

Analysis of the effect of PM10 Concentration,
Temperature, Geographic Location,
Population, and distance to
Electricity-Generation Combustion Points on
the concentration of fine particulate matter
(PM2.5) in North Carolina during the year
2018.

https://github.com/fr55/DataAnalytics_FinalProject

Felipe Raby Amadori

Abstract

Experimental overview. This section should be no longer than 250 words.

Contents

1	Research Question and Rationale	5
2	Dataset Information	6
2.1	EPA PM2.5 Dataset	6
2.1.1	Data Content Information	6
2.2	EPA PM10 Dataset	6
2.2.1	Data Content Information	7
2.3	NOAA Average Temperature Dataset	7
2.3.1	Data Content Information	8
2.4	US Census Bureau US counties shapefile	8
2.4.1	Data Content Information	8
2.5	EPA combustion points for electricity generation in the US Dataset	8
2.5.1	Data Content Information	9
2.6	Analyzed data structure	9
3	Exploratory Data Analysis and Wrangling	11
3.1	EPA PM2.5 and PM10 Datasets	11
3.2	North Carolina Counties Zoning, Geographic information, and Population Data	21
3.3	US Census Bureau US counties shapefile	23
3.4	NOAA Average Temperature Dataset	25
3.5	EPA combustion points for electricity generation in the US Dataset	33
3.6	Additional previsualization of the data	35
4	Analysis	36
5	Summary and Conclusions	37

List of Tables

1	Selections	6
2	Dataset content	7
3	Selections	8
4	Dataset content	9
5	Selections	10
6	Dataset content	10
7	Dataset content	11
8	Summary of data structure	11

List of Figures

1	PM2.5 NC 2018 frequency polygon.	16
2	PM2.5 NC 2018 boxplot.	17
3	PM2.5 NC 2018 scatterplot.	18
4	PM10 NC 2018 scatterplot.	19
5	PM2.5 NC Monitoring Stations Previsualization.	21
6	Counties exploratory map.	25
7	NC Zoning exploratory map.	26
8	Mean Annual Temperature exploratory map.	29
9	Daily Mean Temperature NC 2018 frequency polygon.	30
10	Daily Mean Temperature NC 2018 boxplot.	31
11	Daily Mean Temperature NC 2018 scatterplot.	32
12	Combustion points for electricity generation in the North Carolina.	34

1 Research Question and Rationale

Nowadays air pollution is one of the most relevant health issues in the world. It refers to the contamination of the air by chemicals, biological materials, and other types of pollutants that are harmful to human health. To solve the problem of air pollution, it's necessary to understand the problem, what are the causes, and search for solutions based on the findings.

Particulate matter with a diameter of less than 2.5 micrometers is called PM2.5, and it is a extremely harmful air pollutant because it consists of particles with diameters that are less than or equal to 2.5 microns in size, which can get deeply into the lung, and ultimately impair lung function.

This study focus on trying to understand how PM2.5 concentration in North Carolina vary with temperature, PM10 concentration, zoning (piedmont, coastal, mountain), population, elevation, and distance to combustion points for electricity generation. This last variable was included because according to the EPA combustion for electricity generation is the major point-source sector for PM2.5 in the USA (EPA, 2019).

The research question is: What are the effects of temperature, PM10 concentration, zoning (piedmont, coastal, mountain), population, elevation, distance to combustion points for electricity generation, in PM2.5 concentrations within North Carolina in the year 2018?

2 Dataset Information

For the analysis the following datasets were considered:

2.1 EPA PM2.5 Dataset

This dataset contains data from air quality monitoring of PM2.5 in North Carolina in 2018, and it was obtained using the Download Daily Data Tool in the United States Environmental Protection Agency (EPA) webpage <https://www.epa.gov/outdoor-air-quality-data/download-daily-data> where the options showed in Table 1 were selected:

Option	Selection
Pollutant	PM2.5
Year	2018
Geographic Area	North Carolina
Monitor Site	All Sites
Download	Download CSV (spreadsheet)

Table 1: Selections

The downloaded file was saved in the project folder path `./Data/Raw/` as `EPAair_PM25_NC2018_raw.csv` on 2019-03-31.

2.1.1 Data Content Information

The dataset contains daily mean PM2.5 concentration in $\mu\text{g}/\text{m}^3$ in 2018. Data from 24 stations in 21 different counties of North Carolina with their location in NAD83 lat/long coordinates.

The dataset contains 19 columns, which are shown in Table 2. Column names without description are self-explanatory.

2.2 EPA PM10 Dataset

This dataset contains data from air quality monitoring of PM10 in North Carolina in 2018, and it was obtained using the Download Daily Data Tool in the United States Environmental Protection Agency (EPA) webpage <https://www.epa.gov/outdoor-air-quality-data/download-daily-data> where the options showed in Table 3 were selected:

The downloaded file was saved in the project folder path `./Data/Raw/` as `EPAair_PM10_NC2018_raw.csv` on 2019-03-31.

