(x = hsurvey, y = lookup, by.x = "country",
## by.y = "country.wrong"): object 'hsurvey' not found
```

Let's see if it worked:

```
head(hsurvey)
```

```
## Error in head(hsurvey): object 'hsurvey' not found
```



## 10: 1 and 2-sample Null-Hypothesis tests

### Chapter Goals

1. Learn about hypothesis test objects in R
2. One and two sample tests: Correlations, t-tests and chi-square

Do we get more treasure from chests buried in the sand or at the bottom of the ocean? Is there a relationship between the number of scars a pirate has and how much grogg he can drink? Are pirates with nipple rings more likely to wear bandannas than those without nipple rings? Glad you asked, let's see how we can answer these questions some hypothesis tests.

### Warning about null-hypothesis tests with "frequentist" statistics

Until recently, null-hypothesis testing using frequentist statistics has been the most popular method of conducting inferential statistics. However, it has serious flaws. While I can't go into the details here, I can point out that the main flaw is that frequentist statistics don't give you the information you really want to know. For example, imagine that you are comparing the effectiveness of a cancer drug to a placebo. After conducting a double-blind study, where you give some patients the placebo and some patients the drug, you want to know the probability that that the drug is better than a placebo. Unfortunately frequentist statistics cannot give you this information. They can only tell you the probability of getting a specific result *given* that the null hypothesis (in this case, that the drug is equally as effective as the placebo) is true. If that sounds confusing, it's because it is. A better alternative is Bayesian statistics which *can* give you posterior probability information. Unfortunately, Bayesian statistics can be computationally demanding, so in the past we've lived with frequentist statistics and tried to ignore its fundamental flaws. However, given improvements in processing speed, we can now easily conduct Bayesian alternatives to frequentist tests on modern computers.

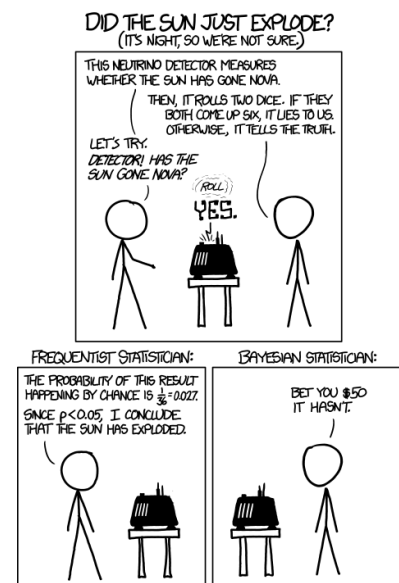


Figure 43: xkcd comic. Currently used without any permission.

We will cover Bayesian statistics in Chapter X and I strongly encourage you to adopt them in your own analyses. However, for the purposes of completeness, I'll show you how to conduct most of the standard frequentist tests here.

### *T-test*

We use t-tests to compare the sample mean of data to some hypothesized mean. In a one sample test, we use one set of observations to test whether or not the population mean is different from a hypothesized value. In a two-sample test, we use two sets of observations to test whether or not the two populations have different means.

The t-test function in R is `t.test()`. The `t.test()` function can take several arguments, here I'll emphasize a few of them. To see them all, check the help menu for `t.test` (`?t.test`).

There are two ways to use the `t.test()` function. The first way is to enter one (or more) vectors as arguments to the function as follows:

`t.test(x, y)`: Conduct either a one sample t-test on a vector `x`, or a two-sample t-test on two vectors `x` and `y`.

#### `t.test(x, y)`

`x, y`

Either one vector of data (`x`) for a one-sample t-test, or two vectors (`x, y`) for a two-sample test.

`alternative`

A character string indicating whether the test is two-tailed or one-tailed (including the direction). Type "t" for two-tailed, "g" for a 'greater than' one-tailed test, or "l" for a "less than" one-tailed test.

`mu`

he population mean under the null hypothesis.

`paired, var.equal`

`paired`: A logical value (either T or F) indicating whether the test is paired (T) or unpaired (F). Only use this for two-sample tests.

`var.equal`: A logical value indicating whether or not you treat the two variances as equal.

Let's do an example using the pirate survey dataset. If you haven't

downloaded the pirate survey dataset, check Chapter 1 for instructions.

### One-Sample t-test

The format for a one-sample t-test is as follows:

```
t.test(x = data, # A vector of data
      mu = 0, # The null hypothesis
      alternative = "t" # Two tailed test (use "l" or "g" for one
      )
```

where  $x$  is a vector of data,  $\mu$  is the population mean under the null hypothesis, and  $\text{alternative}$  is "t" for a two-tailed test, or "l" or "g" for a one-tailed test.

Let's do a one-sample t-test on the age of pirates in our survey. Specifically, let's see if the average age of the pirates is significantly different from 20. In this case, our vector  $x$  is `pirates$age`:

```
test.result <- t.test(x = pirates$age, # Vector of data to test
                     mu = 20, # Null hypothesis is mean = 30
                     alternative = "t" # Two-tailed test
                     )
```

You'll notice that when you assign the `t.test` to an object (in this case we called it `test.result`), you do not see any output. To see the output of the test, you need to tell R to print the object by executing the name of the test object:

```
test.result # Print the results of the t.test

##
## One Sample t-test
##
## data: pirates$age
## t = 2.4556, df = 999, p-value = 0.01424
## alternative hypothesis: true mean is not equal to 20
## 95 percent confidence interval:
## 32.15961 128.91439
## sample estimates:
## mean of x
## 80.537
```

Now, you can see the main output of the test. Looks like we got a test statistic of 2.46 and a resulting p-value that's pretty darn small. But what if you want to access specific values like the test statistic or the p-value? Thankfully, this is easy in R. To see which information you can extract from the t-test object, apply the `names()` function to the test object:

```
par(mar = c(10, 0, 3, 1))
plot(1, xlim = c(-10, 10), ylim = c(0, 1),
     yaxt = "n", main = "One Sample t-test",
     type = "n", bty = "n", ylab = "", xlab = "")

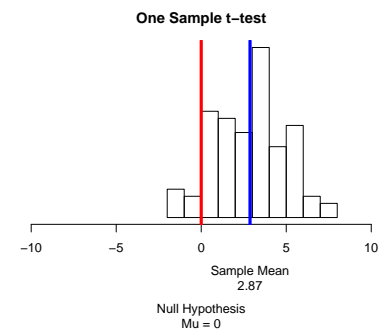
samples <- rnorm(100, mean = 3, sd = 2)

par(new = T)
hist(samples, yaxt = "n", xaxt = "n", xlab = "",
     ylab = "", main = "", xlim = c(-10, 10))

mtext(paste("Sample Mean\n", round(mean(samples), 2), sep = "\n"),
      side = 1, line = 3.5, at = mean(samples))

abline(v = mean(samples), lty = 1, col = "blue", lwd = 4)

mtext("Null Hypothesis\nMu = 0", side = 1, line = 6, at = 0)
abline(v = 0, lty = 1, col = "red", lwd = 4)
```



If you want to see what information is in a test object, just apply `names()` to the object. You can then extract specific information with `$`

```
names(test.result)

## [1] "statistic" "parameter" "p.value" "conf.int" "estimate"
## [6] "null.value" "alternative" "method" "data.name"
```

From this vector of names, I see that I can extract the test statistic with the name `statistic` and the p-value with the name `p.value`. To get these from the test object, use the `$` operator:

```
test.result$statistic # Show me the test statistic

##          t
## 2.455574

test.result$p.value # Show me the p.value

## [1] 0.01423541
```

Being able to quickly extract key numerical information from a test object is huge. For one thing, it allows you to automate the process of running statistical tests over different datasets or simulations. In Chapter XX, we'll see how you can use loops to do this in a snap.

### Two-sample t-test

In a two-sample t-test, we use two sets of observations drawn from two different populations and test whether or not the two populations have the same mean. To conduct a two-sample t-test, we simply enter two vectors as arguments `x` and `y`.

Let's use this convention to compare the ages of pirates who wear headbands and pirates who don't wear headbands. First, we'll create the two vectors `age.headband` and `age.noheadband` containing the age data for pirates who do and do not wear headbands. We'll then enter these vectors as arguments `x` and `y` to `t.test()`:

```
age.headband <- subset(pirates, subset = headband == "yes")$age # Get the first vector
age.noheadband <- subset(pirates, subset = headband == "no")$age # Get the second vector

test.result <- t.test(x = age.headband, # Enter the first vector
                     y = age.noheadband, # Enter the second vector
                     alternative = "two.sided" # Specify a two-tailed test
                     )

test.result

##
## Welch Two Sample t-test
##
## data: age.headband and age.noheadband
## t = -0.7392, df = 96.025, p-value = 0.4616
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -364.7624 166.8143
## sample estimates:
## mean of x mean of y
## 59.33779 158.31183
```

