

# Assignment 4: Data Wrangling

Changxin Yu

## OVERVIEW

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

## Directions

1. Rename this file `<FirstLast>_A03_DataExploration.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.

The completed exercise is due on Friday, Oct7th @ 5:00pm.

## 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, being sure to set string columns to be read in a factors. See the README file for the EPA air datasets for more information (especially if you have not worked with air quality data previously).

```
getwd()

## [1] "E:/things/Duke University/study/2022 Fall/ENVIRON 872/EDA-Fall2022/Assignments"

setwd("E:/things/Duke University/study/2022 Fall/ENVIRON 872/EDA-Fall2022")

library(tidyverse)
library(lubridate)

EPA_03_18 <- read.csv("./Data/Raw/EPAair_03_NC2018_raw.csv", stringsAsFactors = TRUE)
EPA_03_19 <- read.csv("./Data/Raw/EPAair_03_NC2019_raw.csv", stringsAsFactors = TRUE)
EPA_PM25_18 <- read.csv("./Data/Raw/EPAair_PM25_NC2018_raw.csv",
                        stringsAsFactors = TRUE)
EPA_PM25_19 <- read.csv("./Data/Raw/EPAair_PM25_NC2019_raw.csv",
                        stringsAsFactors = TRUE)
```

2. Explore the dimensions, column names, and structure of the datasets.

```

EPAair <- list(EPA_03_18, EPA_03_19, EPA_PM25_18, EPA_PM25_19)
for(df in EPAair){
  print(dim(df))
  print(colnames(df))
  print(str(df))
}

```

```

## [1] 9737    20
## [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"
## '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 37003
## $ 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
## $ 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 4
## $ 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 2
## $ CBSA_NAME
##      : Factor w/ 17 levels "", "Asheville, NC",...: 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
## $ 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 ...
## NULL
## [1] 10592    20
## [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"
## '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 0.038 ...
## $ 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 1 ...
## $ CBSA_CODE : int 25860 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 8 ...
## $ 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/ 30 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 ...
## NULL
## [1] 8983 20
## [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"
## '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 1 ...

```

```

## $ Site.ID : int 370110002 370110002 370110002 370110002 370110002 370110002
## $ POC : int 1 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 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 15
## $ DAILY_OBS_COUNT : int 1 1 1 1 1 1 1 1 1 1 ...
## $ PERCENT_COMPLETE : num 100 100 100 100 100 100 100 100 100 100 ...
## $ AQS_PARAMETER_CODE : int 88502 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 NA ...
## $ CBSA_NAME : Factor w/ 14 levels "", "Asheville, NC",...: 1 1 1 1 1 1 1 1 1 1 ..
## $ 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 11 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 1 ..
## $ SITE_LATITUDE : num 36 36 36 36 36 ...
## $ SITE_LONGITUDE : num -81.9 -81.9 -81.9 -81.9 -81.9 ...
## NULL
## [1] 8581 20
## [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"
## '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 2 ...
## $ Site.ID : int 370110002 370110002 370110002 370110002 370110002 370110002
## $ POC : int 1 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 1 ...
## $ DAILY_AQI_VALUE : int 7 4 5 26 11 5 6 6 15 7 ...
## $ Site.Name : Factor w/ 25 levels "", "Board Of Ed. Bldg.",...: 14 14 14 14 14 14
## $ DAILY_OBS_COUNT : int 1 1 1 1 1 1 1 1 1 1 ...
## $ PERCENT_COMPLETE : num 100 100 100 100 100 100 100 100 100 100 ...
## $ AQS_PARAMETER_CODE : int 88502 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 NA ...
## $ CBSA_NAME : Factor w/ 14 levels "", "Asheville, NC",...: 1 1 1 1 1 1 1 1 1 1 ..
## $ 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 11 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 1 ..
## $ SITE_LATITUDE : num 36 36 36 36 36 ...
## $ SITE_LONGITUDE : num -81.9 -81.9 -81.9 -81.9 -81.9 ...
## NULL

```

## 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
EPAair1=list()
for(df in EPAair){
  df$Date <- as.Date(df$Date, format = "%m/%d/%Y")
  df = list(df)
  EPAair1 <- append(EPAair1, df)
}

#4
EPAair=list()
for(df in EPAair1){
  df <- select(df, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,
              COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
  df = list(df)
  EPAair <- append(EPAair, df)
}

#5
EPA_03_18_processed <- as.data.frame(EPAair[1])
EPA_03_19_processed <- as.data.frame(EPAair[2])
EPA_PM25_18_processed <- as.data.frame(EPAair[3])
EPA_PM25_19_processed <- as.data.frame(EPAair[4])

EPA_PM25_18_processed <- mutate(EPA_PM25_18_processed, AQS_PARAMETER_DESC = "PM2.5")
EPA_PM25_19_processed <- mutate(EPA_PM25_19_processed, AQS_PARAMETER_DESC = "PM2.5")

#6
setwd("E:/things/Duke University/study/2022 Fall/ENVIRON 872/EDA-Fall2022")
# There is always some error of my working directory path.
# It has the default of assignment folder and cannot be changed in the markdown file.
write.csv(EPA_03_18_processed, row.names = FALSE,
          file = "./Data/Processed/EPAair_03_NC2018_processed.csv")
write.csv(EPA_03_19_processed, row.names = FALSE,
          file = "./Data/Processed/EPAair_03_NC2019_processed.csv")
write.csv(EPA_PM25_18_processed, row.names = FALSE,
          file = "./Data/Processed/EPAair_PM25_NC2018_processed.csv")
write.csv(EPA_PM25_19_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 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.
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
EPA_1819 <- rbind(EPA_O3_18_processed, EPA_O3_19_processed,
                  EPA_PM25_18_processed, EPA_PM25_19_processed)

#8
EPA_1819 <- EPA_1819 %>%
  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(meanAQI = mean(DAILY_AQI_VALUE),
            meanLat = mean(SITE_LATITUDE),
            meanLon = mean(SITE_LONGITUDE)) %>%
  mutate(Month = month(Date)) %>%
  mutate(Year = year(Date))

#9
EPA_1819 <- pivot_wider(EPA_1819, names_from = AQS_PARAMETER_DESC,
                       values_from = meanAQI)

#10
dim(EPA_1819)

## [1] 8976    9

#11
setwd("E:/things/Duke University/study/2022 Fall/ENVIRON 872/EDA-Fall2022")
write.csv(EPA_1819, row.names = FALSE,
          file = "./Data/Processed/EPAair_O3_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 a month and year are not available (use the function `drop_na` in your pipe).

13. Call up the dimensions of the summary dataset.

```
#12
EPA_1819_summary <- EPA_1819 %>%
  group_by(Site.Name, Month, Year) %>%
  summarise(meanPM25 = mean(PM2.5),
            meanO3 = mean(Ozone)) %>%
  drop_na(Month | Year)
```

```
#13
dim(EPA_1819_summary)
```

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

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

Answer: Using `na.omit` will remove all rows in a data frame that has one or more NA values. While `drop_na` will remove rows by finding whether there is an NA or not in specified column. Therefore, in this case, `drop_na` is better to focus on Month and Year columns.