

Assignment 4: Data Wrangling

Maia Griffith

OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics on Data Wrangling

Directions

1. Rename this file `<FirstLast>_A04_DataWrangling.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 to **answer the questions** in this assignment document.
5. When you have completed the assignment, **Knit** the text and code into a single PDF file.
6. Ensure that code in code chunks does not extend off the page in the PDF.

The completed exercise is due on Thursday, Sept 28th @ 5:00pm.

Set up your session

- 1a. Load the `tidyverse`, `lubridate`, and `here` packages into your session.
 - 1b. Check your working directory.
 - 1c. Read in all four raw data files associated with the EPA Air dataset, being sure to set string columns to be read in as factors. See the README file for the EPA air datasets for more information (especially if you have not worked with air quality data previously).
2. Apply the `glimpse()` function to reveal the dimensions, column names, and structure of each dataset.

```
#1a Loading the packages
```

```
library(tidyverse)
library(lubridate)
library(here)
```

```
#1b Checking the working directory
```

```
getwd()
```

```
## [1] "/home/guest/ENV872/maia-g/EDE_Fall2023"
```

```
#1c Reading in each file, first the two O3 files and then the two PM25 files
```

```
Air.O3.NC18 <- read.csv(file=here("./Data/Raw/EPAair_O3_NC2018_raw.csv"),
                        stringsAsFactors = TRUE)
```

```
Air.O3.NC19 <- read.csv(file=here("./Data/Raw/EPAair_O3_NC2019_raw.csv"),
```

```

stringsAsFactors = TRUE)

Air.PM25.NC18 <-read.csv(file=here("./Data/Raw/EPAair_PM25_NC2018_raw.csv"),
                        stringsAsFactors = TRUE)
Air.PM25.NC19 <-read.csv(file=here("./Data/Raw/EPAair_PM25_NC2019_raw.csv"),
                        stringsAsFactors = TRUE)

#2 using the `glimpse()` function
glimpse(Air.03.NC18)

## Rows: 9,737
## Columns: 20
## $ Date <fct> 03/01/2018, 03/02/2018, 03/03/201~
## $ Source <fct> AQS, AQS, AQS, AQS, AQS, AQS, AQS~
## $ Site.ID <int> 370030005, 370030005, 370030005, ~
## $ POC <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Max.8.hour.Ozone.Concentration <dbl> 0.043, 0.046, 0.047, 0.049, 0.047~
## $ UNITS <fct> ppm, ppm, ppm, ppm, ppm, ppm, ppm~
## $ DAILY_AQI_VALUE <int> 40, 43, 44, 45, 44, 28, 33, 41, 4~
## $ Site.Name <fct> Taylorsville Liledoun, Taylorsvil~
## $ DAILY_OBS_COUNT <int> 17, 17, 17, 17, 17, 17, 17, 17, 1~
## $ PERCENT_COMPLETE <dbl> 100, 100, 100, 100, 100, 100, 100~
## $ AQS_PARAMETER_CODE <int> 44201, 44201, 44201, 44201, 44201~
## $ AQS_PARAMETER_DESC <fct> Ozone, Ozone, Ozone, Ozone, Ozone~
## $ CBSA_CODE <int> 25860, 25860, 25860, 25860, 25860~
## $ CBSA_NAME <fct> "Hickory-Lenoir-Morganton, NC", "~
## $ STATE_CODE <int> 37, 37, 37, 37, 37, 37, 37, 37, 3~
## $ STATE <fct> North Carolina, North Carolina, N~
## $ COUNTY_CODE <int> 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, ~
## $ COUNTY <fct> Alexander, Alexander, Alexander, ~
## $ SITE_LATITUDE <dbl> 35.9138, 35.9138, 35.9138, 35.913~
## $ SITE_LONGITUDE <dbl> -81.191, -81.191, -81.191, -81.19~

```

```

glimpse(Air.03.NC19)

## Rows: 10,592
## Columns: 20
## $ Date <fct> 01/01/2019, 01/02/2019, 01/03/201~
## $ Source <fct> AirNow, AirNow, AirNow, AirNow, A~
## $ Site.ID <int> 370030005, 370030005, 370030005, ~
## $ POC <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Max.8.hour.Ozone.Concentration <dbl> 0.029, 0.018, 0.016, 0.022, 0.037~
## $ UNITS <fct> ppm, ppm, ppm, ppm, ppm, ppm, ppm~
## $ DAILY_AQI_VALUE <int> 27, 17, 15, 20, 34, 34, 27, 35, 3~
## $ Site.Name <fct> Taylorsville Liledoun, Taylorsvil~
## $ DAILY_OBS_COUNT <int> 24, 24, 24, 24, 24, 24, 24, 24, 2~
## $ PERCENT_COMPLETE <dbl> 100, 100, 100, 100, 100, 100, 100~
## $ AQS_PARAMETER_CODE <int> 44201, 44201, 44201, 44201, 44201~
## $ AQS_PARAMETER_DESC <fct> Ozone, Ozone, Ozone, Ozone, Ozone~
## $ CBSA_CODE <int> 25860, 25860, 25860, 25860, 25860~
## $ CBSA_NAME <fct> "Hickory-Lenoir-Morganton, NC", "~
## $ STATE_CODE <int> 37, 37, 37, 37, 37, 37, 37, 37, 3~

