

Data Visualization in R

Dr. Alexander Mark Weber

Assistant Professor (Partner), Department of Pediatrics,
Division of Neurology, Faculty of Medicine,
Associate Member, Department of Neuroscience,
Associate Member, School of Biomedical Engineering,
Independent Investigator, BC Children's Hospital Research Institute,
University of British Columbia

November 19, 2020

Outline

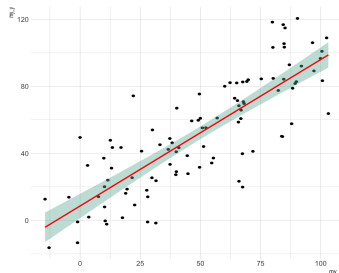
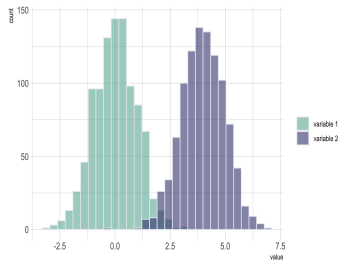
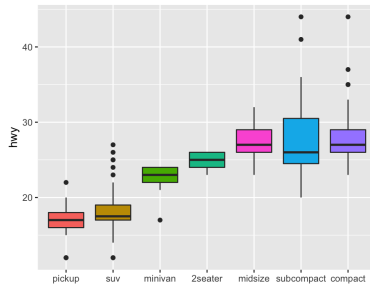
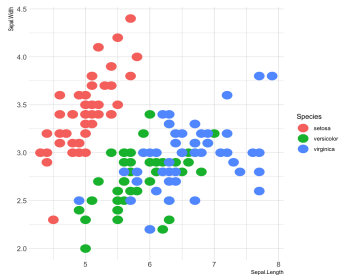
- 1 Learning Objectives
- 2 Basic Introduction to R
- 3 Getting Started
- 4 RStudio
- 5 Loading Data
- 6 Histogram
- 7 Scatterplot
- 8 Boxplot
- 9 Correlation
- 10 Linear Regression
- 11 Summary
- 12 Going Further

Learning Objectives

By the end of this lecture students will be able to...

- explain to a colleague what R is
- install and run R and RStudio
- load data and packages in R and RStudio
- describe what a histogram is and when you would use it
- describe what a scatterplot is and when you would use it
- describe what a boxplot is, what each horizontal line refers to, and when you would use it
- describe and identify outliers
- explain concepts of correlation and simple linear regression
- perform correlation and regression analysis using R and RStudio
- interpret results from correlation and regression
- create and customize histograms, scatterplots, boxplots, and linear models in R

Basic Introduction to R





What is R?



What is R?

- A language and environment for **statistical computing and graphics**



What is R?

- A language and environment for **statistical computing and graphics**
- A GNU Project - 100% free software (other examples include Linux operating systems)



What is R?

- A language and environment for **statistical computing and graphics**
- A GNU Project - 100% free software (other examples include Linux operating systems)
- Widely used among statisticians, data scientists, and scientists alike



What is R?

- A language and environment for **statistical computing and graphics**
- A GNU Project - 100% free software (other examples include Linux operating systems)
- Widely used among statisticians, data scientists, and scientists alike

Why use R?

- Free, well documented, helpful community online, runs on everything



What is R?

- A language and environment for **statistical computing and graphics**
- A GNU Project - 100% free software (other examples include Linux operating systems)
- Widely used among statisticians, data scientists, and scientists alike

Why use R?

- Free, well documented, helpful community online, runs on everything
- Wide variety of statistical techniques: linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering



What is R?

- A language and environment for **statistical computing and graphics**
- A GNU Project - 100% free software (other examples include Linux operating systems)
- Widely used among statisticians, data scientists, and scientists alike

Why use R?

- Free, well documented, helpful community online, runs on everything
- Wide variety of statistical techniques: linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering
- Highly extensible via *packages*



What is R?

- A language and environment for **statistical computing and graphics**
- A GNU Project - 100% free software (other examples include Linux operating systems)
- Widely used among statisticians, data scientists, and scientists alike

Why use R?

- Free, well documented, helpful community online, runs on everything
- Wide variety of statistical techniques: linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering
- Highly extensible via *packages*
- Well-designed publication-quality plots

Getting Started

Where to begin?

- Visit <https://www.r-project.org/>



[\[Home\]](#)

[Download](#)

[CRAN](#)

[R Project](#)

[About R](#)

[Logo](#)

[Contributors](#)

[What's New?](#)

[Reporting Bugs](#)

[Conferences](#)

[Search](#)

[Get Involved: Mailing Lists](#)

[Developer Pages](#)

[R Blog](#)

[R Foundation](#)

[Foundation](#)

[Board](#)

[Members](#)

[Donors](#)

[Donate](#)

[Help With R](#)

The R Project for Statistical Computing

Getting Started

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To [download R](#), please choose your preferred [CRAN mirror](#).

If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

News

- [R version 4.0.3 \(Bunny-Wunnies Freak Out\)](#) has been released on 2020-10-10.
- Thanks to the organisers of useR! 2020 for a successful online conference. Recorded tutorials and talks from the conference are available on the [R Consortium YouTube channel](#).
- [R version 3.6.3 \(Holding the Windsock\)](#) was released on 2020-02-29.
- You can support the R Foundation with a renewable subscription as a [supporting member](#)

News via Twitter



Getting Started

Where to begin?

- Visit <https://www.r-project.org/>



[Home]

Download

CRAN

R Project

About R

Logo

Contributors

What's New?

Reporting Bugs

Conferences

Search

Get Involved: Mailing Lists

Developer Pages

R Blog

R Foundation

Foundation

Board

Members

Donors

Donate

Help With R

The R Project for Statistical Computing

Getting Started

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To **download R**, please choose your preferred **CRAN mirror**.

If you have questions about R like how to download and install the software, or what the license terms are, please read our **answers to frequently asked questions** before you send an email.

News

- **R version 4.0.3 (Bunny-Wunnies Freak Out)** has been released on 2020-10-10.
- Thanks to the organisers of useR! 2020 for a successful online conference. Recorded tutorials and talks from the conference are available on the **R Consortium YouTube channel**.
- **R version 3.6.3 (Holding the Windsock)** was released on 2020-02-29.
- You can support the R Foundation with a renewable subscription as a supporting member

News via Twitter



The R Foundation

@R_Foundation

New R blog entry by Tomas Kalibera and Simon Urbanek: Will R work on Apple Silicon? <https://developer.r-project.org/blog/public/20...>



Nov 2, 2020

CRAN Mirrors

The Comprehensive R Archive Network is available at the following URLs, please choose a location close to you. Some statistics on the status of the mirrors can be found here: [main page](#), [windows release](#), [windows old release](#).

If you want to host a new mirror at your institution, please have a look at the [CRAN Mirror HOWTO](#).

0-Cloud

<https://cloud.r-project.org/>

Automatic redirection to servers worldwide, currently sponsored by Rstudio

Algeria

<https://cran.usbth.dz/>

University of Science and Technology Houari Boumediene

Argentina

<http://mirror.faglp.unlp.edu.ar/CRAN/>

Universidad Nacional de La Plata

Australia

<https://cran.csiro.au/>
<https://mirror.aarnet.edu.au/pub/CRAN/>
<https://cran.ms.unimelb.edu.au/>
<https://cran.curtin.edu.au/>

CSIRO
AARNET
School of Mathematics and Statistics, University of Melbourne
Curtin University

Austria

<https://cran.wu.ac.at/>

Wirtschaftsuniversität Wien

Belgium

<https://www.freestats.org/cran/>
<https://lib.ugent.be/CRAN/>

Patrick Wessa
Ghent University Library

Brazil

<https://ibg.gib.usp.br/mirrors/cran/>
<https://cran-r-csl.ufpb.br/>
<https://cran.fiocruz.br/>
<https://cps.fmvz.usp.br/CRAN/>
<https://biologie.esalq.usp.br/CRAN/>

Computational Biology Center at Universidade Estadual de Santa Cruz
Universidade Federal do Paraná
Oswaldo Cruz Foundation, Rio de Janeiro
University of São Paulo, São Paulo
University of São Paulo, Piracicaba

Bulgaria

<https://ftp.uni-sofia.bg/CRAN/>

Sofia University

Canada

<https://mirror.rcg.sfu.ca/mirror/CRAN/>
<https://mugq.ca/mirror/cran/>

Simon Fraser University, Burnaby
Manitoba Unix User Group

Getting Started

Where to begin?

- Download R for your operating system and install

The Comprehensive R Archive Network

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (2020-10-10, Bunny-Wunnies Freak Out) [R-4.0.3.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).
- Daily snapshots of current patched and development versions are [available here](#). Please read about [new features and bug fixes](#) before filing corresponding feature requests or bug reports.
- Source code of older versions of R is [available here](#).
- Contributed extension [packages](#)

Questions About R

- If you have questions about R like how to download and install the software,

Getting Started

Where to begin?

