

Assignment 4: Data Wrangling (Fall 2024)

Julia Kagiliery

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.

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2    3.5.1      v tibble    3.2.1
## v lubridate  1.9.4      v tidyr     1.3.1
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(lubridate)
library(here)
```

```
## here() starts at /Users/juliakagiliery/Library/Mobile Documents/com~apple~CloudDocs/GitHub Links/EDA
```

- 1b. Check your working directory.

```
print(getwd())
```

```
## [1] "/Users/juliakagiliery/Library/Mobile Documents/com~apple~CloudDocs/GitHub Links/EDAClas2025"
```

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

```
#1a  
#for all below, read in data and make strings factors  
EPAair_03_NC2018_raw <- read.csv(here("Data/Raw/EPAair_03_NC2018_raw.csv"),  
stringsAsFactors = TRUE)  
print(dim(EPAair_03_NC2018_raw))
```

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

```
#1b  
EPAair_03_NC2019_raw <- read.csv(here("Data/Raw/EPAair_03_NC2019_raw.csv"),  
stringsAsFactors = TRUE)  
print(dim(EPAair_03_NC2019_raw))
```

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

```
#1c  
EPAair_PM25_NC2018_raw <- read.csv(here("Data/Raw/EPAair_PM25_NC2018_raw.csv"),  
stringsAsFactors = TRUE)  
print(dim(EPAair_PM25_NC2018_raw))
```

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

```
#2  
EPAair_PM25_NC2019_raw <- read.csv(here("Data/Raw/EPAair_PM25_NC2019_raw.csv"),  
stringsAsFactors = TRUE)  
print(dim(EPAair_PM25_NC2019_raw))
```

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

All four datasets should have the same number of columns but unique record counts (rows). Do your datasets follow this pattern? Yes, this is true. See based on prints that all have 20 columns but anywhere between ~8,500 to ~10,600 rows.

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
EPAair_03_NC2018_processed <- EPAair_03_NC2018_raw |> #make dates actually dates
  mutate(Date = as.Date(Date, format = "%m/%d/%Y")) #need to specify the format

EPAair_03_NC2019_processed <- EPAair_03_NC2019_raw |>
  mutate(Date = as.Date(Date, format = "%m/%d/%Y"))

EPAair_PM25_NC2018_processed <- EPAair_PM25_NC2018_raw |>
  mutate(Date = as.Date(Date, format = "%m/%d/%Y"))

EPAair_PM25_NC2019_processed <- EPAair_PM25_NC2019_raw |>
  mutate(Date = as.Date(Date, format = "%m/%d/%Y"))

#4
EPAair_03_NC2018_processed <- EPAair_03_NC2018_processed |> #only want these specified columns
  select(
    Date,
    DAILY_AQI_VALUE,
    Site.Name,
    AQS_PARAMETER_DESC,
    COUNTY,
    SITE_LATITUDE,
    SITE_LONGITUDE
  )

EPAair_03_NC2019_processed <- EPAair_03_NC2019_processed |>
  select(
    Date,
    DAILY_AQI_VALUE,
    Site.Name,
    AQS_PARAMETER_DESC,
    COUNTY,
    SITE_LATITUDE,
    SITE_LONGITUDE
  )

EPAair_PM25_NC2018_processed <- EPAair_PM25_NC2018_processed |>
  select(
    Date,
    DAILY_AQI_VALUE,
    Site.Name,
    AQS_PARAMETER_DESC,
    COUNTY,
    SITE_LATITUDE,
    SITE_LONGITUDE
  )

EPAair_PM25_NC2019_processed <- EPAair_PM25_NC2019_processed |>
```

```

select(
  Date,
  DAILY_AQI_VALUE,
  Site.Name,
  AQS_PARAMETER_DESC,
  COUNTY,
  SITE_LATITUDE,
  SITE_LONGITUDE
)

#5
EPAair_PM25_NC2018_processed <- EPAair_PM25_NC2018_processed |>
  mutate(AQS_PARAMETER_DESC = "PM2.5") #need the column to be only this value

EPAair_PM25_NC2019_processed <- EPAair_PM25_NC2019_processed |>
  mutate(AQS_PARAMETER_DESC = "PM2.5")

#6
write_csv(
  EPAair_03_NC2018_processed,
  here("Data/Processed/EPAair_03_NC2018_processed.csv") #saving in this relative path with this name
)
write_csv(
  EPAair_03_NC2019_processed,
  here("Data/Processed/EPAair_03_NC2019_processed.csv")
)
write_csv(
  EPAair_PM25_NC2018_processed,
  here("Data/Processed/EPAair_PM25_NC2018_processed.csv")
)
write_csv(
  EPAair_PM25_NC2019_processed,
  here("Data/Processed/EPAair_PM25_NC2019_processed.csv")
)

#I checked the dimensions of all my processed data sets and I get the same dimensions as above (9737, 1

