

Assignment 10: Data Scraping

Yilin Zhong

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(lubridate)

getwd()
```

```
## [1] "C:/Users/victo/Desktop/Spring 2023/EDA/EDA-Spring2023"
```

2. We will be scraping data from the NC DEQs Local Water Supply Planning website, specifically the Durham’s 2022 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.
 - Note the web address: <https://www.ncwater.org/WUDC/app/LWSP/report.php?psid=03-32-010&year=2022>

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

```
#2
the_website<- read_html("https://www.ncwater.org/WUDC/app/LWSP/report.php?pwsid=03-32-010&year=2022")
```

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), with the first value being “27.6400”.

```
#3
water.system.name <- html_nodes(the_website,
                                "div+ table tr:nth-child(1) td:nth-child(2)") %>% html_text()

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

ownership <- html_nodes(the_website,
                        "div+ table tr:nth-child(2) td:nth-child(4)") %>% html_text()

max.withdrawals.mgd <-html_nodes(the_website,
                                "th~ td+ td") %>% html_text()
```

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 withdrawal 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...

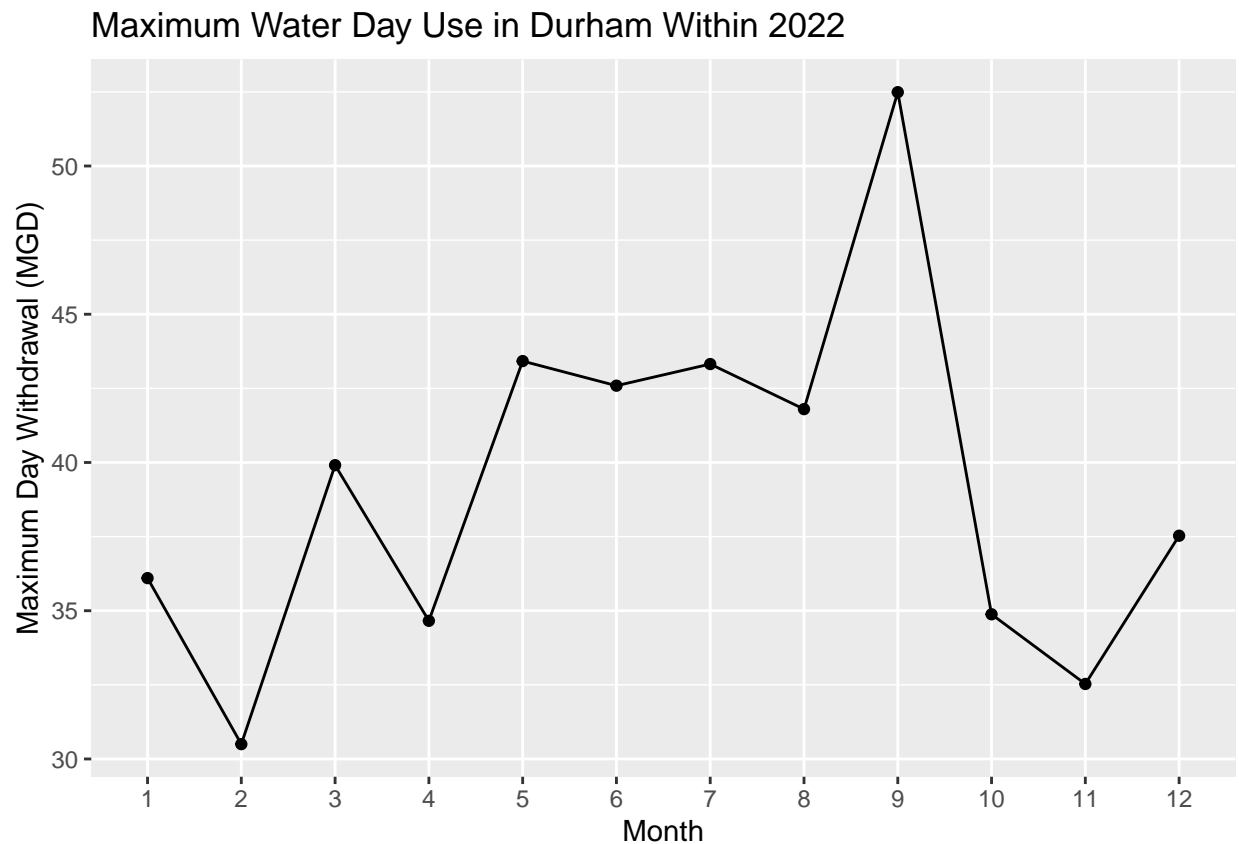
5. Create a line plot of the average daily withdrawals across the months for 2022