Column	Description
Date	mm/dd/YY
Source	AQS (Air Quality System)
Site ID	A unique number identifying the site.
POC	“Parameter Occurrence Code”, distinguishes different instruments that measure the same parameter at the same site.
Daily Mean PM2.5 Concentration	
Units	Concentration Units
DAILY_AQI_VALUE	AQI = Air quality index
Site Name	
DAILY_OBS_COUNT	
PERCENT_COMPLETE	
AQS_PARAMETER_CODE	
AQS_PARAMETER_DESC	
CBSA_CODE	
CBSA_NAME	
STATE_CODE	
COUNTY_CODE	A unique number identifying the County.
COUNTY	
SITE_LATITUDE	NAD83
SITE_LONGITUDE	NAD83

Table 2: Dataset content

2.2.1 Data Content Information

The dataset contains daily mean PM10 concentration in ug/m3 in 2018. Data from 9 stations in 8 different counties of North Carolina with their location in NAD83 lat/long coordinates.

The dataset contains 19 columns, which are shown in Table 4. Column names without description are self-explanatory.

2.3 NOAA Average Temperature Dataset

This dataset contains data from temperature monitoring in North Carolina in 2018, and it was obtained using the Data Search Tool in the National Center for Environmental Information of the National Oceanic and Atmospheric Administration (NOAA). Webpage <https://www.ncdc.noaa.gov/cdo-web>. Options showed in Table 5 were selected: XXXXXXArreglar

The downloaded file was saved in the project folder path ./Data/Raw/ as NOAA_TAVG_NC2018_raw.csv on 2019-03-28.

Option	Selection
Pollutant	PM10
Year	2018
Geographic Area	North Carolina
Monitor Site	All Sites
Download	Download CSV (spreadsheet)

Table 3: Selections

2.3.1 Data Content Information

The dataset contains daily mean air temperature in Farenheit in 2018. Data from 39 stations in North Carolina with their location in NAD83 lat/long coordinates. No county information.

The dataset contains 7 columns, which are shown in Table 6. Column names without description are self-explanatory.

2.4 US Census Bureau US counties shapefile

This dataset contains geographic and geometric information of all the counties of the US. The data is in NAD83 lat/long coordinates. De file was provided by John Fay in the Environmental Data Analytics (ENV 872L) course at Duke University, Spring 2019.

The files containing the information were saved in the project folder path `./Data/Spatial/` as `cb_2017_us_county_20m` on 2019-03-28.

2.4.1 Data Content Information

The dataset contains geographic and geometric information of all the counties of the US in NAD83 lat/long coordinates.

The dataset contains 10 columns, which are shown in Table 7. Column names without description are self-explanatory.

2.5 EPA combustion points for electricity generation in the US Dataset

This dataset contains facility-level locations for combustion points for electricity generation in the US, and it was obtained from the United States Environmental Protection Agency (EPA) webpage <https://www3.epa.gov/air/emissions/where.htm>. The Top PM2.5 emitting sectors link was selected.

The downloaded file was saved in the project folder path `./Data/Raw/` as `EPA_ElecGenComb_US_raw.kml` on 2019-03-31.

Column	Description
Date	mm/dd/YY
Source	AQS (Air Quality System)
Site ID	A unique number identifying the site.
POC	“Parameter Occurrence Code”, distinguishes different instruments that measure the same parameter at the same site.
Daily Mean PM10 Concentration	
Units	Concentration Units
DAILY_AQI_VALUE	AQI = Air quality index
Site Name	
DAILY_OBS_COUNT	
PERCENT_COMPLETE	
AQS_PARAMETER_CODE	
AQS_PARAMETER_DESC	
CBSA_CODE	
CBSA_NAME	
STATE_CODE	A unique number identifying the County.
COUNTY_CODE	A unique number identifying the County.
COUNTY	
SITE_LATITUDE	NAD83
SITE_LONGITUDE	NAD83

Table 4: Dataset content

2.5.1 Data Content Information

The dataset is a kml file that contains combustion points for electricity generation in the US. The data is in WGS84 lat/long coordinates.

All data sets, variable, and files are named according to the following naming convention: *databasename_datatype_details_stage.format*, where:

- databasename refers to the database from where the data originated
- datatype is a description of data
- details are additional descriptive details, particularly important for processed data
- stage refers to the stage in data management pipelines (e.g., raw, cleaned, temp or processed)

2.6 Analyzed data structure

With these datasets an exploratory data analysis was done and for the study. The datasets were wrangled and a file called PM2.5_Full_Elev_utm.shp was created, which has the data

Option	Selection
Pollutant	PM10
Year	2018
Geographic Area	North Carolina
Monitor Site	All Sites
Download	Download CSV (spreadsheet)

Table 5: Selections

Column	Description
STATION	A unique code identifying the site.
NAME	Station Name
Site ID	A unique number identifying the site.
LATITUDE	NAD83
LONGITUDE	NAD83
DATE	dd/mm/YY
TAVG	Daily Average Temperature in °F

Table 6: Dataset content

structure shown in Table 8.