```
par(mar = c(10, 0, 3, 1))
plot(1, xlim = c(-10, 10), ylim = c(0, 1),
     yaxt = "n", main = "Two Sample t-test",
     type = "n", bty = "n", ylab = "", xlab = "")

samples.1 <- rnorm(100, mean = 3, sd = 2)

par(new = T)
hist(samples.1, yaxt = "n", xaxt = "n", xlab = "",
     ylab = "", main = "", xlim = c(-10, 10),
     col = rgb(0, 0, 1, alpha = .1))

mtext(paste("Sample Mean\n", round(mean(samples.1), 2), sep = ""),
     side = 1, line = 3.5, at = mean(samples.1))

abline(v = mean(samples.1), lty = 1,
       col = rgb(0, 0, 1, alpha = 1), lwd = 4)

samples.2 <- rnorm(100, mean = -3, sd = 2)

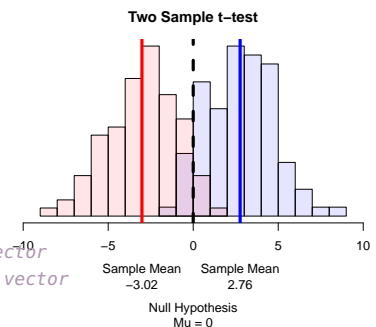
par(new = T)
hist(samples.2, yaxt = "n", xaxt = "n", xlab = "",
     ylab = "", main = "", xlim = c(-10, 10),
     col = rgb(1, 0, 0, alpha = .1))

mtext(paste("Sample Mean\n", round(mean(samples.2), 2), sep = ""),
     side = 1, line = 3.5, at = mean(samples.2))

abline(v = mean(samples.2), lty = 1,
       col = rgb(1, 0, 0, alpha = 1), lwd = 4)

mtext("Null Hypothesis\nMu = 0", side = 1, line = 6, at = 0)

abline(v = 0, lty = 2, col = "black", lwd = 4)
```





Looks like we see a test statistic of  $-0.74$  with a resulting p-value of 0.46. According to null-hypothesis test logic, we fail reject the null hypothesis.

### *Specifying t-tests with formula notation*

The second (and I think better) way to specify arguments to the `t.test()` function is by using the `formula` and `subset` arguments. Using this notation, we specify the dependent and independent variables as a formula in the form `dv ~ iv` where `dv` is the dependent variable, and `iv` is the independent variable with two levels in a dataframe. As you'll see, this convention is a bit nicer to use when working with data in dataframes because we don't need to define two separate vectors prior to the test.

#### `t.test(formula, data)`

##### `formula`

An (optional) formula in the form `dv ~ iv` where `dv` is the dependent variable, and `iv` is the independent variable with two levels in a dataframe. Specify the dataframe in the `data` argument

##### `data`

The dataframe containing the columns specified in `formula`.

##### `subset`

A logical vector indicating a subset of data to use. If the independent variable in the formula specification has more than two values, you'll need to use `subset` to restrict the data to only two values of the `iv`.

##### `alternative, mu, paired, var.equal`

Additional arguments (see previous `t.test()` description)

Using this formulation, we don't have to define separate vectors (`x` and `y`) prior to conducting the test. Instead, we can use the `formula` argument to tell R which columns in a dataframe correspond to the dependent and independent variables. When you use the formula version of `t.test`, the independent variable *must* only have two possible values. If R finds more than two values in the `iv`, it will return an error. To ensure that there are only two values present, include

the appropriate subset argument. For example, if the independent variable is sex and you want to compare males and females, include the argument `subset = sex %in% c("male", "female")`

Let's repeat the previous t-test using a formula. To do this, we'll specify three new arguments:

- `data = pirates`: The columns in formula come from the dataframe `pirates`
- `formula = age ~ headband`: Conduct a test on age as a function of headband.
- `subset = headband %in% c("yes", "no")`: Restrict our analysis to pirates who answered "yes" or "no" to whether or not they wear a headband.

Here's how the alternative notation for the same test looks:

```
test.result <- t.test(formula = age ~ headband, #dv is weight, iv is Diet
                      subset = headband %in% c("yes", "no"), # Only use valid headband values
                      data = pirates, # Dataframe is pirates
                      alternative = "two.sided" # Two-sided test
                      )

test.result

##
## Welch Two Sample t-test
##
## data: age by headband
## t = 0.7392, df = 96.025, p-value = 0.4616
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -166.8143 364.7624
## sample estimates:
## mean in group no mean in group yes
## 158.31183 59.33779
```

As you can see, the results of this test and the prior test are identical. You can use whichever version makes more sense to you. Personally, I like the formula version because all the necessary commands (including the specification that the two diets are 1 and 2) are contained within the `t.test()` function. Of course, you can also specify additional restrictions in the subset argument.

Let's try making the previous test a little more complicated by adding a subset argument. Let's say a pirate tells you "Oh, well there's only a difference between the age of headband and no-headband pirates for those pirates who went to college at Captain

Chunk's Canon Crew." We can test this by adding the additional restriction `college == "Captain Chunk's Canon Crew"`, to the `subset` argument:

```
test.result <- t.test(formula = age ~ headband,
                      subset = headband %in% c("yes", "no") &
                        college == "Captain Chunk's Canon Crew",
                      data = pirates,
                      alternative = "two.sided"
                      )

test.result

##
## Welch Two Sample t-test
##
## data: age by headband
## t = 0.8812, df = 60.233, p-value = 0.3817
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -230.9589 594.7197
## sample estimates:
## mean in group no mean in group yes
##      229.96667      48.08629
```

Looks like we still don't find a significant difference in age between headband and no-headband wearers, even just for pirates who went to Captain Chunk's Canon Crew.

### Correlation test

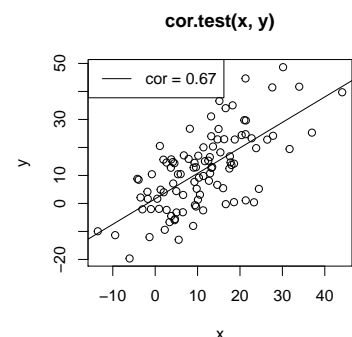
Next we'll cover two-sample correlation tests. Recall that in a correlation test, you are accessing the relationship between two variables on a ratio or interval scale.

To run a correlation test, use the `cor.test(x, y)` function. The test has the following arguments

```
x <- rnorm(n = 100, mean = 10, sd = 10)
y <- x + rnorm(n = 100, mean = 0, sd = 10)

plot(x, y, main = "cor.test(x, y)")
abline(lm(y ~ x))

legend("topleft",
       legend = paste("cor = ", round(cor(x, y), 2), sep = ""),
       lty = 1)
```



`cor.test(x, y)`: Conduct a correlation test between two vectors `x` and `y`.

## cor.test()

**x, y**

Two numeric data vectors of the same length

**alternative**

A string indicating the direction of the test. You can use "t" for two-sided, "l", for less than, and "g" for greater than.

**method**

A string indicating which correlation coefficient to test. You can use "pearson", "kendall", or "spearman". The default is Pearson.

**conf.level**

The confidence level for the Pearson correlation coefficient.

Let's conduct a correlation test on the age of pirates and the number of parrots they've had in their lifetime. We'll set `x = pirates$age`, and `y = pirates$parrots.lifetime`

```
test.result <- cor.test(x = pirates$age,
                       y = pirates$parrots.lifetime
                       )

test.result

##
## Pearson's product-moment correlation
##
## data: pirates$age and pirates$parrots.lifetime
## t = 2.1327, df = 998, p-value = 0.03319
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.005385926 0.128811730
## sample estimates:
##      cor
## 0.06735652
```

Looks like we have a positive correlation of 0.07! To see what information we can extract for this test, let's run the `names()` on the test object:

```
names(test.result)

## [1] "statistic" "parameter" "p.value" "estimate" "null.value"
## [6] "alternative" "method" "data.name" "conf.int"
```

Looks like we've got a lot of information in this test object. As an example, let's look at the confidence interval:

```
test.result$conf.int
## [1] 0.005385926 0.128811730
## attr(,"conf.level")
## [1] 0.95
```

You'll notice that when we tried to access the confidence interval, we got an additional piece of information called `attr("conf.level")`. This means that the result of the command `test.result$conf.int` not only contains the bounds of the confidence interval, but also the level of confidence. This is a good thing because the confidence interval only makes sense in terms of the level of confidence used to calculate the interval.

### *Chi-square test*

Next, we'll cover chi-square tests. In a chi-square test test, we test whether or not there is a relationship between two variables on a nominal scale (like sex, eye color, first name etc.). To conduct a chi-square test, we use the `chi.square()` function.

#### **chisq.test()**

**x, y**

Two vectors (can be numeric, factor, or string) of the same length. Alternatively, you can simply enter a matrix as the x argument and ignore the y argument.

**correct**

a logical indicating whether to apply continuity correction when computing the test statistic for 2 by 2 tables: one half is subtracted from all  $|O - E|$  differences; however, the correction will not be bigger than the differences themselves.

Let's use the `chisq.test()` function to test if there is a relationship between the college a pirate went to and the type of sword he/she uses. We'll use the `tattoo` and `college` vectors in `pirates`

