

Tools for Working with Data - 211

Nicole Sorhagen, Ph.D.

2020-08-11

Contents

1	About this book	5
2	Set up project on Rstudio Cloud	7
3	Introduction	9
3.1	Layout of Rstudio cloud	9
3.2	Importing data into Rstudio cloud	17
4	Packages	25
4.1	Installing packages	25
4.2	Loading Packages	29
4.3	Misc	30

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.

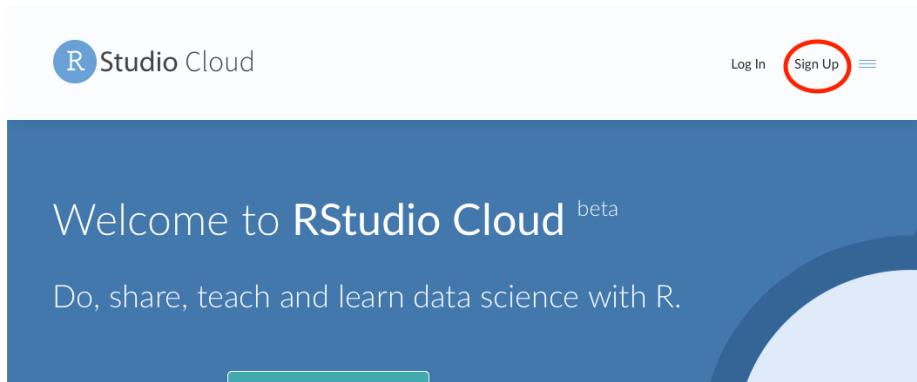
This work is licensed under a Creative Commons Attribution-NoDerivatives 4.0 International License.

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)

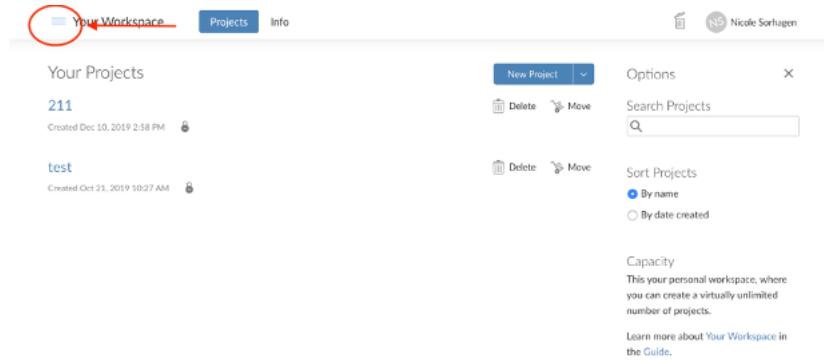


Then join our shared RStudio cloud workspace with the link that I sent you in the email titled 'Rstudio cloud shared workspace'.

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...

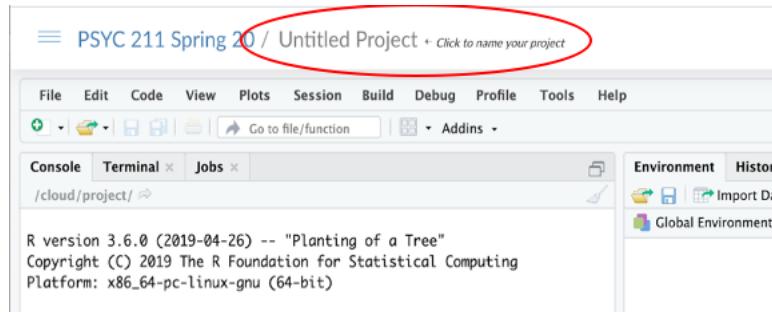
Expand the R studio cloud options by clicking on the 3 lines in the top left corner.



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

Once you are in the shared the classes workspace, open a new project.

