

Instructions to Download R, RStudio, and Packages

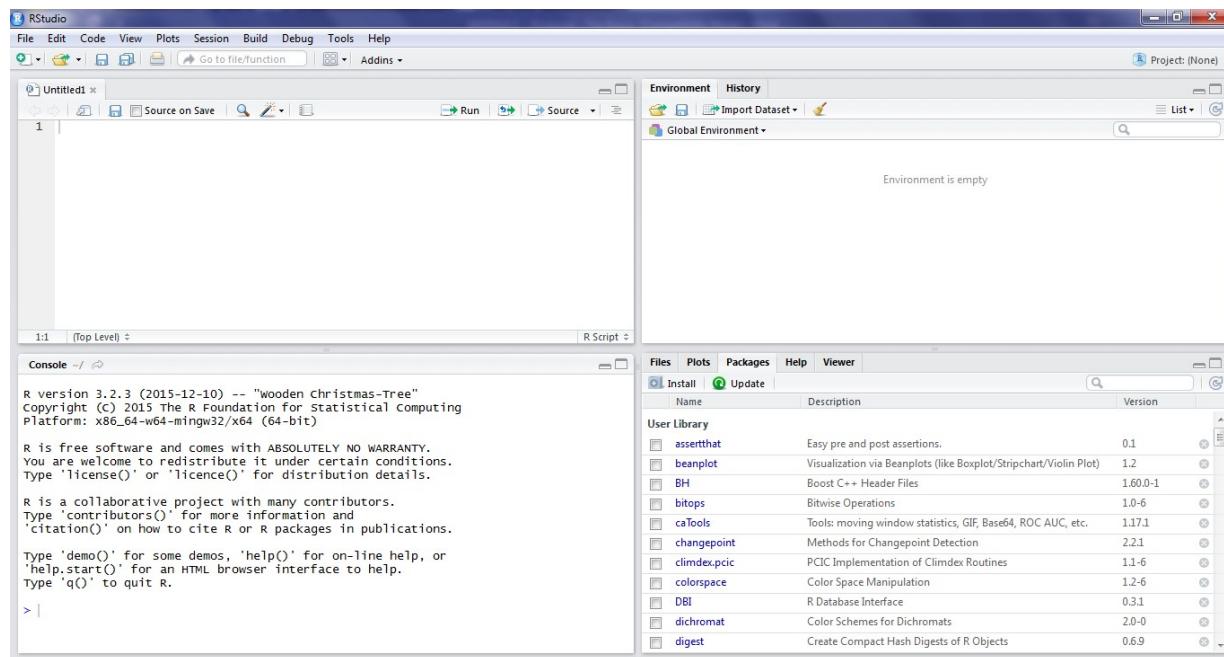
A) You must download base R first, since it must be installed on your computer for RStudio to work.

- 1) Go to this webpage: <https://cloud.r-project.org/>
- 2) Under “Download and Install R,” click the link that corresponds to your computer’s operating system.
- 3) FOR WINDOWS:
 - a) In the first line (after “base”), click the link “install R for the first time”.
 - b) Click “Download R 4 ...” at the top of the page.
- 4) FOR MAC:
 - a) Under “Latest Release”, click on the first file name that ends in “.pkg”.
- 5) After downloading, proceed as you normally would to complete installation of a program.

B) Next, install RStudio – this requires base R to already be installed on your computer.

- 1) Go to this webpage: <https://www.rstudio.com/products/rstudio/download/#download>
- 2) Under “All Installers,” click the link that corresponds to your computer’s operating system.
- 3) After downloading, proceed as you normally would to complete installation of a program.

C) I recommend opening RStudio after installation is complete to ensure it runs correctly. You should see a screen that looks something like this:



D) Install four packages in RStudio that we will use during the workshop.

- 1) Open RStudio.
- 2) Click the “Tools” toolbar at the top of the screen.

Computer Science in Modern Biology

Introduction to R

Instructor:

Rachel Pilla

Nicole Berry

TAs:

Andrew Cannizzaro



1

Workshop Schedule

Day 1 (Monday):

- Introduction to R
- Coding basics
- Working with data in R
- Troubleshooting R scripts

Day 2 (Tuesday):

- Formatting and manipulating data
- Basic analyses and plotting
- Practice group exercises

2

Workshop Logistics

Zoom

- Feel free to unmute yourself to ask questions at any point

Chat

- Post any questions or comments here
- The TAs will be checking the chat to help answer questions

Break-Out Rooms

- If you are really stuck, you can move to a break-out room with a TA to share your screen and get you up to speed
- Break time is a great opportunity for this if needed!

