

# STAT 231: Problem Set 4A

Jett Knight

due by 10 PM on Monday, March 15

In order to most effectively digest the textbook chapter readings – and the new R commands each presents – series A homework assignments are designed to encourage you to read the textbook chapters actively and in line with the textbook’s Prop Tip of page 33:

**“Pro Tip:** If you want to learn how to use a particular command, we highly recommend running the example code on your own”

A more thorough reading and light practice of the textbook chapter prior to class allows us to dive quicker and deeper into the topics and commands during class. Furthermore, learning a programming language is like learning any other language – practice, practice, practice is the key to fluency. By having two assignments each week, I hope to encourage practice throughout the week. A little coding each day will take you a long way!

*Series A assignments are intended to be completed individually.* While most of our work in this class will be collaborative, it is important each individual completes the active readings. The problems should be straightforward based on the textbook readings, but if you have any questions, feel free to ask me!

Steps to proceed:

1. In RStudio, go to File > Open Project, navigate to the folder with the course-content repo, select the course-content project (course-content.Rproj), and click "Open"
2. Pull the course-content repo (e.g. using the blue-ish down arrow in the Git tab in upper right window)
3. Copy ps4A.Rmd from the course repo to your repo (see page 6 of the GitHub Classroom Guide for Stat231 if needed)
4. Close the course-content repo project in RStudio
5. Open YOUR repo project in RStudio
6. In the ps4A.Rmd file in YOUR repo, replace "YOUR NAME HERE" with your name
7. Add in your responses, committing and pushing to YOUR repo in appropriate places along the way
8. Run "Knit PDF"
9. Upload the pdf to Gradescope. Don’t forget to select which of your pages are associated with each problem. *You will not get credit for work on unassigned pages (e.g., if you only selected the first page but your solution spans two pages, you would lose points for any part on the second page that the grader can’t see).*

# 1. Web scraping

a.

In Section 5.5.1, the `rvest` package is used to scrape a Wikipedia page. BUT **WAIT!** While we technically might be able to scrape a webpage, that doesn't necessarily mean we are allowed to. **ETHICS ALERT!** Before scraping a web page, you should always check whether doing so is allowed. If you're unsure of the permissions for a particular domain, you can use the handy `paths_allowed` function within the `robotstxt` package.

Check the permissions for the Wikipedia page using the code below. Uncomment the code. You should get a response "TRUE", indicating that a bot has permissions to access this page.

```
# you may need to install the robotstxt package if you're on your machine and haven't used it before
library(robotstxt)
```

```
## Warning: package 'robotstxt' was built under R version 4.0.4
```

```
paths_allowed("https://en.wikipedia.org/wiki/Mile_run_world_record_progression")
```

```
## en.wikipedia.org
```

```
## [1] TRUE
```

b.

Now, follow the code in the textbook to scrape the tables from the Wikipedia page on "Mile run world record progression" (e.g., pages 118-120). Use `length(tables)` to identify how many tables are in the object you created called `tables`. How many tables are there? (Note: the Wikipedia page has been updated since the first edition of the textbook was printed, so the number should be different than that in the textbook!)

ANSWER: There are 12 tables.

```
library(methods)
url <- "http://en.wikipedia.org/wiki/Mile_run_world_record_progression"
tables <- url %>%
  read_html() %>%
  html_nodes("table")
length(tables)
```

```
## [1] 12
```

c.

Next, look at the Wikipedia page: [https://en.wikipedia.org/wiki/Mile\\_run\\_world\\_record\\_progression](https://en.wikipedia.org/wiki/Mile_run_world_record_progression). The table toward the bottom titled "Women Indoor IAAF era" shows four records – one for Mary Decker, two for Doina Melinte, and one for Genzebe Dibaba.

Suppose we want to work in R with this "Women Indoor IAAF era" table. From your `tables` object created in part b, create a dataframe called `women_indoor` that includes this "Women Indoor IAAF era" table. You can use the same code as used in the textbook to create `Table3` and `Table4`, except you'll need to update

the number within the double brackets to correspond to the correct table. You'll likely need to look at a number of different tables in the `tables` object before finding which one corresponds to the "Women Indoor IAAF" table. Print the table. Who holds the indoor one-mile world record for IAAF women, and what was her time?

