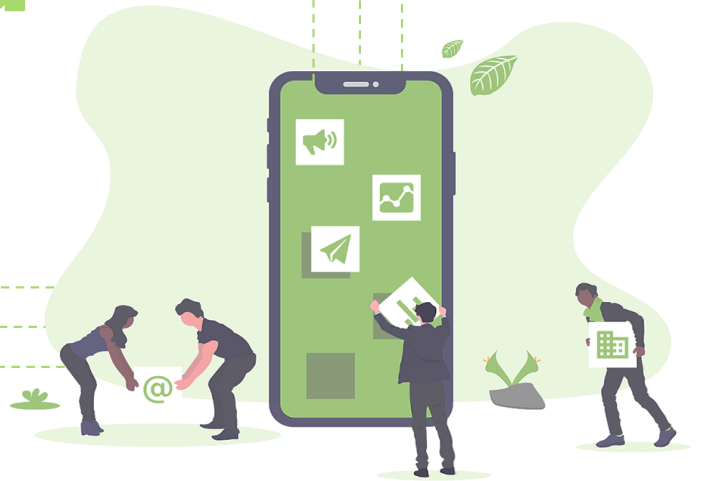


# The Data Science Track



Prepared By: R. Daynalo

1

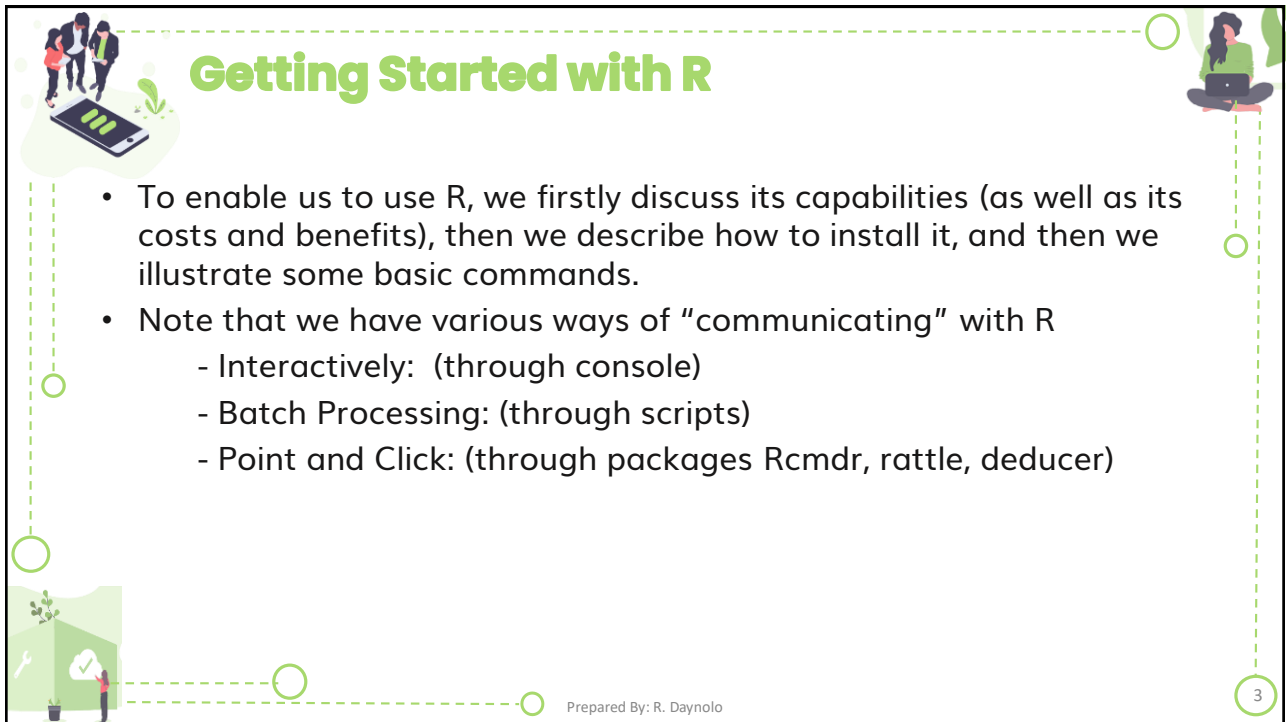
1

## 7. Overview and History of R



2

1



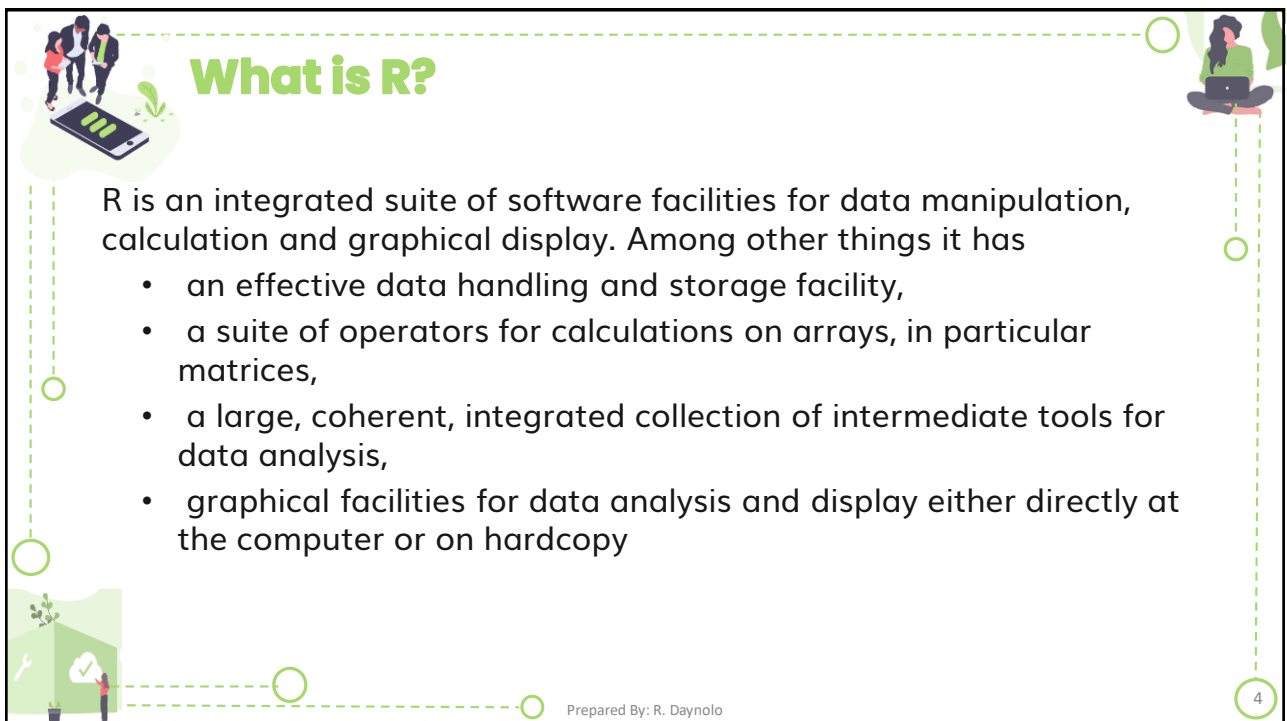
## Getting Started with R

- To enable us to use R, we firstly discuss its capabilities (as well as its costs and benefits), then we describe how to install it, and then we illustrate some basic commands.
- Note that we have various ways of “communicating” with R
  - Interactively: (through console)
  - Batch Processing: (through scripts)
  - Point and Click: (through packages Rcmdr, rattle, deducer)

Prepared By: R. Daynolo

3

3



## What is R?


R is an integrated suite of software facilities for data manipulation, calculation and graphical display. Among other things it has

- an effective data handling and storage facility,
- a suite of operators for calculations on arrays, in particular matrices,
- a large, coherent, integrated collection of intermediate tools for data analysis,
- graphical facilities for data analysis and display either directly at the computer or on hardcopy


Prepared By: R. Daynolo

4


4



## What is R?




- R is a statistical programming environment for performing standard and specialized statistical methods
  - "environment" : a fully planned and coherent system, rather than an incremental accretion of very specific and inflexible tools, as is frequently the case with other data analysis software
- R is a free/open source statistical package
  - based on S language developed at Bell Laboratories




Prepared By: R. Daynolo

5


5



## What is S?




- S is a language that was developed by John Chambers and others at Bell Labs.
- S was initiated in 1976 as an internal statistical analysis environment – originally implemented as Fortran libraries.
- Early versions of the language did not contain functions for statistical modeling.
- In 1988 the system was rewritten in C and began to resemble the system that we have today (this was Version 3 of the language). The book ***Statistical Models in S*** by Chambers and Hastie (the white book) documents the statistical analysis functionality.




Prepared By: R. Daynolo

6


6



## What is S?




- Version 4 of the S language was released in 1998 and is the version we use today. The book ***Programming with Data*** by John Chambers (the green book) documents this version of the language.




Prepared By: R. Daynola