Column	Description
STATEFP	A unique number identifying the State.
COUNTYFP	County Federal Information Processing Standards (FIPS) Code
COUNTYNS	Provides the American National Standards Institute (ANSI) code for the county or equivalent entity, as used by GNIS.
AFFGEOID	AFF Summary Level Code
GEOID	NAD83
NAME	County Name
LSAD	Legal/statistical area description
ALAND	County Land Area in square meters
AWATER	County Water Area in square meters
Geometry	Geometry and geographic information

Table 7: Dataset content

Variable	Units	N.Elements	Range	Source.File
Date	YY-mm-dd	343	From 2018-01-01 to 2018-12-09	EPAair_PM25_NC2018_raw.csv
Site_ID	-	24	-	EPAair_PM25_NC2018_raw.csv
COUNTY	-	21	-	EPAair_PM25_NC2018_raw.csv
Population	People	21	From 5,507 to 1,034,290	https://en.wikipedia.org/
Zone	-	3	Coastal, Piedmont, and Mountains	NC County Maps
PM2_5	ug/m3	6499	From -2.5 to 34.2	EPAair_PM25_NC2018_raw.csv
PM10	ug/m3	926	From 0 to 35	EPAair_PM10_NC2018_raw.csv
TAVG	Fahrenheit	4011	From 11 to 87	NOAA_TAVG_NC2018_raw.csv
Emiss_Dist	meters	24	From 813.5 to 81800.9	Self made
Elevation	meters	24	From 0.04 to 1418.8	Package elevatr

Table 8: Summary of data structure

3 Exploratory Data Analysis and Wrangling

3.1 EPA PM2.5 and PM10 Datasets

Uploading PM2.5 and PM10 2018 raw data files associated with EPA Air dataset and format date column.

```
EPA_AQPM25_NC2018_raw <- read.csv("../Data/Raw/EPAair_PM25_NC2018_raw.csv")
EPA_AQPM10_NC2018_raw <- read.csv("../Data/Raw/EPAair_PM10_NC2018_raw.csv")

#Formatting Dates
EPA_AQPM25_NC2018_raw$Date <- as.Date(EPA_AQPM25_NC2018_raw$Date, format = "%m/%d/%Y")
EPA_AQPM10_NC2018_raw$Date <- as.Date(EPA_AQPM10_NC2018_raw$Date, format = "%m/%d/%Y")
```

Data exploration of the PM2.5 and PM10 2018 raw data files associated with EPA Air dataset.

```
dim(EPA_AQPM25_NC2018_raw)
```

```
## [1] 9644 20
```

```
dim(EPA_AQPM10_NC2018_raw)
```

```
## [1] 2905 20
```

```
str(EPA_AQPM25_NC2018_raw)
```

```
## 'data.frame': 9644 obs. of 20 variables:
## $ Date : Date, format: "2018-01-02" "2018-01-05" ...
## $ Source : Factor w/ 2 levels "AirNow","AQS": 2 2 2 2 2 2 2 2 ...
## $ Site.ID : int 370110002 370110002 370110002 370110002 370110002 ...
## $ 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 ...
## $ 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 ...
## $ AQS_PARAMETER_DESC : Factor w/ 2 levels "Acceptable PM2.5 AQI & Speciat ...
## $ 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 ...
## $ 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 ...
## $ COUNTY_CODE : int 11 11 11 11 11 11 11 11 11 11 ...
## $ COUNTY : Factor w/ 21 levels "Avery", "Buncombe", ...: 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 ...
```

```
str(EPA_AQPM10_NC2018_raw)
```

```
## 'data.frame': 2905 obs. of 20 variables:
## $ Date : Date, format: "2018-01-01" "2018-01-02" ...
## $ Source : Factor w/ 1 level "AQS": 1 1 1 1 1 1 1 1 1 1 ...
## $ Site.ID : int 370510009 370510009 370510009 370510009 370510009 ...
## $ POC : int 3 3 3 3 3 3 3 3 3 3 ...
## $ Daily.Mean.PM10.Concentration: int 3 9 15 12 12 10 9 15 22 15 ...
## $ UNITS : Factor w/ 1 level "ug/m3 SC": 1 1 1 1 1 1 1 1 1 1 ...
## $ DAILY_AQI_VALUE : int 3 8 14 11 11 9 8 14 20 14 ...
## $ Site.Name : Factor w/ 9 levels "Durham Armory", ...: 9 9 9 9 9 9 9 9 ...
## $ 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 81102 81102 81102 81102 81102 81102 81102 81102 ...
## $ AQS_PARAMETER_DESC : Factor w/ 1 level "PM10 Total 0-10um STP": 1 1 1 1 ...
```

```
## $ CBSA_CODE : int 22180 22180 22180 22180 22180 22180 22180 22180 22180
## $ CBSA_NAME : Factor w/ 8 levels "", "Charlotte-Concord-Gastonia,
## $ 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 51 51 51 51 51 51 51 51 51 51 ...
## $ COUNTY : Factor w/ 8 levels "Cumberland", "Durham", ...: 1 1 1
## $ SITE_LATITUDE : num 35 35 35 35 35 ...
## $ SITE_LONGITUDE : num -79 -79 -79 -79 -79 ...
```

```
colnames(EPA_AQPM25_NC2018_raw)
```

```
## [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"
```

```
colnames(EPA_AQPM10_NC2018_raw)
```

```
## [1] "Date" "Source"
## [3] "Site.ID" "POC"
## [5] "Daily.Mean.PM10.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"
```

```
summary(EPA_AQPM25_NC2018_raw)
```

```
##      Date      Source      Site.ID      POC
## Min.   :2018-01-01  AirNow: 873  Min.   :370110002  Min.   :1.000
## 1st Qu.:2018-04-04  AQS    :8771  1st Qu.:370650099  1st Qu.:3.000
## Median :2018-06-27                Median :371190041  Median :3.000
## Mean   :2018-06-30                Mean   :371023866  Mean   :2.948
## 3rd Qu.:2018-09-30                3rd Qu.:371230001  3rd Qu.:3.000
## Max.   :2018-12-31                Max.   :371830021  Max.   :5.000
##
## Daily.Mean.PM2.5.Concentration  UNITS  DAILY_AQI_VALUE
## Min.   : -2.80                  ug/m3 LC:9644  Min.   : 0.00
```

```

## 1st Qu.: 5.00                                1st Qu.:21.00
## Median : 7.20                                Median :30.00
## Mean   : 7.61                                Mean   :31.22
## 3rd Qu.: 9.80                                3rd Qu.:41.00
## Max.   :34.20                                Max.   :97.00
##
##
##      Site.Name      DAILY_OBS_COUNT PERCENT_COMPLETE
## Garinger High School      : 732  Min.      :1      Min.      :100
## Millbrook School          : 722  1st Qu.:1      1st Qu.:100
## Remount                    : 668  Median :1      Median :100
## Montclair Elementary School: 648  Mean    :1      Mean    :100
## Hattie Avenue              : 510  3rd Qu.:1      3rd Qu.:100
## Board Of Ed. Bldg.         : 478  Max.    :1      Max.    :100
## (Other)                    :5886
## AQS_PARAMETER_CODE                                AQS_PARAMETER_DESC
## Min.      :88101      Acceptable PM2.5 AQI & Speciation Mass:2008
## 1st Qu.:88101      PM2.5 - Local Conditions                      :7636
## Median :88101
## Mean    :88184
## 3rd Qu.:88101
## Max.    :88502
##
##      CBSA_CODE                                CBSA_NAME      STATE_CODE
## Min.      :11700      Charlotte-Concord-Gastonia, NC-SC:2048  Min.      :37
## 1st Qu.:16740      Raleigh, NC                                :1418  1st Qu.:37
## Median :24780      Winston-Salem, NC                          :1323  Median :37
## Mean    :29881                                :1165  Mean    :37
## 3rd Qu.:40580      Asheville, NC                              : 532  3rd Qu.:37
## Max.    :49180      Durham-Chapel Hill, NC                     : 469  Max.    :37
## NA's    :1165      (Other)                                    :2689
##
##      STATE      COUNTY_CODE      COUNTY      SITE_LATITUDE
## North Carolina:9644  Min.      : 11.0  Mecklenburg:2048  Min.      :34.36
##                      1st Qu.: 65.0  Wake            :1069  1st Qu.:35.24
##                      Median :119.0  Forsyth         : 876  Median :35.64
##                      Mean    :102.4  Buncombe       : 478  Mean    :35.58
##                      3rd Qu.:123.0  Durham          : 469  3rd Qu.:35.91
##                      Max.    :183.0  Pitt            : 461  Max.    :36.11
##                      (Other)      :4243
## SITE_LONGITUDE
## Min.      :-83.44
## 1st Qu.: -80.87
## Median   : -80.23
## Mean     : -80.03
## 3rd Qu.: -78.82
## Max.     : -76.21