- Download R for your operating system and install

The Comprehensive R Archive Network

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (2020-10-10, Bunny-Wunnies Freak Out) [R-4.0.3.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).
- Daily snapshots of current patched and development versions are [available here](#). Please read about [new features](#) and [bug fixes](#) before filing corresponding feature requests or bug reports.
- Source code of older versions of R is [available here](#).
- Contributed extension [packages](#)

Questions About R

- If you have questions about R like how to download and install the software,

R for Windows

Subdirectories:

[base](#)

Binaries for base distribution. This is what you want to [install R for the first time](#).

[contrib](#)

Binaries of contributed CRAN packages (for R \geq 2.13.x; managed by Uwe Ligges). There is also information on [third party software](#) available for CRAN Windows services and corresponding environment and make variables.

[old contrib](#)

Binaries of contributed CRAN packages for outdated versions of R (for R $<$ 2.13.x; managed by Uwe Ligges).

[Rtools](#)

Tools to build R and R packages. This is what you want to build your own packages on Windows, or to build R itself.

Please do not submit binaries to CRAN. Package developers might want to contact Uwe Ligges directly in case of questions / suggestions related to Windows binaries.

You may also want to read the [R FAQ](#) and [R for Windows FAQ](#).

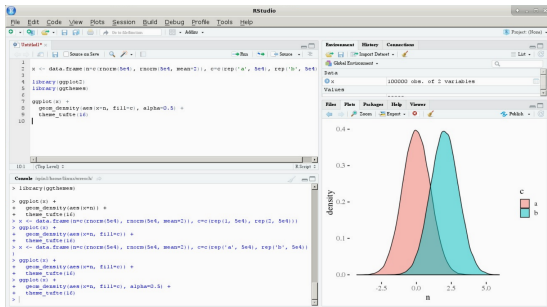
Note: CRAN does some checks on these binaries for viruses, but cannot give guarantees. Use the normal precautions with downloaded executables.

What is RStudio?

- RStudio is an integrated development environment (IDE) for R.

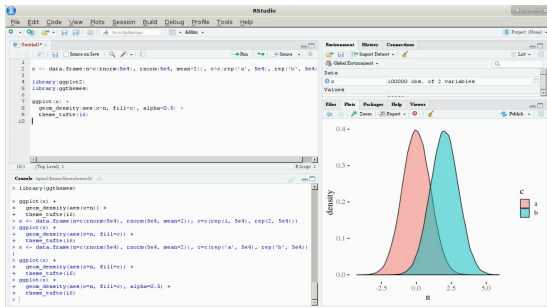
What is RStudio?

- RStudio is an integrated development environment (IDE) for R.
- Includes: console, syntax-highlighting editor, tools for plotting, history, debugging and workspace management



What is RStudio?

- RStudio is an integrated development environment (IDE) for R.
- Includes: console, syntax-highlighting editor, tools for plotting, history, debugging and workspace management
- Otherwise R is run through a command line interface (CLI)



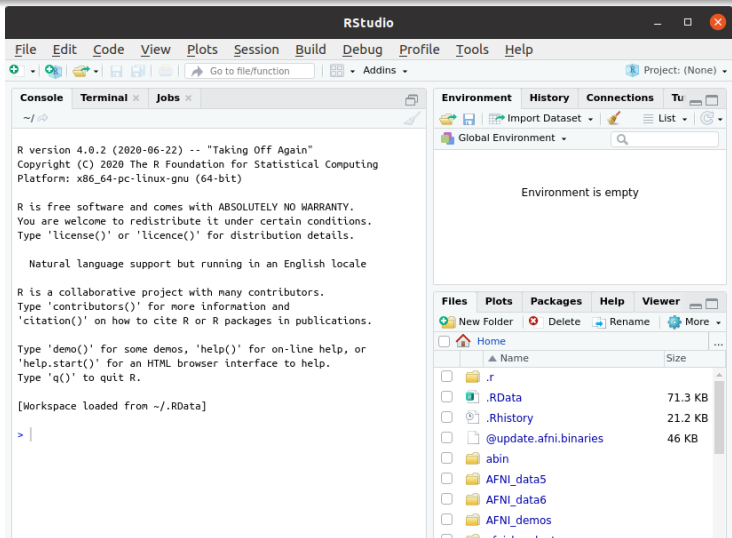
Installation

- Visit <https://rstudio.com/products/rstudio/download/>

	RStudio Desktop Open Source License	RStudio Desktop Pro Commercial License	RStudio Server Open Source License	RStudio Server Pro Commercial License
	Free	\$995 /year	Free	\$4,975 /year (5 Named Users)
	DOWNLOAD	BUY	DOWNLOAD	BUY
	Learn more	Learn more	Learn more	Evaluation Learn more
Integrated Tools for R	✓	✓	✓	✓
Priority Support		✓		✓
Access via Web Browser			✓	✓
RStudio Professional Drivers		✓		✓
Connect to RStudio Server Pro remotely		✓		
Enterprise Security				✓
Project Sharing				✓
Manage Multiple R Sessions & Versions				✓

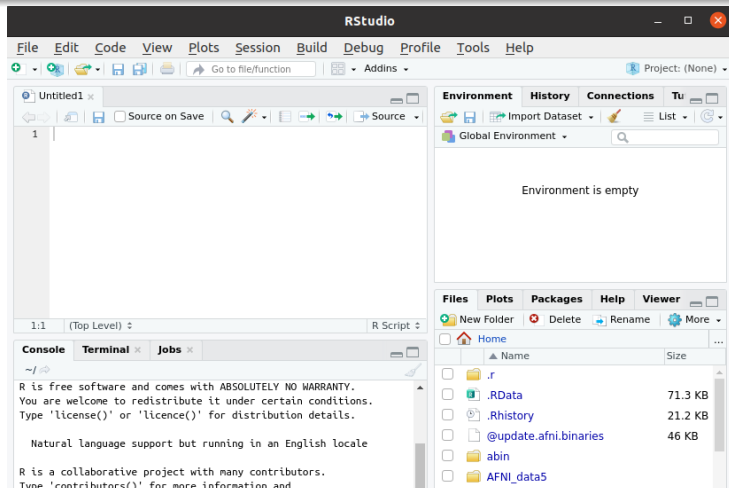
Getting Started

- Let's run Rstudio and take a look

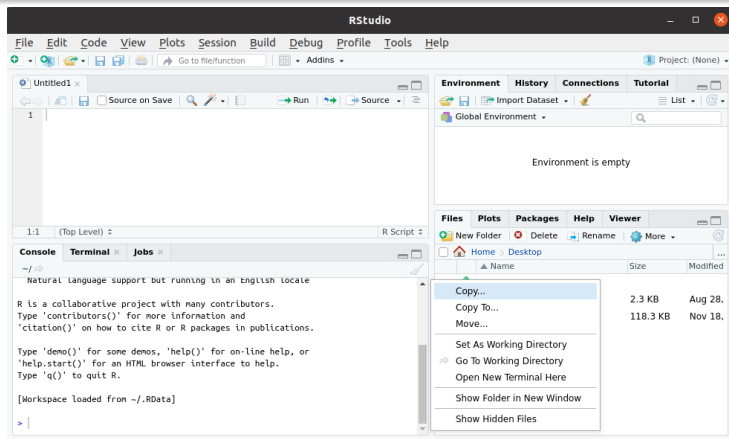


Getting Started

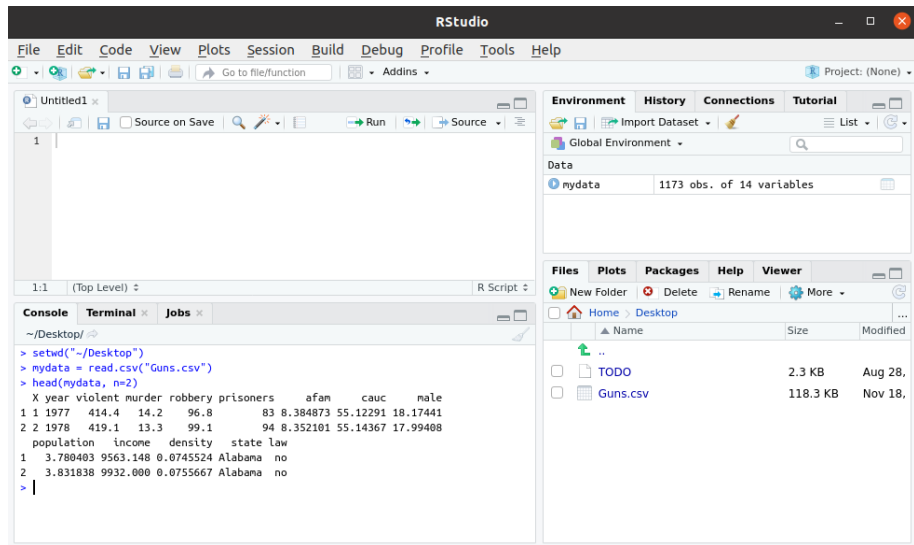
- Let's open up a scripting window
- File > New File > RScript