3

What is R?

- FREE and OPEN SOURCE statistical and computational program
- Simple Google searches can give you effective solutions to questions and troubleshooting issues in R
- Widely used in sciences and rapidly growing in popularity
- Can handle more advanced computations, statistical analyses, and bigger data files than Excel
- Lots of styles to code and get the exact same output or result

4

Introduction & Basics in R

5

R script:

- This is where you write and edit your code and comments
- RStudio color codes your scripts to make it easier to identify names of functions, numbers, and comments
- It also automatically fills in parentheses & quotations

6

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- Left Sidebar:** Untitled1 x, Source on Save, Run, Source.
- Console Area:** Displays the R command-line interface. It shows the R version (3.2.3), copyright information, and a welcome message about the "wooden Christmas-Tree". It also includes help documentation for 'demo()', 'help()', 'help.start()', and 'q()'.
- Environment Area:** Shows the Global Environment pane with the message "Environment is empty".
- Packages Area:** Shows the Packages pane listing various R packages with their descriptions and versions.

R console:

- This is where your commands are run, and where results appear
- The “>” symbol in blue means R is ready to work.
- If you see a “+” after running code, then something did not finish running, and/or you have an error or forgot a parenthesis

7

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- Left Sidebar:** Untitled1 x, Source on Save, Run, Source.
- Console Area:** Displays the R command-line interface, identical to the one in the previous screenshot.
- Environment Area:** Shows the Global Environment pane with the message "Environment is empty".
- Packages Area:** Shows the Packages pane listing various R packages with their descriptions and versions.

Environment:

- This is where you can see which objects you have created
- It also tells you what type of object it is, and its size and dimensions

8

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help, and Addins. The top toolbar has icons for file operations like Open, Save, Run, and Source. The left sidebar has tabs for Environment, History, Global Environment, and a search bar. The main area has two large blue boxes: one labeled "R script" containing the text "R script", and another labeled "Environment" with the message "Environment is empty". Below these is the R Console window, which displays the R startup message and help information. To the right of the console is the "Plots, Packages, & Help" pane, which lists installed packages with their descriptions and versions. The packages listed include assertthat, beamplot, BH, bitops, caTools, changepoint, climex.pic, colorspace, DBI, dichromat, and digest.

Plots/Packages/Help:

- This is where your plots will appear
- You can also see what packages are installed and available
- ***Most importantly,*** this is where you find help files that tell you how to use functions/arguments

9

Packages in R

- **Packages** are bundles of tools and functions that others have developed to be used in R
- They are often grouped to specific types of functions, analyses, or datasets
- `rLakeAnalyzer`, for example, has lots of functions developed by limnologists to help you analyze common types of data collected from lakes

10

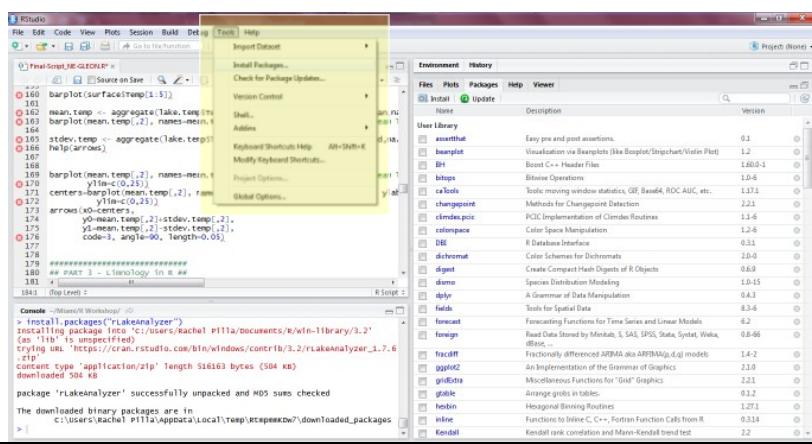
Packages in R

- Currently, there are over 16,000 packages!
- We are going to install 2 packages to use today:
 - 1) dplyr
 - 2) tidyverse

11

Installing a Package

- Tools → Install Packages...