```

```
## $ STATE <fct> North Carolina, North Carolina, N-
## $ COUNTY_CODE <int> 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, ~
## $ COUNTY <fct> Alexander, Alexander, Alexander, ~
## $ SITE_LATITUDE <dbl> 35.9138, 35.9138, 35.9138, 35.913~
## $ SITE_LONGITUDE <dbl> -81.191, -81.191, -81.191, -81.19~
```

`glimpse(Air.PM25.NC18)`

```
## Rows: 8,983
## Columns: 20
## $ Date <fct> 01/02/2018, 01/05/2018, 01/08/2018, 01/~
## $ Source <fct> AQS, AQS, AQS, AQS, AQS, AQS, AQS, AQS, ~
## $ Site.ID <int> 370110002, 370110002, 370110002, 370110~
## $ POC <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Mean.PM2.5.Concentration <dbl> 2.9, 3.7, 5.3, 0.8, 2.5, 4.5, 1.8, 2.5, ~
## $ UNITS <fct> ug/m3 LC, ug/m3 LC, ug/m3 LC, ug/m3 LC, ~
## $ DAILY_AQI_VALUE <int> 12, 15, 22, 3, 10, 19, 8, 10, 18, 7, 24~
## $ Site.Name <fct> Linville Falls, Linville Falls, Linvill~
## $ DAILY_OBS_COUNT <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ PERCENT_COMPLETE <dbl> 100, 100, 100, 100, 100, 100, 100, 100, ~
## $ AQS_PARAMETER_CODE <int> 88502, 88502, 88502, 88502, 88502, 8850~
## $ AQS_PARAMETER_DESC <fct> Acceptable PM2.5 AQI & Speciation Mass, ~
## $ CBSA_CODE <int> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
## $ CBSA_NAME <fct> "", "", "", "", "", "", "", "", "", "", ~
## $ STATE_CODE <int> 37, 37, 37, 37, 37, 37, 37, 37, 37, 37, ~
## $ STATE <fct> North Carolina, North Carolina, North C~
## $ COUNTY_CODE <int> 11, 11, 11, 11, 11, 11, 11, 11, 11, 11, ~
## $ COUNTY <fct> Avery, Avery, Avery, Avery, Avery, Aver~
## $ SITE_LATITUDE <dbl> 35.97235, 35.97235, 35.97235, 35.97235, ~
## $ SITE_LONGITUDE <dbl> -81.93307, -81.93307, -81.93307, -81.93~
```

`glimpse(Air.PM25.NC19)`

```
## Rows: 8,581
## Columns: 20
## $ Date <fct> 01/03/2019, 01/06/2019, 01/09/2019, 01/~
## $ Source <fct> AQS, AQS, AQS, AQS, AQS, AQS, AQS, AQS, ~
## $ Site.ID <int> 370110002, 370110002, 370110002, 370110~
## $ POC <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Mean.PM2.5.Concentration <dbl> 1.6, 1.0, 1.3, 6.3, 2.6, 1.2, 1.5, 1.5, ~
## $ UNITS <fct> ug/m3 LC, ug/m3 LC, ug/m3 LC, ug/m3 LC, ~
## $ DAILY_AQI_VALUE <int> 7, 4, 5, 26, 11, 5, 6, 6, 15, 7, 14, 20~
## $ Site.Name <fct> Linville Falls, Linville Falls, Linvill~
## $ DAILY_OBS_COUNT <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ PERCENT_COMPLETE <dbl> 100, 100, 100, 100, 100, 100, 100, 100, ~
## $ AQS_PARAMETER_CODE <int> 88502, 88502, 88502, 88502, 88502, 8850~
## $ AQS_PARAMETER_DESC <fct> Acceptable PM2.5 AQI & Speciation Mass, ~
## $ CBSA_CODE <int> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
## $ CBSA_NAME <fct> "", "", "", "", "", "", "", "", "", "", ~
## $ STATE_CODE <int> 37, 37, 37, 37, 37, 37, 37, 37, 37, 37, ~
## $ STATE <fct> North Carolina, North Carolina, North C~
## $ COUNTY_CODE <int> 11, 11, 11, 11, 11, 11, 11, 11, 11, 11, ~
## $ COUNTY <fct> Avery, Avery, Avery, Avery, Avery, Aver~
```

```
## $ SITE_LATITUDE      <dbl> 35.97235, 35.97235, 35.97235, 35.97235,~
## $ SITE_LONGITUDE     <dbl> -81.93307, -81.93307, -81.93307, -81.93~
```

Wrangle individual datasets to create processed files.

