

Stat 220 Introduction to Data Science

Deepak Bastola

2024-02-27

Contents

Course overview	7
0.1 Learning Objectives	7
 Set-up Instructions	 11
1 What is R, RStudio, and RMarkdown?	11
1.1 What is RStudio?	11
1.2 R Studio Server	11
1.3 R/RStudio	12
1.4 Installing R/RStudio (not needed if you are using the maize server)	12
1.5 What is RMarkdown?	13
1.6 Install LaTeX (for knitting R Markdown documents to PDF): . .	13
1.7 Updating R/RStudio (not needed if you are using the maize2 server)	14
1.8 Opening a new file	14
1.9 Running codes and knitting .Rmd files:	14
1.10 Few More Instructions	15
1.11 VPN	15
 2 Assignments in Stat 220	 17
2.1 Do's and Don't of collaboration for individual assignments	17
2.2 Format and Content	18
 3 Software in Stat 220	 19
3.1 File organization: Using maize	19
3.2 File organization: Using your own Rstudio	22
3.3 RStudio projects	22
3.4 Best practices (or what not to do)	22
3.5 Git and GitHub	23
3.6 Slack	23
3.7 Acknowledgements	24

4	GitHub Guide for Students in Stat 220	25
4.1	Overview	25
4.2	Getting setup with Git and GitHub	25
4.3	Individual assignments	26
4.4	Group work	29
4.5	Additional resources	30
4.6	Acknowledgements	30
4.7	Reuse	30
5	R Markdown Syntax	31
6	Github Tutorial	35
6.1	Tutorial 1: Creating and cloning a Repository starting from Github to RStudio	36
6.2	Tutorial 2: Creating a new GitHub repository using <code>usethis</code> R package (RStudio to Github) (Works ONLY on local RStudio)	37
	Class Activities	41
7	Class Activity 1	41
7.1	Extras (optional)	45
7.2	Questions	47
8	Class Activity 2	49
9	Class Activity 3	51
9.1	Question 1: data types	51
9.2	Question 2: Subsetting and coercion	52
9.3	Question 3: Lists	54
10	Class Activity 4	57
10.1	Your turn 1	57
11	Class Activity 5	69
11.1	Problem 1: Changing color and shape scales	69
11.2	Problem 2: US maps	73
11.3	Problem 3: Choropleth map	77
12	Class Activity 6	83
12.1	Problem 1: <code>select()</code>	83
12.2	Problem 2: <code>filter()</code>	84
12.3	Problem 3: <code>arrange()</code>	85
12.4	Problem 4: <code>mutate()</code>	86
12.5	Problem 5: <code>summarize()</code> or <code>summarise()</code>	87
12.6	Problem 6	88

13 Class Activity 7	91
13.1 Problem 1: Boolean Operators	91
13.2 Problem 2: Joining artists and bands data	99
13.3 Problem 3: Filtering and counting rows in the data	100
14 Class Activity 8	103
14.1 Your turn 1	103
14.2 Your turn 2	107
15 Class Activity 9	111
15.1 Your Turn 1	111
15.2 Your Turn 2	114
16 Class Activity 10	117
16.1 Your Turn 1	117
16.2 Your Turn 2	120
17 Class Activity 11	125
17.1 Problem 1	125
17.2 Problem 2	128
17.3 Problem 3	129
17.4 Problem 4	130
18 Class Activity 12	133
18.1 Group Activity 1	133
18.2 Group Activity 2	135
18.3 Group Activity 3	137
19 Class Activity 13	141
20 Class Activity 14	143
20.1 Group Activity 1	143
20.2 Group Activity 2	145
21 Class Activity 15	149
21.1 Group Activity 1	149
21.2 Group Activity 2	152
22 Class Activity 16	157
22.1 Group Activity 1	157
22.2 Group Activity 2	158
22.3 Group Activity 3	160
23 Class Activity 17	163
23.1 Group Activity 1	163
23.2 Group Activity 2	165

24 Class Activity 18	169
24.1 Introduction to Shiny Apps	169
25 Class Activity 19	175
25.1 Question 1	178
26 Class Activity 20	181
26.1 Group Activity 1	181
26.2 Group Activity 2	181
27 Class Activity 21	187
27.1 Group Activity 1	187
28 Class Activity 22	193
28.1 Group Activity 1	193
28.2 Group Activity 2	195
28.3 Extra: Code to recreate the plot in the slides for the diabetes dataset.	196
29 Class Activity 23	199
29.1 Group Activity 1	199
29.2 Group Activity 2	202
30 Class Activity 24	205
30.1 Group Activity 1	205
30.2 Group Activity 2	206

Course overview

Greetings and welcome to Introduction to Data Science! In this course, we will delve into the computational aspects of data analysis, covering topics such as data acquisition, management, and visualization tools. Throughout this course, we will emphasize the principles of data-scientific, reproducible research and dynamic programming, utilizing the R/RStudio ecosystem.

If you have taken Stat 120, 230, or 250 at Carleton, you will find yourself well-equipped to handle the material. However, it is important to refresh your R and R-markdown skills before the start of the class. Specifically, I expect all students to be able to load a data set into R, calculate basic summary statistics, and perform basic exploratory data analysis. In the first week of class, we will delve into Git and GitHub version control, though prior exposure to these topics is not necessary.

0.1 Learning Objectives

- Develop research questions that can be answered by data. Import/scrape data into R and reshape it to the form necessary for analysis.
- Manipulate common types of data, including numeric, categorical (factors), text, date-times, geo-location variables in order to provide insight into your data and facilitate analysis.
- Explore data using both graphical and numeric methods to provide insight and uncover relationships/patterns.
- Utilize fundamental programming concepts such as iteration, conditional execution, and functions to streamline your code.
- Build, tune, use, and evaluate basic statistical learning models to uncover clusters and classify observations.
- Draw informed conclusions from your data and communicate your findings using both written and interactive platforms.

Set-up Instructions

Chapter 1

What is R, RStudio, and RMarkdown?

R is a free and open source statistical programming language that facilitates statistical computation. There are a myriad of application that can be done in R, thanks to a huge online support community and dedicated packages. However, R has no graphical user interface and it has to be run by typing commands into a text interface.

1.1 What is RStudio?

RStudio provides graphical interface to R! You can think of RStudio as a graphical front-end to R that that provides extra functionality. The use of the R programming language with the RStudio interface is an essential component of this course.

1.2 R Studio Server

The quickest way to get started is to go to <https://maize2.mathcs.carleton.edu>, which opens an R Studio window in your web browser. Once logged in, I recommend that you do the following:

- Step 1: Create a folder for this course where you can save all of your work. In the Files window, click on New Folder.
- Step 2: Click on Tools -> Global Options -> R Markdown. Then uncheck the box that says “Show output inline...”

(It is also possible to download RStudio on your own laptop. Instructions may be found at the end of this document.)

1.3 R/RStudio

The use of the R programming language with the RStudio interface is an essential component of this course. You have two options for using RStudio:

- The **server version** of RStudio on the web at (<https://maize2.mathcs.carleton.edu>). The advantage of using the server version is that all of your work will be stored in the cloud, where it is automatically saved and backed up. This means that you can access your work from any computer on campus using a web browser. This server may run slow during peak days/hours. I also recommend you to download a local version of R server in your computer in case of rare outages.
- A **local version** of RStudio installed on your machine. This option is highly recommended due to the computational resources this course demands. Using this version you can only store your files in your local machine. Additionally, we can save our work on GitHub. We will learn how to use GitHub in the beginning of the course. Both R and RStudio are free and open-source. Please make sure that you have recently updated both R and RStudio.

1.4 Installing R/RStudio (not needed if you are using the maize server)

Download the latest version of R: <https://cran.r-project.org/> Download the free Rstudio desktop version: <https://www.rstudio.com/products/rstudio/download/>

Use the default download and install options for each. For R, download the “precompiled binary” distribution rather than the source code

Updating R/RStudio (not needed if you are using the maize server)

If you have used a local version of R/RStudio before and it is still installed on your machine, then you should make sure that you have the most recent versions of each program.

- To check your version of R, run the command `getRversion()` and compare your version to the newest version posted on <https://cran.r-project.org/>. If you need an update, then install the newer version using the installation directions above.
- In RStudio, check for updates with the menu option **Help > Check for updates**. Follow directions if an update is needed.

**** Did it work? (A sanity check after your install/update) ****

Do whatever is appropriate for your operating system to launch RStudio. You should get a window similar to the screenshot you see here, but yours will be

more boring because you haven't written any code or made any figures yet!

Put your cursor in the pane labeled *Console*, which is where you interact with the live R process. Create a simple object with code like `x <- 2 * 4` (followed by enter or return). Then inspect the `x` object by typing `x` followed by enter or return. You should see the value 8 printed. If this happened, you've succeeded in installing R and RStudio!

1.5 What is RMarkdown?

An R Markdown file (.Rmd file) combines R commands and written analyses, which are 'knit' together into an HTML, PDF, or Microsoft Word document.

An R Markdown file contains three essential elements:

- Header: The header (top) of the file contains information like the document title, author, date and your preferred output format (pdf_document, word_document, or html_document).
- Written analysis: You write up your analysis after the header and embed R code where needed. The online help below shows ways to add formatting details like bold words, lists, section labels, etc to your final pdf/word/html document. For example, adding `**` before and after a word will bold that word in your compiled document.
- R chunks: R chunks contain the R commands that you want evaluated. You embed these chunks within your written analysis and they are evaluated when you compile the document.

1.6 Install LaTeX (for knitting R Markdown documents to PDF):

You need a Latex compiler to create a pdf document from a R Markdown file. If you use the maize server, you don't need to install anything. If you are using a local RStudio, you should install a Latex compiler. Below are the recommended installers for Windows and Mac:

- MacTeX for Mac (3.2GB)
- MiKTeX for Windows (190MB)
- Alternatively, you can install the `tinytex` R package by running `install.packages("tinytex")` in the console.

1.7 Updating R/RStudio (not needed if you are using the maize2 server)

If you have used a local version of R/RStudio before and it is still installed on your machine, then you should make sure that you have the most recent versions of each program.

- To check your version of R, run the command `getRversion()` and compare your version to the newest version posted on <https://cran.r-project.org/>. If you need an update, then install the newer version using the installation directions above.
- In RStudio, check for updates with the menu option **Help > Check for updates**. Follow directions if an update is needed.

1.8 Opening a new file

If using Rstudio on your computer, using the **File>Open File** menu to find and open this .Rmd file.

If using Maize Rstudio from your browser:

- In the Files tab, select **Upload** and **Choose File** to find the .Rmd that you downloaded. Click *OK* to upload to your course folder/location in the maize server account.
- Click on the .Rmd file in the appropriate folder to open the file.

1.9 Running codes and knitting .Rmd files:

- You can run a line of code by placing your cursor in the line of code and clicking **Run Selected Line(s)**
- You can run an entire chunk by clicking the green triangle on the right side of the code chunk.
- After each small edit or code addition, **Knit** your Markdown. If you wait until the end to Knit, it will be harder to find errors in your work.
- Format output type: You can use any of pdf_document, html_document type, or word_document type.
- **Maize users:** You may also need to allow for “pop-up” in your web browser when knitting documents.

1.10 Few More Instructions

The default setting in Rstudio when you are running chunks is that the “output” (numbers, graphs) are shown **inline** within the Markdown Rmd. If you prefer to have your plots appear on the right of the console and not below the chunk, then change the settings as follows:

1. Select Tools > Global Options.
2. Click the R Markdown section and uncheck (if needed) the option Show output inline for all R Markdown documents.
3. Click OK.

Now try running R chunks in the .Rmd file to see the difference. You can recheck this box if you prefer the default setting.

1.11 VPN

If you plan to do any work off campus this term, you need to install Carleton’s VPN. This will allow you to access the **maize** server (if needed).

Installing the GlobalProtect VPN

Follow the directions here to install VPN.

Chapter 2

Assignments in Stat 220

2.1 Do's and Don't of collaboration for individual assignments

- You *can* discuss homework problems with classmates but you must write up **your own** homework solutions and **do your own work in R (no sharing commands or output)**.
 - **Do not share R commands/code in any way**, including, but not limited to, sending commands via email, slack, text, or showing commands in a shared screen with the intention of showing a classmate your solution to a problem.
 - You **can** share a screen to help troubleshoot a coding problem in R.
- You *can* use the following resources to complete your homework:
 - Carleton faculty (myself, other math or statistics faculty, etc)
 - discussions with classmates (see above) or knowledgeable friends
 - Carleton resources like stats lab assistants
 - student solutions provided in the back of your student textbook or in the student solution manual.
- You *cannot* use any resources other than the ones listed above to complete assignments (homework, reports, etc) for this class. (e.g. you cannot use a friend's old assignments or reports, answers found on the internet, textbook (instructor) solutions manual, etc.)

2.1.1 Examples that violate the academic integrity policy

- sending your .Rmd homework file to another person in the class
- receiving an .Rmd homework file from another person
- sharing a screen and copying code, verbatim, from another person
- sending/receiving R commands

- neglecting to acknowledge classmates with whom you worked with on an assignment

2.2 Format and Content

Submit via GitHub (for most assignments) an organized and correctly ordered assignment.

- Content: Good data scientists need to do more than just write code; they should be able to interpret and explain their analyzes.
 - Provide a **written answer** first, followed by any required R code and output.
 - Use **complete sentences** when answering any problem that requires an explanation or overall problem summary.
- When including code:
 - Be sure to show the natural sequence of work needed to answer the problem.
 - Include brief comments explain your code steps.
 - Do not include typos or unnecessary commands/output.
 - Always include code output.
- At the top of each individual assignment **include the names of classmates that you worked with** on all or part of the assignment (but each person must write up their assignment on their own)

Disability Accommodations: Carleton College is committed to providing equitable access to learning opportunities for all students. The Disability Services office (Henry House, 107 Union Street) is the campus office that collaborates with students who have disabilities to provide and/or arrange reasonable accommodations. If you have, or think you may have, a disability (e.g., mental health, attentional, learning, autism spectrum disorders, chronic health, traumatic brain injury and concussions, vision, hearing, mobility, or speech impairments), please contact disability@carleton.edu or call Sam Thayer ('10), Accessibility Specialist (x4464) or Chris Dallager, Director of Disability Services (x5250) to arrange a confidential discussion regarding equitable access and reasonable accommodations.

Academic Honesty: All work that you turn in under your name must follow Carleton's academic integrity policy. The use of textbook solution manuals (physical or online solutions), homework, reports or exams done by past students are not allowed. Look at the College's Writing Across the Curriculum website for additional guidance on plagiarism and how to avoid plagiarism in their writing.

Chapter 3

Software in Stat 220

You will work with many .Rmd Markdown files in this course. These include class activities, homework template, project helper files etc. To stay organized, I *strongly* suggest you create a **stat220** folder that contains the following sub-folders:

- **stat220** folder
 - **Assignments:** This folder will contain subfolders for each assignment. Each assignment subfolder (e.g. homework1, homework2, ...) will be a Github connected RStudio project that you will create **once an assignment is posted**.
 - **Content:** This folder should be used to save any non-assignment files (e.g. slides, examples) for this class. You will create this subfolder by creating an RStudio project (see step 5 below).

To get started with this organization, follow the steps below.

3.1 File organization: Using maize

The server (online) version of Rstudio is run from a unix server. You can navigate this file system using unix commands, but I assume that most or all of you will just use Rstudio to access your files on this server.

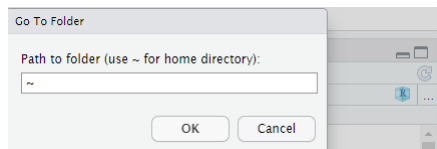
1. In Rstudio, click the **Files** *tab* in the lower right-hand window.

Note: this is **not** the same as the **File** *menu* option.

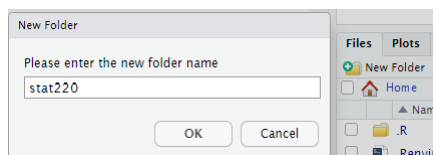


2. Verify that you are in your **HOME** folder (should simply say

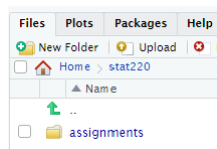
Home right under the New Folder button). To navigate to your Home folder (if somehow you are not in it), click the ... button (far right side of the **Files** tab) and enter a ~ (tilde) symbol



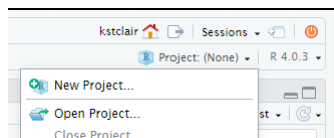
- **3.** Click the **New Folder** button and name the folder **stat220**.



- **4.** Click on this newly created (empty) **stat220** folder. Within the folder create another **New Folder** and name it **assignments**.



- **5.** Within the **stat220** folder, create an **RStudio project** called **content** with the following steps:
 - **a.** Click the **Project** button in the upper righthand corner of your RStudio window and select **New Project....**



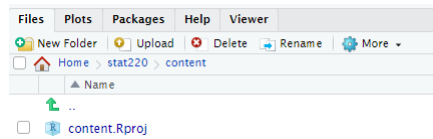
- **b.** Select **New Directory** and then **New Project**



- c. Enter **content** as the **Directory name** and use the **Browse** button to find your **stat220** folder. Then click **Create Project**.



- d. You should now have a new folder called **content** in your **stat220** folder and this folder will contain an RStudio project **.Rproj**. Feel free to add subfolders to this **content** folder (e.g. slides, examples, etc).



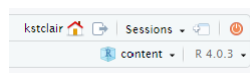
Warning: Do not create an RStudio project in the main `stat220` folder because it is not good practice to have RStudio projects in subfolders of another project (e.g. a project within a project is not recommended).

3.2 File organization: Using your own Rstudio

Create a folder called **stat220** somewhere on your computer. Within this folder create an **assignments** subfolder. Then complete **step 5** from above to create a **content** RStudio project folder.

3.3 RStudio projects

Once you've created a project, your R session should be running within that project folder. You can check which project you are in by checking the project name in the upper righthand part of your RStudio window. Here we see the **content** project is open:



Running R from an RStudio project sets your **working directory** to the project folder:

```
> getwd()
[1] "/Accounts/kstclair/stat220/content"
```

This allows for easy file path access to all files related to this project.

To **start** a project, click on the `.Rproj` file or use the **Open Project...** option shown in step 5 above.

3.4 Best practices (or what not to do)

- Never save files to a lab computer hard drive (e.g. desktop, downloads, etc). They will be erased when you log off.

- Do not use gmail as a file storage system! Avoid emailing yourself files that you created (and saved) on a lab computer. Eventually you will lose work this way.
- Avoid using online versions of google drive and dropbox. Similar to gmail, downloading, editing a doc, then uploading it back to drive/dropbox is another great way to lose work.
- Avoid this and this.

3.5 Git and GitHub

Git is version control software that you install locally on your computer. Git is already installed on the maize RStudio server.

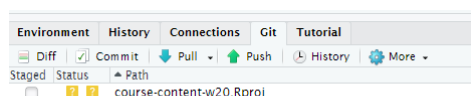
Github is a cloud-based service for hosting git projects. It allows multiple users to share and contribute to projects and it is how you will be submitting homework assignments and projects for this class. More information about Git and Github can also be found in Getting setup with Git and GitHub and Git and Github.

If you are using a local install of R/RStudio, then you will need to install Git.

Installing Git

Directions for both Windows & Mac here: <http://happygitwithr.com/install-git.html>.

1. If you are using **maize**, then there is nothing you need to install.
2. Windows users should follow Option 1 in 6.2.
3. Mac users can follow Option 1 in 6.3 if comfortable, otherwise follow Option 2
4. Linux users can follow 6.4.



3.6 Slack

We will use Slack for all course communication. Sign up for our course Slack team here! You will need to create an account with a username, and log in to read and post. You can download a standalone Slack application to your Mac, Windows, Linux and/or Android/iOS device. You can control whether you receive notifications on new posts by going to Preferences, as well as decide which 'channels' to subscribe to. A 'channel' is a discussion thread, which is used to organize communications into topics. You can learn more about Slack features here.

Several channels have been set up for specific parts of the course. Feel free to ask questions anytime. You can browse the available channels in our team by clicking on “Channels” on the left-hand panel.

3.7 Acknowledgements

This installation guide is based on the guide from Adam Loy and Katie St. Clair.

Chapter 4

GitHub Guide for Students in Stat 220

4.1 Overview

If you are using the maize RStudio server, then you can connect to GitHub without any extra software downloads. If you are using RStudio on your computer, then you will need to download Git software (as directed in Software in Stat 220) to use GitHub connected projects. You will use GitHub to submit homework and collaborate on projects.

4.2 Getting setup with Git and GitHub

If you are **not** working on the maize RStudio server, then make sure that you have installed all of the software mentioned in Software in Stat 220. In addition, you should install the `usethis` and `gitcreds` R packages.

Everyone needs to connect Git and GitHub by doing the following:

1. Register for account on GitHub (<https://github.com/>). I recommend using a username that incorporates your name (e.g., dbastola) and Carleton email address for your Github account.
2. If you haven't done so already, accept the invite to the class organization DataScienceWinter24. This organization is where all course homework files and project repositories will live.
3. Setup options in Git by running the following code chunk in your console:

```
#install.packages("usethis") # uncomment to install  
usethis::use_git_config(user.name = "Jane Doe", user.email = "jane@example.org")
```

changing the first two arguments to your own name and email (this should be the email associated with your GitHub account).

4. In order to push changes to github (i.e. to track changes and submit homework), you will need to prove that you have permission to change a Github repo. This is done with a personal access token (PAT). Note that you will need to install the packages `usethis` and `gitcreds` to do this.

```
usethis::create_github_token()
```

Call ``gitcreds::gitcreds_set()`` to register this token in the local Git credentials. It is also a great idea to store this token in any password-management software. Opening URL `'https://github.com/settings/tokens/new?scopes=repo,user,gist,workf'`

“Generate token” and store your tokens somewhere safe in your local computer as you will need this again in the future. You can additionally add PAT to your `.Renviro`n file as well. Copy it and paste it into your `.Renviro`n file as system variable `GITHUB_PAT` using

```
usethis::edit_r_environ()
```

Add to the file and save. You can also set the PAT token in R using the following.

```
#install.packages("gitcreds") # uncomment to install
gitcreds::gitcreds_set()
```

You can check that you’ve stored a credential with `gitcreds_get()`:

```
gitcreds::gitcreds_get()
```

You should get something like this:

```
...
#> <gitcreds>
#> protocol: https
#> host      : github.com
#> username: PersonalAccessToken
#> password: <-- hidden -->
...
```

Treat your PAT token like a password! For details, follow the step in Section 9.1 on this page to do this: <https://happygitwithr.com/https-pat.html>.

4.3 Individual assignments

If you followed the suggestions in the File organization in RStudio page, then you should already have an assignments folder on your computer or maize account.

Each new assignment/project will be posted as a repository on GitHub and

added directly to your account (within the Stat220 organization). This repository will contain assignment details (README, .Rmd).

4.3.1 Creating an individual assignment repo and project

1. Go to our course GitHub organization page (DataScienceWinter24) and find your homework repo, such as `hw1-username` (where your username is attached).
2. Enter the online assignment repository on GitHub. Click the green “Code” button. Most of you should just use the default setting which is to “clone” (copy) using HTTPS. Click the clipboard to the right of the URL to copy the repo location.
3. Now open up RStudio and create a project as follows:

- Click the **Project** button in the upper right corner of your RStudio

window and select **New Project...**



- Select **Version Control** and then **New Project**



- Paste the link you just copied into the Repository URL box. Leave

the Project directory name blank (or keep the auto-filled name). Use the **Browse** button to find your **assignments** folder, then click **Create Project**



4.3.2 Working on your assignment

An RStudio project should now open, which will allow you to start working on your homework assignment. You should see the project assignment name in the top right side of RStudio. You will probably see a blank console screen when you open a new project. Look in the **Files** tab for your homework .Rmd file. Click on whatever file you want to edit (probably the .Rmd file) and edit away. Make sure that your current assignment's project is the one open and showing in the upper right project name. To **open** a project, click on the **.Rproj** file or use the **Open Project...** option available in the upper right project link.

4.3.2.1 Commits

After you make changes to the homework assignment, commit them. What are commits you ask? Commits are essentially taking a snapshot of your projects. Commits save this snapshot to your local version of Git (located on your hard drive or the maize server). For example, if I make changes to a code so that it prints "Hello world", and then commit them with an informative message, I can look at the history of my commits and view the code that I wrote at that time. If I made some more changes to the function that resulted in an error, I could go back to the commit where the code was originally working. This prevents you from creating several versions of your homework (homework-v1, homework-v2, ...) or from trying to remember what your code originally looked like.

You can make commits in the Git tab in RStudio.



Click the **Commit** button in the Git tab. Check the boxes of the files that you want to commit, enter your commit message (briefly state what changes have been made), then hit **Commit**. You can read how to do this in RStudio in more detail here: <http://r-pkgs.had.co.nz/git.html#git-commit>.

Two things about committing.

- You should **commit somewhat frequently**. At minimum, if you're doing a homework assignment, you should make a commit each time that you've finished a question.
- Leave **informative commit messages**. "Added stuff" will not help you if you're looking at your commit history in a year. A message like "Added initial version of hello-world function" will be more useful.

4.3.2.2 Pushing changes to Github

At some point you'll want to get the updated version of the assignment back onto GitHub, either so that we can help you with your code or so that it can be graded. You will also want to push work frequently when you have a shared GitHub repo for project collaborations (i.e. more than one person is working on a project and code). If you are ready to push, you can again click on the "Up" **Push** arrow in the Git tab or in the Commit pop-up window or in the Git tab (shown above).

To "turn in" an assignment, all you need to do is push all your relevant files to Github by the deadline.

4.4 Group work

Collaborative Github assignments are pretty similar to individual assignments.

4.4.1 Creating a group/partner assignment repo and project

Go to our course GitHub organization page(DataScienceWinter24) and find the repo for your group, for example if your group name is "team01" the you might find the `mp1-team01` repo. Clone this repo to your computer/maize account using the same steps done for an individual assignment (see steps 2-3).

4.4.1.1 Working with collaborative repos

For group homework, I suggest that only the *recorder* edit the group-homework-x.Rmd file to avoid merge conflicts! Other group members can create a new Markdown doc to run and save commands. Only the recorder needs to **push** changes (answers) to the Github repo and all others can then **pull** these changes (i.e. the final answers) after the HW is submitted.

When you are working together on a Github project, you should commit and push your modifications frequently. You will also need to frequently **pull** updates from Github down to your local version of RStudio. These updates are changes that your teammates have made since your last pull. To pull in changes, click the “Down” **Pull** arrow in the Git tab (shown above).

If you get an error about conflict after pulling or pushing, don’t freak out! This can happen if you edit a file (usually an .Rmd or .R file) in a location that was also changed by a teammate. When this happens you should attempt to fix the **merge conflict**. Take a look at this resource site and try to fix the merge conflict in Rstudio.