```

##

summary(EPA_AQPM10_NC2018_raw)

```
##          Date          Source      Site.ID          POC
##  Min.    :2018-01-01    AQS:2905  Min.    :370510009  Min.    :1.000
##  1st Qu.:2018-04-07          1st Qu.:370670022  1st Qu.:3.000
##  Median :2018-07-02          Median :371170001  Median :3.000
##  Mean   :2018-07-03          Mean   :371072712  Mean   :3.172
##  3rd Qu.:2018-10-04          3rd Qu.:371190042  3rd Qu.:4.000
##  Max.   :2018-12-31          Max.   :371830014  Max.   :5.000
##
##  Daily.Mean.PM10.Concentration  UNITS      DAILY_AQI_VALUE
##  Min.    : 0.00                ug/m3 SC:2905  Min.    : 0.00
##  1st Qu.:10.00                1st Qu.: 9.00
##  Median :13.00                Median :12.00
##  Mean   :13.72                Mean   :12.67
##  3rd Qu.:17.00                3rd Qu.:16.00
##  Max.   :64.00                Max.   :55.00
##
##          Site.Name      DAILY_OBS_COUNT PERCENT_COMPLETE
##  Millbrook School      :588  Min.    :1      Min.    :100
##  Garinger High School  :351  1st Qu.:1      1st Qu.:100
##  Montclair Elementary School:344  Median :1      Median :100
##  Hattie Avenue         :342  Mean   :1      Mean   :100
##  Durham Armory          :335  3rd Qu.:1      3rd Qu.:100
##  William Owen School   :321  Max.   :1      Max.   :100
##  (Other)                :624
##  AQS_PARAMETER_CODE      AQS_PARAMETER_DESC      CBSA_CODE
##  Min.    :81102          PM10 Total 0-10um STP:2905  Min.    :16740
##  1st Qu.:81102          1st Qu.:16740
##  Median :81102          Median :22180
##  Mean   :81102          Mean   :28310
##  3rd Qu.:81102          3rd Qu.:39580
##  Max.   :81102          Max.   :49180
##                          NA's    :247
##          CBSA_NAME      STATE_CODE
##  Charlotte-Concord-Gastonia, NC-SC:695  Min.    :37
##  Raleigh, NC                    :588  1st Qu.:37
##  Winston-Salem, NC              :342  Median :37
##  Durham-Chapel Hill, NC         :335  Mean   :37
##  Fayetteville, NC               :321  3rd Qu.:37
##  Greensboro-High Point, NC      :320  Max.   :37
##  (Other)                        :304
##          STATE      COUNTY_CODE      COUNTY      SITE_LATITUDE
```