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



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 cloud

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

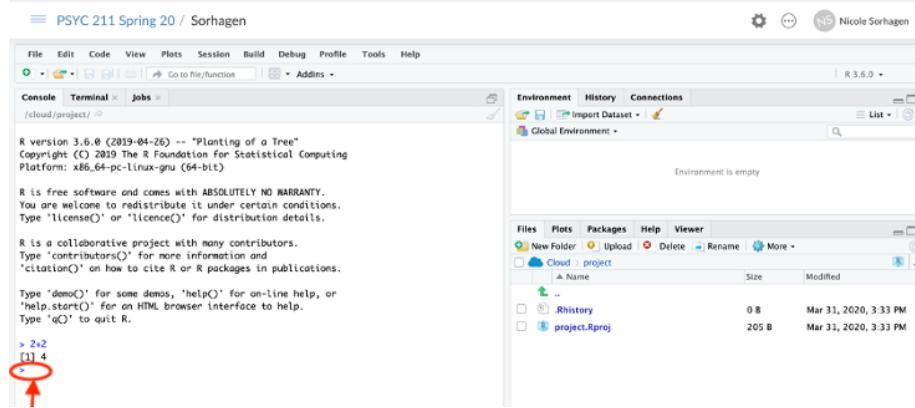
The console panel of R studio is where you can type commands and where you will see the output of commands.

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.

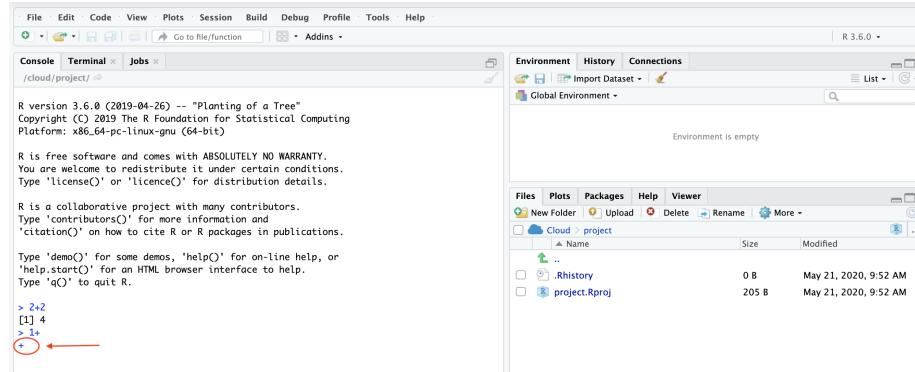
The `>` in the last line of the console means that the console is ready for a command (see red circle in the picture above).



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.



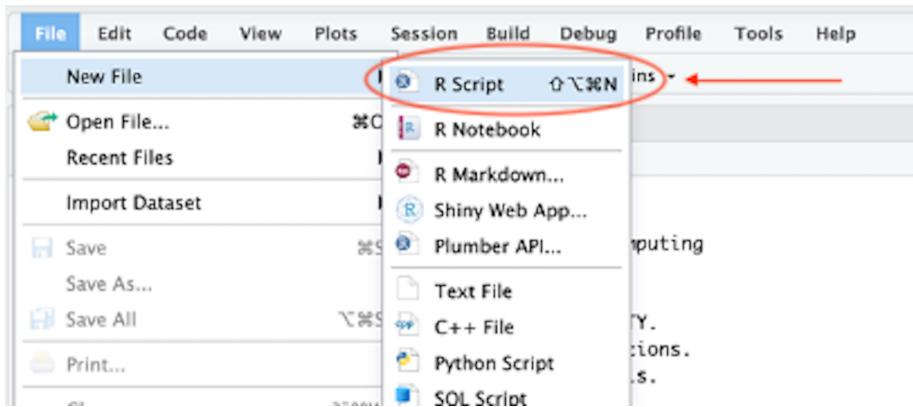
Push the **ESC** button on your keyboard to get back to the command prompt.

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. It also lets you share with others – from professors, to collaborators, to peer reviewers.

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.

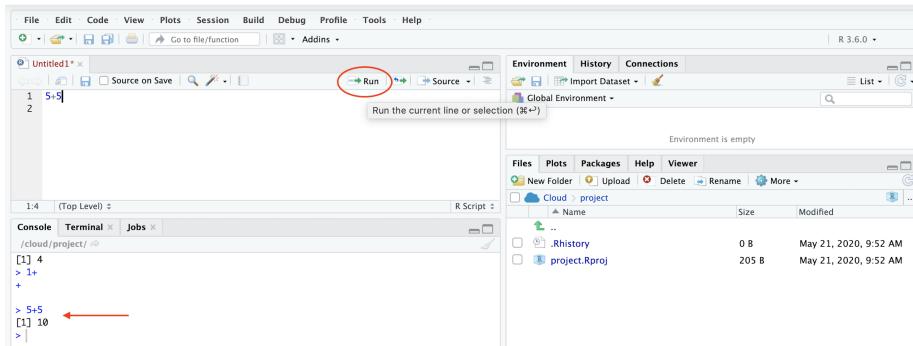
You should run code from scripts

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. 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).

