

Tools for Working with Data - 211

Nicole Sorhagen, Ph.D.

2020-09-01

Contents

1 About this book	5
2 Set up project on Rstudio Cloud	7
3 Introduction	9
3.1 Layout of Rstudio	9
3.2 Importing data into Rstudio cloud	18
4 Packages	25
4.1 Installing packages	25
4.2 Loading Packages	29
4.3 Misc	30
5 Picturing Data	33
5.1 Histograms	33
5.2 Scatterplots	38
5.3 Additional resources	42
6 Descriptive Statistics	43
6.1 Descriptive statistics using Tidyverse and Psych packages.	43
6.2 Descriptive statistics using base R	46

Chapter 1

About this book

This book describes how to use R as a tool to work with data.

R statistics is becoming increasingly popular for data management and analysis due to its accessibility and versatility. For example, R can produce records of data analyses, which is consistent with the growing move towards reproducible and open science within the field of psychology. R statistics is also known for making elegant graphs, which can help develop data visualization skills. Because it is open-sourced it is extremely flexible - people create and share packages that make certain aspects of data analysis easy.

R is a programming language. Although learning a programming language can seem a bit intimidating, there are many benefits to trying to figure it out. Mastering the basics of R could be useful for your future coursework, as well as for data management and analysis needs outside the classroom (independent research, future employment, etc.). That is to say, learning the basics of a programming language is a highly transferable skill.

The R programming language can be used within the R software as well as other programs. RStudio is a IDE (integrated development environment) and was designed to make the use of the R programming language more user friendly.

R, and its companion program RStudio, are free and available in PC, Mac, and Linux versions, so students can have it on their own computer - eliminating the need to visit computer labs or to buy student versions of expensive software. R can be downloaded from the CRAN (Comprehensive R Archive Network) (<https://www.r-project.org/>). Rstudio can be found here: <https://rstudio.com/>.

While you are welcome to download R and Rstudio on your personal computer, you do not have to for this course. We will be using Rstudio on a website called Rstudio cloud for class work (this is discussed in more detail in the next chapter). So I am not going to go into detail on downloading the programs on

to your computer here. Please email me if you are interested in this and are having a hard time figuring it out.

Finally, please note that I will be updating this book over the course of the semester.

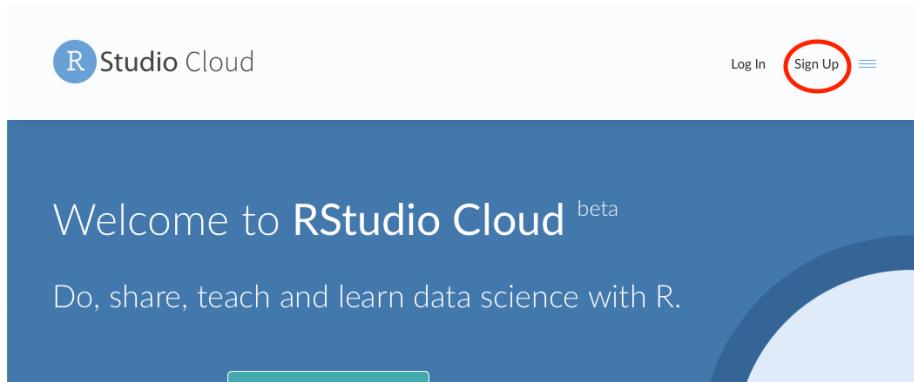
This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License.

Chapter 2

Set up project on Rstudio Cloud

We will use Rstudio cloud on this website: <https://rstudio.cloud>.

You must first make an Rstudio account by clicking the sign up button in the top right corner. (this is free)

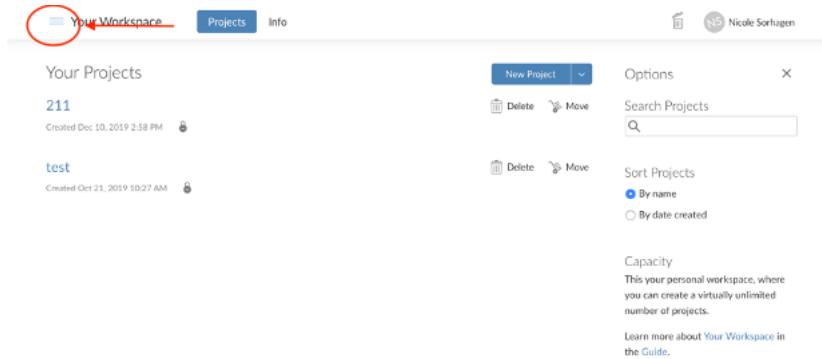


Then join our shared RStudio cloud workspace with the link that I sent you in the email titled 'Rstudio cloud shared workspace'.

You MUST join our shared workspace. I will be checking your work through this shared RStudio cloud workspace. Within this shared workspace, I will be able to see everyone's project, but you will only be able to see your project and my project.

Once you are in your Rstudio Cloud account...

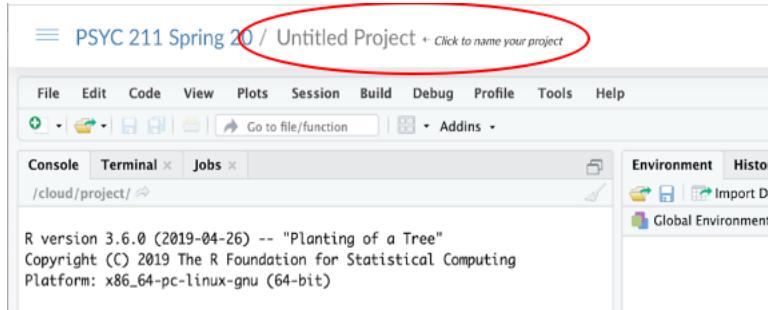
Expand the R studio cloud options by clicking on the 3 lines in the top left corner.



Then select our course (which will be titled the name of course). If you cannot see this option – then you have not been added to our shared workspace.

Once you are in the Class's shared workspace, open a new project. The new project button is on the top right.

Call this project your last name by clicking on the box that says 'Untitled Project' and typing your last name.



Congrats! You have successfully set up your project on our class RStudio Cloud workplace! You will use this project for all of your R work for the rest of the semester. Next let's learn a bit about the software environment and the R programming language.

Chapter 3

Introduction

This chapter introduces the RStudio cloud environment and describes how to import data into the RStudio cloud.

R cannot handle typos and is case sensitive ('Gender' is not the same as 'gender'). If your code will not run check for typos and caps.

Related to this point, do not be afraid to copy and paste with using R. I often copy and paste code and replace variable or dataset names as needed. (This is one of the few times in education where copy and paste is OK!)

3.1 Layout of Rstudio

Rstudio has four panes: the console panel, the script panel, the environment and history panel, and the files and plots panel. Each will be describe in turn next.

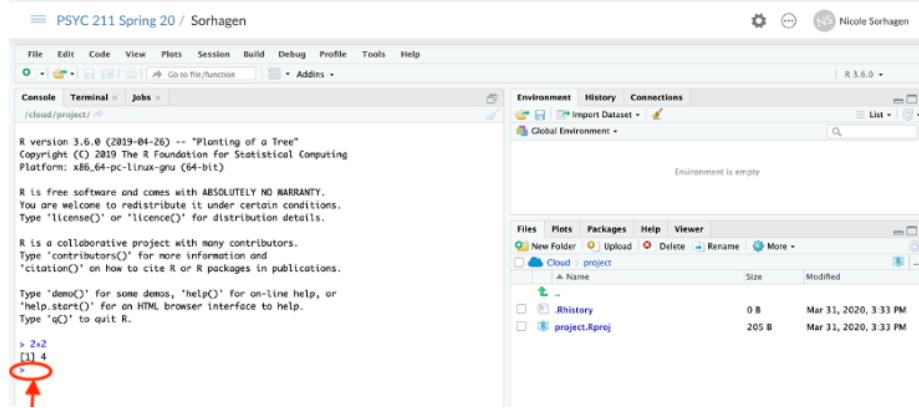
3.1.1 Console

The console panel of R studio is where you can type commands (or R code) and where you will see the output of commands.

In its most basic form, you can think of R as a fancy calculator.

For example: In the console type `2+2` and then press RETURN on your keyboard. The answer '`4`' will appear on the next line.

The `>` in the last line of the console means that the console is ready for a command (see red circle in the picture above).



```
R version 3.6.0 (2019-04-26) -- "Planting of a Tree"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

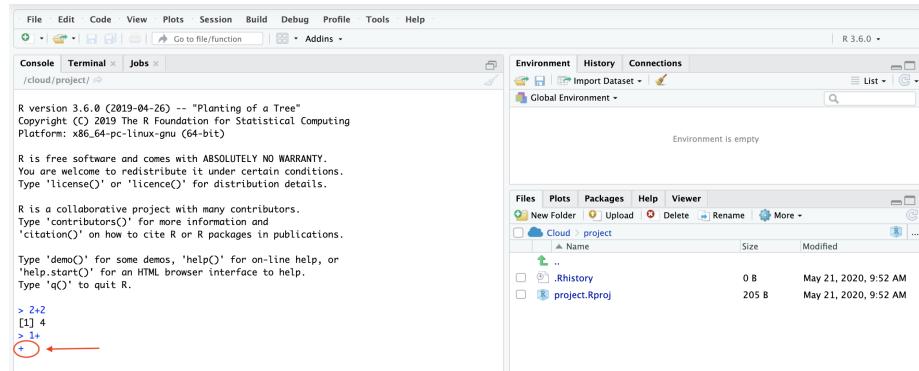
Type 'demo()' for some demos, 'help()' for on-line help,
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> 2+2
[1] 4
```

If **>** is missing from the last line, it means that R is waiting for you to complete a command.

For example, type **1+** in the console and then hit enter.

The plus sign means the command is incomplete.



```
R version 3.6.0 (2019-04-26) -- "Planting of a Tree"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help,
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> 2+2
[1] 4
> 1+
+-->
```

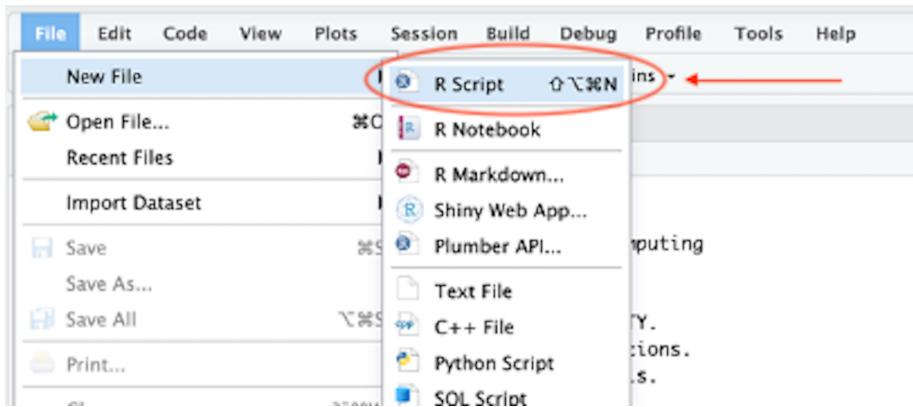
Push the **ESC** button on your keyboard to get back to the command prompt (**>**).

