

0.Welcome to PSTAT 10!

Principles of Data Science with R

Dr. Uma Ravat

PSTAT 10

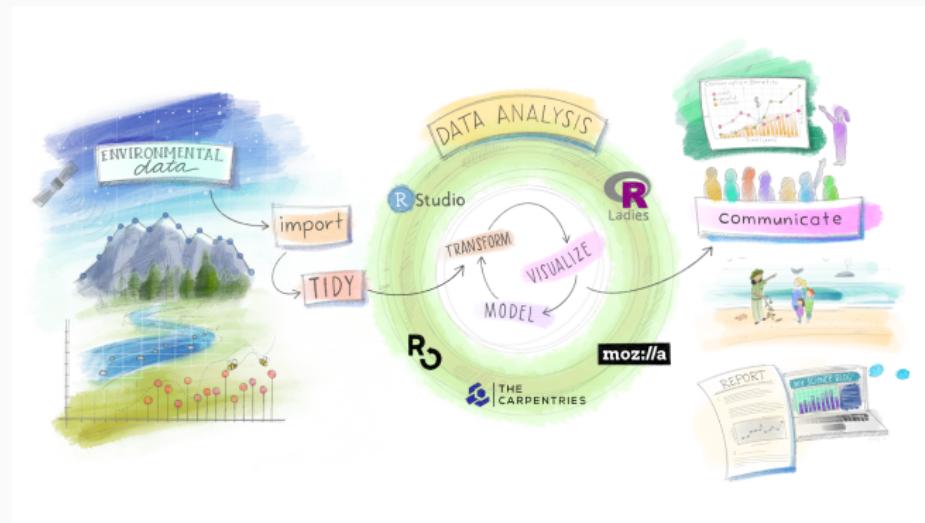
0.1 Welcome!

Plan for today

1. What is Data Science?
2. Your Turn at Data Science
3. **Course Overview** Read syllabus at length after lecture
4. **(Course Toolkit)** Rmarkdown - finish in section 1, homework 1

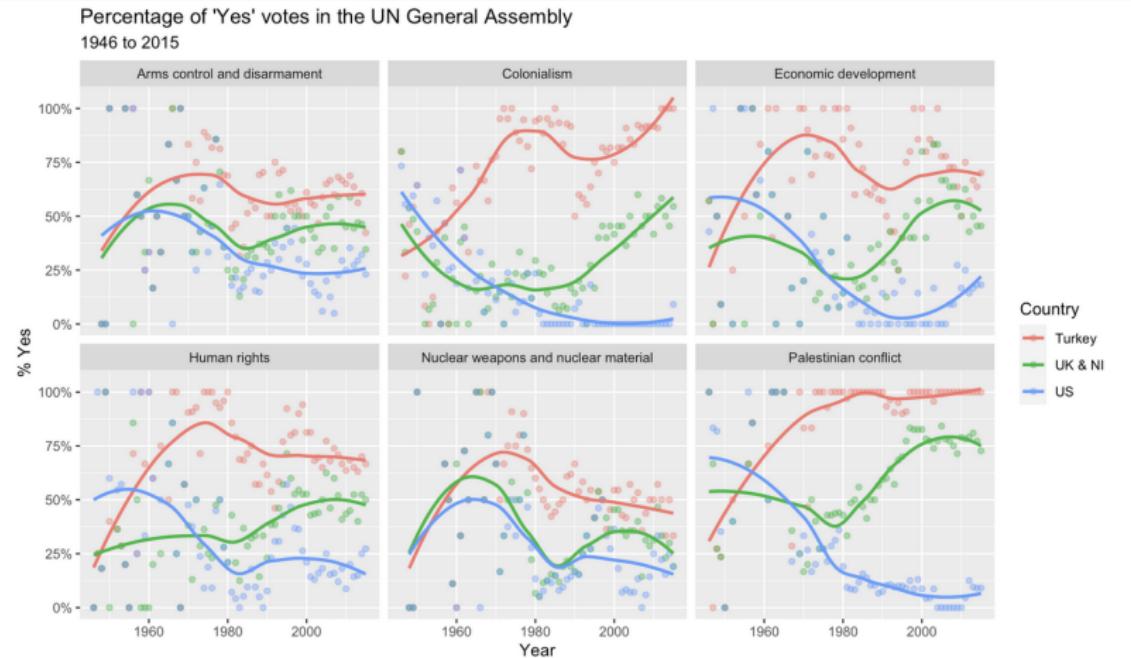
What is Data Science?