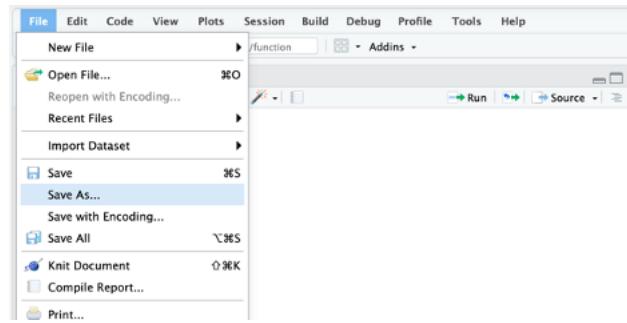
After the code is run, the results will automatically appear the in console (see red arrow in the picture below).



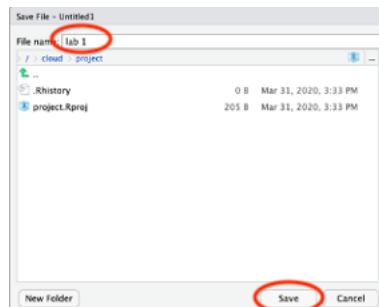
In order to use the script again you must **save** it.

From the drop-down menu select:

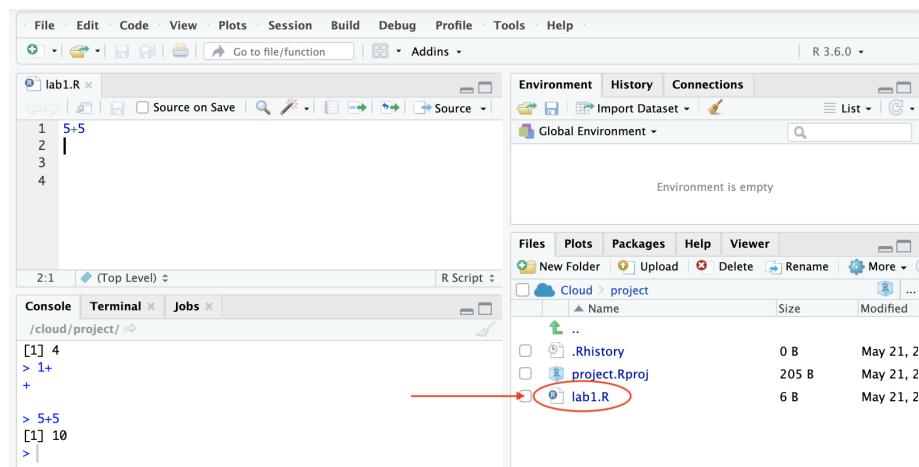
FILE -> SAVE AS



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.

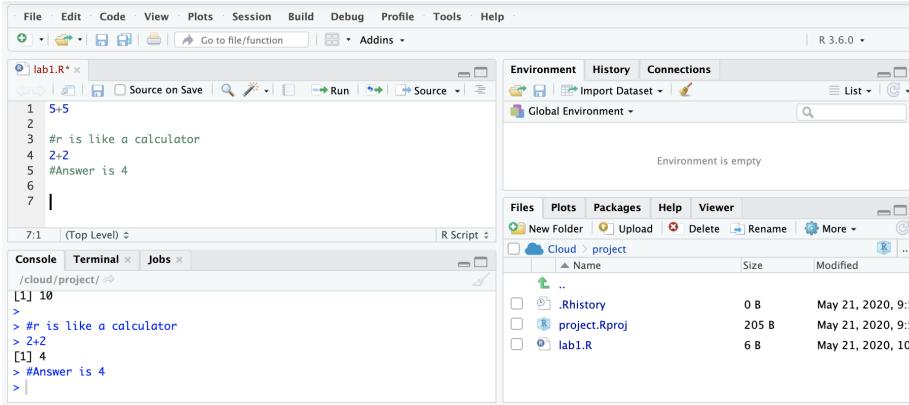


This script file is a record of your work and is how you will be graded for this lab. Make sure you saved this file and complete the rest of this lab in your 'lab1' script.

Within a script you should include comments to yourself and others using **#**. 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 are green in the script.

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.

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.