3.1.2 Script

One of the benefits of using R is that you can save a record of your work using scripts. Records of your work allow you to easily start and stop an assignment or research project. You can pick up where you left off whether it is 20 minutes later or 2 years later. These script files also let you share your data analysis with others – from professors, to collaborators, to peer reviewers (ensure that your work is reproducible).

To create a new script, go to the top bar menu:

FILE -> NEW FILE -> R SCRIPT



A new script will open in the top left of the RStudio platform.

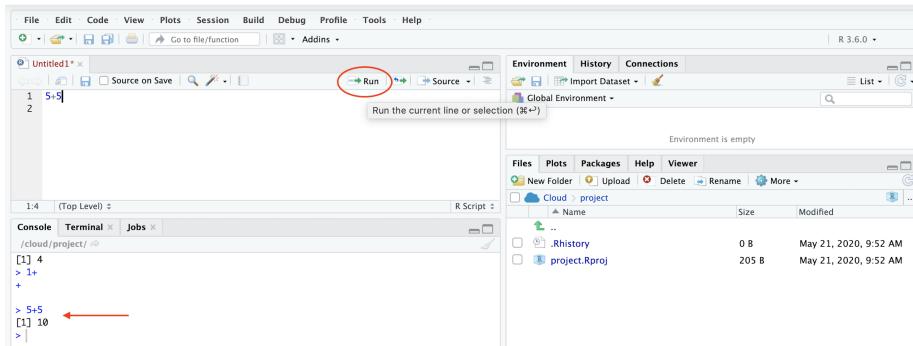
You should run code from scripts

Scripts are similar to running command in the console (this is what you did in the last section).

For example, type `5+5` in the script panel.

In order to run command in a script you should click the run button while the cursor is in the code or the code is selected. Instead on clicking on the run button, you can also run the code by pressing the COMMAND and RETURN keys on your keyboard at the same time (the ALT and RETURN key on a pc).

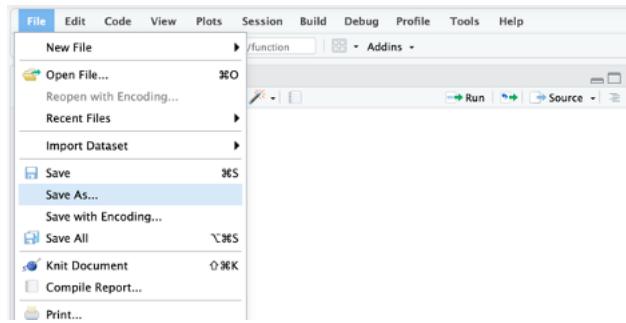
After the code is run, the results will automatically appear the in console (see red arrow in the picture below).



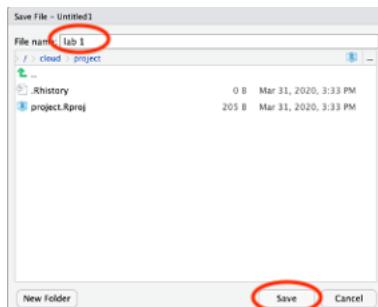
In order to use the script again you must **save** it.

From the drop-down menu select:

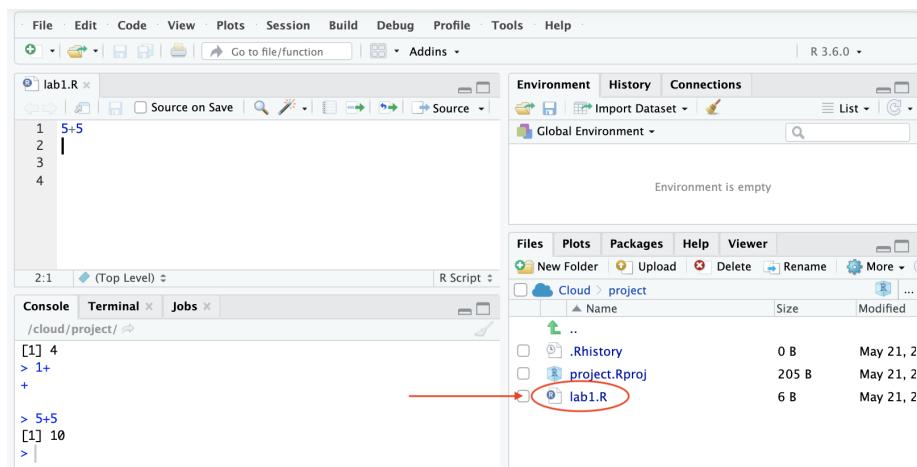
FILE -> SAVE AS



Type **lab 1** into the file name box. And then click the **SAVE** button.



Your file should now be listed in the files window in the bottom right.



This script file is a record of your work and is how you will be graded for the ‘reproducing the assigned chapters’ assignment. Make sure you saved this file and complete the rest of this lab in your ‘lab1’ script.

Within a script you should include comments to yourself and others using a **#**. Anything with a **#** in front of it will not run. **These comments and**

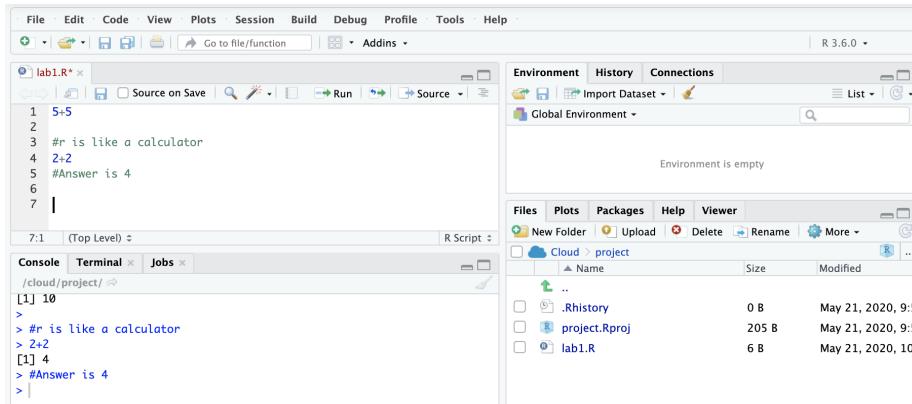
explanations are an important part of an R script.

For example, type the following in to the script and then run it.

```
#r is like a calculator
```

```
2+2
```

```
#Answer is 4
```



Note that the comments appear green in the script.

You can make notes to me this way (if needed). For example:

```
#This will not run. Dr S, I need help
```

I will leave comments to you this way. When I do, I will start my comment with my initials. For example:

```
#NS: try loading the tidyverse package first
```

Scripts should be organized so that someone else can follow what you did.

3.1.3 Environment and history

In the top right corner of RStudio is the environment and history window. The **history tab** shows every line of code that has been run in the current session. I do not use this much.

The **environment tab** is where all active **objects** are listed. An object is something can hold information for later use. The information can be data, values, output, or functions. Normally my panel is set on this.

Objects are assigned using `<-`. Values on the right side of `<-` will be assigned to the object on the left side.

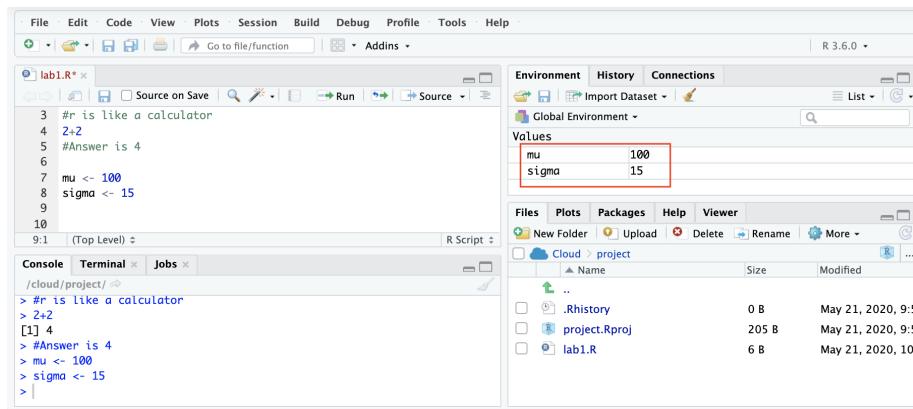
For example, let's tell R that the population mean of IQ scores is 100 and the population standard deviation is 15.

To do this use the following code:

```
mu <- 100
```

```
sigma <- 15
```

After you run these commands, the objects will now be listed in the environment panel in the top left.



The shortcut for making `<-` is the ALT and `-` key together. (or OPTION and `-` on a mac)

3.1.3.1 Vectors

It is possible to store more than one number in an object. One way to do this is to use a **vector**. Assign a set of numbers a vector with the **combine** function: `c()`, by typing all of the numbers you want to store within the parentheses in a comma separated list.

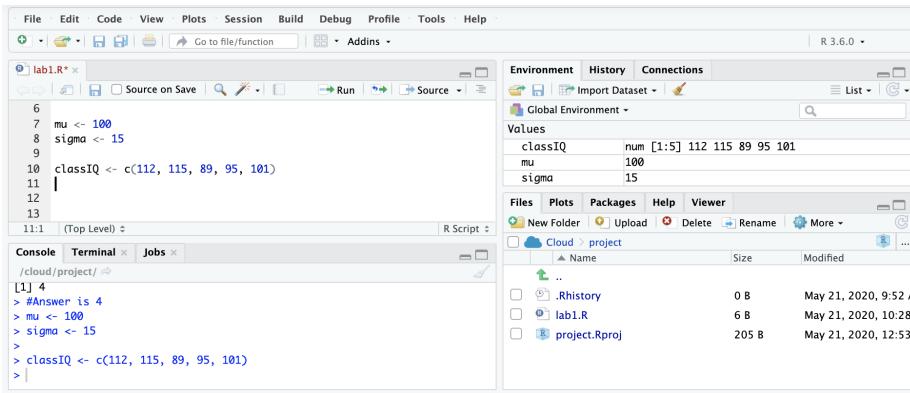
For example, let's enter the IQ scores of students in a small class.

To do this use the following code:

```
classIQ <- c(112, 115, 89, 95, 101)
```

After you run this code, the `classIQ` vector should appear in the environment.