```

#4
Month<-c("1", "5", "9", "2", "6", "10", "3", "7","11", "4", "8", "12")
durham.df<-data.frame(
  "Month" = as.numeric(Month),
  "Year" = rep(2022),
  "max.withdrawals.mgd" = as.numeric(max.withdrawals.mgd)
) %>%
  mutate(
    Water.System.Name = !!water.system.name,
    PWSID = !!PWSID,
    Ownership = !!ownership,
    Date = my(paste(Month,"-",Year))
  )

#5
ggplot(durham.df, aes(x = as.factor(Month), y =max.withdrawals.mgd, group=1 )) +
  geom_line() +
  geom_point()+
  labs(x = "Month",
       y = "Maximum Day Withdrawal (MGD)",
       title = "Maximum Water Day Use in Durham Within 2022")

```



- 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. Be sure to modify the code to reflect the year and site (pwsid) scraped.

#6.

```
scrape.it <- function(the_year, the_pwsid){
  function_website <- read_html(paste0('https://www.ncwater.org/WUDC/app/LWSP/report.php?',
    'pwsid=', the_pwsid, '&', 'year=', the_year))

  water.system.name_tag <- 'div+ table tr:nth-child(1) td:nth-child(2)'
  PWSID_tag <- 'td tr:nth-child(1) td:nth-child(5)'
  ownership_tag <- 'div+ table tr:nth-child(2) td:nth-child(4)'
  max.withdrawals.mgd_tag <- 'th~ td+ td'

  function_water.system.name<- function_website %>%
    html_nodes(water.system.name_tag) %>% html_text()

  function_the_PWSID<- function_website %>%
    html_nodes(PWSID_tag) %>% html_text()

  function_the_ownership<- function_website %>%
    html_nodes(ownership_tag) %>% html_text()

  function_the_max.withdrawals.mgd<- function_website %>%
    html_nodes(max.withdrawals.mgd_tag) %>% html_text()

  Month<-c("1", "5", "9", "2", "6", "10", "3", "7", "11", "4", "8", "12")

  water.df<-data.frame(
    "Month" = as.numeric(Month),
    "Year" = rep(the_year),
    "max.withdrawals.mgd" = as.numeric(function_the_max.withdrawals.mgd)
  ) %>%
  mutate(
    Water.System.Name = !!function_water.system.name,
    PWSID = !!function_the_PWSID,
    Ownership = !!function_the_ownership,
  )

  water.df<-water.df %>%
    mutate(Date = my(paste(Month, "-", Year)))

  return(water.df)
}
```

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

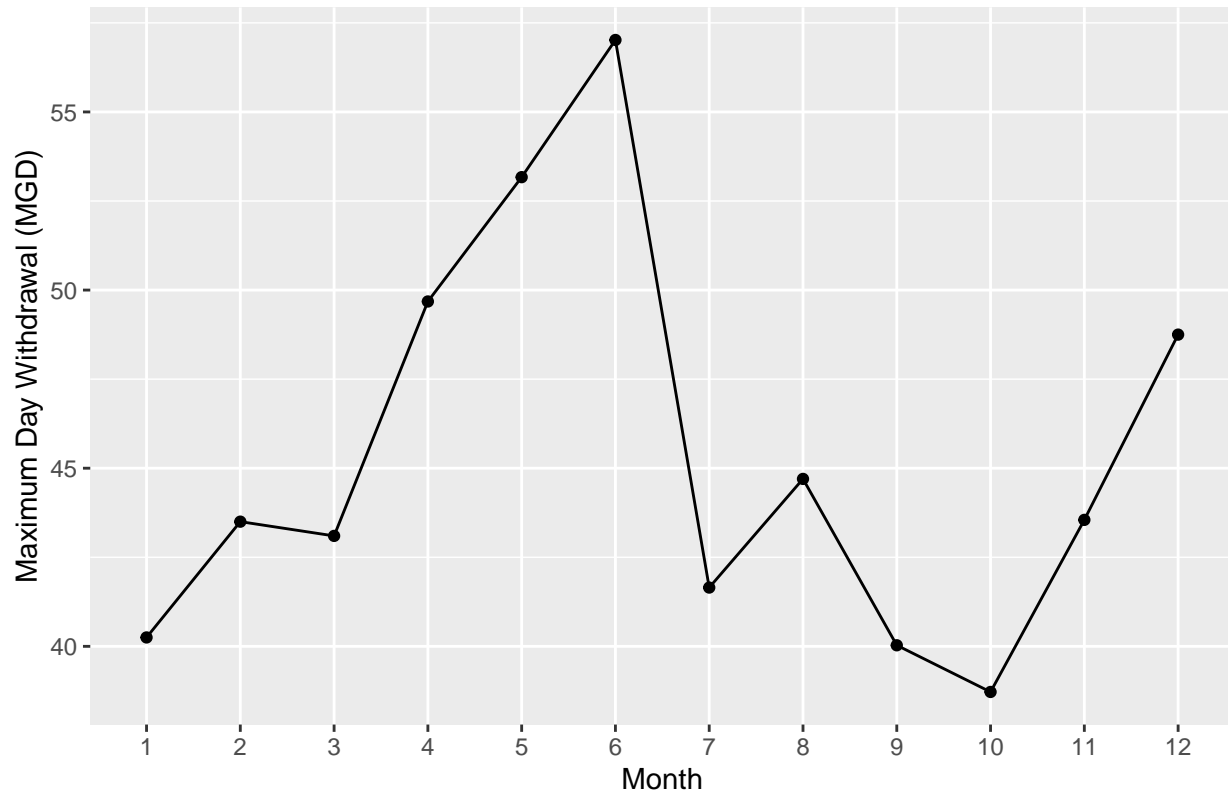
#7

