

Assignment 4: Data Wrangling

Azura Liu

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., “Fay_A04_DataWrangling.Rmd”) prior to submission.

The completed exercise is due on Monday, Feb 7 @ 7:00pm.

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. 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 set up
getwd()

## [1] "C:/Users/Idae/Desktop/ENV872/Environmental_Data_Analytics_2022/Assignments"

library("tidyverse")
library("lubridate")
PM25.18<-read.csv("../Data/Raw/EPAair_PM25_NC2018_raw.csv")
PM25.19<-read.csv("../Data/Raw/EPAair_PM25_NC2019_raw.csv")
O3.19<-read.csv("../Data/Raw/EPAair_O3_NC2019_raw.csv")
O3.18<-read.csv("../Data/Raw/EPAair_O3_NC2018_raw.csv")

#2 explore datasets
dim(PM25.18)

## [1] 8983    20

colnames(PM25.18)

## [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] "AQ5_PARAMETER_CODE" "AQ5_PARAMETER_DESC"
```

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

```
str(PM25.18)
```

```
## 'data.frame': 8983 obs. of 20 variables:
## $ Date : chr "01/02/2018" "01/05/2018" "01/08/2018" "01/11/2018" ...
## $ Source : chr "AQS" "AQS" "AQS" "AQS" ...
## $ 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 : chr "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 ...
## $ Site.Name : chr "Linville Falls" "Linville Falls" "Linville Falls" "Linville Falls" ...
## $ 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 : chr "Acceptable PM2.5 AQI & Speciation Mass" "Acceptable PM2.5 AQI & Speciation Mass" ...
## $ CBSA_CODE : int NA NA NA NA NA NA NA NA NA NA ...
## $ CBSA_NAME : chr "" "" "" "" "" "" "" "" "" "" ...
## $ STATE_CODE : int 37 37 37 37 37 37 37 37 37 37 ...
## $ STATE : chr "North Carolina" "North Carolina" "North Carolina" "North Carolina" ...
## $ COUNTY_CODE : int 11 11 11 11 11 11 11 11 11 11 ...
## $ COUNTY : chr "Avery" "Avery" "Avery" "Avery" ...
## $ SITE_LATITUDE : num 36 36 36 36 36 ...
## $ SITE_LONGITUDE : num -81.9 -81.9 -81.9 -81.9 -81.9 ...
```

```
dim(PM25.19)
```

```
## [1] 8581 20
```

```
colnames(PM25.19)
```

```
## [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"
```

```
str(PM25.19)
```

```
## 'data.frame': 8581 obs. of 20 variables:
## $ Date : chr "01/03/2019" "01/06/2019" "01/09/2019" "01/12/2019" ...
## $ Source : chr "AQS" "AQS" "AQS" "AQS" ...
## $ 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 : chr "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 ...
## $ Site.Name : chr "Linville Falls" "Linville Falls" "Linville Falls" "Linville Falls" ...
```

```
## $ 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   : chr   "Acceptable PM2.5 AQI & Speciation Mass" "Acceptable PM2.5 AQI & Speciation Mass"
## $ CBSA_CODE            : int  NA NA NA NA NA NA NA NA NA NA ...
## $ CBSA_NAME            : chr   "" "" "" "" "" "" "" "" "" "" ...
## $ STATE_CODE           : int  37 37 37 37 37 37 37 37 37 37 ...
## $ STATE                : chr   "North Carolina" "North Carolina" "North Carolina" "North Carolina"
## $ COUNTY_CODE          : int  11 11 11 11 11 11 11 11 11 11 ...
## $ COUNTY               : chr   "Avery" "Avery" "Avery" "Avery" ...
## $ SITE_LATITUDE        : num  36 36 36 36 36 ...
## $ SITE_LONGITUDE       : num  -81.9 -81.9 -81.9 -81.9 -81.9 ...
```

```
dim(O3.18)
```

```
## [1] 9737 20
```

```
colnames(O3.18)
```

```
## [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"
```

```
str(O3.18)
```

```
## 'data.frame': 9737 obs. of 20 variables:
## $ Date      : chr  "03/01/2018" "03/02/2018" "03/03/2018" "03/04/2018" ...
## $ Source    : chr  "AQS" "AQS" "AQS" "AQS" ...
## $ Site.ID   : int  370030005 370030005 370030005 370030005 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 0.049
## $ UNITS     : chr  "ppm" "ppm" "ppm" "ppm" ...
## $ DAILY_AQI_VALUE : int  40 43 44 45 44 28 33 41 45 40 ...
## $ Site.Name  : chr  "Taylorsville Liledoun" "Taylorsville Liledoun" "Taylorsville Liledoun"
## $ 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 : chr  "Ozone" "Ozone" "Ozone" "Ozone" ...
## $ CBSA_CODE   : int  25860 25860 25860 25860 25860 25860 25860 25860 25860 25860
```