7


7



## Historical Notes




- 1991: Created in New Zealand by Ross Ihaka and Robert Gentleman of University of Auckland. Their experience developing R is documented in a 1996 *JCGS* paper.
- 1993: First announcement of R to the public.
- 1995: Martin Mächler convinces Ross and Robert to use the GNU General Public License to make R free software.
- 1996: A public mailing list is created (R-help and R-devel)
- 1997: The R Core Group is formed (containing some people associated with S-PLUS). The core group controls the source code for R.
- 2000: R version 1.0.0 is released.
- 2018: R Version 3.5.0 is released.



Prepared By: R. Daynola


8

8




## Back to R

- Although R is a programming language, i.e. generating computer code to complete tasks is required, there are now Graphical User Interface (GUI) Add Ons like R Commander, and rattle, which allow users to "point and click".
- Initially developed by Robert Gentleman and Ross Ihaka and now maintained by the "R core development team"
- Cross platform compatibility: Windows, MacOS, Linux




Prepared By: R. Daynola

9





## Costs and Benefits of R

ADVANTAGES	DISADVANTAGES
• Fast and free.	• Not user friendly, i.e. steep learning curve with minimal GUI.
• State of the art: Statistical researchers provide their methods as R packages.	• No commercial support; figuring out correct methods or how to use a function on your own can be frustrating.
• Some packages, such as Mx, WinBugs, and other programs use or will use R.	• Easy to make mistakes and not know.
• Active user community	• Some users complain about hostility on the R listserve
• Excellent for computer intensive analyses, etc., and Interfaces with database storage software (SQL)	• Working with large datasets is limited by RAM
• Forces you to think about your analysis.	• Data prep & cleaning can be messier & more mistake prone in R vs. SPSS or SAS



Prepared By: R. Daynola

10





## Costs and Benefits of R

Free Software!

- The freedom to run the program, for any purpose (freedom 0).
- The freedom to study how the program works, and adapt it to your needs (freedom 1). Access to the source code is a precondition for this.
- The freedom to redistribute copies so you can help your neighbor (freedom 2).
- The freedom to improve the program, and release your improvements to the public, so that the whole community benefits (freedom 3). Access to the source code is a precondition for this.



