

# Assignment 10: Data Scraping

Student Name

## 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(lubridate);library(viridis);library(here)
here()
```

```
## [1] "/home/guest/EDE_Fall2024"
```

```
#install.packages("rvest")
library(rvest)

#install.packages("dataRetrieval")
library(dataRetrieval)

#install.packages("tidycensus")
library(tidycensus)

# Set theme
mytheme <- theme_gray() +
  theme(axis.text = element_text(angle = 45, color = "purple"),
        legend.position = "top")
theme_set(mytheme)
```

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.
- Note the web address: <https://www.ncwater.org/WUDC/app/LWSP/report.php?pwsid=03-32-010&year=2023>

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

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

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
water_system <- the_URL %>%
  html_nodes('div+ table tr:nth-child(1) td:nth-child(2)') %>% html_text()
the_PWSID<- the_URL %>%
  html_nodes('td tr:nth-child(1) td:nth-child(5)') %>% html_text()
the_ownership <- the_URL %>%
  html_nodes('div+ table tr:nth-child(2) td:nth-child(4)') %>% html_text()

the_MGD <- the_URL %>%
  html_nodes('th~ td+ td') %>% html_text()
the_MGD = as.numeric(the_MGD)
```

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 maximum daily withdrawals across the months for 2023, making sure the months are presented in proper sequence.

```
#4
Month <- the_URL %>%
  html_nodes('.fancy-table:nth-child(31) tr+ tr th') %>% html_text()

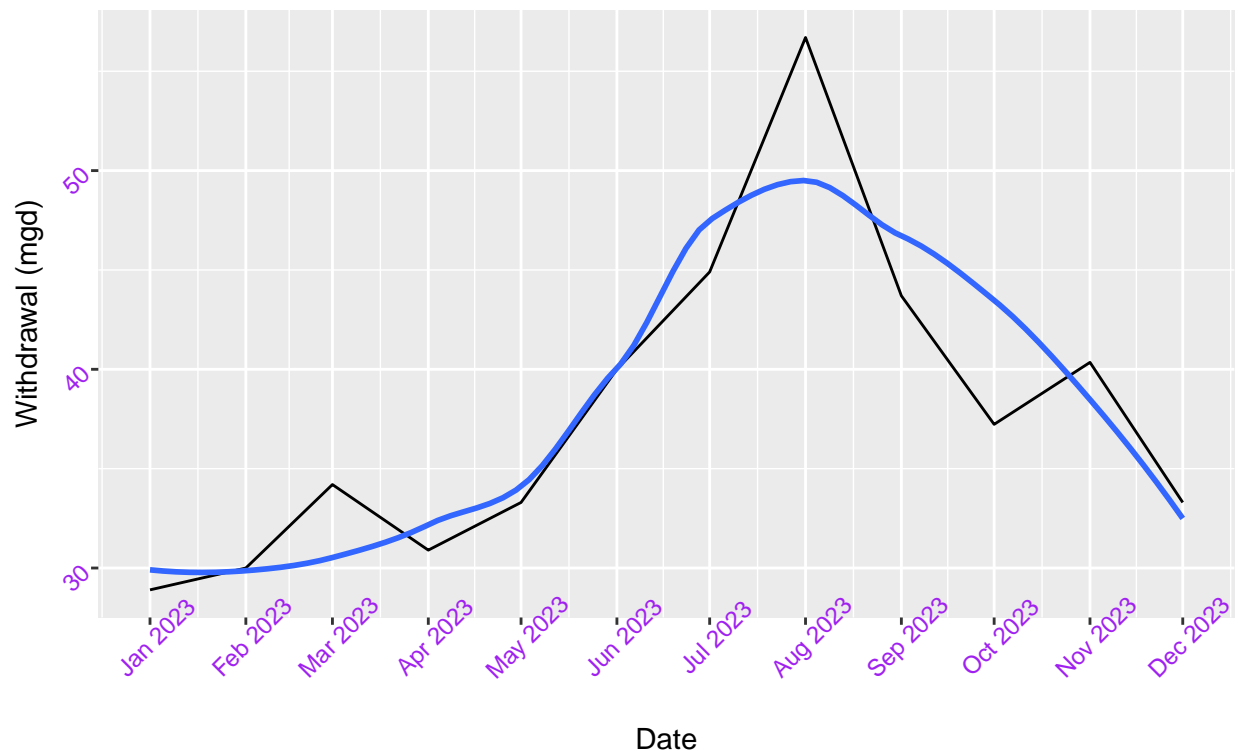
withdrawal_dataframe <- data.frame(
  "Water_System_Name"= water_system,
  "PWSID" = the_PWSID,
  "Ownership" = the_ownership,
  "Maximum_Day_Use" = the_MGD,
  "Month"= Month,
  "Year" = 2023)

withdrawal_dataframe <- withdrawal_dataframe %>%
  mutate(Month = factor(Month, levels =
    c("Jan", "Feb", "Mar", "Apr", "May", "Jun", "Jul", "Aug",
      "Sep", "Oct", "Nov", "Dec"))) %>%
  arrange(Month) %>%
  mutate(Date = my(paste0(Month, "-", Year)))

#5
ggplot(withdrawal_dataframe, aes(x=Date, y= Maximum_Day_Use)) +
  geom_line() +
  geom_smooth(method="loess", se=FALSE) +
  labs(title = paste("2022 Water Usage Data for", water_system),
    subtitle = the_ownership,
    y="Withdrawal (mgd)",
    x="Date") +
  scale_x_date(date_labels = "%b %Y", date_breaks = "1 month")

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

## 2022 Water Usage Data for Durham Municipality



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.**