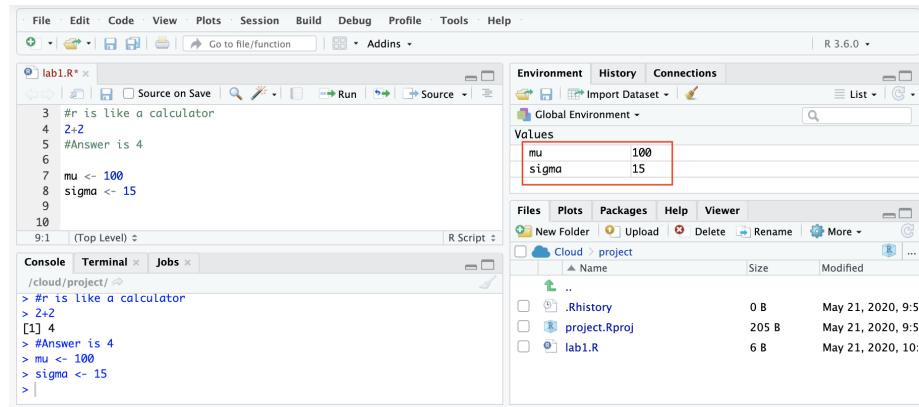
Objects are assigned using `<-`. Values on the right side of `<-` will be assigned to the object on the left side.

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()`. Do use this, type all the numbers you want to store within the parentheses in a comma separated list.

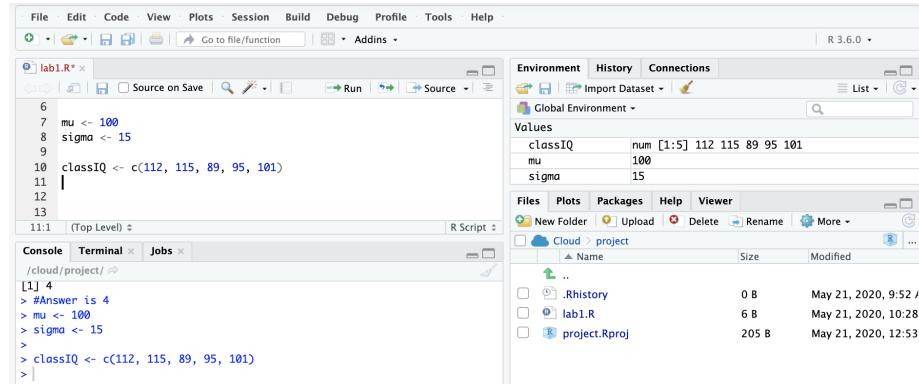
For example, let's enter IQ scores of students in a small class.

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)

The screenshot shows the RStudio Cloud interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. The version R 3.6.0 is displayed. The left pane contains an R script named 'lab1.R' with the following code:

```
7 mu <- 100
8 sigma <- 15
9
10 classIQ <- c(112, 115, 89, 95, 101)
11
12 #get z-scores
13 (classIQ-mu)/sigma
14
```

The right pane shows the Environment and Files panes. The Environment pane displays the Global Environment with variables: classIQ (num [1:5] 112 115 89 95 101), mu (100), and sigma (15). The Files pane shows a project directory with files: .Rhistory (0 B), lab1.R (108 B), and project.Rproj (205 B).

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:

The screenshot shows the RStudio Cloud interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. The version R 3.6.0 is displayed. The left pane contains an R script named 'lab1.R' with the following code:

```
9
10 classIQ <- c(112, 115, 89, 95, 101)
11
12 #get z-scores
13 (classIQ-mu)/sigma
14
15 zscores <- (classIQ-mu)/sigma
16
17
```

The right pane shows the Environment and Files panes. The Environment pane displays the Global Environment with variables: 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). The Files pane shows a project directory with files: .Rhistory (0 B), lab1.R (108 B), and project.Rproj (205 B).

