

# Assignment 10: Data Scraping

Susanna Jenkins

## 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)
library(dplyr)

getwd()
```

```
## [1] "/home/guest/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?pwsid=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

websiteURL <-
  read_html('https://www.ncwater.org/WUDC/app/LWSP/report.php?pwsid=03-32-010&year=2022')

websiteURL

## {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= ...
```

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 <- websiteURL %>%
  html_nodes("div+ table tr:nth-child(1) td:nth-child(2)") %>%
  html_text()

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

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

max.withdrawals.mgd <- websiteURL %>%
  html_nodes('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 (\*\*MAX) daily withdrawals across the months for 2022

#4

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

df.withdrawals.Durham
```

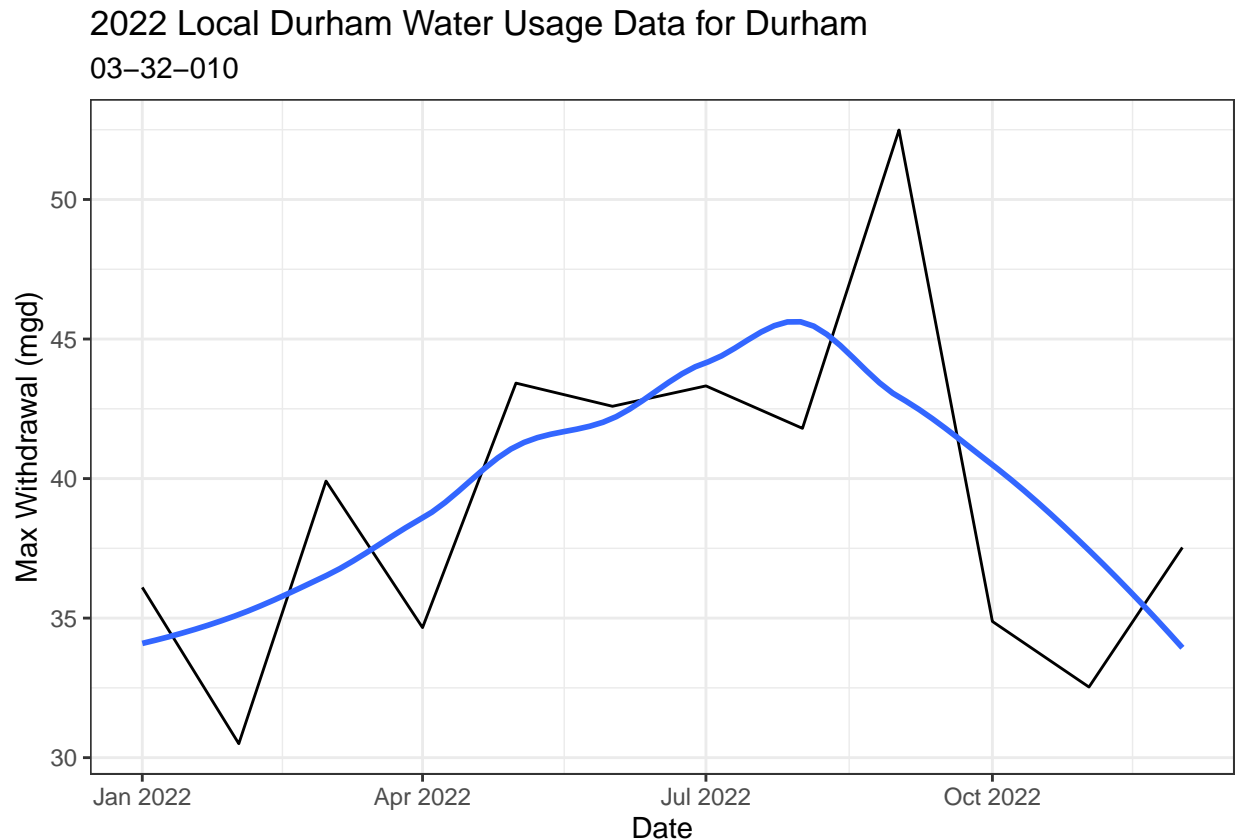
##	Month	Year	Max-Withdrawals_mgd	WaterSystemName	PWSID	Ownership
## 1	1	2022	36.10	Durham	03-32-010	Municipality
## 2	5	2022	43.42	Durham	03-32-010	Municipality
## 3	9	2022	52.49	Durham	03-32-010	Municipality
## 4	2	2022	30.50	Durham	03-32-010	Municipality
## 5	6	2022	42.59	Durham	03-32-010	Municipality
## 6	10	2022	34.88	Durham	03-32-010	Municipality
## 7	3	2022	39.91	Durham	03-32-010	Municipality
## 8	7	2022	43.32	Durham	03-32-010	Municipality
## 9	11	2022	32.53	Durham	03-32-010	Municipality
## 10	4	2022	34.66	Durham	03-32-010	Municipality
## 11	8	2022	41.80	Durham	03-32-010	Municipality
## 12	12	2022	37.53	Durham	03-32-010	Municipality
##	Date					
## 1	2022-01-01					
## 2	2022-05-01					
## 3	2022-09-01					
## 4	2022-02-01					
## 5	2022-06-01					
## 6	2022-10-01					
## 7	2022-03-01					
## 8	2022-07-01					
## 9	2022-11-01					
## 10	2022-04-01					
## 11	2022-08-01					
## 12	2022-12-01					