<https://www.fsf.org/>



Prepared By: R. Daynalo

11


11

## Design of the R System

The R system is divided into 2 conceptual parts:


1. The "base" R system that you download from CRAN
2. Everything else




Prepared By: R. Daynalo

12

12




## Design of the R System



R functionality is divided into a number of packages:


- The “base” R system contains, among other things, the **base** package which is required to run R and contains the most fundamental functions.
- The other packages contained in the “base” system include **utils**, **stats**, **datasets**, **graphics**, **grDevices**, **grid**, **methods**, **tools**, **parallel**, **compiler**, **splines**, **tcltk**, **stats4**.
- There are also “Recommend” packages: **boot**, **class**, **cluster**, **codetools**, **foreign**, **KernSmooth**, **lattice**, **mgcv**, **nlme**, **rpart**, **survival**, **MASS**, **spatial**, **nnet**, **Matrix**.




Prepared By: R. Daynola

13


13



## R Packages



- When you download R from the Comprehensive R Archive Network (CRAN), you get that “base” R system
- The base R system comes with basic functionality; implements the R language
- One reason R is so useful is the large collection of packages that extend the basic functionality of R
- R packages are developed and published by the larger R community



Prepared By: R. Daynola

14

14

## Obtaining R Packages

- The primary location for obtaining R packages is **CRAN**.
- For biological applications, many packages are available from the **Bioconductor Project**.
- You can obtain information about the available packages on CRAN with the `available.packages()` function.

```
a <- available.packages()
head(rownames(a), 3) ## Show the names of the first few packages
```

```
## [1] "A3"      "abc"     "abcdeFBA"
```

- Currently, the CRAN package repository features **13,626** available packages, covering a wide range of topics.

Prepared By: R. Daynalo

15

15

## Installing R Packages

- Packages can be installed with the `install.packages()` function in R.
- To install a single package, pass the name of the package to the `install.packages()` function as its first argument.
- The following code installs the **ggplot2** package from CRAN

```
install.packages("ggplot2")
```

- The command downloads the **ggplot2** package from CRAN and installs it on your computer
- Any packages on which this package depends will also be downloaded and installed.

Prepared By: R. Daynalo

16

16



## Installing R Packages

- You can install multiple R packages at once with a single call to `install.packages()`
- Place the names of the R packages in a character vector

```
install.packages(c("ggplot2", "lattice", "devtools"))
```

Prepared By: R. Daynola

17

17

## Installing R Packages in RStudio

The screenshot shows the RStudio interface with the 'User Library' pane open. The pane displays a list of installed and available packages. The 'User Library' pane is divided into two sections: 'User Library' and 'Available Packages'. The 'User Library' section lists the following packages:

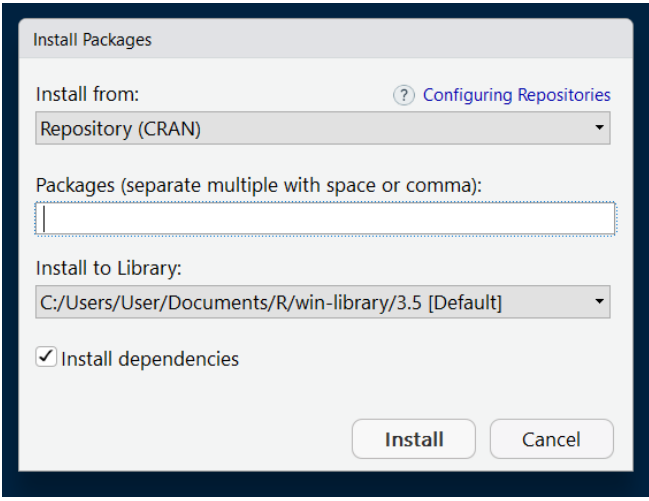
Package Name	Description	Version
abind	Combine Multidimensional Arrays	1.4-5
acepack	ACE and AAS for Selecting Multiple Regression Transformations	1.4-1
AppliedPredictiveMo...	Functions and Data Sets for 'Applied Predictive Modeling'	1.1-7
assertthat	Easy Pre and Post Assertions	0.2.0
backports	Reimplementations of Functions Introduced Since R-3.0.0	1.1.2
base64enc	Tools for base64 encoding	0.1-3
BayesFactor	Computation of Bayes Factors for Common Designs	0.9.12-4.2
BH	Boost C++ Header Files	1.66.0-1
bindr	Parameterized Active Bindings	0.1.1
bindrcpp	An 'Rcpp' Interface to Active Bindings	0.2.2
bitops	Bitwise Operations	1.0-6
broom	Convert Statistical Analysis Objects into Tidy Tibbles	0.5.0
devtools	Tools for Managing R Packages	1.11.0

Prepared By: R. Daynola

18

18

## Installing R Packages in RStudio



The image shows the 'Install Packages' dialog box in RStudio. It has a title bar 'Install Packages'. Below the title bar, there is a section 'Install from:' with a dropdown menu set to 'Repository (CRAN)' and a link '? Configuring Repositories'. Below that is a text input field for 'Packages (separate multiple with space or comma):'. Then there is a section 'Install to Library:' with a dropdown menu set to 'C:/Users/User/Documents/R/win-library/3.5 [Default]'. At the bottom left, there is a checked checkbox 'Install dependencies'. At the bottom right, there are two buttons: 'Install' and 'Cancel'.