```
test.result <- with(pirates,
  chisq.test(x = college,
             y = sword.type)
)

test.result
```

```
##
## Pearson's Chi-squared test
##
## data: college and sword.type
## X-squared = 2.088, df = 3, p-value = 0.5543
```

Looks like we got a test-statistic of 2.09 and a p-value of 0.55. According to this test, we do not have sufficient data to reject the null hypothesis that there is no relationship between the two variables.

Let's see what other information we can get from the chi-square test object:

```
names(test.result)

## [1] "statistic" "parameter" "p.value" "method" "data.name" "observed"
## [7] "expected" "residuals" "stdres"
```

I encourage you to run the `names()` function on statistical objects. You never know what interesting things you'll discover!

We've got some interesting new options here. Let's look at the value of `observed`, the observed frequencies in the data

```
test.result$observed

##                                sword.type
## college                        banana cutlass sabre
## Captain Chunk's Canon Crew           25    549    44
## Jack Sparrow's School of Fashion and Piratry    17    268    20
##                                sword.type
## college                        scimitar
## Captain Chunk's Canon Crew                41
## Jack Sparrow's School of Fashion and Piratry    26
```

Cool. It looks like R stores a table of the observed frequencies and the expected frequencies under the null-hypothesis. Thanks R!

# 11: Regression and ANOVA

## Chapter Goals

1. Learn about regression
2. ANOVA

## The Linear Model

The linear model is easily the most famous and widely used model in all of statistics. Why? Because it can apply to so many interesting research questions where you are trying to predict a variable of interest (the *response* or *dependent variable*) on the basis of one or more other variables (the *predictor* or *independent variables*).

The linear model takes the following form:

$$y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots\beta_nx_n$$

where the  $x$  values represent the predictors, while the beta values represent weights.

To use the linear model in R, we use the `lm()` function:

### `lm()`

#### function

A function in a form  $y \sim x_1 + x_2 + \dots$ , where  $y$  is the dependent variable, and  $x_1, x_2, \dots$  are the independent variables.

#### data

The dataframe containing the columns specified in the formula.

#### subset

An optional vector specifying a subset of observations to be used in the fitting process. For example `subset = age > 50`

Let's try an example with some made-up data. For this example, I'll create three independent variables  $x_1$ ,  $x_2$  and  $x_3$  from normal distributions. I'll then create a dependent variable  $y$  as a linear function of the three independent variables (with a little error thrown in)<sup>3</sup>. Finally, I'll run the linear model and see if we can recover the true beta values:

```
# Step 1: Create a dataframe of predictors

random.data <- data.frame(
  "x1" = rnorm(100, mean = 0, sd = 1),
  "x2" = rnorm(100, mean = 4, sd = 5),
  "x3" = rnorm(100, mean = -2, sd = 2)
)

# Step 2: Create the DV with beta values 0, 1, 2, 3
random.data$y <- with(random.data, 0 + 1 * x1 + 2 * x2 + 3 * x3)

# Step 3: Add some random noise to the DV
random.data$y <- random.data$y + rnorm(100, mean = 0, sd = 1)

# Step 4: Run the model then print the result
result <- lm(y ~ x1 + x2 + x3, data = random.data)
result

##
## Call:
## lm(formula = y ~ x1 + x2 + x3, data = random.data)
##
## Coefficients:
## (Intercept)      x1      x2      x3
##   -0.1173    1.0113    2.0116    3.0148
```

Looks like our estimated beta values of  $-0.1173$ ,  $1.0113$ ,  $2.0116$  and  $3.0148$  are pretty close to the true values of  $0$ ,  $1$ ,  $2$ ,  $3$ ! This means the model did a pretty good job of estimating the true beta values from the (noisy) data.

We can get lots of other information from the linear model object. Here are three of them (to see all of them, run `names(result)`):

- `coefficients`: A vector of the estimated beta values
- `residuals`: A vector of the differences between the true response values and the fitted response values.
- `fitted.values`: A vector of the fitted values.

These attributes let us easily calculate some interesting model diagnostics. For example, let's see how far the model fits are on average from the true values:

<sup>3</sup> If there was no error in the model, then the response variable would be a perfect linear combination of the predictor variables. When this happens, the model will always perfectly fit the data and no stats are necessary (or even possible!)