- I'm going to grab some public data from <https://vincentarelbundock.github.io/Rdatasets/datasets.html>
- Place file in Desktop, and change the working directory



Loading Data

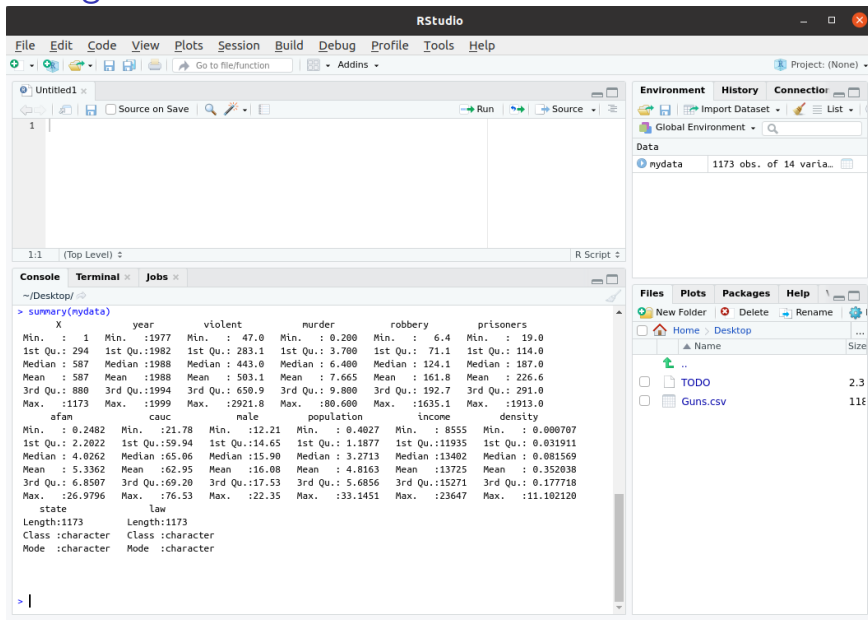


The screenshot shows the RStudio interface with the following components:

- Menu Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Toolbar:** Includes icons for opening files, saving, running code, and other standard RStudio functions.
- Source Editor:** Contains a script named 'Untitled1.R' with the following R code:

```
1  
  
1:1 (Top Level) ± R Script ±  
  
Console Terminal Jobs  
~/Desktop/ ↗  
> setwd("~/Desktop")  
> mydata = read.csv("Guns.csv")  
> head(mydata, n=2)  
  X year violent murder robbery prisoners      afam      cauc      male  
1 1 1977   414.4    14.2    96.8          83 8.384873 55.12291 18.17441  
2 2 1978   419.1    13.3    99.1          94 8.352101 55.14367 17.99408  
  population income density state law  
1 3.780403 9563.148 0.0745524 Alabama no  
2 3.831838 9932.000 0.0755667 Alabama no  
> |
```
- Environment Panel:** Shows the 'Global Environment' with a data object 'mydata' containing 1173 observations and 14 variables.
- Files Panel:** Displays the file structure of the current directory (~/Desktop/), showing files like 'TODO' (2.3 KB) and 'Guns.csv' (118.3 KB).

Loading Data



The screenshot shows the RStudio environment with the following components:

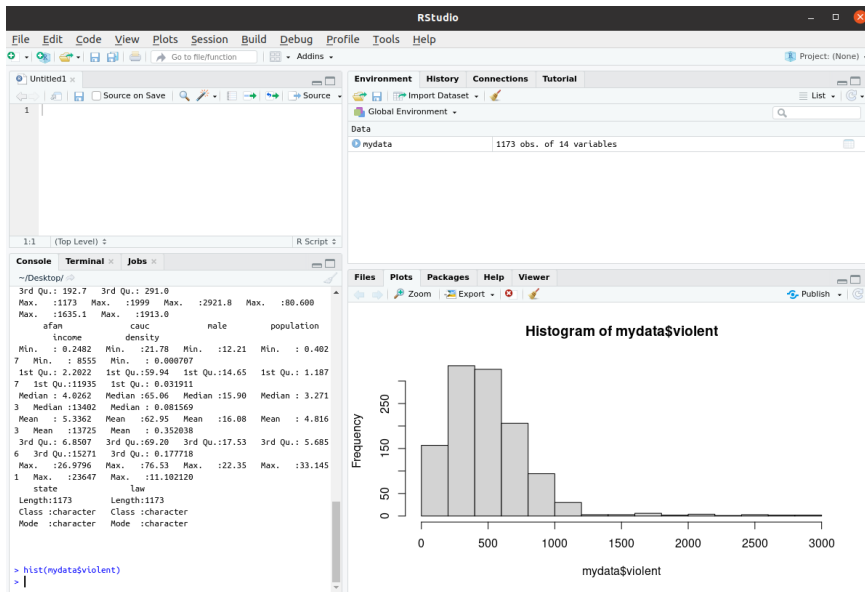
- Source Editor:** Contains a single line of code: `1`.
- Console:** Displays the output of the `summary(mydata)` command. The output shows summary statistics for variables: X, year, violent, murder, robbery, prisoners, afan, cauc, male, population, income, density, state, and law. Each variable has its minimum, 1st quartile, median, mean, 3rd quartile, and maximum value displayed.
- Environment:** Shows the Global Environment with a variable named `nydata` containing 1173 observations of 14 variables.
- Files:** Shows the file explorer with a folder named `TODO` and a file named `Guns.csv`.

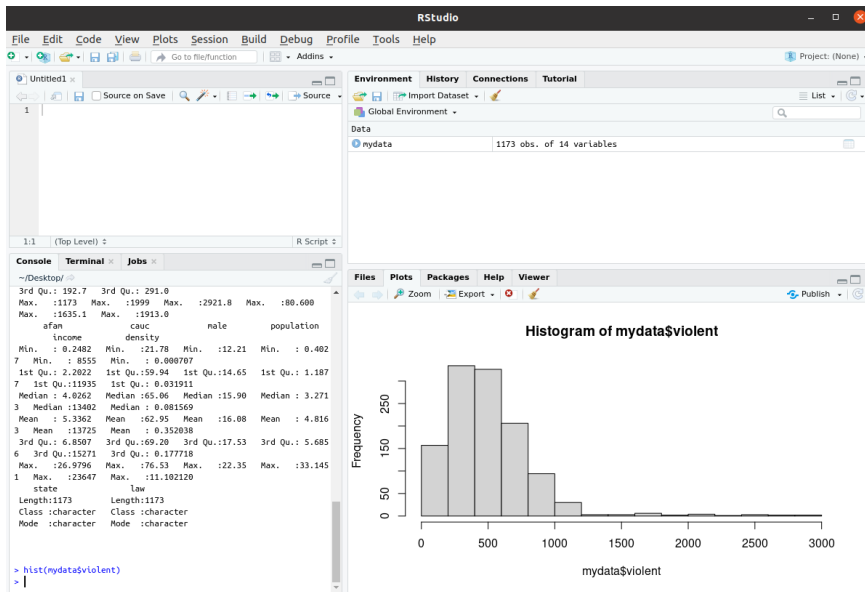
Console Output:

```
> summary(mydata)
  X      year      violent      murder      robbery      prisoners
Min. : 1   Min. :1977   Min. : 47.0   Min. : 0.200   Min. : 6.4   Min. : 19.0
1st Qu.: 294 1st Qu.:1982 1st Qu.: 283.1 1st Qu.: 3.700 1st Qu.: 71.1 1st Qu.: 114.0
Median : 587 Median :1988 Median : 443.0 Median : 6.400 Median : 124.1 Median : 187.0
Mean : 587 Mean :1988 Mean : 503.1 Mean : 7.665 Mean : 161.8 Mean : 226.6
3rd Qu.: 880 3rd Qu.:1994 3rd Qu.: 650.9 3rd Qu.: 9.800 3rd Qu.: 192.7 3rd Qu.: 291.0
Max. :1173 Max. :1999 Max. :2921.8 Max. :80.600 Max. :1635.1 Max. :1913.0

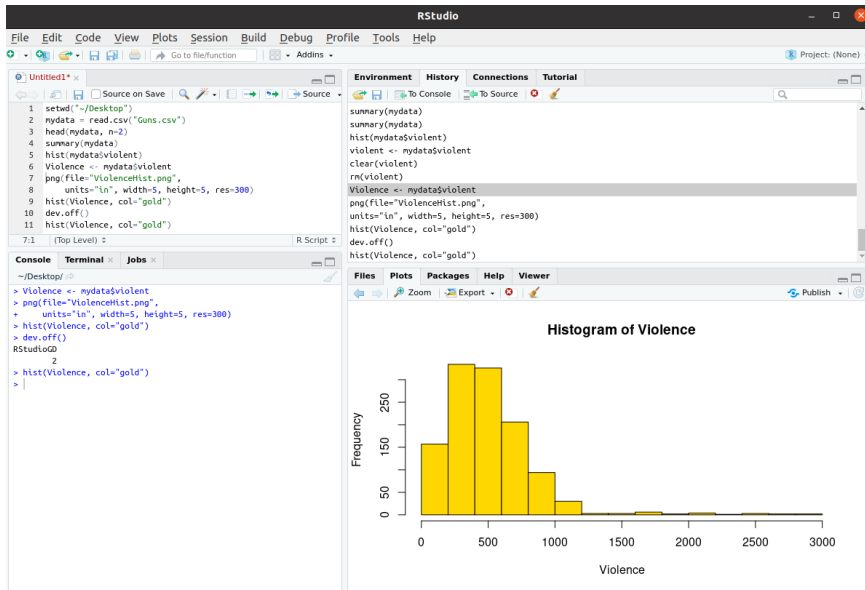
  afan      cauc      male      population      income      density
Min. : 0.2482 Min. :21.78 Min. :12.21 Min. : 0.4027 Min. : 8555 Min. : 0.000707
1st Qu.: 2.2022 1st Qu.:59.94 1st Qu.:14.65 1st Qu.: 1.1877 1st Qu.:11935 1st Qu.: 0.031911
Median : 4.0262 Median :65.06 Median :15.90 Median : 3.2713 Median :13402 Median : 0.081569
Mean : 5.3362 Mean :62.95 Mean :16.08 Mean : 4.8163 Mean :13725 Mean : 0.352038
3rd Qu.: 6.8507 3rd Qu.:69.20 3rd Qu.:17.53 3rd Qu.: 5.6856 3rd Qu.:15271 3rd Qu.: 0.177718
Max. :26.9796 Max. :76.53 Max. :22.35 Max. :33.1451 Max. :23647 Max. :11.102120

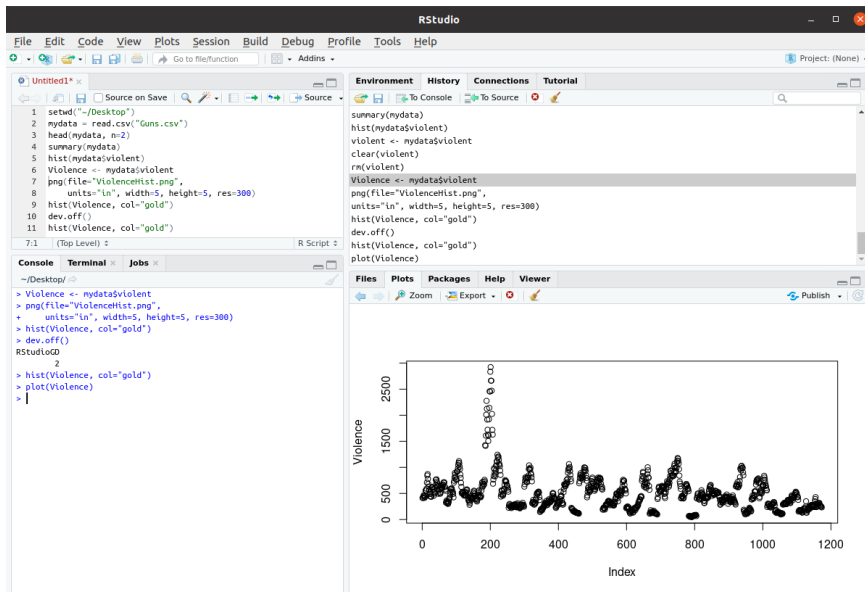
  state      law
Length:1173 Length:1173
Class :character Class :character
Mode :character Mode :character
```

