

# Assignment 4: Data Wrangling

Thomas Hancock

## OVERVIEW

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

## Directions

1. Change “Student Name” on line 3 (above) with your name.
2. Work through the steps, **creating code and output** that fulfill each instruction.
3. Be sure to **answer the questions** in this assignment document.
4. When you have completed the assignment, **Knit** the text and code into a single PDF file.
5. After Knitting, submit the completed exercise (PDF file) to the dropbox in Sakai. Add your last name into the file name (e.g., “Salk\_A04\_DataWrangling.Rmd”) prior to submission.

The completed exercise is due on Tuesday, February 4 at 1:00 pm.

## Set up your session

1. Check your working directory, load the `tidyverse` and `lubridate` packages, and upload all four raw data files associated with the EPA Air dataset. See the README file for the EPA air datasets for more information (especially if you have not worked with air quality data previously).
2. Explore the dimensions, column names, and structure of the datasets.

```
#1 Load necessary packages and raw data  
getwd()
```

```
## [1] "C:/Users/thoma/Thomas/2018 Grad School/Duke MEM/ENV 872/Environmental_Data_Analytics_2020/Assignments/Assignment 4/Assignment 4.Rmd"
```

```
library(tidyverse)  
library(lubridate)  
EPAair_03_NC2018 <- read.csv("../Data/Raw/EPAair_03_NC2018_raw.csv")  
EPAair_03_NC2019 <- read.csv("../Data/Raw/EPAair_03_NC2019_raw.csv")  
EPAair_PM25_NC2018 <- read.csv("../Data/Raw/EPAair_PM25_NC2018_raw.csv")  
EPAair_PM25_NC2019 <- read.csv("../Data/Raw/EPAair_PM25_NC2019_raw.csv")
```

```
#2 Explore data  
colnames(EPAair_03_NC2018) # Report column names
```

```
## [1] "Date"  
## [2] "Source"  
## [3] "Site.ID"  
## [4] "POC"  
## [5] "Daily.Max.8.hour.Ozone.Concentration"  
## [6] "UNITS"  
## [7] "DAILY_AQI_VALUE"  
## [8] "Site.Name"  
## [9] "DAILY_OBS_COUNT"
```

```
## [10] "PERCENT_COMPLETE"
## [11] "AQS_PARAMETER_CODE"
## [12] "AQS_PARAMETER_DESC"
## [13] "CBSA_CODE"
## [14] "CBSA_NAME"
## [15] "STATE_CODE"
## [16] "STATE"
## [17] "COUNTY_CODE"
## [18] "COUNTY"
## [19] "SITE_LATITUDE"
## [20] "SITE_LONGITUDE"
```

```
colnames(EPAair_03_NC2019)
```

```
## [1] "Date"
## [2] "Source"
## [3] "Site.ID"
## [4] "POC"
## [5] "Daily.Max.8.hour.Ozone.Concentration"
## [6] "UNITS"
## [7] "DAILY_AQI_VALUE"
## [8] "Site.Name"
## [9] "DAILY_OBS_COUNT"
## [10] "PERCENT_COMPLETE"
## [11] "AQS_PARAMETER_CODE"
## [12] "AQS_PARAMETER_DESC"
## [13] "CBSA_CODE"
## [14] "CBSA_NAME"
## [15] "STATE_CODE"
## [16] "STATE"
## [17] "COUNTY_CODE"
## [18] "COUNTY"
## [19] "SITE_LATITUDE"
## [20] "SITE_LONGITUDE"
```

```
colnames(EPAair_PM25_NC2018)
```

```
## [1] "Date" "Source"
## [3] "Site.ID" "POC"
## [5] "Daily.Mean.PM2.5.Concentration" "UNITS"
## [7] "DAILY_AQI_VALUE" "Site.Name"
## [9] "DAILY_OBS_COUNT" "PERCENT_COMPLETE"
## [11] "AQS_PARAMETER_CODE" "AQS_PARAMETER_DESC"
## [13] "CBSA_CODE" "CBSA_NAME"
## [15] "STATE_CODE" "STATE"
## [17] "COUNTY_CODE" "COUNTY"
## [19] "SITE_LATITUDE" "SITE_LONGITUDE"
```