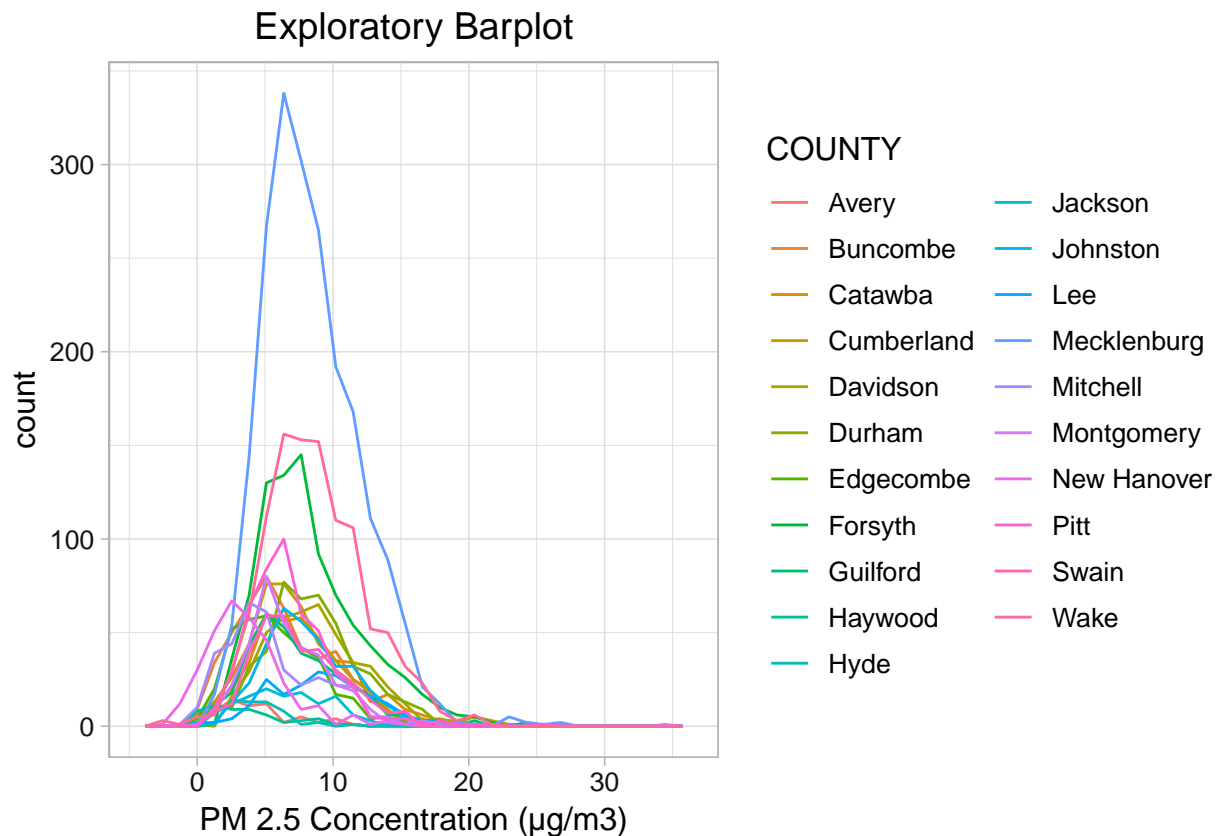


Figure 1: PM2.5 NC 2018 frequency polygon.

```
## North Carolina:2905    Min.    : 51.0    Mecklenburg:695    Min.    :35.04
##                        1st Qu.: 67.0    Wake       :588    1st Qu.:35.24
##                        Median :117.0    Forsyth    :342    Median :35.86
##                        Mean   :107.3    Durham     :335    Mean   :35.67
##                        3rd Qu.:119.0    Cumberland :321    3rd Qu.:36.03
##                        Max.   :183.0    Guilford   :320    Max.   :36.11
##                                (Other)  :304
##
## SITE_LONGITUDE
## Min.    :-80.87
## 1st Qu. :-80.23
## Median  :-78.95
## Mean    :-79.36
## 3rd Qu. :-78.57
## Max.    :-76.91
##
```

Visual data exploration of the PM2.5 2018 raw data file in Figure 1, Figure 2, and Figure 3.

Visual data exploration of the PM10 2018 raw data file in Figure 4.