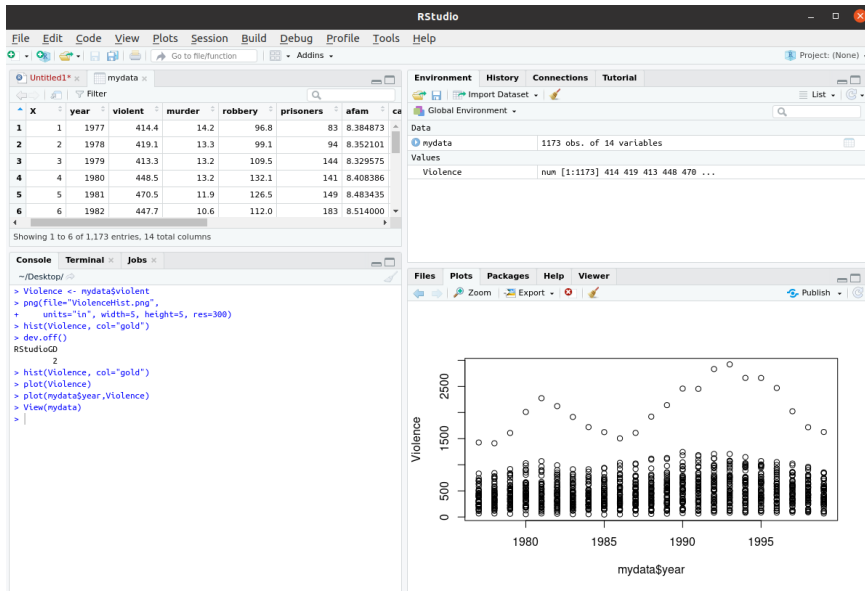


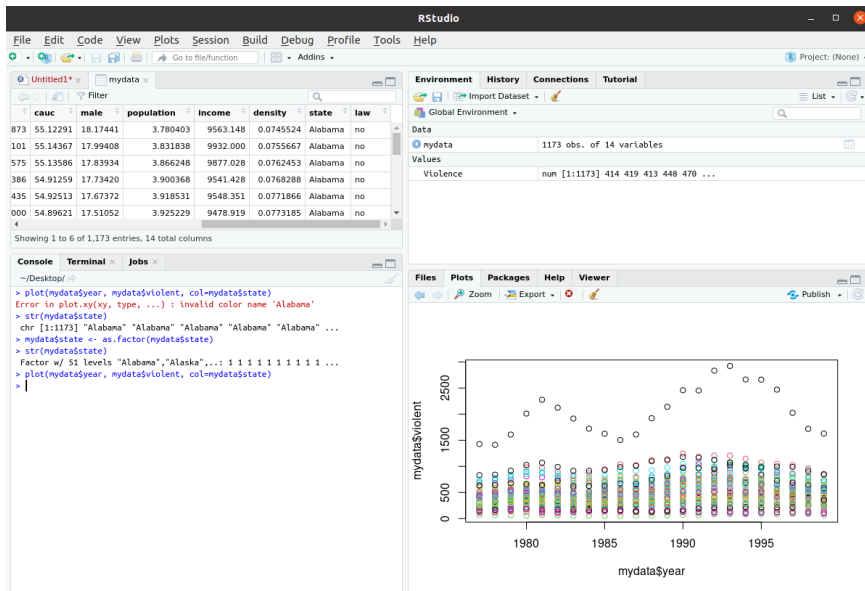


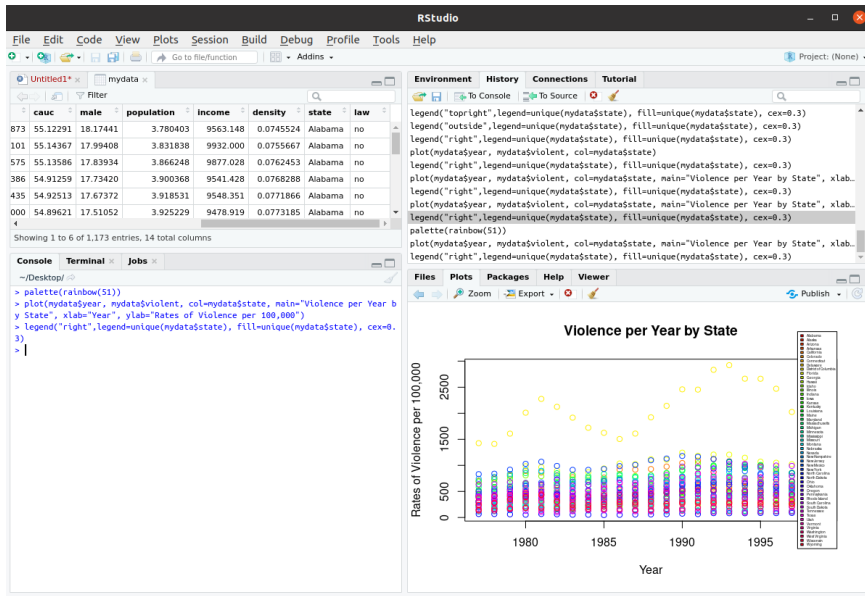
Histogram

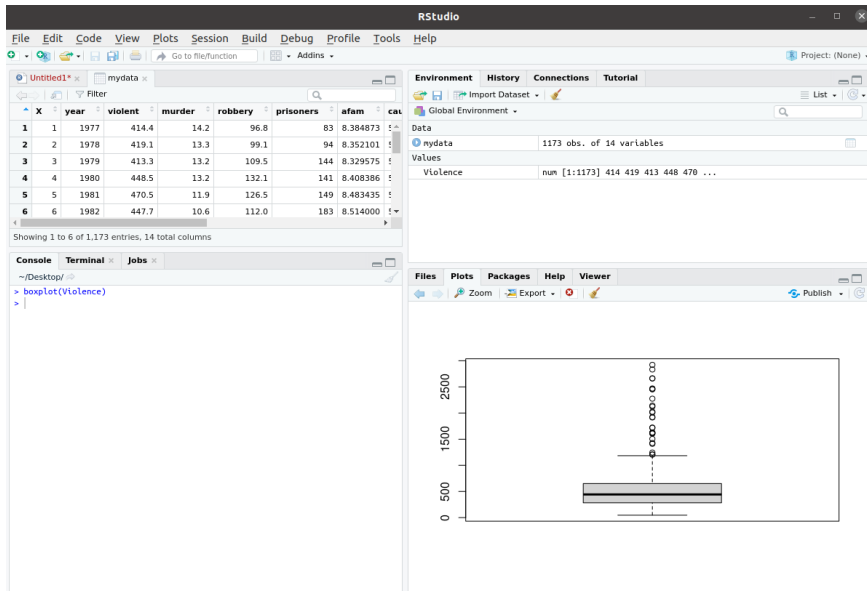








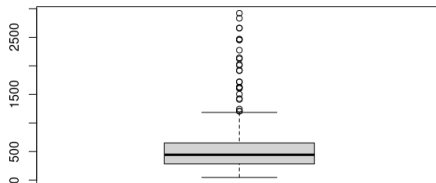




Boxplot

Boxplot Elements

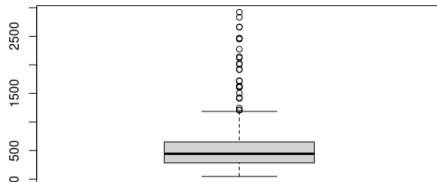
- Median (Q_2): middle value of the dataset



Boxplot

Boxplot Elements

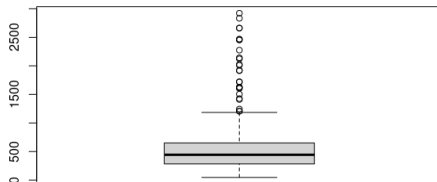
- Median (Q_2): middle value of the dataset
- First/lower quartile (Q_1): median of the lower half of the values



Boxplot

Boxplot Elements

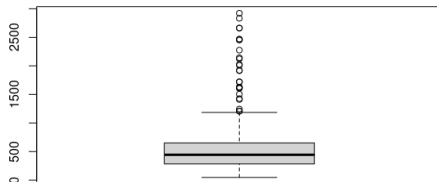
- Median (Q_2): middle value of the dataset
- First/lower quartile (Q_1): median of the lower half of the values
- Third / upper quartile (Q_3): median of the upper half of the values



Boxplot

Boxplot Elements

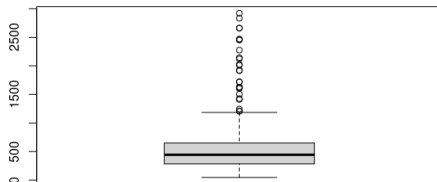
- Median (Q_2): middle value of the dataset
- First/lower quartile (Q_1): median of the lower half of the values
- Third / upper quartile (Q_3): median of the upper half of the values
- The interquartile range, or IQR, is $Q_3 - Q_1$



Boxplot

Boxplot Elements

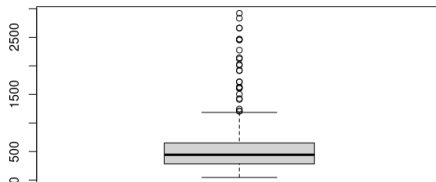
- Median (Q_2): middle value of the dataset
- First/lower quartile (Q_1): median of the lower half of the values
- Third / upper quartile (Q_3): median of the upper half of the values
- The interquartile range, or IQR, is $Q_3 - Q_1$
- Whiskers are the max value that lies within $1.5 \times \text{IQR}$ of first and third quartiles

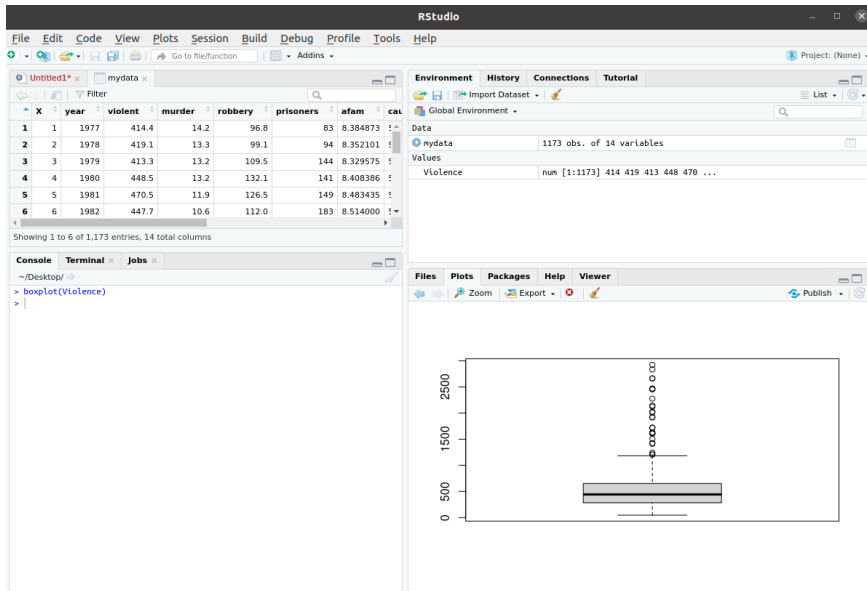


Boxplot

Boxplot Elements

- Median (Q_2): middle value of the dataset
- First/lower quartile (Q_1): median of the lower half of the values
- Third / upper quartile (Q_3): median of the upper half of the values