4.5 Additional resources

- Happy Git and GitHub for the useR
- Rstudio, Git and GitHub
- Interactive learning guide for Git
- GitHub Guides
- Git setup for Windows (video)
- Git setup for Mac (video)
- How to clone, edit, and push homework assignments with GitHub Classroom (video)

4.6 Acknowledgements

Most of this content in this guide was taken from <https://github.com/jfiksels/github-classroom-for-students>, edited for our classroom use by Katie St. Clair.

4.7 Reuse

This guide is licensed under the CC BY-NC 3.0 Creative Commons License.

Chapter 5

R Markdown Syntax

Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

5.0.1 Lists in R Markdown:

You can use asterisk mark to provide emphasis, such as **italics** or ****bold****. You can create lists with a dash:

```
- Item 1
- Item 2
- Item 3
  + Subitem 1
* Item 4
```

to produce

- Item 1
- Item 2
- Item 3
 - Subitem 1
- Item 4

You can embed Latex equations in-line, $\frac{1}{n} \sum_{i=1}^n x_i$ to produce $\frac{1}{n} \sum_{i=1}^n x_i$ or in a new line as $\text{Var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ to produce

$$\text{Var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

5.0.2 Embed an R code chunk:

Use the following

```
```r
Use back ticks to
create a block of code
```
```

to produce:

```
Use back ticks to
create a block of code
```

You can also evaluate and display the results of R code. Each task can be accomplished in a suitably labeled chunk like the following:

```
summary(cars)
```

| speed | dist |
|---------------|----------------|
| Min. : 4.0 | Min. : 2.00 |
| 1st Qu.: 12.0 | 1st Qu.: 26.00 |
| Median : 15.0 | Median : 36.00 |
| Mean : 15.4 | Mean : 42.98 |
| 3rd Qu.: 19.0 | 3rd Qu.: 56.00 |
| Max. : 25.0 | Max. : 120.00 |

```
fit <- lm(dist ~ speed, data = cars)
fit
```

Call:

```
lm(formula = dist ~ speed, data = cars)
```

Coefficients:

| | |
|-------------|-------|
| (Intercept) | speed |
| -17.579 | 3.932 |

5.0.3 Including Plots:

You can also embed plots. See Figure 5.1 for example:

```
par(mar = c(0, 1, 0, 1))
pie(
  c(280, 60, 20),
  c('Sky', 'Sunny side of pyramid', 'Shady side of pyramid'),
  col = c('#0292D8', '#F7EA39', '#C4B632'),
  init.angle = -50, border = NA
)
```

(Credit: Yihui Xie)



Figure 5.1: A fancy pie chart.

5.0.4 Read in data files:

```
simple_data <- read.csv("https://deepbas.io/data/simple-1.dat", )
summary(simple_data)
```

```

      initials      state      age
Length:3      Length:3      Min.   :45.0
Class :character Class :character 1st Qu.:47.5
Mode  :character Mode  :character Median :50.0
                                   Mean  :52.0
                                   3rd Qu.:55.5
                                   Max.   :61.0

      time
Length:3
Class :character
Mode  :character
```

```
knitr::kable(simple_data)
```

| initials | state | age | time |
|----------|-------|-----|------|
| vib | MA | 61 | 6:01 |
| adc | TX | 45 | 5:45 |
| kme | CT | 50 | 4:19 |

5.0.5 Hide the code:

If we enter the `echo = FALSE` option in the R chunk (see the .Rmd file). This prevents the R code from being printed to your document; you just see the results.

| initials | state | age | time |
|----------|-------|-----|------|
| vib | MA | 61 | 6:01 |
| adc | TX | 45 | 5:45 |
| kme | CT | 50 | 4:19 |

Chapter 6

Github Tutorial

```
# load the required libraries
library(credentials) # to help with PAT access
library(gitcreds)
library(usethis)

# STEPS INVOLVED TO ESTABLISH GIT CREDENTIALS / PAT

# Step 1

# usethis::use_git_config(user.name = "deepbas", user.email = "deepbas99@gmail.com")

# Step 2

# usethis::create_github_token()

# Step 3

# if this is the second/subsequent iteration start from here

# gitcreds::gitcreds_set()

# Verify

# gitcreds::gitcreds_get()
```

In this worksheet, you will practice creating a GitHub repository using the `usethis::use_github()` function and cloning it back to your local machine using RStudio's menu options.

6.1 Tutorial 1: Creating and cloning a Repository starting from Github to RStudio

1. Visit the GitHub website at <https://github.com> and sign in using your GitHub account. If you don't have an account yet, you can create one for free.
2. Once logged in, click on the “+” icon in the top right corner of the web-page, then click on “New repository”.
3. Enter a name for your new repository in the “Repository name” field. You may also provide an optional description.
4. Choose the visibility of your repository by selecting either “Public” or “Private”. Public repositories are visible to anyone, while private repositories are only visible to you and any collaborators you invite.
5. (Optional) Check the box to initialize the repository with a README file.
6. Click on the “Create repository” button to create your new repository.

This will create a new GitHub repository on your Github account. Follow further to clone the repository to your local folder using RStudio.

1. Go to your GitHub repository webpage and click on the green “Code” button. This will display a dropdown menu with a URL for your repository. Click on the clipboard icon to copy the URL to your clipboard.
2. Open RStudio, and from the “File” menu, select “New Project”.
3. In the “New Project” dialog, choose “Version Control”.
4. Select “Git” as the version control system.
5. In the “Repository URL” field, paste the URL that you copied from your GitHub repository webpage.
6. Choose a local directory where you want to clone the repository by clicking on the “Browse” button and navigating to the desired folder on your computer.
7. Click on “Create Project” to clone the GitHub repository to your local computer.

6.2 Tutorial 2: Creating a new GitHub repository using usethis R package (RStudio to Github) (Works ONLY on local RStudio)

6.2.1 Prerequisites

1. Install the usethis package if you haven't already: `install.packages("usethis")`
2. Make sure you have a GitHub account, and you are logged in.
3. Configure Git with your name and email address if you haven't already. Run the following commands in the R console, replacing "Your Name" and "youremail@example.com" with your information:

```
usethis::use_git_config(user.name = "Your Name", user.email = "youremail@example.com")
```

4. Create a new R project in RStudio by clicking on "File" > "New Project" > "New Directory" > "New Project." Give your project a name and choose a location on your computer to save it. Click "Create Project."
5. Make a new file or copy and paste a .Rmd file that you want to have in your repo and save it to your requirement.
6. In the R console, load the usethis package:

```
library(usethis)
```

7. Initialize a Git repository for your project by running:

```
usethis::use_git()
```

8. Now, let's create a new GitHub repository using the `usethis::use_github()` function. Run the following command:

```
usethis::use_github()
```

9. Follow the instructions in the R console, and your GitHub repository will be created. Note the repository URL, as you will need it in the next activity.

Class Activities

Chapter 7

Class Activity 1

The R package `babynames` provides data about the popularity of individual baby names from the US Social Security Administration. Data includes all names used at least 5 times in a year beginning in 1880.

```
#install.packages("babynames") # uncomment to install  
library(babynames)
```

Below is the list for first few cases of baby names.

```
head(babynames)
```

```
# A tibble: 6 x 5  
  year sex  name      n  prop  
  <dbl> <chr> <chr>   <int> <dbl>  
1  1880 F    Mary    7065 0.0724  
2  1880 F    Anna    2604 0.0267  
3  1880 F    Emma    2003 0.0205  
4  1880 F  Elizabeth 1939 0.0199  
5  1880 F    Minnie   1746 0.0179  
6  1880 F  Margaret 1578 0.0162
```

1. How many cases and variables are in the dataset `babynames`?

Answer:

```
dim(babynames)
```

```
[1] 1924665      5
```

There are 1924665 cases and 5 variables in the dataset `babynames`.

Let's use the package `tidyverse` to do some exploratory data analysis.

```
#install.packages("tidyverse") # uncomment to install
library(tidyverse)
babynames %>% filter(name=='Aimee')
```

```
# A tibble: 150 x 5
  year sex  name      n      prop
  <dbl> <chr> <chr> <int>   <dbl>
1  1880 F    Aimee    13 0.000133
2  1881 F    Aimee    11 0.000111
3  1882 F    Aimee    13 0.000112
4  1883 F    Aimee    11 0.0000916
5  1884 F    Aimee    15 0.000109
6  1885 F    Aimee    17 0.000120
7  1886 F    Aimee    17 0.000111
8  1887 F    Aimee    18 0.000116
9  1888 F    Aimee    12 0.0000633
10 1889 F    Aimee    16 0.0000846
# i 140 more rows
```

```
filtered_names <- babynames %>% filter(name=='Aimee')
```

```
#install.packages("ggplot2") # uncomment to install
library(ggplot2)
```

```
ggplot(data=filtered_names, aes(x=year, y=prop)) +
  geom_line(aes(colour=sex)) +
  xlab('Year') +
  ylab('Prop. of Babies Named Aimee')
```

2. What do you see in the Figure 1? Explain in a few sentences.

Click for answer

Answer:

In Figure 1, we can see the proportion of babies named Aimee by year for both males and females. We notice that the name Aimee has been more popular among females than males throughout the years. There is a peak in popularity around the 1970s for female babies, and then the popularity declines.

3. Repeat question 2 to infer how does the proportion of babies with your first name trend over time. Examine the generated plot and describe the trend of your name's popularity over time. Consider the following points:

Has the popularity of your name increased, decreased, or remained stable over the years? Is there a noticeable difference in popularity between sexes? Are there any interesting patterns or trends, such as sudden increases or decreases in popularity?

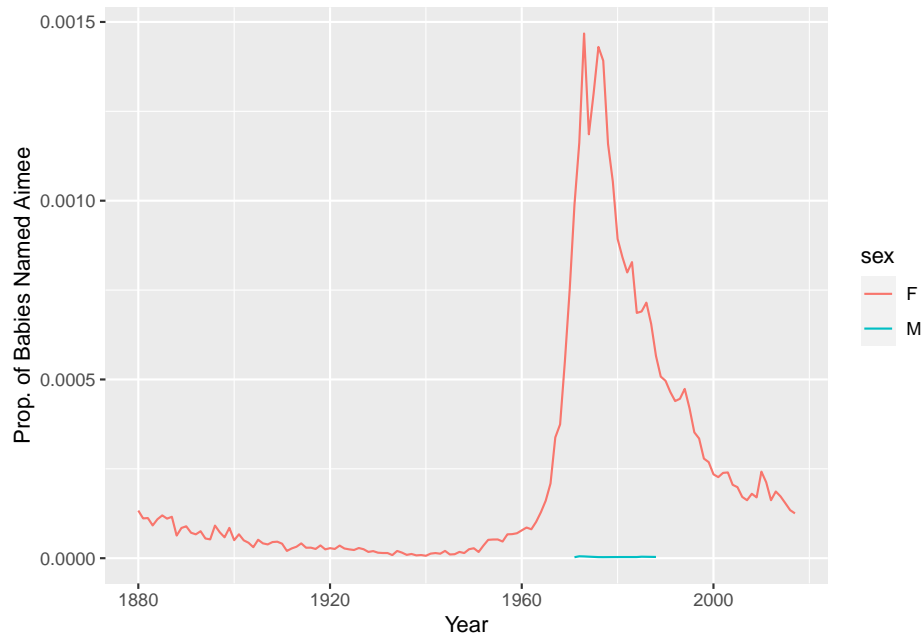


Figure 7.1: A trend chart

Answer: Answers will vary.

```
# Replace 'YourName' with your first name
your_name <- "Dee"

your_name_data <- babynames %>% filter(name == your_name)

ggplot(data=your_name_data, aes(x=year, y=prop)) +
  geom_line(aes(colour=sex)) +
  xlab('Year') +
  ylab(paste('Prop. of Babies Named', your_name))
```



4 (Optional). Compare the popularity of your first name with a randomly chosen name from the dataset. Examine the generated plot and compare the popularity of your first name with the randomly chosen name. Consider the following points:

Are there differences in popularity trends between the two names? Is one name consistently more popular than the other, or do their popularity levels change over time? Are there any interesting patterns or trends in the data, such as periods of rapid increase or decrease in popularity?

Answer Answers will vary

```
# Replace 'YourName' with your first name
your_name_data <- babynames %>% filter(name == 'Dee')

# Replace 'RandomName' with a randomly chosen name from the dataset
random_name_data <- babynames %>% filter(name == 'Max')

# Combine the two datasets
combined_data <- bind_rows(your_name_data, random_name_data)

# Plot the data
ggplot(data=combined_data, aes(x=year, y=prop)) +
  geom_line(aes(colour=sex, linetype=name)) +
  xlab('Year')
```



7.1 Extras (optional)

7.1.1 Part 1: Setting Working Directory and Loading Data

1. Set your working directory to a folder on your computer where you would like to save your R scripts and data files.

```
# Replace 'your_directory_path' with the path to your desired folder
# setwd("your_directory_path")
```

2. Load the mtcars dataset which comes preloaded with R. This dataset consists of various car features and their corresponding miles per gallon (mpg) values.

```
data(mtcars)
head(mtcars)
```

| | mpg | cyl | disp | hp | drat | wt | qsec | vs | am |
|-------------------|------|-----|------|-----|------|-------|-------|----|----|
| Mazda RX4 | 21.0 | 6 | 160 | 110 | 3.90 | 2.620 | 16.46 | 0 | 1 |
| Mazda RX4 Wag | 21.0 | 6 | 160 | 110 | 3.90 | 2.875 | 17.02 | 0 | 1 |
| Datsun 710 | 22.8 | 4 | 108 | 93 | 3.85 | 2.320 | 18.61 | 1 | 1 |
| Hornet 4 Drive | 21.4 | 6 | 258 | 110 | 3.08 | 3.215 | 19.44 | 1 | 0 |
| Hornet Sportabout | 18.7 | 8 | 360 | 175 | 3.15 | 3.440 | 17.02 | 0 | 0 |
| Valiant | 18.1 | 6 | 225 | 105 | 2.76 | 3.460 | 20.22 | 1 | 0 |

gear carb

| | | |
|-------------------|---|---|
| Mazda RX4 | 4 | 4 |
| Mazda RX4 Wag | 4 | 4 |
| Datsun 710 | 4 | 1 |
| Hornet 4 Drive | 3 | 1 |
| Hornet Sportabout | 3 | 2 |
| Valiant | 3 | 1 |

7.1.2 Part 2: Downloading Packages

1. Install the “tidyverse” package, which is a collection of useful R packages for data manipulation, exploration, and visualization.

```
# Uncomment the line below to install the package  
# install.packages("tidyverse")
```

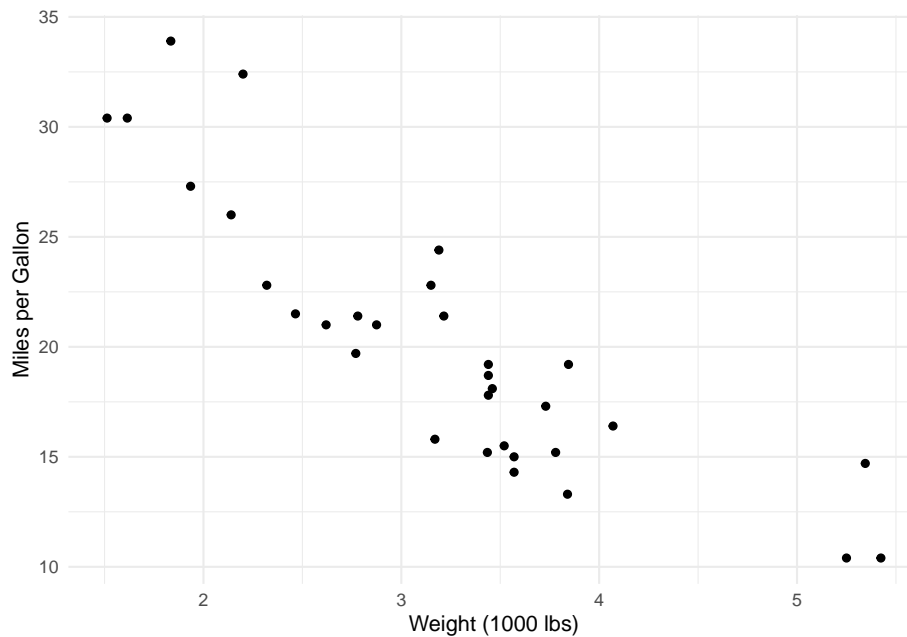
2. Load the “tidyverse” package into your R session.

```
library(tidyverse)
```

7.1.3 Part 3: Creating and Compiling an R Markdown File

1. Create a new R Markdown file in RStudio by clicking on “File” > “New File” > “R Markdown...”. Save the file in your working directory.
2. Add the following code to your R Markdown file to create a scatter plot of the mtcars dataset, showing the relationship between miles per gallon (mpg) and the weight of the car (wt).

```
# Create a scatter plot  
ggplot(data = mtcars, aes(x = wt, y = mpg)) +  
  geom_point() +  
  xlab("Weight (1000 lbs)") +  
  ylab("Miles per Gallon") +  
  theme_minimal()
```



3. Knit your R Markdown file to create an output document. Click the “Knit” button at the top of the RStudio script editor, and choose the output format you prefer (e.g., HTML, PDF, or Word).

7.2 Questions

7.2.1 1. How does the weight of a car (wt) affect its miles per gallon (mpg) based on the scatter plot you created?

[Click for answer](#)

Answer:

Based on the scatter plot, there appears to be a negative relationship between the weight of a car (wt) and its miles per gallon (mpg). As the weight of a car increases, its fuel efficiency (mpg) tends to decrease.

7.2.2 2. What is the importance of setting a working directory in R?

Click for answer

Answer:

Setting a working directory in R is important because it determines the default location where R will read from or write to when loading or saving files. This makes it easier to keep your files organized and ensures that your R scripts can access the necessary files without needing to specify the full file paths. It also simplifies sharing your R projects with others since the file paths within your scripts will be relative to the working directory.

7.2.3 3. Explain the role of R Markdown in creating reproducible research documents.

Click for answer

Answer:

R Markdown plays a crucial role in creating reproducible research documents by allowing you to combine text, code, and output (e.g., tables, figures) within a single document. This integration of narrative, data, and results makes it easier to document your data analysis process, ensuring that others can easily understand, reproduce, and build upon your work. R Markdown also supports various output formats (e.g., HTML, PDF, Word) to make it easy to share your research findings with others.

Chapter 8

Class Activity 2

Let's practice some common data assignments and manipulations in R.

- a. Create a vector of all integers from 4 to 10, and save it as **a1**.

Click for answer

```
a1 <- 4:10  
a1
```

```
[1] 4 5 6 7 8 9 10
```

- b. Create a vector of *even* integers from 4 to 10, and save it as **a2**.

Click for answer

```
a2 <- seq(4, 10, by=2)  
a2
```

```
[1] 4 6 8 10
```

- c. What do you get when you add **a1** to **a2**?

Click for answer

```
a1_plus_a2 <- a1 + a2  
a1_plus_a2
```

```
[1] 8 11 14 17 12 15 18
```

Answer: When you add **a1** to **a2**, you get a vector containing the element-wise sum: 8, 11, 14, 17, 12, 15, 18.

- d. What does the command **sum(a1)** do?

Click for answer

```
sum_a1 <- sum(a1)
sum_a1
```

```
[1] 49
```

Answer: The command `sum(a1)` calculates the sum of all elements in the vector `a1`. In this case, it returns 49.

e. What does the command `length(a1)` do?

Click for answer

```
length_a1 <- length(a1)
length_a1
```

```
[1] 7
```

Answer: The command `length(a1)` returns the number of elements in the vector `a1`. In this case, there are 7 elements.

f. Use the `sum` and `length` commands to calculate the average of the values in `a1`.

Click for answer

```
average_a1 <- sum(a1) / length(a1)
average_a1
```

```
[1] 7
```

Answer: The average of the values in `a1` is 7.

Chapter 9

Class Activity 3

```
# some interesting data objects
x <- c(3,6,9,5,10)
x.mat <- cbind(x, 2*x)
x.df <- data.frame(x=x,double.x=x*2)
my.list <- list(myVec=x, myDf=x.df, myString=c("hi","bye"))
```

9.1 Question 1: data types

- What data type is x?

Click for answer

Answer:

```
# code
typeof(x)
```

```
[1] "double"
```

- What data type is c(x, x/2)?

Click for answer

Answer:

```
# code
typeof(c(x, x/2))
```

```
[1] "double"
```

- What data type is c(x,NA)? What data type is c(x,"NA")?

Click for answer

Answer:

```
# code
typeof(c(x, NA))

[1] "double"

typeof(c(x, "NA"))

[1] "character"
```

9.2 Question 2: Subsetting and coercion

- How can we reverse the order of entries in `x`?

Click for answer

Answer:

```
# code
rev(x)

[1] 10  5  9  6  3

x[length(x):1]

[1] 10  5  9  6  3
```

- What does `which(x < 5)` equal?

Click for answer

Answer:

```
# code
which(x<5)

[1] 1
```

- Extract the element of `x` that corresponds to the location in the preceding question.

Click for answer

Answer:

```
# code
x[which(x<5)]

[1] 3
```

- What does `sum(c(TRUE,FALSE,TRUE,FALSE))` equal?

Click for answer

Answer:

```
# code  
sum(c(TRUE,FALSE,TRUE,FALSE))
```

```
[1] 2
```

- What does `sum(x[c(TRUE,FALSE,TRUE,FALSE)])` equal?

Click for answer

Answer:

```
# code  
sum(x[c(TRUE,FALSE,TRUE,FALSE, TRUE)])
```

```
[1] 22
```

- What does `sum(x < 5)` equal?

Click for answer

Answer:

```
# code  
sum(x < 5)
```

```
[1] 1
```

- What does `sum(x[x < 5])` equal?

Click for answer

Answer:

```
# code  
sum(x[x < 5])
```

```
[1] 3
```

- Why `dim(x.mat[1:2,1])` return NULL while `dim(x.mat[1:2,1:2])` returns a dimension?

Click for answer

Answer:

```
# code  
dim(x.mat[1:2,1])
```

```
NULL
```

```
dim(x.mat[1:2,1:2])
```

```
[1] 2 2
```

9.3 Question 3: Lists

- Using `my.list`, show three ways to write one command that gives the 3rd entry of variable `x` in data frame `myDf`

Click for answer

Answer:

```
# code  
my.list[[1]][3]
```

```
[1] 9
```

```
my.list[["myVec"]][3]
```

```
[1] 9
```

```
my.list[1]$myVec[3]
```

```
[1] 9
```

```
my.list$myVec[3]
```

```
[1] 9
```

- What class of object does the command `my.list[3]` return?

Click for answer

Answer:

```
# code  
class(my.list[3])
```

```
[1] "list"
```

- What class of object does the command `my.list[[3]]` return?

Click for answer

Answer:

```
# code  
class(my.list[[3]])
```

```
[1] "character"
```

- What class of object does the command `unlist(my.list)` return? Why are all the entries **characters**?

Click for answer

Answer:

```
# code  
class(unlist(my.list))
```

```
[1] "character"
```


Chapter 10

Class Activity 4

```
# Load the required libraries  
library(tidyverse)  
library(ggplot2)  
library(datasauRus)
```

10.1 Your turn 1

This worksheet will guide you through creating various plots using the `ggplot2` package in R. We will be using the `datasaurus_dozen` dataset from the `datasauRus` package for demonstration purposes. The dataset contains 13 different datasets, and we'll use them to create a variety of plots.

10.1.1 Scatterplot

a. Run the following code.

```
ggplot(data = dino_data, mapping = aes(x = x, y = y)) +  
  geom_point()
```



- b. You *must* remember to put the aesthetic mappings in the `aes()` function! What happens if you forget?

[Click for answer](#)

Answer:

If you forget to put the aesthetic mappings inside the `aes()` function, `ggplot2` will not be able to map the variables to the aesthetics correctly, and you might encounter an error or unexpected behavior in your plot.

```
# Add a layer and see what happens  
ggplot(data = dino_data , x = x , y = y)
```

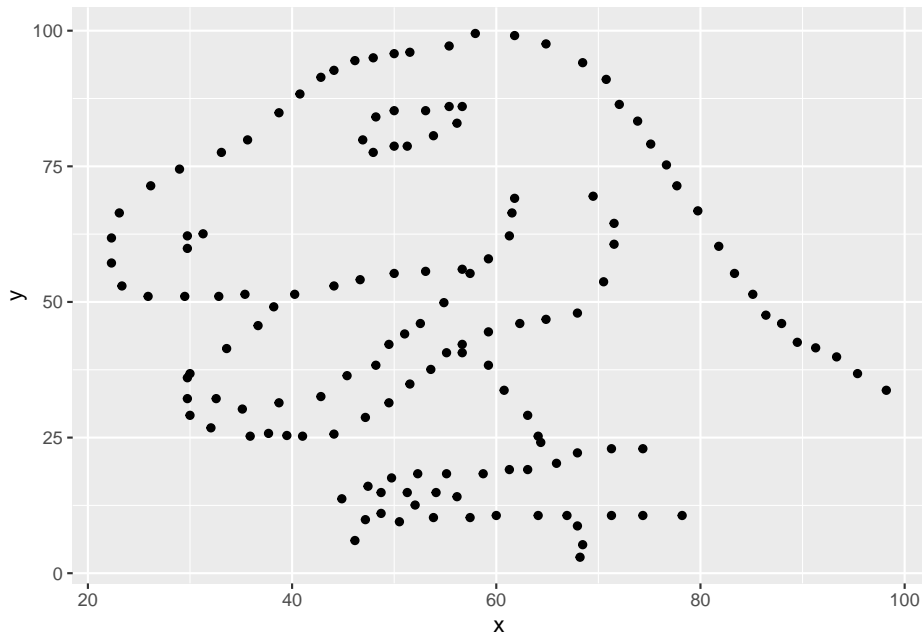


- c. The aesthetic mappings can be specified in the geom layer if you prefer, instead of the main `ggplot()` call. Give it a try:

Click for answer

Answer:

```
# Rebuild the scatterplot with your aesthetic mapping in the geom layer
ggplot(data = dino_data) +
  geom_point(aes(x = x, y = y))
```



10.1.2 Bar Plot

In this problem, we'll explore creating a bar plot using the `datasaurus_dozen` dataset.

- Create a new data frame containing the count of observations in each dataset.

Click for answer

Answer:

```
dataset_counts <- datasaurus_dozen %>%
  group_by(dataset) %>%
  summarise(count = n())
```

- Create a bar plot showing the number of observations in each dataset.