3.1.3.2 Data frames

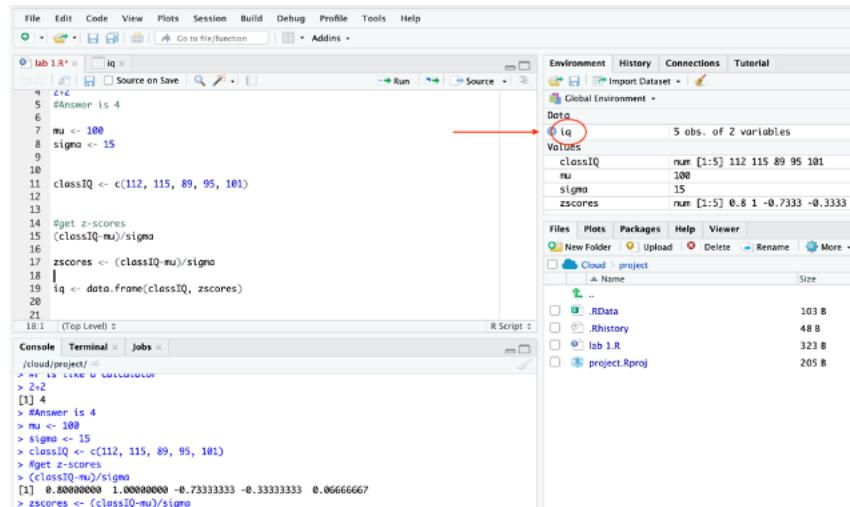
Right now the IQ scores and the z-scores are in separate objects. Variables often need to be in a single object in order to do some basic analyses. You can combine the classIQ and the zscores variable using the **data.frame** command.

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.

```

File Edit Code View Plots Session Build Debug Profile Tools Help
lab 1.R * iq *
Filter
classIQ zscores
1 112 0.80000000
2 115 1.00000000
3 89 -0.73333333
4 95 -0.33333333
5 101 0.06666667

Showing 1 to 5 of 5 entries, 2 total columns
Console Terminal Jobs
/cloud/project/ >
[1] 4
> #Answe is 4
> mu <- 100
> sigma <- 15
> classIQ <- c(112, 115, 89, 95, 101)
> #get z-scores
> (classIQ-mu)/sigma
[1] 0.8000000 1.0000000 -0.7333333 -0.3333333 0.06666667
> zscores <- (classIQ-mu)/sigma
> iq <- data.frame(classIQ, zscores)
> View(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

Finally, 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 nice feature!