Prepared By: R. Daynola

19

## Installing an R Package from Bioconductor Project

- To install core packages, type the following in an R command window:
 


```
if (!requireNamespace("BiocManager"))
  install.packages("BiocManager")
BiocManager::install()
```
- To install specific packages, e.g., "GenomicFeatures" and "AnnotationDbi", with
 

```
BiocManager::install(c("GenomicFeatures", "AnnotationDbi"))
```

<http://bioconductor.org/install/>

Prepared By: R. Daynola

20





## Installing an R Package from Bioconductor Project

- Installing a package does not make it immediately available to you in R; you must load the package
- The `library()` function is used to load packages into R
- The following code is used to load the `ggplot3` package into R

```
library(ggplot2)
```


Note: Do not put the package name in quotes!

Prepared By: R. Daynola

21



21



## Some R Resources

Available from CRAN (<https://cran.r-project.org/>)


- An Introduction to R
- Writing R Extensions
- R Data Import/Export
- R Installation and Administration (mostly for building R from sources)
- R Internals (not for the faint of heart)


Prepared By: R. Daynola

22

22



## Some Useful Books on S/R



**Standard Texts**

- Chambers (2008). *Software for Data Analysis*, Springer. (your textbook)
- Chambers (1998). *Programming with Data*, Springer.
- Venables & Ripley (2002). *Modern Applied Statistics with S*, Springer.
- Venables & Ripley (2000). *S Programming*, Springer.
- Pinheiro & Bates (2000). *Mixed-Effects Models in S and S-PLUS*, Springer.
- Murrell (2005). *R Graphics*, Chapman & Hall/CRC Press.


**Other Resources:**

- A longer list of books is at <http://www.r-project.org/doc/bib/R-books.html>


Prepared By: R. Daynola

23

23



## R Environment

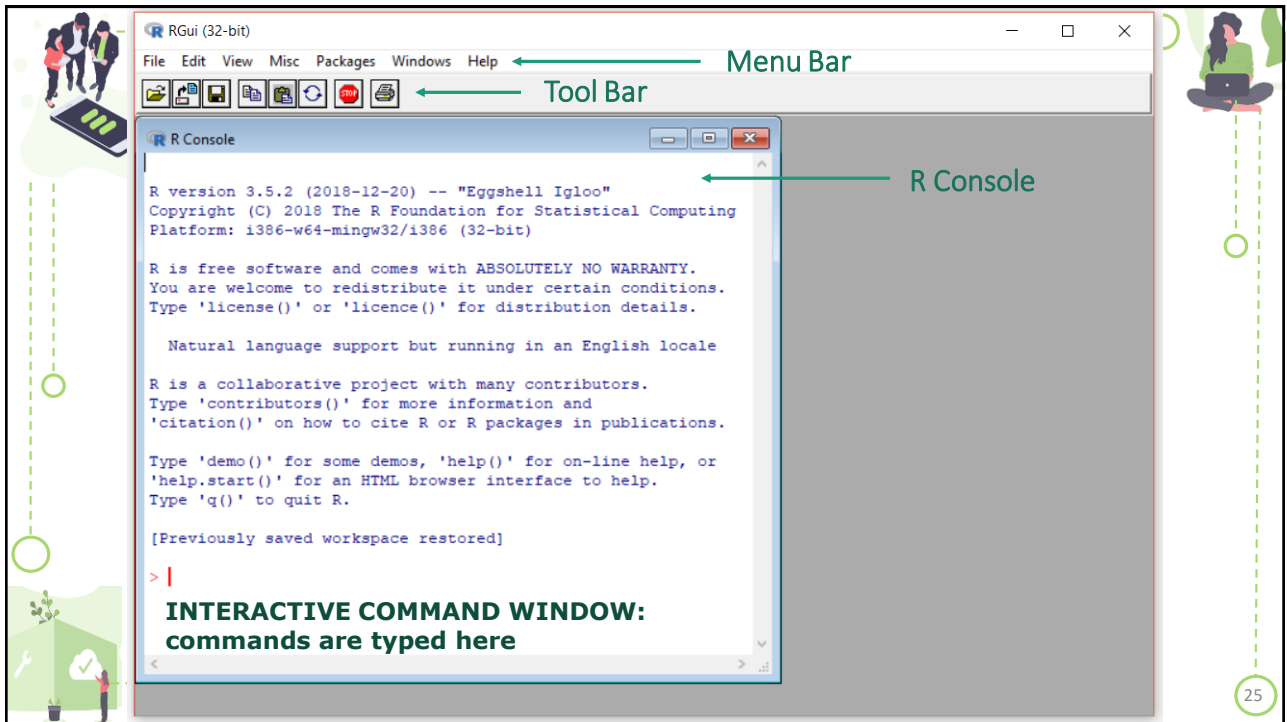


- Start-up window
  - Menu bar
  - Tool bar
  - Command window or R console
- Graphics device
- R codes
  - Expression language with simple syntax
  - Case-sensitive
  - Use of text editor (e.g. R editor, or Notepad, but not MS Word) for ease of management of codes
  - Always save R codes in a text editor rather than saving your work space

Prepared By: R. Daynola

24

24




25

## Setting up R


- To change the default working directory:  
`setwd(file="drive:/path/folder/")`
- To access help files for a function:  
`?function`
- To access help files for a text:  
`??text` or `??"text phrase"`
- To view an example:  
`example(function)`
- To add comments:  
`#commentline`

Prepared By: R. Daynola


26



## Interfaces




- For Windows and MAC OS, the standard R download comes with an R GUI, which is adequate for simple tasks
- **RStudio**. Very popular, with a nice interface and well thought out, especially for more advanced usage: can be a bit buggy, so make sure you update it regularly. Available in all platforms.