```
the_df <- scrape.it(2015, '03-32-010')
view(the_df)

ggplot(the_df, aes(x = as.factor(Month), y =max.withdrawals.mgd, group=1 )) +
  geom_line() +
  geom_point()+
  labs(x = "Month",
```

```
y = "Maximum Day Withdrawal (MGD)",
title = "Maximum Water Day Use in Durham Within 2015")
```

Maximum Water Day Use in Durham Within 2015



- 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.

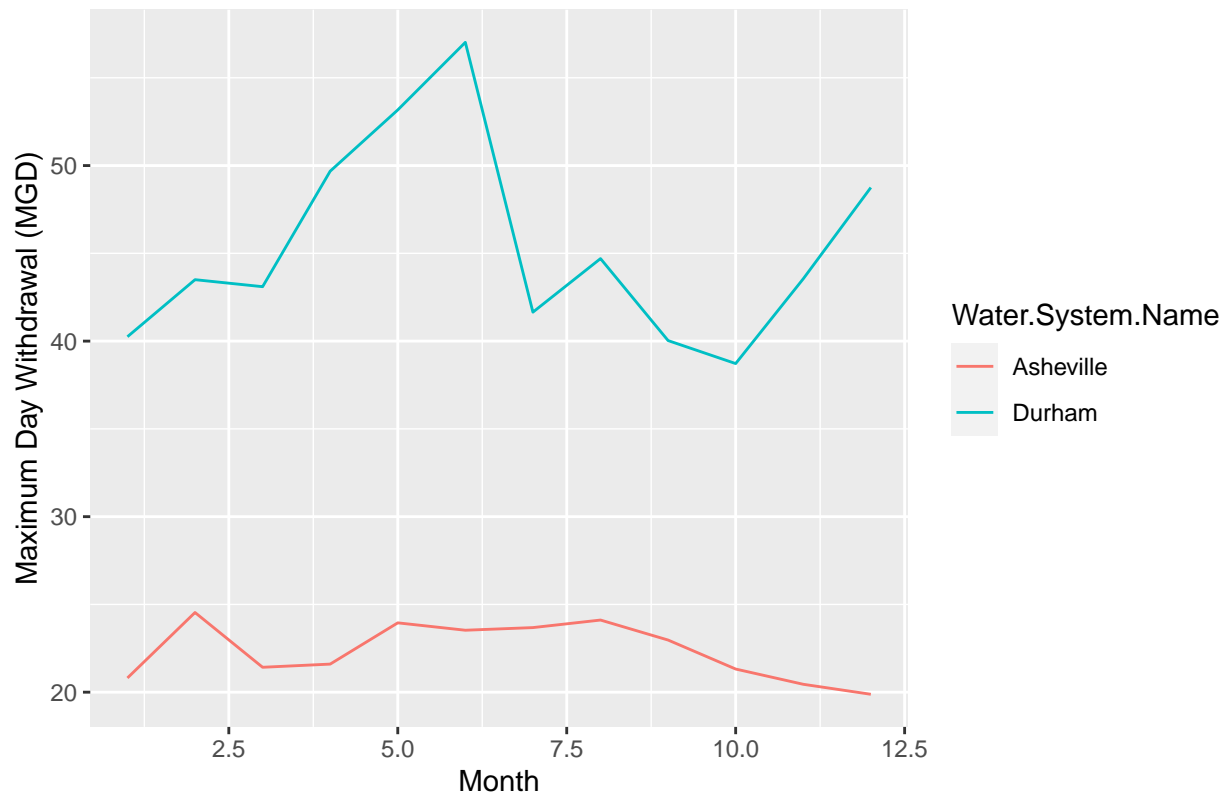
```
#8
asheville.df<-scrape.it(2015, '01-11-010')
view(asheville.df)

asheville.durham.df<-full_join(asheville.df, the_df)