```
abs.resid <- abs(result$residuals) # Calculate the absolute value of the residuals
mean(abs.resid) # Calculate the mean

## [1] 0.7175083
```

So it looks like, on average, our model fits are 0.72 off from the true values (This value is driven by the standard deviation in errors, which we set to 1).

The linear model assumes that there should be no relationship between the predicted values and the distribution of errors. We can easily check this with a scatterplot with the predicted values on the x-axis, and residuals on the y-axis (see Figure 44 in the margin).

### Generalized Linear Model (GLM)

In the Generalized Linear Model (GLM), we take the original linear model, but apply a link function that wraps around the linear combination of predictors

```
plot(x = result$fitted.values,
     y = result$residuals,
     main = "Model Diagnostics")
```

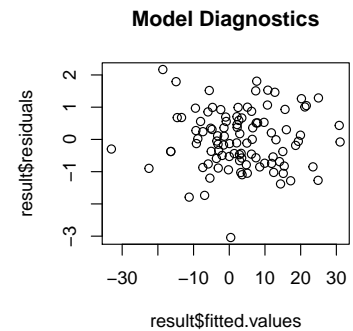


Figure 44: We see no evidence for a relationship between predicted values and residuals. This is good!

### glm()

function, data, subset

The same arguments as in `lm()`

family

One of the following strings, indicating the link function for the general linear model

- "binomial": Binary logistic regression, useful when the response is either 0 or 1.
- "gaussian": Standard linear regression. Using this family will give you the same result as `lm()`
- "Gamma": Gamma regression, useful for exponential response data
- "inverse.gaussian": Inverse-Gaussian regression, useful when the dv is strictly positive and skewed to the right.
- "poisson": Poisson regression, useful for count data. For example, "How many parrots has a pirate owned over his/her lifetime?"

The key new argument in `glm()` compared to `lm()` is the `family` argument. This argument tells R which link function to use. To see more information about the families, look at help under `?family`.

### Binary Logistic Regression

Probably the most common non-Normal family you will use is binomial which corresponds to binary logistic regression. In binary logistic regression, we predict a binary outcome variable (containing 0s and 1s) as the logit transformation of a linear combination of a set of predictors. Formally:

$$p(y = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$

To conduct binary logistic regression, we use the `family = "binomial"` argument. Let's try an example with some more made-up data. We'll set the true beta values to 0, 1, -2 and 3

```
# Step 1: Create a dataframe of predictors
random.data <- data.frame(
  "x1" = rnorm(500, mean = 0, sd = 2),
  "x2" = rnorm(500, mean = 4, sd = 2),
  "x3" = rnorm(500, mean = -2, sd = 2)
)

# Step 2: Create the DV with beta values 0, 1, -2, 3
random.data$y <- with(random.data, 0 + 1 * x1 - 2 * x2 + 0 * x3)

# Step 3: Add some random noise to the DV
random.data$y <- random.data$y + rnorm(500, mean = 0, sd = 1)

# Step 4: Convert to probability
random.data$y.prob <- with(random.data, 1 / (1 + exp(-y)))

# Step 5: Create binary response
random.data$y.bin <- round(random.data$y.prob, 0)

# Step 4: Run the model then print the result
result <- glm(y.bin ~ x1 + x2 + x3,
  data = random.data, family = "binomial")

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

result

##
## Call: glm(formula = y.bin ~ x1 + x2 + x3, family = "binomial", data = random.data)
##
## Coefficients:
## (Intercept)      x1      x2      x3
##   -2.2764    4.4889   -8.1380   -0.6464
##
## Degrees of Freedom: 499 Total (i.e. Null); 496 Residual
## Null Deviance:      155
## Residual Deviance: 13.06 AIC: 21.06
```

```
# Logit
logit.fun <- function(x) {1 / (1 + exp(-x))}

curve(logit.fun,
  from = -3,
  to = 3,
  lwd = 2,
  main = "Logit",
  ylab = "p(y = 1)",
  xlab = expression("b_{0} + b_{1}x_{1} + b_{2}x_{2} + ... b_{n}x_{n}")
)

abline(h = .5, lty = 2)
abline(v = 0, lty = 1)
```

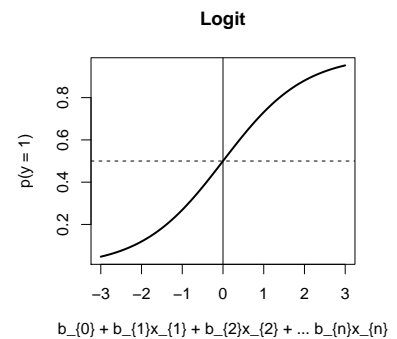


Figure 45: The logit function used in binary logistic regression