Prepared By: R. Daynola

27


27



## RStudio



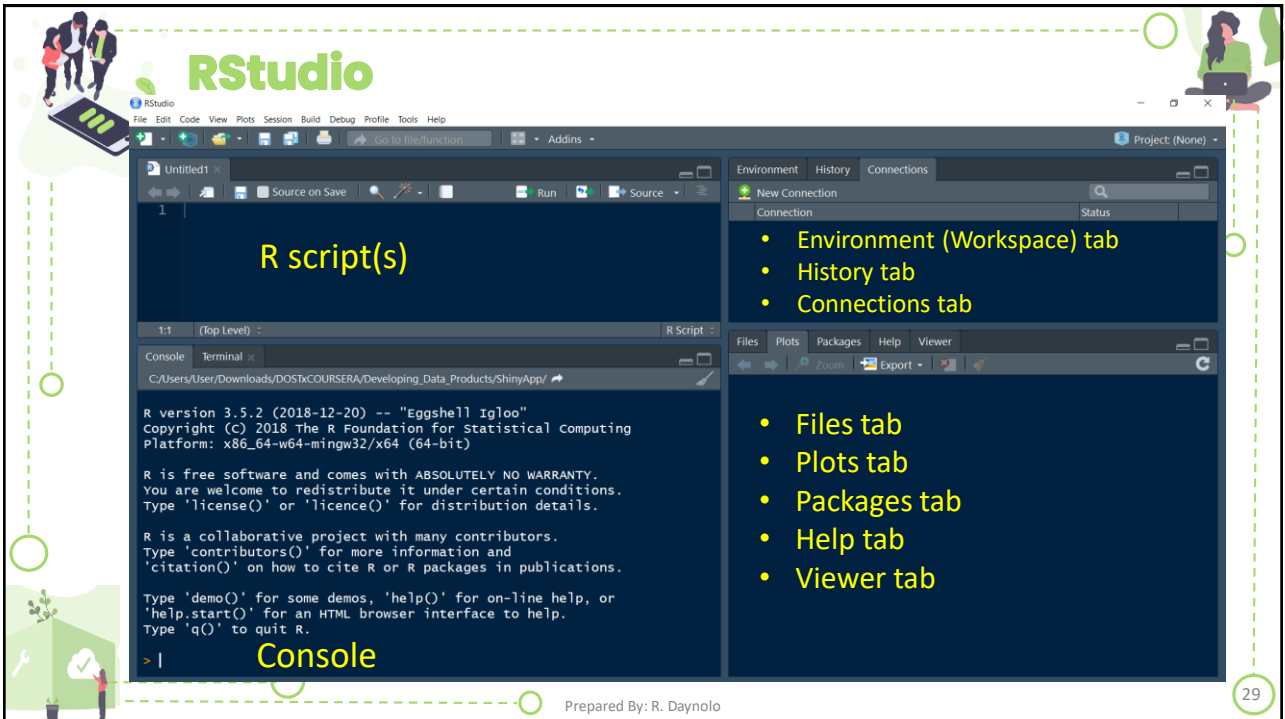
- RStudio allows the user to run R in a more user-friendly environment. It is open-source (i.e. free) and available at <http://www.rstudio.com/>



Prepared By: R. Daynola

28

28



**RStudio**

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Untitled1 x

R script(s)

1:1 (Top Level) R Script

Console

```
C:/Users/User/Downloads/DOSTxCOURSE/Developing_Data_Products/ShinyApp/
```

R version 3.5.2 (2018-12-20) -- "Eggshell Igloo"  
copyright (C) 2018 The R Foundation for Statistical Computing  
Platform: x86\_64-w64-mingw32/x64 (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, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

> |

Environment History Connections

New Connection

Connection Status

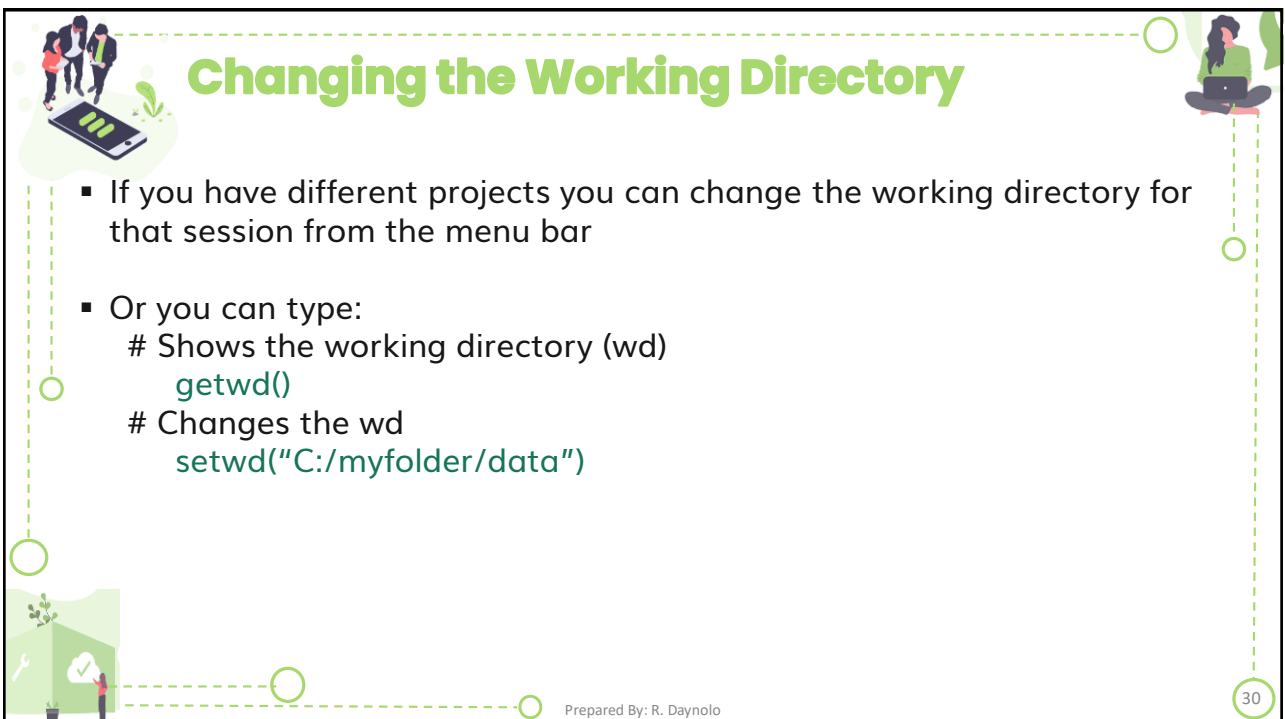
- Environment (Workspace) tab
- History tab
- Connections tab