```
colnames(EPAair_PM25_NC2019)
```

```
## [1] "Date" "Source"
## [3] "Site.ID" "POC"
## [5] "Daily.Mean.PM2.5.Concentration" "UNITS"
## [7] "DAILY_AQI_VALUE" "Site.Name"
## [9] "DAILY_OBS_COUNT" "PERCENT_COMPLETE"
## [11] "AQS_PARAMETER_CODE" "AQS_PARAMETER_DESC"
```

```

## [13] "CBSA_CODE"           "CBSA_NAME"
## [15] "STATE_CODE"          "STATE"
## [17] "COUNTY_CODE"        "COUNTY"
## [19] "SITE_LATITUDE"       "SITE_LONGITUDE"

dim(EPAair_03_NC2018) # Report dimensions

## [1] 9737 20
dim(EPAair_03_NC2019)

## [1] 10592 20
dim(EPAair_PM25_NC2018)

## [1] 8983 20
dim(EPAair_PM25_NC2019)

## [1] 8581 20
str(EPAair_03_NC2018) # Show structure

## 'data.frame': 9737 obs. of 20 variables:
## $ Date : Factor w/ 364 levels "01/01/2018","01/02/2018",...: 60 61 62
## $ Source : Factor w/ 1 level "AQS": 1 1 1 1 1 1 1 1 1 1 ...
## $ Site.ID : int 370030005 370030005 370030005 370030005 370030005 370030005 ...
## $ POC : int 1 1 1 1 1 1 1 1 1 1 ...
## $ Daily.Max.8.hour.Ozone.Concentration: num 0.043 0.046 0.047 0.049 0.047 0.03 0.036 0.044 0.049 0 ...
## $ UNITS : Factor w/ 1 level "ppm": 1 1 1 1 1 1 1 1 1 1 ...
## $ DAILY_AQI_VALUE : int 40 43 44 45 44 28 33 41 45 40 ...
## $ Site.Name : Factor w/ 40 levels "", "Beaufort",...: 35 35 35 35 35 35 35 35 35 35 ...
## $ DAILY_OBS_COUNT : int 17 17 17 17 17 17 17 17 17 17 ...
## $ PERCENT_COMPLETE : num 100 100 100 100 100 100 100 100 100 100 ...
## $ AQS_PARAMETER_CODE : int 44201 44201 44201 44201 44201 44201 44201 44201 44201 44201 ...
## $ AQS_PARAMETER_DESC : Factor w/ 1 level "Ozone": 1 1 1 1 1 1 1 1 1 1 ...
## $ CBSA_CODE : int 25860 25860 25860 25860 25860 25860 25860 25860 25860 25860 ...
## $ CBSA_NAME : Factor w/ 17 levels "", "Asheville, NC",...: 9 9 9 9 9 9 9 9 9 9 ...
## $ STATE_CODE : int 37 37 37 37 37 37 37 37 37 37 ...
## $ STATE : Factor w/ 1 level "North Carolina": 1 1 1 1 1 1 1 1 1 1 ...
## $ COUNTY_CODE : int 3 3 3 3 3 3 3 3 3 3 ...
## $ COUNTY : Factor w/ 32 levels "Alexander", "Avery",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ SITE_LATITUDE : num 35.9 35.9 35.9 35.9 35.9 ...
## $ SITE_LONGITUDE : num -81.2 -81.2 -81.2 -81.2 -81.2 ...

str(EPAair_03_NC2019)

## 'data.frame': 10592 obs. of 20 variables:
## $ Date : Factor w/ 365 levels "01/01/2019","01/02/2019",...: 1 2 3 4 5
## $ Source : Factor w/ 2 levels "AirNow", "AQS": 1 1 1 1 1 1 1 1 1 1 ...
## $ Site.ID : int 370030005 370030005 370030005 370030005 370030005 370030005 ...
## $ POC : int 1 1 1 1 1 1 1 1 1 1 ...
## $ Daily.Max.8.hour.Ozone.Concentration: num 0.029 0.018 0.016 0.022 0.037 0.037 0.029 0.038 0.038 0 ...
## $ UNITS : Factor w/ 1 level "ppm": 1 1 1 1 1 1 1 1 1 1 ...
## $ DAILY_AQI_VALUE : int 27 17 15 20 34 34 27 35 35 28 ...
## $ Site.Name : Factor w/ 38 levels "", "Beaufort",...: 33 33 33 33 33 33 33 33 33 33 ...
## $ DAILY_OBS_COUNT : int 24 24 24 24 24 24 24 24 24 24 ...
## $ PERCENT_COMPLETE : num 100 100 100 100 100 100 100 100 100 100 ...
## $ AQS_PARAMETER_CODE : int 44201 44201 44201 44201 44201 44201 44201 44201 44201 44201 ...