Here is a picture of what your screen should look like:



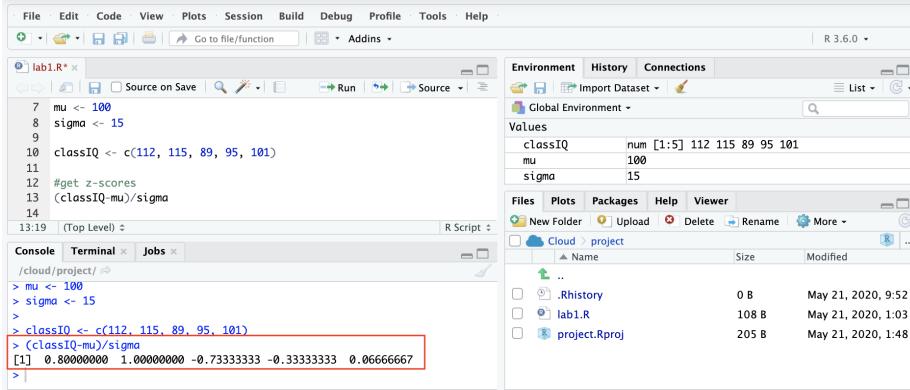
Calculations with vectors apply to all data points.

For example, let's calculate the z-scores for each of the IQ scores.

To do this use the following code:

```
#get z-scores
(classIQ-mu)/sigma
```

The results will appear in the console (See the red box in the picture below)



It is possible to save these answers as a vector using the `<-` function.

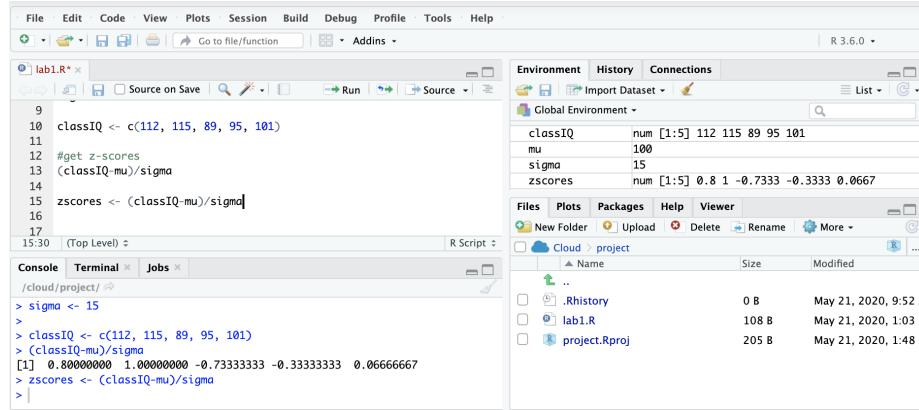
For example, let's save those zscores in a vector called `zscores`.

To do this use the following code:

```
zscores <- (classIQ-mu)/sigma
```

There should now be a vector in the environment called `zscores`.

Here is a picture so that you can check your progress:



3.1.3.2 Data frames

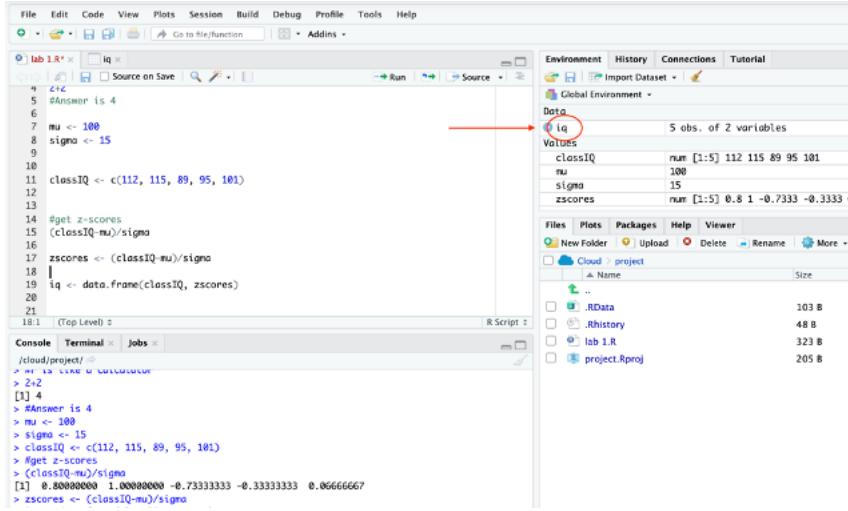
Right now the IQ scores and the z-scores are in separate objects. Variables often need to be in a single object in order to do some basic analyses. You can combine the classIQ and the zscores variable using the **data.frame** command.

To do this use the following code:

```
iq <- data.frame(classIQ, zscores)
```

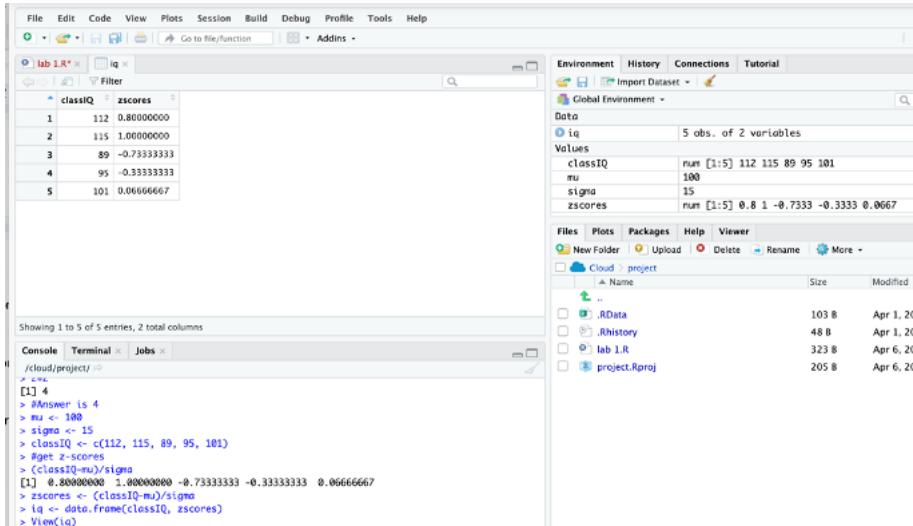
- This command takes the form of DatasetName <- data.frame(Variable1, Variable2, etc)
- The dataset name can be anything you want that you have not already used
 - the name must be one word (there cannot be spaces in the name)

This object will be listed under data instead of values in the environment panel.



Double-click on the word ‘iq’ in the environment panel to look at the dataset that you just created (it is circled in red in the picture above).

A new tab will open with a spreadsheet view of the dataset. When you are done viewing the data, you can close it by click on the ‘x’ next to the name iq.



Note that after looking at the dataset this way, the command `view(iq)` appeared in the console. You can look at the dataset with the `view` command as well

When typing the code to create the data frame, you may have noticed that RStudio uses **predictive text**. This means that RStudio will suggest functions and objects as you type. You should take advantage of this feature!

The screenshot shows the RStudio Cloud interface. On the left, the code editor displays R script code for calculating z-scores. In the center, the Environment pane shows the global environment with variables like classIQ, rm, sigma, and zscores. On the right, the Files pane shows a project directory with files like .RData, .history, lab.R, and project.Rproj.

```

File Edit Code View Plots Session Build Debug Profile Tools Help
lab 1.R * Source on Save Go to file/function Run Source Addins
12
13
14 #get z-scores
15 (classIQ-mu)/sigma
16
17 zscores <- (classIQ-mu)/sigma
18
19 iq <- data[,
20   data.frame (utils) data.frame(..., row.names = NULL, check.rows = FALSE,
21   check.names = TRUE, fix.empty.names = TRUE)
22   data.class (base) check.names = TRUE, fix.empty.names = TRUE)
23   data.entry (utils) setNames = TRUE, check.names = TRUE)
24   data.frames The function data.frames() creates data frames, tightly coupled
25   data.matrix (base) collections of variables which share many of the properties of
26   datasets (utils) matrices and of lists, used as the fundamental data structure by
27   most of R's modeling software.
28   Press F1 for additional help
28 iq <- data.frame(classIQ, zscores)
19:11 [Top Level] R Script
Console Terminal Jobs
~/cloud/project
> rm -r .history
> 2>
[1] 4
> Answer <- 4
> x <- 100
> signs <- 15
> classIQ <- c(112, 115, 89, 95, 181)
> #get z-scores
> (classIQ-mu)/sigma
[1] 0.80000000 1.00000000 -0.73333333 -0.33333333 0.06666666
> zscores <- (classIQ-mu)/sigma
[1] 4

```

3.2 Importing data into Rstudio cloud

In the last section you learned how to assign data to a **vector** using the **combine** function.

Another way to assign data to an object is by first entering the data into a spreadsheet (like google sheets or excel) and then import the data into RStudio. This will be our preferred method.

First download the exam2.csv file from d2l.

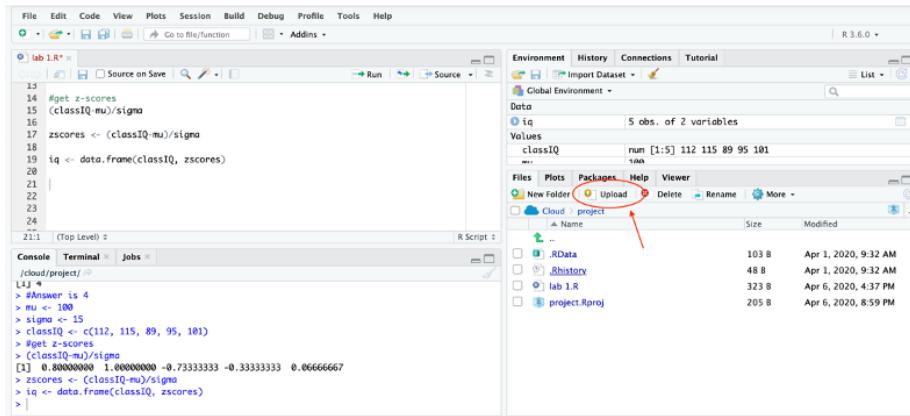
This file contains data on the exam 2 scores from a previous class of mine. It also has some made up data on students' studying, eating, and drinking behaviors before the exam.

Make sure the file remains a .cvs file after you download it. If you open it in excel, excel might change the file to an .xls file.

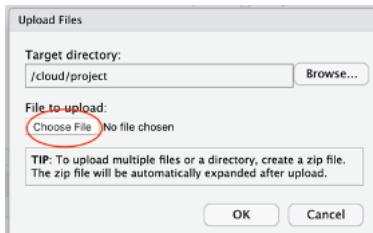
When using RStudio in the cloud, you will first need to upload the data into your project and then import it into the global environment.