The screenshot shows the RStudio interface with the 'Packages' tab selected in the top navigation bar. The 'Install' tab is active. A list of packages is displayed, including:

Name	Description	Version
assertthat	Easy pre and post assertions.	0.1
broom	Visualization via Broomplots (like BroomPlot/Stripchart/Vislin Plot)	1.2
BH	Boost C++ Header Files	1.60.0-5
bitops	Bitwise Operations	1.0-6
cxtools	Toxic moving window statistics, GBR, Base64, ROC AUC, etc.	1.171
changepoint	Methods for Changepoint Detection	2.21
clm2	PGC Implementation of Clm2dss Routines	1.1-6
colorspace	Color Space Manipulation	2.4
DBI	R Database Interface	0.3.1
dichromat	Color Schemes for Dichiromats	16.0
digest	Create Compact Hash Objects of R Objects	0.6.3
domc	Species Distribution Modeling	0.10.15
dplyr	A Grammar of Data Manipulation	0.4.3
fields	Tools for Spatial Data	8.3-6
forecast	Forecasting Functions for Time Series and Linear Models	6.2
foreign	Read Data Stored by Minitab, S, SAS, SPSS, Stata, Systat, WinBUGS, etc.	0.8-66
frcdiff	Fractionally differenced ARIMA(s,d,B)IMA(q,d,q) models	1.4-2
gridExtra	An Implementation of the 'gridExtra' R package	2.2.0
gridSVG	Miscellaneous Functions for "Grid" Graphics	2.2.1
grid	Arrange grobs in tables.	0.3.2
hexbin	Hexagonal Binning Routines	1.271
inline	Functions to Inline C, C++, Fortran Function Calls from R	0.3.14
Kendall	Kendall rank correlation and Mann-Kendall trend test	2.2

The console at the bottom shows the command `install.packages("LakeAnalyzer")` being run, and the output indicating the package has been successfully unpacked and NAMESPACE checked.

12

Installing a Package

- Type the packages you want to install:
dplyr, tidyverse

13

Installing a Package

- Console will tell you when they have been installed

14

Installing a Package

- Packages only need to be installed ONCE
- But, each time you re-open R, you need to load the package(s) you want to use so R can make all the functions available:

```
library(package name)
```

 Hadley Wickham ✅
@hadleywickham

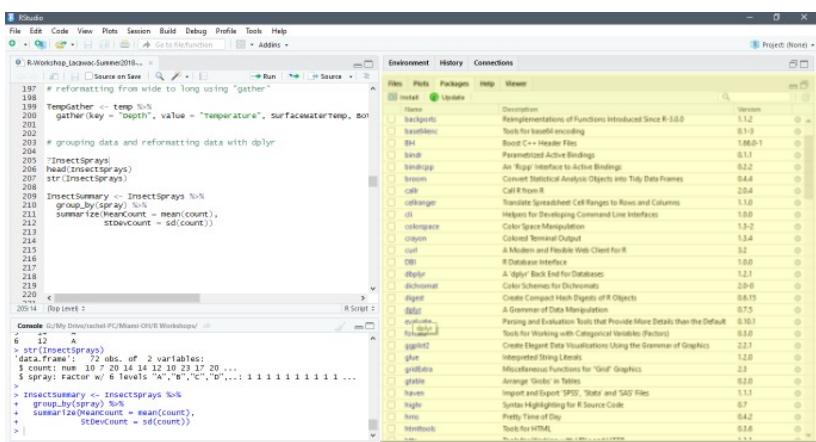
Replies to @ijlytle
@ijlytle a package is like a book, a library is like a library; you use library() to check a package out of the library #rsats