```
#6.
scrape.it <- function(el_PWSID, the_year){

  the_URL <-
    read_html(paste0(
      'https://www.ncwater.org/WUDC/app/LWSP/report.php?pwsid=',
      el_PWSID, '&year=', the_year))

  #Set the element address variables (determined in the previous step)
  water_system_tag <- 'div+ table tr:nth-child(1) td:nth-child(2)'
  the_PWSID_tag <- 'td tr:nth-child(1) td:nth-child(5)'
  the_ownership_tag <- 'div+ table tr:nth-child(2) td:nth-child(4)'
  the_MGD_tag <- 'th~ td+ td'

  #Scrape the data items
  water_system <- the_URL %>%
    html_nodes(water_system_tag) %>% html_text()
  the_PWSID <- the_URL %>%
    html_nodes(the_PWSID_tag) %>% html_text()
  the_ownership <- the_URL %>%
```

```

  html_nodes(the_ownership_tag)%>% html_text()
the_MGD <- the_URL %>%
  html_nodes(the_MGD_tag)%>% html_text()

#Convert to a dataframe
PWSID_dataframe <- data.frame(
  "Water_System_Name"= water_system,
  "PWSID" = the_PWSID,
  "Ownership" = the_ownership,
  "Maximum_Day_Use" = as.numeric(the_MGD),
  "Month"= as.factor(Month),
  "Year" = the_year)

PWSID_dataframe <-PWSID_dataframe %>%
mutate(Month = factor(Month, levels =
                      c("Jan", "Feb", "Mar", "Apr", "May", "Jun", "Jul", "Aug",
                        "Sep", "Oct", "Nov", "Dec"))) %>%

arrange(Month) %>%
mutate(Date = my(paste0(Month, "-", Year)))

return(PWSID_dataframe)
}

```

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