```

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, 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
AirQualityO3PMNC <- rbind(EPAair_O3_NC2018_processed, EPAair_O3_NC2019_processed, EPAair_PM25_NC2018_processed, EPAair_PM25_NC2019_processed)
print(dim(AirQualityO3PMNC)) # per Dr. Fay's advice, this should be 37,893 x 7 which I do have.

## [1] 37893      7
```

```
#8
#the below code finds where the 4 data sets have common sites then I manually exclude blanks (this is good)
common_sites <- Reduce(
  intersect,
  list(
    EPAair_O3_NC2018_processed$Site.Name,
    EPAair_O3_NC2019_processed$Site.Name,
    EPAair_PM25_NC2018_processed$Site.Name,
    EPAair_PM25_NC2019_processed$Site.Name
  )
)

common_sites <- common_sites[common_sites != ""]
print(common_sites) #just want to make sure this works and it did!
```

```
## [1] "Linville Falls"      "Durham Armory"      "Leggett"
## [4] "Hattie Avenue"       "Clemmons Middle"   "Mendenhall School"
## [7] "Frying Pan Mountain" "West Johnston Co."  "Garinger High School"
## [10] "Castle Hayne"        "Pitt Agri. Center"  "Bryson City"
## [13] "Millbrook School"
```

```
AirQualityO3PMNC <- AirQualityO3PMNC |>
  filter(Site.Name %in% common_sites) #keeping only sites in common

print(dim(AirQualityO3PMNC)) # per Dr. Fay, this should be 16,510 observations in 7 columns which I have

## [1] 16510      7
```

```

AirQuality03PMNC <- AirQuality03PMNC |> #creating the summary sites with AQI, longitude, and latitude m
  group_by(Date, Site.Name, AQS_PARAMETER_DESC, COUNTY) |>
  summarise(
    MEAN_DAILY_AQI_VALUE = mean(DAILY_AQI_VALUE, na.rm = TRUE),
    MEAN_SITE_LONGITUDE = mean(SITE_LONGITUDE, na.rm = TRUE),
    MEAN_SITE_LATITUDE = mean(SITE_LATITUDE, na.rm = TRUE),
    .groups = "drop"
  )

AirQuality03PMNC <- AirQuality03PMNC |> # add the two new rows requested, year and month (need proper f
  mutate(
    Year = year(Date),
    Month = month(Date)
  )

print(dim(AirQuality03PMNC)) # per the instructions, this should be 14,752 x 9, correct!

```

```
## [1] 14752      9
```

```

#9
AirQuality03PMNCWider <- AirQuality03PMNC |> #now have columns for ozone and PM so the row shows both v
  pivot_wider(
    names_from = AQS_PARAMETER_DESC, # Column that will become new column names #name of my new column
    values_from = MEAN_DAILY_AQI_VALUE # Column that fills in the new columns
  )

#10
print(dim(AirQuality03PMNCWider))

```

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

```

#11
write_csv(AirQuality03PMNCWider, here("Data/Processed/EPAair_03_PM25_NC1819_Processed.csv")) # save the

```

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
AirQuality03PMNCWiderSummary <- AirQuality03PMNCWider |> #new columns are mean PM2.5 and mean ozone (as
  group_by(Site.Name, Month, Year) |>
  summarise(
    meanPM2.5 = mean(PM2.5),
    meanozone = mean(Ozone),
    .groups = "drop"
  ) |>

```

```
drop_na(meanozone) #drop NA mean ozone but we will still have NA PM 2.5  
  
#13  
print(dim(AirQuality03PMNCWiderSummary))
```

```
## [1] 182 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: `na.omit` will remove the whole column when there is an NA value present (so we lose the whole ozone column) but `drop_na` only removes rows (not columns) where we have an NA value which lets us maintain a column full of the actual ozone values.