

1.Welcome to PSTAT 194TR!

Transfer exploration seminar: Statistics and Data Science

Dr. Uma Ravat

PSTAT 194TR

0.1 Welcome!

Plan for today

1. What is Data Science?
2. Your Turn at Data Science
3. (**R Toolkit**) Rstudio, Rmarkdown

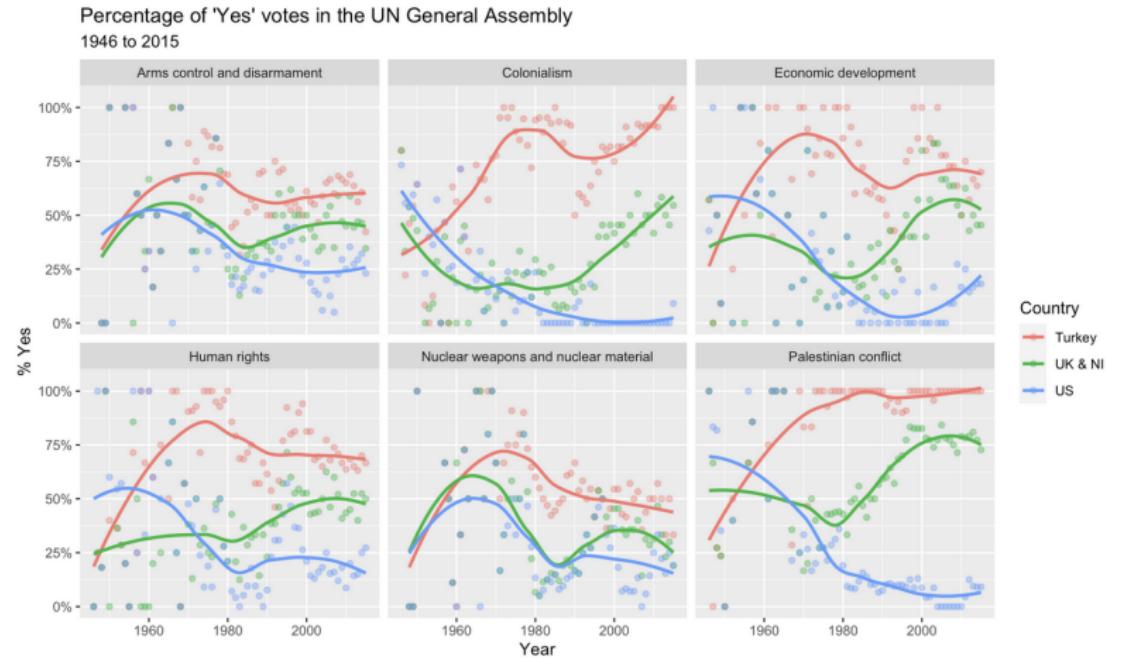
What is Data Science?

- Data science is an exciting discipline that allows you to look at raw data and transform it into understanding, insight, and knowledge that you can act on.



Your data science tasks result in...

1. A Visualization



2. Report/Website

UN Votes

Your name here
2023-08-24

Introduction

How do nations countries vote in the United Nations General Assembly. How have their voting patterns evolved throughout time, and how can we identify, on any given year, which countries are they more certain about? Answering these questions at a high level is the focus of this analysis.

Packages

We will use the `dplyr`, `tidyverse`, `tidytext`, and `gridExtra` packages for the data wrangling and visualization, and the `grid` package for interactive display of tabular output.

`library(tidyverse)`
`library(dplyr)`
`library(tidytext)`
`library(grid)`

Data

The data we're using originally come from the `unvotes` package, but it's been modified a bit by joining the various data frames provided in the package to help you get started with the analysis.

```
country <- read_csv("https://raw.githubusercontent.com/rfordatascience/tidytuesday/master/2023/2023-08-01/unvotes.csv")
```

UN voting patterns

Let's create a data visualization that depicts how the voting record of the UK & US changed over time on a variety of issues, and compared it to two other countries: US and Turkey.