Click for answer

Answer:

```
ggplot(data = dataset_counts, aes(x = dataset, y = count)) +
  geom_bar(stat = "identity") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



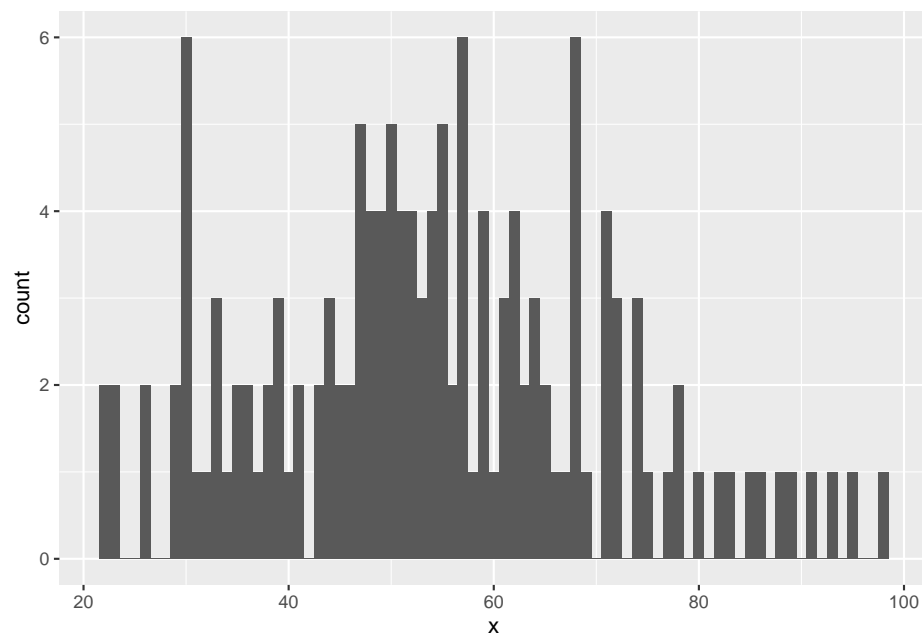
10.1.3 Histogram

- a. Create a histogram of the `x` variable for the `dino` dataset.

Click for answer

Answer:

```
ggplot(data = dino_data, aes(x = x)) +  
  geom_histogram(binwidth = 1)
```



b. Overlay a density curve on the histogram.

[Click for answer](#)

Answer:

```
ggplot(data = dino_data, aes(x = x)) +  
  geom_histogram(aes(y = after_stat(density)), binwidth = 2, fill = "lightblue") +  
  geom_density(color = "red")
```



10.1.4 Boxplot

- a. Create a boxplot of the x variable for each dataset in `datasaurus_dozen`.

Click for answer

Answer:

```
ggplot(data = datasaurus_dozen, aes(x = dataset, y = x)) +  
  geom_boxplot() +  
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



10.1.5 Faceting

Click for answer

Answer:

- Create a scatterplot of x vs. y for each dataset in `datasaurus_dozen` using `facet_wrap()`.

```
ggplot(data = datasaurus_dozen, aes(x = x, y = y)) +
  geom_point() +
  facet_wrap(~ dataset) +
  theme_minimal()
```




10.1.6 Variable Transformation

- a. The scatterplot of the `dino` dataset without any transformations is given below.

Click for answer

Answer:

```
ggplot(data = dino_data, aes(x = x, y = y)) +  
  geom_point() +  
  theme_minimal() -> p1
```

- b. Now, apply the square root transformation to both the x and y axes using the `scale_x_sqrt()` and `scale_y_sqrt()` functions in the `dino` dataset.

Click for answer

Answer:

```
ggplot(data = dino_data, aes(x = x, y = y)) +  
  geom_point() +  
  scale_x_sqrt() +  
  scale_y_sqrt() +  
  theme_minimal() -> p2
```

- c. Finally, use `grid.arrange()` function from `gridExtra` package to plot the above two plots side-by-side. Which plot do you prefer and why?

Click for answer

Answer: The second plot is more revealing of a dinosaur than the first plot.

```
library(gridExtra)
grid.arrange(p1, p2, nrow = 1)
```



10.1.7 (Optional) Lne plot

- a. Create a line plot of the `x` variable over the `y` variable for the `dino` dataset. To make it more interesting, let's first calculate the rolling mean of the `y` variable.

Click for answer

Answer:

```
dino_data <- dino_data %>%
  arrange(x) %>%
  mutate(rolling_mean_y = zoo::rollmean(y, k = 5, fill = NA))

# Line plot
ggplot(data = dino_data, aes(x = x, y = rolling_mean_y)) +
  geom_line(color = "blue") +
  theme_minimal()
```



Chapter 11

Class Activity 5

```
# Load the required libraries
library(tidyverse)
library(ggplot2)
library(ggthemes)
```

11.1 Problem 1: Changing color and shape scales

In this problem, you will learn about the effects of changing colors, scales, and shapes in `ggplot2` for both gradient and discrete color choices. You will be given a series of questions and examples to enhance your understanding. Consider the following scatter plot

```
# Generate sample data
set.seed(42)
data <- data.frame(
  Category = factor(sample(1:3, 50, replace = TRUE), labels = c("A", "B", "C")),
  X = 10 ^ rnorm(50, mean = 2, sd = 1),
  Y = rnorm(50, mean = 0, sd = 1)
)

p <- ggplot(data, aes(x = X, y = Y, color = Category)) +
  geom_point(size = 3)

p
```

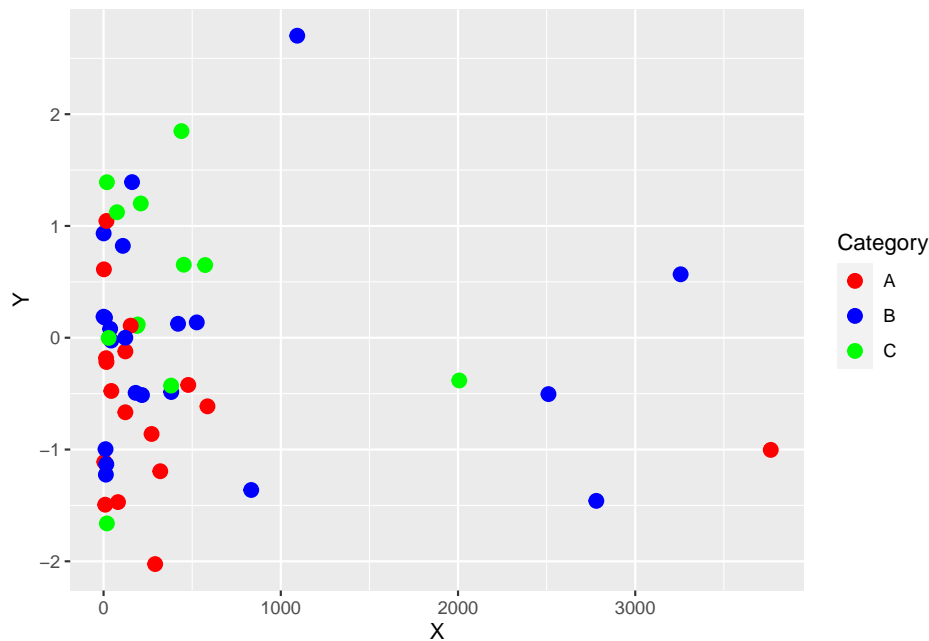


- a. Modify the scatter plot to use custom colors for each category using `scale_color_manual()`. What is the effect of changing the colors on the plot's readability?

[Click for answer](#)

Answer: Changing colors using `scale_color_manual()` allows for better distinction between categories and enhances the plot's readability.

```
p <- ggplot(data, aes(x = X, y = Y, color = Category, group = Category)) +  
  geom_point(size = 3) +  
  scale_color_manual(values = c("A" = "red", "B" = "blue", "C" = "green"))  
p
```



- b. Modify the scatter plot to use custom shapes for each category using `scale_shape_manual()`. What is the effect of changing the shapes on the plot's readability?

Click for answer

Answer: Changing the shapes using `scale_shape_manual()` helps to distinguish between categories and improves the plot's readability

```
p <- ggplot(data, aes(x = X, y = Y, shape = Category, group = Category)) +
  geom_point(size = 3) +
  scale_shape_manual(values = c("A" = 16, "B" = 17, "C" = 18))
```

p



- c. Try modifying the plot by combining color, shape, and theme customizations. Additionally, try using `geom_smooth()` to add trend lines for each category. Pay attention to how each element affects the overall readability and interpretability of the plot.

Click for answer

Answer:

```
# Base plot
p <- ggplot(data, aes(x = X, y = Y)) +
  geom_point(aes(color = Category, shape = Category), size = 3) + # Assign color and shape to Category
  geom_smooth(aes(group = Category, color = Category), method = "lm", se = FALSE) + # Add trend lines
  scale_shape_manual(values = c("A" = 19, "B" = 8, "C" = 24)) + # Customize shapes for each category
  scale_color_brewer(palette = "Dark2") + # Customize color palette
  ggthemes::theme_tufte() +
  labs(title = "Separate Trend Lines for Each Category")
```

p



11.2 Problem 2: US maps

Now, let's learn about the effect of changing various coordinate systems in `ggplot2` using a map example from the `usmap` package. We will explore the different types of coordinate systems available in `ggplot2` and how they can be applied to the map visualization.

```
#install.packages("usmap") #uncomment to install
library(usmap)
```

11.2.1 a. Plot a simple map of the United States using `ggplot2` and the `usmap` package.

Click for answer

Answer:

```
us <- plot_usmap()
us
```

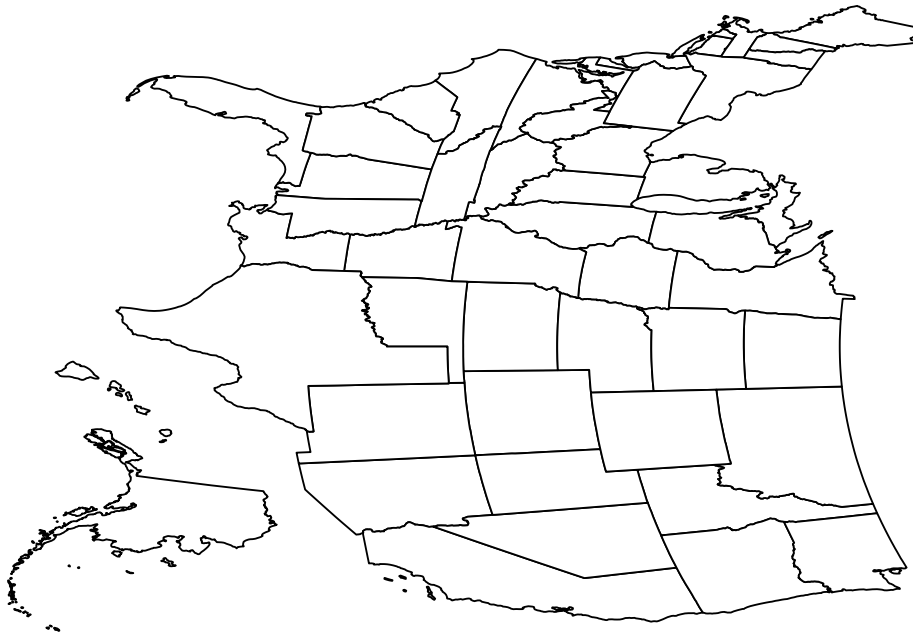


11.2.2 b. Apply the `coord_flip()` function to the map to flip the x and y axes.

[Click for answer](#)

Answer:

```
us_flipped <- us + coord_flip()  
us_flipped
```



11.2.3 c. Apply the `coord_polar()` function to the map to transform the plot to a polar coordinate system

[Click for answer](#)

Answer:

```
us_polar <- us + coord_polar()  
us_polar
```



11.2.4 d. Apply the `coord_quickmap()` function to the map to provide an approximation for a map projection.

[Click for answer](#)

Answer:

```
us_quickmap <- us + coord_quickmap()  
us_quickmap
```



11.3 Problem 3: Choropleth map

In today's class we created choropleth maps of states in the US based on ACS data.

```
states <- map_data("state")
ACS <- read.csv("https://raw.githubusercontent.com/deepbas/statdatasets/main/ACS.csv")
ACS <- dplyr::filter(ACS, !(region %in% c("Alaska", "Hawaii"))) # only 48+D.C.
ACS$region <- tolower(ACS$region) # lower case (match states regions)
```

11.3.1 (a) Mapping median income

Create a choropleth plot that uses color to create a MedianIncome map of the US.

Click for answer

Answer:

```
# map median income
ggplot(data=ACS) + coord_map() +
  geom_map(aes(map_id = region, fill = MedianIncome), map = states) +
  expand_limits(x=states$long, y=states$lat) + ggtitle("Median Income")
```



11.3.2 (b) Mapping deviations from national median income

The median income in the US in 2016 was estimated to be \$27,000. Redraw your map in (a) to visualize each state's deviation from national median income.

Click for answer

Answer:

```
# compare state income to national income
ggplot(data=ACS) + coord_map() +
  geom_map(aes(map_id = region, fill = MedianIncome - 27000), map = states) +
  expand_limits(x=states$long, y=states$lat) + ggtitle("Deviation from national median income")
```



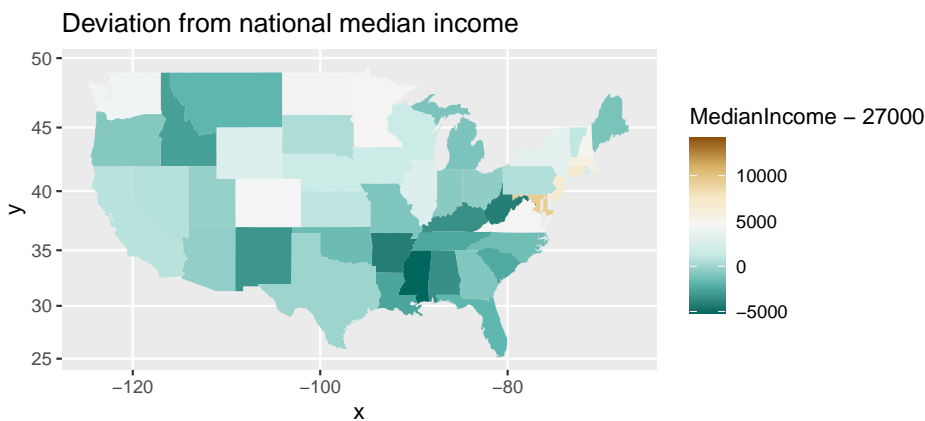
11.3.3 (c) Changing numerically scaled color

You should use a *diverging* color for (b) to highlight larger deviations from the national median. Add `scale_fill_distiller` to the map from (b) and select a diverging palette.

Click for answer

Answer:

```
# change to a diverging color
ggplot(data=ACS) + coord_map() +
  geom_map(aes(map_id = region, fill = MedianIncome - 27000), map = states) +
  expand_limits(x=states$long, y=states$lat) + ggtitle("Deviation from national median income") +
  scale_fill_distiller(type = "div")
```



11.3.4 (d) Fixing a midpoint on a diverging scale

Use `scale_fill_gradient2` to fix a midpoint scale value at white color, with diverging colors for larger positive and negative values. Apply these colors to your map in (b) and fix the midpoint at an appropriate value.

Click for answer

Answer:

```
# change to a gradient fill color
ggplot(data=ACS) + coord_map() +
  geom_map(aes(map_id = region, fill = MedianIncome - 27000), map = states) +
  expand_limits(x=states$long, y=states$lat) + ggtitle("Deviation from national median income") +
  scale_fill_gradient2(
    low = "lightblue",    # Set the low color to red
    mid = "white",        # Set the mid color to yellow
    high = "maroon",      # Set the high color to green
    midpoint = 0
  )
```



11.3.5 (e) Polygon map

```
# Merge income data with geographic information
income_data <- left_join(states, ACS, by = c("region" = "region"))
```

Next, we will use this merged data to create a polygon map that focuses on the boundaries and shapes of each state, colored by median income.

11.3.5.1 Understanding Mercator Projection

The Mercator projection is a cylindrical map projection that was widely used for navigation charts because it represents lines of constant course, known as rhumb lines, as straight segments. However, this projection distorts the size of objects as the latitude increases from the Equator to the poles. For example, Greenland appears larger than Africa on a Mercator projection map, while in reality, Africa is about 14 times larger.

For this task, you will create a polygon map to visualize the `MedianIncome` across different states using the Mercator projection. Pay attention to the shapes and sizes of states as depicted on the map.

Click for answer

```
library(sf)

ggplot(data = income_data) +
  geom_polygon(aes(x = long, y = lat, group = group, fill = MedianIncome), color = "white",
    coord_sf(crs = st_crs("+proj=merc +lon_0=0 +k=1 +x_0=0 +y_0=0 +datum=WGS84 +units=m +no_defs"),
    labs(fill = "Median Income", title = "Median Income by State") +
    theme_minimal() +
    scale_fill_viridis_c())
```


Median Income by State



11.3.6 (f) Visualizing Relative Income Deviation with Robinson Projection

The Robinson projection is a map projection of a world map which shows the entire globe as if it were flat. It was specifically created in an attempt to find a good compromise to the problem of readily showing the whole globe as a flat image. The projection is neither equal-area nor conformal, abandoning both for a compromise. The Robinson projection is widely used for thematic and educational maps due to its visually pleasing representation of the Earth.

For this task, you will visualize the relative income deviation across states using the Robinson projection. Consider how the projection's compromise between size and shape affects the presentation of income data.

Click for answer

```
# Calculate income deviation as a percentage
national_median <- 27000

# Merge the updated income data with geographic information
ACS$IncomeDeviationPercent <- ((ACS$MedianIncome - national_median) / national_median) * 100
income_data <- left_join(states, ACS, by = c("region" = "region"))

# Define the CRS for Robinson projection
robinson_crs <- st_crs("+proj=robin +lon_0=0 +x_0=0 +y_0=0 +datum=WGS84 +units=m +no_defs")

# Plot the income deviation using Robinson projection with geom_polygon
ggplot(data = income_data) +
  geom_polygon(aes(x = long, y = lat, group = group, fill = IncomeDeviationPercent), color = "white") +
  coord_sf(crs = robinson_crs, datum = NA) +
  labs(fill = "Income Deviation (%)", title = "Income Deviation from National Median by State (%)") +
  theme_minimal() +
  scale_fill_distiller(palette = "Spectral", name = "Deviation (%)")
```

Income Deviation from National Median by State (%) (Robinson Projection)



Chapter 12

Class Activity 6

```
# load the necessary libraries
library(dplyr)
library(ggplot2)
library(babynames)
```

We will work with the `babynames` dataset again in this class activity. The header of the dataset looks like this:

```
knitr::kable(head(babynames))
```

| year | sex | name | n | prop |
|------|-----|-----------|------|-----------|
| 1880 | F | Mary | 7065 | 0.0723836 |
| 1880 | F | Anna | 2604 | 0.0266790 |
| 1880 | F | Emma | 2003 | 0.0205215 |
| 1880 | F | Elizabeth | 1939 | 0.0198658 |
| 1880 | F | Minnie | 1746 | 0.0178884 |
| 1880 | F | Margaret | 1578 | 0.0161672 |

In this tutorial, we will learn about the five main verbs of `dplyr` and how to use them to manipulate data:

- `select()`: Choose columns from a data frame
- `filter()`: Choose rows based on a condition
- `arrange()`: Sort the rows of a data frame
- `mutate()`: Add new columns based on existing columns
- `summarise()`: Aggregate data and compute summary statistics

12.1 Problem 1: `select()`

Which of these is NOT a way to select the `name` and `n` columns together?

```
select(babynames, -c(year, sex, prop)) #1
select(babynames, name:n) #2
select(babynames, starts_with("n")) #3
select(babynames, ends_with("n")) #4
```

Click for answer

Answer: 4 is not the way to select the `name` and `n` columns together

12.2 Problem 2: filter()

Use `filter()` with the logical operators to extract:

12.2.1 a. All of the names where `prop` is greater than or equal to 0.08

Click for answer

```
filter(babynames, prop >= 0.08)
```

```
# A tibble: 3 x 5
  year sex  name      n  prop
  <dbl> <chr> <chr>  <int> <dbl>
1  1880 M    John   9655 0.0815
2  1880 M   William 9532 0.0805
3  1881 M    John   8769 0.0810
```

```
# alternate
babynames %>% filter(prop >= 0.08)
```

```
# A tibble: 3 x 5
  year sex  name      n  prop
  <dbl> <chr> <chr>  <int> <dbl>
1  1880 M    John   9655 0.0815
2  1880 M   William 9532 0.0805
3  1881 M    John   8769 0.0810
```

12.2.2 b. All of the babies named “Rose”

Click for answer

```
babynames %>% filter(name == "Rose")
```

```
# A tibble: 247 x 5
  year sex  name      n  prop
  <dbl> <chr> <chr>  <int> <dbl>
1  1880 F    Rose    700 0.00717
2  1880 M    Rose     7 0.0000591
```

```

3 1881 F      Rose      734 0.00743
4 1882 F      Rose      886 0.00766
5 1883 F      Rose      877 0.00730
6 1883 M      Rose         5 0.0000445
7 1884 F      Rose     1060 0.00770
8 1884 M      Rose         5 0.0000407
9 1885 F      Rose     1164 0.00820
10 1885 M      Rose         9 0.0000776
# i 237 more rows

```

12.2.3 c. Use filter() to choose all rows where name is “John” and sex is “M”.

Click for answer

```
babynames %>% filter(name == "John", sex == "M")
```

```

# A tibble: 138 x 5
   year sex   name      n  prop
<dbl> <chr> <chr> <int> <dbl>
1  1880 M     John   9655 0.0815
2  1881 M     John   8769 0.0810
3  1882 M     John   9557 0.0783
4  1883 M     John   8894 0.0791
5  1884 M     John   9388 0.0765
6  1885 M     John   8756 0.0755
7  1886 M     John   9026 0.0758
8  1887 M     John   8110 0.0742
9  1888 M     John   9247 0.0712
10 1889 M     John   8548 0.0718
# i 128 more rows

```

12.3 Problem 3: arrange()

12.3.1 a. Use arrange() to sort the babynames dataset by the prop column in descending order.

Click for answer

```
babynames %>% arrange(desc(prop))
```

```

# A tibble: 1,924,665 x 5
   year sex   name      n  prop
<dbl> <chr> <chr> <int> <dbl>
1  1880 M     John   9655 0.0815
2  1881 M     John   8769 0.0810
3  1880 M   William  9532 0.0805

```

```

4 1883 M John 8894 0.0791
5 1881 M William 8524 0.0787
6 1882 M John 9557 0.0783
7 1884 M John 9388 0.0765
8 1882 M William 9298 0.0762
9 1886 M John 9026 0.0758
10 1885 M John 8756 0.0755
# i 1,924,655 more rows

```

12.3.2 b. Use `arrange()` to sort the babynames dataset by year (ascending) and then by prop (descending).

Click for answer

```
babynames %>% arrange(year, desc(prop))
```

```

# A tibble: 1,924,665 x 5
   year sex  name      n  prop
  <dbl> <chr> <chr>  <int> <dbl>
1  1880 M   John   9655 0.0815
2  1880 M   William 9532 0.0805
3  1880 F   Mary   7065 0.0724
4  1880 M   James  5927 0.0501
5  1880 M   Charles 5348 0.0452
6  1880 M   George 5126 0.0433
7  1880 M   Frank  3242 0.0274
8  1880 F   Anna   2604 0.0267
9  1880 M   Joseph 2632 0.0222
10 1880 M   Thomas 2534 0.0214
# i 1,924,655 more rows

```

12.4 Problem 4: `mutate()`

12.4.1 a. Use `mutate()` to create a new column called decade which contains the decade the record is in (e.g., 1990 for the years 1990-1999).

Click for answer

```
babynames %>% mutate(decade = (year %/% 10) * 10)
```

```

# A tibble: 1,924,665 x 6
   year sex  name      n  prop decade
  <dbl> <chr> <chr>  <int> <dbl> <dbl>
1  1880 F   Mary   7065 0.0724  1880
2  1880 F   Anna   2604 0.0267  1880

```

```

3 1880 F      Emma      2003 0.0205 1880
4 1880 F      Elizabeth 1939 0.0199 1880
5 1880 F      Minnie    1746 0.0179 1880
6 1880 F      Margaret  1578 0.0162 1880
7 1880 F      Ida       1472 0.0151 1880
8 1880 F      Alice     1414 0.0145 1880
9 1880 F      Bertha    1320 0.0135 1880
10 1880 F     Sarah      1288 0.0132 1880
# i 1,924,655 more rows

```

12.5 Problem 5: summarize() or summarise()

Use the codes mentioned so far to compute three statistics:

- the total number of children who ever had your name
- the maximum number of children given your name in a single year
- the mean number of children given your name per year/decade (optional)

Click for answer

```

babynames %>%
  filter(name == "Dee", sex == "M")

```

```

# A tibble: 136 x 5
   year sex  name      n    prop
  <dbl> <chr> <chr> <int>  <dbl>
1  1880 M    Dee      20 0.000169
2  1881 M    Dee      32 0.000296
3  1882 M    Dee      23 0.000188
4  1883 M    Dee      22 0.000196
5  1884 M    Dee      27 0.000220
6  1885 M    Dee      28 0.000241
7  1886 M    Dee      26 0.000218
8  1887 M    Dee      39 0.000357
9  1888 M    Dee      35 0.000269
10 1889 M    Dee      24 0.000202
# i 126 more rows

```

```

babynames %>%
  filter(name == "Dee", sex == "M") %>%
  summarise(max_number = max(n))

```

```

# A tibble: 1 x 1
  max_number
  <int>
1       125

```

```
babynames %>%
  filter(name == "Dee", sex == "M") %>%
  mutate(decade = (year %/% 10) * 10) %>%
  group_by(decade) %>%
  summarise(total = sum(n),
            max = max(n),
            mean = mean(n))
```

```
# A tibble: 14 x 4
  decade total    max  mean
  <dbl> <int> <int> <dbl>
1  1880   276    39  27.6
2  1890   271    43  27.1
3  1900   302    38  30.2
4  1910   818   125  81.8
5  1920  1090   125 109
6  1930  1010   118 101
7  1940   967   120  96.7
8  1950   957   118  95.7
9  1960   683   102  68.3
10 1970   380    57   38
11 1980   217    30  21.7
12 1990   130    17   13
13 2000    87    13   9.67
14 2010    52    12   7.43
```

12.6 Problem 6

12.6.1 a. Use `min_rank()` and `mutate()` to rank each row in `babynames` from largest prop to smallest prop.

Click for answer

```
babynames %>% mutate(rank = min_rank(desc(prop))) %>% arrange(rank)
```

```
# A tibble: 1,924,665 x 6
  year sex  name    n  prop  rank
  <dbl> <chr> <chr> <int> <dbl> <int>
1  1880 M    John  9655 0.0815     1
2  1881 M    John  8769 0.0810     2
3  1880 M   William 9532 0.0805     3
4  1883 M    John  8894 0.0791     4
5  1881 M   William 8524 0.0787     5
6  1882 M    John  9557 0.0783     6
7  1884 M    John  9388 0.0765     7
8  1882 M   William 9298 0.0762     8
```



```

9 1886 M John 9026 0.0758 9
10 1885 M John 8756 0.0755 10
# i 1,924,655 more rows

```

12.6.2 b. Compute each name's rank within its year and sex.