```
summary(result)

##
## Call:
## glm(formula = y.bin ~ x1 + x2 + x3, family = "binomial", data = random.data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.59351  -0.00003   0.00000   0.00000   2.09097
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -2.2764      1.6807  -1.354  0.17560
## x1           4.4889      1.7209   2.609  0.00909 **
## x2          -8.1380      3.0519  -2.667  0.00766 **
## x3          -0.6464      0.5467  -1.182  0.23710
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 155.017  on 499  degrees of freedom
## Residual deviance:  13.059  on 496  degrees of freedom
## AIC: 21.059
##
## Number of Fisher Scoring iterations: 13
```

## ANOVA

Once you've calculated a regression object, you can easily create an ANOVA table based on the regression analysis using the `anova()` function.

### `anova(mod)`

To use the `anova()` function, you apply it to an existing linear model object. For example, let's apply it to the previous model:

```
anova(result)

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Analysis of Deviance Table
##
## Model: binomial, link: logit
##
## Response: y.bin
##
## Terms added sequentially (first to last)
##
##      Df Deviance Resid. Df Resid. Dev
## NULL              499      155.017
## x1      1    18.042      498      136.975
## x2      1   122.251      497      14.724
## x3      1     1.665      496       13.059
```

Looks like we have significant effects for all predictors

## *12: Writing your own functions (Coming Soon!)*

*The basic structure of a function*

*Tips and tricks for complex functions*

*Storing and loading your functions to and from a function file*



## *13: Loops and Simulations (Coming Soon!)*

*When and when not to use loops*

*The list object*

*Simple loops over one index*

*Loops over multiple indices*

*Printing and saving temporary results*

*Parallel computing with snowfall()*





## *14: Bayesian Inference (Coming Soon!)*

*What are Bayesian statistics?*

*Bayesian one and two sample tests*

*Bayesian general linear model*



## *15: Model fitting (Coming Soon!)*

*What is a model?*

*What is a loss function?*

*Minimizing loss functions with optimization routines*

*A worked example: Prospect Theory*



## *16: Writing and sharing your work (Coming Soon!)*

*RMarkdown*

*Shiny*

*Sweave (R and Latex)*



## *Appendix*

	white	blue1	chocolate	gold	darkgreen	darkblue	gray30	gray56	gray82		gray26	gray52	gray78	onyx	darkblue	lightgreen	yellow	pumpkin	red2	teal	turquoise	purple	indigo	brown	gray	thislavender	smoke
	aliceblue	blue2	chocolate	darkgray	seafoam	darkblue	gray31	gray57	gray83		gray27	gray53	gray79	onyx	darkblue	lightgreen	yellow	pumpkin	red3	teal	turquoise	purple	seafoam	gray	thislavender	yellow	
artique	whitewash	blue3	chocolate	pink	slate	brick	gray32	gray58	gray84		gray28	gray54	gray80	hopkin	wine	lightpink	lightcyan	pumpkin	red4	teal	turquoise	purple	seafoam	gray	thislavender	yellow1	
antique	whitewash	blue4	chocolate	darkgray	slate	brick	gray33	gray59	gray85		gray29	gray55	gray81	hopkin	nonchill	lightpink	lightcyan	yellow	seafoam	orange	violet	purple	seafoam	gray	thislavender	yellow2	
antique	whitewash	blue5	coral	larkspur	slate	brick	gray34	gray60	gray86		gray30	gray56	gray82	hopkin	nonchill	lightpink	lightcyan	yellow	malabar	green	violet	purple	seafoam	gray	thislavender	yellow3	
antique	whitewash	brown	coral	hollyhock	slate	brick	gray35	gray61	gray87		gray31	gray57	gray83	hopkin	nonchill	lightpink	lightcyan	green	spring	green	violet	purple	seafoam	gray	tomato	yellow4	