ANSWER: Genzebe Dibaba from Ethiopia holds the indoor one-mile world record for IAAF women. Her time was 4:13.31.

```
women_indoor <- html_table(tables[[10]])
women_indoor
```

```
##      Time      Athlete  Nationality      Date
## 1  4:20.5    Mary Decker United States February 19, 1982
## 2  4:18.86  Doina Melinte      Romania February 13, 1988
## 3  4:17.14  Doina Melinte      Romania February 9, 1990
## 4  4:13.31  Genzebe Dibaba     Ethiopia February 17, 2016
##
##              Venue
## 1      San Diego United States
## 2 East Rutherford United States
## 3 East Rutherford United States
## 4              Stockholm Sweden
```

d. Lastly:

- create a dataframe called `women_outdoor` that contains the table for "Women's IAAF era" (starting with Anne Smith's record and ending with Sifan Hassan's record)
- combine `women_indoor` and `women_outdoor` into one dataframe called `women` using the `bind_rows()` function. Include a variable called `Type` in this new dataframe to indicate whether a particular observation corresponds to an indoor record or an outdoor record. (Hint: create `Type` separately in each dataframe before combining)
- arrange `women` by ascending time, and identify the fastest world record

Is the fastest record from an indoor or outdoor event?

ANSWER: The fastest record is from an outdoor event.

```
women_outdoor <- html_table(tables[[8]])
women_indoor$Type <- "Indoor"
women_outdoor$Type <- "Outdoor"
women <- bind_rows(women_indoor, women_outdoor)
women %>%
  arrange(Time)
```

```
##      Time      Athlete  Nationality      Date
## 1  4:12.33    Sifan Hassan Netherlands 12 July 2019
## 2  4:12.56  Svetlana Masterkova      Russia 14 August 1996[9]
## 3  4:13.31    Genzebe Dibaba     Ethiopia February 17, 2016
## 4  4:15.61    Paula Ivan      Romania 10 July 1989[9]
## 5  4:16.71  Mary Decker-Slaney United States 21 August 1985[9]
## 6  4:17.14    Doina Melinte      Romania February 9, 1990
## 7  4:17.44    Maricica Puica      Romania 9 September 1982[9]
```

## 8	4:18.08	Mary Decker-Tabb	United States	9 July 1982[9]
## 9	4:18.86	Doina Melinte	Romania	February 13, 1988
## 10	4:20.5	Mary Decker	United States	February 19, 1982
## 11	4:20.89	Lyudmila Veselkova	Soviet Union	12 September 1981[9]
## 12	4:21.7	Mary Decker	United States	26 January 1980[9]
## 13	4:22.1	Natalia Mara<U+0219>escu	Romania	27 January 1979[9]
## 14	4:23.8	Natalia Mara<U+0219>escu	Romania	21 May 1977[9]
## 15	4:29.5	Paola Pigni	Italy	8 August 1973[9]
## 16	4:35.3	Ellen Tittel	West Germany	20 August 1971[9]
## 17	4:36.8	Maria Gommers	Netherlands	14 June 1969[9]
## 18	4:37.0	Anne Smith	United Kingdom	3 June 1967[9]
##		Venue	Type	Auto
## 1		Monaco	Outdoor	
## 2		Zürich	Outdoor	
## 3		Stockholm	Sweden Indoor	<NA>
## 4		Nice	Outdoor	
## 5		Zürich	Outdoor	
## 6	East Rutherford	United States	Indoor	<NA>
## 7		Rieti	Outdoor	
## 8		Paris	Outdoor	
## 9	East Rutherford	United States	Indoor	<NA>
## 10	San Diego	United States	Indoor	<NA>
## 11		Bologna	Outdoor	
## 12		Auckland	Outdoor	4:21.68
## 13		Auckland	Outdoor	4:22.09
## 14		Bucharest	Outdoor	
## 15		Viareggio	Outdoor	
## 16		Sittard	Outdoor	
## 17		Leicester	Outdoor	
## 18		London	Outdoor	