#5

```
ggplot(df.withdrawals.Durham, aes(x = Date, y = Max-Withdrawals_mgd)) +
  geom_line() +
```

```
geom_smooth(method="loess",se=FALSE) +
labs(title = paste("2022 Local Durham Water Usage Data for",water.system.name),
      subtitle = PWSID,
      y = "Max Withdrawal (mgd)",
      x = "Date") +
theme_bw()
```

## 'geom\_smooth()' using formula = 'y ~ x'



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

#6.

```
base.URL <-
  'https://www.ncwater.org/WUDC/app/LWSP/report.php?pwsid='
PWSID.scrape <- '03-32-010'
year.scrape <- 2022

scrape.MGD <- function(PWSID.scrape, year.scrape){
  the.website <- read_html(paste0(base.URL, PWSID.scrape, '&year=', year.scrape))
```

```

water.system.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.tag <- 'th~ td+ td'

the.water.system.name <-
  the.website %>% html_nodes(water.system.tag) %>% html_text()

the.PWSID <-
  the.website %>% html_nodes(PWSID.tag) %>% html_text()

the.ownership <-
  the.website %>% html_nodes(ownership.tag) %>% html_text()

the.max.withdrawals <-
  the.website %>% html_nodes(max.withdrawals.tag) %>% html_text()

df.withdrawals <-
  data.frame("Month" = c(1,5,9,2,6,10,3,7,11,4,8,12),
             "Year" = rep(year.scrape,12),
             "Max.Withdrawals.mgd" = as.numeric(the.max.withdrawals)) %>%
  mutate("WaterSystemName" = !!the.water.system.name,
         "PWSID" = !!the.PWSID,
         "Ownership" = !!the.ownership,
         Date = my(paste(Month,"-",Year)))

return(df.withdrawals)
}

```

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

```

#7

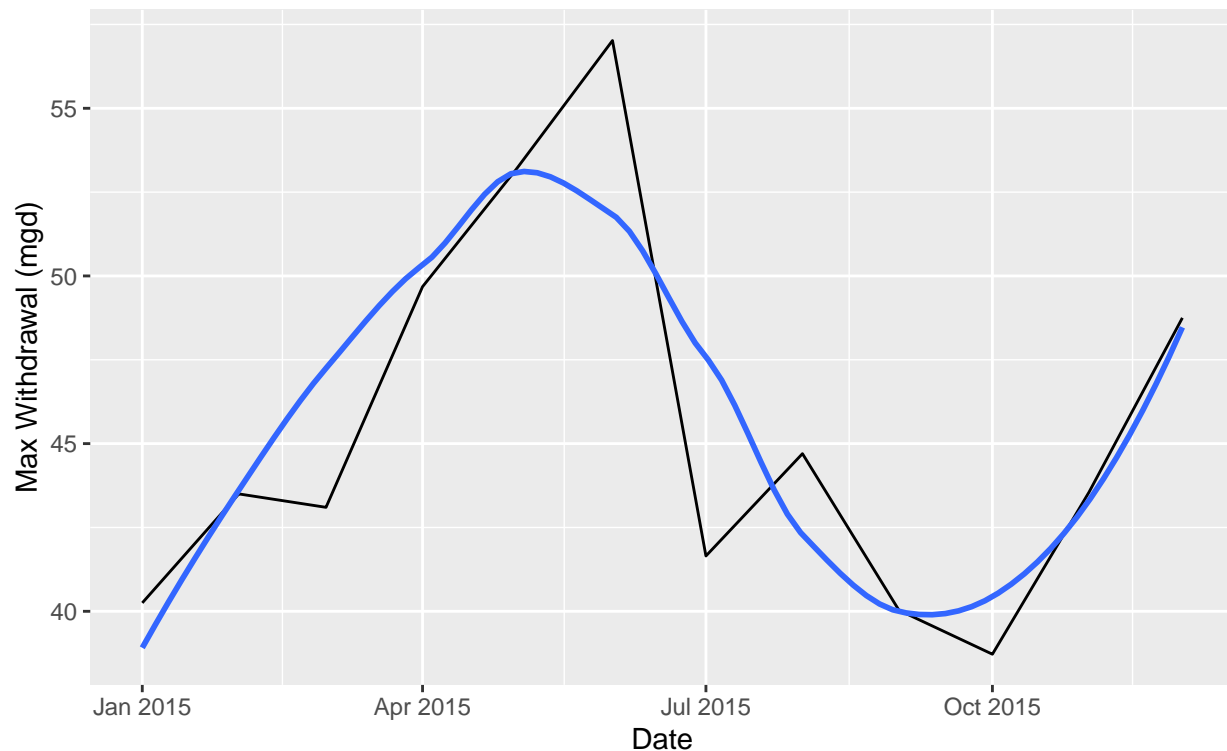
withdrawal.2015.df <- scrape.MGD('03-32-010', 2015)
view(withdrawal.2015.df)

ggplot(withdrawal.2015.df, aes(x = Date, y = Max.Withdrawals.mgd)) +
  geom_line() +
  geom_smooth(se = FALSE) +
  labs(title = paste(year.scrape,"Local Durham Water Usage Data for Durham"),
       subtitle = PWSID.scrape,
       x = "Date",
       y = "Max Withdrawal (mgd)")