Files Plots Packages Help Viewer

- Files tab
- Plots tab
- Packages tab
- Help tab
- Viewer tab

Prepared By: R. Daynola

29



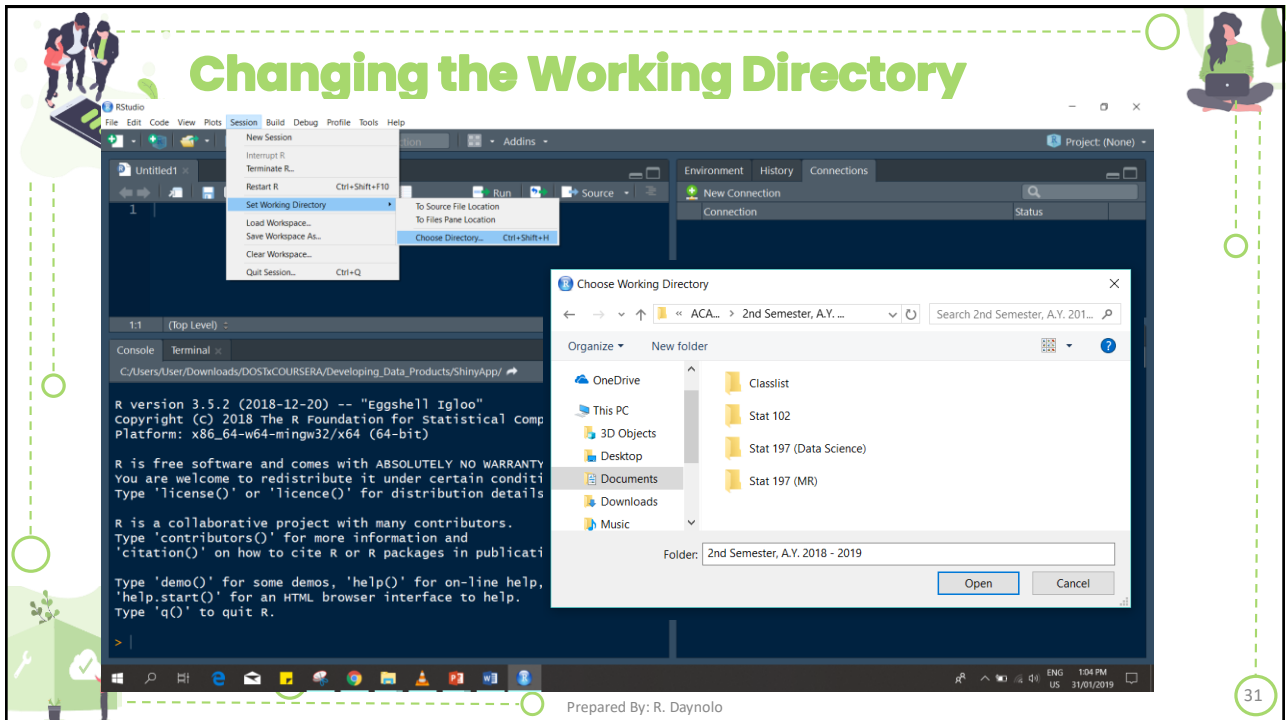
## Changing the Working Directory

- If you have different projects you can change the working directory for that session from the menu bar
- Or you can type:
  - # Shows the working directory (wd)  
`getwd()`
  - # Changes the wd  
`setwd("C:/myfolder/data")`

