

# Assignment 4: Data Wrangling

Jessica Citrola

## 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
library(tidyverse)
library(lubridate)

EPA.O3.2018 <- read.csv("../Data/Raw/EPAair_O3_NC2018_raw.csv", stringsAsFactors = TRUE)
EPA.O3.2019 <- read.csv("../Data/Raw/EPAair_O3_NC2019_raw.csv", stringsAsFactors = TRUE)
EPA.PM.2018 <- read.csv("../Data/Raw/EPAair_PM25_NC2018_raw.csv", stringsAsFactors = TRUE)
EPA.PM.2019 <- read.csv("../Data/Raw/EPAair_PM25_NC2019_raw.csv", stringsAsFactors = TRUE)

#2

dim(EPA.O3.2018)

## [1] 9737 20
```

```
colnames(EPA.03.2018)
```

```
## [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"
```

```
head(EPA.03.2018)
```

```
##      Date Source   Site.ID POC Daily.Max.8.hour.Ozone.Concentration UNITS
## 1 03/01/2018   AQS 370030005   1                        0.043   ppm
## 2 03/02/2018   AQS 370030005   1                        0.046   ppm
## 3 03/03/2018   AQS 370030005   1                        0.047   ppm
## 4 03/04/2018   AQS 370030005   1                        0.049   ppm
## 5 03/05/2018   AQS 370030005   1                        0.047   ppm
## 6 03/06/2018   AQS 370030005   1                        0.030   ppm
##   DAILY_AQI_VALUE      Site.Name DAILY_OBS_COUNT PERCENT_COMPLETE
## 1              40 Taylorsville Liledoun           17           100
## 2              43 Taylorsville Liledoun           17           100
## 3              44 Taylorsville Liledoun           17           100
## 4              45 Taylorsville Liledoun           17           100
## 5              44 Taylorsville Liledoun           17           100
## 6              28 Taylorsville Liledoun           17           100
##   AQS_PARAMETER_CODE AQS_PARAMETER_DESC CBSA_CODE      CBSA_NAME
## 1              44201              Ozone   25860 Hickory-Lenoir-Morganton, NC
## 2              44201              Ozone   25860 Hickory-Lenoir-Morganton, NC
## 3              44201              Ozone   25860 Hickory-Lenoir-Morganton, NC
## 4              44201              Ozone   25860 Hickory-Lenoir-Morganton, NC
## 5              44201              Ozone   25860 Hickory-Lenoir-Morganton, NC
## 6              44201              Ozone   25860 Hickory-Lenoir-Morganton, NC
##   STATE_CODE      STATE COUNTY_CODE   COUNTY SITE_LATITUDE SITE_LONGITUDE
## 1          37 North Carolina          3 Alexander      35.9138      -81.191
## 2          37 North Carolina          3 Alexander      35.9138      -81.191
## 3          37 North Carolina          3 Alexander      35.9138      -81.191
## 4          37 North Carolina          3 Alexander      35.9138      -81.191
## 5          37 North Carolina          3 Alexander      35.9138      -81.191
## 6          37 North Carolina          3 Alexander      35.9138      -81.191
```

```
str(EPA.03.2018)
```

```
## '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 ...
```

```
dim(EPA.03.2019)
```

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

```
colnames(EPA.03.2019)
```

```
## [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"
```

```
head(EPA.03.2019)
```

```
##           Date Source   Site.ID POC Daily.Max.8.hour.Ozone.Concentration UNITS
## 1 01/01/2019 AirNow 370030005 1 0.029 ppm
## 2 01/02/2019 AirNow 370030005 1 0.018 ppm
## 3 01/03/2019 AirNow 370030005 1 0.016 ppm
## 4 01/04/2019 AirNow 370030005 1 0.022 ppm
## 5 01/05/2019 AirNow 370030005 1 0.037 ppm
## 6 01/06/2019 AirNow 370030005 1 0.037 ppm
##   DAILY_AQI_VALUE Site.Name DAILY_OBS_COUNT PERCENT_COMPLETE
## 1              27 Taylorsville Liledoun          24          100
## 2              17 Taylorsville Liledoun          24          100
## 3              15 Taylorsville Liledoun          24          100
## 4              20 Taylorsville Liledoun          24          100
## 5              34 Taylorsville Liledoun          24          100
## 6              34 Taylorsville Liledoun          24          100
##   AQS_PARAMETER_CODE AQS_PARAMETER_DESC CBSA_CODE CBSA_NAME
## 1              44201 Ozone 25860 Hickory-Lenoir-Morganton, NC
## 2              44201 Ozone 25860 Hickory-Lenoir-Morganton, NC
## 3              44201 Ozone 25860 Hickory-Lenoir-Morganton, NC
## 4              44201 Ozone 25860 Hickory-Lenoir-Morganton, NC
## 5              44201 Ozone 25860 Hickory-Lenoir-Morganton, NC
## 6              44201 Ozone 25860 Hickory-Lenoir-Morganton, NC
##   STATE_CODE STATE COUNTY_CODE COUNTY SITE_LATITUDE SITE_LONGITUDE
## 1          37 North Carolina 3 Alexander 35.9138 -81.191
## 2          37 North Carolina 3 Alexander 35.9138 -81.191
## 3          37 North Carolina 3 Alexander 35.9138 -81.191
## 4          37 North Carolina 3 Alexander 35.9138 -81.191
## 5          37 North Carolina 3 Alexander 35.9138 -81.191
## 6          37 North Carolina 3 Alexander 35.9138 -81.191
```