```

```
## $ AQS_PARAMETER_DESC      : Factor w/ 1 level "Ozone": 1 1 1 1 1 1 1 1 1 ...
## $ CBSA_CODE                : int   25860 25860 25860 25860 25860 25860 25860 25860 25860 ...
## $ CBSA_NAME                : Factor w/ 15 levels "", "Asheville, NC",...: 8 8 8 8 8 8 8 8 8 ...
## $ STATE_CODE               : int    37 37 37 37 37 37 37 37 37 ...
## $ STATE                    : Factor w/ 1 level "North Carolina": 1 1 1 1 1 1 1 1 1 ...
## $ COUNTY_CODE              : int     3 3 3 3 3 3 3 3 3 ...
## $ COUNTY                   : Factor w/ 30 levels "Alexander", "Avery",...: 1 1 1 1 1 1 1 1 1 ...
## $ SITE_LATITUDE            : num   35.9 35.9 35.9 35.9 35.9 ...
## $ SITE_LONGITUDE           : num  -81.2 -81.2 -81.2 -81.2 -81.2 ...
```

```
str(EPAair_PM25_NC2018)
```

```
## 'data.frame':      8983 obs. of  20 variables:
## $ Date                : Factor w/ 365 levels "01/01/2018", "01/02/2018",...: 2 5 8 11 14 17 ...
## $ Source               : Factor w/ 1 level "AQS": 1 1 1 1 1 1 1 1 1 ...
## $ Site.ID              : int   370110002 370110002 370110002 370110002 370110002 370110002 ...
## $ POC                  : int    1 1 1 1 1 1 1 1 1 ...
## $ Daily.Mean.PM2.5.Concentration: num   2.9 3.7 5.3 0.8 2.5 4.5 1.8 2.5 4.2 1.7 ...
## $ UNITS                 : Factor w/ 1 level "ug/m3 LC": 1 1 1 1 1 1 1 1 1 ...
## $ DAILY_AQI_VALUE       : int   12 15 22 3 10 19 8 10 18 7 ...
## $ Site.Name             : Factor w/ 25 levels "", "Blackstone",...: 15 15 15 15 15 15 15 15 15 ...
## $ DAILY_OBS_COUNT       : int    1 1 1 1 1 1 1 1 1 ...
## $ PERCENT_COMPLETE      : num   100 100 100 100 100 100 100 100 100 ...
## $ AQS_PARAMETER_CODE    : int   88502 88502 88502 88502 88502 88502 88502 88502 88502 ...
## $ AQS_PARAMETER_DESC    : Factor w/ 2 levels "Acceptable PM2.5 AQI & Speciation Mass",...: 1 ...
## $ CBSA_CODE             : int    NA NA NA NA NA NA NA NA NA ...
## $ CBSA_NAME             : Factor w/ 14 levels "", "Asheville, NC",...: 1 1 1 1 1 1 1 1 1 ...
## $ STATE_CODE            : int    37 37 37 37 37 37 37 37 37 ...