3.2.1 Upload the data into Rstudio Cloud

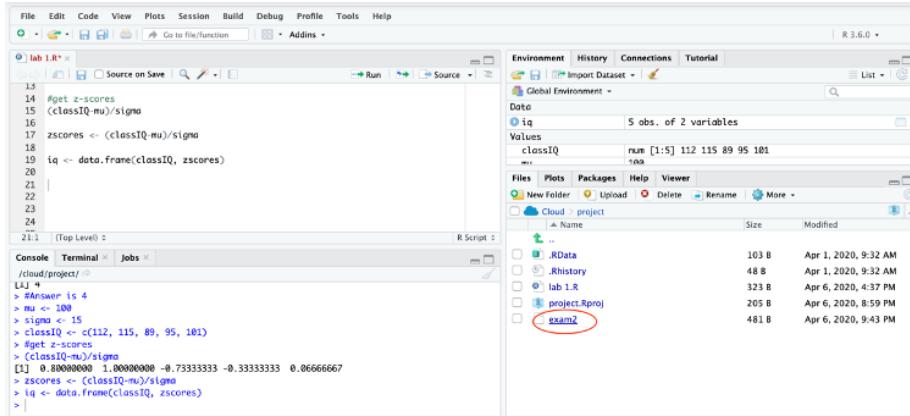
Then select the UPLOAD button in the files window.



In the Upload Files window, click the CHOOSE FILE button and then navigate to the exam2.csv file on your computer. Then click the OK button.



The data file should now be listed in the files section of RStudio.

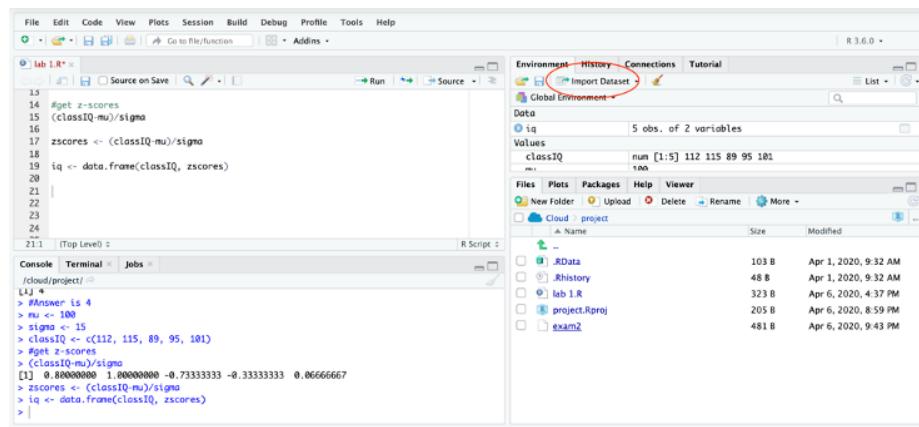


3.2.2 Import data

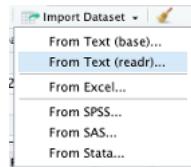
Then you need to import the data into the environment (i.e. assign the data to an object). This can be done through using point and click options or with code.

3.2.2.1 Point and click

First click on the **IMPORT DATASET** button in the environment panel.



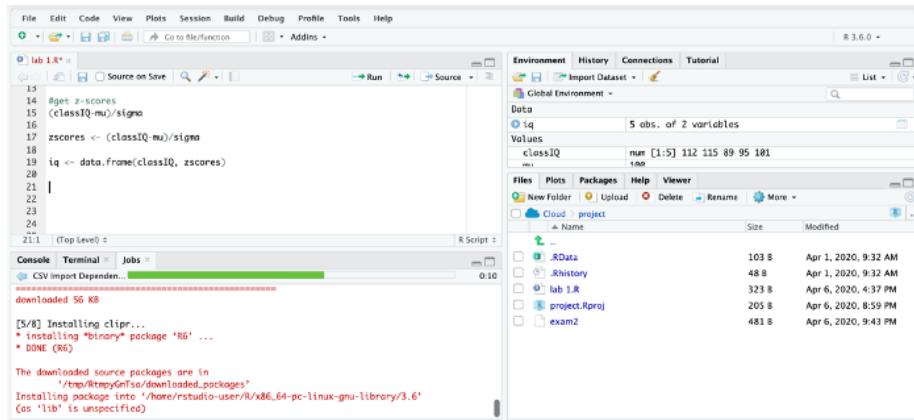
Then select the '**FROM TEXT (READR)**'



The first time you select this – the following window will appear asking if you would like to install the `readr` package. Select YES. I will introduce packages in the next section.

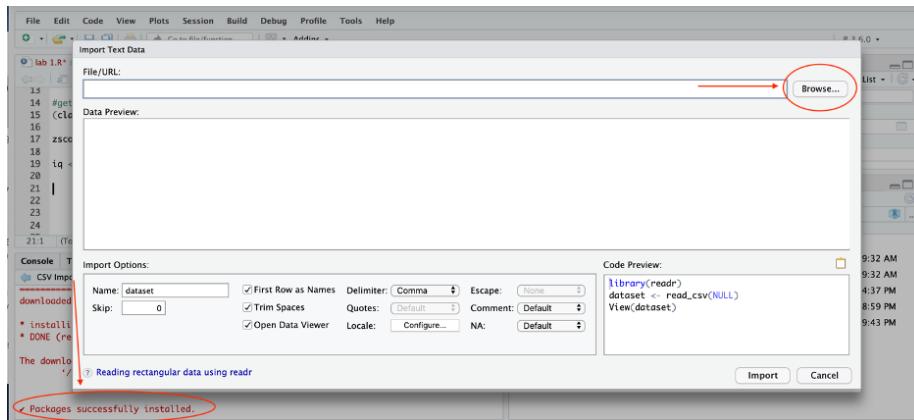


After you select yes, R will begin downloading the package. This can take a few minutes and will look something like this:



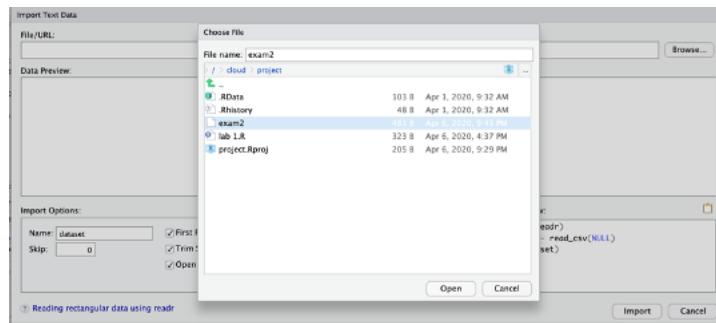
It is important to be *patient* here and let the package download completely before you move on to the next step.

When the download is complete, your screen should look like this:



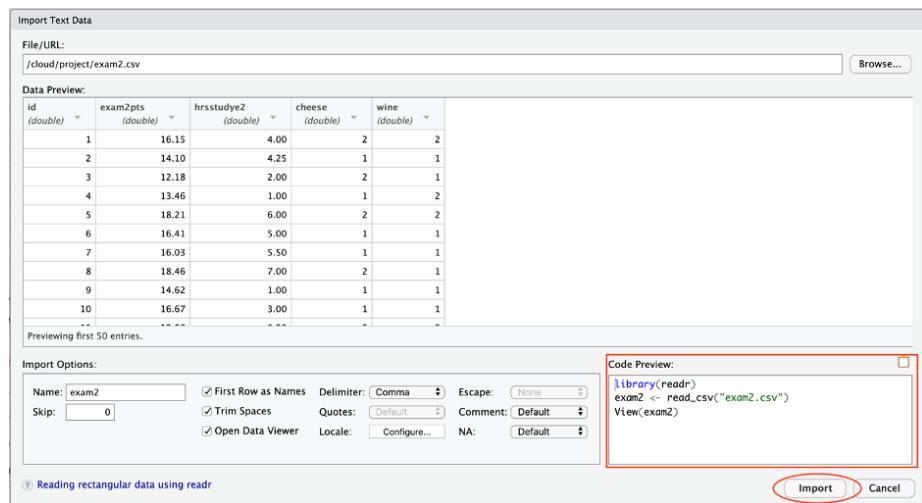
Note that you can see that the package was successfully installed in the console in the bottom left. The next time you use `readr` to import data – you will not have to download the package first.

Next select the BROWSE button in the top left corner of the import data window.



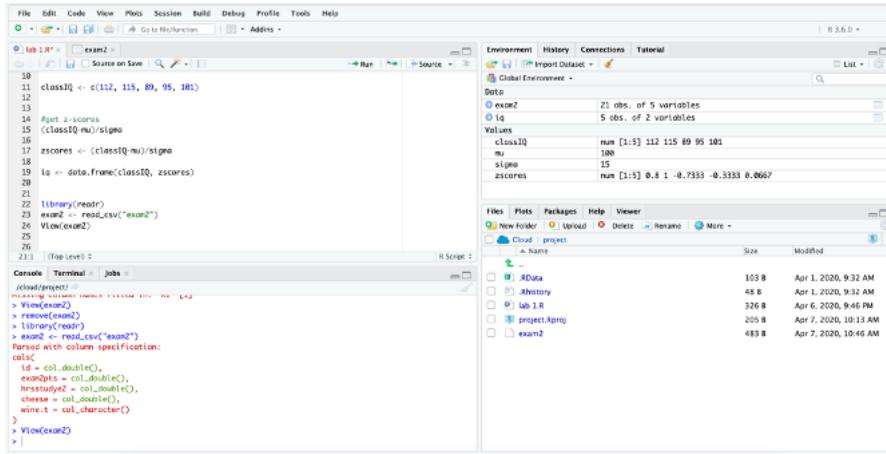
In the choose file window, select ‘exam2’. And then select OPEN.

The next window should look like this:



From here you should click the IMPORT button.

But first note the **Code Preview** box. This is the code you could use to import data (instead of clicking through all these windows). It's a good idea to copy this code before you click the import button and then paste it into your script for your records and in case you need to assign the file to an object again (because it is faster with code).



