Assignment 10: Data Scraping

Aidan Power

OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics on data scraping.

Directions

- 1. Rename this file <FirstLast>_A10_DataScraping.Rmd (replacing <FirstLast> with your first and last name).
- 2. Change "Student Name" on line 3 (above) with your name.
- 3. Work through the steps, **creating code and output** that fulfill each instruction.
- 4. Be sure your code is tidy; use line breaks to ensure your code fits in the knitted output.
- 5. Be sure to **answer the questions** in this assignment document.
- 6. When you have completed the assignment, **Knit** the text and code into a single PDF file.

Set up

- 1. Set up your session:
- Load the packages tidyverse, rvest, and any others you end up using.
- Check your working directory

```
#1
library(tidyverse)
library(rvest)
library(here)
library(dplyr)
library(ggplot2)
```

[1] "/home/guest/EDE_Fall2024"

- 2. We will be scraping data from the NC DEQs Local Water Supply Planning website, specifically the Durham's 2023 Municipal Local Water Supply Plan (LWSP):
- Navigate to https://www.ncwater.org/WUDC/app/LWSP/search.php
- Scroll down and select the LWSP link next to Durham Municipality.

Indicate this website as the as the URL to be scraped. (In other words, read the contents into an rvest webpage object.)

```
#2
webpage <- read_html('https://www.ncwater.org/WUDC/app/LWSP/report.php?pwsid=03-32-010&year=2023')
webpage

## {html_document}
## <html xmlns="http://www.w3.org/1999/xhtml" lang="en" xml:lang="en">
## [1] <head>\n<title>DWR :: Local Water Supply Planning</title>\n<meta http-equ ...
## [2] <body id="plan">\r\n<!--<div id="division-header">\r\n<a name="top" href= ...</pre>
class(webpage)
```

- ## [1] "xml_document" "xml_node"
 - 3. The data we want to collect are listed below:
 - From the "1. System Information" section:
 - Water system name
 - PWSID
 - Ownership
 - From the "3. Water Supply Sources" section:
 - Maximum Day Use (MGD) for each month

In the code chunk below scrape these values, assigning them to four separate variables.

HINT: The first value should be "Durham", the second "03-32-010", the third "Municipality", and the last should be a vector of 12 numeric values (represented as strings)".

```
#3
#Extract Variables
Water_system_name <- webpage %>%
    html_nodes("div+ table tr:nth-child(1) td:nth-child(2)") %>%
    html_text()

PWSID <- webpage %>%
    html_nodes("td tr:nth-child(1) td:nth-child(5)") %>%
    html_text()

Ownership <- webpage %>%
    html_nodes("div+ table tr:nth-child(2) td:nth-child(4)") %>%
    html_text()

Maximum_day_use <- webpage %>%
    html_nodes("th~ td+ td") %>%
    html_text()

Maximum_day_use #Made sure there are 12 values
```

```
## [1] "28.9000" "33.3000" "43.7000" "30.0000" "40.0000" "37.2300" "34.2000"
## [8] "44.9000" "40.3500" "30.9000" "56.7000" "33.3000"
```

4. Convert your scraped data into a dataframe. This dataframe should have a column for each of the 4 variables scraped and a row for the month corresponding to the withdrawal data. Also add a Date column that includes your month and year in data format. (Feel free to add a Year column too, if you wish.)

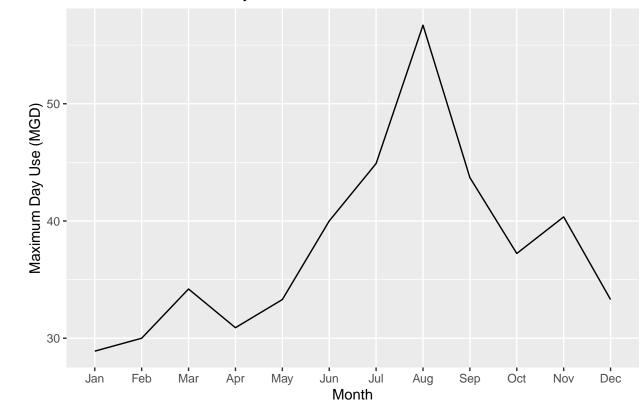
TIP: Use rep() to repeat a value when creating a dataframe.

NOTE: It's likely you won't be able to scrape the monthly widthrawal data in chronological order. You can overcome this by creating a month column manually assigning values in the order the data are scraped: "Jan", "May", "Sept", "Feb", etc... Or, you could scrape month values from the web page...

5. Create a line plot of the maximum daily withdrawals across the months for 2023, making sure, the months are presented in proper sequence.

```
#4
Month <- c("Jan", "May", "Sep", "Feb", "Jun", "Oct",</pre>
           "Mar", "Jul", "Nov", "Apr", "Aug", "Dec")
m.num <- c(01, 05, 09, 02, 06, 10, 03, 07, 11, 04, 08, 12)
Year \leftarrow rep(2023, 12)
df_watersupply <- data.frame("Month" = Month,</pre>
                               "Year" = Year,
                               "Date" = make_date(Year,m.num),
                               "Water System Name" = Water system name,
                               "PWSID" = PWSID,
                               "Ownership" = Ownership,
                               "Maximum Day Use (MGD)" = as.numeric(Maximum_day_use))
#5
df watersupply$Month <- factor(df watersupply$Month, levels=c("Jan","Feb","Mar","Apr","May",
                                                                  "Jun", "Jul", "Aug", "Sep", "Oct",
                                                                  "Nov", "Dec"))
ggplot(df_watersupply, aes(x=Month,y=Maximum.Day.Use..MGD., group=1)) +
  geom_line() +
  labs(title = "Durham Maximum Daily Withdrawals in 2023", x="Month",
       y="Maximum Day Use (MGD)")
```