```

File Edit Code View Plots Session Build Debug Profile Tools Help
lab 1.R * iq *
Source on Save Run Source
R 3.6.0
12
13
14 #get z-scores
15 (classIQ-mu)/sigma
16
17 zscores <- (classIQ-mu)/sigma
18
19 iq <- data[[
20   data [Completion list]
21   data.class [base]
22   data.entry [utils]
23   data.frame [base]
24   data.matrix [base]
25   data.table [utils]
26   datasets:: [base]
27
28 iq <- data.frame(classIQ, zscores)
Press F1 for additional help
R Script
19.11 (Top level) 2
Console Terminal Jobs
/cloud/project/ > 2> 2
[1] 4
> #Answe is 4
> mu <- 100
> sigma <- 15
> classIQ <- c(112, 115, 89, 95, 101)
> #get z-scores
> (classIQ-mu)/sigma
[1] 0.8000000 1.0000000 -0.7333333 -0.3333333 0.06666667
> zscores <- (classIQ-mu)/sigma

```

3.2 Importing data into Rstudio cloud

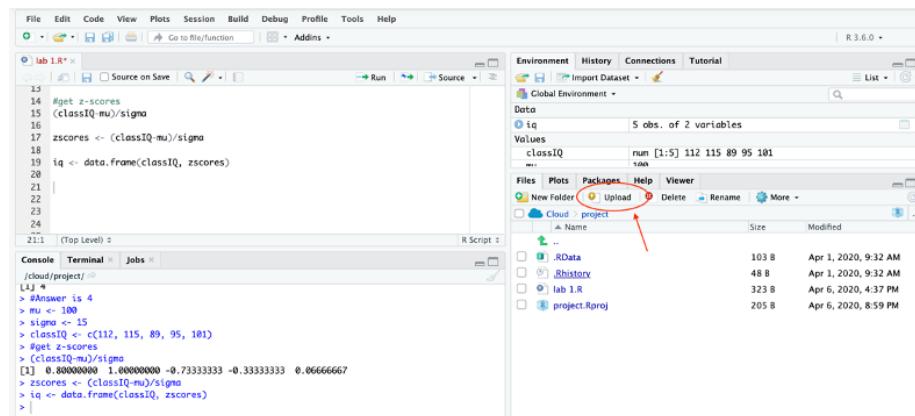
In the Introduction section you learned how to assign data to a **vector** using the **combine** function.

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.

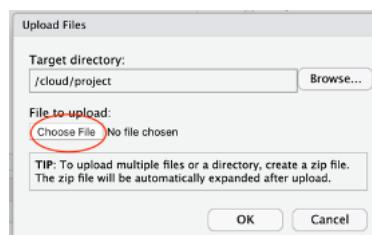
First download the exam2.csv file from d2l.

3.2.1 Upload the data into Rstudio Cloud

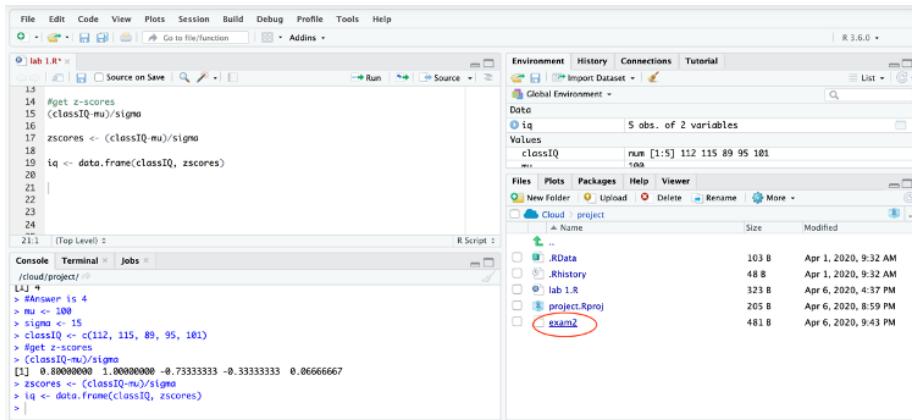
Then select the UPLOAD button in the files window.



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.

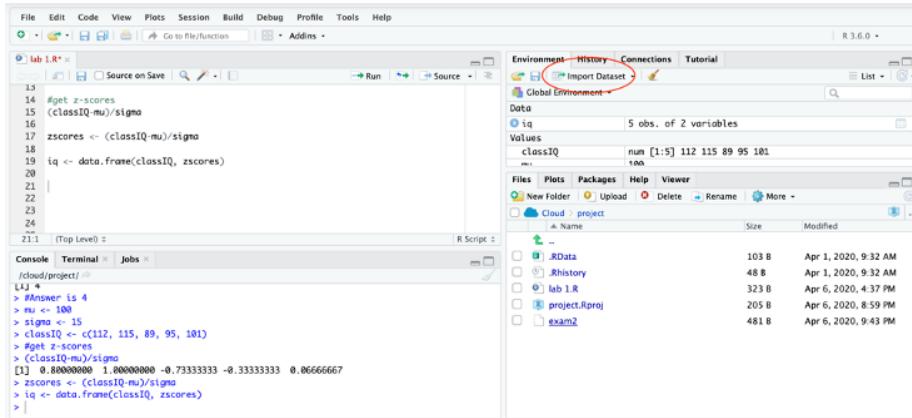


3.2.2 Import data

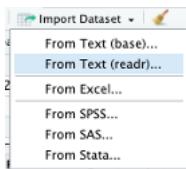
Then you need to import the data into the environment (i.e. assign the data to an object). This can be done through using point and click options or with code.

3.2.2.1 Point and click

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



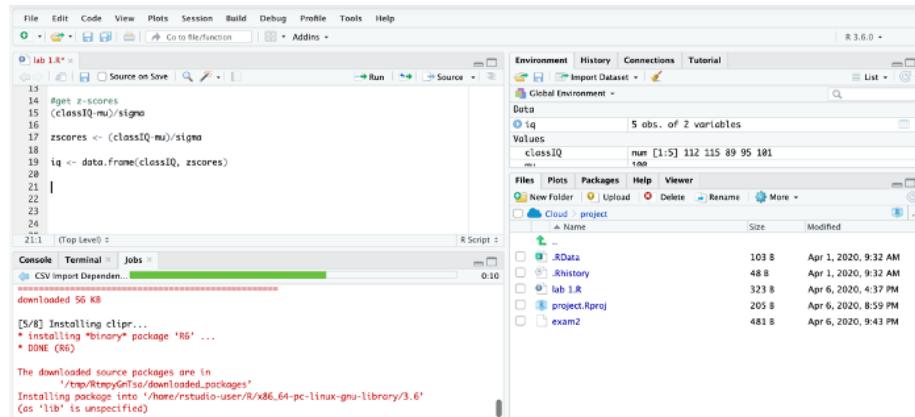