Click for answer

```

babynames %>% group_by(year, sex) %>% mutate(rank = min_rank(desc(prop)))

# A tibble: 1,924,665 x 6
# Groups:   year, sex [276]
   year sex name n prop rank
  <dbl> <chr> <chr> <int> <dbl> <int>
1 1880 F Mary 7065 0.0724 1
2 1880 F Anna 2604 0.0267 2
3 1880 F Emma 2003 0.0205 3
4 1880 F Elizabeth 1939 0.0199 4
5 1880 F Minnie 1746 0.0179 5
6 1880 F Margaret 1578 0.0162 6
7 1880 F Ida 1472 0.0151 7
8 1880 F Alice 1414 0.0145 8
9 1880 F Bertha 1320 0.0135 9
10 1880 F Sarah 1288 0.0132 10
# i 1,924,655 more rows

```

12.6.3 c. Then compute the median rank for each combination of name and sex, and arrange the results from highest median rank to lowest.

Click for answer

```

babynames %>%
  group_by(year, sex) %>%
  mutate(rank = min_rank(desc(prop))) %>%
  group_by(name, sex) %>%
  summarize(score = median(rank)) %>%
  arrange(score)

# A tibble: 107,973 x 3
# Groups:   name [97,310]
   name sex score
  <chr> <chr> <dbl>
1 Mary F 1
2 James M 3
3 John M 3

```

| | | | |
|----|-----------|---|-----|
| 4 | William | M | 4 |
| 5 | Robert | M | 6 |
| 6 | Michael | M | 7.5 |
| 7 | Charles | M | 9 |
| 8 | Elizabeth | F | 10 |
| 9 | Joseph | M | 10 |
| 10 | Thomas | M | 11 |

i 107,963 more rows

Chapter 13

Class Activity 7

```
# load the necessary libraries  
library(tidyverse)  
library(babynames)
```

13.1 Problem 1: Boolean Operators

Use Boolean operators to alter the code below to return only the rows that contain:

13.1.1 a. Girls named Rhea

Click for answer

```
filter(babynames, name == "Rhea", sex == "F")
```

```
# A tibble: 136 x 5  
  year sex   name     n   prop  
  <dbl> <chr> <chr> <int> <dbl>  
1  1882 F     Rhea     7 0.0000605  
2  1883 F     Rhea     8 0.0000666  
3  1884 F     Rhea    13 0.0000945  
4  1885 F     Rhea    11 0.0000775  
5  1886 F     Rhea    13 0.0000846  
6  1887 F     Rhea    14 0.0000901  
7  1888 F     Rhea    20 0.000106  
8  1889 F     Rhea    31 0.000164  
9  1890 F     Rhea    39 0.000193  
10 1891 F     Rhea    24 0.000122  
# i 126 more rows
```

```
babynames %>% filter(name == "Rhea", sex == "F")
```

```
# A tibble: 136 x 5
  year sex   name     n     prop
  <dbl> <chr> <chr> <int>   <dbl>
1  1882 F     Rhea      7 0.0000605
2  1883 F     Rhea      8 0.0000666
3  1884 F     Rhea     13 0.0000945
4  1885 F     Rhea     11 0.0000775
5  1886 F     Rhea     13 0.0000846
6  1887 F     Rhea     14 0.0000901
7  1888 F     Rhea     20 0.000106
8  1889 F     Rhea     31 0.000164
9  1890 F     Rhea     39 0.000193
10 1891 F     Rhea     24 0.000122
# i 126 more rows
```

13.1.2 b. Names that were used by exactly 5 or 6 children in 1990

Click for answer

```
filter(babynames, year == 1990, n == 5 | n == 6)
```

```
# A tibble: 6,144 x 5
  year sex   name     n     prop
  <dbl> <chr> <chr> <int>   <dbl>
1  1990 F   Aariel      6 0.00000292
2  1990 F  Aarion      6 0.00000292
3  1990 F Abagael      6 0.00000292
4  1990 F  Abbye      6 0.00000292
5  1990 F  Abiola      6 0.00000292
6  1990 F Abreanna      6 0.00000292
7  1990 F  Abygail      6 0.00000292
8  1990 F  Acadia      6 0.00000292
9  1990 F Adilenne      6 0.00000292
10 1990 F  Adriena      6 0.00000292
# i 6,134 more rows
```

```
babynames %>% filter(year == "1990", n == 5 | n == 6)
```

```
# A tibble: 6,144 x 5
  year sex   name     n     prop
  <dbl> <chr> <chr> <int>   <dbl>
1  1990 F   Aariel      6 0.00000292
2  1990 F  Aarion      6 0.00000292
3  1990 F Abagael      6 0.00000292
```

```

4 1990 F      Abbye      6 0.00000292
5 1990 F      Abiola     6 0.00000292
6 1990 F      Abreanna   6 0.00000292
7 1990 F      Abygail    6 0.00000292
8 1990 F      Acadia     6 0.00000292
9 1990 F      Adilenne   6 0.00000292
10 1990 F     Adriana    6 0.00000292
# i 6,134 more rows

```

13.1.3 c. Names that are one of Apple, Yoroi, Ada

Click for answer

```
filter(babynames, name == "Apple" | name == "Yoroi" | name == "Ada")
```

```

# A tibble: 200 x 5
  year sex  name      n    prop
  <dbl> <chr> <chr> <int>  <dbl>
1  1880 F    Ada    652 0.00668
2  1881 F    Ada    628 0.00635
3  1882 F    Ada    689 0.00596
4  1883 F    Ada    778 0.00648
5  1884 F    Ada    854 0.00621
6  1885 F    Ada    876 0.00617
7  1885 M    Ada      5 0.0000431
8  1886 F    Ada    915 0.00595
9  1886 M    Ada      6 0.0000504
10 1887 F    Ada    910 0.00586
# i 190 more rows

```

13.1.4 d. Store the data tibble in part c into a new tibble and change all the character columns to upper case. Also, rename the n variable to count.

Click for answer

```

aya <- babynames %>% filter(name == "Apple" | name == "Yoroi" | name == "Ada")
aya %>% mutate_if(is.character, toupper)

```

```

# A tibble: 200 x 5
  year sex  name      n    prop
  <dbl> <chr> <chr> <int>  <dbl>
1  1880 F    ADA    652 0.00668
2  1881 F    ADA    628 0.00635
3  1882 F    ADA    689 0.00596
4  1883 F    ADA    778 0.00648
5  1884 F    ADA    854 0.00621

```

```

6 1885 F    ADA    876 0.00617
7 1885 M    ADA     5 0.0000431
8 1886 F    ADA    915 0.00595
9 1886 M    ADA     6 0.0000504
10 1887 F   ADA    910 0.00586
# i 190 more rows

```

```
aya %>% mutate_at(vars(name), toupper)
```

```

# A tibble: 200 x 5
  year sex   name     n     prop
  <dbl> <chr> <chr> <int>   <dbl>
1 1880 F    ADA    652 0.00668
2 1881 F    ADA    628 0.00635
3 1882 F    ADA    689 0.00596
4 1883 F    ADA    778 0.00648
5 1884 F    ADA    854 0.00621
6 1885 F    ADA    876 0.00617
7 1885 M    ADA     5 0.0000431
8 1886 F    ADA    915 0.00595
9 1886 M    ADA     6 0.0000504
10 1887 F   ADA    910 0.00586
# i 190 more rows

```

```
aya %>% rename(count = n)
```

```

# A tibble: 200 x 5
  year sex   name count     prop
  <dbl> <chr> <chr> <int>   <dbl>
1 1880 F    Ada    652 0.00668
2 1881 F    Ada    628 0.00635
3 1882 F    Ada    689 0.00596
4 1883 F    Ada    778 0.00648
5 1884 F    Ada    854 0.00621
6 1885 F    Ada    876 0.00617
7 1885 M    Ada     5 0.0000431
8 1886 F    Ada    915 0.00595
9 1886 M    Ada     6 0.0000504
10 1887 F   Ada    910 0.00586
# i 190 more rows

```

13.1.5 e. Change all the column names to upper case in the previous problem.

[Click for answer](#)

```
aya %>% rename_at(vars(year:prop), toupper)
```

```
# A tibble: 200 x 5
  YEAR SEX  NAME      N      PROP
  <dbl> <chr> <chr> <int>   <dbl>
1  1880 F    Ada     652 0.00668
2  1881 F    Ada     628 0.00635
3  1882 F    Ada     689 0.00596
4  1883 F    Ada     778 0.00648
5  1884 F    Ada     854 0.00621
6  1885 F    Ada     876 0.00617
7  1885 M    Ada       5 0.0000431
8  1886 F    Ada     915 0.00595
9  1886 M    Ada       6 0.0000504
10 1887 F    Ada     910 0.00586
# i 190 more rows
```

13.1.6 f. What do these commands do?

```
polluted_cities <- tribble(
  ~city, ~size, ~amount,
  "New York", "large", 23,
  "New York", "small", 14,
  "London", "large", 22,
  "London", "small", 16,
  "Beijing", "large", 121,
  "Beijing", "small", 56
)
```

```
polluted_cities
```

```
# A tibble: 6 x 3
  city      size amount
  <chr>    <chr> <dbl>
1 New York large    23
2 New York small    14
3 London  large    22
4 London  small    16
5 Beijing large   121
6 Beijing small    56
```

```
polluted_cities %>% select_if(is.numeric) #1
polluted_cities %>% rename_all(toupper) #2
polluted_cities %>% rename_if(is.character, toupper) #3
polluted_cities %>% rename_at(vars(contains("it")), toupper) #4
```

Click for answer

answer:

1. Selects all numeric columns from the polluted_cities dataset.
2. Renames all column names in the polluted_cities dataset to uppercase.
3. Renames column names with character data type in the polluted_cities dataset to uppercase.
4. Renames column names containing “it” in the polluted_cities dataset to uppercase.

```
polluted_cities %>% select_if(is.numeric) #1
```

```
# A tibble: 6 x 1
  amount
  <dbl>
1     23
2     14
3     22
4     16
5    121
6     56
```

```
polluted_cities %>% rename_all(toupper) #2
```

```
# A tibble: 6 x 3
  CITY      SIZE AMOUNT
  <chr>    <chr> <dbl>
1 New York large    23
2 New York small    14
3 London   large    22
4 London   small    16
5 Beijing  large   121
6 Beijing  small    56
```

```
polluted_cities %>% rename_if(is.character, toupper) #3
```

```
# A tibble: 6 x 3
  CITY      SIZE amount
  <chr>    <chr> <dbl>
1 New York large    23
2 New York small    14
3 London   large    22
4 London   small    16
5 Beijing  large   121
6 Beijing  small    56
```

```
polluted_cities %>% rename_at(vars(contains("it")), toupper) #4
```



```
# A tibble: 6 x 3
  CITY      size amount
  <chr>    <chr> <dbl>
1 New York large    23
2 New York small    14
3 London   large    22
4 London   small    16
5 Beijing  large   121
6 Beijing  small    56
```

Let's look at an interesting example on how to join related information on various artists, bands, songs, and their labels.

```
artists <- tibble(first = c("Jimmy", "George", "Mick", "Tom", "Davy", "John",
                           "Paul", "Jimmy", "Joe", "Elvis", "Keith", "Paul",
                           "Ringo", "Joe", "Brian", "Nancy"),
                  last = c("Buffett", "Harrison", "Jagger", "Jones", "Jones",
                           "Lennon", "McCartney", "Page", "Perry", "Presley",
                           "Richards", "Simon", "Starr", "Walsh", "Wilson", "Wilson"),
                  instrument = c("Guitar", "Guitar", "Vocals", "Vocals", "Vocals",
                                 "Guitar", "Bass", "Guitar", "Guitar", "Vocals", "Guitar",
                                 "Guitar", "Drums", "Guitar", "Vocals", "Vocals"))

bands <- tibble(first = c("John", "John Paul", "Jimmy", "Robert", "George", "John",
                           "Paul", "Ringo", "Jimmy", "Mick", "Keith", "Charlie", "Ronnie"),
                last = c("Bonham", "Jones", "Page", "Plant", "Harrison", "Lennon",
                           "McCartney", "Starr", "Buffett", "Jagger", "Richards", "Watts", "Wood"),
                band = c("Led Zeppelin", "Led Zeppelin", "Led Zeppelin", "Led Zeppelin",
                           "The Beatles", "The Beatles", "The Beatles", "The Beatles",
                           "The Coral Reefers", "The Rolling Stones", "The Rolling Stones",
                           "The Rolling Stones", "The Rolling Stones"))

albums <- tibble(album = c("A Hard Day's Night", "Magical Mystery Tour", "Beggar's Banquet",
                           "Abbey Road", "Led Zeppelin IV", "The Dark Side of the Moon", "Aerosmith",
                           "Rumours", "Hotel California"),
                 band = c("The Beatles", "The Beatles", "The Rolling Stones", "The Beatles",
                           "Led Zeppelin", "Pink Floyd", "Aerosmith", "Fleetwood Mac", "Eagles"),
                 year = c(1964, 1967, 1968, 1969, 1971, 1973, 1973, 1977, 1982))

songs <- tibble(song = c("Come Together", "Dream On", "Hello, Goodbye", "It's Not Unusual"),
                album = c("Abbey Road", "Aerosmith", "Magical Mystery Tour", "Along Came Jones"),
                first = c("John", "Steven", "Paul", "Tom"),
                last = c("Lennon", "Tyler", "McCartney", "Jones"))
```

```
labels <- tibble(album = c("Abbey Road", "A Hard Days Night", "Magical Mystery Tour",
                           "Led Zeppelin IV", "The Dark Side of the Moon", "Hotel California",
                           "Rumours", "Aerosmith", "Beggar's Banquet"),
                 label = c("Apple", "Parlophone", "Parlophone", "Atlantic", "Harvest",
                           "Asylum", "Warner Brothers", "Columbia", "Decca"))
```

Let's take a glimpse of the tibbles `artists` and `bands`. Notice that there are different number of rows in the dataset.

```
glimpse(artists)
```

```
Rows: 16
Columns: 3
$ first      <chr> "Jimmy", "George", "Mick", "Tom", "Davy~
$ last       <chr> "Buffett", "Harrison", "Jagger", "Jones~
$ instrument <chr> "Guitar", "Guitar", "Vocals", "Vocals",~
```

```
glimpse(bands)
```

```
Rows: 13
Columns: 3
$ first <chr> "John", "John Paul", "Jimmy", "Robert", "Geo~
$ last  <chr> "Bonham", "Jones", "Page", "Plant", "Harriso~
$ band  <chr> "Led Zeppelin", "Led Zeppelin", "Led Zeppeli~
```

```
glimpse(albums)
```

```
Rows: 9
Columns: 3
$ album <chr> "A Hard Day's Night", "Magical Mystery Tour"~
$ band  <chr> "The Beatles", "The Beatles", "The Rolling S~
$ year  <dbl> 1964, 1967, 1968, 1969, 1971, 1973, 1973, 19~
```

```
glimpse(songs)
```

```
Rows: 4
Columns: 4
$ song <chr> "Come Together", "Dream On", "Hello, Goodbye~
$ album <chr> "Abbey Road", "Aerosmith", "Magical Mystery ~
$ first <chr> "John", "Steven", "Paul", "Tom"
$ last  <chr> "Lennon", "Tyler", "McCartney", "Jones"
```

```
glimpse(labels)
```

```
Rows: 9
Columns: 2
$ album <chr> "Abbey Road", "A Hard Days Night", "Magical ~
```

```
$ label <chr> "Apple", "Parlophone", "Parlophone", "Atlant~
```

13.2 Problem 2: Joining artists and bands data

- 13.2.1 a. Join the artists and bands tibbles using `left_join()`, `right_join()`, and `full_join()`. Verify that the datasets obtained from `left_join()` and `right_join()` are the same using `setequal()`.

Click for answer

```
bands2 <- left_join(bands, artists, by = c("first", "last"))
bands3 <- right_join(artists, bands, by = c("first", "last"))
full_bands <- full_join(artists, bands, by = c("first", "last"))

# Check if the datasets are the same
setequal(bands2, bands3)
```

```
[1] TRUE
```

- 13.2.2 b. Use the pipe operator, `%>%`, to create one table that combines all information from artists, bands, albums, songs, and labels.

Click for answer

```
all_info <- artists %>%
  full_join(bands, by = c("first", "last")) %>%
  full_join(songs, by = c("first", "last")) %>%
  full_join(albums, by = c("album", "band")) %>%
  full_join(labels, by = c("album"))
all_info
```

A tibble: 30 x 8

| | first | last | instrument | band | song | album | year | label |
|----|--------|-----------|------------|-------|-------|-------|-------|-------|
| | <chr> | <chr> | <chr> | <chr> | <chr> | <chr> | <dbl> | <chr> |
| 1 | Jimmy | Buffett | Guitar | The ~ | <NA> | <NA> | NA | <NA> |
| 2 | George | Harrison | Guitar | The ~ | <NA> | <NA> | NA | <NA> |
| 3 | Mick | Jagger | Vocals | The ~ | <NA> | <NA> | NA | <NA> |
| 4 | Tom | Jones | Vocals | <NA> | It's~ | Alon~ | NA | <NA> |
| 5 | Davy | Jones | Vocals | <NA> | <NA> | <NA> | NA | <NA> |
| 6 | John | Lennon | Guitar | The ~ | Come~ | Abbe~ | 1969 | Apple |
| 7 | Paul | McCartney | Bass | The ~ | Hell~ | Magi~ | 1967 | Parl~ |
| 8 | Jimmy | Page | Guitar | Led ~ | <NA> | <NA> | NA | <NA> |
| 9 | Joe | Perry | Guitar | <NA> | <NA> | <NA> | NA | <NA> |
| 10 | Elvis | Presley | Vocals | <NA> | <NA> | <NA> | NA | <NA> |

```
# i 20 more rows
```

13.3 Problem 3: Filtering and counting rows in the data

13.3.1 a. Collect artists that have songs provided, and return rows of artists that don't have bands info.

Click for answer

```
# Artists with songs
artists_with_songs <- artists %>%
  semi_join(songs, by = c("first", "last"))

# Artists without bands info
artists_without_bands <- artists %>%
  anti_join(bands, by = c("first", "last"))

artists_with_songs
```

```
# A tibble: 3 x 3
  first last    instrument
  <chr> <chr>    <chr>
1 Tom   Jones    Vocals
2 John  Lennon   Guitar
3 Paul  McCartney Bass
artists_without_bands
```

```
# A tibble: 8 x 3
  first last    instrument
  <chr> <chr>    <chr>
1 Tom   Jones    Vocals
2 Davy  Jones    Vocals
3 Joe   Perry    Guitar
4 Elvis Presley Vocals
5 Paul  Simon    Guitar
6 Joe   Walsh    Guitar
7 Brian Wilson Vocals
8 Nancy Wilson Vocals
```

13.3. PROBLEM 3: FILTERING AND COUNTING ROWS IN THE DATA101

13.3.2 b. Collect the albums made by a band, count the number of rows, find the rows of songs that match a row in labels, and count the number of rows.

Click for answer

```
# Albums made by a band
albums_by_band <- albums %>% semi_join(bands, by = "band")
n_albums_by_band <- nrow(albums_by_band)

# Rows of songs that match a row in labels
songs_with_labels <- songs %>% semi_join(labels, by = "album")
n_songs_with_labels <- nrow(songs_with_labels)

n_albums_by_band
```

```
[1] 5
```

```
n_songs_with_labels
```

```
[1] 3
```


Chapter 14

Class Activity 8

```
# load the necessary libraries
library(tidyverse)
library(lubridate)
```

14.1 Your turn 1

In the provided R code, we start with two datasets, `DBP_wide` and `BP_wide`, representing blood pressure measurements in a wide format. We then demonstrate how to transform `BP_wide` into a long format using `pivot_longer()`.

```
DBP_wide <- tibble(id = letters[1:4],
  sex = c("F", "M", "M", "F"),
  v1.DBP = c(88, 84, 102, 70),
  v2.DBP = c(78, 78, 96, 76),
  v3.DBP = c(94, 82, 94, 74),
  age=c(23, 56, 41, 38)
)
DBP_wide
```

```
# A tibble: 4 x 6
  id    sex  v1.DBP v2.DBP v3.DBP  age
<chr> <chr>  <dbl>  <dbl>  <dbl> <dbl>
1 a     F      88     78     94    23
2 b     M      84     78     82    56
3 c     M     102     96     94    41
4 d     F      70     76     74    38
```

```
BP_wide <- tibble(id = letters[1:4],
  sex = c("F", "M", "M", "F"),
```

```

      SBP_v1 = c(130, 120, 130, 119),
      SBP_v2 = c(110, 116, 136, 106),
      SBP_v3 = c(112, 122, 138, 118))
BP_wide

# A tibble: 4 x 5
  id    sex  SBP_v1 SBP_v2 SBP_v3
<chr> <chr>  <dbl>  <dbl> <dbl>
1 a     F      130    110   112
2 b     M      120    116   122
3 c     M      130    136   138
4 d     F      119    106   118

BP_long <- BP_wide %>%
  pivot_longer(names_to = "visit", values_to = "SBP", SBP_v1:SBP_v3) %>%
  mutate(visit = parse_number(visit))
BP_long

# A tibble: 12 x 4
  id    sex  visit  SBP
<chr> <chr>  <dbl> <dbl>
1 a     F      1    130
2 a     F      2    110
3 a     F      3    112
4 b     M      1    120
5 b     M      2    116
6 b     M      3    122
7 c     M      1    130
8 c     M      2    136
9 c     M      3    138
10 d    F      1    119
11 d    F      2    106
12 d    F      3    118

```

14.1.1 a. Create a long dataframe from DBP_wide based on the repeated DBP columns and save it as DBP_long.

Click for answer

Answer:

```

DBP_long <- DBP_wide %>%
  pivot_longer(names_to = "visit",
               values_to = "DBP",
               cols = v1.DBP:v3.DBP)
DBP_long

```



```
# A tibble: 12 x 5
  id    sex    age visit    DBP
  <chr> <chr> <dbl> <chr>  <dbl>
1 a     F      23 v1.DBP   88
2 a     F      23 v2.DBP   78
3 a     F      23 v3.DBP   94
4 b     M      56 v1.DBP   84
5 b     M      56 v2.DBP   78
6 b     M      56 v3.DBP   82
7 c     M      41 v1.DBP  102
8 c     M      41 v2.DBP   96
9 c     M      41 v3.DBP   94
10 d    F      38 v1.DBP   70
11 d    F      38 v2.DBP   76
12 d    F      38 v3.DBP   74
```

14.1.2 b. Clean up the visit column of DBP_long so that the values are 1, 2, 3, and save it as DBP_long.

Click for answer

Answer:

```
DBP_long <- DBP_long %>%
  mutate(visit = parse_number(visit))
DBP_long
```

```
# A tibble: 12 x 5
  id    sex    age visit    DBP
  <chr> <chr> <dbl> <dbl>  <dbl>
1 a     F      23     1   88
2 a     F      23     2   78
3 a     F      23     3   94
4 b     M      56     1   84
5 b     M      56     2   78
6 b     M      56     3   82
7 c     M      41     1  102
8 c     M      41     2   96
9 c     M      41     3   94
10 d    F      38     1   70
11 d    F      38     2   76
12 d    F      38     3   74
```

14.1.3 c. Make DBP_long wide with column names visit.1, visit.2, visit.3 for the DBP values, and save it as DBP_wide2

Click for answer

Answer:

```
DBP_wide2 <- DBP_long %>%
  pivot_wider(names_from = "visit",
              values_from = "DBP",
              names_prefix = "visit.")
DBP_wide2
```

```
# A tibble: 4 x 6
  id    sex    age visit.1 visit.2 visit.3
<chr> <chr> <dbl>   <dbl>   <dbl>   <dbl>
1 a     F      23      88      78      94
2 b     M      56      84      78      82
3 c     M      41     102      96      94
4 d     F      38      70      76      74
```

14.1.4 d. Join DBP_long with BP_long2 to create a single data frame with columns id, sex, visit, SBP, DBP, and age. Save this as BP_both_long.

Click for answer

Answer:

```
BP_both_long <- left_join(BP_long, DBP_long, by = c("id", "sex", "visit"))
BP_both_long
```

```
# A tibble: 12 x 6
  id    sex  visit  SBP  age  DBP
<chr> <chr> <dbl> <dbl> <dbl> <dbl>
1 a     F      1    130   23   88
2 a     F      2    110   23   78
3 a     F      3    112   23   94
4 b     M      1    120   56   84
5 b     M      2    116   56   78
6 b     M      3    122   56   82
7 c     M      1    130   41  102
8 c     M      2    136   41   96
9 c     M      3    138   41   94
10 d    F      1    119   38   70
11 d    F      2    106   38   76
12 d    F      3    118   38   74
```

14.1.5 e. Calculate the mean SBP and DBP for each visit and save the result as `mean_BP_by_visit`.

Click for answer

Answer:

```
mean_BP_by_visit <- BP_both_long %>%
  group_by(visit) %>%
  summarize(mean_SBP = mean(SBP),
             mean_DBP = mean(DBP))
mean_BP_by_visit
```

```
# A tibble: 3 x 3
  visit mean_SBP mean_DBP
  <dbl>   <dbl>   <dbl>
1     1     125.     86
2     2     117     82
3     3     122     86
```

14.2 Your turn 2

14.2.1 a. Parsing Complex Dates: Use `dmy_hms()` to parse the following date-time string: “25-Dec-2020 17:30:00”

Click for answer

Answer:

```
parsed_date <- dmy_hms("25-Dec-2020 17:30:00")
parsed_date
```

```
[1] "2020-12-25 17:30:00 UTC"
```

14.2.2 b. Advanced Date Arithmetic: Calculate the exact age in years for someone born on “1995-05-15 09:30:00”.

Click for answer

Answer:

```
dob <- ymd_hms("1995-05-15 09:30:00")
exact_age <- as.duration(interval(dob, now())) / dyears(1)
exact_age
```

```
[1] 28.79053
```

14.2.3 c. Creating Date-Time Objects: Create a date-time object for March 15, 2020, 13:30:00 using `make_datetime()`.

Click for answer

Answer:

```
new_date_time <- make_datetime(2020, 3, 15, 13, 30, 0)
new_date_time
```

```
[1] "2020-03-15 13:30:00 UTC"
```

14.2.4 d. Extracting Components from Date-Time Objects: Extract the year, month (as a number), day, hour, and minute from “2022-07-01 14:45:00”.

Click for answer

Answer:

```
example_date_time <- ymd_hms("2022-07-01 14:45:00")
extracted_components <- tibble(
  year = year(example_date_time),
  month = month(example_date_time),
  day = day(example_date_time),
  hour = hour(example_date_time),
  minute = minute(example_date_time)
)
extracted_components
```

```
# A tibble: 1 x 5
  year month  day  hour minute
<dbl> <dbl> <int> <int> <int>
1  2022     7    1    14     45
```

14.2.5 e. Advanced Date-Time Arithmetic with Periods: Add 2 months and 15 days to “2021-08-01”.