8:34 AM · Dec 8, 2014 · Echofon

15

Package Information

- Under the “Packages” tab, click on **dplyr**

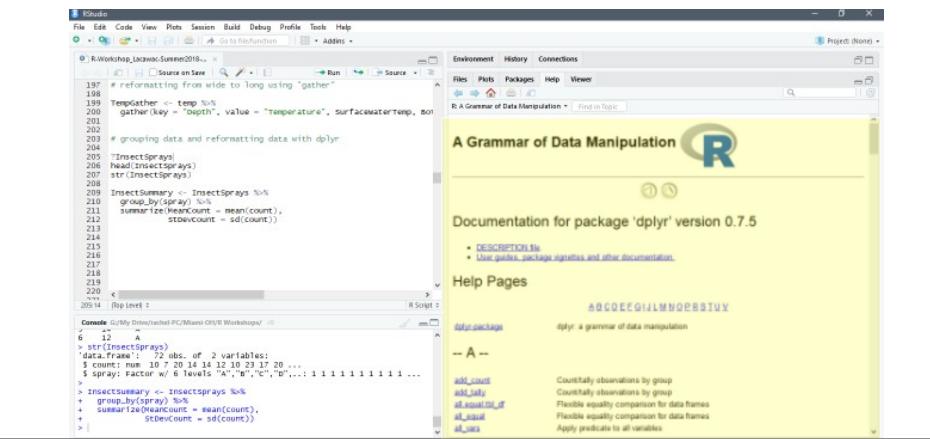


Name	Description	Version
backports	Reimplementations of Functions Introduced Since R-3.0.0	1.1.2
gridExtra	Extensible Grid Arranging Function	0.1.3
Rcpp	Boost C++ Interface	1.0.6-1
grid	Parametrized Active Bindings	0.1.1
gridExtra	An 'Rcpp' Interface to Active Bindings	0.2.2
grid	Convert Statistical Analysis Objects into Tidy Data Frames	0.4.4
grid	Call R from R	2.0.4
gridExtra	Translate Between Cell Ranges to Rows and Columns	1.0.0
grid	Tools for Overriding Command Line Interfaces	1.1.0
grid	Color Space Manipulation	1.3.2
grid	Colored Terminal Output	1.3.4
grid	A Modern and Flexible Web Client for R	3.2
DBI	R Database Interface	1.0.0
dplyr	A 'dplyr' Back End for Databases	1.2.1
grid	Grid Based Drawing	2.3.0
grid	Create Compact Hash Digests of R Objects	0.8.75
grid	Parsing and Evaluation Tools that Provide More Details than the Default	0.10.3
gridExtra	Tools for Working with Categorical Variables (Factors)	0.1.0
grid	Create Elegance Data Visualizations Using the Grammar of Graphics	2.2.1
grid	Integrated String Lengths	1.2.0
grid	Interactive Grids for 'Grid' Graphics	2.1
grid	Arrange 'Grids' in Tables	0.3.0
haven	Import and Export 'SPSS', 'Stata' and 'SAS' Files	1.3.1
highr	Syntax Highlighting for R Source Code	0.7
htmltools	Pretty Print of Day	0.4.2
htmltools	Tools for HTML	0.1.6

16

Package Information

- This will list all the available functions
- Clicking on a function takes you to the help file



17

Using an R Script

- Write commands in the script (top left pane)
 - Save it, edit it, revisit it later, etc.
- Code NOT automatically run when you hit Enter
- To run code:
 - “Run” button in upper right corner of R Script
 - “CTRL + Enter” shortcut on PC
 - “Command + Return” shortcut on Mac

18

Using an R Script

- Add in **comments** using “#”
 - Comments can be any information in an R script that will NOT be run, but can be used for notes/explanations/etc.
 - Anything behind “#” will be ignored by R
- R is case sensitive
 - “Mean” ≠ “mean”
- R doesn’t care too much about spaces (or lack thereof)

19

Using an R Script

- Be careful to close all parentheses & quotations (and in the correct spots!)
 - A majority of your errors will come from missing or misplaced parentheses or quotations – but they’re usually easy fixes!
- “>” at the beginning of Console means it’s ready to go!
 - “+” means code is still running or it is not complete
 - Often missing parenthesis, bracket, or quotation

20

Key Components of Code

- **Functions** allow you to manipulate data, apply calculations, run statistical analysis, and more
- **Arguments** are the defining information for functions, and allow you to “customize” the function to have it do what you want it to do
- **Objects** are pieces of data saved in R that can be called up, reused, or manipulated

21

Key Components of Code

```
x <- seq(1, 10)
```

name of the **function** to
create a **sequence**

22

Key Components of Code

x <- seq(1,10)

the arguments to define the function, to create a sequence from 1 through 10

23

Key Components of Code

x <- seq(1,10)

name of the object that saves the results of the function in R

24

Key Components of Code

x <- seq(1, 10)

names of objects must begin with a letter,
but can include uppercase or lowercase
letters, numbers, underscores, and periods –
make your object names clear and useful!