Prepared By: R. Daynola

30

# Changing the Working Directory

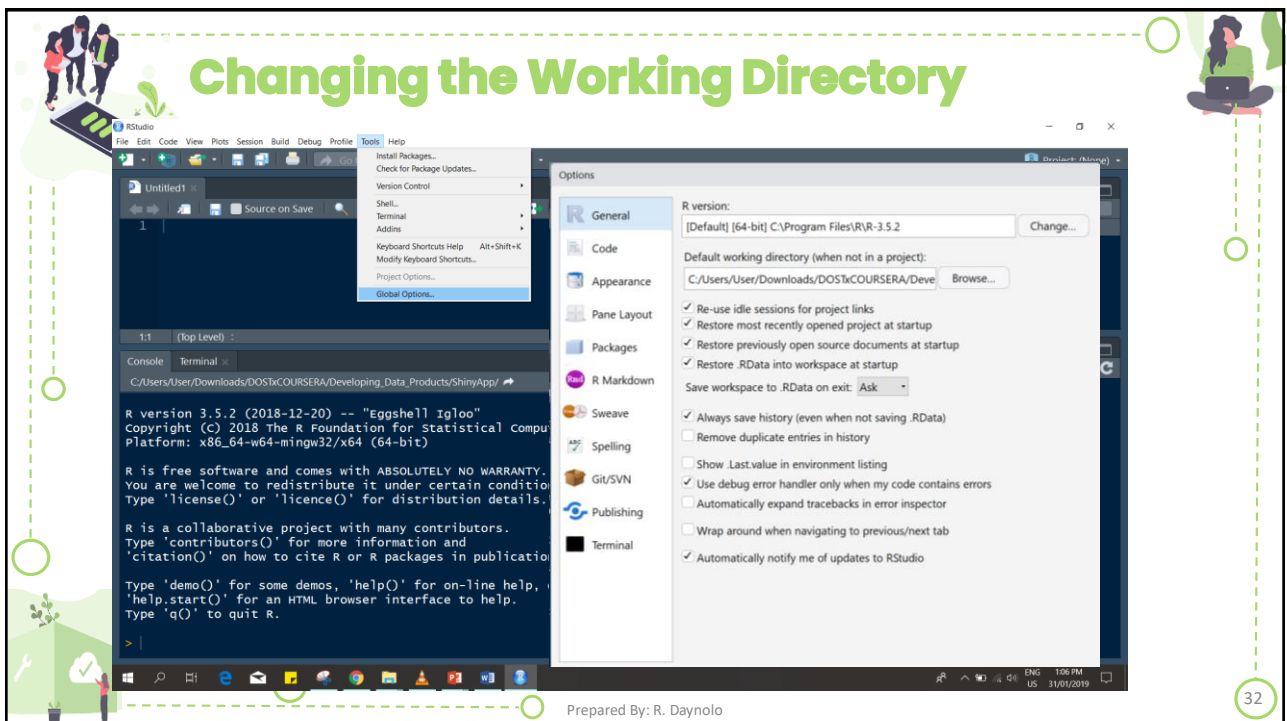


The screenshot shows the RStudio interface. The 'Session' menu is open, and 'Set Working Directory' is selected. The 'Choose Working Directory' dialog box is displayed, showing the current directory as 'C:\Users\User\Downloads\DOSTxCOURSE\Developing\_Data\_Products\ShinyApp\'. The 'Folder' field is set to '2nd Semester, A.Y. 2018 - 2019'. The 'Open' button is highlighted.

Prepared By: R. Daynola

31

# Changing the Working Directory



The screenshot shows the RStudio interface with the 'Global Options' dialog box open. The 'General' tab is selected, showing the 'R version' as '[Default] [64-bit] C:\Program Files\R\R-3.5.2'. The 'Default working directory (when not in a project):' is set to 'C:\Users\User\Downloads\DOSTxCOURSE\Developing\_Data\_Products\ShinyApp\'. The 'Save workspace to .RData on exit' is set to 'Ask'.

Prepared By: R. Daynola

32



## Package Tab

- The package tab shows the list of add-ons included in the installation of RStudio. If checked, the package is loaded in R, if not, any command related to that package won't work, you will need to select it. You may also install other add-ons by clicking on the **Install Packages** icon.
- Another way to activate a package is by typing, for example, `library(foreign)`. This will automatically check the `-foreign` package (it helps bring data from proprietary formats like Stata, SAS, or SPSS)

Prepared By: R. Daynola

33

33

## Package Tab

The screenshot shows the RStudio interface with the Package Tab selected. The console displays the R version 3.5.2 (2018-12-20) and the user's location. The Package Tab shows a list of installed and available packages.

Name	Description	Version	Status
abind	Combine Multidimensional Arrays	1.4-5	●
acepack	ACE and AVAS for Selecting Multiple Regression Transformations	1.4.1	●
AppliedPredictive...	Functions and Data Sets for 'Applied Predictive Modeling'	1.1-7	●
assertthat	Easy Pre and Post Assertions	0.2.0	●
backports	Reimplementations of Functions Introduced Since R-3.0.0	1.1.2	●
base64enc	Tools for base64 encoding	0.1-3	●
BayesFactor	Computation of Bayes Factors for Common Designs	0.9.12-4.2	●
BH	Boost C++ Header Files	1.66.0-1	●
bindr	Parametrized Active Bindings	0.1.1	●

Prepared By: R. Daynola

34

34

## Plots Tab

- The plots tab will display the graphs.
- To extract the graph, click on **Export** where you can save the file as an image (PNG, JPG, etc.) or as PDF, these options are useful when you only want to share the graph. Probably the easiest way to export a graph is by copying it to the clipboard and then paste it directly into your Word document.

Prepared By: R. Daynola

35

35

## Plots Tab

```

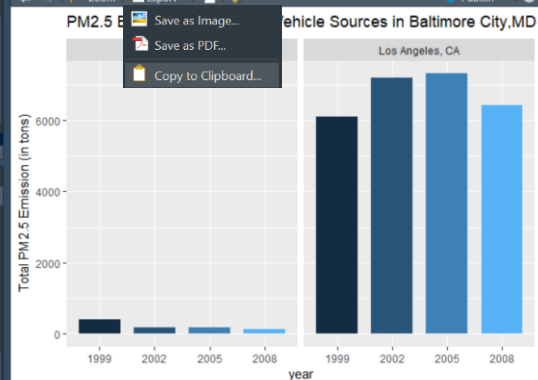
20
21 #Merge Baltimore and Los Angeles data
22 dataNEI <- rbind(baltimoreNEI, losAngelesNEI)
23
24 #Plot and save in PNG format
25 #png("plot6.png",width=640,height=480,units="px")
26
27 library(ggplot2)
28
29 g <- ggplot(dataNEI, aes(x=factor(year), y=Emissions, fill=city)) +
30   geom_bar(aes(fill=year), stat="identity", width=0.75) +
31   facet_grid(.~city) +
32   guides(fill=FALSE) + theme_grey() +
33   labs(x="year", y=expression("Total PM2.5 Emission (in tons)")) +
34   labs(title=expression("PM2.5 Emissions from Motor Vehicle Sources in Baltimore City, MD and Los Angeles, CA (1999-2008)"))
35   print(g)
36
37
37:1 (Top Level) :
R Script