- The interquartile range, or IQR, is $Q_3 - Q_1$
- Whiskers are the max value that lies within $1.5 \times \text{IQR}$ of first and third quartiles
- Outliers are circles that don't lie within the box or whiskers





Packages

What are packages?

- Packages are sets of *functions* developed by the community

Packages

What are packages?

- Packages are sets of *functions* developed by the community
- They increase the power of R by improving existing base R functionalities, or by adding new ones

Packages

What are packages?

- Packages are sets of *functions* developed by the community
- They increase the power of R by improving existing base R functionalities, or by adding new ones
- Can be loaded with `library(packagename)`

Packages

What are packages?

- Packages are sets of *functions* developed by the community
- They increase the power of R by improving existing base R functionalities, or by adding new ones
- Can be loaded with `library(packagename)`
- If you don't have the package yet, you can install it with `install.packages("packagename")`

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function
Addins
Project: (None)

Untitled1*

mydata

stats

stats

Filter

Item	group1	vars	n	mean	sd	median
X16	6	Colorado	1	23	482.78696	63.53352
X17	7	Connecticut	1	23	420.67391	65.36586
X18	8	Delaware	1	23	564.35652	122.20878
X19	9	District of Columbia	1	23	2048.97826	466.43853
X110	10	Florida	1	23	999.23478	151.54791
X111	11	Georgia	1	23	595.57826	99.63111

Showing 6 to 11 of 51 entries, 15 total columns

Environment History Connections Tutorial

Global Environment

Data	
mydata	1173 obs. of 14 variables
stats	51 obs. of 15 variables

Values	
outliers	2048.97826086957
Violence	num [1:1173] 414 419 413 448 470 ...