3.2.2.2 Code

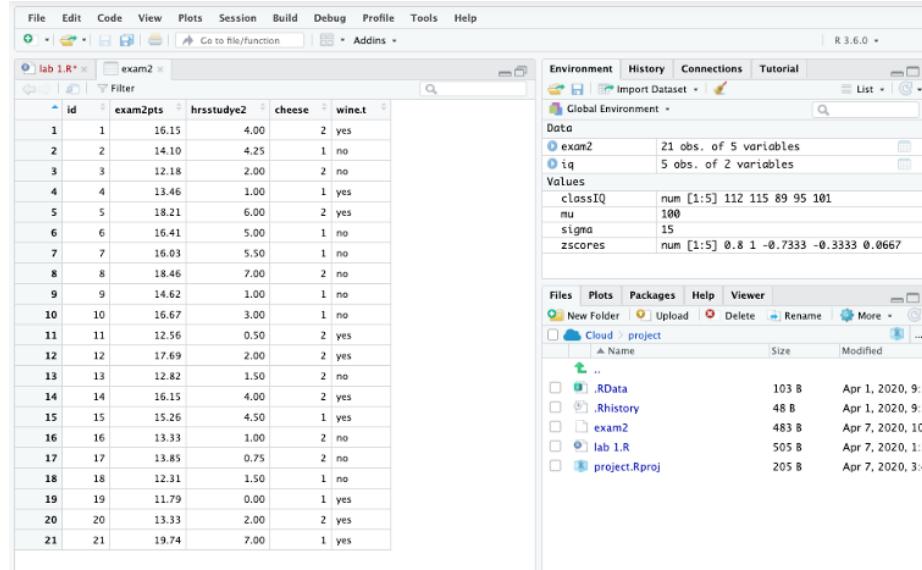
Alternatively you could have typed the code that you copy and pasted (and not gone through all of the point and click windows).

```
library(readr)
exam2 <- read_csv("exam2.csv")
```

- This command takes the form of DatasetName <- read_csv("FILENAME.csv")
- The dataset name can be anything that you have not already used
 - the name must be one word (there cannot be spaces in the name)
- If you have not installed the readr package, you will have to do so first (see the packages chapter for more information)

3.2.2.3 View data

Double click on the word exam2 in the environment panel to look at the dataset.



Each column is a different variable. Each row is a different participant (in this example a student).

- The first column is an arbitrary student ID number – so that the students' identity is protected.
- The second column is exam points earned by the students out of 20 (this is real data from a Fall 2019 class).
- The third column is the number of hours the students studied for the exam (this is made up data).
- The fourth column is whether or not students ate cheese the night before the exam (1 = no; 2 = yes... also made up data).
- The last column is data on whether or not students drank wine the night before the exam (1 = no; 2 = yes... also made up data).

Chapter 4

Packages

NOTE: Please open a new script and call it lab 2 (or week 2) for the replication of this chapter and the picturing data chapter assignment.

Base R refers to the functions that automatically come with R. But many people build on top of Base R to make R better. The way they do this is through **packages**, which contain new R functions. There are thousands of packages available that can do fancy things like quickly compute descriptive statistics and create APA style tables (and much much more).

The first time you use a package, you need to install it. Once a package is installed, you will need to tell R that you want to use it by loading it. You will need to load any packages you want to use each time you open the R program. (I am not exactly sure how this works in the RStudio cloud because it does not seem to shut down when you close out of the RStudio cloud website. See the Restarting R section below for a work around.)

That is, you only have to install a package once. You will have to load a package every time you want to use it.

4.1 Installing packages

The first time you use a package, you need to install it. We actually did this once already while importing data! This time let's learn more about the process.

In RStudio, packages can be installed through point and click (GUI) or with code.

4.1.1 Installing packages using point and click (GUI)

Let's first install a package called **Tidyverse**. Tidyverse was created by Hadley Wickham and his team with the aim of making various aspects of data analysis in R easier. It is actually collection of packages that include a lot of functions (e.g., subsetting, transforming, visualizing) that many people think of as essential for data analysis. (See the tidyverse website for additional information: <https://www.tidyverse.org>).

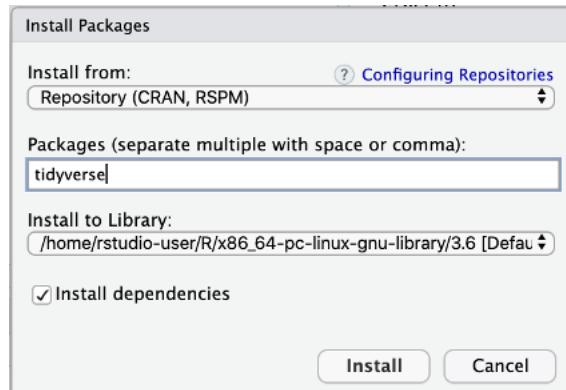
To install a package with GUI go to the top bar menu:

TOOLS -> INSTALL PACKAGES



In the install packages window, type the name of the package you would like to install. For example, type `tidyverse` in the packages box.

Then click INSTALL.



Again, installing a package can be a little slow on the RStudio cloud. Please be patient (maybe this is a good time to stretch your legs, refill your beverage, let the dog out, etc.)

Your screen should look like this when it is starting to install:

The screenshot shows the RStudio interface with the following details:

- Console:** Displays R code being run, including `library(readr)` and `install.packages("tidyverse")`.
- Environment:** Shows the Global Environment pane with variables `exam2`, `iq`, and `zscores`. It also shows the `Values` pane for these variables.
- Packages:** Shows the Packages tab in the environment pane, indicating that `tidyverse` is being installed.
- File Explorer:** Shows a project named "project" with files like `RData`, `Rhistory`, `lab 1.R`, and `exam2`.

It should look like this when it is in the process of installing:

The screenshot shows the RStudio interface with the following details:

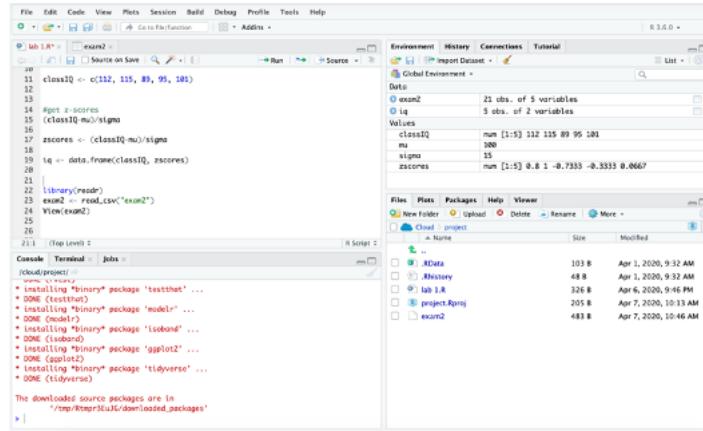
- Console:** Displays R code being run, including `library(readr)` and `install.packages("broom")`.
- Environment:** Shows the Global Environment pane with variables `exam2`, `iq`, and `zscores`.
- Packages:** Shows the Packages tab in the environment pane, indicating that `broom` is being installed.
- File Explorer:** Shows a project named "project" with files like `RData`, `Rhistory`, `lab 1.R`, and `exam2`.
- Output:** Shows the progress of the package download and extraction in the console.

```

10 library(readr)
11 exam2 <- read_csv("exam2")
12 View(exam2)
13
14 #get z-scores
15 (classIQ-mu)/sigma
16
17 zscores <- (classIQ-mu)/sigma
18
19 iq <- data.frame(classIQ, zscores)
20
21
22 library(readr)
23 exam2 <- read_csv("exam2")
24 View(exam2)
25
26
21.1 (Top Level) R Script 2
Console Terminal Jobs
/coud/project
downloaded 2.2 MB
trying URL 'http://package-proxy/src/contrib/broom_0.5.5.tar.gz'
Content type 'application/x-tar' length 1965899 bytes (1.9 MB)
downloaded 575 KB
trying URL 'http://package-proxy/src/contrib/dplyr_1.4.2.tar.gz'
Content type 'application/x-tar' length 589414 bytes (575 KB)
downloaded 575 KB
trying URL 'http://package-proxy/src/contrib/dplyr_0.8.5.tar.gz'

```

And then this when the installation is complete:



Do not proceed until the console says the package has been installed.

4.1.2 Installing packages using code

You can also install a package using this code:

```
install.packages()
```

To install tidyverse, for example, you would use this code:

```
install.packages("tidyverse")
```

- Note that the word tidyverse is in quotes

But do not run this code – as you have already installed it with GUI.

Instead, let's install a package called psych using the `install.packages` command. The **psych package** is a package for personality, psychometric, and psychological research. It has been developed at Northwestern University (maintained by William Revelle) to include useful functions for personality and psychological research.

To install this package, use following command:

```
install.packages("psych")
```

Your screen should look like this when the package is completely installed:

```

File Edit Code View Plots Session Build Debug Profile Tools Help
lab 1.R * exam2.R
Source on Save Run Source
14 #get z-scores
15 (classIQ-mu)/sigma
16
17 zscores <- (classIQ-mu)/sigma
18 iq <- data.frame(classIQ, zscores)
19
20 library(readr)
21 exam2 <- read_csv("exam2")
22 View(exam2)
23
24 library(tidyverse)
25
26 install.packages("psych")
27
28 library(psych)
29
30
31
32
33
34
28.1 (Top Level) R Script
Console Terminal Jobs
/ccloud/project/
downloaded 3.6 MB
* installing *binary* package 'moment' ...
* DONE (moment)
* installing *binary* package 'psych' ...
* DONE (psych)

The downloaded source packages are in
  '/tmp/RtmpQWzul/downloaded_packages'
> |

```

Environment History Connections Tutorial
Import Dataset List
Global Environment
Data
exam2 21 obs. of 5 variables
iq 5 obs. of 2 variables
Values
classIQ num [1:5] 112 115 89 95 101
mu 100
sigma 15
zscores num [1:5] 0.8 1 -0.7333 -0.3333 0.0667

Files Plots Packages Help Viewer
New Folder Upload Delete Rename More
Cloud project
Name Size Modified
RData 103 B Apr 1, 2020, 9:32 AM
Rhistory 48 B Apr 1, 2020, 9:32 AM
exam2 483 B Apr 7, 2020, 10:46 AM
lab 1.R 326 B Apr 6, 2020, 9:46 PM
project.Rproj 205 B Apr 7, 2020, 12:32 PM

Remember that installing packages is the first step to using them and they only have to be installed once.

Next let's learn how to load packages, so that you can use thier functions.

4.2 Loading Packages

Installing a package is only the first step.

In order to use a package, it must be loaded first.

Packages can only be loaded with code. Packages need to be loaded every time you open the RStudio program. Most people's R scripts begin with the code that load packages.

When you have the Rstudio program installed on your computer this is straight forward (either the program is open or closed). This is less clear with Rstudio cloud because it does not seem to always shut down when you close the web browser site. (Please see the section on restarting Rstudio in the misc section below for a work around.)