25

Key Components of Code

x <- seq(1, 10)

assignment operator tells R to
save the result of the function
as the named object

26

Types of Data

An object can hold data that is:

- Numeric
- Integer
- Characters (strings)
- Logical (TRUE/FALSE)
- Complex ($1+4i$)

27

Types of Data Structures

- An object can be in the form of:
 - Vector
 - 1-D object with the same data type
 - Matrix
 - 2-D object with the same data type
 - Data frame
 - 2-D object with different data types for each column if desired
 - Generally the most useful for biological and ecological data

28

Live Coding

open RStudio and code along!

29

Computer Science in Modern Biology

Introduction



Instructor:

Rachel Pilla

Nicole Berry

TAs:

Andrew Cannizzarro



30

Review of Day 1

31

Workshop Logistics

Zoom

- Feel free to unmute yourself to ask questions at any point

Chat

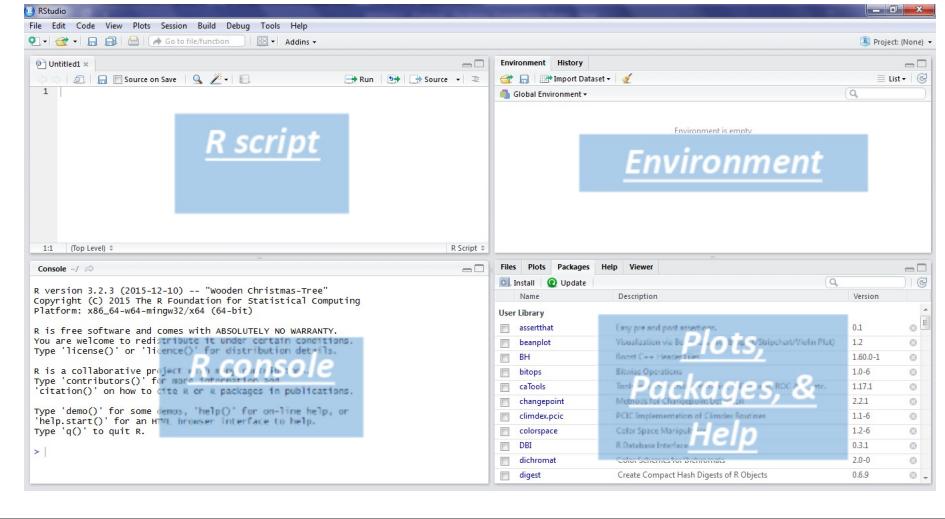
- Post any questions or comments here
- The TAs will be checking the chat to help answer questions

Break-Out Rooms

- If you are really stuck, you can move to a break-out room with a TA to share your screen and get you up to speed
- Break time is a great opportunity for this if needed!

32

Layout of RStudio



33

Packages

- Packages only need to be installed ONCE
- But, each time you re-open R (today!), you need to load the packages you want to use so R can make all the functions available:

```
library(package name)
```



Hadley Wickham
 @hadleywickham

Replies to [@ijlytle](#)

[@ijlytle](#) a package is like a book, a library is like a library; you use library() to check a package out of the library #rsats

8:34 AM · Dec 8, 2014 · [Echofon](#)

34

Using an R Script

- To run code:
 - “Run” button in upper right corner of R Script
 - “CTRL + Enter” shortcut on PC
 - “Command + Return” shortcut on Mac
- “>” at the beginning of Console means it’s ready to go!
 - “+” means code is still running or it is not complete
 - Often missing parenthesis, bracket, or quotation

35

Key Components of Code

- **Functions** allow you to manipulate data, apply calculations, run statistical analysis, and more
- **Arguments** are the defining information for functions, and allow you to “customize” the function to have it do what you want it to do
- **Objects** are pieces of data saved in R that can be called up, reused, or manipulated

36

Types of Data

An object can hold data that is:

- Numeric
- Integer
- Characters (strings)
- Logical (TRUE/FALSE)
- Complex ($1+4i$)

37

Types of Data Structures

- An object can be in the form of:
 - Vector
 - 1-D object with the same data type
 - Matrix
 - 2-D object with the same data type
 - Data frame
 - 2-D object with different data types for each column if desired
 - Generally the most useful for biological and ecological data

38

Live Coding

open RStudio and code along!

