

Assignment 4: Data Wrangling (Fall 2024)

Camber Vincent

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

```
#1b
getwd()
```

```
## [1] "/Users/cambervincent/EDA_Fall_2024/Assignments"
```

```
#1c
pm25_2019<-read.csv(here("Data/Raw/EPAair_PM25_NC2019_raw.csv"),
                    stringsAsFactors = TRUE)
pm25_2018<-read.csv(here("Data/Raw/EPAair_PM25_NC2018_raw.csv"),
                    stringsAsFactors = TRUE)
```

```
o3_2019<-read.csv(here("Data/Raw/EPAair_O3_NC2019_raw.csv"),
  stringsAsFactors = TRUE)
o3_2018<-read.csv(here("Data/Raw/EPAair_O3_NC2018_raw.csv"),
  stringsAsFactors = TRUE)
```

```
#2
dim(o3_2018)
```

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

```
dim(o3_2019)
```

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

```
dim(pm25_2018)
```

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

```
dim(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? All four of the datasets have 20 columns, but unique record counts. The O3_2018 dataset has 9737 records, the O3_2019 dataset has 10592 records, the PM25_2018 dataset has 8983 records, and the PM25_2019 dataset has 8581 records.

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
#Date conversion, format specified to properly read in all values
o3_2018$Date<-as.Date(o3_2018$Date,format="%m/%d/%Y")
o3_2019$Date<-as.Date(o3_2019$Date,format="%m/%d/%Y")
pm25_2018$Date<-as.Date(pm25_2018$Date,format="%m/%d/%Y")
pm25_2019$Date<-as.Date(pm25_2019$Date,format="%m/%d/%Y")
```

```
#4
#Selected columns piped to new data sets
o3_2018_selected<-o3_2018%>%
```

```

    select(Date,DAILY_AQI_VALUE,Site.Name,AQS_PARAMETER_DESC,COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
o3_2019_selected<-o3_2019%>%
    select(Date,DAILY_AQI_VALUE,Site.Name,AQS_PARAMETER_DESC,COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
pm25_2018_selected<-pm25_2018%>%
    select(Date,DAILY_AQI_VALUE,Site.Name,AQS_PARAMETER_DESC,COUNTY, SITE_LATITUDE, SITE_LONGITUDE)
pm25_2019_selected<-pm25_2019%>%
    select(Date,DAILY_AQI_VALUE,Site.Name,AQS_PARAMETER_DESC,COUNTY, SITE_LATITUDE, SITE_LONGITUDE)

#5
#Used dplyr package language to mutate column data, rather than base r for readability
pm25_2018_selected<-pm25_2018_selected%>%
    mutate(AQS_PARAMETER_DESC="PM2.5")
pm25_2019_selected<-pm25_2019_selected%>%
    mutate(AQS_PARAMETER_DESC="PM2.5")

#6
#Saving csv files to the processed data folder
write.csv(pm25_2019_selected,
          file="/Users/cambervincent/EDA_Fall_2024/Data/Processed/EPAair_PM25_NC2019_processed.csv",
          row.names=FALSE)
write.csv(pm25_2018_selected,
          file="/Users/cambervincent/EDA_Fall_2024/Data/Processed/EPAair_PM25_NC2018_processed.csv",
          row.names=FALSE)
write.csv(o3_2019_selected,
          file="/Users/cambervincent/EDA_Fall_2024/Data/Processed/EPAair_O3_NC2019_processed.csv",
          row.names=FALSE)
write.csv(o3_2018_selected,
          file="/Users/cambervincent/EDA_Fall_2024/Data/Processed/EPAair_O3_NC2018_processed.csv",
          row.names=FALSE)

```

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
o3_pm25_1819<-rbind(o3_2018_selected,
                    o3_2019_selected,
                    pm25_2018_selected,
                    pm25_2019_selected) #rbind run

#8
o3_pm25_1819_wrangled<-o3_pm25_1819%>%
  filter(Site.Name%in%
         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"))%>% #filter to site locations

  group_by(Date,Site.Name,AQS_PARAMETER_DESC,COUNTY)%>% #group_by step of split-combine-apply
  summarize(AQI_mean=mean(DAILY_AQI_VALUE),
            lat_mean=mean(SITE_LATITUDE),
            lon_mean=mean(SITE_LONGITUDE), #mean value calculation of split-combine-apply
            .groups="drop")%>% #ungroups previously grouped data for further analysis

  mutate(Month=month(Date),Year=year(Date)) #adding month and year columns

#9
o3_pm25_1819_tidy<-pivot_wider(o3_pm25_1819_wrangled,
                               names_from=AQS_PARAMETER_DESC,
                               values_from=AQI_mean) #spreading out data based on parameter

#10
dim(o3_pm25_1819_tidy) #call dimensions

## [1] 8976    9

#11
write.csv(o3_pm25_1819_tidy,
          file="/Users/cambervincent/EDA_Fall_2024/Data/Processed/EPAair_O3_PM25_NC1819_Processed.csv",
          row.names=FALSE) #saving csv file
```

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
o3_pm25_1819_summary<-o3_pm25_1819_tidy%>%
  group_by(Site.Name,Month,Year)%>% #group by site, month, and year per instructions

  summarize(
    PM2.5=mean(PM2.5,na.rm=TRUE),
    Ozone=mean(Ozone,na.rm=TRUE), #remove n/a values from mean calculation
    .groups="drop"#prevents "groups" message from displaying in final Rmd knit
                    #ungrouping is not needed for further analysis, just for visualization
  )%>%

  drop_na(Ozone) #removing instances where Ozone values are unavailable

#13
dim(o3_pm25_1819_summary) #call dimensions

```

```
## [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: When `drop_na` was used, the final data frame had the dimensions of 239 x 5. The `na.omit` function led to a final summary data frame with dimensions of 223 x 5. 16 rows of data were additionally omitted when the `na.omit` function was used. `Drop_na` only removes the rows of the dataframe where Ozone values are not available, whereas `na.omit` removed the rows where any column (including PM2.5) had an n/a or NaN value.