Plotting the data is a bit of a challenge because it's ordered by changing which countries the code above `filter`s to. Note that the country name should be specified and capitalized exactly the same way as it appears in the data. See the [Appendix](#) for a list of the countries in the data.

```
unvotes %>%  
  filter(country %in% c("US", "UK", "Turkey")) %>%  
  group_by(issue, year) %>%  
  summarise(vote = sum(vote), n = n(),  
    country_percent = yes / max(yearly_vote, na.rm = TRUE),  
    yearly_vote = sum(vote) * 100 / n(),  
    issue_percent = yes / max(issue_vote, na.rm = TRUE),  
    issue_percent_label = paste0(issue, " (% of total votes)")) %>%  
  mutate(issue_percent_label = percent)  
  # Add labels to the plot  
  # Create a "Percentage of 'Yes' votes in the UN General Assembly"  
  # substitute a "%(n) to 100%"  
  # n = "Year"  
  # else = "Country"  
  #
```

Percentage of "Yes" votes in the UN General Assembly

References

1. Niles Gavaghan (2019). `unvotes`: A tidy dataset of international votes from the UN General Assembly. *GitHub*. GitHub version 0.2.0.
2. Niles Gavaghan (2019). `unvotes`: International votes from the UN General Assembly. *R package distribution*.
3. United Nations. (2019). `unvotes`: International votes from the UN General Assembly. *Machine-readable handbook of International Organization*, added by Eric Rennert (published May 21, 2019).
4. Much of the analysis has been modeled on the examples presented in the `unvotes` package vignettes.

Appendix

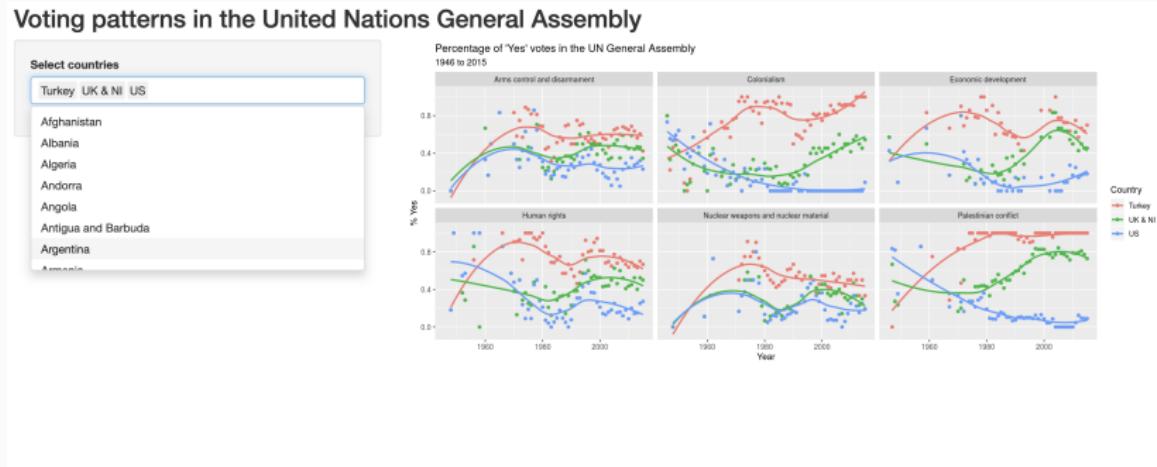
Below is a list of countries in the dataset.

Show 40 Search:

country
1 Afghanistan
2 Albania
3 Algeria
4 Armenia
5 Angola
6 Antigua and Barbuda
7 Argentina
8 Armenia
9 Australia
10 Austria

Showing 1 to 10 of 780 entities 1 2 3 4 5 ... 20 200

3. a data product: RShiny App



Data Science tasks require

1. Computer Programming

- A computer program consists of instructions that tell the computer what to do (with some data).
- PSTAT 10

2. Statistical Analysis

- the **science** of collecting, organizing, exploring, interpreting, and presenting **data** to uncover patterns, trends, making predictions based on the data.
- Statistical Analysis is grounded in Probability Theory

3. Probability

- is the study of uncertainty or randomness and its consequences in the world around us.
 - PSTAT 120A.
-
- In the first part of this course, we will focus on R programming language and introduction to probability using the R software.
 - In the second part of the course, you will be using the Python software.