3. Change the Date columns to be date objects.
4. Select the following columns: Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY, SITE_LATITUDE, SITE_LONGITUDE
5. For the PM2.5 datasets, fill all cells in AQS_PARAMETER_DESC with “PM2.5” (all cells in this column should be identical).
6. Save all four processed datasets in the Processed folder. Use the same file names as the raw files but replace “raw” with “processed”.

```
#3 Date col is in 2 digit month, 2 digit day, and 4 digit year format
# Looking at the class for each to show they are currently factors.
class(Air.03.NC18$Date)
```

```
## [1] "factor"
```

```
class(Air.03.NC19$Date)
```

```
## [1] "factor"
```

```
class(Air.PM25.NC18$Date)
```

```
## [1] "factor"
```

```
class(Air.PM25.NC19$Date)
```

```
## [1] "factor"
```

```
# Using lubridate functions to change date columns in each dataset.
Air.03.NC18$Date <- mdy(Air.03.NC18$Date)
Air.03.NC19$Date <- mdy(Air.03.NC19$Date)

Air.PM25.NC18$Date <- mdy(Air.PM25.NC18$Date)
Air.PM25.NC19$Date <- mdy(Air.PM25.NC19$Date)

# Looking at the class for each to show they are now date objects.
class(Air.03.NC18$Date)
```

```
## [1] "Date"
```

```
class(Air.03.NC19$Date)
```

```
## [1] "Date"
```

```
class(Air.PM25.NC18$Date)
```

```
## [1] "Date"
```

```
class(Air.PM25.NC19$Date)
```

```
## [1] "Date"
```

```
#4 Selecting only certain columns. I used the select() function.
```

```
Air.03.NC18.processed <-
```

```
  Air.03.NC18 %>%
```

```
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,  
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

```
Air.03.NC19.processed <-
```

```
  Air.03.NC19 %>%
```

```
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,  
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

```
Air.PM25.NC18.processed <-
```

```
  Air.PM25.NC18 %>%
```

```
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,  
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

```
Air.PM25.NC19.processed <-
```

```
  Air.PM25.NC19 %>%
```

```
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,  
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

```
#5 Changing column values to PM2.5 in both datasets
```

```
Air.PM25.NC18.processed$AQS_PARAMETER_DESC <- "PM2.5"
```

```
Air.PM25.NC19.processed$AQS_PARAMETER_DESC <- "PM2.5"
```

```
#6 Saving processed data
```

```
write.csv(Air.03.NC18.processed, row.names = FALSE,  
         file = "./Data/Processed/EPAair_03_NC2018_processed.csv")
```

```
write.csv(Air.03.NC19.processed, row.names = FALSE,  
         file = "./Data/Processed/EPAair_03_NC2019_processed.csv")
```

```
write.csv(Air.PM25.NC18.processed, row.names = FALSE,  
         file = "./Data/Processed/EPAair_PM25_NC2018_processed.csv")
```

```
write.csv(Air.PM25.NC19.processed, row.names = FALSE,  
         file = "./Data/Processed/EPAair_PM25_NC2019_processed.csv")
```

Combine datasets

7. Combine the four datasets with `rbind`. Make sure your column names are identical prior to running this code.
8. Wrangle your new dataset with a pipe function (`%>%`) so that it fills the following conditions:

- Include only sites that the four data frames have in common: “Linville Falls”, “Durham Armory”, “Leggett”, “Hattie Avenue”, “Clemmons Middle”, “Mendenhall School”, “Frying Pan Mountain”, “West Johnston Co.”, “Garinger High School”, “Castle Hayne”, “Pitt Agri. Center”, “Bryson City”, “Millbrook School” (the function `intersect` can figure out common factor levels - but it will include sites with missing site information, which you don’t want...)
 - Some sites have multiple measurements per day. Use the split-apply-combine strategy to generate daily means: group by date, site name, AQS parameter, and county. Take the mean of the AQI value, latitude, and longitude.
 - Add columns for “Month” and “Year” by parsing your “Date” column (hint: `lubridate` package)
 - Hint: the dimensions of this dataset should be 14,752 x 9.
9. Spread your datasets such that AQI values for ozone and PM2.5 are in separate columns. Each location on a specific date should now occupy only one row.
 10. Call up the dimensions of your new tidy dataset.
 11. Save your processed dataset with the following file name: “EPAair_O3_PM25_NC1819_Processed.csv”