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



Your data science tasks result in...

1. A Visualization



2. Report/Website

Introduction

U.S. election patterns

References

Appendix

UN Votes

Your name here

2022-08-24

Introduction

How do various countries vote in the United Nations General Assembly, how have their voting patterns evolved throughout time, and/or differently do they view certain issues? Answering these questions (at a high level) is the focus of this analysis.

Packages

We will use the `tidyverse`, `shiny`, and `shinyWidgets` packages for data cleaning and visualization, and the `gapminder` package for obtaining UN data.

```
library(tidyverse)
library(shiny)
library(shinyWidgets)
library(gapminder)
```

Data

The data we're using originally came from the `anewdemos` 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.

```
library("readr")
library("anewdemos", "gapminder", "shiny")
```

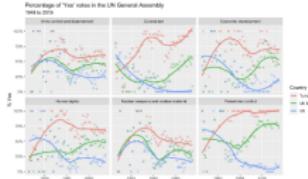
UN voting patterns

The following code shows the voting record of the UN's 194 member states over time on a variety of issues, and compares it to two other countries: US and Turkey.

```
# We can use the gapminder package to get data by changing which countries the code above looks for and to see if our country should be split and categorized exactly the way it appears in the data. See the FAQ for a lot of details.
```

```
countries %>
  filter(name %in% c("USA", "TR", "Turkey")) %>
  group_by(year, country, year_start, year_end) %>
  summarise(votes = sum(vote), p_percent_vote = sum(vote)/n()) %>
  mutate(p_percent_vote = ifelse(is.na(p_percent_vote), 0, p_percent_vote)) %>
  mutate(is_na = ifelse(is.na(votes), 1, 0)) %>
  mutate(is_na = ifelse(is_na == 1, TRUE, FALSE)) %>
  mutate(is_na = ifelse(is_na == 0, FALSE, TRUE)) %>
  mutate(is_na = ifelse(is_na == 1, TRUE, FALSE)) %>
  mutate(is_na = ifelse(is_na == 0, FALSE, TRUE))
```

title = "Percentage of 'yes' votes in the UN General Assembly",
x_label = "Year",
y_label = "Country",
color = "#0072BD"



References

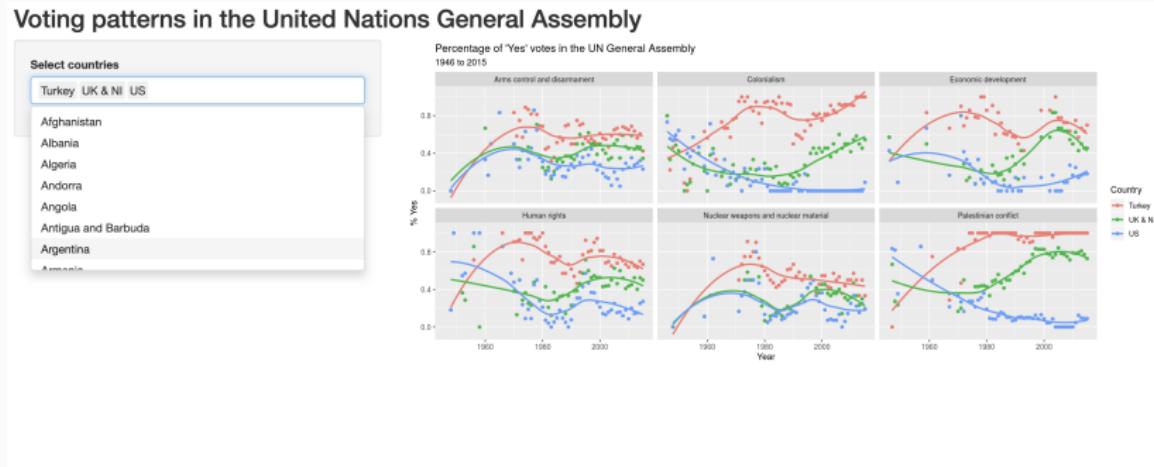
6. *Uline-CetraGroup, Rundt, Introduction to Data science course*
1. David Robinson (2011). *UNIVote*: United Nations General Assembly Voting Data. R package version 0.2.3.
2. Eric Weisstein. *Data and Analyses of "Voting in the UN General Assembly"*. Routledge Handbooks of International Organizations, edited by Bob Reiterman (published May 27, 2008).
3. Much of the analysis has been modified on the exercises presented in the *univote* package vignettes.