Click for answer

Answer:

```
initial_date <- ymd("2021-08-01")
new_date <- initial_date + months(2) + days(15)
new_date
```

```
[1] "2021-10-16"
```

14.2.6 f. Duration and Time Differences: Calculate the duration in days, weeks, months, and years between “2019-04-01” and “2022-04-01”.

Click for answer

Answer:

```
start_date <- ymd("2019-04-01")
end_date <- ymd("2022-04-01")
time_diff <- end_date - start_date
duration_days <- as.duration(time_diff)
duration_weeks <- duration_days / dweeks(1)
duration_months <- duration_days / dmonths(1)
duration_years <- duration_days / dyears(1)

duration_results <- tibble(
  days = duration_days,
  weeks = duration_weeks,
  months = duration_months,
  years = duration_years
)
duration_results
```

```
# A tibble: 1 x 4
  days                weeks months years
<Duration>          <dbl>  <dbl> <dbl>
1 94694400s (~3 years)  157.    36.0  3.00
```


Chapter 15

Class Activity 9

```
# load the necessary libraries
library(tidyverse)
```

15.1 Your Turn 1

15.1.1 a) read_csv()

Use `read_csv()` to import the `desserts` data set from <https://raw.githubusercontent.com/deepbas/statdatasets/main/desserts.csv>. Use `glimpse` to see if the data import is alright.

```
url <- "https://raw.githubusercontent.com/deepbas/statdatasets/main/desserts.csv"
desserts <- read_csv(url)
glimpse(desserts)
```

```
Rows: 549
Columns: 16
$ series          <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ episode         <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ baker           <chr> "Annetha", "David", "Edd", "~
$ technical       <chr> "2nd", "3rd", "1st", "N/A", ~
$ result          <chr> "IN", "IN", "IN", "IN", "IN"~
$ uk_airdate      <chr> "17 August 2010", "17 August~
$ us_season       <dbl> NA, NA, NA, NA, NA, NA, NA, ~
$ us_airdate      <date> NA, NA, NA, NA, NA, NA, NA, ~
$ showstopper_chocolate <chr> "chocolate", "chocolate", "n~
$ showstopper_dessert <chr> "other", "other", "other", "~
$ showstopper_fruit <chr> "no fruit", "no fruit", "no ~
$ showstopper_nut  <chr> "no nut", "no nut", "no nut"~
```

```
$ signature_chocolate <chr> "no chocolate", "chocolate",~
$ signature_dessert    <chr> "cake", "cake", "cake", "cak~
$ signature_fruit      <chr> "no fruit", "fruit", "fruit"~
$ signature_nut        <chr> "no nut", "no nut", "no nut"~
```

15.1.2 b) Are there any issues with the data import? If so, what are they?

Click for answer

Answer: Based on the output of glimpse, we can see that the ‘technical’ column should be a numeric column and the ‘uk_airstate’ column should be a date column. We can also identify any issues with missing values.

your r-code

```
desserts <- read_csv(url,
  col_types = list(
    technical = col_number(),
    uk_airstate = col_date()
  )
)
```

```
problems(desserts)
```

```
# A tibble: 556 x 5
   row   col expected      actual      file
  <int> <int> <chr>         <chr>    <chr>
1     2     6 date in ISO8601 17 August 2010 ""
2     3     6 date in ISO8601 17 August 2010 ""
3     4     6 date in ISO8601 17 August 2010 ""
4     5     4 a number      N/A      ""
5     5     6 date in ISO8601 17 August 2010 ""
6     6     6 date in ISO8601 17 August 2010 ""
7     7     4 a number      N/A      ""
8     7     6 date in ISO8601 17 August 2010 ""
9     8     6 date in ISO8601 17 August 2010 ""
10    9     4 a number      N/A      ""
# i 546 more rows
```

15.1.3 c) Import the dataset with correct data types, if needed. Fix the problems, if any.

Click for answer

```
desserts <- read_csv(url,
  col_types = list(
```



```

    technical = col_number(),
    uk_airdate = col_date()
  )
)

problems(desserts)

```

```

# A tibble: 556 x 5
   row   col expected      actual      file
  <int> <int> <chr>      <chr>      <chr>
1     2     6 date in ISO8601 17 August 2010 ""
2     3     6 date in ISO8601 17 August 2010 ""
3     4     6 date in ISO8601 17 August 2010 ""
4     5     4 a number      N/A      ""
5     5     6 date in ISO8601 17 August 2010 ""
6     6     6 date in ISO8601 17 August 2010 ""
7     7     4 a number      N/A      ""
8     7     6 date in ISO8601 17 August 2010 ""
9     8     6 date in ISO8601 17 August 2010 ""
10    9     4 a number      N/A      ""
# i 546 more rows

```

```

desserts <- read_csv(url,
  col_types = list(
    technical = col_number(),
    uk_airdate = col_date(format = "%d %B %Y")
  )
)

problems(desserts)

```

```

# A tibble: 7 x 5
   row   col expected actual file
  <int> <int> <chr>      <chr> <chr>
1     5     4 a number N/A    ""
2     7     4 a number N/A    ""
3     9     4 a number N/A    ""
4    11     4 a number N/A    ""
5    35     4 a number N/A    ""
6    36     4 a number N/A    ""
7    37     4 a number N/A    ""

```

```

desserts <- read_csv(url,
  col_types = list(
    technical = col_number(),
    uk_airdate = col_date(format = "%d %B %Y")
  )
)

```

```
),
  na = c("", "NA", "N/A")
)
```

```
problems(desserts)
```

```
# A tibble: 0 x 5
# i 5 variables: row <int>, col <int>, expected <chr>,
#   actual <chr>, file <chr>
```

15.2 Your Turn 2

Use the appropriate `read_<type>()` function to import the following data sets:

- <https://deepbas.io/data/simple-1.dat>
- <https://deepbas.io/data/mild-1.csv>
- <https://deepbas.io/data/tricky-1.csv>
- <https://deepbas.io/data/tricky-2.csv>

Identify and fix any issues you encounter.

15.2.1 a) Importing simple data:

Click for answer

```
simple1 <- readr::read_csv("https://deepbas.io/data/simple-1.dat")
problems(simple1)
```

```
# A tibble: 0 x 5
# i 5 variables: row <int>, col <int>, expected <chr>,
#   actual <chr>, file <chr>
```

15.2.2 b) Importing mildly tricky data:

Click for answer

```
mild1 <- readr::read_delim("https://deepbas.io/data/mild-1.csv", delim = "|")
problems(mild1)
```

```
# A tibble: 0 x 5
# i 5 variables: row <int>, col <int>, expected <chr>,
#   actual <chr>, file <chr>
```

15.2.3 c) Importing tricky data 1:

Click for answer

```
tricky1 <- read_csv("https://deepbas.io/data/tricky-1.csv")
problems(tricky1)
```

```
# A tibble: 2 x 5
  row   col expected actual   file
<int> <int> <chr>    <chr>   <chr>
1     4     4 5 columns 4 columns ""
2     7     4 5 columns 4 columns ""
```

```
# Fix missing values
```

```
tricky1[3, ] <- c(tricky1[3, 1:2], NA, tricky1[3, 3:4])
tricky1[6, ] <- c(tricky1[4, 1], NA, tricky1[4, 3:5])
```

15.2.4 d) Importing tricky data 2:

Click for answer

```
tricky2 <- read_csv("https://deepbas.io/data/tricky-2.csv")
problems(tricky2)
```

```
# A tibble: 0 x 5
# i 5 variables: row <int>, col <int>, expected <chr>,
#   actual <chr>, file <chr>
```

```
# Fix missing values
```

```
tricky2_part1 <- read_csv("https://deepbas.io/data/tricky-2.csv", n_max = 7) %>%
  separate(city, c("city", "state"), sep = ",") %>%
  select(-c(7))
```

```
tricky2_part2 <- read_csv(
  "https://deepbas.io/data/tricky-2.csv",
  skip = 8,
  col_names = c("iata", "airport", "city", "state", "latitude", "longitude")
)
```

```
# Combine parts
```

```
data_combined <- full_join(tricky2_part1, tricky2_part2)
```


Chapter 16

Class Activity 10

```
# load the necessary libraries
library(tidyverse)
library(tidyr)
```

16.1 Your Turn 1

```
students <- tibble(
  id = 1:24,
  grade = sample(c("9th", "10th", "11th"), 24, replace = TRUE),
  region = sample(c("North America", "Europe", "Asia", "South America", "Middle East", "Africa"),
  score = round(runif(24,50, 100))
)
```

- 16.1.1 a. Create a new column `grade_fac` by converting the `grade` column into a factor. Reorder the levels of `grade_fac` to be “9th”, “10th”, and “11th”. Sort the dataset based on the `grade_fac` column.

Click for answer

Answer:

```
students_a <- students %>%
  mutate(grade_fac = factor(grade)) %>%
  mutate(grade_fac = fct_relevel(grade_fac, c("9th", "10th", "11th"))) %>%
  arrange(grade_fac)
print(students_a, n = 24)
```

```
# A tibble: 24 x 5
      id grade region      score grade_fac
  <int> <chr> <chr>      <dbl> <fct>
1     1  9th  South America    73 9th
2     2  9th  Europe          73 9th
3     3  9th  Middle East     67 9th
4     4  9th  Middle East     64 9th
5     5 14th  South America    90 9th
6     6 16th  Middle East     90 9th
7     7 19th  Middle East     94 9th
8     8  8th  Europe          85 10th
9     9 18th  South America    79 10th
10    10 20th  Middle East     60 10th
11    11 21th  Africa          78 10th
12    12  2th  Africa          90 11th
13    13  3th  Europe          86 11th
14    14  4th  South America   100 11th
15    15  9th  Africa          60 11th
16    16 10th  Africa          65 11th
17    17 11th  Europe          57 11th
18    18 12th  Europe          70 11th
19    19 13th  Europe          62 11th
20    20 15th  Middle East     90 11th
21    21 17th  North America    83 11th
22    22 22th  Middle East     76 11th
23    23 23th  Europe          92 11th
24    24 24th  Asia           51 11th
```

16.1.2 b. Create a new column `region_fac` by converting the `region` column into a factor. Collapse the `region_fac` levels into three categories: “Americas”, “EMEA” and “Asia”. Count the number of students in each collapsed region category.

Click for answer

Answer:

```
students_b <- students_a %>%
  mutate(region_fac = factor(region)) %>%
  mutate(region_collapsed = fct_collapse(region_fac,
                                          Americas = c("North America", "South America"),
                                          EMEA = c("Europe", "Middle East", "Africa"),
                                          Asia = "Asia")) %>%

  count(region_collapsed)
print(students_b)
```

```
# A tibble: 3 x 2
  region_collapsed    n
  <fct>             <int>
1 EMEA                18
2 Asia                 1
3 Americas             5
```

16.1.3 c. Create a new column `grade_infreq` that is a copy of the `grade_fac` column. Reorder the levels of `grade_infreq` based on their frequency in the dataset. Print the levels of `grade_infreq` to check the ordering.

Click for answer

Answer:

```
students_c <- students_a %>%
  mutate(grade_infreq = grade_fac) %>%
  mutate(grade_infreq = fct_infreq(grade_infreq))

levels(students_c$grade_infreq)
```

```
[1] "11th" "9th"  "10th"
```

16.1.4 d. Create a new column `grade_lumped` by lumping the least frequent level of the `grade_fac` column into an ‘Others’ category. Count the number of students in each of the categories of the `grade_lumped` column.

Click for answer

Answer:

```
students_d <- students_a %>%
  mutate(grade_lumped = fct_lump(grade_fac, n = 1, other_level = "Others")) %>%
  count(grade_lumped)
students_d
```

```
# A tibble: 2 x 2
  grade_lumped    n
  <fct>         <int>
1 11th          13
2 Others        11
```

16.2 Your Turn 2

Lets import the `gss_cat` dataset from the `forcats` library. This datast contains a sample of categorical variables from the General Social survey.

```
# import gss_cat dataset from forcats library
forcats::gss_cat
```

```
# A tibble: 21,483 x 9
   year marital      age race  rincome partyid relig denom
   <int> <fct>      <int> <fct> <fct>   <fct>   <fct> <fct>
1  2000 Never marr~    26 White $8000 ~ Ind,ne~ Prot~ Sout~
2  2000 Divorced      48 White $8000 ~ Not st~ Prot~ Bapt~
3  2000 Widowed      67 White Not ap~ Indepe~ Prot~ No d~
4  2000 Never marr~    39 White Not ap~ Ind,ne~ Orth~ Not ~
5  2000 Divorced      25 White Not ap~ Not st~ None  Not ~
6  2000 Married      25 White $20000~ Strong~ Prot~ Sout~
7  2000 Never marr~    36 White $25000~ Not st~ Chri~ Not ~
8  2000 Divorced      44 White $7000 ~ Ind,ne~ Prot~ Luth~
9  2000 Married      44 White $25000~ Not st~ Prot~ Other
10 2000 Married      47 White $25000~ Strong~ Prot~ Sout~
# i 21,473 more rows
# i 1 more variable: tvhours <int>
```

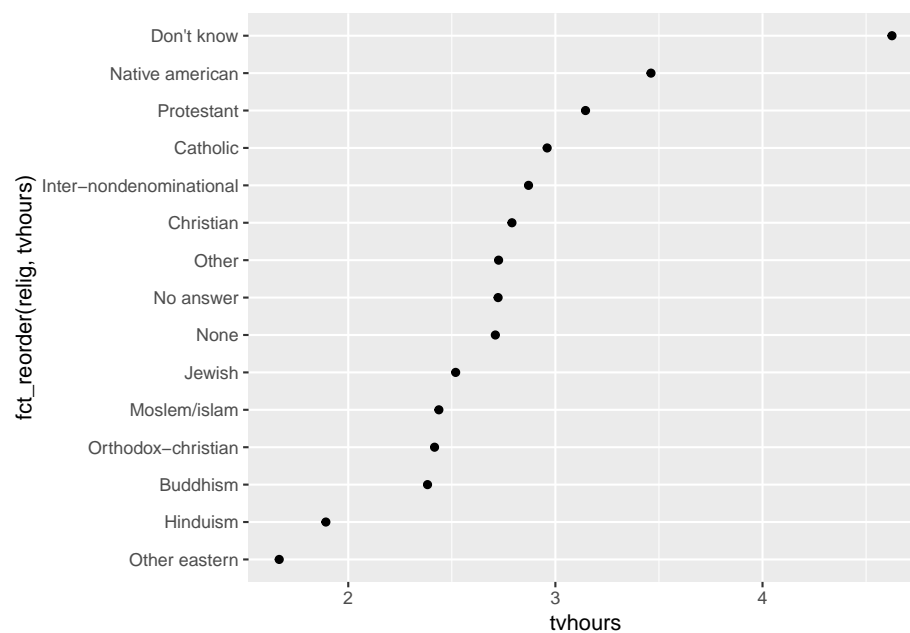
Use `gss_cat` to answer the following questions.

16.2.1 a. Which religions watch the least TV?

Click for answer

Answer:

```
# your r-code
gss_cat %>%
  drop_na(tvhours) %>%
  group_by(relig) %>%
  summarize(tvhours = mean(tvhours)) %>%
  ggplot(aes(tvhours, fct_reorder(relig, tvhours))) +
  geom_point()
```

16.2.2 b. Do married people watch more or less TV than single people?

[Click for answer](#)

Answer:

```
# your r-code
gss_cat %>%
  drop_na(tvhours) %>%
  group_by(marital) %>%
  summarize(tvhours = mean(tvhours)) %>%
  ggplot(aes(tvhours, fct_reorder(marital, tvhours))) +
  geom_point()
```



16.2.3 c. Collapse the marital variable to have levels married, not_married, and no_answer. Include "Never married", "Divorced", and "Widowed" in not_married

Click for answer

Answer:

```
# your r-code
gss_cat %>%
  drop_na(tvhours) %>%
  select(marital, tvhours) %>%
  mutate(
    maritalStatus =
      fct_collapse(
        marital,
        married = c("Married",
                     "Separated"),
        not_married = c("Never married",
                        "Divorced",
                        "Widowed"),
        no_answer = c("No answer"))
  ) -> marital_c

levels(marital_c$maritalStatus)
```

```
[1] "no_answer"  "not_married" "married"
```


Chapter 17

Class Activity 11

```
# load the necessary libraries
library(tidyverse)
library(stringr)
```

17.1 Problem 1

Let's learn about combining strings with different separators first.

```
place <- "Central Park"
activity <- "jogging"
activities <- c("jogging", "picnicking", "boating")
my_sentence <- str_c(place, " is great for ", activity, ".", sep = "")
my_sentence
```

```
[1] "Central Park is great for jogging."
```

- a. What happens when a `str_c` entry is a vector?

Click for answer

Answer: When an entry in `str_c` is a vector, it will combine the strings with each element of the vector, creating multiple combined strings.

```
my_sentences <- str_c(place, " is great for ", activities, ".", sep = "")
my_sentences
```

```
[1] "Central Park is great for jogging."
[2] "Central Park is great for picnicking."
[3] "Central Park is great for boating."
```

- b. How do you combine strings with `str_glue`?

Click for answer

Answer: You can combine strings with `str_glue` using curly braces `{}` to insert variables directly into the string.

```
my_sentence <- str_glue("{place} is great for {activity}.")
my_sentence
```

Central Park is great for jogging.

```
my_sentences1 <- str_glue("{place} is great for {activities}.")
my_sentences1
```

Central Park is great for jogging.

Central Park is great for picnicking.

Central Park is great for boating.

c. What does `str_flatten` do?

Click for answer

Answer: `str_flatten` collapses a character vector into a single string by concatenating the elements with a specified separator.

```
str_flatten(my_sentences, collapse = " and ")
```

```
[1] "Central Park is great for jogging. and Central Park is great for picnicking. and C
```

d. What will using a `\n` separator do in the command below?

Click for answer

Answer: Using a `\n` separator in the command will insert a newline character between the strings being combined, making them display on separate lines when printed.

```
p <- str_c(place, " is great for ", activity, sep = "\n")
writeLines(p)
```

Central Park

is great for

jogging

e. Does `str_length` count spaces and special characters??

Click for answer

Answer: Yes, `str_length` counts spaces and special characters as part of the string's length.

```
p
```

```
[1] "Central Park\n is great for \njogging"
```

```
str_length(p)
```

```
[1] 35
```

f. How do you count the number of e's in a string?

Click for answer

Answer: You can count the number of e's in a string using `str_count` with a pattern that matches the character 'e'.

```
text <- "The quick brown fox jumps over the lazy dog."
pattern <- "e"
vowel_count <- str_count(text, pattern)
vowel_count
```

```
[1] 3
```

g. What happens with negative positions?

Click for answer

Answer: Negative positions in `str_sub` count the positions from the end of the string rather than from the beginning.

```
str_sub(my_sentence, start = -3, end = -1)
```

```
[1] "ng."
```

h. How do you extract a substring with positive and negative positions?

Click for answer

Answer: You can extract a substring with positive and negative positions using `str_sub` and specifying the start and end positions with either positive or negative numbers.

```
my_sentence <- "Central Park is great for jogging."
positive_substr <- str_sub(my_sentence, start = 1, end = 12)
negative_substr <- str_sub(my_sentence, start = -8, end = -1)
positive_substr
```

```
[1] "Central Park"
```

```
negative_substr
```

```
[1] "jogging."
```

i. With a vector of positions?

Click for answer

Answer: Using a vector of positions with `str_sub` will extract substrings starting and ending at the specified positions in the vector.

```
str_sub(my_sentence, start = c(1, 9), end = c(4, 15))
```

```
[1] "Cent"      "Park is"
```

j. How do you extract multiple **substrings** using a vector of positions?

Click for answer

Answer: You can extract multiple **substrings** using a vector of positions with **str_sub** by specifying the start and end positions in separate vectors.

```
my_sentence <- "Central Park is great for jogging."
substrs <- str_sub(my_sentence, start = c(1, 14, 24), end = c(12, 19, 30))
substrs
```

```
[1] "Central Park" "is gre"      "or jogg"
```

17.2 Problem 2

a. Use the string parsing functions that you learned today to do tasks described in the comments below:

```
s1 <- "12%" # remove %
s2 <- "New Jersey_*" # remove _*
s3 <- "2,150" # remove comma(,)
s4 <- "Learning #datascience is fun!" # extract #datascience
s5 <- "123 Main St, Springfield, MA, 01101" # separate info
```

Click for answer

```
# Cleaning steps
s1_clean <- str_replace(s1, "%", "")
s2_clean <- str_replace(s2, "_\\*", "")
s3_clean <- str_replace(s3, ",", "")
s4_clean <- str_extract(s4, "#\\w+")
s5_clean <- str_split(s5, ",\\s?")
```

```
# Print cleaned strings
s1_clean
```

```
[1] "12"
```

```
s2_clean
```

```
[1] "New Jersey"
```

```
s3_clean
```

```
[1] "2150"
```



```
s4_clean
```

```
[1] "#datascience"
```

```
s5_clean
```

```
[[1]]
```

```
[1] "123 Main St" "Springfield" "MA" "01101"
```

17.3 Problem 3

- a. Use the string parsing functions that you learned today to do tasks described in the comments below:

```
s1 <- "25%" # remove %
s2 <- "Los Angeles_" # remove _
s3 <- "1,250" # remove comma(,)
s4 <- "Discover #machinelearning today!" # extract #machinelearning
s5 <- "456 Main St, San Francisco, CA, 94107" # separate info
```

Click for answer

```
# Cleaning steps
s1_clean <- str_replace(s1, "%", "")
s2_clean <- str_replace(s2, "_", "")
s3_clean <- str_replace(s3, ",", "")
s4_clean <- str_extract(s4, "#\\w+")
s5_clean <- str_split(s5, "\\s?")
```

```
# Print cleaned strings
```

```
s1_clean
```

```
[1] "25"
```

```
s2_clean
```

```
[1] "Los Angeles"
```

```
s3_clean
```

```
[1] "1250"
```

```
s4_clean
```

```
[1] "#machinelearning"
```

```
s5_clean
```

```
[[1]]
```

```
[1] "456 Main St" "San Francisco" "CA"
```

```
[4] "94107"
```

17.4 Problem 4

- a. Let's look at the following dataset containing information about movies and their release years. We'll extract the release year from the movie title, create a new column with decades, and count the number of movies in each decade.

```
# Sample dataset
movies <- tibble(
  title = c(
    "The Godfather (1972)", "Pulp Fiction (1994)", "The Dark Knight (2008)",
    "Forrest Gump (1994)", "The Shawshank Redemption (1994)", "The Matrix (1999)",
    "Inception (2010)", "Interstellar (2014)", "Parasite (2019)", "Fight Club (1999)"
  )
)
movies
```

```
# A tibble: 10 x 1
  title
  <chr>
1 The Godfather (1972)
2 Pulp Fiction (1994)
3 The Dark Knight (2008)
4 Forrest Gump (1994)
5 The Shawshank Redemption (1994)
6 The Matrix (1999)
7 Inception (2010)
8 Interstellar (2014)
9 Parasite (2019)
10 Fight Club (1999)
```

Click for answer

```
# Processing the dataset
movies_processed <- movies %>%
  mutate(
    release_year = as.integer(str_extract(title, "\\d{4}")),
    decade = floor(release_year / 10) * 10
  ) %>%
  count(decade) %>%
  rename(num_movies = n)

# Print the processed dataset
movies_processed
```

```
# A tibble: 4 x 2
  decade num_movies
  <dbl>      <int>
1 1970         1
2 1990         5
3 2000         1
4 2010         3
```

| | | |
|---|------|---|
| 1 | 1970 | 1 |
| 2 | 1990 | 5 |
| 3 | 2000 | 1 |
| 4 | 2010 | 3 |

Chapter 18

Class Activity 12

```
# load the necessary libraries
library(stringr)
library(dplyr)
library(readr)
```

In this tutorial, we will learn about string manipulations using regular expressions and the `stringr` library in R. We will cover different examples and use cases to help you understand the concepts and functions related to string manipulation.

18.1 Group Activity 1

```
x <- "My SSN is 593-29-9502 and my age is 55"
y <- "My phone number is 612-643-1539"
z <- "My old SSN number is 39532 9423."
out <- str_flatten(c(x,y,z), collapse = ". ")
```

18.1.1 a. What characters in `x` will `str_view_all(x, "-.-")` find?

Click for answer

answer:

The pattern searches for a dash, followed by any two characters, followed by another dash. In `x`, it finds “-29-” which is a part of the SSN.

```
str_view_all(x, "-.-")
```

```
[1] | My SSN is 593<-29->9502 and my age is 55
```

18.1.2 b. What pattern will `str_view_all(x, "-\\d{2}-")` find?

Click for answer

answer:

The pattern searches for a dash, followed by two digits, followed by another dash. In `x`, it finds the same “-29-” as in the previous example, which is a part of the SSN.

```
str_view_all(x, "-\\d{2}-") # "-" then 2 digits then "-"
```

```
[1] | My SSN is 593<-29->9502 and my age is 55
```

18.1.3 c. What pattern will `str_view_all(out, "\\d{2}\\.*")` find?

Click for answer

answer:

The pattern searches for two digits followed by an optional period. In `out`, it finds “55” and “55.”, which represent the age in the first sentence.

```
str_view_all(out, "\\s\\d{2}\\.*") # 2 digits then "."
```

```
[1] | My SSN is 593-29-9502 and my age is< 55.> My phone number is 612-643-1539. My ol
```

18.1.4 d. Use `str_view_all` to determine the correct regex pattern to identify all SSN in `out`

We can get the SSN with the usual format (###-##-####) with a regex that has 3, 2, and 4 digits separated by a dash.

```
str_view_all(out, "[0-8]\\d{2})-(\\d{2})-(\\d{4})")
```

```
[1] | My SSN is <593-29-9502> and my age is 55. My phone number is 612-643-1539. My ol
```

This misses the oddly formatted SSN in the third entry. Rather than use a dash, we can specify the divider as `[-\\s]?` which allows either 0 or 1 occurrences of either a dash or space divider:

```
str_view_all(out, "[0-8]\\d{2})[-\\s]?((\\d{2})[-\\s]?((\\d{4}))")
```

```
[1] | My SSN is <593-29-9502> and my age is 55. My phone number is 612-643-1539. My ol
```

Click for answer

answer:

The first pattern finds the SSNs in the standard format (###-##-####) by searching for 3 digits, a dash, 2 digits, another dash, and 4 digits. The second

pattern does the same but allows for a space instead of a dash as a divider. It finds all SSNs in out, including the oddly formatted one in the third sentence.

18.1.5 e. Write a regular expression to extract dates in the format YYYY-MM-DD from a given text.

```
date_pattern <- "\\d{4}-\\d{2}-\\d{2}"
text <- "The event will take place on 2023-07-20 and end on 2023-07-22."
str_extract_all(text, date_pattern)
```

```
[[1]]
[1] "2023-07-20" "2023-07-22"
```

Click for answer

Answer: The pattern searches for 4 digits, a dash, 2 digits, another dash, and 2 digits. In the given text, it finds the dates “2023-07-20” and “2023-07-22”.