Durham Maximum Daily Withdrawals in 2023

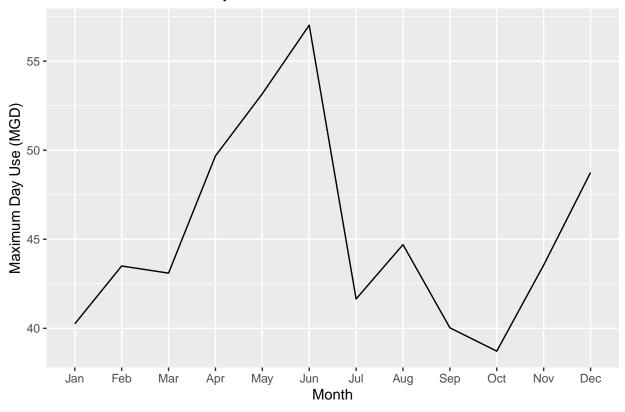


6. Note that the PWSID and the year appear in the web address for the page we scraped. Construct a function using your code above that can scrape data for any PWSID and year for which the NC DEQ has data, returning a dataframe. Be sure to modify the code to reflect the year and site (pwsid) scraped.

```
html_nodes("th~ td+ td") %>%
  html_text()
Month <- c("Jan", "May", "Sep", "Feb", "Jun", "Oct",</pre>
            "Mar", "Jul", "Nov", "Apr", "Aug", "Dec")
m.num <- c(01, 05, 09, 02, 06, 10, 03, 07, 11, 04, 08, 12)
Year.1 <- rep(Year,12)</pre>
df_watersupply2 <- data.frame("Month" = Month,</pre>
"Year" = rep(Year, 12),
"Date" = make_date(Year.1,m.num),
"Water System Name" = Water_system_name,
"PWSID" = PWSID scraped,
"Ownership" = Ownership,
"Maximum Day Use (MGD)" = as.numeric(Maximum_day_use))
return(df_watersupply2)
the_df <- scrape.it('03-32-010','2015')
view(the_df)
```

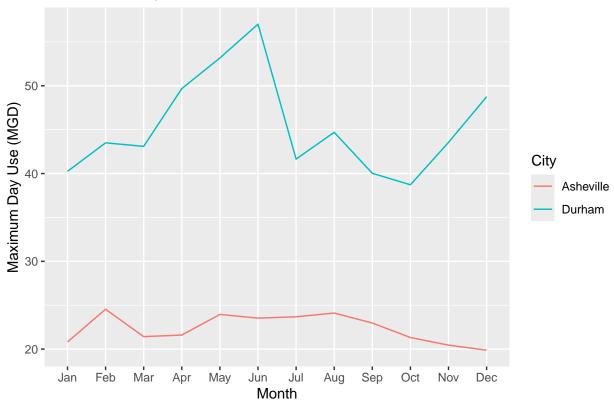
7. Use the function above to extract and plot max daily withdrawals for Durham (PWSID='03-32-010') for each month in 2015

Durham Maximum Daily Withdrawals in 2015



8. Use the function above to extract data for Asheville (PWSID = 01-11-010) in 2015. Combine this data with the Durham data collected above and create a plot that compares Asheville's to Durham's water withdrawals.

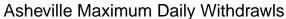
Maximum Daily Withdrawals in 2015

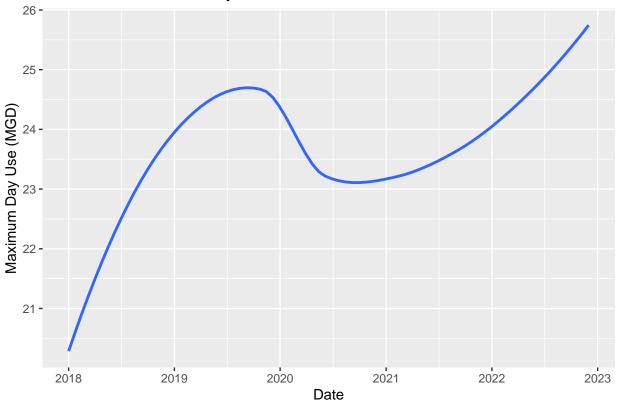


9. Use the code & function you created above to plot Asheville's max daily withdrawal by months for the years 2018 thru 2022.Add a smoothed line to the plot (method = 'loess').

TIP: See Section 3.2 in the "10_Data_Scraping.Rmd" where we apply "map2()" to iteratively run a function over two inputs. Pipe the output of the map2() function to bindrows() to combine the dataframes into a single one.

'geom_smooth()' using formula = 'y ~ x'





Question: Just by looking at the plot (i.e. not running statistics), does Asheville have a trend in water usage over time? > Answer: Water usage appeared to increase until about 2020 where it dropped and then did not increase until the beginning of 2021. With the exception of this decrease, water usage in Asheville has been consistently increasing over time. >