39

Additional Tips & Useful Resources

40

“NA” Values

- R recognizes “NA” or blank cells from Excel as missing values
- If you have another format for indicating “NA” values, you can add an argument when importing your data file:
`read.csv(... , na.strings = "n/a")`
- Many functions will automatically ignore/skip over “NA” values
 - Be sure you read the associated help files to ensure the function is treating “NA” values as you are expecting, or adjust the argument accordingly!

41

Assorted Useful Functions

- | | |
|-----------------|---------------|
| • ncol(...) | • cbind(...) |
| • nrow(...) | • rbind(...) |
| • colnames(...) | • rep(...) |
| • rownames(...) | • seq(...) |
| • length(...) | • range(...) |
| • dim(...) | • sum(...) |
| • unique(...) | • sd(...) |
| • paste(...) | • var(...) |
| • sort(...) | • diff(...) |
| • order(...) | • round(...) |
| • t(...) | • ifelse(...) |

42

Swirl

- Package for self-guided learning of R in the R environment
- Learn at your own pace and select topics most useful to you
 - Introductory material
 - Regressions
 - Data wrangling
 - Statistical inference
 - Advanced programming
- <https://swirlstats.com/>

43

RStudio Cheat Sheets

- Short, handy guides to help with many packages and programming needs in RStudio
- Available free for download in multiple language options
- <https://rstudio.com/resources/cheatsheets/>

44

ggplot2 Cheat Sheet

Data Visualization with ggplot2 : CHEAT SHEET

This cheat sheet provides a quick reference for ggplot2 functions across various categories:

- Basics:** Explains the grammar of graphics, components of a ggplot, and how to display values.
- Graphical Primitives:** Shows examples for points, lines, polygons, rectangles, and text.
- Two Variables:** Focuses on continuous x and discrete x variables.
- One Variable:** Focuses on continuous y variables.
- Three Variables:** Focuses on continuous x, continuous y, and discrete z variables.
- Continuous Functions:** Focuses on density, area, and lines.
- Continuous Bivariate Distribution:** Focuses on joint distributions.
- Visualizing Error:** Focuses on error bars.
- Maps:** Focuses on choropleth maps.
- Workflow:** Provides tips for building a ggplot.

The R Studio logo is at the bottom left.

45

ggplot2 Cheat Sheet

This cheat sheet covers various aspects of ggplot2:

- Stats:** An alternative way to build a layer.
- Scales:** Maps values to the visual values of an aesthetic. Includes general purpose scales, coordinate systems, and position adjustments.
- Coordinate Systems:** Focuses on polar coordinates, facets, and themes.
- Faceting:** Divides a plot into subplots based on third variables.
- Position Adjustments:** Positions elements relative to one another.
- Labels:** Adds labels to plots.
- Themes:** Themes for ggplot2.
- Legends:** Legend options for ggplot2.
- Zooming:** Zooming features for ggplot2.

The R Studio logo is at the bottom left.

46

Tidyverse

- Packages to help organize data and streamline R coding, specifically for data science
- Functions in the packages share similar syntax and coding structure
- “dplyr”, “tidyr”, and “ggplot2” are part of tidyverse
- <https://www.tidyverse.org/>

47

Lubridate

- Package for more effectively working with dates between Excel and R
- Can recognize date/times with minimal formatting issues
- Also allows for date/time differencing, among other useful functions
- Included in tidyverse
- <https://lubridate.tidyverse.org/>

48

49

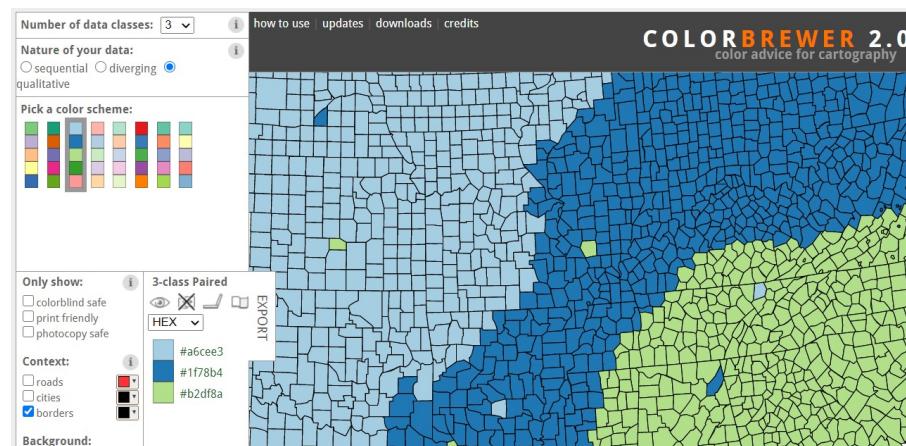
50

ColorBrewer

- Pre-made color schemes and palettes
- Effective graphic color selections for maximum perception
 - Qualitative, sequential, and diverging
- Includes options for photo-copy friendly, color-blind friendly, etc.
- <http://colorbrewer2.org/>