The command to load a package is:

`library()`

For example, load the tidyverse package with this:

`library(tidyverse)`

After you run this code, your screen should look like this:

The screenshot shows the RStudio interface with the following details:

- Environment View:** Shows objects `exam2` (21 obs. of 5 variables), `iq` (5 obs. of 2 variables), and `zscores` (num [1:5] 0.8 1 -0.7333 -0.3333 0.0667).
- Global Environment View:** Shows objects `classIQ` (num [1:5] 112 115 89 95 101), `mu` (100), `sigma` (15), and `zscores` (num [1:5] 0.8 1 -0.7333 -0.3333 0.0667).
- Console View:** Shows the command `library(tidyverse)` highlighted with a red oval.
- Output View:** Shows the results of loading tidyverse, including conflicts with dplyr and purrr.
- File View:** Shows a file tree with files like `RData`, `Rhistory`, `exam2`, `lab 1.R`, and `project.Rproj`.

The console shows that the Tidyverse package has been loaded (don't worry about the conflicts for now).

Next let's load the psych package using this command:

```
library(psych)
```

The screenshot shows the RStudio interface with the following details:

- Environment View:** Shows objects `exam2` (21 obs. of 5 variables), `iq` (5 obs. of 2 variables), and `zscores` (num [1:5] 0.8 1 -0.7333 -0.3333 0.0667).
- Global Environment View:** Shows objects `classIQ` (num [1:5] 112 115 89 95 101), `mu` (100), `sigma` (15), and `zscores` (num [1:5] 0.8 1 -0.7333 -0.3333 0.0667).
- Console View:** Shows the command `library(psych)`.
- Output View:** Shows the output of `install.packages("psych")` and the successful attachment of the `psych` package.
- File View:** Shows a file tree with files like `RData`, `Rhistory`, `exam2`, `lab 1.R`, and `project.Rproj`.

Again,

don't worry about the warning about masked functions for now.

4.3 Misc

You can get additional information using the `help()` function and `? help` operator in R. They both provide access to documentation pages for all functions and packages.

For example, use the following code to get more information about Tidyverse:
`?tidyverse`

Or this command to get more information about Psych:

```
help(psych)
```

4.3.1 Restarting R

Because it is unclear whether Rstudio completely turns off when you close the website, you could restart the R session to simulate the act of closing and reopening the Rstudio program (like you could if it were installed on your computer).

To do this, in the drop down menu go to:

SESSION -> RESTART R



When you restart the R session, everything in the script, environment, console, and files will remain.

All packages that were loaded will be cleared, so you will have to reload them if you want to use them.

If something is not working like it is suppose to (and you have checked for type-os), try restarting the R session. It could be that the functions of one package conflict or mask the functions of another package.

Chapter 5

Picturing Data

“The simple graph has brought more information to the data analyst’s mind than any other device.” — John Tukey

This chapter focuses on how to make graphs and figures in R. Data visualization is useful for descriptive statistics, data analysis, and communicating results.

5.1 Histograms

Here you will learn how to make a histogram. Histograms plot the frequency of each score in a set of data. Thus, they are essentially a graphic of a frequency distribution. They are useful for checking the shape of a distribution (many statistical tests assume data is approximately normally distributed), checking for coding errors, and checking for outliers.

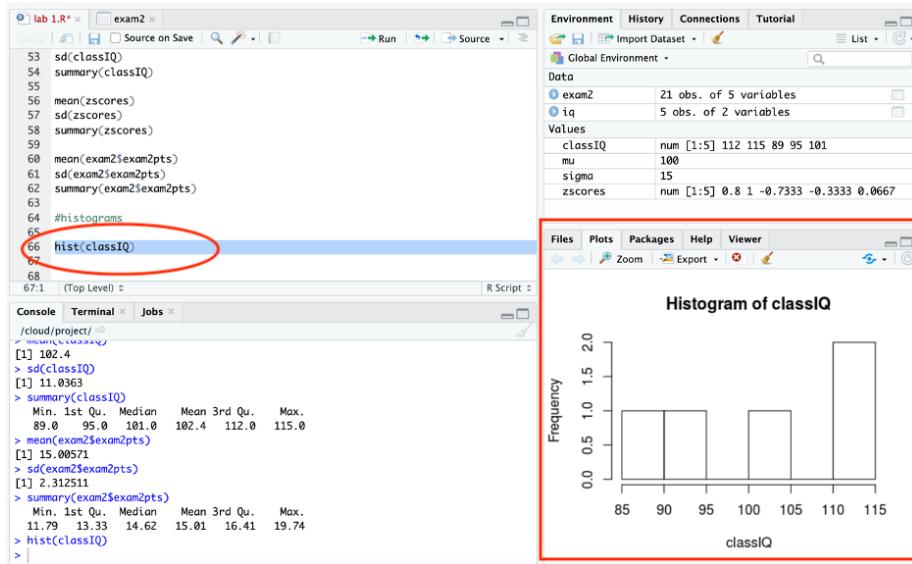
5.1.1 Histograms with base R

Let’s first make histogram with base R by using the `hist()` function.

A **function** in R is any kind of operation. For example, the `hist()` function will create a histogram. An **argument** is what a function acts on.

For example, `hist(classIQ)` will return the histogram of the IQ scores in the `classIQ` vector. This code applies the function `hist` to the variable `classIQ`.

After you run this command, your screen should look similar to this:



I circled and boxed what should match here. (Please excuse a few differences between this screenshot and your screen, like the name of the script, the code in the script before the histogram, and the results in the console. I had presented the material in a different order the last time I taught it.)

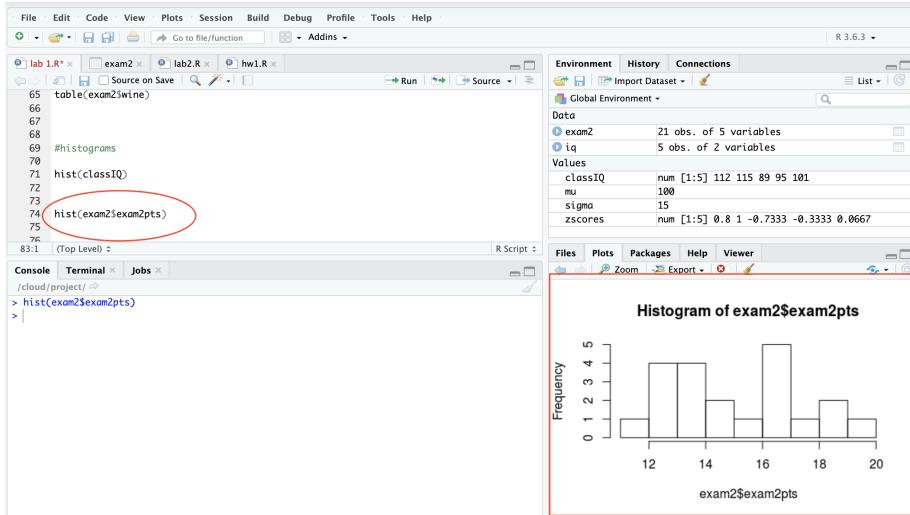
In order to use a base R function with a variable within a data frame you have to tell R to first look in the data frame in order to find the variable. You do this with the dollar sign (\$). Place the \$ between the name of the data frame and the name of the variable.

For example, to use the `hist()` function to create a histogram of the exam 2 points variable in the exam 2 dataset, use this code:

```
hist(exam2$exam2pts)
```

- `exam2$exam2pts` is telling R to first go to the `exam2` dataset and then use the `exam2pts` variable

Here is a picture:



The data looks somewhat normally distributed, with a slight positive skew.

5.1.2 Histograms with tidyverse

The base R option is quick and easy. But it is not customizable. Because of this – many people prefer to use **ggplot** (of the Tidyverse package - so tidyverse needs to be loaded).

Ggplot is typically taught with the analogy of a globe that is built one layer at a time. You start with a world of only ocean (no land). Then you progressively add “layers” of land, colors, terrain, legends, etc. This system is based on the grammar of graphics: statistical graphics map **data** onto perceptible **aesthetic attributes** (e.g., position, color, shape, size, line type) of **geometric objects** (e.g., points, bars, lines). Code can also be added to ggplots to make graphs in APA style.

With ggplot, you build plots step-by-step, layer-by-layer using the following steps:

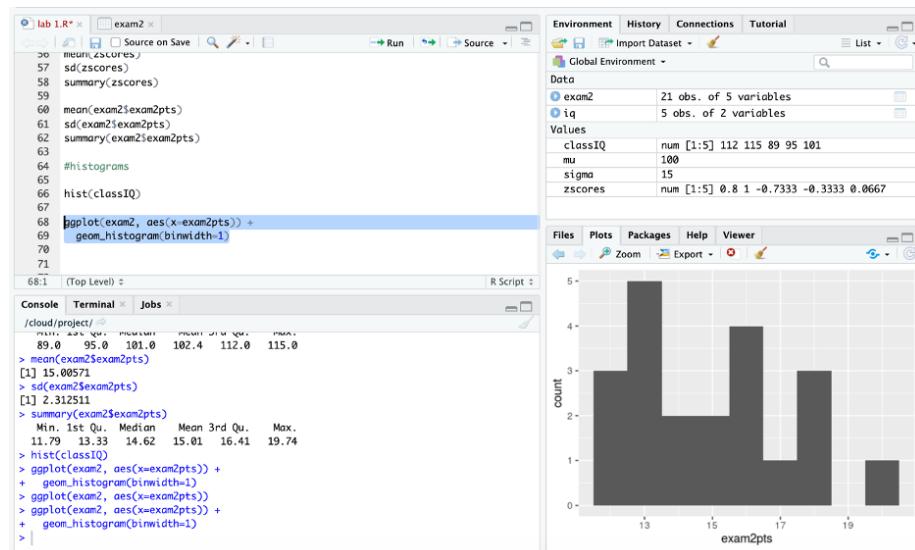
1. Start with **ggplot()**
 2. Supply a dataset and aesthetic mapping, **aes()**
 3. Add on . . .
- + **Layers**, like **geom_point()** or **geom_histogram()**
 - + **Scales**, like **scale_colour_brewer()**
 - + **Faceting Specifications**, like **facet_wrap()**
 - + **Coordinate Systems**, like **coord_flip()**

The code for a histogram of the exam 2 points is:

```
ggplot(exam2, aes(x=exam2pts)) +
  geom_histogram(binwidth=1)
```

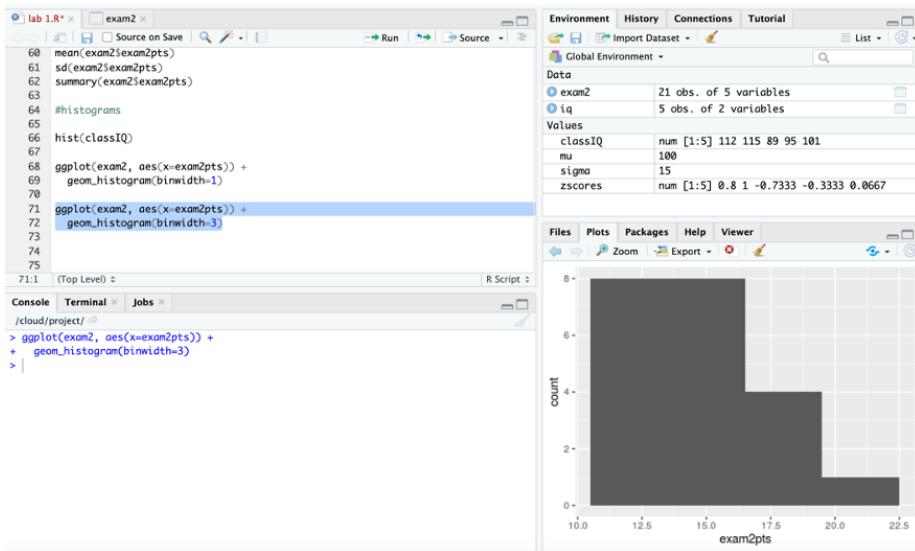
- The first line starts with ggplot and then supplies a dataset and aesthetic mapping, `aes()`
 - Because you supply the dataset this way - you do not need to use \$ to tell R where the variable is
- The second line adds the layer of a histogram

After you run this code your screen should look like this:



Note that the histogram here is more detailed than the one you produced with base R. This is because base R used 5-points bins, while ggplot used 1-point bins (because you told R to). Here it is easier to see that the data is slightly skewed right.