```
str(EPA.03.2019)
```

```
## 'data.frame': 10592 obs. of 20 variables:
## $ Date : Factor w/ 365 levels "01/01/2019","01/02/2019",...: 1 2 3 4 ...
## $ Source : Factor w/ 2 levels "AirNow","AQS": 1 1 1 1 1 1 1 1 1 ...
## $ 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 0.029 0.018 0.016 0.022 0.037 0.037 0.029 0.038 0.038 ...
## $ UNITS : Factor w/ 1 level "ppm": 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 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 ...
## $ 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 ...
## $ 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 ...
```

```
dim(EPA.PM.2018)
```

```
## [1] 8983 20
```

```
colnames(EPA.PM.2018)
```

```
## [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"
```

```
head(EPA.PM.2018)
```

```
##      Date Source  Site.ID POC Daily.Mean.PM2.5.Concentration  UNITS
## 1 01/02/2018  AQS 370110002  1                2.9 ug/m3 LC
## 2 01/05/2018  AQS 370110002  1                3.7 ug/m3 LC
## 3 01/08/2018  AQS 370110002  1                5.3 ug/m3 LC
## 4 01/11/2018  AQS 370110002  1                0.8 ug/m3 LC
## 5 01/14/2018  AQS 370110002  1                2.5 ug/m3 LC
## 6 01/17/2018  AQS 370110002  1                4.5 ug/m3 LC
##      DAILY_AQI_VALUE      Site.Name DAILY_OBS_COUNT PERCENT_COMPLETE
## 1                12 Linville Falls                1                100
## 2                15 Linville Falls                1                100
## 3                22 Linville Falls                1                100
## 4                 3 Linville Falls                1                100
## 5                10 Linville Falls                1                100
## 6                19 Linville Falls                1                100
##      AQS_PARAMETER_CODE      AQS_PARAMETER_DESC CBSA_CODE CBSA_NAME
## 1                88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 2                88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 3                88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 4                88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 5                88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 6                88502 Acceptable PM2.5 AQI & Speciation Mass      NA
##      STATE_CODE      STATE COUNTY_CODE COUNTY SITE_LATITUDE SITE_LONGITUDE
## 1                37 North Carolina                11 Avery      35.97235      -81.93307
## 2                37 North Carolina                11 Avery      35.97235      -81.93307
## 3                37 North Carolina                11 Avery      35.97235      -81.93307
## 4                37 North Carolina                11 Avery      35.97235      -81.93307
## 5                37 North Carolina                11 Avery      35.97235      -81.93307
## 6                37 North Carolina                11 Avery      35.97235      -81.93307
```

```
str(EPA.PM.2018)
```

```
## '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 3
## $ 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 ...
```

```
dim(EPA.PM.2019)
```

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

```
colnames(EPA.PM.2019)
```

```
## [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"
```

```
head(EPA.PM.2019)
```

```
##      Date Source Site.ID POC Daily.Mean.PM2.5.Concentration UNITS
## 1 01/03/2019 AQS 370110002 1 1.6 ug/m3 LC
## 2 01/06/2019 AQS 370110002 1 1.0 ug/m3 LC
## 3 01/09/2019 AQS 370110002 1 1.3 ug/m3 LC
## 4 01/12/2019 AQS 370110002 1 6.3 ug/m3 LC
## 5 01/15/2019 AQS 370110002 1 2.6 ug/m3 LC
## 6 01/18/2019 AQS 370110002 1 1.2 ug/m3 LC
## DAILY_AQI_VALUE Site.Name DAILY_OBS_COUNT PERCENT_COMPLETE
```