Files Plots Packages Help Viewer

Zoom Export

Publish

Console Terminal Jobs

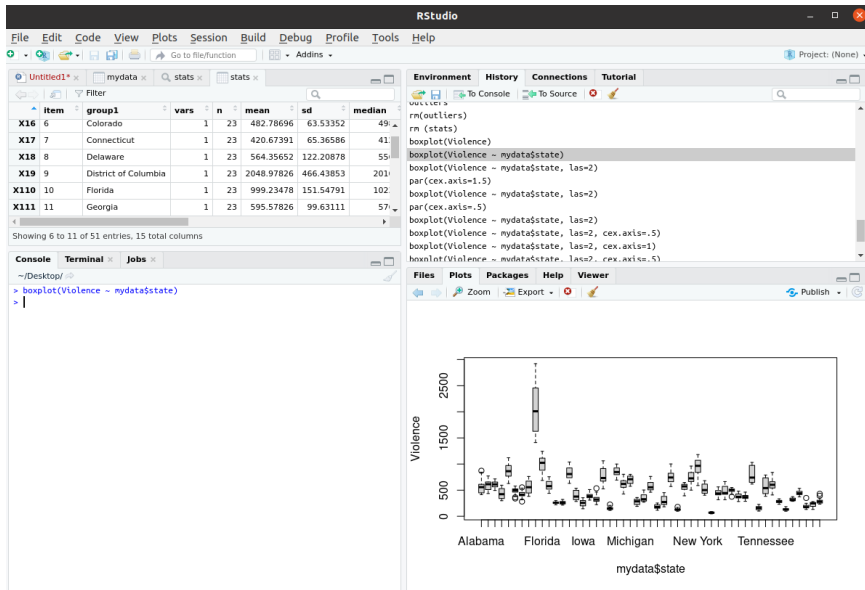
```

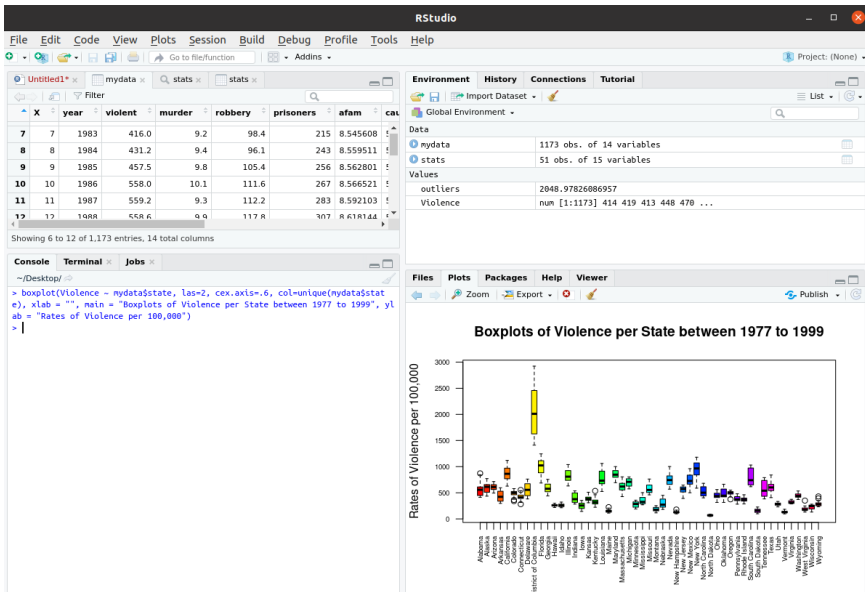
~/Desktop/
> library(psych)
> stats <- describeBy(mydata$violent, mydata$state, mat=TRUE)
> stats$mean
[1] 558.17391 596.79565 608.09565 438.40000 877.48261 482.78696
[7] 428.67391 564.35652 2048.97826 999.23478 595.57826 259.43478
[13] 262.64348 828.16522 399.86957 247.60000 391.70870 334.69565
[19] 778.43913 154.04783 853.63478 616.71739 695.21384 277.26522
[25] 355.52609 582.65217 183.47826 300.73478 755.03478 130.84783
[31] 552.82174 743.86957 941.31739 526.39565 68.00435 445.45217
[37] 491.58696 496.33913 383.45652 369.40000 800.26522 157.36087
[43] 580.68261 612.65652 283.50870 133.20870 322.95217 447.72174
[49] 188.76087 225.28261 287.50000

> boxplot(stats$mean)
> outliers = boxplot(stats$mean, plot=FALSE)$out
> stats[stats$mean %in% outliers,]
  item group1 vars n mean sd median trimmed
X19 9 District of Columbia 1 23 2048.978 466.4385 2010.6 2028.089
mad min max range skew kurtosis se
X19 593.4848 1411.7 2921.8 1510.1 0.3306733 -1.260257 97.25916

> View(stats)
>

```





tidyverse and ggplot2

tidyverse

- tidyverse is a collection of packages for tidy data

tidyverse and ggplot2

tidyverse

- tidyverse is a collection of packages for tidy data
- The core packages are ggplot2, dplyr, tidyr, readr, purrr, tibble, stringr, and forcats

tidyverse and ggplot2

tidyverse

- tidyverse is a collection of packages for tidy data
- The core packages are ggplot2, dplyr, tidyr, readr, purrr, tibble, stringr, and forcats
- As of November 2018, the packages make up 5 out of the top 10 most downloaded R packages

tidyverse and ggplot2

tidyverse

- tidyverse is a collection of packages for tidy data
- The core packages are ggplot2, dplyr, tidyr, readr, purrr, tibble, stringr, and forcats
- As of November 2018, the packages make up 5 out of the top 10 most downloaded R packages
- `install.packages("tidyverse")`

tidyverse and ggplot2

tidyverse

- tidyverse is a collection of packages for tidy data
- The core packages are ggplot2, dplyr, tidyr, readr, purrr, tibble, stringr, and forcats
- As of November 2018, the packages make up 5 out of the top 10 most downloaded R packages
- `install.packages("tidyverse")`

ggplot2

- ggplot2 is a data visualization package

tidyverse and ggplot2

tidyverse

- tidyverse is a collection of packages for tidy data
- The core packages are ggplot2, dplyr, tidyr, readr, purrr, tibble, stringr, and forcats
- As of November 2018, the packages make up 5 out of the top 10 most downloaded R packages
- `install.packages("tidyverse")`

ggplot2

- ggplot2 is a data visualization package
- follows 'Grammar of Graphics': a general scheme for data visualization which breaks up graphs into semantic components such as scales and layers

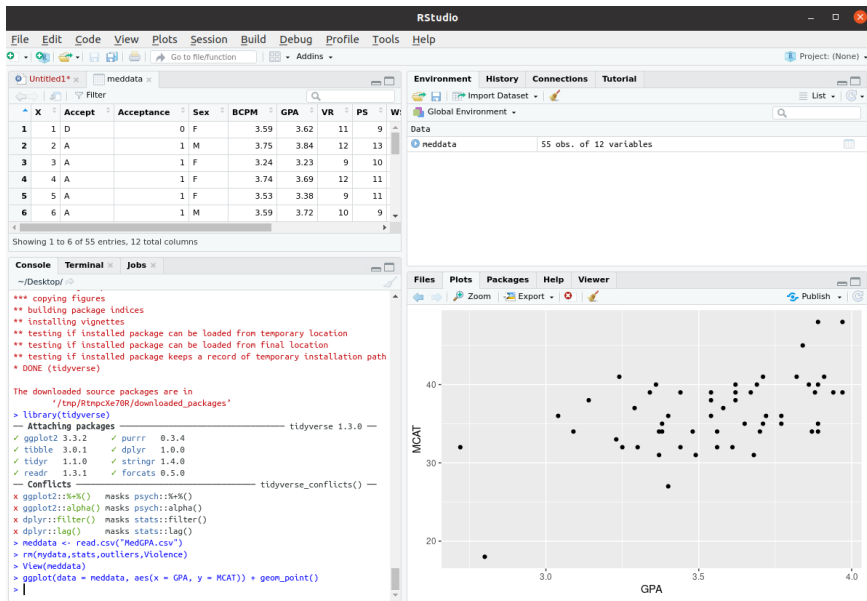
tidyverse and ggplot2

tidyverse

- tidyverse is a collection of packages for tidy data
- The core packages are ggplot2, dplyr, tidyr, readr, purrr, tibble, stringr, and forcats
- As of November 2018, the packages make up 5 out of the top 10 most downloaded R packages
- `install.packages("tidyverse")`

ggplot2

- ggplot2 is a data visualization package
- follows 'Grammar of Graphics': a general scheme for data visualization which breaks up graphs into semantic components such as scales and layers
- `library(tidyverse)`



Correlation

What is correlation?

- Correlation is a statistical measure that suggests the level of linear dependence between two variables, that occur in pair

Correlation

What is correlation?

- Correlation is a statistical measure that suggests the level of linear dependence between two variables, that occur in pair
- Correlation can take values between -1 to $+1$

Correlation

What is correlation?

- Correlation is a statistical measure that suggests the level of linear dependence between two variables, that occur in pair
- Correlation can take values between -1 to $+1$
- If for every time an independent variable increases and the dependent variable increases, we say this is a positive correlation, and will have a value closer to 1

Correlation

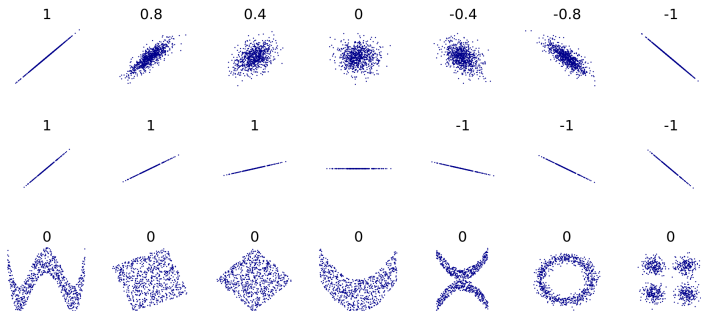
What is correlation?