```
## $ CBSA_NAME           : chr  "Hickory-Lenoir-Morganton, NC" "Hickory-Lenoir-Morganton, NC" ...
## $ STATE_CODE          : int   37 37 37 37 37 37 37 37 37 37 ...
## $ STATE               : chr   "North Carolina" "North Carolina" "North Carolina" "North Carolina" ...
## $ COUNTY_CODE         : int   3 3 3 3 3 3 3 3 3 3 ...
## $ COUNTY              : chr   "Alexander" "Alexander" "Alexander" "Alexander" ...
## $ SITE_LATITUDE       : num   35.9 35.9 35.9 35.9 35.9 35.9 ...
## $ SITE_LONGITUDE      : num   -81.2 -81.2 -81.2 -81.2 -81.2 -81.2 ...
```

```
dim(O3.19)
```

```
## [1] 10592    20
```

```
colnames(O3.19)
```

```
## [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] "AQ5_PARAMETER_CODE"
## [12] "AQ5_PARAMETER_DESC"
## [13] "CBSA_CODE"
## [14] "CBSA_NAME"
## [15] "STATE_CODE"
## [16] "STATE"
## [17] "COUNTY_CODE"
## [18] "COUNTY"
## [19] "SITE_LATITUDE"
## [20] "SITE_LONGITUDE"
```

```
str(O3.19)
```

```
## 'data.frame':    10592 obs. of  20 variables:
## $ Date              : chr   "01/01/2019" "01/02/2019" "01/03/2019" "01/04/2019" ...
## $ Source            : chr   "AirNow" "AirNow" "AirNow" "AirNow" ...
## $ 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 ...
## $ UNITS             : chr    "ppm" "ppm" "ppm" "ppm" ...
## $ DAILY_AQI_VALUE   : int    27 17 15 20 34 34 27 35 35 28 ...
## $ Site.Name         : chr    "Taylorsville Liledoun" "Taylorsville Liledoun" "Taylorsville Liledoun" ...
## $ 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 ...
## $ AQ5_PARAMETER_CODE: int    44201 44201 44201 44201 44201 44201 44201 44201 44201 44201 ...
## $ AQ5_PARAMETER_DESC: chr    "Ozone" "Ozone" "Ozone" "Ozone" ...
## $ CBSA_CODE         : int    25860 25860 25860 25860 25860 25860 25860 25860 25860 25860 ...
## $ CBSA_NAME         : chr    "Hickory-Lenoir-Morganton, NC" "Hickory-Lenoir-Morganton, NC" ...
## $ STATE_CODE        : int     37 37 37 37 37 37 37 37 37 37 ...
## $ STATE             : chr    "North Carolina" "North Carolina" "North Carolina" "North Carolina" ...
## $ COUNTY_CODE       : int     3 3 3 3 3 3 3 3 3 3 ...
## $ COUNTY            : chr    "Alexander" "Alexander" "Alexander" "Alexander" ...
```

```
## $ 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 ...
```

Wrangle individual datasets to create processed files.

3. Change date to a date object
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 set date

```
class(PM25.18$Date)
```

```
## [1] "character"
```

```
PM25.18$Date<-as.Date(PM25.18$Date, format = "%m/%d/%Y")
```

```
class(PM25.18$Date)
```

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

```
class(PM25.19$Date)
```

```
## [1] "character"
```

```
PM25.19$Date<-as.Date(PM25.19$Date, format = "%m/%d/%Y")
```

```
class(PM25.19$Date)
```

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

```
class(O3.18$Date)
```

```
## [1] "character"
```

```
O3.18$Date<-as.Date(O3.18$Date, format = "%m/%d/%Y")
```

```
class(O3.18$Date)
```

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

```
class(O3.19$Date)
```

```
## [1] "character"
```

```
O3.19$Date<-as.Date(O3.19$Date, format = "%m/%d/%Y")
```

```
class(O3.19$Date)
```

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

#4 subsetting

```
PM25.18.sub<-select(PM25.18, Date, DAILY_AQI_VALUE, Site.Name,
                    AQS_PARAMETER_DESC, COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

```
PM25.19.sub<-select(PM25.19, Date, DAILY_AQI_VALUE, Site.Name,
                    AQS_PARAMETER_DESC, COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

```
O3.18.sub<-select(O3.18, Date, DAILY_AQI_VALUE, Site.Name,
                  AQS_PARAMETER_DESC, COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

```
O3.19.sub<-select(O3.19, Date, DAILY_AQI_VALUE, Site.Name,
                  AQS_PARAMETER_DESC, COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
```

#5 fill a column

```

PM25.18.sub$AQS_PARAMETER_DESC<-"PM2.5")
PM25.19.sub$AQS_PARAMETER_DESC<-"PM2.5")

#6 save work
write.csv(PM25.18.sub, row.names = FALSE, file = "../Data/Processed/EPAair_PM25_NC2018_processed.csv")
write.csv(PM25.19.sub, row.names = FALSE, file = "../Data/Processed/EPAair_PM25_NC2019_processed.csv")
write.csv(O3.18.sub, row.names = FALSE, file = "../Data/Processed/EPAair_O3_NC2019_processed.csv")
write.csv(O3.19.sub, row.names = FALSE, file = "../Data/Processed/EPAair_O3_NC2018_processed.csv")

