

Assignment 4: Data Wrangling

Laura Brockington

Spring 2023

OVERVIEW

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

Directions

1. Rename this file `<FirstLast>_A04_DataWrangling.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, Feb 20th @ 5:00pm.

Set up your session

- 1a. Load the `tidyverse`, `lubridate`, and `here` packages into your session.
 - 1b. Check your working directory.
 - 1c. Read in all four raw data files associated with the EPA Air dataset, being sure to set string columns to be read in as factors. See the README file for the EPA air datasets for more information (especially if you have not worked with air quality data previously).
2. Apply the `glimpse()` function to reveal the dimensions, column names, and structure of each dataset.

```
# the formatR package nor the width cutoff worked to keep my code from running
# off the page when I knit, so I added this to the header 'header-includes:
# \usepackage{fveextra}
# \DefineVerbatimEnvironment{Highlighting}{Verbatim}{breaklines,commandchars=\\\{\}\}'
# 1a
library(tidyverse) #loading necessary packages
library(lubridate)
library(here)
# 1b
getwd() #checking working directory
```

```
## [1] "/Users/laura/Desktop/EDA/EDA"
```

```

# 1c
o3_2018 <- read.csv(file = here("./Data/Raw/EPAair_O3_NC2018_raw.csv"), header = TRUE,
  stringsAsFactors = TRUE) #reading in datasets
o3_2019 <- read.csv(file = here("./Data/Raw/EPAair_O3_NC2019_raw.csv"), header = TRUE,
  stringsAsFactors = TRUE)
pm25_2018 <- read.csv(file = here("./Data/Raw/EPAair_PM25_NC2018_raw.csv"), header =
  TRUE,
  stringsAsFactors = TRUE)
pm25_2019 <- read.csv(file = here("./Data/Raw/EPAair_PM25_NC2019_raw.csv"), header =
  TRUE,
  stringsAsFactors = TRUE)
# 2
glimpse(o3_2018) #revealing the structure and columns of each dataset

```

```

## Rows: 9,737
## Columns: 20
## $ Date <fct> 03/01/2018, 03/02/2018, 03/03/201~
## $ Source <fct> AQS, AQS, AQS, AQS, AQS, AQS, AQS~
## $ Site.ID <int> 370030005, 370030005, 370030005, ~
## $ POC <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Max.8.hour.Ozone.Concentration <dbl> 0.043, 0.046, 0.047, 0.049, 0.047~
## $ UNITS <fct> ppm, ppm, ppm, ppm, ppm, ppm, ppm~
## $ DAILY_AQI_VALUE <int> 40, 43, 44, 45, 44, 28, 33, 41, 4~
## $ Site.Name <fct> Taylorsville Liledoun, Taylorsvil~
## $ DAILY_OBS_COUNT <int> 17, 17, 17, 17, 17, 17, 17, 17, 1~
## $ PERCENT_COMPLETE <dbl> 100, 100, 100, 100, 100, 100, 100~
## $ AQS_PARAMETER_CODE <int> 44201, 44201, 44201, 44201, 44201~
## $ AQS_PARAMETER_DESC <fct> Ozone, Ozone, Ozone, Ozone, Ozone~
## $ CBSA_CODE <int> 25860, 25860, 25860, 25860, 25860~
## $ CBSA_NAME <fct> "Hickory-Lenoir-Morganton, NC", "~
## $ STATE_CODE <int> 37, 37, 37, 37, 37, 37, 37, 37, 3~
## $ STATE <fct> North Carolina, North Carolina, N~
## $ COUNTY_CODE <int> 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, ~
## $ COUNTY <fct> Alexander, Alexander, Alexander, ~
## $ SITE_LATITUDE <dbl> 35.9138, 35.9138, 35.9138, 35.913~
## $ SITE_LONGITUDE <dbl> -81.191, -81.191, -81.191, -81.19~