In ggplot, it is easy to change the amount of points per bin by changing the number after the binwidth. For example, here I change the number to 3:



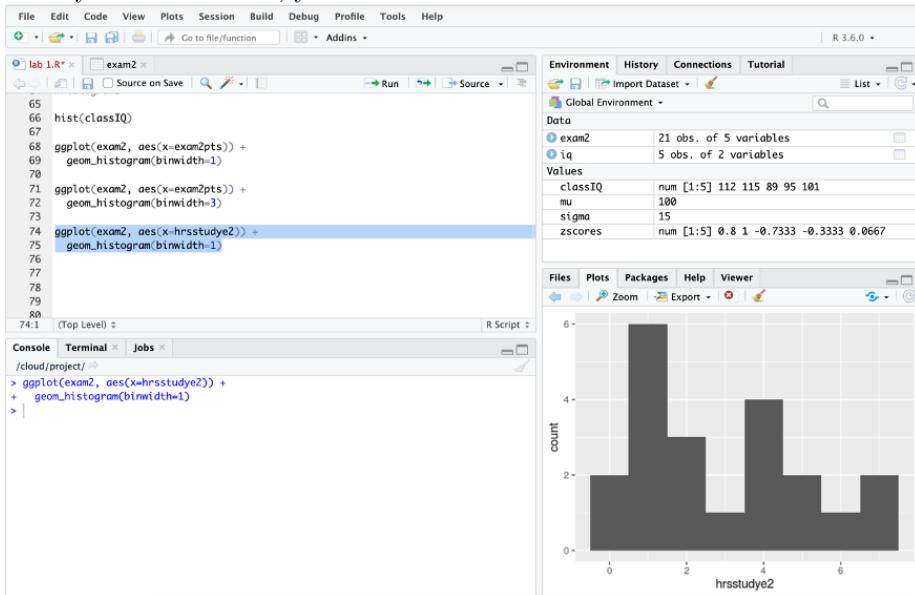
Some say that 10 bins in a histogram is a good rule of thumb (There are more precise equations for determining the “right” number of bins as well).

Let's look at a histogram of the number of hours studied for exam 2 next.

Here is the code:

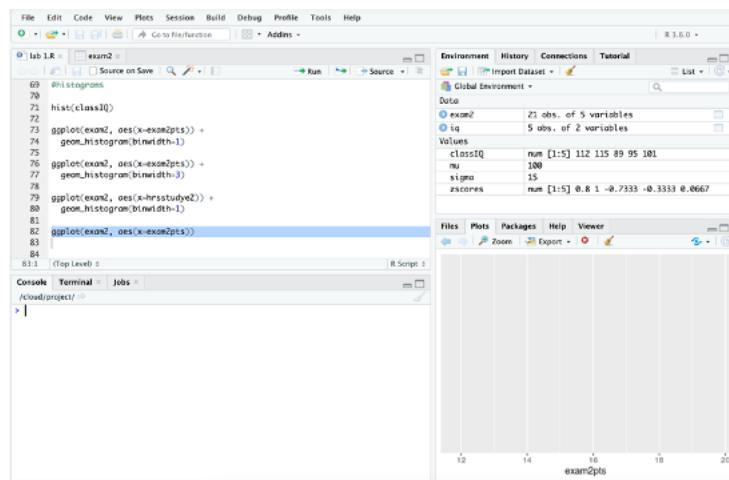
```
ggplot(exam2, aes(x=hrsstudy2)) +
  geom_histogram(binwidth=1)
```

After you run this code, your screen should look like this:



The histogram show that data approximates the normal distribution and is roughly mound shape.

You do not need to run this – but I just want to show you that if I run only the first line of the histogram code, the figure would look like this:



...so this is the world as only ocean – without land. The second line of the ggplot code (i.e. `geom_histogram(binwidth=1)`) adds the “land”.

Please note that I am going to start providing less screen shots of the whole Rstudio window from this point forward. When I include R code know that I mean that the code should be typed into a script.

5.2 Scatterplots

A **Scatterplot** is a graph where one variable is plotted on the y-axis and the other is plotted on the x-axis. Each dot represents one participant, measured on two variables.

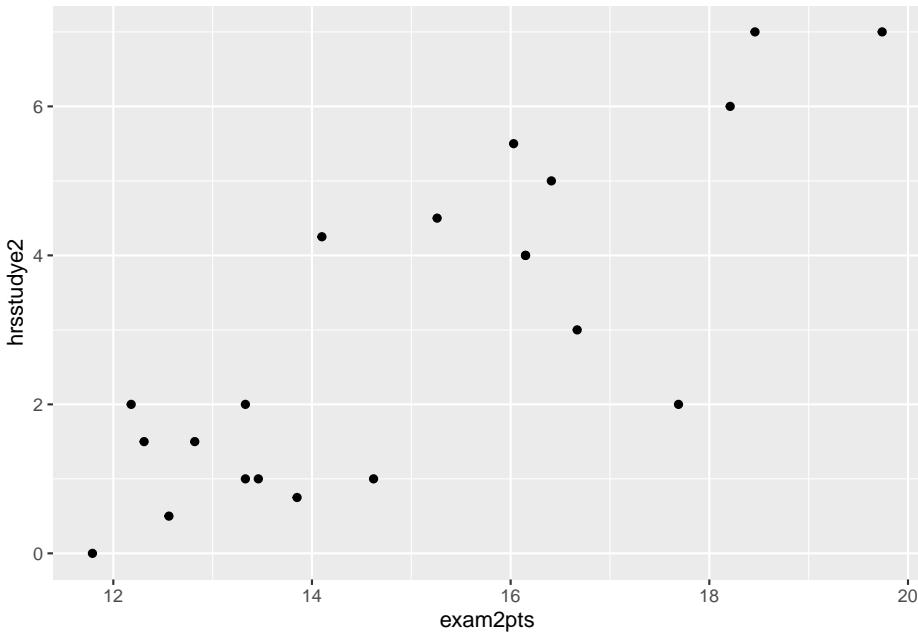
We are going to focus on using ggplots to create scatterplots because it is the more powerful data visualization tool in R.

5.2.1 Two continuous variables

Using the exam 2 dataset, let's say we hypothesized that there is a positive association between exam 2 scores and the number of hours studied for the exam. One of the first steps of exploring this association is to create a scatterplot.

Here is the code and resulting graph:

```
ggplot(exam2, aes(x=exam2pts, y=hrsstudye2)) +
  geom_point()
```



The data points are trending upward, suggesting a positive relation between exam 2 scores and the number of hours. The students who studied longer for the exam received higher grades; While those students who studied for less time received lower grades.

5.2.2 One continuous and one categorical variable

Let's say you were interested in the relation between cheese eating and exam 2 scores. You hypothesized that exam scores will be lower for students who ate cheese the night before the exam because cheese gives nightmares.

Traditionally psychology likes to visualize the relation between a continuous and categorical variable using a bar graph. However, bar graphs can be misleading about the true nature of the data. Because of this, I prefer to continue to use a scatterplot to look at the association between a continuous and categorical variable - with some alterations to show the mean and variability (which is important to show with group data).

In ggplots you can alter the scatterplot to include the mean and variability, in addition to the actual data points, by including the `stat_summary()` function in the ggplot code. The `stat_summary()` function adds statistics to a ggplot.

To use the `stat_summary()` function you need to install the Hmisc package. It does not need to be loaded. *This is a rare exception on how packages in R normally work - The package does not need to be loaded in order for R to use it.*

Here is the code to install Hmisc:

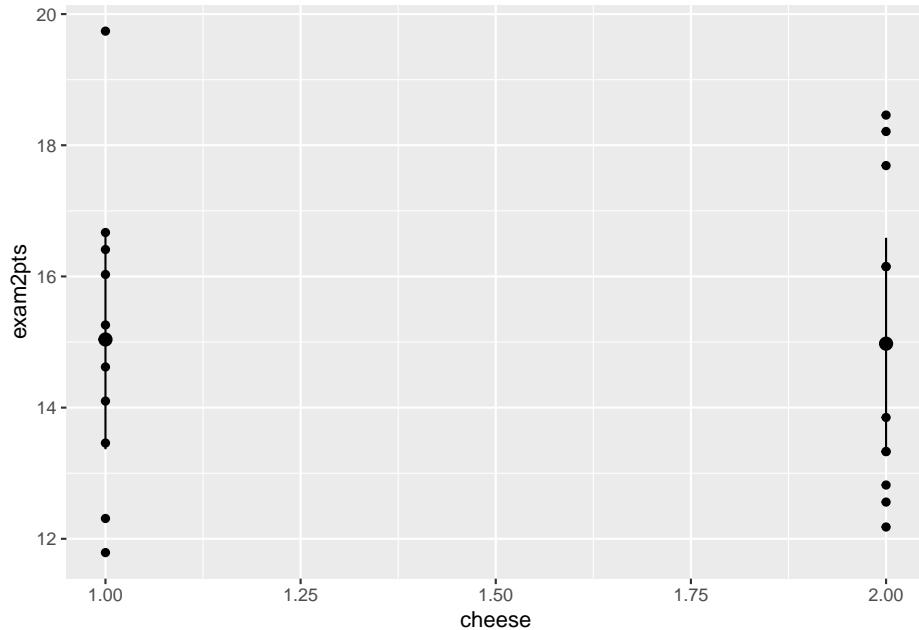
```
install.packages("Hmisc")
- remember to be patient and wait until the package is completely installed.
```

Then create the scatterplot with the following ggplot code:

```
ggplot(exam2, aes(x = cheese, y = exam2pts)) +
  geom_point() +
  stat_summary(fun.data = mean_cl_normal)
```

- x is the categorical variable
- y is the continuous variable
- `geom_point()` includes the data points
- The `fun.data = mean_cl_normal` within the `stat_summary()` function adds the mean and the confidence interval around the mean.