```

```
## 'geom_smooth()' using method = 'loess' and formula = 'y ~ x'
```

## 2022 Local Durham Water Usage Data for Durham 03-32-010



- 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.

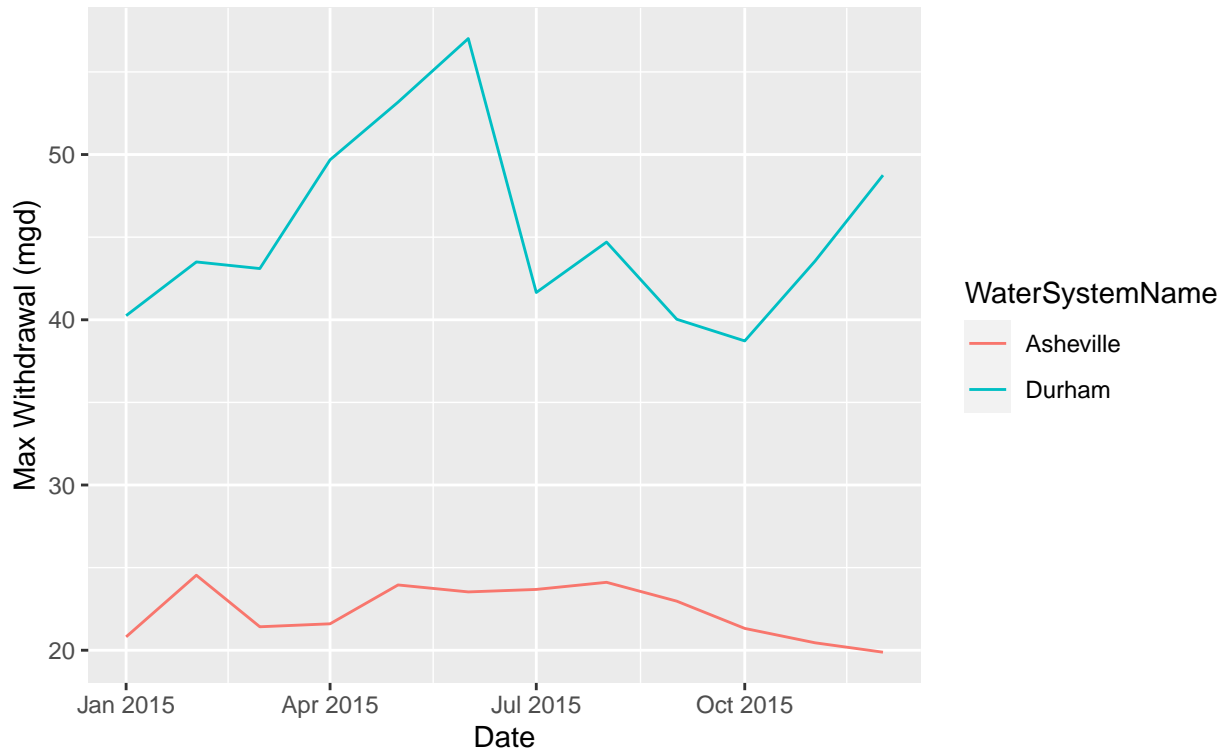
#8

```
withdrawal.Asheville.2015.df <- scrape.MGD('01-11-010', 2015)
view(withdrawal.Asheville.2015.df)

Asheville.Durham <- rbind(withdrawal.Asheville.2015.df, withdrawal.2015.df)

ggplot(Asheville.Durham, aes(x = Date, y = Max.Withdrawals.mgd,
                             color=WaterSystemName)) +
  geom_line() +
  labs(title = "2015 Local Water Usage Data for Durham & Asheville",
        subtitle = PWSID.scrape,
        x = "Date",
        y = "Max Withdrawal (mgd)")
```