## $ STATE                 : Factor w/ 1 level "North Carolina": 1 1 1 1 1 1 1 1 1 ...
## $ COUNTY_CODE           : int   11 11 11 11 11 11 11 11 11 ...
## $ COUNTY                : Factor w/ 21 levels "Avery", "Buncombe",...: 1 1 1 1 1 1 1 1 1 ...
## $ SITE_LATITUDE         : num   36 36 36 36 36 ...
## $ SITE_LONGITUDE        : num  -81.9 -81.9 -81.9 -81.9 -81.9 ...
```

```
str(EPAair_PM25_NC2019)
```

```
## 'data.frame':      8581 obs. of  20 variables:
## $ Date                : Factor w/ 365 levels "01/01/2019", "01/02/2019",...: 3 6 9 12 15 18 ...
## $ Source               : Factor w/ 2 levels "AirNow", "AQS": 2 2 2 2 2 2 2 2 2 ...
## $ Site.ID              : int   370110002 370110002 370110002 370110002 370110002 370110002 ...
## $ POC                  : int    1 1 1 1 1 1 1 1 1 ...
## $ Daily.Mean.PM2.5.Concentration: num   1.6 1 1.3 6.3 2.6 1.2 1.5 1.5 3.7 1.6 ...
## $ UNITS                 : Factor w/ 1 level "ug/m3 LC": 1 1 1 1 1 1 1 1 1 ...
## $ DAILY_AQI_VALUE       : int    7 4 5 26 11 5 6 15 7 ...
## $ Site.Name             : Factor w/ 25 levels "", "Board Of Ed. Bldg.",...: 14 14 14 14 14 14 14 14 14 ...
## $ DAILY_OBS_COUNT       : int    1 1 1 1 1 1 1 1 1 ...
## $ PERCENT_COMPLETE      : num   100 100 100 100 100 100 100 100 100 ...
## $ AQS_PARAMETER_CODE    : int   88502 88502 88502 88502 88502 88502 88502 88502 88502 ...
## $ AQS_PARAMETER_DESC    : Factor w/ 2 levels "Acceptable PM2.5 AQI & Speciation Mass",...: 1 ...
## $ CBSA_CODE             : int    NA NA NA NA NA NA NA NA NA ...
## $ CBSA_NAME             : Factor w/ 14 levels "", "Asheville, NC",...: 1 1 1 1 1 1 1 1 1 ...
## $ STATE_CODE            : int    37 37 37 37 37 37 37 37 37 ...
## $ STATE                 : Factor w/ 1 level "North Carolina": 1 1 1 1 1 1 1 1 1 ...
## $ COUNTY_CODE           : int   11 11 11 11 11 11 11 11 11 ...
## $ COUNTY                : Factor w/ 21 levels "Avery", "Buncombe",...: 1 1 1 1 1 1 1 1 1 ...
```

```
## $ SITE_LATITUDE          : num  36 36 36 36 36 ...
## $ SITE_LONGITUDE         : num -81.9 -81.9 -81.9 -81.9 -81.9 ...
```