```

Combine datasets

- Combine the four datasets with `rbind`. Make sure your column names are identical prior to running this code.
- Wrangle your new dataset with a pipe function (`%>%`) so that it fills the following conditions:
 - Filter records to include just the 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 `intersect` function can figure out common factor levels if we didn’t give you this list…)
 - 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_NC2122_Processed.csv”

```

#7 combine datasets
EPA_Air<-rbind(PM25.18.sub,PM25.19.sub,O3.18.sub,O3.19.sub)
dim (EPA_Air) #this is correct

```

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

```
summary(EPA_Air)
```

```

##      Date      DAILY_AQI_VALUE  Site.Name      AQS_PARAMETER_DESC
##  Min.   :2018-01-01  Min.    : 0.00  Length:37893  Length:37893
## 1st Qu.:2018-06-27  1st Qu.: 27.00  Class :character  Class :character
## Median :2019-01-06  Median : 36.00  Mode  :character  Mode  :character
## Mean   :2018-12-26  Mean    : 36.27
## 3rd Qu.:2019-06-23  3rd Qu.: 45.00
## Max.   :2019-12-31  Max.    :136.00
##      COUNTY      SITE_LATITUDE  SITE_LONGITUDE
## Length:37893    Min.    :34.36  Min.    :-83.80
## Class :character 1st Qu.:35.26  1st Qu.: -81.37
## Mode  :character Median :35.64  Median : -80.23
##                  Mean    :35.62  Mean    :-80.21
##                  3rd Qu.:35.99  3rd Qu.: -78.77
##                  Max.    :36.51  Max.    :-76.21

```

```

#8 piping
EPA_Air_Piped<-EPA_Air %>%

```

```

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) %>%
summarize(DAILY_AQI_VALUE=mean(DAILY_AQI_VALUE),
           SITE_LONGITUDE = mean(SITE_LONGITUDE),
           SITE_LATITUDE = mean(SITE_LATITUDE))>%
mutate(Month = month (Date),
       Year = year (Date))

```

```

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

```

```

## [1] 14752      9
summary(EPA_Air_Piped)

```

```

##      Date      Site.Name      AQS_PARAMETER_DESC      COUNTY
## Min.   :2018-01-01 Length:14752 Length:14752 Length:14752
## 1st Qu.:2018-07-01 Class :character Class :character Class :character
## Median :2019-01-08 Mode  :character Mode  :character Mode  :character
## Mean   :2018-12-30
## 3rd Qu.:2019-06-28
## Max.   :2019-12-31
## DAILY_AQI_VALUE SITE_LONGITUDE SITE_LATITUDE      Month
## Min.   : 0.00 Min.   :-83.44 Min.   :34.36 Min.   : 1.000
## 1st Qu.: 25.00 1st Qu.: -80.79 1st Qu.:35.43 1st Qu.: 4.000
## Median : 35.00 Median : -79.80 Median :35.86 Median : 6.000
## Mean   : 35.19 Mean   : -79.67 Mean   :35.68 Mean   : 6.402
## 3rd Qu.: 44.00 3rd Qu.: -78.46 3rd Qu.:36.03 3rd Qu.: 9.000
## Max.   :129.00 Max.   : -77.36 Max.   :36.11 Max.   :12.000
##      Year
## Min.   :2018
## 1st Qu.:2018
## Median :2019
## Mean   :2019
## 3rd Qu.:2019
## Max.   :2019

```

```

#9 Spread AQI
EPA_Air_Wider <- pivot_wider(EPA_Air_Piped, names_from = AQS_PARAMETER_DESC, values_from = DAILY_AQI_VALUE)

#10 check dimension
dim(EPA_Air_Wider)

```

```

## [1] 8976      9
#11 save work
write.csv(EPA_Air_Wider, row.names = FALSE, file = "../Data/Processed/EPAair_03_PM25_NC2122_Processed.csv")

```

Generate summary tables

12a. Use the split-apply-combine strategy to generate a summary data frame from your results from Step 9 above. Data should be grouped by site, month, and year. Generate the mean AQI values for ozone and PM2.5 for each group.

12b. BONUS: Add a piped statement to 12a that removes rows where both mean ozone and mean PM2.5 have missing values.

13. Call up the dimensions of the summary dataset.

```
#12(a,b) summary table
```

```
EPA_Air_Summary <-  
  EPA_Air_Wider %>%  
  group_by(Site.Name, Month, Year) %>%  
  summarise(mean.PM25 = mean(PM2.5),  
            mean.O3 = mean(Ozone)) %>%  
  filter(is.na(mean.PM25) == F & is.na(mean.O3) == F  
         | is.na(mean.PM25) == T & is.na(mean.O3) == F  
         | is.na(mean.PM25) == F & is.na(mean.O3) == T )
```

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

```
#cannot figure out how to drop only when both are null with drop_na  
#nor removing the filtered results within pipe...I'm sure there is an easier way
```

```
#13 check dim  
dim(EPA_Air_Summary)
```

```
## [1] 292 5
```

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

Answer: “drop_na” is from the package “dplyr” and allows more data manipulation. We can specify the columns we want to apply the function to. “na.omit” removes the entire row if missing data is presented at all.