Then select the '**FROM TEXT (READR)**'



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.

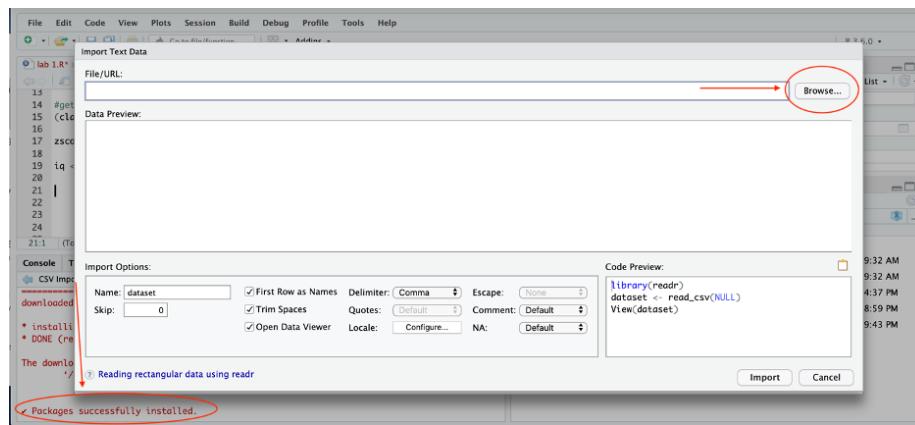


After you select yes, R will begin downloading the package. This can take a few minutes and will look something like this:



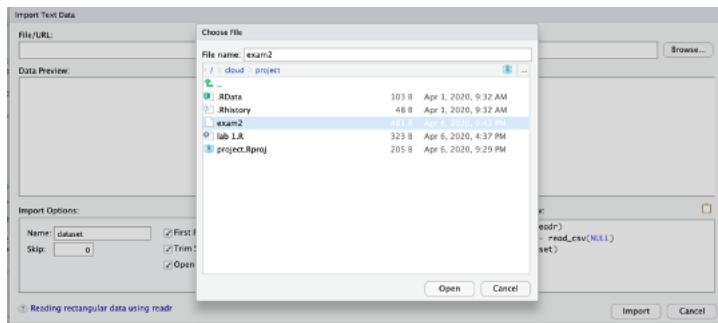
It is important to be *patient* here and let the package download completely before you move on to the next step.

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.



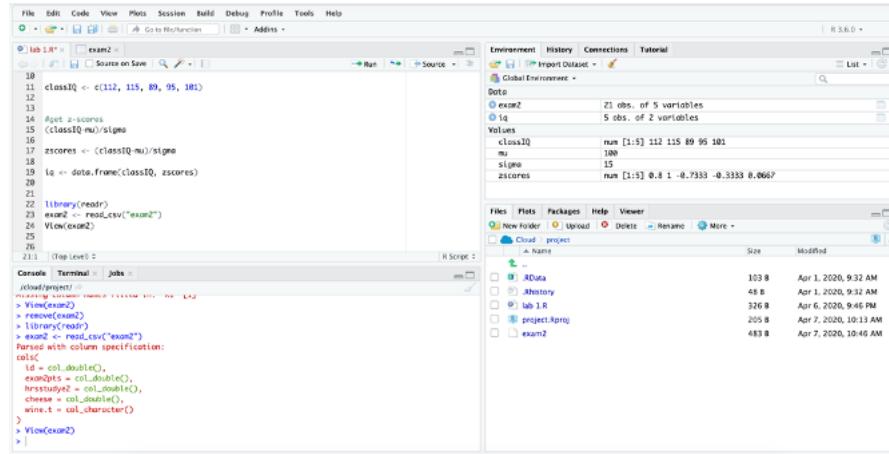
In the choose file window, select ‘exam2’. And then select OPEN.

The next window should look like this:

id (double)	exam2pts (double)	hrsstudy2 (double)	cheese (double)	wine (double)
1	16.15	4.00	2	2
2	14.10	4.25	1	1
3	12.18	2.00	2	1
4	13.46	1.00	1	2
5	18.21	6.00	2	2
6	16.41	5.00	1	1
7	16.03	5.50	1	1
8	18.46	7.00	2	1
9	14.62	1.00	1	1
10	16.67	3.00	1	1
..

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). Copy this code before I 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.

	id	exam2pts	hrsstudy	cheese	wine
1	1	16.15	4.00	2	yes
2	2	14.10	4.25	1	no
3	3	12.18	2.00	2	no
4	4	13.46	1.00	1	yes
5	5	18.21	6.00	2	yes
6	6	16.41	5.00	1	no
7	7	16.03	5.50	1	no
8	8	18.46	7.00	2	no
9	9	14.62	1.00	1	no
10	10	16.67	3.00	1	no
11	11	12.56	0.50	2	yes
12	12	17.69	2.00	2	yes
13	13	12.82	1.50	2	no
14	14	16.15	4.00	2	yes
15	15	15.26	4.50	1	yes
16	16	13.33	1.00	2	no
17	17	13.85	0.75	2	no
18	18	12.31	1.50	1	no
19	19	11.79	0.00	1	yes
20	20	13.33	2.00	2	yes
21	21	19.74	7.00	1	yes

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>).

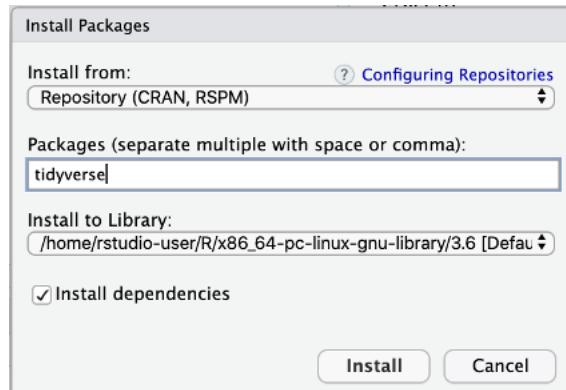
To install a package with GUI go to the top bar menu:

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:

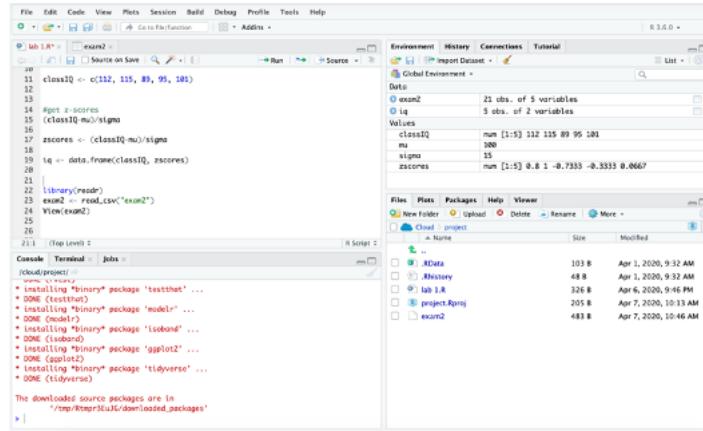
- Console:** Displays R code being run, including `library(readr)` and `install.packages("tidyverse")`.
- Environment:** Shows the Global Environment pane with variables `exam2`, `iq`, `mu`, `sigma`, and `zscores` defined.
- Packages:** Shows the Packages pane where `tidyverse` is listed as being installed.
- File Explorer:** Shows a project named "Cloud" with files like .RData, .Rhistory, lab 1.R, project.Rproj, 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 the progress of package installation, including:
 - Download of broom_0.5.5.tar.gz (2.2 MB)
 - Content type application/x-tar length 1965899 bytes (1.9 MB)
 - Downloaded 575 KB
 - trying URL 'http://package-proxy/src/contrib/broom_0.5.5.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'
- Environment:** Shows the Global Environment pane with variables `exam2`, `iq`, `mu`, `sigma`, and `zscores` defined.
- Packages:** Shows the Packages pane where `tidyverse` is listed as being installed.
- File Explorer:** Shows a project named "Cloud" with files like .RData, .Rhistory, lab 1.R, project.Rproj, and exam2.

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:

```

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
35.1 (Top Level) <--> R Script
36 Console Terminal <--> Jobs
37 /cloud/project/
38
39 downloaded 3.6 MB
40
41 * installing binary package 'moment' ...
42 DONE (moment)
43 * installing binary package 'psych' ...
44 DONE (psych)
45
46 The downloaded source packages are in
47   '/tmp/RtmpQWzul/downloaded_packages'
48
49

```

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.