51

ColorBrewer



52

Computer Science in Modern Biology

Thank you for joining!

to share feedback or ask questions,
please feel free to contact me:

 pillarm@miamioh.edu
 [@rmpilla](https://twitter.com/rmpilla)



- 3) Click “Install Packages...”.
- 4) Make sure “Repository (CRAN)” is selected in the “Install from:” box.
- 5) Type the names of packages to install in the “Packages” box:
`dplyr, tidyverse`
- 6) Make sure “Install dependencies” is checked.
- 7) Click “Install”.
- 8) If successfully installed, you should see something like this in the “Console” panel:



The screenshot shows the RStudio Console window with the title "Console G:/My Drive/rachel-PC/Miami-OH/". The console output displays the installation of several R packages from CRAN. It shows the URL being tried, the content type, length, and download time for each package. After the download, it lists the packages successfully unpacked and MD5 sums checked. Finally, it indicates the location where the downloaded binary packages are stored.

```
Console G:/My Drive/rachel-PC/Miami-OH/
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.0/dplyr_1.0.0.zip'
Content type 'application/zip' length 1304095 bytes (1.2 MB)
downloaded 1.2 MB

trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.0/tidyr_1.1.0.zip'
Content type 'application/zip' length 1514564 bytes (1.4 MB)
downloaded 1.4 MB

trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.0/ggplot2_3.3.2.zip'
Content type 'application/zip' length 4066720 bytes (3.9 MB)
downloaded 3.9 MB

trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.0/lubridate_1.7.9.zip'
Content type 'application/zip' length 1748550 bytes (1.7 MB)
downloaded 1.7 MB

package 'rlang' successfully unpacked and MD5 sums checked
package 'tidyselect' successfully unpacked and MD5 sums checked
package 'vctrs' successfully unpacked and MD5 sums checked
package 'dplyr' successfully unpacked and MD5 sums checked
package 'tidyverse' successfully unpacked and MD5 sums checked
package 'ggplot2' successfully unpacked and MD5 sums checked
package 'lubridate' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:\Users\rache\AppData\Local\Temp\Rtmp63z4Ek\downloaded_packages
```

*If you have any problems during the installation,
please email me (pillarm@miamioh.edu) and I'll do my best to help you out!*

Group 1: Mammals (load “ggplot2”, and use built-in dataset “msleep”)

- 1) How many different species are included in this data set? How many orders do they represent? Provide numerical results.
- 2) Is there a correlation between total sleep and REM sleep? Provide statistical results and a visual.
- 3) Is there a difference in time spent awake based on eating habit? Provide statistical results and a visual.
- 4) Can body weight explain brain weight of these species? Is this relationship improved when removing outliers? Provide statistical results and visuals.

Group 2: Diamonds (load “ggplot2”, and use built-in dataset “diamonds”)

- 1) What is the average length and width of each cut of diamond? Provide numerical results and a visual(s).
- 2) Is the carat of these diamonds significantly different than 1 (average for center stone)? Provide statistical results.
- 3) Is there a difference in price based on color? If so, which color(s) have the highest prices? Provide statistical results and a visual.
- 4) Can price be explained by depth? Provide statistical results and a visual.

Group 3: Housing (load “ggplot2”, and use built-in dataset “txhousing”)

- 1) What is the total number of housing sales for each city from 2000-2015 combined? Provide numerical results and a visual.
- 2) Is median sale price across all cities significantly different across months? Provide statistical results and a visual.
- 3) Is there a change in average yearly inventory for Dallas or for Houston since 2000? Provide statistical results and a visual(s).
- 4) Does total number of listings explain total number of sales by city, summed across all years? Which cities are the two highest? Provide statistical results and a visual.

Notes

