

# Assignment 4: Data Wrangling

Joshua Frear

## 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 Tuesday, Feb 16 @ 11:59pm.

## 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
getwd()

## [1] "D:/Data_872/Environmental_Data_Analytics_2021/Assignments"

library(tidyverse)
library(lubridate)

Ozone_NC2018 <- read.csv("../Data/Raw/EPAair_03_NC2018_raw.csv",
                        stringsAsFactors = TRUE)
Ozone_NC2019 <- read.csv("../Data/Raw/EPAair_03_NC2019_raw.csv",
                        stringsAsFactors = TRUE)
PM25_NC2018 <- read.csv("../Data/Raw/EPAair_PM25_NC2018_raw.csv",
                        stringsAsFactors = TRUE)
PM25_NC2019 <- read.csv("../Data/Raw/EPAair_PM25_NC2019_raw.csv",
                        stringsAsFactors = TRUE)

#2
#create a function to do the three exploratory actions for each dataset
explore <- function(x) {
  print(dim(x))
  print(colnames(x))
  print(str(x))
}
```

```
explore(Ozone_NC2018)
```

```
## [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 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 ...
## NULL
```

```
explore(Ozone_NC2019)
```

```
## [1] 10592    20
## [1] "Date"
## [2] "Source"
## [3] "Site.ID"
## [4] "POC"
## [5] "Daily.Max.8.hour.Ozone.Concentration"
```



```
## $ 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
```

```
explore(PM25_NC2019)
```

```
## [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
Ozone_NC2018$Date <- as.Date(Ozone_NC2018$Date, format = "%m/%d/%Y")
Ozone_NC2019$Date <- as.Date(Ozone_NC2019$Date, format = "%m/%d/%Y")
PM25_NC2018$Date <- as.Date(PM25_NC2018$Date, format = "%m/%d/%Y")
PM25_NC2019$Date <- as.Date(PM25_NC2019$Date, format = "%m/%d/%Y")

#4
Ozone_NC2018.processed <-
  Ozone_NC2018 %>%
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
Ozone_NC2019.processed <-
  Ozone_NC2019 %>%
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
PM25_NC2018.processed <-
  PM25_NC2018 %>%
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
PM25_NC2019.processed <-
  PM25_NC2019 %>%
  select(Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,
         COUNTY, SITE_LATITUDE, SITE_LONGITUDE)

#5
PM25_NC2018.processed$AQS_PARAMETER_DESC = "PM2.5"
PM25_NC2019.processed$AQS_PARAMETER_DESC = "PM2.5"

#6
write.csv(Ozone_NC2018.processed,
         file = "../Data/Processed/EPAair_O3_NC2018_processed.csv")
write.csv(Ozone_NC2019.processed,
         file = "../Data/Processed/EPAair_O3_NC2019_processed.csv")
write.csv(PM25_NC2018.processed,
         file = "../Data/Processed/EPAair_PM25_NC2018_processed.csv")
write.csv(PM25_NC2019.processed,
         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\_NC1718\_Processed.csv”

```
#7
EPAair_combined <-
  rbind(Ozone_NC2018.processed, Ozone_NC2019.processed,
        PM25_NC2018.processed, PM25_NC2019.processed)

summary(EPAair_combined)
```

##	Date	DAILY_AQI_VALUE	Site.Name
##	Min. :2018-01-01	Min. : 0.00	Millbrook School : 2169
##	1st Qu.:2018-06-27	1st Qu.: 27.00	Garinger High School: 1818
##	Median :2019-01-06	Median : 36.00	Hattie Avenue : 1432
##	Mean :2018-12-26	Mean : 36.27	Durham Armory : 1405
##	3rd Qu.:2019-06-23	3rd Qu.: 45.00	Pitt Agri. Center : 1303
##	Max. :2019-12-31	Max. :136.00	Clemmons Middle : 1261
##			(Other) :28505
##	AQS_PARAMETER_DESC	COUNTY	SITE_LATITUDE SITE_LONGITUDE
##	Ozone:20329	Mecklenburg: 3903	Min. :34.36 Min. : -83.80
##	PM2.5:17564	Forsyth : 3175	1st Qu.:35.26 1st Qu.: -81.37
##		Wake : 2846	Median :35.64 Median : -80.23
##		Cumberland : 1795	Mean :35.62 Mean : -80.21
##		Haywood : 1672	3rd Qu.:35.99 3rd Qu.: -78.77
##		Swain : 1628	Max. :36.51 Max. : -76.21
##		(Other) :22874	

```
str(EPAair_combined)
```

```
## 'data.frame': 37893 obs. of 7 variables:
## $ Date : Date, format: "2018-03-01" "2018-03-02" ...
## $ DAILY_AQI_VALUE : int 40 43 44 45 44 28 33 41 45 40 ...
## $ Site.Name : Factor w/ 51 levels "", "Beaufort",...: 35 35 35 35 35 35 35 35 35 35 ...
## $ AQS_PARAMETER_DESC: Factor w/ 2 levels "Ozone", "PM2.5": 1 1 1 1 1 1 1 1 1 1 ...
## $ COUNTY : Factor w/ 37 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 ...
```

```
#8
common_sites <- intersect(intersect(PM25_NC2018.processed$Site.Name,
                                     PM25_NC2019.processed$Site.Name),
                           intersect(Ozone_NC2018.processed$Site.Name,
                                     Ozone_NC2019.processed$Site.Name))

EPAair_tidy <-
```

```

EPAair_combined %>%
  filter(Site.Name %in% common_sites) %>%
  filter(Site.Name != "") %>%
  group_by(Date, Site.Name, AQS_PARAMETER_DESC, COUNTY) %>%
  summarise(meanAQI = mean(DAILY_AQI_VALUE),
            meanlatitude = mean(SITE_LATITUDE),
            meanlongitude = 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

```

#9
EPAair_tidy <-
  EPAair_tidy %>%
  pivot_wider(names_from = AQS_PARAMETER_DESC, values_from = meanAQI)

```

```

#10
dim(EPAair_tidy)

```

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

```

#11
write.csv(EPAair_tidy, "../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 a month and year are not available (use the function `drop_na` in your pipe).

13. Call up the dimensions of the summary dataset.

```

#12a
EPAair_summary <-
  EPAair_tidy %>%
  group_by(Site.Name, Month, Year) %>%
  summarise(mean_PM25_AQI = mean(PM2.5),
            mean_03_AQI = mean(Ozone))

```

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

```

#12b
EPAair_summary <-
  EPAair_summary %>%
  drop_na(Month, Year)

```

```

#13
dim(EPAair_summary)

```

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

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

Answer: `na.omit`, by default, removes any row with an NA in any column, and you cannot easily force it to only look at specific columns, like `drop_na` does.