- Correlation is a statistical measure that suggests the level of linear dependence between two variables, that occur in pair
- Correlation can take values between -1 to +1
- If for every time an independent variable increases and the dependent variable increases, we say this is a positive correlation, and will have a value closer to 1
- If for every time an independent variable increases and the dependent variable decreases, we say this is a negative correlation, and will have a value closer to -1

Correlation

Pearson Correlation:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 (y_i - \bar{y})^2}}$$



Correlation

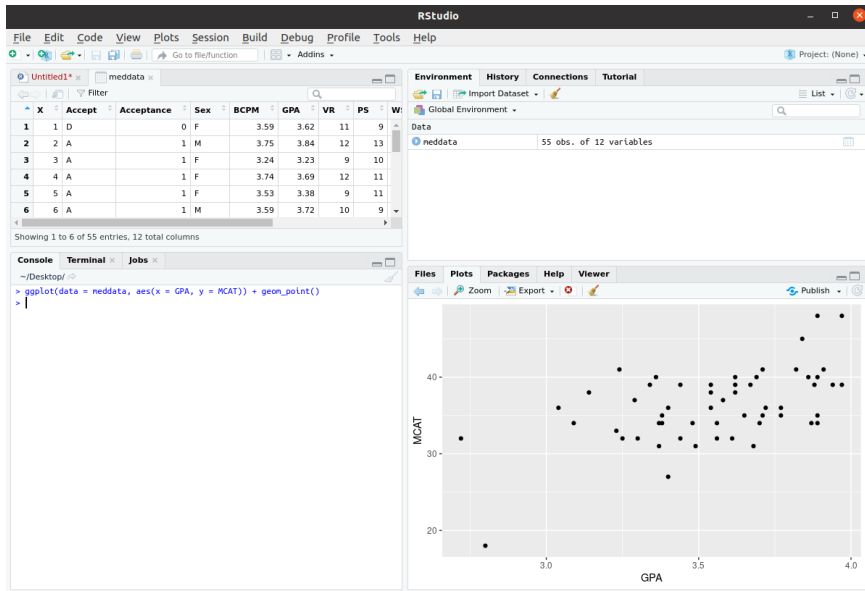
Two Assumptions

- The two variables are normally distributed. We can test this assumption using a histogram

Correlation

Two Assumptions

- The two variables are normally distributed. We can test this assumption using a histogram
- The relationship between the two variables is linear. If this relationship is found to be curved, etc. we need to use another correlation test. We can test this assumption by examining the scatterplot between the two variables.



RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins Project: (None)

Untitled1* meddata

X	Accept	Acceptance	Sex	BCPM	GPA	VR	PS	W
1	1	D	0	F	3.59	3.62	11	9
2	2	A	1	M	3.75	3.84	12	13
3	3	A	1	F	3.24	3.23	9	10
4	4	A	1	F	3.74	3.69	12	11
5	5	A	1	F	3.53	3.38	9	11
6	6	A	1	M	3.59	3.72	10	9

Showing 1 to 6 of 55 entries, 12 total columns

Environment History Connections Tutorial

Global Environment

Data

- meddata 55 obs. of 12 variables
- plot.data num [1:2, 1:55] 3.62 38 3.84 45 3.23 33 3.69 40 3.38 35 ...
- plot1 List of 9
- plot2 List of 9

Files Plots Packages Help Viewer

Zoom Export Publish

MCAT

GPA

X

X

```

~/Desktop/
** R
** inst
** byte-compile and prepare package for lazy loading
** help
*** installing help indices
** building package indices
** installing vignettes
** testing if installed package can be loaded from temporary location
** testing if installed package can be loaded from final location
** testing if installed package keeps a record of temporary installation path
* DONE (gridExtra)

The downloaded source packages are in
  '/tmp/RtmpKc76R/downloaded_packages'
> library(gridExtra)

Attaching package: 'gridExtra'

The following object is masked from 'package:dplyr':

  combine

> plot1 <- ggplot(data = meddata, aes(x = "", y = MCAT)) + geom_boxplot()
> plot2 <- ggplot(data = meddata, aes(x = "", y = GPA)) + geom_boxplot()
> grid.arrange(plot1, plot2, ncol=2)
> |

```

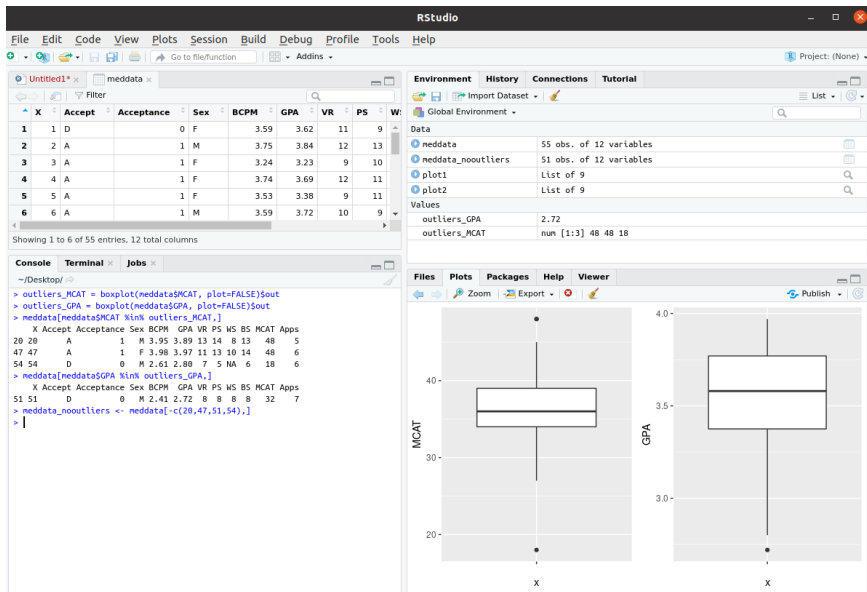

A Word of Warning

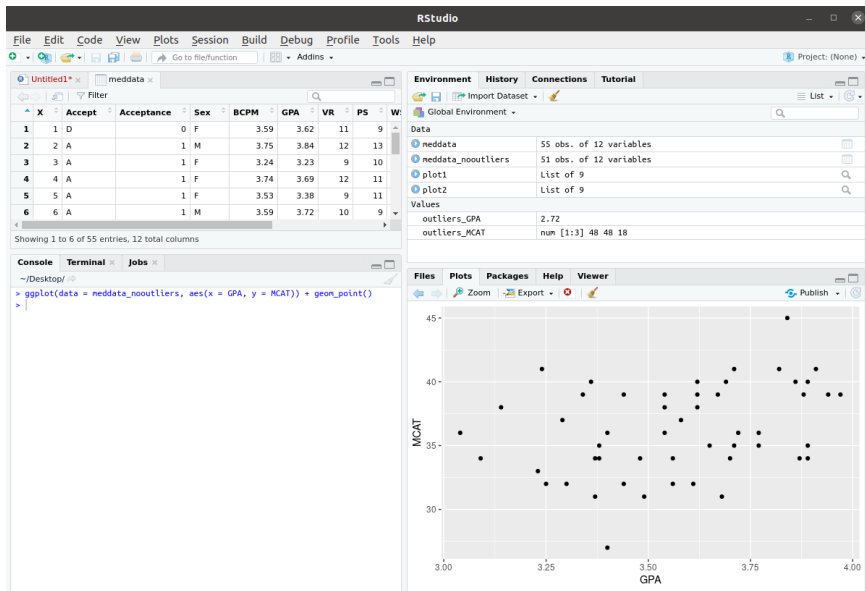
To Remove or Keep Outliers?

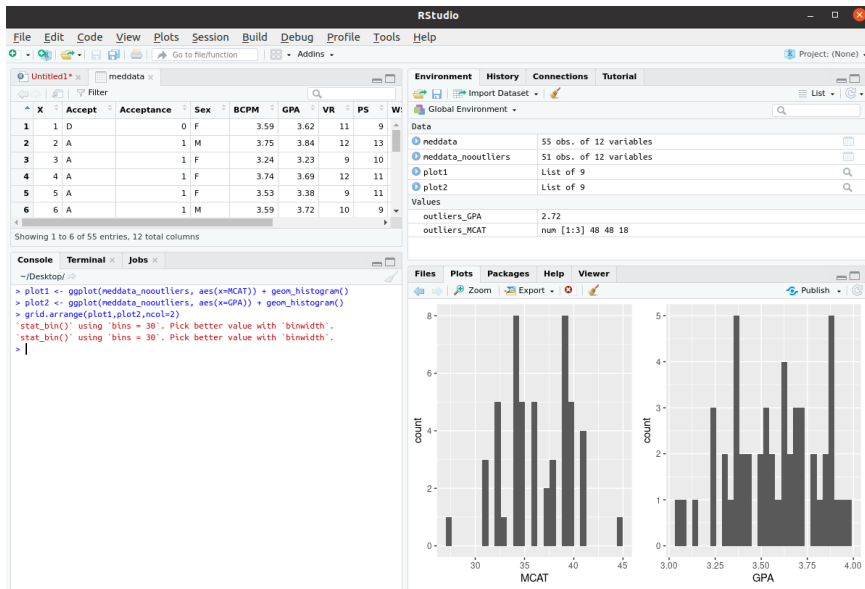
There are generally three reasons why we may have an outlier:

- ❶ There was a recording error
- ❷ Someone made a fundamental mistake collecting the data
- ❸ The data point is perfectly valid, in which case the model we are trying to use cannot account for the behavior.