antique	whitewash	brown1	coral	hollyhock	slate	brick	gray36	gray62	gray88		gray32	gray58	gray84	hopkin	nonchill	lightpink	lightcyan	lime	dumuri	green	violet	purple	seashell	snow	zinnia	flowergreen	
aquamarine	brown2	coral	hollyhock	slate	brick	white	gray37	gray63	gray89		gray33	gray59	gray85	indian	nonchill	lightcyan	salmon	green	dumuri	green	violet	purple	seashell	snow1	tomato2		
aquamarine	brown3	coral	hollyhock	slate	brick	gray	gray38	gray64	gray90		gray34	gray60	gray86	dian	lightblue	lightcyan	salmon	green	light	green	green	violet	purple	seashell	snow2	zinnia3	
aquamarine	brown4	flow	hollyhock	slate	brick	gray	gray39	gray65	gray91		gray35	gray61	gray87	dian	lightblue	lightcyan	salmon	green	light	green	green	violet	purple	seashell	snow3	zinnia4	
aquamarine	brown5	coral	hollyhock	slate	brick	mostwhite	gray40	gray66	gray92		gray36	gray62	gray88	dian	lightblue	lightcyan	salmon	green	light	green	green	violet	purple	seashell	snow4	turquoise	
aquamarine	brown6	coral	hollyhock	slate	brick	gold	gray41	gray67	gray93		gray37	gray63	gray89	dian	lightblue	lightcyan	salmon	green	light	green	green	violet	purple	seashell	snow5	turquoise1	
azure	whitewash	blue2	coral	hollyhock	slate	pink	gold1	gray42	gray68	gray94		gray38	gray64	gray90	ivory	lightblue	lightcyan	seafoam	green	light	green	green	violet	purple	seashell	snow6	
azure	whitewash	blue3	coral	hollyhock	slate	pink	gold2	gray43	gray69	gray95		gray39	gray65	gray91	ivory1	lightblue	lightcyan	skyblue	green	light	green	green	violet	purple	seashell	snow7	
azure	whitewash	blue4	coral	hollyhock	slate	pink	gold3	gray44	gray70	gray96		gray40	gray66	gray92	ivory2	lightblue	lightcyan	skyblue	green	light	green	green	violet	purple	seashell	snow8	
azure	whitewash	blue5	cyan	hollyhock	slate	pink	gold4	gray45	gray71	gray97		gray41	gray67	gray93	ivory3	lightblue	lightcyan	skyblue	green	light	green	green	violet	purple	seashell	snow9	
azure	whitewash	blue6	cyan1	hollyhock	slate	pink	denro	gray46	gray72	gray98		gray42	gray68	gray94	ivory4	lightblue	lightcyan	skyblue	green	light	green	green	violet	purple	seashell	snow10	
beige	whitewash	blue7	cyan2	hollyhock	slate	pink	denro	gray47	gray73	gray99		gray43	gray69	gray95	khaki	lightblue	lightcyan	skyblue	green	light	green	green	violet	purple	seashell	snow11	
bisque	whitewash	blue8	cyan3	hollyhock	slate	pink	denro	gray48	gray74	gray100		gray44	gray70	gray96	khaki1	lightblue	lightcyan	slateblue	green	light	green	green	violet	purple	seashell	snow12	
bisque	whitewash	blue9	cyan4	hollyhock	slate	pink	denro	gray49	gray75	green	gray19	gray45	gray71	gray97	khaki2	lightblue	lightcyan	goldenrod	green	light	green	green	violet	purple	seashell	snow13	
bisque	whitewash	blue10	cyan5	hollyhock	slate	pink	denro	gray50	gray76	green1	gray20	gray46	gray72	gray98	khaki3	lightblue	lightcyan	goldenrod	green	light	green	green	violet	purple	seashell		

Figure 46: The colors stored in `colors()`.







# Index

`[]`, 41  
`%in%`, 44

`a:b`, 26  
`abline()`, 84  
`aggregate()`, 115

`beanplot()`, 81  
`boxplot()`, 79

`c()`, 24  
`cbind()`, 54  
`chisq.test()`, 133  
`cor.test()`, 132  
correlation, 131  
`curve()`, 87  
`cut()`, 113

`data.frame()`, 56  
`dplyr()`, 119

General Linear Model  
    `anova()`, 139  
    `glm()`, 137

`head()`, 57  
`hist()`, 78

`legend()`, 88  
`length()`, 25  
license, 2  
Linear Model, 135  
`lm()`, 135

`matrix()`, 55  
`merge()`, 121

`paste()`, 86  
`plot()`, 76  
`points()`, 83

`rbind()`, 54  
`readltab.e()`, 59  
`rep()`, 27  
`rgb()`, 96  
`rnorm()`, 32  
`runif()`, 33

Sammy Davis Jr., 75  
`sample()`, 34  
`seq()`, 27  
`subset()`, 70

t-test, 126  
    `t.test()`, 126, 129  
`text()`, 85

`View()`, 58

`with()`, 67