```

```
glimpse(o3_2019)
```

```

## Rows: 10,592
## Columns: 20
## $ Date <fct> 01/01/2019, 01/02/2019, 01/03/201~
## $ Source <fct> AirNow, AirNow, AirNow, AirNow, A~
## $ Site.ID <int> 370030005, 370030005, 370030005, ~
## $ POC <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Max.8.hour.Ozone.Concentration <dbl> 0.029, 0.018, 0.016, 0.022, 0.037~
## $ UNITS <fct> ppm, ppm, ppm, ppm, ppm, ppm, ppm~
## $ DAILY_AQI_VALUE <int> 27, 17, 15, 20, 34, 34, 27, 35, 3~
## $ Site.Name <fct> Taylorsville Liledoun, Taylorsvil~
## $ DAILY_OBS_COUNT <int> 24, 24, 24, 24, 24, 24, 24, 24, 2~
## $ PERCENT_COMPLETE <dbl> 100, 100, 100, 100, 100, 100, 100~
## $ AQS_PARAMETER_CODE <int> 44201, 44201, 44201, 44201, 44201~

```

```
## $ AQS_PARAMETER_DESC      <fct> Ozone, Ozone, Ozone, Ozone, Ozone~
## $ CBSA_CODE                <int> 25860, 25860, 25860, 25860, 25860~
## $ CBSA_NAME                <fct> "Hickory-Lenoir-Morganton, NC", "~
## $ STATE_CODE               <int> 37, 37, 37, 37, 37, 37, 37, 37, 3~
## $ STATE                   <fct> North Carolina, North Carolina, N~
## $ COUNTY_CODE              <int> 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, ~
## $ COUNTY                  <fct> Alexander, Alexander, Alexander, ~
## $ SITE_LATITUDE            <dbl> 35.9138, 35.9138, 35.9138, 35.913~
## $ SITE_LONGITUDE           <dbl> -81.191, -81.191, -81.191, -81.19~
```

```
glimpse(pm25_2018)
```

```
## Rows: 8,983
## Columns: 20
## $ Date                  <fct> 01/02/2018, 01/05/2018, 01/08/2018, 01/~
## $ Source                <fct> AQS, AQS, AQS, AQS, AQS, AQS, AQS, AQS,~
## $ Site.ID               <int> 370110002, 370110002, 370110002, 370110~
## $ POC                   <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Mean.PM2.5.Concentration <dbl> 2.9, 3.7, 5.3, 0.8, 2.5, 4.5, 1.8, 2.5,~
## $ UNITS                 <fct> 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, 24~
## $ Site.Name             <fct> Linville Falls, Linville Falls, Linvill~
## $ DAILY_OBS_COUNT       <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ PERCENT_COMPLETE      <dbl> 100, 100, 100, 100, 100, 100, 100, 100,~
## $ AQS_PARAMETER_CODE    <int> 88502, 88502, 88502, 88502, 88502, 8850~
## $ AQS_PARAMETER_DESC    <fct> Acceptable PM2.5 AQI & Speciation Mass,~
## $ CBSA_CODE             <int> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA,~
## $ CBSA_NAME             <fct> "", "", "", "", "", "", "", "", "", "",~
## $ STATE_CODE            <int> 37, 37, 37, 37, 37, 37, 37, 37, 37, 37,~
## $ STATE                 <fct> North Carolina, North Carolina, North C~
## $ COUNTY_CODE           <int> 11, 11, 11, 11, 11, 11, 11, 11, 11, 11,~
## $ COUNTY                <fct> Avery, Avery, Avery, Avery, Avery, Aver~
## $ SITE_LATITUDE         <dbl> 35.97235, 35.97235, 35.97235, 35.97235,~
## $ SITE_LONGITUDE        <dbl> -81.93307, -81.93307, -81.93307, -81.93~
```

```
glimpse(pm25_2019)
```

```
## Rows: 8,581
## Columns: 20
## $ Date                  <fct> 01/03/2019, 01/06/2019, 01/09/2019, 01/~
## $ Source                <fct> AQS, AQS, AQS, AQS, AQS, AQS, AQS, AQS,~
## $ Site.ID               <int> 370110002, 370110002, 370110002, 370110~
## $ POC                   <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ Daily.Mean.PM2.5.Concentration <dbl> 1.6, 1.0, 1.3, 6.3, 2.6, 1.2, 1.5, 1.5,~
## $ UNITS                 <fct> ug/m3 LC, ug/m3 LC, ug/m3 LC, ug/m3 LC,~
## $ DAILY_AQI_VALUE       <int> 7, 4, 5, 26, 11, 5, 6, 15, 7, 14, 20~
## $ Site.Name             <fct> Linville Falls, Linville Falls, Linvill~
## $ DAILY_OBS_COUNT       <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ PERCENT_COMPLETE      <dbl> 100, 100, 100, 100, 100, 100, 100, 100,~
## $ AQS_PARAMETER_CODE    <int> 88502, 88502, 88502, 88502, 88502, 8850~
## $ AQS_PARAMETER_DESC    <fct> Acceptable PM2.5 AQI & Speciation Mass,~
## $ CBSA_CODE             <int> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA,~
```

```
## $ CBSA_NAME          <fct> "", "", "", "", "", "", "", "", "", "", "~
## $ STATE_CODE         <int> 37, 37, 37, 37, 37, 37, 37, 37, 37, 37, ~
## $ STATE              <fct> North Carolina, North Carolina, North C~
## $ COUNTY_CODE        <int> 11, 11, 11, 11, 11, 11, 11, 11, 11, 11, ~
## $ COUNTY             <fct> Avery, Avery, Avery, Avery, Avery, Aver~
## $ SITE_LATITUDE       <dbl> 35.97235, 35.97235, 35.97235, 35.97235, ~
## $ SITE_LONGITUDE      <dbl> -81.93307, -81.93307, -81.93307, -81.93~
```