```

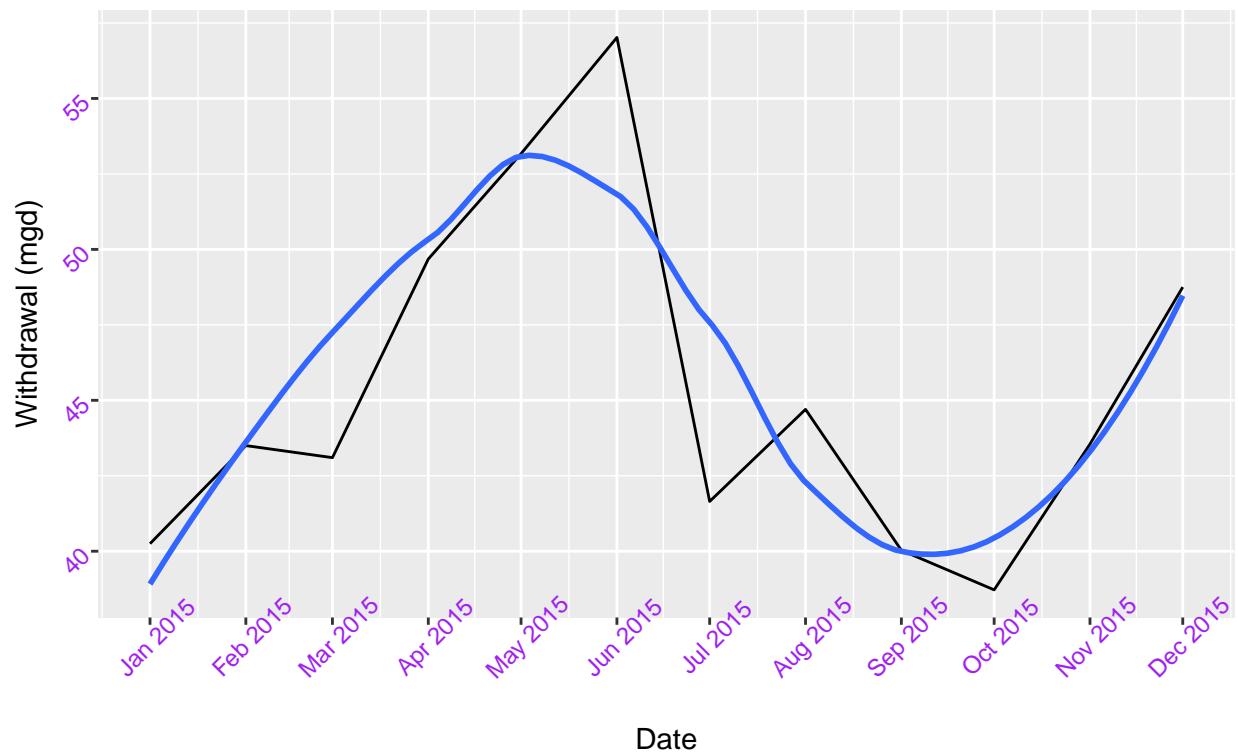
#7
df_2015 <- scrape.it('03-32-010', 2015)