“Sometimes we throw out perfectly good data when we should be throwing out questionable models.”







```
> cor(meddata_nooutliers$GPA,meddata_nooutliers$MCAT, method="pearson")  
[1] 0.3350624  
> cor.test(meddata_nooutliers$GPA,meddata_nooutliers$MCAT, method="pearson")
```

Pearson's product-moment correlation

```
data: meddata_nooutliers$GPA and meddata_nooutliers$MCAT  
t = 2.4893, df = 49, p-value = 0.01624  
alternative hypothesis: true correlation is not equal to 0  
95 percent confidence interval:  
 0.06552957 0.55902680  
sample estimates:  
      cor  
0.3350624
```

```
> |
```

Linear Regression

What is linear regression?

- linear regression is a linear approach to modelling the relationship between a dependent variable and an explanatory variable

Linear Regression

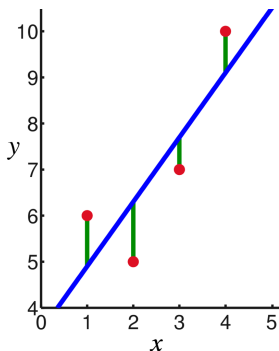
What is linear regression?

- linear regression is a linear approach to modelling the relationship between a dependent variable and an explanatory variable
- can be used to predict values of one variable based on the other

Linear Regression

What is linear regression?

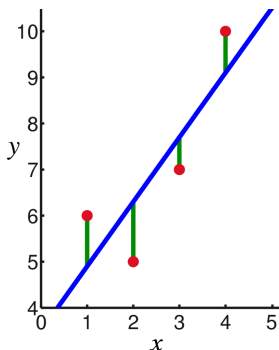
- linear regression is a linear approach to modelling the relationship between a dependent variable and an explanatory variable
- can be used to predict values of one variable based on the other
- the observations (red) are assumed to be the result of random deviations (green) from an underlying relationship (blue)



Linear Regression

What is linear regression?

- linear regression is a linear approach to modelling the relationship between a dependent variable and an explanatory variable
- can be used to predict values of one variable based on the other
- the observations (red) are assumed to be the result of random deviations (green) from an underlying relationship (blue)



Ordinary Least Squares

- minimize the sum of the squares of the differences (green) between the observed dependent variable (red) in the given dataset and those predicted by the linear function (blue).

$$y_i = \alpha + \beta x_i + \epsilon$$

Linear Regression

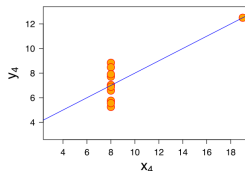
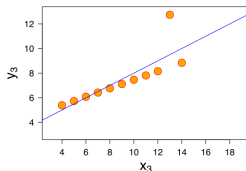
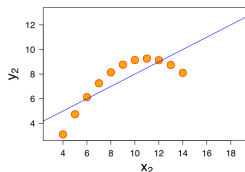
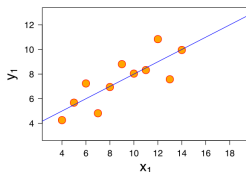
Warning

- just because you can fit a data, doesn't mean you should

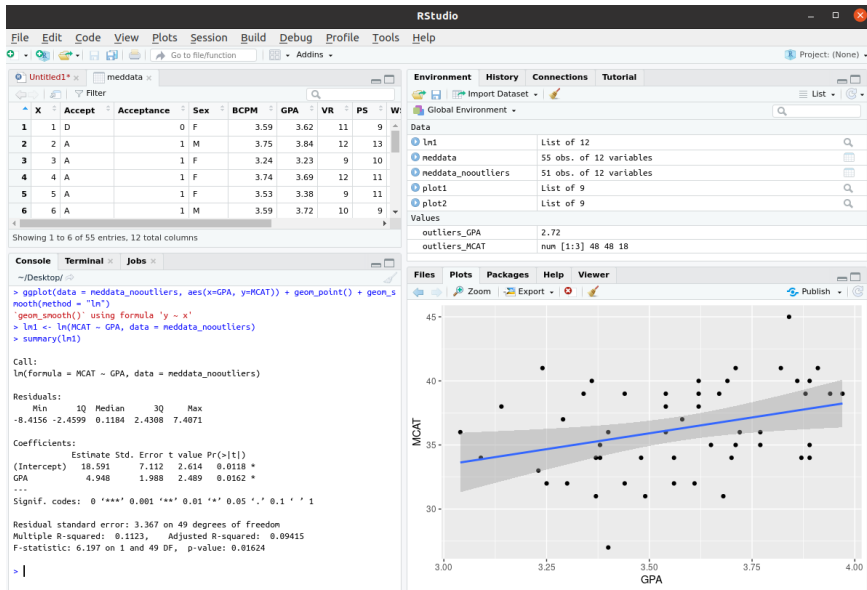
Linear Regression

Warning

- just because you can fit a data, doesn't mean you should
- the data below have approximately the same linear regression line (as well as nearly identical means, standard deviations, and correlations)



Linear Regression



Correlation and Linear Regression

Summary

- Pearson's correlation gives a value r , which tells us the level of linear dependence between two variables (-1 to $+1$)

Correlation and Linear Regression

Summary

- Pearson's correlation gives a value r , which tells us the level of linear dependence between two variables (-1 to $+1$)
- Based on two assumptions: data are normally distributed without any outliers; the relationship between the variables is linear

Correlation and Linear Regression

Summary

- Pearson's correlation gives a value r , which tells us the level of linear dependence between two variables (-1 to $+1$)
- Based on two assumptions: data are normally distributed without any outliers; the relationship between the variables is linear
- Linear Regression models the linear correlation between an independent variable and a dependent one

Correlation and Linear Regression

Summary

- Pearson's correlation gives a value r , which tells us the level of linear dependence between two variables (-1 to $+1$)
- Based on two assumptions: data are normally distributed without any outliers; the relationship between the variables is linear
- Linear Regression models the linear correlation between an independent variable and a dependent one
- Ordinary Least Squares is one method of determining the linear regression, where it minimizes the sum of the squares of the differences between the observed dependent variable and those predicted by the linear function.

Correlation and Linear Regression

Summary

- Ultimately we end up with a linear model:

$$y_i = \alpha + \beta x_i + \epsilon$$

Correlation and Linear Regression

Summary

- Ultimately we end up with a linear model:

$$y_i = \alpha + \beta x_i + \epsilon$$

- And an r-squared, which tells us how much of the variation of the response variable is explained by our predictors, and not by error

Summary

Today we learned...

- what R and RStudio are, how to download and install them, and how to load data and packages

Summary

Today we learned...

- what R and RStudio are, how to download and install them, and how to load data and packages
- what a histogram is and how to use it to get a visual idea of how a single set of our data values are distributed

Summary

Today we learned...

- what R and RStudio are, how to download and install them, and how to load data and packages
- what a histogram is and how to use it to get a visual idea of how a single set of our data values are distributed
- what a scatterplot is and how to use it to visualize two sets of our data and visually see how they are potentially linked

Summary

Today we learned...

- what R and RStudio are, how to download and install them, and how to load data and packages
- what a histogram is and how to use it to get a visual idea of how a single set of our data values are distributed
- what a scatterplot is and how to use it to visualize two sets of our data and visually see how they are potentially linked
- what a boxplot is and how to use it to visualize our data

Summary

Today we learned...

- what R and RStudio are, how to download and install them, and how to load data and packages
- what a histogram is and how to use it to get a visual idea of how a single set of our data values are distributed
- what a scatterplot is and how to use it to visualize two sets of our data and visually see how they are potentially linked
- what a boxplot is and how to use it to visualize our data
- what an outlier is and how to identify and remove it in R
- about Pearson correlation and r

Summary

Today we learned...

- what R and RStudio are, how to download and install them, and how to load data and packages
- what a histogram is and how to use it to get a visual idea of how a single set of our data values are distributed
- what a scatterplot is and how to use it to visualize two sets of our data and visually see how they are potentially linked
- what a boxplot is and how to use it to visualize our data
- what an outlier is and how to identify and remove it in R
- about Pearson correlation and r
- about linear regression, ordinary least squares, a linear regression model, and r^2
- how to plot all of these using base R packages, and with ggplot2

Going Further

Additional Resources

- Getting Started with R and RStudio
<https://www.youtube.com/watch?v=1VKMsawJu8w>
- Programming with R - Tutorial
<http://swcarpentry.github.io/r-novice-inflammation/>
- R for Reproducible Scientific Analysis - Tutorial
<http://swcarpentry.github.io/r-novice-gapminder/>
- Scatterplots, Boxplots, Histograms in R
<http://becomingvisual.com/rfundamentals/statistical-graphs.html>
- Linear Regression in R
<https://www.datacamp.com/community/tutorials/linear-regression-R>