Wrangle individual datasets to create processed files.

3. Change date columns to be date objects.
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
o3_2018$Date <- mdy(o3_2018$Date) #changing date class
o3_2019$Date <- mdy(o3_2019$Date)
pm25_2018$Date <- mdy(pm25_2018$Date)
pm25_2019$Date <- mdy(pm25_2019$Date)

# 4
o3_2018_subset <- select(o3_2018, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,
  COUNTY, SITE_LATITUDE, SITE_LONGITUDE) #selecting certain columns
o3_2019_subset <- select(o3_2019, Date, DAILY_AQI_VALUE, Site.Name, AQS_PARAMETER_DESC,
  COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
pm25_2018_subset <- select(pm25_2018, Date, DAILY_AQI_VALUE, Site.Name,
  ↪ AQS_PARAMETER_DESC,
  COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
pm25_2019_subset <- select(pm25_2019, Date, DAILY_AQI_VALUE, Site.Name,
  ↪ AQS_PARAMETER_DESC,
  COUNTY, SITE_LATITUDE, SITE_LONGITUDE)

# 5
pm25_2018_subset$AQS_PARAMETER_DESC <- "PM2.5" #setting AQS_PARAMETER_DESC as PM2.5
pm25_2019_subset$AQS_PARAMETER_DESC <- "PM2.5"