R Course Objectives

1. Learn basics of reproducible research using Rmarkdown.
2. Learn basics of R programming language for Data science.
3. Learn introductory Statistics and Probability for Data science using R

Student learning outcomes

By the end of the R portion, students will

- Be comfortable using Rmarkdown to write reproducible reports.
- Have experience with introductory programming and data analysis tasks in R,
- Have a basic understanding of the fundamental concepts in probability theory.

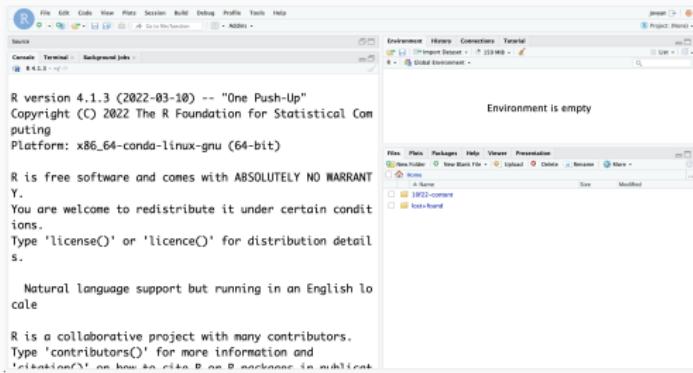
To succeed . . .

- Arrive on time
- Stay engaged
 - Avoid distractions (phone, email, social media, messaging apps turned off)
 - Actively participate
 - Ask questions
- Keep your files organized.
- Take half an hour after each class to review the materials and revise your notes.
- Take fifteen minutes before next class to skim material for upcoming class.

0.2 YOUR TURN 1 (10 minutes)

0.2 YOUR TURN 1 (10 minutes)

1. Login to our course Rstudio server using your **ucsb credentials** (link on canvas)

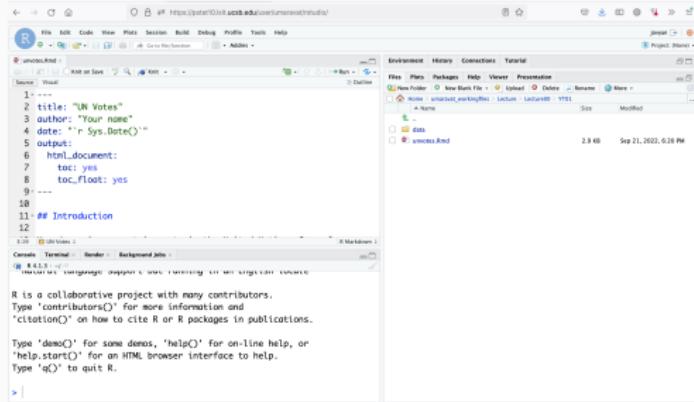


Setup My Rstudio Correctly

2. In the Files pane (bottom right corner), navigate to

yournetid_workingfiles -> Lecture00 -> 3YT01

and spot the file called unvotes.Rmd.



The screenshot shows the RStudio interface. The Source tab displays the Rmd code:

```
1 ---
2 title: "Un Votes"
3 author: "Your name"
4 date: "r Sys.Date()"
5 output:
6   html_document:
7     toc: yes
8     toc_float: yes
9 ---
10
11 ## Introduction
12
```

The Files pane in the bottom right shows the project structure:

- unvotes.Rmd
- data

The unvotes.Rmd file was modified on Sep 21, 2022, at 6:26 PM.

The R console at the bottom shows standard R startup messages and help information.