18.1.6 f. Write a regular expression to extract all words that start with a capital letter in a given text.

```
capital_pattern <- "\\b[A-Z][a-zA-Z]*\\b"
text <- "Alice and Bob went to the Market to buy some Groceries."
str_extract_all(text, capital_pattern)
```

```
[[1]]
[1] "Alice"      "Bob"        "Market"     "Groceries"
```

Click for answer

Answer: The pattern searches for a word boundary, followed by an uppercase letter, and then any sequence of letters. In the given text, it finds the words “Alice”, “Bob”, “Market”, and “Groceries”.

18.2 Group Activity 2

18.2.1 a. Let’s deal with a number string that is longer than 9 digits.

```
ssn <- "([0-8]\\d{2})[-\\s]?\\d{2})[-\\s]?\\d{4}"
test <- c("123-45-67890", "1123 45 6789")
str_view_all(test, ssn)
```

```
[1] | <123-45-6789>0
[2] | 1<123 45 6789>
```

This example captures a 9-digit string as an SSN, but these strings are longer than 9 digits and may not represent an SSN. One way to deal with this is to use the negative lookbehind `?<!` and negative lookahead `?!` operators to ensure that the identified 9-digit string does not have a leading 0 or does not contain more digits.

If we “look behind” from the start of the SSN, we should not see another digit:

```
str_view_all(test, "(?<!\d)([0-8]\d{2})[-\.\s]?(\d{2})[-\.\s]?(\d{4})")
```

```
[1] | <123-45-6789>0
[2] | 1123 45 6789
```

And if we “look ahead” from the end of the SSN, we should not see another digit:

```
str_view_all(test, "(?<!\d)([0-8]\d{2})[-\.\s]?(\d{2})[-\.\s]?(\d{4})(?!\\d)")
```

```
[1] | 123-45-67890
[2] | 1123 45 6789
```

For parts b and c, consider the following string.

```
string1 <- "100 dollars 100 pesos"
```

18.2.2 b. Explain why the following matches the first 100 and not the second.

Click for answer

answer: It looks for one or more digits followed by a space and `dollars`

```
str_view(string1, "\\d+(?= dollars)")
```

```
[1] | <100> dollars 100 pesos
```

18.2.3 c. Explain why the following matches the second 100 and not the first.

Click for answer

answer: It looks for one or more digits not followed by either a digit or space followed by `dollars`

```
str_view(string1, "\\d+(?!\\d| dollars)")
```

```
[1] | 100 dollars <100> pesos
```

For parts d and e, please take a look at `string2`.

```
string2 <- "USD100 PES0100"
```


18.2.4 d. Explain why the following matches the first 100 and not the second.

Click for answer

answer: It looks for exactly 3 digits preceded by USD

```
str_view(string2, "(?<=USD)\\d{3}")
```

```
[1] | USD<100> PES0100
```

18.2.5 e. Explain why the following matches the second 100 and not the first.

Click for answer

answer: It looks for exactly 3 digits that is not preceded by USD

```
str_view(string2, "(?!USD)\\d{3}")
```

```
[1] | USD100 PES0<100>
```

18.3 Group Activity 3

```
tweets<- read_csv("https://raw.githubusercontent.com/deepbas/statdatasets/main/TrumpTweetData.csv")
```

18.3.1 a. What proportion of tweets (text) mention “America”?

```
tweets %>%  
  summarize(prop = mean(str_detect(str_to_title(text), "America")))
```

```
# A tibble: 1 x 1  
  prop  
  <dbl>  
1 0.0926
```

Click for answer

Answer: About 10% of tweets mention “America”.

18.3.2 b. What proportion of these tweets include “great”?

```
tweets %>% filter(str_detect(str_to_title(text), "America")) %>%  
  summarize(prop = mean(str_detect(str_to_lower(text), "great")))
```

```
# A tibble: 1 x 1
  prop
<dbl>
1 0.4
```

Click for answer

Answer: About 40% of tweets mention “great”.

18.3.3 c. What proportion of the tweets mention @?

```
tweets %>% mutate(ct = str_count(text, "@")) %>%
  select(text, ct) %>%
  summarize(prop = mean(ct>0))
```

```
# A tibble: 1 x 1
  prop
<dbl>
1 0.317
```

Click for answer

Answer: About 32% of tweets mention @.

18.3.4 d. Remove the tweets having mentions @.

```
Mentions <- c("@[^\\s]+")
```

```
tw_noMentions <- tweets %>% mutate(textNoMentions = str_replace_all(text, Mentions, ""))
tw_noMentions$text[38]
```

```
[1] "My daughter @IvankaTrump will be on @Greta tonight at 7pm. Enjoy! https://t.co/QySC5PLFMMy"
tw_noMentions$textNoMentions[38]
```

```
[1] "My daughter will be on tonight at 7pm. Enjoy! https://t.co/QySC5PLFMMy"
```

Click for answer

Answer: @: This part of the pattern matches the “@” symbol, which usually indicates the beginning of a mention in a tweet. `[^\\s]+`: This part of the pattern matches one or more characters that are NOT whitespaces. The `^` inside the square brackets `[]` negates the character class (meaning it matches any character that is NOT in the specified class). The double backslash `\\` is used to escape the backslash in the R string, so the pattern `\\s` represents the whitespace character class `\\s`. Finally, the `+` indicates that the pattern should match one or more occurrences of the non-whitespace characters. Together, this

regular expression pattern `@[^\s]+` matches any mention in a tweet, which usually starts with “@” followed by one or more non-whitespace characters.

18.3.5 e. What poportion of tweets originated from an iPhone?

```
tweets %>% group_by(source) %>% summarize(count = n()) %>%  
  mutate(prop = count / sum(count)) %>% filter(source == "iPhone")
```

```
# A tibble: 1 x 3  
  source count  prop  
  <chr>   <int> <dbl>  
1 iPhone    628 0.415
```

Click for answer

Answer: About 42% of the tweets originated from an iPhone.

Chapter 19

Class Activity 13

In-class Midterm 1 !!

Chapter 20

Class Activity 14

```
# load the necessary libraries
library(wordcloud)
library(reshape2)
library(tidyverse)
library(tidyr)
library(tidytext)
library(dplyr)
```

20.1 Group Activity 1

20.1.1 a. Variance and Skewness

The variance of a random variable x is defined as:

$$\text{Var}(x) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

where $x_i = (\sum_i^n x_i)/n$ is the sample mean. Also, the skewness of the random variable x is defined as:

$$\text{Skew}(x) = \frac{\frac{1}{n-2} \left(\sum_{i=1}^n (x_i - \bar{x})^3 \right)}{\text{Var}(x)^{3/2}}$$

Please write functions to calculate the variance and skewness of $\{12, 45, 54, 34, 56, 30, 67, \text{NA}\}$.

```
x <- c(12, 45, 54, 34, 56, 30, 67, NA)
```

Click for answer

```
# function to calculate the variance of a vector
var <- function(x){
  x <- na.omit(x) # omit NA values
  sum((x - mean(x)) ^ 2) / (length(x) - 1)
}
```

```
var(x)
```

```
[1] 346.619
```

```
# function to calculate the skewness of a vector
skewness <- function(x) {
  x <- na.omit(x) # omit NA values
  sum((x - mean(x)) ^ 3) / ((length(x) - 2) * var(x) ^ (3 / 2))
}
```

```
skewness(x)
```

```
[1] -0.3930586
```

20.1.2 b. (Optional) Conditions and breaks

Create a function that iterates through a numeric vector and stops when it encounters the first negative number, returning the position of that negative number. If there are no negative numbers in the vector, the function should return a message stating that there are no negative numbers.

Click for answer

```
find_first_negative <- function(x) {
  negative_positions <- which(x < 0)

  if (length(negative_positions) > 0) {
    return(paste("The first negative number is at position", negative_positions[1]))
  } else {
    return("There are no negative numbers in the vector")
  }
}
```

```
test_vector <- c(5, 12, -7, 20, 15)
find_first_negative(test_vector)
```

```
[1] "The first negative number is at position 3"
```

20.2 Group Activity 2

```
musical_instr_reviews <- read_csv("https://raw.githubusercontent.com/deepbas/statdatasets/main/musical_instr_reviews.csv")
  rename(ratingOverall = overall)
glimpse(musical_instr_reviews)
```

```
Rows: 10,261
Columns: 3
$ reviewerName <chr> "cassandra tu \"Yeah, well, that's j~
$ reviewText <chr> "not much to write about here but it~
$ ratingOverall <dbl> 5, 5, 5, 5, 5, 5, 5, 3, 5, 5, 5, 4, ~
```

20.2.1 a. Write a function to filter the dataset based on the provided rating:

Click for answer

```
filter_reviews_by_rating <- function(data, rating) {
  data %>% filter(ratingOverall == rating)
}
```

20.2.2 b. Write a function to process the text and remove stop words:

Click for answer

```
process_text <- function(data) {
  data %>%
    select(reviewText) %>%
    unnest_tokens(output = word, input = reviewText) %>%
    anti_join(stop_words)
}
```

20.2.3 c. Write a function to join the processed text with sentiment data and create a word count table.

Click for answer

```
create_word_count_table <- function(data) {
  data %>%
    inner_join(get_sentiments("bing")) %>%
    count(word, sentiment, sort = TRUE) %>%
    reshape2::acast(word ~ sentiment, value.var = "n", fill = 0)
}
```

- 20.2.4 d. Create the final function that takes the rating and number of words as input arguments and returns a word cloud plot.

Click for answer

```
word_cloud <- function(rating, num.words) {  
  rating <- as.numeric(rating)  
  num.words <- as.numeric(num.words)  
  
  if (rating >= 1 & rating <= 5) {  
    filtered_reviews <- filter_reviews_by_rating(musical_instr_reviews, rating)  
    processed_reviews <- process_text(filtered_reviews)  
    word_count_table <- create_word_count_table(processed_reviews)  
  
    comparison.cloud(  
      word_count_table,  
      colors = c("blue", "purple"),  
      scale = c(2, 0.5),  
      max.words = num.words,  
      title.size = 2  
    )  
  } else {  
    warning(" Please enter a rating from 1 to 5")  
  }  
}  
  
word_cloud(rating = "4", num.words = 300)
```



Chapter 21

Class Activity 15

```
# load the necessary libraries
library(tidyverse)
library(dplyr)
library(stringr)

energy <- read_csv("https://raw.githubusercontent.com/deepbas/statdatasets/main/energy.csv",
  col_type = cols(
    .default = col_double(),
    Timestamp = col_datetime(format = ""),
    dayWeek = col_factor(levels=c("Mon", "Tues", "Wed", "Thurs", "Fri", "Sat", "Sun"))
  ))
```

21.1 Group Activity 1

21.1.1 a. if and for loop

Write a for loop to iterate over the columns of the 'energy' dataset and print the names of all columns containing the string "House". Please use the function `colnames()` to extract the column names and store the results in a list.

Click for answer

Answer:

```
# Create an empty list to store the column names
house_columns <- list()

# Iterate over the columns of the 'energy' dataset
for (i in seq_along(colnames(energy))) {
```

```
col_name <- colnames(energy)[i]

# Check if the column name contains the string "House"
if (str_detect(col_name, "House")) {
  # Add the column name to the list
  house_columns[[length(house_columns) + 1]] <- col_name
}
}

# Print the list of house columns
house_columns <- unlist(house_columns)
house_columns
```

```
[1] "Allen_House"
[2] "Alumni_Guest_House/Johnson_House"
[3] "Benton_House"
[4] "Berg_House"
[5] "Bird_House"
[6] "Chaney_House"
[7] "Clader_House"
[8] "Dacie_Moses_House"
[9] "Douglas_House"
[10] "Farm_House"
[11] "Geffert_House"
[12] "Headley_House"
[13] "Henrickson_House"
[14] "Henry_House"
[15] "Hill_House"
[16] "Hilton_House"
[17] "Hoppin_House_(Alumni)"
[18] "Huntington_House"
[19] "Jewett_House"
[20] "Jones_House"
[21] "Nutting_House"
[22] "Page_House_West"
[23] "Parish_House_"
[24] "Parr_House"
[25] "Pollock_House"
[26] "Prentice_House"
[27] "Rayment_House"
[28] "Rice_House"
[29] "Rogers_House"
[30] "Ryberg_House"
[31] "Seccombe_House"
[32] "Sperry_House"
```

```
[33] "Stimson_House"  
[34] "Strong_House"  
[35] "Whittier_House"  
[36] "Wilson_House"
```

21.1.2 b. for loop and mean

Using a for loop, calculate and print the mean energy consumption of houses you identified in part a.

[Click for answer](#)

Answer:

```
# Assuming the house_columns vector from the previous step  
  
# Create an empty numeric vector to store the mean energy consumption  
mean_energy_consumption <- numeric()  
  
# Iterate over the house_columns vector  
for (house_col in house_columns) {  
  # Calculate the mean energy consumption for the current house column  
  mean_val <- mean(energy[[house_col]], na.rm = TRUE)  
  
  # Add the mean energy consumption to the vector  
  mean_energy_consumption <- c(mean_energy_consumption, mean_val)  
}  
  
# Combine the house names and mean energy consumption into a dataframe  
house_mean_energy <- bind_cols(House = house_columns, MeanEnergyConsumption = mean_energy_consumption)  
  
# Print the dataframe  
house_mean_energy %>% knitr::kable()
```

| House | MeanEnergyConsumption |
|----------------------------------|-----------------------|
| Allen_House | 0.9821865 |
| Alumni_Guest_House/Johnson_House | 20.2631152 |
| Benton_House | 1.8849290 |
| Berg_House | 1.3174340 |
| Bird_House | 2.3222680 |
| Chaney_House | 1.0715123 |
| Clader_House | 0.4646776 |
| Dacie_Moses_House | 1.2776465 |
| Douglas_House | 0.7219500 |
| Farm_House | 5.0599020 |
| Geffert_House | 0.9360400 |
| Headley_House | 1.4555605 |
| Henrickson_House | 3.4407858 |
| Henry_House | 1.3639619 |
| Hill_House | 1.4735884 |
| Hilton_House | 0.4248030 |
| Hoppin_House_(Alumni) | 1.8760474 |
| Huntington_House | 1.2395238 |
| Jewett_House | 0.8987697 |
| Jones_House | 0.8680271 |
| Nutting_House | 4.3967234 |
| Page_House_West | 1.8923490 |
| Parish_House__ | 12.6793378 |
| Parr_House | 9.7210618 |
| Pollock_House | 1.1831426 |
| Prentice_House | 0.9089497 |
| Rayment_House | 0.8005664 |
| Rice_House | 1.1568457 |
| Rogers_House | 0.5634289 |
| Ryberg_House | 1.0729988 |
| Seccombe_House | 2.6874199 |
| Sperry_House | 0.7052983 |
| Stimson_House | 2.0659904 |
| Strong_House | 2.5410595 |
| Whittier_House | 1.0424369 |
| Wilson_House | 1.0435830 |

21.2 Group Activity 2

1. Make a data frame of quantiles for `energy` buildings in columns 9-90 (you will need `na.rm = TRUE`)

Click for answer

Answer:

```
qdf <- energy %>% select(9:90) %>%
  map_dfc(quantile, probs = seq(.1,.9,.1), na.rm = TRUE)
qdf

# A tibble: 9 x 82
  `100_Nevada_Street` `104_Maple_St.` `106_Winona_St.`
    <dbl>           <dbl>           <dbl>
1      0.0972        1.04           0.601
2      0.120         1.11           0.632
3      0.183         1.18           0.673
4      0.461         1.18           0.681
5      0.710         1.42           0.692
6      0.795         1.42           0.865
7      0.915         1.54           1.10
8      1.11          1.56           1.20
9      1.24          1.67           1.27
# i 79 more variables: Allen_House <dbl>,
#   `Alumni_Guest_House/Johnson_House` <dbl>,
#   Arboretum_Office <dbl>, Art_Studios <dbl>,
#   Benton_House <dbl>, Berg_House <dbl>, Bird_House <dbl>,
#   Boliou_Memorial_Art_Bldg. <dbl>, Burton_Hall <dbl>,
#   `Cassat_Hall/_James_Hall` <dbl>,
#   `Center_for_Mathematics_&_Computing` <dbl>, ...
```

2. Add a variable to identify the quantile

Click for answer

Answer:

```
qdf <- energy %>% select(9:90) %>%
  map_dfc(quantile, probs = seq(.1,.9,.1), na.rm = TRUE) %>%
  mutate(stat = str_c("quantile_", seq(10,90,10)))
qdf

# A tibble: 9 x 83
  `100_Nevada_Street` `104_Maple_St.` `106_Winona_St.`
    <dbl>           <dbl>           <dbl>
1      0.0972        1.04           0.601
2      0.120         1.11           0.632
3      0.183         1.18           0.673
4      0.461         1.18           0.681
5      0.710         1.42           0.692
6      0.795         1.42           0.865
7      0.915         1.54           1.10
```

```

8           1.11           1.56           1.20
9           1.24           1.67           1.27
# i 80 more variables: Allen_House <dbl>,
#   `Alumni_Guest_House/Johnson_House` <dbl>,
#   Arboretum_Office <dbl>, Art_Studios <dbl>,
#   Benton_House <dbl>, Berg_House <dbl>, Bird_House <dbl>,
#   Boliou_Memorial_Art_Bldg. <dbl>, Burton_Hall <dbl>,
#   `Cassat_Hall/_James_Hall` <dbl>,
#   `Center_for_Mathematics_&Computing` <dbl>, ...

```

3. Reshape the data frame to make variables `stat` (describing the quantile), `building` and `quant` (quantile value)

Click for answer

Answer:

```

qdf <- energy %>% select(9:90) %>%
  map_dfc(quantile, probs = seq(.1,.9,.1), na.rm = TRUE) %>%
  mutate(stat = str_c("quantile_", seq(10,90,10))) %>%
  pivot_longer(names_to = "building", values_to = "quantiles", 1:82)
qdf

```

```

# A tibble: 738 x 3
   stat      building      quantiles
  <chr>    <chr>         <dbl>
1 quantile_10 100_Nevada_Street  0.0972
2 quantile_10 104_Maple_St.    1.04
3 quantile_10 106_Winona_St.    0.601
4 quantile_10 Allen_House 0.756
5 quantile_10 Alumni_Guest_House/Johnson_House 17.0
6 quantile_10 Arboretum_Office 0.13
7 quantile_10 Art_Studios 0.23
8 quantile_10 Benton_House 1.59
9 quantile_10 Berg_House 1.06
10 quantile_10 Bird_House 1.42
# i 728 more rows

```

4. Plot the KWH value for each quantile on the x-axis for the buildings Sayles-Hill, Language_&_Dining_Center, Olin_Hall_of_Science

Click for answer

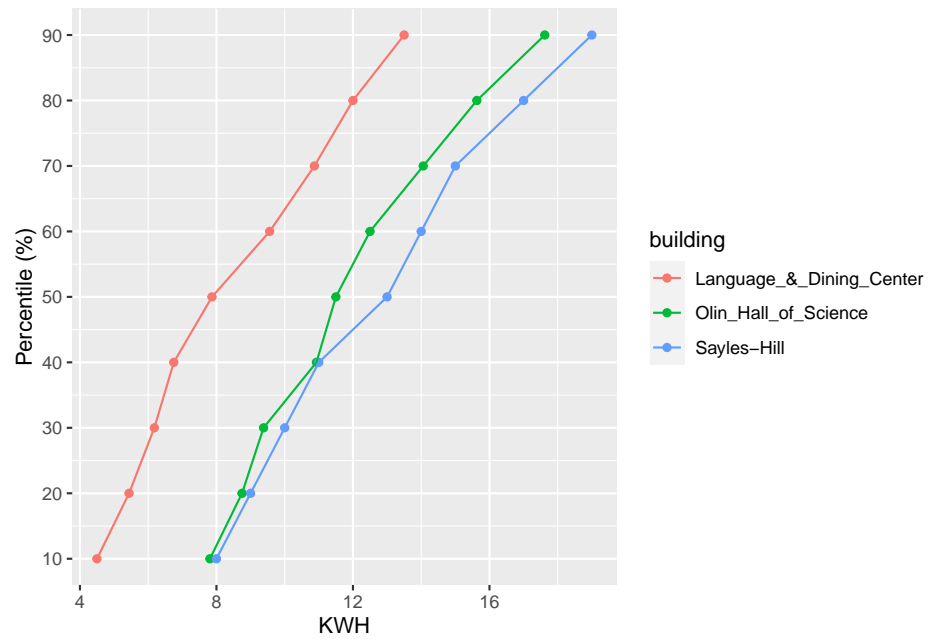
Answer:

```

qdf %>%
  filter(building %in% c("Sayles-Hill", "Language_&_Dining_Center", "Olin_Hall_of_Science")) +
  ggplot(aes(x=quantiles, y=parse_number(stat), color=building)) +
  geom_point() +
  geom_line(aes(group=building)) +

```

```
labs(y="Percentile (%)",x="KWH") +  
scale_y_continuous(breaks=seq(10,90,by=10))
```



Chapter 22

Class Activity 16

```
# load the necessary libraries
library(tidyverse)
library(stringr)
library(polite)
library(rvest)
```

22.1 Group Activity 1

- 22.1.1 a. Scrape the first table in List_of_NASA_missions wiki page. Additionally, use `janitor::clean_names()` to clean the column names and store the resulting table as `NASA_missions.csv` in your working folder.

Click for answer

```
wiki_NASA <- "https://en.wikipedia.org/wiki/List_of_NASA_missions"

# Scrape the data and write the first table to a CSV file
bow(wiki_NASA) %>%
  scrape() %>%
  html_nodes("table") %>%
  .[[1]] %>%
  html_table(fill = TRUE) %>%
  janitor::clean_names() %>%
  write_csv("NASA_missions.csv")
```

- 22.1.2 b. Now, write a code snippet to scrape all the URLs from the anchor tags (`a`) on a given Wikipedia page, convert the relative URLs to absolute URLs, and store the results in a tibble and save it as `NASA_missions_urls.csv` in your working folder.

Click for answer

```
# Scrape the data and write the URLs to a CSV file
bow(wiki_NASA) %>%
  scrape() %>%
  html_nodes("a") %>%
  html_attr("href") %>%
  url_absolute("https://en.wikipedia.org/") %>%
  data.frame(url = .) %>%
  write_csv("NASA_missions_urls.csv")
```

22.2 Group Activity 2

- 22.2.1 a. How do you scrape a table from a web page using `rvest`, clean the column names with `janitor`, and prepare the data for analysis in R?

Click for answer

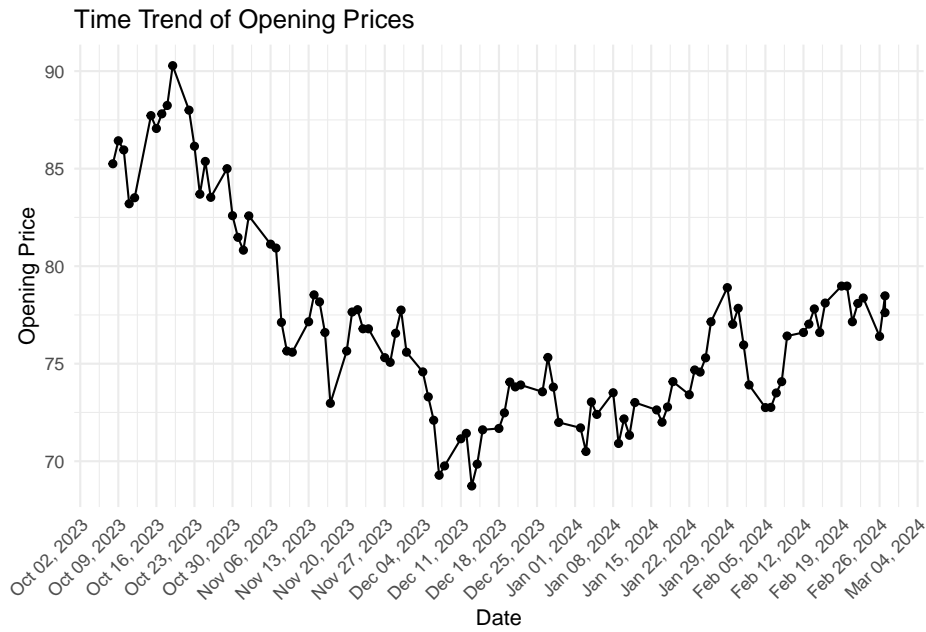
```
yf <- "https://finance.yahoo.com/quote/CL%3DF/history?p=CL%3DF"
bow(yf) %>% scrape() %>%
  html_nodes("table") %>% .[[1]] %>%
  html_table() %>% janitor::clean_names() %>%
  slice(-n()) %>%
  mutate(date = lubridate::mdy(date)) %>%
  mutate_at(vars(open:adj_close), as.numeric) -> ticker
```

- 22.2.2 b. Write the R code to create a time trend plot of opening prices from the scraped data using `ggplot2`.

Click for answer

```
ggplot(ticker, aes(x = date, y = open)) +
  geom_line() + # Plot lines
  geom_point() + # Add points
  scale_x_date(date_labels = "%b %d, %Y", date_breaks = "1 week") +
  labs(title = "Time Trend of Opening Prices", x = "Date", y = "Opening Price") +
```

```
theme_minimal() +
theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



22.2.3 c. How can you transform the data into a long format suitable for plotting multiple price types with ggplot2?

Click for answer

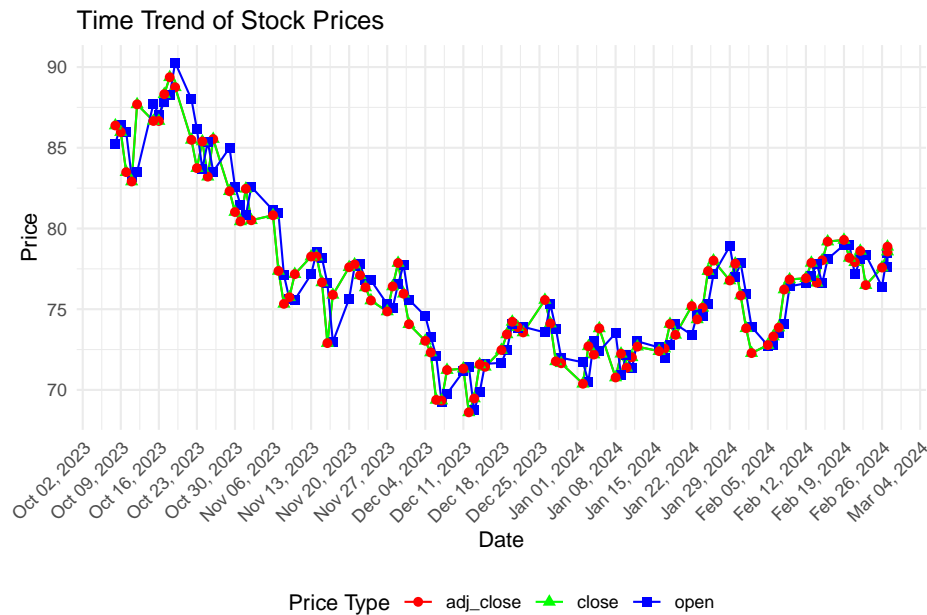
```
ticker_long <- ticker %>%
  pivot_longer(cols = c(open, close, adj_close), names_to = "PriceType", values_to = "Price")
```

22.2.4 d. Show how to create a ggplot2 visualization that includes lines and points, with different colors and shapes for each price type, and make the x-axis dates legible.

Click for answer

```
ggplot(ticker_long, aes(x = date, y = Price, color = PriceType)) +
  geom_line() +
  geom_point(aes(shape = PriceType), size = 2) + # Different shapes for each price type
  scale_color_manual(values = c("open" = "blue", "close" = "green", "adj_close" = "red")) +
  scale_x_date(date_labels = "%b %d, %Y", date_breaks = "1 week") +
  labs(title = "Time Trend of Stock Prices", x = "Date", y = "Price") +
```

```
theme_minimal() +
theme(
  axis.text.x = element_text(angle = 45, hjust = 1),
  legend.position = "bottom"
) +
guides(shape = guide_legend(title = "Price Type"), color = guide_legend(title = "Pri
```



22.3 Group Activity 3

In this activity, you'll scrape web data using `rvest` and tidy up the results into a well-formatted table. Start by extracting job titles from a given URL, then gather the associated company names, and trim any leading or trailing whitespace from the location data. Next, retrieve the posting dates and the URLs for the full job descriptions. Finally, combine all these elements into a single dataframe, ensuring that each piece of information aligns correctly. Your task is to produce a clean and informative table that could be useful for job seekers. To facilitate the selection of the correct CSS selectors, you may find the `SelectorGadget` Chrome extension particularly useful.