```

```


> g <- ggplot(dataNEI, aes(x=factor(year), y=Emissions, fill=city)) +
+   geom_bar(aes(fill=year), stat="identity", width=0.75) +
+   facet_grid(.~city) +
+   guides(fill=FALSE) + theme_grey() +
+   labs(x="year", y=expression("Total PM2.5 Emission (in tons)")) +
+   labs(title=expression("PM2.5 Emissions from Motor Vehicle Sources in Baltimore City, MD and Los Angeles, CA (1999-2008)"))
> print(g)
> #dev.off()
> print(g)

```




36


36



## Help and Documentation




- There is a large amount of (free) documentation and help available. Some help is automatically installed. Typing in the console window the command  
`> help(rnorm)`
- gives help on the `rnorm()` function. It gives a description of the function, possible arguments and the values that are used as default for optional arguments. Typing  
`> example(rnorm)`
- gives some examples of how the function can be used




Prepared By: R. Daynolo

37


37



## R Scripts




- R is an interpreter that uses a command line based environment.
- This means that you have to type commands, rather than use the mouse and menus.
- This has an advantage that you do not always have to retype all commands.
- You can store your commands in files, the so-called scripts.
- These scripts have typically file names with extension `.R`, e.g. `foo.R`.




Prepared By: R. Daynolo

38

38




## R Scripts



- You can run (send to the console window) part of the code by selecting the lines and pressing CTRL + ENTER or click Run in the editor window.
- If you do not select anything, R will run the line your cursor is on.
- You can always run the whole script with the console command source, so e.g. for the script in the file foo.R you type:  


```
> source("foo.R")
```
- You can also click Run all in the editor window or type CTRL+SHIFT+S to run the whole script at once.



Prepared By: R. Daynalo

39


39



## R Scripts



- To create a new R script you can either go to **File > New R Script**, or click on the icon with the "+" sign and select "R Script", or simply press **CTRL+SHIFT+N**.
- Make sure to save the script.

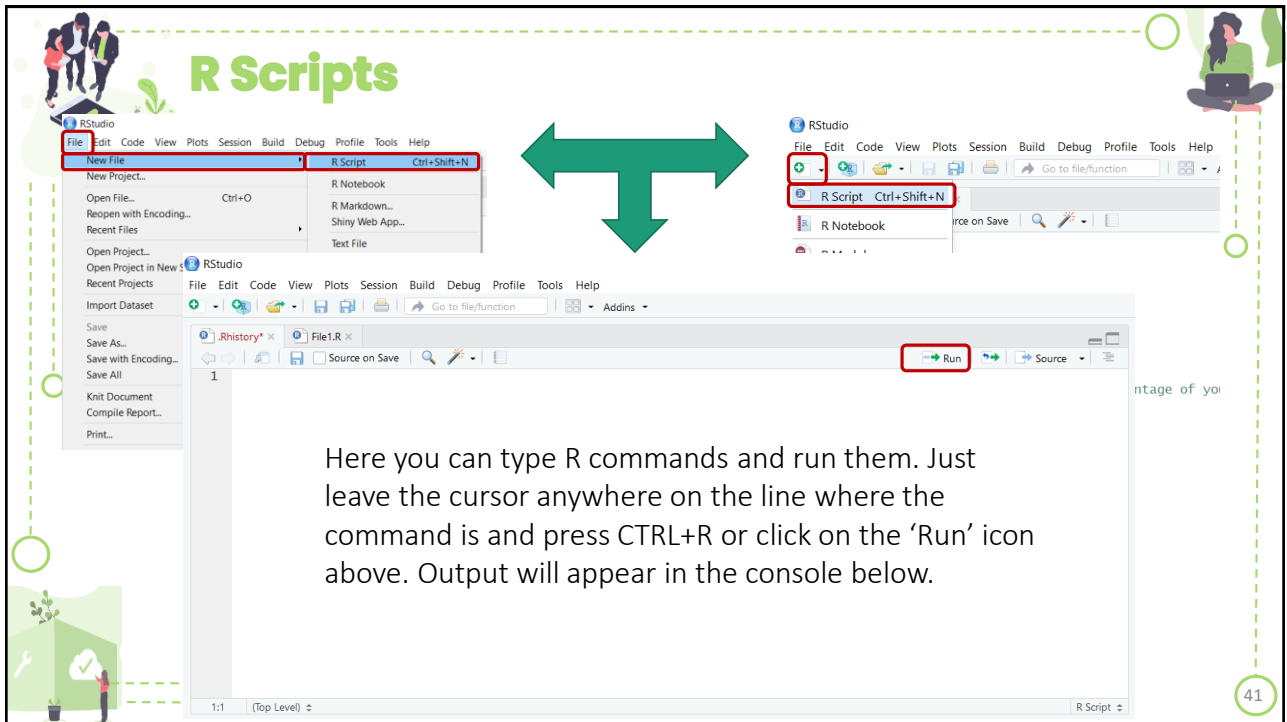


Prepared By: R. Daynalo

40

40

## R Scripts



Here you can type R commands and run them. Just leave the cursor anywhere on the line where the command is and press CTRL+R or click on the 'Run' icon above. Output will appear in the console below.

41

41

## Three Ways of Quitting from R Session

1. Enter in Command Window:  
`> quit()`
2. Click on  
**File** ► **Exit**
3. Click on Close button (X at upper right hand corner of R console window).

Prepared By: R. Daynola

42

42