Knitting

5. Open it and click “Knit”.



6. Go back to the file `unvotes.Rmd` and change your name on top (in the `yaml` – we'll talk about what this means later) and knit again.
7. Change the country names to those you're interested in. Spelling and capitalization must match how the countries appear in the data, so take a peek at the Appendix to confirm spelling.
8. Knit again. Voila, your first data visualization!

You should now be familiar with the computing environment we will use in this course.

0.3 R Markdown documents (and newer quarto documents)

0.3 R Markdown documents (and newer quarto documents)

- with .Rmd extension
- **rmarkdown** is an R package
 - write code and prose in **reproducible** computational documents

rmarkdown.rstudio.com



Take a look at the **Rmarkdown gallery**

R Markdown syntax

- Code goes in code chunks, defined by three backticks
- narrative goes outside of chunks
- Simple markdown syntax for text

Tour: R Markdown

R Markdown help in RStudio

R Markdown Cheat Sheet Help -> Cheatsheets

R Markdown :: CHEAT SHEET

What is R Markdown?



.Rmd files - An R Markdown (.Rmd) file is a record of your research. It contains the code that a scientist needs to reproduce your work along with the narration that a reader needs to understand your work.



Reproducible Research - At the click of a button, or the type of a command, you can rerun the code in an R Markdown file to reproduce your work and export the results as a finished report.



Dynamic Documents - You can choose to export the finished report in a variety of formats, including html, pdf, MS Word, or RTF documents; html or pdf based slides, Notebooks, and more.

Workflow



- ① Open a new .Rmd file at File ► New File ► R Markdown. Use the wizard that opens to pre-populate the file with a template

- ② Write document by editing template

- ③ Knit document to create report; use knit button or render() to knit

- ④ Preview Output in IDE window

- ⑤ Publish (optional) to web server

- ⑥ Examine build log in R Markdown console

- ⑦ Use output file that is saved along side .Rmd

The screenshot shows the RStudio environment. In the center, the R Markdown editor displays the following code:

```
1 # ...
2 # title: "R Markdown"
3 # author: "RStudio"
4 # output: 2
5 # html_document
6 # toc: TRUE
7 ...
8
9 ---[if !setup, !include(FALSE)]
10 knitr::opts_chunk$set(echo = TRUE)
11 ...
12
13 ## R Markdown
14
15 This is on R Markdown document.
16 Markdown is a simple formating
17 Syntax for authoring HTML, PDF,
18 and MS Word documents.
19
20 ---{r cars}
21 summary(cars)
22
23
24 For more details on using R Markdown
25 see http://rmarkdown.rstudio.com.
```

Annotations with arrows point to various UI elements:

- set preview location
- insert code chunk
- run code chunk(s)
- go to code chunk
- publish
- show outline
- run all previous chunks
- modify chunk options
- run current chunk

The preview pane on the right shows the rendered HTML output:

R Markdown

This is an R Markdown formatting syntax for documents.

summary(cars)

	speed
#&	Min.
#&	1st Qu.
#&	Median
#&	Mean
#&	3rd Qu.
#&	Max.

For more details on u
<http://rmarkdown.rsti>

The file browser at the bottom shows a file named "report.html".

render

Use `markdown::render()` to render/knit at cmd line. Important args:

`input` - file to render
`output_format`

`output_options` -
List of render
options (as in YAML)