```
url <- "https://realpython.github.io/fake-jobs/"
```

Click for answer


```

title <- bow(url) %>% scrape() %>% html_elements(css = ".is-5") %>% html_text() # part 1
company <- bow(url) %>% scrape() %>% html_elements(css = ".company") %>% html_text() # part 2
location <- bow(url) %>% scrape() %>% html_elements(css = ".location") %>% html_text() %>% str_trim()
time <- bow(url) %>% scrape() %>% html_elements(css = "time") %>% html_text() # part 4
html <- bow(url) %>% scrape() %>% html_element(css = ".card-footer-item+ .card-footer-item") %>%
  html_text()

# Create a dataframe
tibble(title = title, company = company, location = location, time = time, html = html) # port 6

# A tibble: 100 x 5
  title                company      location time  html
  <chr>                <chr>      <chr>    <chr> <chr>
1 Senior Python Developer Payne, Ro~ Stewart~ 2021~ http~
2 Energy engineer       Vasquez-D~ Christo~ 2021~ http~
3 Legal executive       Jackson, ~ Port Er~ 2021~ http~
4 Fitness centre manager Savage-Br~ East Se~ 2021~ http~
5 Product manager       Ramirez I~ North J~ 2021~ http~
6 Medical technical officer Rogers-Ya~ Davidvi~ 2021~ http~
7 Physiological scientist Kramer-Kl~ South C~ 2021~ http~
8 Textile designer      Meyers-Jo~ Port Jo~ 2021~ http~
9 Television floor manager Hughes-Wi~ Osborne~ 2021~ http~
10 Waste management officer Jones, Wi~ Scottto~ 2021~ http~
# i 90 more rows

```


Chapter 23

Class Activity 17

```
# load the necessary libraries
library(tidyverse)
library(stringr)
library(purrr)
library(ggthemes)
library(rvest)
library(polite)
```

23.1 Group Activity 1

1. Go to the the numbers webpage and extract the table on the front page.

Click for answer

```
session1 <- read_html("https://www.the-numbers.com/movie/budgets/all") %>%
  html_nodes(css = "table") %>%
  html_table()

table_base <- session1 %>% .[[1]]
```

2. Find out the number of pages that contain the movie table, while looking for the changes in the url in the address bar. How does the url changes when you go to the next page?

Click for answer

Answer: The starting count of the movie gets concatenated to the url in increments of 100.

3. Write a for loop to store all the data in multiple pages to a single data frame. Please do the same using `purrr::map_df()` as well.

Click for answer

```
library(tidyverse)
library(rvest)

new_urls <- "https://www.the-numbers.com/movie/budgets/all/"

# Create an empty data frame
df1 <- list()

# Generate a vector of indices
index <- seq(1, 6301, 100)

# Loop through indices, scrape data, and bind the resulting data frames
start_time <- proc.time() # Capture start time
for (i in 1:length(index)) {
  url <- str_glue("{new_urls}{index[i]}")
  webpage <- read_html(url)
  table_new <- html_table(webpage)[[1]] %>%
    janitor::clean_names() %>%
    mutate(across(everything(), as.character))
  df1[[i]] <- table_new
}
end_time <- proc.time() # Capture end time
end_time - start_time # Calculate duration

      user system elapsed
3.885   0.110  28.335

df1_final <- do.call(rbind, df1)
df1_final1 <- reduce(df1, dplyr::bind_rows)

# alternate using map_df()
start_time <- proc.time() # Capture start time

urls <- map(index, function(i) str_glue({new_urls}, {index[i]}))
urls <- map(index, ~str_glue({new_urls}, {.x}))

library(tidyverse)
library(rvest)
library(glue)
library(janitor)

# Assuming 'urls' is already defined
movies_data <- map_df(urls, ~read_html(.x) %>%
  html_table() %>%
  .[[1]] %>%
```

```

      janitor::clean_names() %>%
      mutate(across(everything(), as.character)))
end_time <- proc.time() # Capture end time
end_time - start_time # Calculate duration

      user  system elapsed
3.871    0.097   28.088

movies_data %>% slice_head(n=6)

# A tibble: 6 x 6
  x      release_date movie production_budget domestic_gross
<chr> <chr>         <chr>         <chr>          <chr>
1 1      Dec 9, 2022 Avata~ $460,000,000    $684,075,767
2 2      Apr 23, 2019 Aveng~ $400,000,000    $858,373,000
3 3      May 20, 2011 Pirat~ $379,000,000    $241,071,802
4 4      Apr 22, 2015 Aveng~ $365,000,000    $459,005,868
5 5      May 17, 2023 Fast X $340,000,000    $146,126,015
6 6      Dec 16, 2015 Star ~ $306,000,000    $936,662,225
# i 1 more variable: worldwide_gross <chr>

```

23.2 Group Activity 2

1. Go to the scrapethissite and extract the table on the front page.

Click for answer

```

session1 <- read_html("https://www.scrapethissite.com/pages/forms/") %>%
  html_nodes(css = "table") %>%
  html_table()

table_base <- session1 %>% .[[1]]

```

2. Find out the number of pages that contain the movie table, while looking for the changes in the url in the address bar. How does the url changes when you go to the next page?

Click for answer

Answer: The url field has ?page_num= added with the number of pages running from 1 to 24.

3. Write a for loop to store all the data in multiple pages to a single data frame. Please do the same using purrr::map_df() as well.

Click for answer

```

library(tidyverse)
library(rvest)

```

```

new_urls <- "http://scrapethissite.com/pages/forms/?page_num="

# Generate a vector of indices
index <- seq(1, 24)

df2 <- list()
start_time <- proc.time() # Capture start time

for (i in index) {
  url <- str_glue("{new_urls}{i}")
  webpage <- read_html(url)
  table_new <- html_table(webpage)[[1]] %>%
    janitor::clean_names() %>%
    #set_names(~ifelse(is.na(.) | . == "", paste("V", seq_along(.), sep=""), .)) %>%
    mutate(across(everything(), as.character))
  df2[[i]] <- table_new
}
end_time <- proc.time() # Capture end time
end_time - start_time # Calculate duration

      user system elapsed
1.612   0.059   9.937

df2_final <- bind_rows(df2)
df2_final

# A tibble: 582 x 9
  team_name      year wins losses ot_losses win_percent
  <chr>         <chr> <chr> <chr>   <chr>         <chr>
1 Boston Bruins  1990  44   24    <NA>         0.55
2 Buffalo Sabres 1990  31   30    <NA>         0.388
3 Calgary Flames 1990  46   26    <NA>         0.575
4 Chicago Blackha~ 1990  49   23    <NA>         0.613
5 Detroit Red Win~ 1990  34   38    <NA>         0.425
6 Edmonton Oilers 1990  37   37    <NA>         0.463
7 Hartford Whalers 1990  31   38    <NA>         0.388
8 Los Angeles Kin~ 1990  46   24    <NA>         0.575
9 Minnesota North~ 1990  27   39    <NA>         0.338
10 Montreal Canadi~ 1990  39   30    <NA>         0.487
# i 572 more rows
# i 3 more variables: goals_for_gf <chr>,
#   goals_against_ga <chr>, x <chr>

# alternate using map
urls <- map(index, function(i) str_glue({new_urls}, {i}))
urls <- map(index, ~str_glue("{new_urls}{.x}"))

```

```

start_time <- proc.time() # Capture start time
sports_data <- map_df(urls, ~read_html(.x) %>%
  html_table() %>%
  .[[1]] %>%
  janitor::clean_names() %>%
  mutate(across(everything(), as.character)))

end_time <- proc.time() # Capture end time
end_time - start_time # Calculate duration

```

```

      user system elapsed
1.527   0.094   8.805

```

```
sports_data %>% slice_head(n=7)
```

```
# A tibble: 7 x 9
```

| | team_name | year | wins | losses | ot_losses | win_percent |
|---|-------------------|-------|-------|--------|-----------|-------------|
| | <chr> | <chr> | <chr> | <chr> | <chr> | <chr> |
| 1 | Boston Bruins | 1990 | 44 | 24 | <NA> | 0.55 |
| 2 | Buffalo Sabres | 1990 | 31 | 30 | <NA> | 0.388 |
| 3 | Calgary Flames | 1990 | 46 | 26 | <NA> | 0.575 |
| 4 | Chicago Blackhaw~ | 1990 | 49 | 23 | <NA> | 0.613 |
| 5 | Detroit Red Wings | 1990 | 34 | 38 | <NA> | 0.425 |
| 6 | Edmonton Oilers | 1990 | 37 | 37 | <NA> | 0.463 |
| 7 | Hartford Whalers | 1990 | 31 | 38 | <NA> | 0.388 |

```

# i 3 more variables: goals_for_gf <chr>,
#   goals_against_ga <chr>, x <chr>

```


Chapter 24

Class Activity 18

```
# load the necessary libraries
library(tidyverse)
library(shiny)
library(readr)
library(janitor)
library(purrr)
library(lubridate)
library(DT)
library(ggthemes)
library(rvest)
library(polite)
```

24.1 Introduction to Shiny Apps

Shiny is a powerful R package that enables the creation of interactive web applications directly from R. A Shiny app consists of two main components:

- *User Interface (UI)*: Defines the layout and appearance of your app.
- *Server Function*: Contains the logic to manage user inputs and generate outputs.

24.1.0.1 Step 1: Setting Up a Basic Shiny App

Start with a simple “Hello World!” app to understand the basic structure.

```
ui <- fluidPage("Hello World!")
server <- function(input, output) {}
shinyApp(ui = ui, server = server, options = list(height = 200))
```

24.1.0.2 Step 2: Enhancing the UI with More Information

Add more elements to the UI to make it informative.

```
ui <- fluidPage(
  titlePanel("Tracking Covid in Minnesota"),
  h1("A Detailed Overview"),
  p("elements1", "elements2"),
  br(),
  p("Detailed Information:", strong("Critical Data Points"))
)

server <- function(input, output) {}
shinyApp(ui = ui, server = server, options = list(height = 200))
```

24.1.0.3 Step 3: Structuring the App with a Sidebar Layout

Organize the app's layout using sidebarLayout to separate inputs from results.

```
ui <- fluidPage(
  titlePanel("Tracking Covid in Minnesota"),
  sidebarLayout(
    sidebarPanel(p("Inputs will be placed here")),
    mainPanel(p("Results will be displayed here"))
  )
)

server <- function(input, output) {}
shinyApp(ui = ui, server = server, options = list(height = 200))
```

24.1.0.4 Step 4: Reading and Preparing Data

Fetch and prepare Covid data for Minnesota for analysis.

```
table_usafacts <- bow(url = "https://usafacts.org/visualizations/coronavirus-covid-19-")
  scrape() %>% html_nodes("a") %>% # find all links
  html_attr("href") %>% # get the url
  str_subset(".csv") # find those that end in csv

covid_data <- read_csv(table_usafacts[2]) %>%
  filter(State == "MN") %>%
  select(-countyFIPS, -StateFIPS, -State) %>%
  pivot_longer(-1, names_to = "Dates", values_to = "Cases") %>%
  clean_names() %>%
  mutate(dates = ymd(dates),
         counties = factor(str_remove(county_name, " County"))) %>%
  select(-county_name)
```

```
head(covid_data, 10) %>% knitr::kable()
```

24.1.0.5 Step 5: Incorporating User Inputs

Enhance the UI with inputs for filtering data based on user selection.

```
ui <- fluidPage(
  titlePanel("Tracking Covid in Minnesota"),
  sidebarLayout(
    sidebarPanel(
      sliderInput("monthInput", "Select Month Range", min = 1, max = 12, value = c(3, 6)),
      selectInput("countyInput", "Select County", choices = unique(covid_data$counties))
    ),
    mainPanel(
      plotOutput("covidPlot"),
      tableOutput("covidTable")
    )
  )
)

server <- function(input, output) {}

shinyApp(ui = ui, server = server)
```

24.1.0.6 Step 6: Developing the Server Logic

Implement server logic to respond to user inputs and generate outputs.

```
server <- function(input, output) {
  # Filter data based on inputs
  filtered_data <- reactive({
    covid_data %>%
      filter(counties == input$countyInput,
             month(dates) >= input$monthInput[1],
             month(dates) <= input$monthInput[2])
  })

  # Generate plot
  output$covidPlot <- renderPlot({
    ggplot(filtered_data(), aes(x = dates, y = cases, color = counties)) +
      geom_line() +
      theme_minimal() +
      labs(title = "Covid Cases Over Time", y = "Number of Cases")
  })
}
```

```

# Generate table
output$covidTable <- renderTable({
  filtered_data()
})
}

# Launch the app
shinyApp(ui = ui, server = server)

```

24.1.1 Implementation 2

```

ui1 <- fluidPage(
  titlePanel("Tracking Covid in Minnesota"),
  sidebarLayout(
    sidebarPanel(
      sliderInput("monthInput", "Month", 0, 12, c(3, 6)),
      radioButtons("yearInput", "Year",
        choices = c("2020", "2021", "2022", "2023"),
        selected = "2020"),
      selectInput(inputId = "dv", label = "County",
        choices = levels(covid_data$counties),
        selected = c("Aitkin"))
    ),
    mainPanel(
      plotOutput(outputId = "plot"), br(),
      DT::dataTableOutput(outputId = "table")
    )
  )
)

server1 <- function(input, output) {
  filtered_data <- reactive({
    subset(covid_data,
      counties %in% input$dv &
      month >= input$monthInput[1] & month <= input$monthInput[2] &
      year == input$yearInput) })

  output$plot <- renderPlot({
    ggplot(filtered_data(), aes(x=dates, y=cases, color=counties)) + theme_economist() +
      geom_point(alpha=0.5, color = "blue") + theme(legend.position = "none") +
      ylab("Number of Cases") + xlab("Date"))

  output$table <- DT::renderDataTable({
    filtered_data()})
}

```

```
}
```

```
app1 <- shinyApp(ui = ui1, server = server1, options = list(height = 1200))  
app1
```


Chapter 25

Class Activity 19

```
# List of required packages
library(tidyverse)
library(shiny)
library(dplyr)
library(babynames)
```

Explore the magic firsthand! Simply paste the code provided in the slides into the field below and run it to witness the output in action.

```
# Try the Shiny Code Yourself! 1
library(shiny)
ui <- fluidPage(
  titlePanel("Updating Plot Based on User Input"),
  sidebarLayout(
    sidebarPanel(
      numericInput("obs", "Number of observations:",
        value = 100, min = 1)
    ),
    mainPanel( plotOutput("distPlot"))
  )
)
server <- function(input, output) {
  output$distPlot <- renderPlot({
    hist(rnorm(input$obs))
  })
}
shinyApp(ui, server)
```

```
# Try the Shiny Code Yourself! 2
ui <- fluidPage(
```

```

    mainPanel(
      actionButton("addButton", "Add 1"),
      textOutput("result")
    )
  )
  server <- function(input, output) {
    sum_so_far <- eventReactive(input$addButton, {
      if (is.null(input$addButton)) {
        0
      } else {
        input$addButton
      }
    })
    output$result <- renderText({
      paste("Sum so far:", sum_so_far())
    })
  }
  shinyApp(ui, server)

```

Try the Shiny Code Yourself! 3

```

ui <- fluidPage(
  mainPanel(
    numericInput("obs", "Number of observations", value = 100, min = 1),
    actionButton("updateButton", "Update Plot"),
    plotOutput("distPlot")
  )
)
server <- function(input, output) {
  data_to_plot <- eventReactive(input$updateButton, {
    rnorm(input$obs)
  })
  output$distPlot <- renderPlot({
    hist(data_to_plot())
  })
}
shinyApp(ui, server)

```

Try the Shiny Code Yourself! 4

```

ui <- fluidPage(
  titlePanel("Simple Counter"),
  numericInput("adjustValue", "Adjust by:", 1),
  actionButton("incrementButton", "Increment"),
  actionButton("decrementButton", "Decrement"),
  br(),
  textOutput("currentValue")
)

```



```

server <- function(input, output) {
  counter <- reactiveValues(value = 0)
  observeEvent(input$incrementButton, {
    counter$value <- counter$value + input$adjustValue
  })
  observeEvent(input$decrementButton, {
    counter$value <- counter$value - input$adjustValue
  })
  output$currentValue <- renderText({
    paste("Current Value:", counter$value)
  })
}
shinyApp(ui, server)

```

```

# Try the Shiny Code Yourself! 5
ui <- fluidPage(
  mainPanel(
    actionButton("showAlert", "Show Alert"),
    textOutput("alertCount")
  )
)
server <- function(input, output) {
  alert_counter <- reactiveValues(count = 0)
  observeEvent(input$showAlert, {
    showModal(modalDialog(
      title = "Alert",
      "This is an alert message!"
    ))
    alert_counter$count <- alert_counter$count + 1
  })
  output$alertCount <- renderText({
    paste("Number of alerts shown:", alert_counter$count)
  })
}
shinyApp(ui, server)

```

```

# Try the Shiny Code Yourself! 6
ui <- fluidPage(
  textInput("textInput", "Enter some text"),
  actionButton("submitButton", "Submit"),
  textOutput("outputText"),
  textOutput("counter")
)
server <- function(input, output) {
  counter <- reactiveValues(clicks = 0)
  observeEvent(input$submitButton, {

```

```

counter$clicks <- counter$clicks + 1
output$outputText <- renderText({
  # Using isolate() to prevent reactivity on textInput changes
  paste("You submitted:", isolate(input$textInput))
})
output$counter <- renderText({
  paste("Submit button clicked", counter$clicks, "times")
})
})
}
shinyApp(ui, server)

```

25.1 Question 1

We are utilizing the `babynames` dataset to explore and visualize how the popularity of specific names has evolved over time in a Shiny app.

25.1.1 (a) Create a user interface (UI) by using `textInput` for the baby name input, `selectInput` for gender selection, and adding an `actionButton` to initiate

```

# fill in the portions marked with # fill #
ui <- fluidPage(
  titlePanel("Baby Names Trend"),
  sidebarLayout(
    sidebarPanel(
      textAreaInput("name", "Enter a Baby Name:", "Emma"),
      selectizeInput("gender", "Select Gender:", choices = c("Male" = "M", "Female" = "F")),
      actionButton("goButton", "Show Trend")
    ),
    mainPanel(plotOutput("nameTrend"))
  )
)

```

- 25.1.2 (b) Filter the `babynames` dataset based on user inputs for name and gender by utilizing `eventReactive` to respond to the action button press, and apply filter within this reactive context, using `isolate` to access the inputs without establishing reactive dependencies.

```
# fill in the portions marked with # fill #

server <- function(input, output) {
  nameData <- eventReactive(input$goButton, {
    req(input$name) # Ensure the name input is not empty
    babynames %>%
      filter(name == isolate(input$name), sex == isolate(input$gender))
  })
}
```

- 25.1.3 (c) Create a line plot to visualize the trend of a selected baby name over the years by using `renderPlot`, showing the number of occurrences (n) of the name across different years (year) in addition to what you did in part b.

```
server <- function(input, output) {
  nameData <- eventReactive(input$goButton, {
    req(input$name) # Ensure the name input is not empty
    babynames %>%
      filter(name == isolate(input$name), sex == isolate(input$gender))
  })

  output$nameTrend <- renderPlot({
    req(nameData())
    ggplot(nameData(), aes(x = year, y = n)) +
      geom_line() +
      labs(title = paste("Trend for name", isolate(input$name)),
           x = "Year", y = "Number of Babies") +
      theme_minimal()
  })
}
```

Remember, without `isolate`, any change to the inputs (`input$name` and `input$gender`) would immediately trigger the reactive expressions (`nameData` and the plotting logic within `renderPlot`) to re-execute. This means the plot

would update every time the user types a letter in the name input field or changes the gender, rather than waiting for the user to press the `goButton`.

25.1.4 (d) Run the Shiny app by combining the UI and server components using `shinyApp`.

```
shinyApp(ui, server, options = list(height = 800))
```

Chapter 26

Class Activity 20

```
# load the necessary libraries
library(tidyverse)
library(shiny)
library(rvest)
library(polite)
library(leaflet)
library(sp)
library(maptools)
library(rgeos)
library(stringr)
library(RColorBrewer)
library(terra)
library(raster)
```

26.1 Group Activity 1

Explore the magic of `leaflet` firsthand! Simply paste the code provided in the slides into the field below and run it to witness the output in action.

```
# paste the code here
```

26.2 Group Activity 2

Explore COVID-19 vaccination rates across the Midwest with this R script, which scrapes data, processes it, and creates an interactive, state-specific leaflet map for clear visualization.

26.2.1 Q1: How do we obtain and prepare the COVID-19 vaccination data?

To analyze COVID-19 vaccination rates, we start by scraping the necessary data from a webpage. This involves selecting the specific HTML elements that contain the data, parsing it into a readable format, and cleaning the column names for easier manipulation. We also filter the data to include only certain states for focused analysis.

```
covid_final <- read_html("https://usafacts.org/visualizations/covid-vaccine-tracker-sta-
  html_elements(css = "table") %>%
  html_table() %>%
  .[[1]] %>%
  janitor::clean_names() %>%
  mutate_at(2:4, parse_number) %>%
  mutate(state = str_to_lower(state)) %>%
  filter(state %in% c("ohio", "indiana", "michigan", "illinois", "missouri", "wisconsin",
    "iowa", "kansas", "nebraska", "south dakota", "north dakota"))
```

26.2.2 Q2: How do we represent geographical data for mapping?

Next, we create a spatial representation of the Midwest states. This involves mapping state boundaries and converting them into a spatial polygons format, which Leaflet can use to plot the data geographically.

```
USA_Midwest <- maps::map("state",
  regions = c("ohio", "indiana", "michigan", "illinois", "missouri",
    "iowa", "kansas", "nebraska", "south dakota", "north dakota"),
  plot = FALSE, fill = TRUE) %>%
  map2SpatialPolygons(IDs = str_remove(.$names, "(?=.)+."))
```

26.2.3 Q3: How do we integrate the vaccination data with the geographical map?

Merging the COVID-19 vaccination data with the geographical map allows us to visualize vaccination rates across different regions. This step combines the spatial polygons with the vaccination data based on state names.

```
map <- SpatialPolygonsDataFrame(USA_Midwest, covid_final, match.ID = FALSE)
```

26.2.4 Q4: How do we categorize vaccination rates for visualization?

To visually differentiate regions based on their vaccination rates, we create bins of vaccination percentages and assign a color palette. This helps in categorizing states into different levels of vaccination coverage.

```
bins <- seq(min(map$percent_fully_vaccinated), max(map$percent_fully_vaccinated), length.out = 6)
pal <- colorBin("viridis", domain = map$percent_fully_vaccinated, bins = bins)
```

26.2.5 Q5: How do we add informative labels to the map?

Labels enhance the map's interactivity by providing detailed information on hover or click. Here, we format the labels to display the state name and its observed vaccination rate.

```
labels <- sprintf("<strong> %s </strong> <br/> Observed: %s", str_to_upper(map$state), map$percent_fully_vaccinated)
lapply(htmltools::HTML)
```

26.2.6 Q6: How do we plot the finalized map with Leaflet?

```
leaflet(map) %>%
  addTiles() %>%
  setView(lng = -93.1616, lat = 44.4583, zoom = 4) %>%
  addPolygons(
    color = "grey",
    weight = 1,
    fillColor = ~pal(percent_fully_vaccinated),
    fillOpacity = 0.7,
    highlightOptions = highlightOptions(weight = 5),
    label = labels
  ) %>%
  addLegend(
    pal = pal,
    values = ~percent_fully_vaccinated,
    opacity = 0.5,
    title = "Percent Vaccinated",
    position = "bottomright"
  )
```

26.2.7 Q7: Think about how would you include your favorite state beside these midwestern states to the leaflet plot. How would you go about including all the states? You may exclude Alaska and Hawaii from the analysis.

```
covid_final <- read_html("https://usafacts.org/visualizations/covid-vaccine-tracker-states/state/percent_fully_vaccinated")
html_elements(css = "table") %>%
  html_table() %>%
  .[[1]] %>%
  janitor::clean_names() %>%
```

```

mutate_at(2:4, parse_number) %>%
mutate(state = str_to_lower(state)) %>%
filter(state %in% setdiff(state, c("hawaii", "alaska")))

all_states <- covid_final$state %>% unique()

USA_all <- maps::map("state" ,
                    plot = FALSE, fill = TRUE) %>%
  map2SpatialPolygons(IDs = str_remove(.$names, "(?:).+"))

map <- SpatialPolygonsDataFrame(USA_all, covid_final, match.ID = FALSE)

bins <- seq(min(map$percent_fully_vaccinated), max(map$percent_fully_vaccinated), length.out = 10)
pal <- colorBin("viridis", domain = map$percent_fully_vaccinated, bins = bins)

labels <- sprintf("<strong> %s </strong> <br/> Observed: %s", str_to_upper(map$state),
  lapply(htmltools::HTML))

leaflet(map) %>%
  addTiles() %>%
  setView(lng = -93.1616, lat = 44.4583, zoom = 4) %>%
  addPolygons(
    color = "grey",
    weight = 1,
    fillColor = ~pal(percent_fully_vaccinated),
    fillOpacity = 0.7,
    highlightOptions = highlightOptions(weight = 5),
    label = labels
  ) %>%
  addLegend(
    pal = pal,
    values = ~percent_fully_vaccinated,
    opacity = 0.5,
    title = "Percent Vaccinated",
    position = "bottomright"
  )

```

26.2.8 Q8. Deploying a Shiny App to shinyapps.io

Open an account on shinyapps.io and follow the steps to deploy one of the apps to shinyapps.io.

26.2.8.1 Step 1: Create Your Shiny App

Create a file named `app.R` in a new folder. This folder will contain all the necessary files for your app.

26.2.8.2 Step 2: Organize Your App Folder

Ensure your `app.R` file is saved in a dedicated folder for your app. This folder can also include any additional data files, data folder for data, or other resources your app needs.

26.2.8.3 Step 3: Deploy Your App to shinyapps.io

Before deploying, you'll need to set up an account on `shinyapps.io`. Use the `rsconnect` package to authenticate your R session with your `shinyapps.io` account. You will need the API key from your `shinyapps.io` account dashboard.

