

Assignment 4: Data Wrangling (Fall 2024)

Lauren Xu

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.
6. Ensure that code in code chunks does not extend off the page in the PDF.

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 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).
2. Add the appropriate code to reveal the dimensions of the four datasets.

```
#1a  
#Install packages  
library(tidyverse)  
library(lubridate)  
library(here)
```

```
#1b  
#Check your working directory  
here()
```

```
## [1] "/home/guest/EDE_Spring2025"
```

```
#1c  
# Read in all four datasets separately  
EPA_03_2018 <- read.csv(here("/home/guest/EDE_Spring2025/Data/Raw/EPAair_03_NC2018_raw.csv"),
```

```

stringsAsFactors = TRUE)

EPA_O3_2019 <- read.csv(here("/home/guest/EDE_Spring2025/Data/Raw/EPAair_O3_NC2019_raw.csv"),
  stringsAsFactors = TRUE)

EPA_PM25_2018 <- read.csv(here("/home/guest/EDE_Spring2025/Data/Raw/EPAair_PM25_NC2018_raw.csv"),
  stringsAsFactors = TRUE)

EPA_PM25_2019 <- read.csv(here("/home/guest/EDE_Spring2025/Data/Raw/EPAair_PM25_NC2019_raw.csv"),
  stringsAsFactors = TRUE)

#2
# Reveal the dimensions of each dataset
dim(EPA_O3_2018)

## [1] 9737    20

dim(EPA_O3_2019)

## [1] 10592    20

dim(EPA_PM25_2018)

## [1] 8983    20

dim(EPA_PM25_2019)

## [1] 8581    20

```

All four datasets should have the same number of columns but unique record counts (rows). Do your datasets follow this pattern? Answer: Yes. `dim(EPA_O3_2018)`
`[1] 9737 20` `dim(EPA_O3_2019)`
`[1] 10592 20` `dim(EPA_PM25_2018)` `[1] 8983 20` `dim(EPA_PM25_2019)` `[1] 8581 20`

Wrangle individual datasets to create processed files.

3. Change the 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
# Convert date columns to be date objects
EPA_03_2018$Date <- as.Date(EPA_03_2018$Date, format="%m/%d/%Y")
EPA_03_2019$Date <- as.Date(EPA_03_2019$Date, format="%m/%d/%Y")
EPA_PM25_2018$Date <- as.Date(EPA_PM25_2018$Date, format="%m/%d/%Y")
EPA_PM25_2019$Date <- as.Date(EPA_PM25_2019$Date, format="%m/%d/%Y")

#4
# Select the required columns
selected_cols <- c("Date", "DAILY_AQI_VALUE", "Site.Name", "AQ5_PARAMETER_DESC", "COUNTY", "SITE_LATITUDE")

EPA_03_2018 <- EPA_03_2018 %>% select(all_of(selected_cols))
EPA_03_2019 <- EPA_03_2019 %>% select(all_of(selected_cols))
EPA_PM25_2018 <- EPA_PM25_2018 %>% select(all_of(selected_cols))
EPA_PM25_2019 <- EPA_PM25_2019 %>% select(all_of(selected_cols))

#5
# Update AQ5_PARAMETER_DESC for PM2.5 Datasets
EPA_PM25_2018$AQ5_PARAMETER_DESC <- "PM2.5"
EPA_PM25_2019$AQ5_PARAMETER_DESC <- "PM2.5"

#6
# Define save paths
processed_folder_A4 <- here("/home/guest/EDE_Spring2025/Data/Processed/")

# Write the processed files
write.csv(EPA_03_2018, file=paste0(processed_folder_A4, "EPAair_03_NC2018_processed.csv"), row.names = F)
write.csv(EPA_03_2019, file=paste0(processed_folder_A4, "EPAair_03_NC2019_processed.csv"), row.names = F)
write.csv(EPA_PM25_2018, file=paste0(processed_folder_A4, "EPAair_PM25_NC2018_processed.csv"), row.names = F)
write.csv(EPA_PM25_2019, file=paste0(processed_folder_A4, "EPAair_PM25_NC2019_processed.csv"), row.names = F)

```

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 only 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, which you don’t want...)

- Some sites have multiple measurements per day. Use the split-apply-combine strategy to generate daily means: group by date, site name, AQ5 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
# Combine all datasets into one
EPA_combined <- rbind(EPA_O3_2018, EPA_O3_2019, EPA_PM25_2018, EPA_PM25_2019)

#8
# Define the common sites
common_sites <- 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")

# Wrangle the dataset
EPA_cleaned <- EPA_combined %>%
  filter(Site.Name %in% common_sites) %>%
  mutate(
    Date = as.Date(Date, format="%Y-%m-%d"),
    Month = month(Date),
    Year = year(Date)
  ) %>%
  group_by(Date, Month, Year, Site.Name, AQS_PARAMETER_DESC, COUNTY) %>%
  summarise(
    DAILY_AQI_VALUE = mean(DAILY_AQI_VALUE, na.rm = TRUE),
    SITE_LATITUDE = mean(SITE_LATITUDE, na.rm = TRUE),
    SITE_LONGITUDE = mean(SITE_LONGITUDE, na.rm = TRUE)
  ) %>%
  ungroup()
```

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

```
#9
#Spread the datasets
EPA_wide <- EPA_cleaned %>%
  pivot_wider(
    names_from = AQS_PARAMETER_DESC,
    values_from = DAILY_AQI_VALUE)

#10
#Call the dimensions of new tidy dataset
dim(EPA_wide)
```

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

```
#11
#Save the new tidy dataset
write.csv(EPA_wide, file = "/home/guest/EDE_Spring2025/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 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
EPA_summary <- EPA_wide %>%
  mutate(Month = month(Date), Year = year(Date)) %>%
  group_by(Site.Name, Month, Year) %>%
  summarise(
    mean_ozone_AQI = mean(Ozone, na.rm = TRUE),
    mean_PM25_AQI = mean(PM2.5, na.rm = TRUE)
  ) %>%
  drop_na(mean_ozone_AQI) %>%
  ungroup()
```

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

```
#13
dim(EPA_summary)
```

```
## [1] 239 5
```

14. Why did we use the function `drop_na` rather than `na.omit`? Hint: replace `drop_na` with `na.omit` in part 12 and observe what happens with the dimensions of the summary data frame.

Answer: (1) Using `na.omit`: `dim(EPA_summary)=223 5`; (2) Using `drop_na`: `dim(EPA_summary)=239 5`. The `na.omit` removes entire rows if any column has an NA value, while the `drop_na` removes rows only where `mean_ozone_AQI` is NA but keeps other rows. Thus, using `drop_na` can save more valuable data.