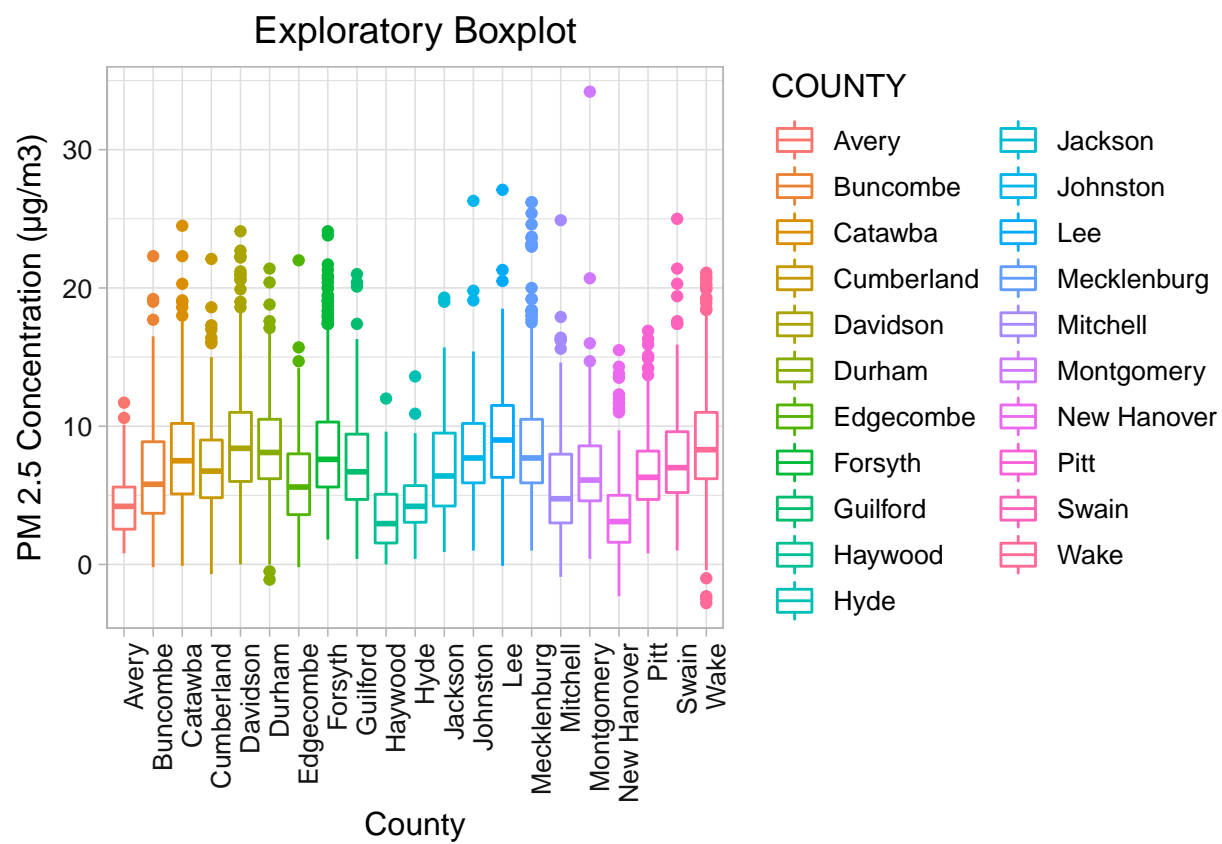


Figure 2: PM2.5 NC 2018 boxplot.

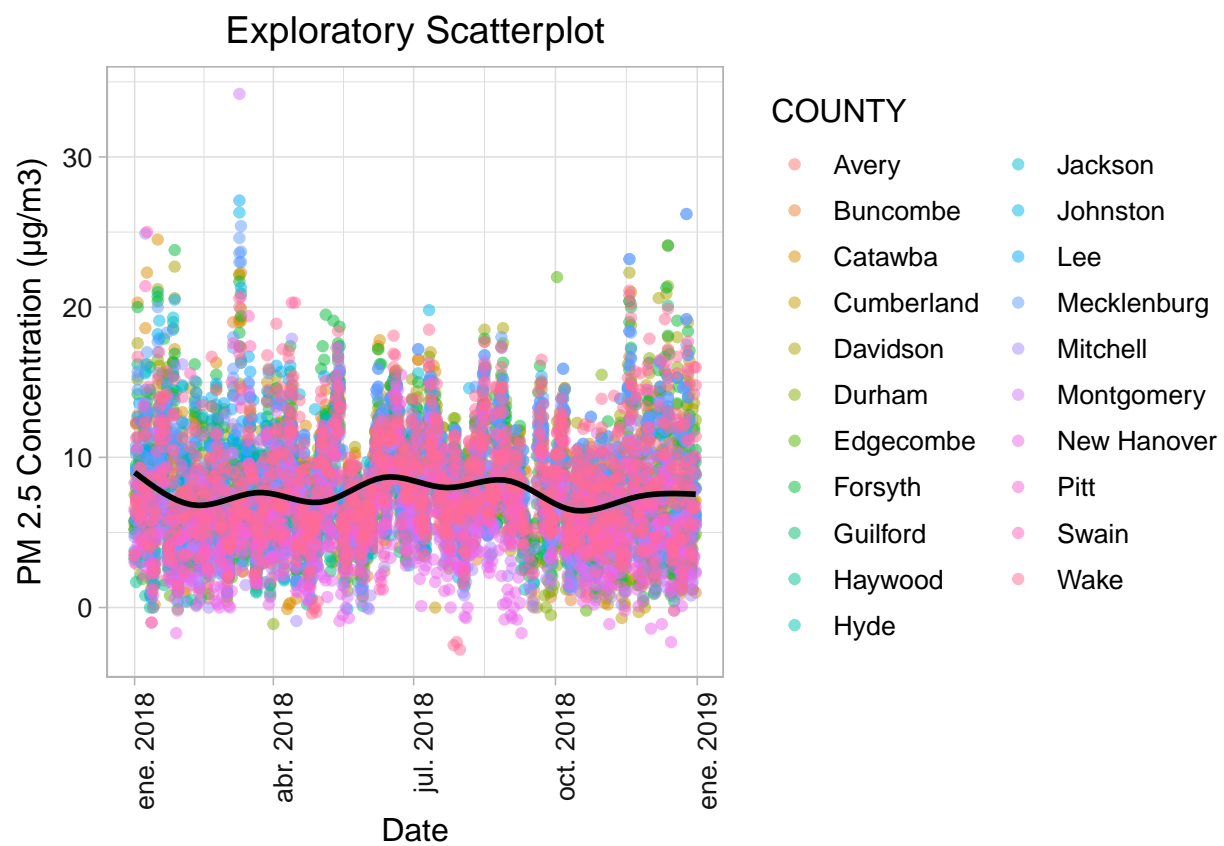


Figure 3: PM_{2.5} NC 2018 scatterplot.