# 6
write.csv(o3_2018_subset, row.names = FALSE, file =
  ↪ "./Data/Processed/EPAair_03_NC2018_processed.csv") #saving subsetted datasets to
  ↪ processed folder
write.csv(o3_2019_subset, row.names = FALSE, file =
  ↪ "./Data/Processed/EPAair_03_NC2019_processed.csv")
write.csv(pm25_2018_subset, row.names = FALSE, file =
  ↪ "./Data/Processed/EPAair_PM25_NC2018_processed.csv")
write.csv(pm25_2019_subset, 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 - but it will include sites with missing site information...)
 - Some sites have multiple measurements per day. Use the split-apply-combine strategy to generate daily means: group by date, site name, 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
o3_2018$Daily.Mean.PM2.5.Concentration <- NA #adding columns with NA so the dataframes
  ↳ match
o3_2019$Daily.Mean.PM2.5.Concentration <- NA
pm25_2018$Daily.Max.8.hour.Ozone.Concentration <- NA
pm25_2019$Daily.Max.8.hour.Ozone.Concentration <- NA
EPAair_combine <- rbind(o3_2018, o3_2019, pm25_2018, pm25_2019) #combining the 4 datasets
#8
EPAair_combine_processed <- #wrangling dataset with a pipe function to fulfill the above
  ↳ requirements
EPAair_combine %>%
  mutate_at(c('Site.Name'), ~na_if(., "")) %>%
  drop_na(Site.Name) %>%
  filter(Site.Name == "Linville Falls" | Site.Name == "Durham Armory" | Site.Name ==
  ↳ "Leggett" | Site.Name == "Hattie Avenue" | Site.Name == "Clemmons Middle" |
  ↳ Site.Name == "Mendenhall School" | Site.Name == "Frying Pan Mountain" | Site.Name
  ↳ == "West Johnston Co." | Site.Name == "Garinger High School" | Site.Name == "Castle
  ↳ Hayne" | Site.Name == "Pitt Agri. Center" | Site.Name == "Bryson City" | Site.Name
  ↳ == "Millbrook School") %>%
  group_by(Date, Site.Name, AQS_PARAMETER_CODE, COUNTY) %>%
  summarize(mean_AQI = mean(DAILY_AQI_VALUE),
            mean_latitude = mean(SITE_LATITUDE),
            mean_longitude = mean(SITE_LONGITUDE)) %>%
  mutate(month = month(Date)) %>%
  mutate(year = year(Date))
```

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

```

#9
EPAair_by_AQI_type <- EPAair_combine_processed #duplicating dataset
EPAair_by_AQI_type <- pivot_wider(EPAair_by_AQI_type, names_from = AQS_PARAMETER_CODE,
  ↪ values_from = mean_AQI) #spreading the AQI_type column
colnames(EPAair_by_AQI_type)[8] = "PM2.5a" #changing column names
colnames(EPAair_by_AQI_type)[9] = "o3"
colnames(EPAair_by_AQI_type)[10] = "PM2.5b"
EPAair_by_AQI_type["PM2.5a"][is.na(EPAair_by_AQI_type["PM2.5a"])] <- 0 #setting NAs to 0
EPAair_by_AQI_type["PM2.5b"][is.na(EPAair_by_AQI_type["PM2.5b"])] <- 0
EPAair_by_AQI_type$PM2.5 <- (EPAair_by_AQI_type$PM2.5a + EPAair_by_AQI_type$PM2.5b)
  ↪ #summing the two PM2.5 columns and adding to new column
EPAair_by_AQI_type <- EPAair_by_AQI_type[, -8] #removing extra columns
EPAair_by_AQI_type <- EPAair_by_AQI_type[, -9]
EPAair_by_AQI_type["PM2.5"][EPAair_by_AQI_type["PM2.5"] == 0] <- NA #changing 0 back to
  ↪ NAs
#10
dim(EPAair_by_AQI_type) #calling dimensions of dataset

```

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

```

#11
write.csv(EPAair_by_AQI_type, row.names = FALSE,
  file = "./Data/Processed/EPAair_03_PM25_NC1819_Processed.csv") #saving dataset

```

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 mean **ozone** values are not available (use the function `drop_na` in your pipe). It's ok to have missing mean PM2.5 values in this result.

13. Call up the dimensions of the summary dataset.

```

#12
EPAair_summary <- #generating summary data frame grouped by site, month, and year with
  ↪ mean AQI values for ozone and PM2.5, then dropping rows with NAs as mean AQI for
  ↪ ozone or o3
EPAair_by_AQI_type %>%
  group_by(Site.Name, month, year) %>%
  summarize(mean_AQI_o3 = mean(o3),
    mean_AQI_PM2.5 = mean(PM2.5)) %>%
  drop_na(mean_AQI_o3)

```

```

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

```

```

#13
dim(EPAair_summary) #calling dimensions of dataset

```

```
## [1] 182    5
```

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

Answer: We used `drop_na()` because it eliminates any row that has a missing (i.e. NA) value from the specified column, while `na.omit()` eliminates any row with NA anywhere in the dataframe. In this case, we only wanted to omit those rows with NAs in the ozone column, so `drop_na` was a better choice.