ggplot(df_2015,aes(x=Date,y= Maximum_Day_Use)) +
  geom_line() +
  geom_smooth(method="loess",se=FALSE) +
  labs(title = paste(df_2015$Year, "Water Usage Data for", df_2015$Water_System_Name),
       subtitle = df_2015$Ownership,
       y="Withdrawal (mgd)",
       x="Date")+
  scale_x_date(date_labels = "%b %Y", date_breaks = "1 month")

```

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

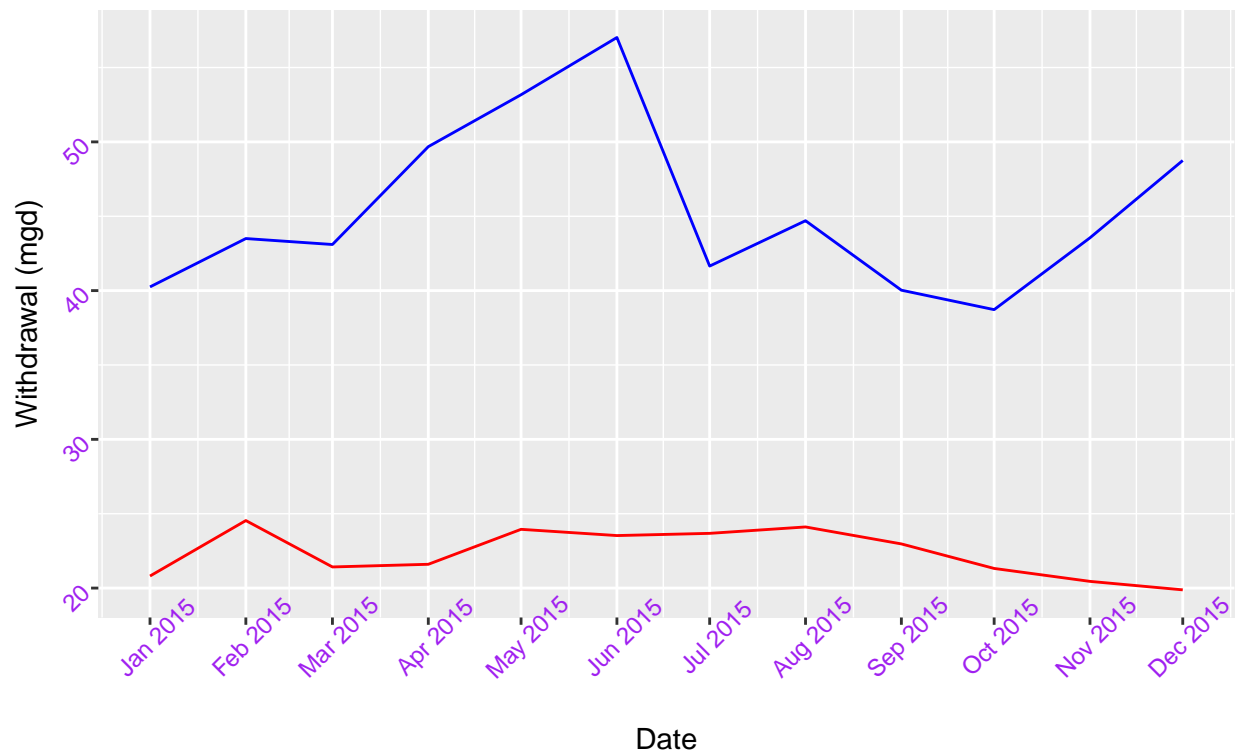
## 2015 Water Usage Data for Durham Municipality



- 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
df_Ashville_2015 <- scrape.it("01-11-010", 2015)
ggplot() +
  geom_line(data = df_2015,
            aes(x=Date,y= Maximum_Day_Use), color="blue") +
  geom_line(data = df_Ashville_2015,
            aes(x=Date,y= Maximum_Day_Use), color="red") +
  labs(title = paste("2015 Water Usage Data for Asheville and Durham"),
        subtitle = "Municipalities",
        y="Withdrawal (mgd)",
        x="Date")+
  scale_x_date(date_labels = "%b %Y", date_breaks = "1 month")
```

## 2015 Water Usage Data for Asheville and Durham Municipalities



- 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.

```
#9
df_Ashville_2018 <- scrape.it("01-11-010", 2018)
df_Ashville_2019 <- scrape.it("01-11-010", 2019)
df_Ashville_2020 <- scrape.it("01-11-010", 2020)
df_Ashville_2021 <- scrape.it("01-11-010", 2021)
df_Ashville_2022 <- scrape.it("01-11-010", 2022)

the_facility_id <- "01-11-010"
the_years <- c(2018, 2019, 2020, 2021, 2022)
dfs_Ashville <- map2(the_facility_id, the_years, scrape.it)
df_Ashville <- bind_rows(dfs_Ashville)

ggplot(df_Ashville, aes(y=Maximum_Day_Use, x=Month, group=1)) +
  geom_line(aes(color = Year, group = Year)) +
  geom_smooth(method = "loess", se=FALSE, color = "red", size = 1.2) +
  labs(title = paste("2018-2022 Water usage data for Asheville"),
```

```

subtitle = "Municipality",
y="Withdrawal (MGD)",
x="Month")