## Joining, by = c("Month", "Year", "max.withdrawals.mgd", "Water.System.Name",
## "PWSID", "Ownership", "Date")

ggplot(asheville.durham.df, aes(x = Month, y =max.withdrawals.mgd, color=Water.System.Name)) +
  geom_line() +
  labs(x = "Month",
       y = "Maximum Day Withdrawal (MGD)",
       title = "Maximum Water Day Use in Asheville and Durham Within 2015")
```

Maximum Water Day Use in Asheville and Durham Within 2015



- Use the code & function you created above to plot Asheville's max daily withdrawal by months for the years 2010 thru 2021. Add a smoothed line to the plot (method = 'loess').

TIP: See Section 3.2 in the "09_Data_Scraping.Rmd" where we apply "map2()" to iteratively run a function over two inputs. Pipe the output of the map2() function to bind_rows() to combine the dataframes into a single one.

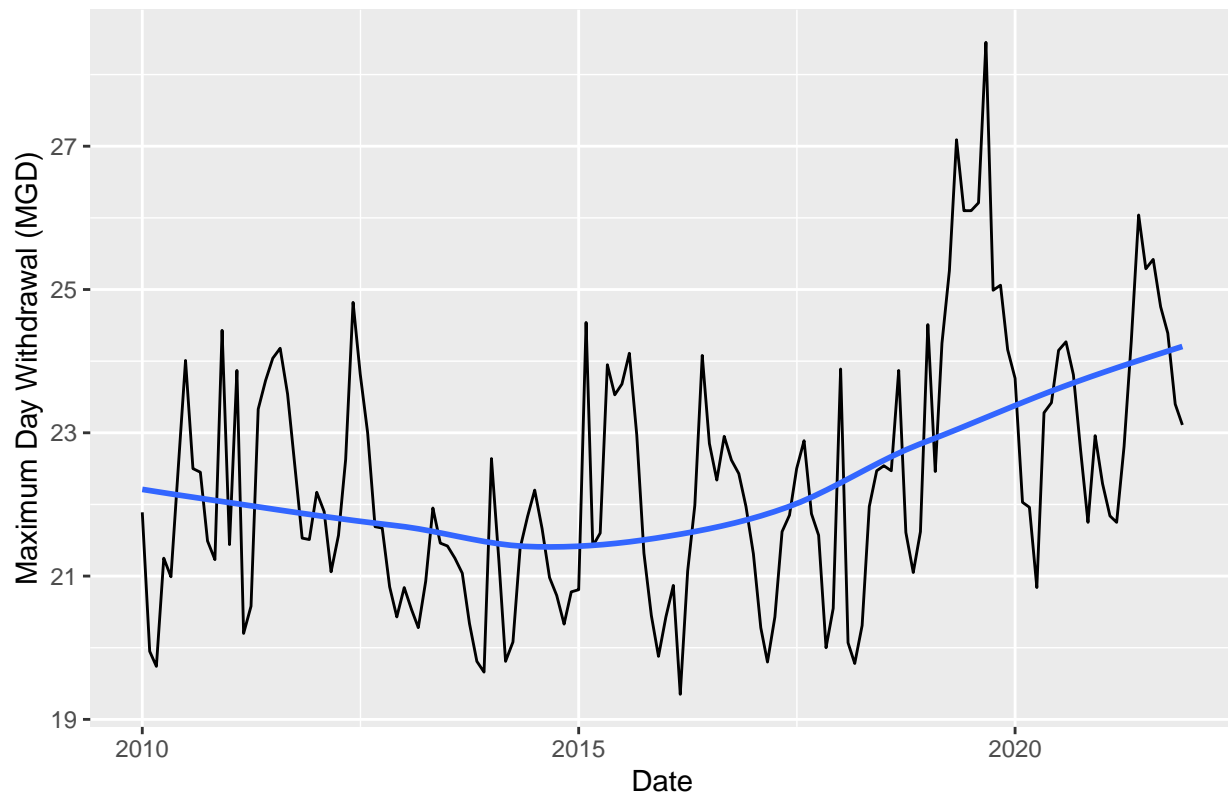
```
#9
the_years = rep(2010:2021)
my_pwsid = '01-11-010'
ash.year.dfs <- lapply(X = the_years,
                      FUN = scrape.it,
                      the_pwsid=my_pwsid)

ash.year.df <- bind_rows(ash.year.dfs)

ggplot(ash.year.df, aes(x=Date, y=max.withdrawals.mgd)) +
  geom_line() +
  geom_smooth(method="loess", se=FALSE) +
  labs(title = paste("Asheville Water usage data from 2010-2021"),
       y="Maximum Day Withdrawal (MGD)",
       x="Date")
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

Asheville Water usage data from 2010–2021



Question: Just by looking at the plot (i.e. not running statistics), does Asheville have a trend in water usage over time?

Answer: Yes, by looking at the plot, the maximum daily water withdrawal for Asheville decreases from 2010 to late 2014. On the other hand, starting from 2015, there is an increasing trend until 2021.