#7 Combining datasets

```
# Making sure all the datasets have the same column names (they do)
colnames(Air.O3.NC18.processed)
```

```
## [1] "Date"           "DAILY_AQI_VALUE" "Site.Name"
## [4] "AQS_PARAMETER_DESC" "COUNTY"         "SITE_LATITUDE"
## [7] "SITE_LONGITUDE"
```

```
colnames(Air.O3.NC19.processed)
```

```
## [1] "Date"           "DAILY_AQI_VALUE" "Site.Name"
## [4] "AQS_PARAMETER_DESC" "COUNTY"         "SITE_LATITUDE"
## [7] "SITE_LONGITUDE"
```

```
colnames(Air.PM25.NC18.processed)
```

```
## [1] "Date"           "DAILY_AQI_VALUE" "Site.Name"
## [4] "AQS_PARAMETER_DESC" "COUNTY"         "SITE_LATITUDE"
## [7] "SITE_LONGITUDE"
```

```
colnames(Air.PM25.NC19.processed)
```

```
## [1] "Date"           "DAILY_AQI_VALUE" "Site.Name"
## [4] "AQS_PARAMETER_DESC" "COUNTY"         "SITE_LATITUDE"
## [7] "SITE_LONGITUDE"
```

```
# Using rbind to combine.
Air.Combo <- rbind(Air.03.NC18.processed,
                  Air.03.NC19.processed,
                  Air.PM25.NC18.processed,
                  Air.PM25.NC19.processed)

# Checking the dimensions to make sure it combined right.
dim(Air.Combo)
```

```
## [1] 37893      7
```

```
#8 Making a pipe.
Air.Combo.Subset <-
  Air.Combo %>%
  filter(Site.Name %in% c("Linville Falls", "Durham Armory", "Leggett",
                        "Hattie Avenue", "Clemmons Middle",
                        "Mendenhall School", "Frying Pan Mountain",
                        "West Johnston Co.", "Garinger High School",
                        "Castle Hayne", "Pitt Agri. Center",
                        "Bryson City", "Millbrook School")) %>%
  group_by(Date, Site.Name, AQS_PARAMETER_DESC, COUNTY) %>%
  summarise(Mean_AQI = mean(DAILY_AQI_VALUE),
            Mean_Lat = mean(SITE_LATITUDE),
            Mean_Long = mean(SITE_LONGITUDE)) %>%
  mutate(Month = month(Date),
         Year = year(Date))
```

```
## 'summarise()' has grouped output by 'Date', 'Site.Name', 'AQS_PARAMETER_DESC'.
## You can override using the '.groups' argument.
```

```
dim(Air.Combo.Subset) # Showing dimensions match the goal
```

```
## [1] 14752      9
```

```
#9 Spreading the data so PM2.5 and Ozone have their own column
Air.Combo.Subset.Spread <-
  pivot_wider(Air.Combo.Subset,
              names_from = AQS_PARAMETER_DESC,
              values_from = Mean_AQI)

#10
dim(Air.Combo.Subset.Spread)
```

```
## [1] 8976      9
```

```
#11
write.csv(Air.Combo.Subset.Spread, row.names = FALSE,
          file = "./Data/Processed/EPAair_03_PM25_NC1819_Processed.csv")
```

Generate summary tables

12. Use the split-apply-combine strategy to generate a summary data frame. Data should be grouped by site, month, and year. Generate the mean AQI values for ozone and PM2.5 for each group. Then, add a pipe to remove instances where mean **ozone** values are not available (use the function `drop_na` in your pipe). It's ok to have missing mean PM2.5 values in this result.
13. Call up the dimensions of the summary dataset.

#12

```
Air.Summary <-  
  Air.Combo.Subset.Spread %>%  
  group_by(Site.Name, Month, Year) %>%  
  summarise(Mean_Ozone = mean(Ozone),  
            Mean_PM2.5 = mean(PM2.5)) %>%  
  drop_na(Mean_Ozone)
```

```
## 'summarise()' has grouped output by 'Site.Name', 'Month'. You can override  
## using the '.groups' argument.
```

```
#na.omit(Mean_Ozone) Used this to test answer to number 14
```

#13

```
dim(Air.Summary)
```

```
## [1] 182    5
```

#14

14. Why did we use the function `drop_na` rather than `na.omit`?

Answer: `drop_na` will look into a specific column for NA values and remove only rows where there are NAs in that specific column, whereas `na.omit` will remove rows in which ANY column has NA values. I tested this in my pipe by running it with `na.omit` and then `drop_na` and `na.omit` deleted rows where PM2.5 had NA values even though I specified the column I wanted it to look at.