```
library(rsconnect)
rsconnect::setAccountInfo(name='<account-name>',
                           token='<api-token>',
                           secret='<api-secret>')
```

Replace `<account-name>`, `<api-token>`, and `<api-secret>` with your actual account name and API key details. Navigate to your app folder in RStudio, or set the working directory to your app folder in R:

```
setwd("/path/to/your/app")
```

Deploy your app using the `deployApp` function:

```
rsconnect::deployApp()
```

Follow the prompts in the R console. Once the deployment process completes, you will receive a URL to access your app hosted on `shinyapps.io`.

Chapter 27

Class Activity 21

```
# load the necessary libraries
library(tidyverse)
library(tidymodels)
library(mlbench)
library(janitor)
library(parsnip)
library(kknn)
library(paletteer)
library(corr)
library(forcats)
library(ggthemes)

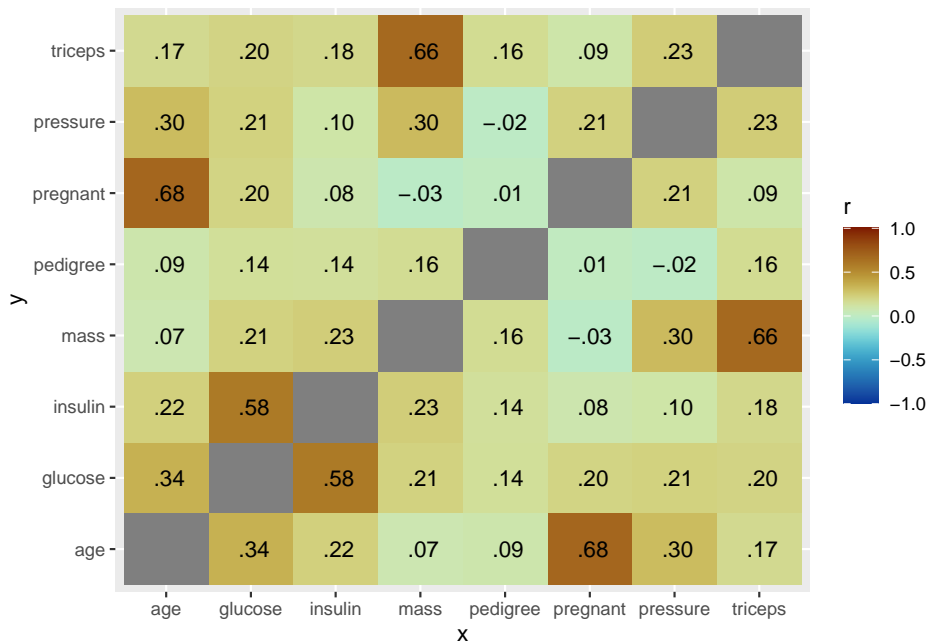
# function that standardizes
standardize <- function(x, na.rm = FALSE) {
  (x - mean(x, na.rm = na.rm)) /
    sd(x, na.rm = na.rm)
}
```

27.1 Group Activity 1

```
# Load the data
data(PimaIndiansDiabetes2)
db <- PimaIndiansDiabetes2 %>% drop_na()

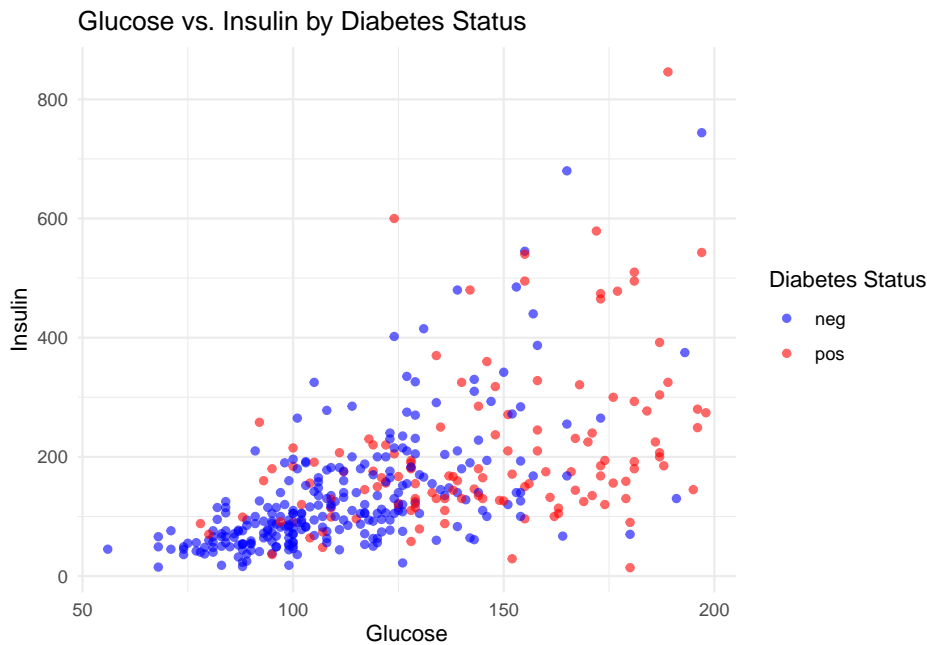
# correlation plot of the variables
db %>%
  select(-diabetes) %>% # only numerical variables
  correlate() %>%
```

```
stretch() %>%
  ggplot(aes(x, y, fill = r)) +
  geom_tile() +
  geom_text(aes(label = as.character(fashion(r)))) +
  scale_fill_paletteer_c("scico::roma", limits = c(-1, 1), direction = -1)
```



- a. Create a scatter plot using ggplot2 to visualize the classification of diabetes status based on glucose and insulin levels, color-coding negative cases in blue and positive cases in red.

```
ggplot(db, aes(x = glucose, y = insulin, color = diabetes)) +
  geom_point(alpha = 0.6) +
  theme_minimal() +
  labs(title = "Glucose vs. Insulin by Diabetes Status",
       x = "Glucose",
       y = "Insulin",
       color = "Diabetes Status") +
  scale_color_manual(values = c("neg" = "blue", "pos" = "red"))
```



- b. Using the provided standardization function, apply it to both the glucose and insulin columns of your dataset to create new standardized columns, then plot these standardized values to analyze diabetes status.

```
db_std <- db %>%
  mutate(glucose_std = standardize(glucose),
         insulin_std = standardize(insulin))
```

- c. Let's perform all the steps involved in classifying whether a patient with certain glucose and insulin would have diabetes or not.

Click for answer

Answer:

```
# 1 Prepare raw data
db_raw <- db %>% select(glucose, insulin, diabetes)

# 2 Create a recipe for data pre-processing
db_recipe <- recipe(diabetes ~ glucose + insulin, data = db_raw) %>%
  step_scale(all_predictors()) %>%
  step_center(all_predictors()) %>%
  prep()

# 3 Apply the recipe to the data set
db_scaled <- bake(db_recipe, db_raw)
```

```
# 4 Create a model specification
knn_spec <- nearest_neighbor(mode = "classification",
                             engine = "kknn",
                             neighbors = 5)

# 5 Fit the model on the pre-processed data
knn_fit <- knn_spec %>%
  fit(diabetes ~ insulin + diabetes, data = db_scaled)

# 6 Classify
# These are standardized value!!
new_observations <- tibble(glucose = c(1, 2), insulin = c(-1, 1))
predict(knn_fit, new_data = new_observations)

# A tibble: 2 x 1
#   .pred_class
#   <fct>
# 1 neg
# 2 neg
```

- d. We already know the labels of some of the patients in the dataset. How well does the model predict their diabetes status? We will see more of this in the coming lectures, but for now try to compare the results for the first 10 cases in the dataset.

Click for answer

Answer:

```
scaled_observations <- db_scaled[1:50,]
predictions <- predict(knn_fit, new_data = scaled_observations)
bind_cols(scaled_observations, predictions) -> predict_data
```

What is the accuracy percentage?

Answer:

```
sum(predictions == db_raw %>% select(diabetes) %>% slice(1:50))/50
```

```
[1] 0.9
```

```
# alternate
accuracy_percentage <- predict_data %>%
  mutate(correct_prediction = diabetes == .pred_class) %>%
  summarize(accuracy = mean(correct_prediction, na.rm = TRUE)) %>%
  pull(accuracy) * 100

accuracy_percentage
```

```
[1] 90
```

- e. Repeat part d. with a different model fitted with different number of neighbors. See if the accuracy percentage change in this new setting.

Chapter 28

Class Activity 22

```
# load the necessary libraries
library(tidyverse)
library(tidymodels)
library(mlbench)      # for PimaIndiansDiabetes2 dataset
library(janitor)
library(parsnip)
library(kknn)
library(ggthemes)
library(purrr)
library(forcats)
```

28.1 Group Activity 1

Load the mlbench package to get PimaIndiansDiabetes2 dataset.

```
# Load the data - diabetes
data(PimaIndiansDiabetes2)
db <- PimaIndiansDiabetes2
db <- db %>% drop_na() %>% mutate(diabetes = fct_rev(factor(diabetes)))
db_raw <- db %>% select(glucose, insulin, diabetes)
```

- Split the data 75-25 into training and test set using the following code.

Click for answer

Answer:

```
set.seed(123)

db_split <- initial_split(db, prop = 0.75)
```

```
# Create training data
db_train <- db_split %>% training()

# Create testing data
db_test <- db_split %>% testing()
```

b. Follow the steps to train a 7-NN classifier using the `tidymodels` toolkit

Click for answer

Answer:

```
# define recipe and preprocess the data
db_recipe <- recipe(diabetes ~ ., data = db_raw) %>%
  step_scale(all_predictors()) %>%
  step_center(all_predictors()) %>%
  prep()

# specify the model
db_knn_spec7 <- nearest_neighbor(mode = "classification",
                                engine = "knn",
                                weight_func = "rectangular",
                                neighbors = 7)

# define the workflow
db_workflow <- workflow() %>%
  add_recipe(db_recipe) %>%
  add_model(db_knn_spec7)

# fit the model
db_fit <- fit(db_workflow, data = db_train)
```

c. Classify the penguins in the `test` data frame.

Click for answer

Answer:

```
test_features <- db_test %>% select(glucose, insulin)
db_pred <- predict(db_fit, test_features, type = "raw")

db_results <- db_test %>%
  select(glucose, insulin, diabetes) %>%
  bind_cols(predicted = db_pred)

head(db_results, 6)
```

| | glucose | insulin | diabetes | predicted |
|---|---------|---------|----------|-----------|
| 4 | 89 | 94 | neg | neg |

| | | | | |
|----|-----|-----|-----|-----|
| 7 | 78 | 88 | pos | neg |
| 15 | 166 | 175 | pos | pos |
| 19 | 103 | 83 | neg | neg |
| 32 | 158 | 245 | pos | pos |
| 36 | 103 | 192 | neg | neg |

28.2 Group Activity 2

Calculate the accuracy, sensitivity, specificity, and positive predictive value by hand using the following confusion matrix.

```
conf_mat(db_results, truth = diabetes, estimate = predicted)
```

```

      Truth
Prediction pos neg
      pos  17   8
      neg  12  61

```

Click for answer

Answer:

```
accuracy(db_results, truth = diabetes,
         estimate = predicted)
```

```

# A tibble: 1 x 3
  .metric .estimator .estimate
  <chr>   <chr>       <dbl>
1 accuracy binary      0.796

```

```
sens(db_results, truth = diabetes,
     estimate = predicted)
```

```

# A tibble: 1 x 3
  .metric .estimator .estimate
  <chr>   <chr>       <dbl>
1 sens    binary      0.586

```

```
spec(db_results, truth = diabetes,
     estimate = predicted)
```

```

# A tibble: 1 x 3
  .metric .estimator .estimate
  <chr>   <chr>       <dbl>
1 spec    binary      0.884

```

```
ppv(db_results, truth = diabetes,
    estimate = predicted)
```

```
# A tibble: 1 x 3
  .metric .estimator .estimate
  <chr>   <chr>       <dbl>
1 ppv    binary         0.68
```

28.3 Extra: Code to recreate the plot in the slides for the diabetes dataset.

Click for answer

Answer:

```
metrics_for_k <- function(k, db_train, db_test){
  db_knn_spec <- nearest_neighbor(mode = "classification",
                                engine = "knn",
                                weight_func = "rectangular",
                                neighbors = k)

  db_knn_wkflow <- workflow() %>%
    add_recipe(db_recipe) %>%
    add_model(db_knn_spec)

  db_knn_fit <- fit(db_knn_wkflow, data = db_train)
  test_features <- db_test %>% select(glucose, insulin)
  nn1_pred <- predict(db_knn_fit, test_features, type = "raw")

  db_results <- db_test %>%
    select(diabetes) %>%
    bind_cols(predicted = nn1_pred)
  custom_metrics <- metric_set(accuracy, sens, spec, ppv)

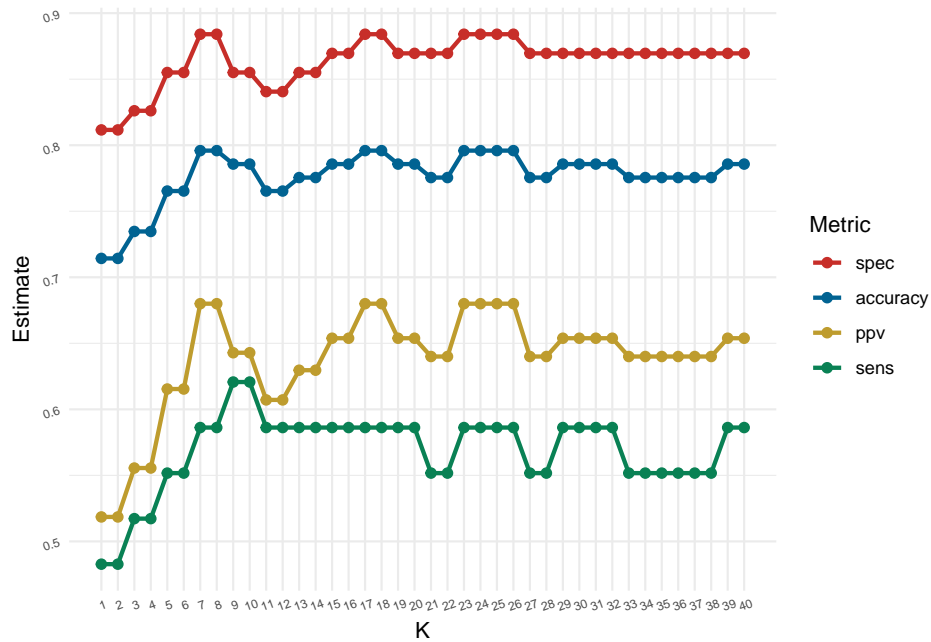
  metrics <- custom_metrics(db_results,
                           truth = diabetes,
                           estimate = predicted)
  metrics <- metrics %>% select(-.estimator) %>% mutate(k = rep(k,4))

  return(list = metrics)
}

optim.results %>%
  ggplot(aes(x = k, y = .estimate, color = forcats::fct_reorder2(.metric, k, .estimate)
  geom_line(size = 1) +
  geom_point(size = 2) +
  theme_minimal() +
  ggthemes::scale_color_wsj() +
  scale_x_continuous(breaks = k) +
```

28.3. EXTRA: CODE TO RECREATE THE PLOT IN THE SLIDES FOR THE DIABETES DATASET.197

```
theme(panel.grid.minor.x = element_blank(),
      axis.text=element_text(size=6, angle = 20))+
labs(color='Metric', y = "Estimate", x = "K")
```



Chapter 29

Class Activity 23

```
# load the necessary libraries
library(tidyverse)
library(tidymodels)
library(mlbench)      # for PimaIndiansDiabetes2 dataset
library(janitor)
library(yardstick) # extra package for getting metrics
library(parsnip) # tidy interface to models
library(ggthemes)
library(forcats)
library(probably)
library(yardstick)
```

29.1 Group Activity 1

Load the mlbench package to get PimaIndiansDiabetes2 dataset.

```
# Load the data - diabetes
data(PimaIndiansDiabetes2)
db <- PimaIndiansDiabetes2
db <- db %>% drop_na()
db_raw <- db %>% select(glucose, insulin, diabetes)

db_split <- initial_split(db_raw, prop = 0.80)
# Create training data
db_train <- db_split %>% training()
# Create testing data
db_test <- db_split %>% testing()
```

- 29.1.1 a. *Creating the Recipe:*** Construct a recipe for the model by normalizing glucose and insulin predictors to predict diabetes status on the training set, ensuring data scales are comparable.

Click for answer

Answer:

```
db_recipe <- recipe(diabetes ~ glucose + insulin, data = db_train) %>%  
  step_scale(all_predictors()) %>%  
  step_center(all_predictors()) %>%  
  prep()
```

- 29.1.2 b. *Model Specification:*** Define the KNN model using a flexible `tune()` placeholder for the number of neighbors, specifying a classification task.

Click for answer

Answer:

```
knn_spec <- nearest_neighbor(weight_func = "rectangular",  
                             engine = "knn",  
                             mode = "classification",  
                             neighbors = tune())
```

- 29.1.3 c. *Creating Folds:*** Divide the training data into 10 stratified folds based on the diabetes outcome to prepare for cross-validation, ensuring representation.

Click for answer

Answer:

```
db_vfold <- vfold_cv(db_train, v = 10, strata = diabetes)
```

- 29.1.4 d. *Cross-Validation Grid:*** Generate a sequence of K values to test with 10-fold cross-validation, evaluating model performance across a range of neighbors.

Click for answer

Answer:


```

k_vals <- tibble(neighbors = seq(from = 1, to = 40, by = 1))

knn_fit <- workflow() %>%
  add_recipe(db_recipe) %>%
  add_model(knn_spec) %>%
  tune_grid(
    resamples = db_vfold,
    grid = k_vals,
    metrics = metric_set(yardstick::ppv, yardstick::accuracy, sens, spec),
    control = control_resamples(save_pred = TRUE))

cv_metrics <- collect_metrics(knn_fit)
cv_metrics %>% group_by(.metric) %>% slice_max(mean)

# A tibble: 12 x 7
# Groups:   .metric [4]
  neighbors .metric .estimator mean      n std_err .config
    <dbl> <chr>    <chr>    <dbl> <int>   <dbl> <chr>
1      39 accuracy binary    0.738    10  0.0273 Prepro~
2      40 accuracy binary    0.738    10  0.0273 Prepro~
3      25 ppv     binary    0.778    10  0.0225 Prepro~
4      26 ppv     binary    0.778    10  0.0225 Prepro~
5      31 sens    binary    0.834    10  0.0261 Prepro~
6      32 sens    binary    0.834    10  0.0261 Prepro~
7      33 sens    binary    0.834    10  0.0226 Prepro~
8      34 sens    binary    0.834    10  0.0226 Prepro~
9      35 sens    binary    0.834    10  0.0226 Prepro~
10     36 sens    binary    0.834    10  0.0226 Prepro~
11     25 spec    binary    0.582    10  0.0487 Prepro~
12     26 spec    binary    0.582    10  0.0487 Prepro~

```

29.1.5 e. Visualization: Plot the cross-validation results to determine the optimal K value, comparing different performance metrics visually.

Click for answer

Answer:

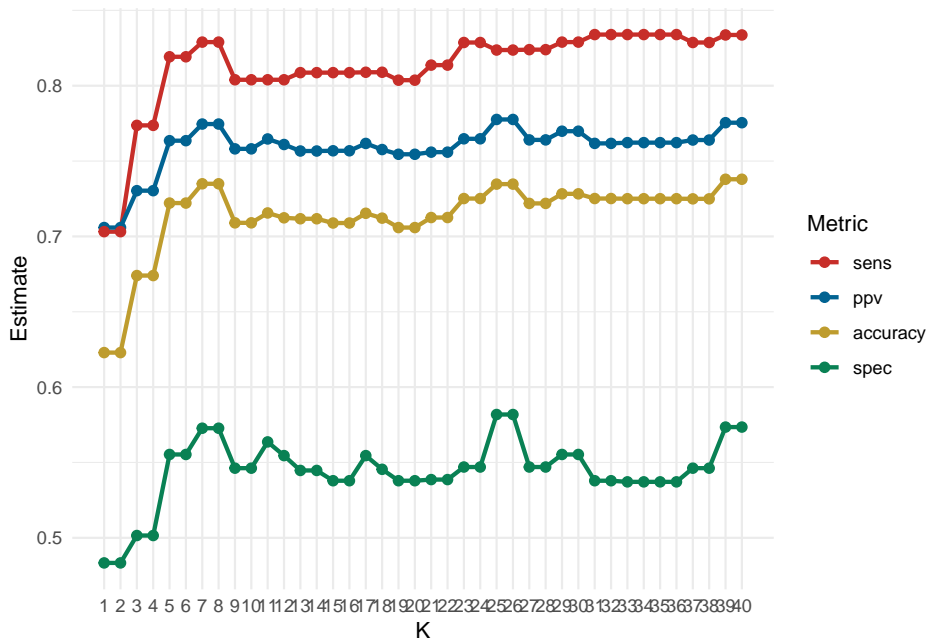
```

final.results <- cv_metrics %>% mutate(.metric = as.factor(.metric)) %>%
  select(neighbors, .metric, mean)

final.results %>%
  ggplot(aes(x = neighbors, y = mean, color = forcats::fct_reorder2(.metric, neighbors, mean))) +
  geom_line(size = 1) +
  geom_point(size = 2) +

```

```
theme_minimal() +
scale_color_wsj() +
scale_x_continuous(breaks = k_vals[[1]]) +
theme(panel.grid.minor.x = element_blank()) +
labs(color='Metric', y = "Estimate", x = "K")
```



29.2 Group Activity 2

29.2.1 a. Data Preparation and Train-Test Split

Load the `mlbench` package and `tidymodels` framework, select relevant features for predicting `glucose`, and split the data into training and test sets. For this activity, use `mass` and `insulin` as your features.

Click for answer

Answer:

```
library(mlbench)
library(tidymodels)
library(dplyr)

data(PimaIndiansDiabetes2)
db <- PimaIndiansDiabetes2 %>%
```

```
drop_na() %>%
select(glucose, mass, insulin)

# Splitting the data
set.seed(2056)
db_split <- initial_split(db, prop = 0.75, strata = glucose)
db_train <- training(db_split)
db_test <- testing(db_split)
```

29.2.2 b. Model Specification

Define a linear regression model for predicting `glucose` as a function of `mass` and `insulin`.

29.2.3 c. Fit the Model

Fit the linear model to the training data, predicting `glucose` based on `mass` and `insulin`.

29.2.4 d. Predict on Test Data and Evaluate the Model

Use the fitted model to predict `glucose` levels on the test set and evaluate the model's accuracy with RMSE and R-squared metrics.

29.2.5 (Bonus) Create a scatter plot to visualize the actual vs. predicted glucose levels, including a regression line for reference.

Chapter 30

Class Activity 24

```
# load the necessary libraries
library(tidyverse)
library(ggthemes)
library(janitor)
library(broom)
library(mlbench)
library(tidymodels)
library(probably)

select <- dplyr::select
theme_set(theme_stata(base_size = 10))

data(PimaIndiansDiabetes2)
db <- PimaIndiansDiabetes2
db <- db %>% drop_na() %>%
  mutate(diabetes = fct_relevel(diabetes, c("neg", "pos"))) # Relevels 'diabetes' factor to ensure
```

30.1 Group Activity 1

In this activity, we will calculate the probability of diabetes for a glucose level of 150 mg/dL using the logistic regression coefficients $\beta_0 = -5.61$ and $\beta_1 = 0.0392$.

30.1.1 a. Calculate Log Odds

First, calculate the log odds for a glucose level of 150 mg/dL.

Click for answer

```
log_odds <- -5.61 + (0.0392 * 150)
log_odds
```

```
[1] 0.27
```

30.1.2 b. Convert Log Odds to Odds

Click for answer

```
odds <- exp(log_odds)
odds
```

```
[1] 1.309964
```

30.1.3 c. Convert Odds to Probability

Click for answer

Finally, convert the odds to probability.

```
probability <- odds / (1 + odds)
probability
```

```
[1] 0.5670929
```

The probability of having diabetes at a glucose level of 150 mg/dL is calculated to be 0.5670929.

30.2 Group Activity 2

a. Let's fit the logistic regression model.

```
set.seed(12345)
db_single <- db %>% select(diabetes, glucose)
db_split <- initial_split(db_single, prop = 0.80)

# Create training data
db_train <- db_split %>% training()

# Create testing data
db_test <- db_split %>% testing()

fitted_logistic_model <- logistic_reg() %>% # Call the model function
  # Set the engine/family of the model
  set_engine("glm") %>%
  # Set the mode
  set_mode("classification") %>%
```

```
# Fit the model
fit(diabetes~., data = db_train)

tidy(fitted_logistic_model)
```

A tibble: 2 x 5

| | term | estimate | std.error | statistic | p.value |
|---|-------------|----------|-----------|-----------|----------|
| | <chr> | <dbl> | <dbl> | <dbl> | <dbl> |
| 1 | (Intercept) | -5.61 | 0.678 | -8.28 | 1.20e-16 |
| 2 | glucose | 0.0392 | 0.00514 | 7.62 | 2.55e-14 |

- b. We are interested in predicting the diabetes status of patients depending on the amount of glucose. Verify that the glucose value of 143.11 gives the probability of having diabetes as $1/2$.
- c. What value of glucose is needed to have a probability of diabetes of 0.5?
- d. Make a classifier that classifies the diabetes status of new patients with a threshold of 0.5, i.e, a new patient is classified as negative if the estimated class probability is less than 0.75. Also, create a confusion matrix of the resulting predictions. Evaluate the model based on accuracy, sensitivity, specificity, and ppv.
- e. Evaluate the performance of a diabetes prediction model at different classification thresholds and visualize how various metrics such as accuracy, sensitivity, and PPV change across these thresholds. Use a sequence of threshold values, apply each one to classify test data, calculate the performance metrics for each classification, and then create a line plot to illustrate the results.