## Wrangle individual datasets to create processed files.

3. Change date to date
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 Change the Date column to read as a date
EPAair_03_NC2018$Date <- as.Date(EPAair_03_NC2018$Date, format = "%m/%d/%Y")
EPAair_03_NC2019$Date <- as.Date(EPAair_03_NC2019$Date, format = "%m/%d/%Y")
EPAair_PM25_NC2018$Date <- as.Date(EPAair_PM25_NC2018$Date, format = "%m/%d/%Y")
EPAair_PM25_NC2019$Date <- as.Date(EPAair_PM25_NC2019$Date, format = "%m/%d/%Y")
class(EPAair_03_NC2018$Date) # make sure date format worked

## [1] "Date"

#4 Select specific columns to include
EPAair_03_NC2018 <- select(EPAair_03_NC2018, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY)
EPAair_03_NC2019 <- select(EPAair_03_NC2019, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY)
EPAair_PM25_NC2018 <- select(EPAair_PM25_NC2018, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY)
EPAair_PM25_NC2019 <- select(EPAair_PM25_NC2019, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY)

#5 Set parameter description to "PM2.5" for the PM2.5 data sets (all entries set to this value)
EPAair_PM25_NC2018$AQS_PARAMETER_DESC <- as.factor("PM2.5")
EPAair_PM25_NC2019$AQS_PARAMETER_DESC <- as.factor("PM2.5")

#6 Save each dataframe
write.csv(EPAair_03_NC2018, row.names = FALSE,
          file = "../Data/Processed/EPAair_03_NC2018_processed.csv")
write.csv(EPAair_03_NC2019, row.names = FALSE,
          file = "../Data/Processed/EPAair_03_NC2019_processed.csv")
write.csv(EPAair_PM25_NC2018, row.names = FALSE,
          file = "../Data/Processed/EPAair_PM25_NC2018_processed.csv")
write.csv(EPAair_PM25_NC2019, 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 all 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)
  - Some sites have multiple measurements per day. Use the split-apply-combine strategy to generate daily means: group by date, site, 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.
- 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.
  - Call up the dimensions of your new tidy dataset.
  - Save your processed dataset with the following file name: "EPAair\_O3\_PM25\_NC1718\_Processed.csv"

```
#7 Rbind the 4 dataframes into a single long data frame
EPAair_combined <- rbind(EPAair_O3_NC2018, EPAair_O3_NC2019, EPAair_PM25_NC2018, EPAair_PM25_NC2019)

#8 Process data
common <- intersect(EPAair_PM25_NC2019$Site.Name, intersect(EPAair_PM25_NC2018$Site.Name, intersect(EPAair_O3_NC2018$Site.Name, EPAair_O3_NC2019$Site.Name)))

EPAair_combined_daily <-
  EPAair_combined %>%
  filter(Site.Name %in% common & Site.Name != "") %>% # Filter to only include common sites (without any)
  group_by(Date, Site.Name, AQS_PARAMETER_DESC, COUNTY) %>% # Group entries
  summarise(mean_AQI = mean(DAILY_AQI_VALUE), # Create averages of AQI value and lat/long
            mean_lat = mean(SITE_LATITUDE),
            mean_lon = mean(SITE_LONGITUDE)) %>%
  mutate(Month = month(Date)) %>% # Create a column with just the month
  mutate(Year = year(Date)) # Create a column with just the year

#9 Spread the dataframe so Ozone and PM2.5 each have their own column
EPAair_combined_Processed <- spread(EPAair_combined_daily, AQS_PARAMETER_DESC, mean_AQI)

#10 Find the dimensions of the processed dataframe
dim(EPAair_combined_Processed)

## [1] 8976    9

#11 Save the file (note: I changed the name to ...NC1819... instead of NC1718 since we have years 2018 and 2019)
write.csv(EPAair_combined_Processed, file = "../Data/Processed/EPAair_O3_PM25_NC1819_Processed.csv")
```

## Generate summary tables

- 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 a month and year are not available (use the function `drop_na` in your pipe).
- Call up the dimensions of the summary dataset.

```
#12a Create summary
EPAair_combined_summary <-
  EPAair_combined_Processed %>%
  group_by(Site.Name, Month, Year) %>% # Group by site, month, and year
  summarise(MeanOzone = mean(Ozone), # Create summaries of ozone and pm2.5 levels for groups
            MeanPM25 = mean(PM2.5))

#12b Remove entries with NA as month or year
EPAair_combined_summary <-
  EPAair_combined_summary %>%
  drop_na(Month, Year) # Drop rows with NA in these columns

#13 Call up dimensions of summary dataset
dim(EPAair_combined_summary)
```

```
## [1] 308 5
```

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

Answer: The `'drop_na'` function allows us to only drop rows with NA in specified columns (in this case, Year and Month). If we used the `'na.omit'` function, it would drop all entries that have an NA in any of the columns, including the Ozone and PM2.5 columns (which is most of them).