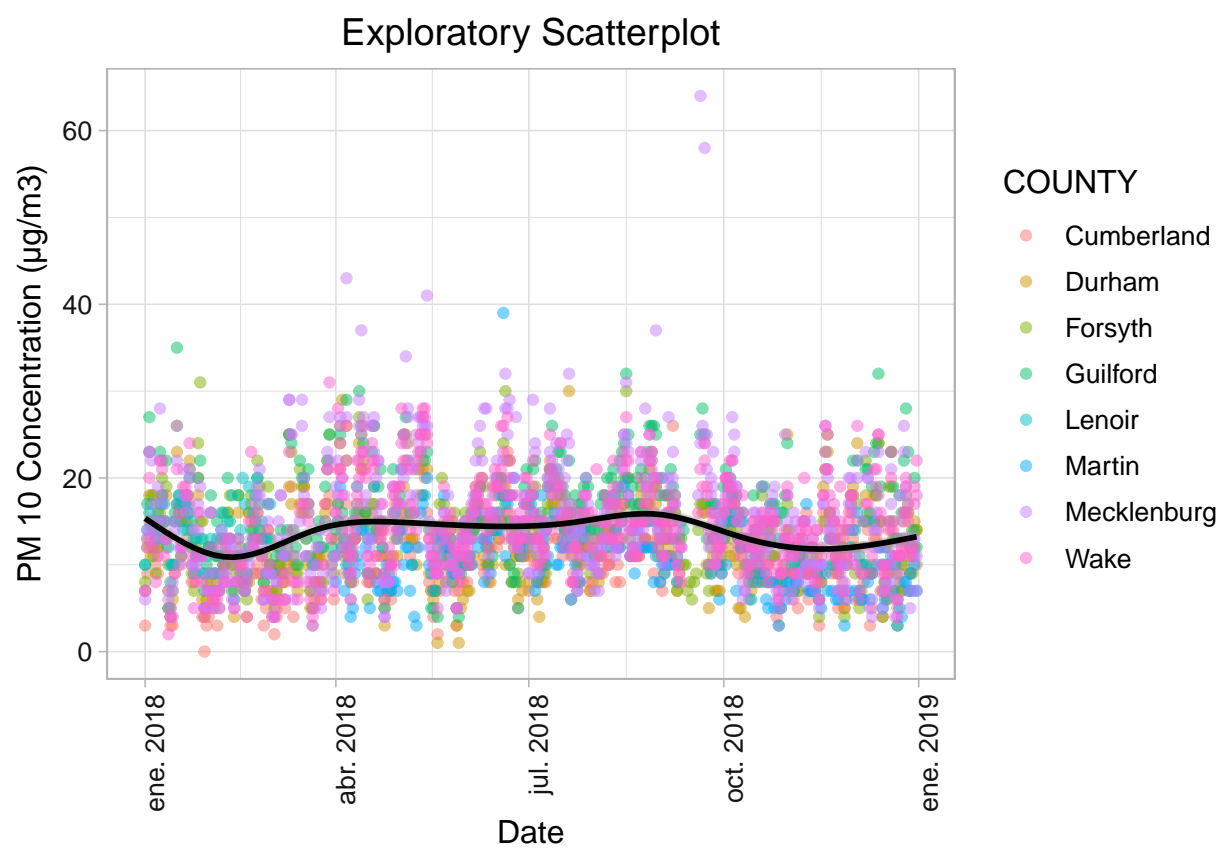


Figure 4: PM10 NC 2018 scatterplot.