## 2015 Local Water Usage Data for Durham & Asheville 03-32-010



- 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

years = rep(2010:2021)
PWSID.scrape = '01-11-010'

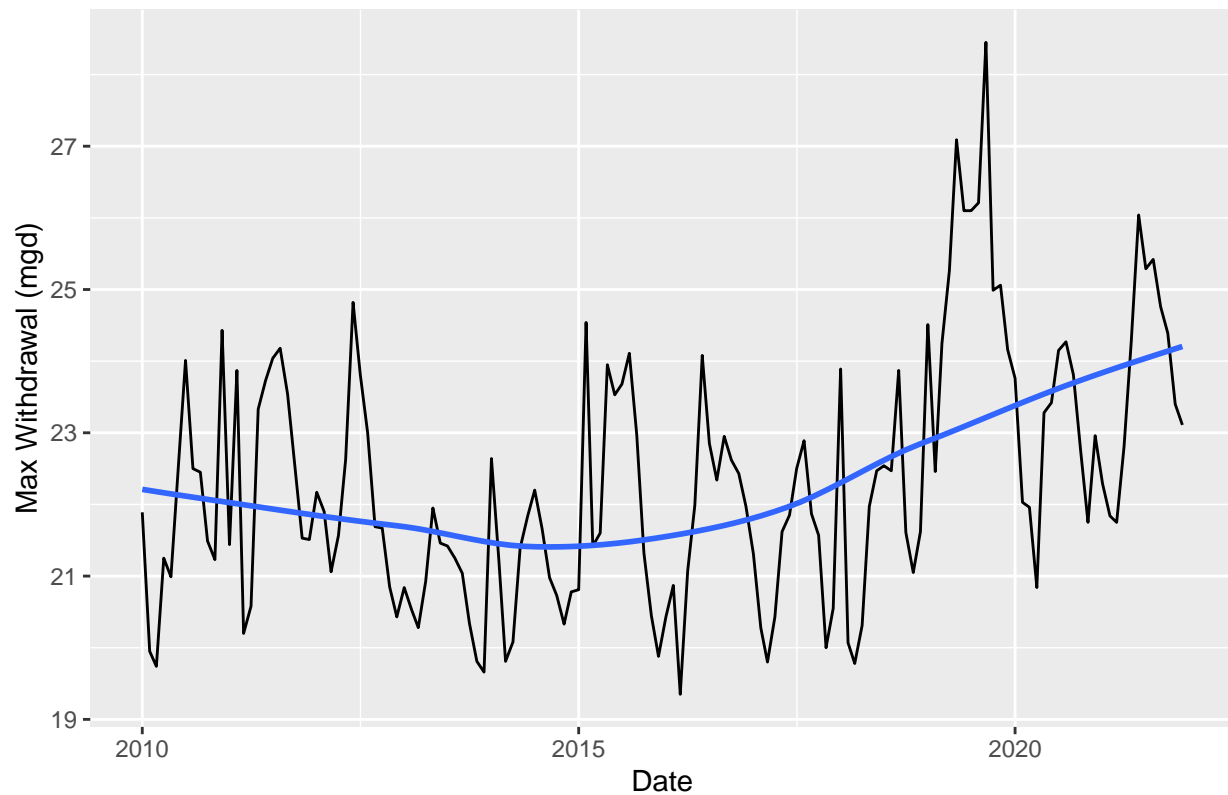
Asheville.2010.2021 <- map2(PWSID.scrape, years, scrape.MGD)

Asheville.10.21 <- bind_rows(Asheville.2010.2021)

ggplot(Asheville.10.21, aes(x=Date, y=Max.Withdrawals.mgd)) +
  geom_line() +
  geom_smooth(method="loess", se=FALSE) +
  labs(title = "2010-2021 Water usage data for Asheville",
       y="Max Withdrawal (mgd)",
       x="Date")
```

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

2010–2021 Water usage data for Asheville



Question: Just by looking at the plot (i.e. not running statistics), does Asheville have a trend in water usage over time? Over time, it looks like it is slightly increasing over time. It appears there is seasonality in the data but generally it looks to be increasing overall from 2010 to 2021.