The graph should look like this:



The means are represented by the large dots. The lines represent the 95% confidence intervals, which shows the certainty around the mean and is based on the sample mean, standard deviation, and n.

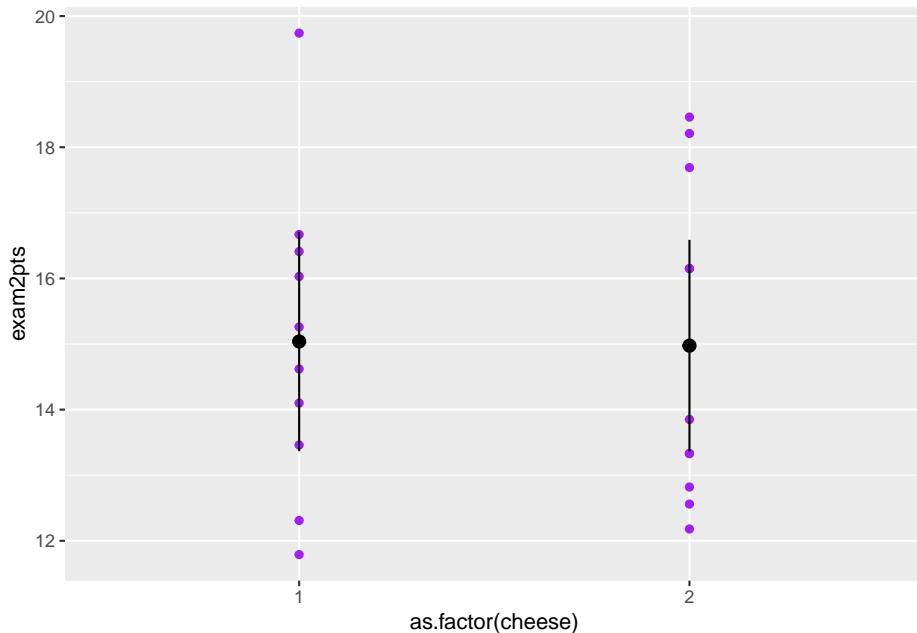
The small dots are the data points representing the participants cheese eating and exam grades.

Here you can see that the mean exam 2 scores are pretty similar for students who did and did not eat cheese the night before the exam. The spread of the scores is also similar. (Remember from the introduction chapter 1 = no and 2 = yes).

I like to make a few alterations to the previous code for aesthetics...

```
ggplot(exam2, aes(x = as.factor(cheese), y = exam2pts)) +
  geom_point(color = "purple") +
  stat_summary(fun.data = mean_cl_normal)
```

- The `color = "purple"` in the `geom_point()` function changes the color of the data points making the graph easier to read
- The `as.factor(cheese)` tells R to treat the cheese variable as a factor, which makes the x-axis more visually appealing



5.3 Additional resources

<https://rstudio.cloud/learn/primers/1.1>

Chapter 6

Descriptive Statistics

NOTE: Please open a new script and save it as lab 3 (or week 3) for your reproduction of this chapter.

This section focuses on functions that find descriptive statistics. **Descriptive statistics** refer to measures of central tendency (mean, median, and mode) and measures of variability (standard deviation, variance, range, etc.).

There are several functions that find descriptive statistics within R. My preferred method uses the Tidyverse and Psych packages, which I describe first. Next I will show you how to find descriptive statistics using base R.

6.1 Descriptive statistics using Tidyverse and Psych packages.

First load the tidyverse and psych packages (if they are not already loaded)

```
library(tidyverse)  
library(psych)
```

The `describe()` function of the Psych package was made to produce the most frequently requested stats in psychology research in an easy to read data frame. I pair this with tidyverse styled code (because of the piping - which I explain next).

Here is the code to get descriptive statistics for the exam2 dataset we made:

```
exam2 %>%  
  describe()
```

- The `describe()` function applies to all of the variables in the dataset (here `exam2`)
- The `%>%` in this code is called a **pipe**
- Pipes are part of the Tidyverse package
- The shortcut to write a pipe (`%>%`) is `shift + command + M` (`shift + alt + M` on a pc)
- Pipes are a way to write strings of functions more easily
- You can think of it as a “THEN”
- So this code would be read as “use the `exam2` dataset” THEN “compute descriptive statistics with the `describe` function”

The twitter handle WeAreRladies uses this example to show the sequential nature of a pipe (`%>%`):

I woke up %>% showered %>% dressed %>% glammed up %>% took breakfast %>% showed up to work

Let's look at the results

```
##          vars   n  mean    sd median trimmed  mad   min   max range skew
## id           1 21 11.00  6.20  11.00  11.00 7.41  1.00 21.00 20.00  0.00
## exam2pts     2 21 15.01  2.31  14.62  14.88 2.65 11.79 19.74  7.95  0.37
## hrsstudye2   3 21  3.02  2.20   2.00   2.88 2.22  0.00  7.00  7.00  0.41
## cheese        4 21  1.52  0.51   2.00   1.53 0.00  1.00  2.00  1.00 -0.09
## wine          5 21  1.48  0.51   1.00   1.47 0.00  1.00  2.00  1.00  0.09
##                  kurtosis   se
## id            -1.37  1.35
## exam2pts      -1.13  0.50
## hrsstudye2   -1.26  0.48
## cheese         -2.08  0.11
## wine          -2.08  0.11
```

The results show that the average exam points was 15.01 (out of 20), with a standard deviation of 2.31 points. Students studied for the exam for an average of 3.02 hours (SD = 2.20).

As the cheese and wine variables are nominal, the mean and standard deviation are not particularly meaningful. Also, any statistics on the ID numbers are meaningless.

This is a good place to mention that it is vital that you as a researcher understand what the numbers you are looking at are and the assumptions that they carry. R (or any computer program) will not tell you if what you asked for does not make sense or is not appropriate.

Rather than getting meaningless results that you have to ignore, you could add the `select()` function to the command above to select certain variables within

6.1. DESCRIPTIVE STATISTICS USING TIDYVERSE AND PSYCH PACKAGES.45

a dataset. Note that you must include more than one variable for the `select()` function. Use the `pull()` function if you want to select only one variable.

For example, to select the `exam2pts` and `hrsstudy2` variables use the following code:

```
exam2 %>%
  select(exam2pts, hrsstudy2) %>%
  describe()

##           vars   n   mean     sd median trimmed   mad    min    max range skew
## exam2pts      1 21 15.01  2.31  14.62  14.88  2.65 11.79 19.74  7.95 0.37
## hrsstudy2      2 21  3.02  2.20   2.00    2.88  2.22  0.00  7.00  7.00 0.41
##                  kurtosis   se
## exam2pts      -1.13 0.50
## hrsstudy2     -1.26 0.48
```

You should create frequency tables for nominal data. Do this with the `count()` function.

For example, create a frequency table for the `cheese` variable with this code:

```
exam2 %>%
  count(cheese)

## # A tibble: 2 x 2
##   cheese     n
##   <dbl> <int>
## 1     1     10
## 2     2     11
```

- this code is saying to "use the exam 2 dataset and then count the cheese variable.

The results show that 10 students did not eat cheese the night before Exam 2 (see Introduction for codebook – or what the 1 and 2 mean) and 11 students did eat cheese the night before the exam.

Finally, often we want to know descriptive statistics by group. For example, say you were interested in relation between cheese eating and exam 2 scores. You would want to know the descriptive statistics of the exam 2 scores for the students who did and did not eat cheese.

To do this I use the `describeBy()` function of the `psych` package, which reports basic summary statistics by a grouping variable. You have to tell R where to find the grouping variable by first including the dataset, followed by a `$` and the

variable name. In this example: `exam2$cheese` (I can't figure out how to avoid the \$ here - I will give extra credit if you can.)

Use the `pull()` function to select the `exam2pts` variable.

```
exam2 %>%
  pull(exam2pts) %>%
  describeBy(exam2$cheese)

##
## Descriptive statistics by group
## group: 1
##   vars n  mean   sd median trimmed  mad   min   max range skew kurtosis   se
## X1    1 10 15.04 2.34 14.94  14.86 2.19 11.79 19.74  7.95 0.41 -0.75 0.74
## -----
## group: 2
##   vars n  mean   sd median trimmed  mad   min   max range skew kurtosis   se
## X1    1 11 14.98 2.4  13.85   14.9 2.48 12.18 18.46  6.28 0.28 -1.78 0.72
```

The results show that the average exam 2 score for students who ate cheese was 15.05 ($SD = 2.34$) and the average exam 2 score for students who did not eat cheese was 14.98 ($SD = 2.4$). Other statistics that you might find useful are the group's n, median, minimum and maximum scores, range, and standard error (se).

6.2 Descriptive statistics using base R

Descriptive statistics can also be computed using base R.

When a variable is stored directly in an object, you can apply the mean and standard deviation functions to the object.

For example:

```
mean(classIQ)
```

```
## [1] 102.4
```

```
sd(classIQ)
```

```
## [1] 11.0363
```

The `summary` function provides the range and median as well:

```
summary(classIQ)

##      Min. 1st Qu. Median     Mean 3rd Qu.     Max.
##    89.0    95.0   101.0   102.4   112.0   115.0
```

Remember that if a variable is in a data frame, you have to tell R to first look in the data frame in order to find the variable. You do this with the dollar sign (\$). Place the \$ between the name of the data frame and the name of the variable.

For example, to find the average points earned on Exam 2 use the following code:

```
mean(exam2$exam2pts)
```

```
## [1] 15.00571
```

Note that when typing this code RStudio will provide a list of the variables in the exam2 after you type the \$. It is very convenient.

The `table()` function of base R performs categorical tabulations of data, frequency tables, and cross tabulations.

For the present example, the code is:

```
table(exam2$cheese)
```

```
##
##  1  2
## 10 11
```

Note that the table is laid out differently than the Tidyverse one above. But you can still easily see that 10 students did not eat cheese the night before Exam 2 and 11 students did eat cheese the night before the exam.

Some people think that the Tidyverse and Psych packages make computing descriptive statistics a bit easier/more direct/better/easier to understand than base R. You should decide which you prefer. (I tend to prefer Tidyverse and Psych). Another package that compute descriptive statistics is skimr