```
## 1      7 Linville Falls      1      100
## 2      4 Linville Falls      1      100
## 3      5 Linville Falls      1      100
## 4     26 Linville Falls      1      100
## 5     11 Linville Falls      1      100
## 6      5 Linville Falls      1      100
##   AQS_PARAMETER_CODE      AQS_PARAMETER_DESC CBSA_CODE CBSA_NAME
## 1      88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 2      88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 3      88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 4      88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 5      88502 Acceptable PM2.5 AQI & Speciation Mass      NA
## 6      88502 Acceptable PM2.5 AQI & Speciation Mass      NA
##   STATE_CODE      STATE COUNTY_CODE COUNTY SITE_LATITUDE SITE_LONGITUDE
## 1      37 North Carolina      11 Avery      35.97235      -81.93307
## 2      37 North Carolina      11 Avery      35.97235      -81.93307
## 3      37 North Carolina      11 Avery      35.97235      -81.93307
## 4      37 North Carolina      11 Avery      35.97235      -81.93307
## 5      37 North Carolina      11 Avery      35.97235      -81.93307
## 6      37 North Carolina      11 Avery      35.97235      -81.93307
```

```
str(EPA.PM.2019)
```

```
## '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 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 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 1 1 1 1 1 1 1 1 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 ...
```

## 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

```
class(EPA.03.2018$Date)
```

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

```
EPA.03.2018$Date <- mdy(EPA.03.2018$Date)
class(EPA.03.2018$Date)
```

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

```
class(EPA.03.2019$Date)
```

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

```
EPA.03.2019$Date <- mdy(EPA.03.2019$Date)
class(EPA.03.2019$Date)
```

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

```
class(EPA.PM.2018$Date)
```

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

```
EPA.PM.2018$Date <- mdy(EPA.PM.2018$Date)
class(EPA.PM.2018$Date)
```

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

```
class(EPA.PM.2019$Date)
```

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

```
EPA.PM.2019$Date <- mdy(EPA.PM.2019$Date)
class(EPA.PM.2019$Date)
```

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

#4

```
EPA.03.2018.select <- select(EPA.03.2018, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY,
EPA.03.2019.select <- select(EPA.03.2019, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY,
EPA.PM.2018.select <- select(EPA.PM.2018, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY,
```



```

EPA.PM.2019.select <- select(EPA.PM.2019, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC, COUNTY,
#5

EPA.PM.2018.select <- mutate(EPA.PM.2018.select, AQS_PARAMETER_DESC = ("PM2.5"))

EPA.PM.2019.select <- mutate(EPA.PM.2019.select, AQS_PARAMETER_DESC = ("PM2.5"))

#6

write.csv(EPA.O3.2018.select, row.names = FALSE,
          file = "/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Process
write.csv(EPA.O3.2019.select, row.names = FALSE,
          file = "/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Process
write.csv(EPA.PM.2018.select, row.names = FALSE,
          file = "/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Process
write.csv(EPA.PM.2019.select, row.names = FALSE,
          file = "/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Process

```

## 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:
  - 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.
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\_NC2122\_Processed.csv”

```

#7

#O3.2018 <- read.csv("/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Pro
#O3.2019 <- read.csv("/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Pro
#PM.2018 <- read.csv("/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Pro

```

```
#PM.2019 <- read.csv("/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Pro
03.PM.Combined <- rbind(EPA.O3.2018.select , EPA.O3.2019.select, EPA.PM.2018.select, EPA.PM.2019.select,
#8
is.POSIXlt(03.PM.Combined$Date)
```

```
## [1] FALSE
```

```
03.PM.Updated <-
  03.PM.Combined %>%
  filter(Site.Name == "Linville Falls" | Site.Name == "Durham Armory" | Site.Name == "Leggett" | Site.N
  group_by(Date, Site.Name, AQS_PARAMETER_DESC, COUNTY) %>%
  summarise(meanAQI = mean(DAILY_AQI_VALUE),
             mean(SITE_LATITUDE),
             mean(SITE_LONGITUDE)) %>%
  mutate(Month = month(Date)) %>%
  mutate(Year = year(Date))
```

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

```
#9
03.PM.Updated.Spread <- pivot_wider(03.PM.Updated, names_from = AQS_PARAMETER_DESC, values_from = meanA
#10
dim(03.PM.Updated.Spread)
```

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

```
#11
write.csv(03.PM.Updated.Spread, row.names = FALSE,
          file = "/Users/jessicacitrola/Documents/ENV872/Environmental_Data_Analytics_2022/Data/Process
```

## 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)
```

```
03.PM.summary <-  
  03.PM.Updated.Spread %>%  
  group_by(Site.Name, Month, Year) %>%  
  filter(!is.na(PM2.5) & !is.na(Ozone)) %>%  
  summarise(meanAQI.PM = mean(PM2.5),  
            meanAQI.Ozone = mean(Ozone))
```

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

```
#13
```

```
dim(03.PM.summary)
```

```
## [1] 223 5
```

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

Answer:

#Not being graded, used is.na.