Appendix

Below is a list of countries in the dataset.

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

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).

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. (for large amount of data) **Working with Databases** using **SQL** (Structured Query Language) to communicate with all the data stored in a database.

- This course is an introduction to these components of data science using the software R

0.2 YOUR TURN 1 (10 minutes)

0.2 YOUR TURN 1 (10 minutes)

1. Login to our PSTAT 10 Rstudio computing environment using your **ucsb credentials** at the Rstudio Server link on left navigation menu on Canvas

2. In the Files pane (bottom right corner), navigate to **yournetid_workingfiles -> Lecture -> Lecture00 -> YT01** and spot the file called unvotes.Rmd.

The screenshot shows the RStudio interface. The Source pane 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: # Introduction
11:
12:
```

The Files pane shows the project structure:

- unvotes.Rmd
- data

The bottom console pane shows the R environment:

```
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.
```

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.

Review Your Turn 1

0.3 Course and syllabus overview

Course Objectives

Using the R programming language to introduce the basics of data science

Part 1. Learn the basics of programming using the R programming language.

Part 2. Learn Introductory Statistics and Probability for Data science using R

Part 3. Learn Relational Databases and SQL by interfacing with R

Student learning outcomes

By the end of the course, students should have a solid foundation in

- programming with R,
- a basic understanding of probability and simulations and
- the ability to work with databases using SQLite.

This course provides a stepping stone for further studies in data science and related fields.

Administrative items

1. No crashing - wait till you are officially enrolled
2. No section switching
3. On Canvas, see syllabus and read about office hours, email policy, other policies.

Course success

Show up, don't fall behind.

1. Syllabus **Read it**
2. Lectures **Show up and participate**
3. Labs **Don't miss**
4. Homeworks **Do them on time, come to Office Hours for help**
5. Homework Reflection Surveys **Don't miss**
6. Quizzes **Schedule a time to take it in Week 2,3,5,7,10**
7. Final exam **Don't miss**