```

```

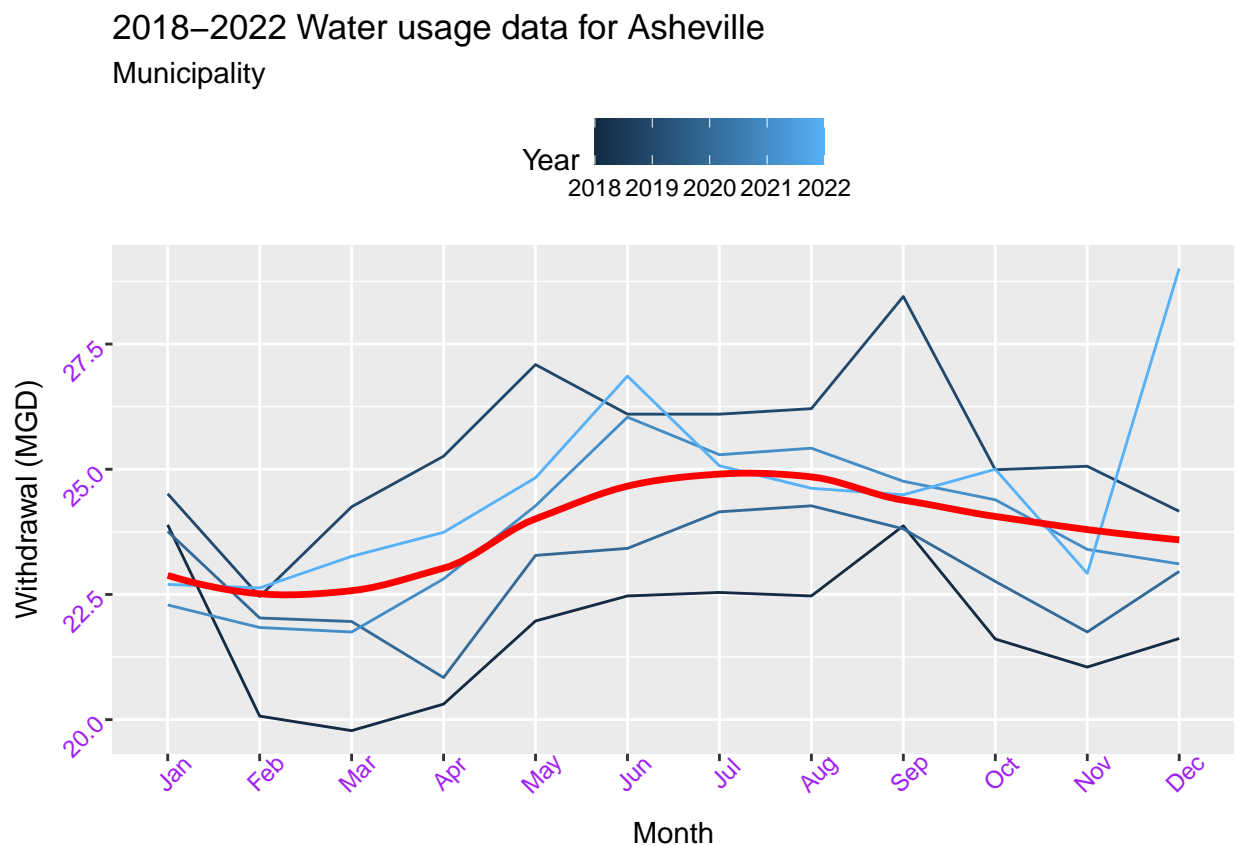
## Warning: Using 'size' aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use 'linewidth' instead.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.

```

```

## '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: Asheville's trend increases in the summer months which makes sense due to the increased heat, and then declines in the winter months. Overall water trends seem to be pretty similar throughout the years, except 2022 seems to have a steep increase in November for unknown reasons.