`output_file`
`output_dir`

`params` - list of
params to use

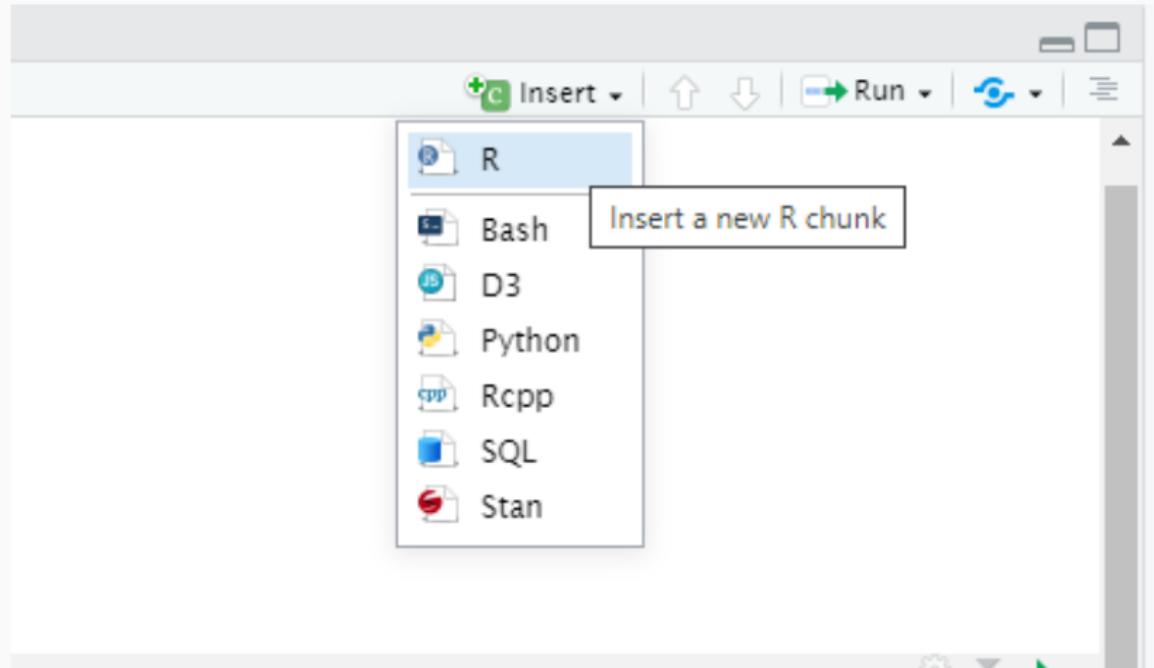
How will we use R Markdown?

- You'll always have a template R Markdown document to start with.
- You can also create your own .Rmd file (File -> New File -> R markdown ... -> Knit and modify text or code as necessary.)

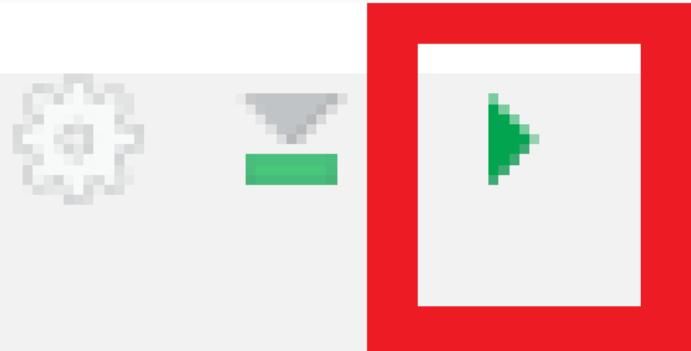
(Optional) Your Turn: Rmarkdown basics

- Open `yournetid_workingfiles/L00/P01_optional/ -> P01_netid.Rmd`
- This is optional but highly recommended.
- Get help on Week 2 Monday office hours.

Demo - Adding Chunks



Demo - Run Code



- Run entire chunk using the “green triangle”
- Run any part(one or more lines) of your code by selecting it and pressing Ctrl+Enter / Cmd+Enter

Lecture 0 Summary

- Core elements of Data Science project life-cycle
 - Programming
 - Statistics
 - Probability
- Accessing Rstudio server instance for the course
- Created a Data Science project report for UN votes.
- Rmarkdown essentials.(Complete it but do not turn it in .)

Post Lecture 0 to-do for you

- Read syllabus carefully
- Note down important dates,
- Get familiar with Course site on Canvas
- Visit Office hours
 - Get help with lecture material if you struggled in lecture today.
 - Practice will make it perfect for you!

Have a great start to the course! See you next lecture!