8. Office hours are helpful - over 20 hours each week. Make use of them!

Don't share my work with anyone or on sites like Coursehero, Chegg, AI etc.

Preferred order of Methods of communication

1. Office Hours

- Come and get real time answers to your questions or just say hi!
- ULA office hours
- TA office hours
- instructor office hours
- HW clinics
- over 20 hours each week. **Please make use of these help hours**

3. Email for **truly private matters** and follow **email policy**

- Grade discrepancies do not need an email - Fill in the form
- your TA is your first point of contact for any issues regarding anything else about section, homework, quiz
- Read syllabus for DSP accomodations and conflicts before contacting Instructor during office hours for any DSP accommodation and conflicts
- instructor for private matters that you do not want any other course staff to know about

Collaboration and sharing your/my work

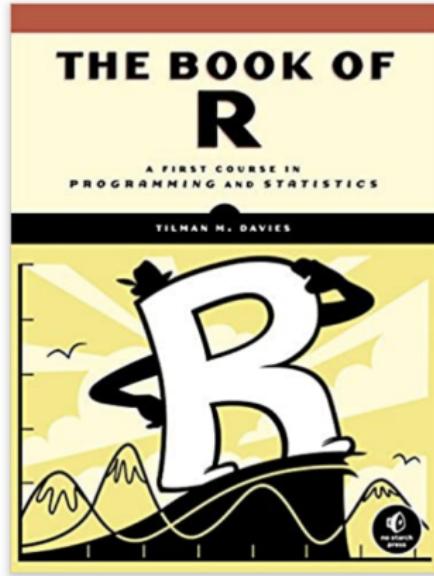
Asking questions and helping others is encouraged

- Discussing course topics with other is welcome!
- Discuss with peers sitting next to you.
- Ask questions, answer questions in class
- Don't judge people or questions
- share ideas liberally.

Limits of collaboration

- Don't share solutions with each other
- You should never see or have possession of anyone else's solutions, including from previous versions of this course
- See syllabus about Heuristics to keep in mind while sharing

TEXTBOOK



You do not need a textbook for this class. Slides will be provided and starter .Rmd will be provided.

Title: The Book of R: A First course in Programming and Statistics

Author: Tilman M. Davies

Publisher: No Starch Press

Summary of Technology/Platforms

Personal PC, Mac or Laptop, Phone : To access course sites

Canvas Course site: Canvas for submitting assignments, grades, reviewing solutions

Rstudio Instance: for code files for lecture, lab, homework

Any of the many Office hours: for questions

Email: Use sparingly for truly private matters for lab TA, headTA or instructor

Do not email course staff separately for the same issue - Use Emails' CC feature while emailing multiple course staff about the same issue.

Most importantly...

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

0.4 R Markdown documents

0.4 R Markdown 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 interface with several windows open:

- Code Editor:** Displays an R Markdown script (report.Rmd) with code like `knitr::opts_chunk\$set(echo = TRUE)` and `summary(cars)`. Callouts point to buttons for "set preview location", "insert code chunk", "run code chunk(s)", "run all previous chunks", "modify chunk/options", and "run current chunk".
- Console:** Shows the command `library(rmarkdown)` and the result of `render("report.Rmd", output_file = "report.html")`.
- Output pane:** Displays the rendered HTML output of the R Markdown document.
- Help pane:** Shows the "R Markdown" help page from the cheatsheet.
- File browser:** Shows the file structure with "report.html" and "report.Rmd" files.

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

input	file to render	output_options	output_file	params
		- List of render options (as in YAML)	output_dir	- list of params to use

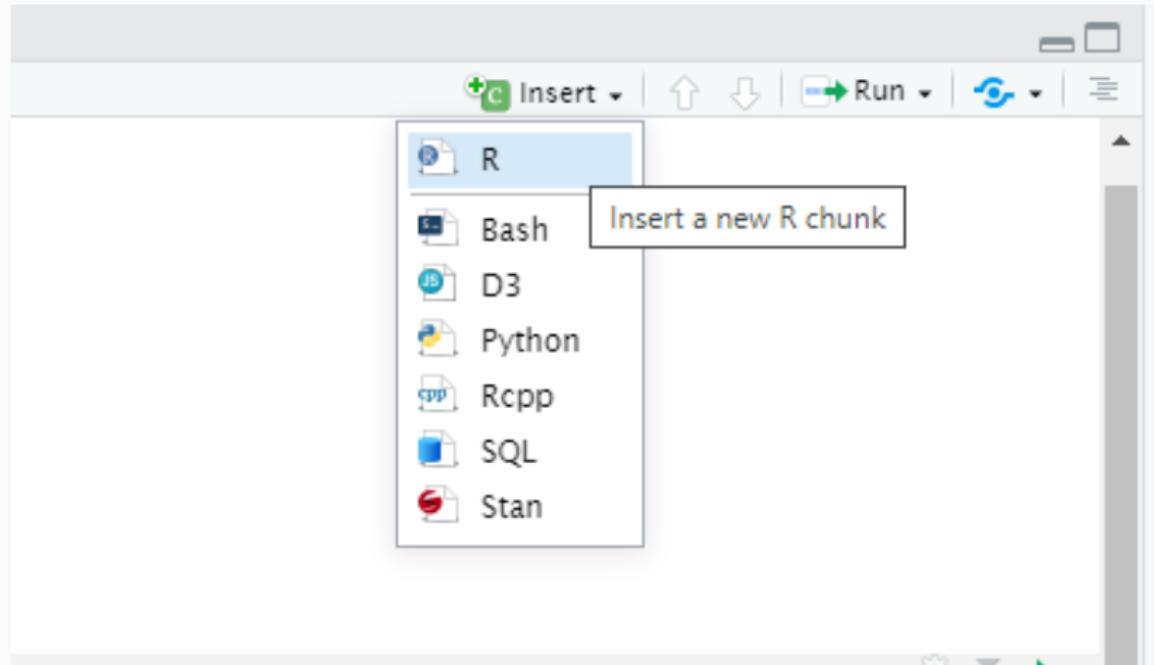
How will we use R Markdown?

- Every worksheet / homework etc. is an R Markdown document
- You'll always have a template R Markdown document to start with
- The amount of scaffolding in the template will decrease over the quarter
- You can also create your own .Rmd file (File -> New File -> R markdown .. -> Knit and modify text or code as necessary.)

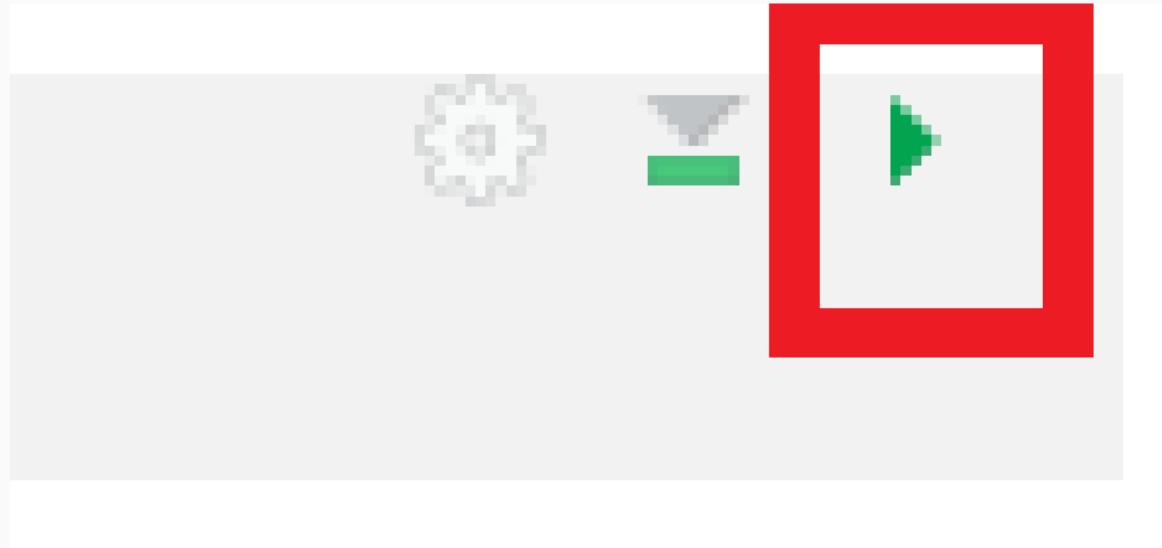
Your Turn: Rmarkdown basics

- Open `yournetid_workingfiles/Homework/HW01/ -> Homework01.Rmd`
- You will start here, continue in Section, complete it and submit it as your first homework on Wednesday on canvas.
- HW Submission link will be released in Week 1 module on Canvas
- Use office hours to practice, review ask questions about today's lecture and these exercises.

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

Summary: Introduction

- Core elements of Data Science project life-cycle
 - Programming
 - Statistics and Probability
 - Databases
- Accessing Rstudio instance for the course
- created a Data Science project report for UN votes.
- Course overview and Brief Syllabus walk through
- Rmarkdown essentials.(Complete it in section 1)

Post Lecture 0 to-do for you

- Read syllabus carefully
- Note down important dates, final exam
- Get familiar with Course site on Canvas
- Go to both Sections each week and ask questions when you are stuck.
- Complete Homework 1 and submit on time.
- 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 quarter! See you next lecture!