Data wrangling of the PM2.5 and PM10 2018 raw data files.

```
#Selecting Columns
EPA_AQ_PM25_NC2018_Temp <- select(EPA_AQPM25_NC2018_raw, Date, Site.ID,
                                Daily.Mean.PM2.5.Concentration, AQS_PARAMETER_DESC,
                                COUNTY:SITE_LONGITUDE)

#Changing column name
colnames(EPA_AQ_PM25_NC2018_Temp)[colnames(EPA_AQ_PM25_NC2018_Temp)
                                == "Daily.Mean.PM2.5.Concentration"] <- "Daily.Mean.Concentr"

#Selecting Columns
EPA_AQ_PM10_NC2018_Temp <- select(EPA_AQPM10_NC2018_raw, Date, Site.ID,
                                Daily.Mean.PM10.Concentration, AQS_PARAMETER_DESC,
                                COUNTY:SITE_LONGITUDE)

#Changing column name
colnames(EPA_AQ_PM10_NC2018_Temp)[colnames(EPA_AQ_PM10_NC2018_Temp)
                                == "Daily.Mean.PM10.Concentration"] <- "Daily.Mean.Concentr"

#Create AQS_PARAMETER_DESC Column with Contaminant description.
EPA_AQ_PM25_NC2018_Temp$AQS_PARAMETER_DESC <- "PM2.5"
EPA_AQ_PM10_NC2018_Temp$AQS_PARAMETER_DESC <- "PM10"

#Eliminates duplicate dates
EPA_AQ_PM25_NC2018_Cleaned <- EPA_AQ_PM25_NC2018_Temp [!duplicated(EPA_AQ_PM25_NC2018_Temp)]

EPA_AQ_PM10_NC2018_Cleaned <- EPA_AQ_PM10_NC2018_Temp [!duplicated(EPA_AQ_PM10_NC2018_Temp)]

# Combine the data.
EPA_AQ_PM2.5PM10_NC2018_Cleaned <- rbind(EPA_AQ_PM25_NC2018_Cleaned, EPA_AQ_PM10_NC2018_Cleaned)

#Save the data in the processed folder
write.csv(EPA_AQ_PM2.5PM10_NC2018_Cleaned,
          "./Data/Processed/EPA_AQ_PM2.5PM10_NC2018_Cleaned.csv")

#Spread PM2.5 and PM10
EPA_AQ_PM2.5PM10_NC2018_Spread <-
  EPA_AQ_PM2.5PM10_NC2018_Cleaned %>%
  spread(AQS_PARAMETER_DESC, Daily.Mean.Concentration)

#Remove rows without PM2.5 data
EPA_AQ_PM2.5PM10_NC2018_Spread <- EPA_AQ_PM2.5PM10_NC2018_Spread[!is.na(EPA_AQ_PM2.5PM10_NC2018_Spread$Daily.Mean.PM2.5.Concentration)]

#Convert the dataset to a spatially enabled "sf" data frame
```

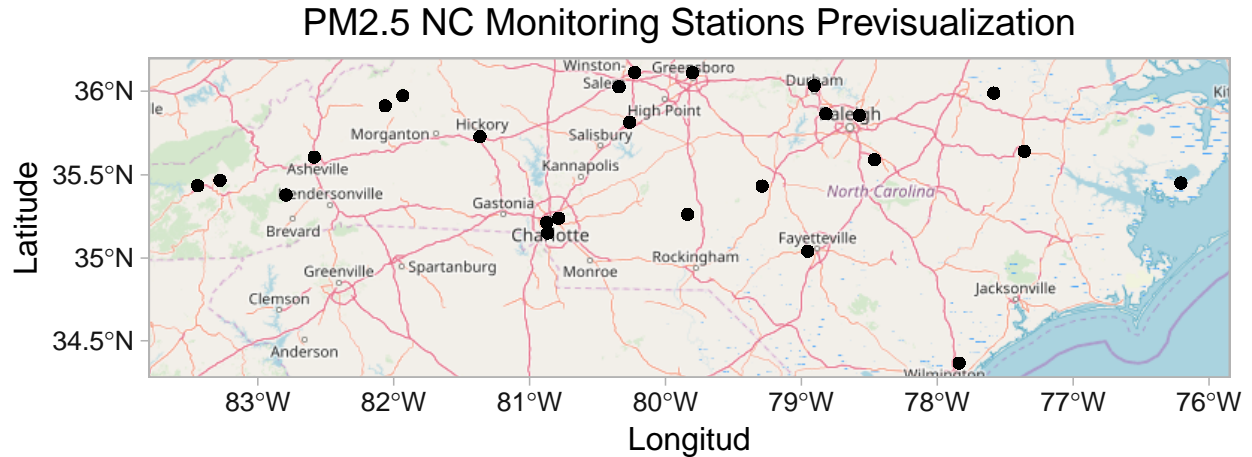


Figure 5: PM2.5 NC Monitoring Stations Previsualization.

```
PM2.5_PM10_sf <- st_as_sf(EPA_AQ_PM2.5PM10_NC2018_Spread, coords = c('SITE_LONGITUDE', 'SITE_LATITUDE'))
#Convert all to UTM Zone 17 (crs = 26917)
PM2.5_PM10_sf_utm <- st_transform(PM2.5_PM10_sf, c=26917)
```

In Figure 5 is presented the locations of the PM2.5 monitoring stations.

3.2 North Carolina Counties Zoning, Geographic information, and Population Data

Downloading the list of North Carolina Counties and Population from a Wikipedia URL.

```
#North Carolina Counties
url <- "https://en.wikipedia.org/wiki/List_of_counties_in_North_Carolina"
webpage <- read_html(url)

County_Name <- webpage %>% html_nodes("th:nth-child(1)") %>% html_text()
County_Population <- webpage %>% html_nodes("tr:nth-child(7)") %>% html_text()
```

```

#Remove unwanted info and characters
County_Info <- data_frame(County = County_Name[9:108])
County_Info$County <- str_replace(County_Info$County, " County", "")
County_Info$County <- str_replace(County_Info$County, "\n", "")

Population <- data_frame(Population=County_Population[2:101])

County_Info <- cbind(County_Info, Population)

County_Info$Population <- str_replace(County_Info$Population,",", "")
County_Info$Population <- str_replace(County_Info$Population,",", "")

County_Info$Population <- as.numeric(County_Info$Population)

```

Assigning the corresponding zone to each county. Info from: Rudersdorf, Amy. 2010. "NC County Maps." Government & Heritage Library, State Library of North Carolina.

```

#North Carolina Zones
County_Info$Zone<-ifelse(County_Info$County == 'Ashe'|County_Info$County =='Alleghany'|C
                        ifelse(County_Info$County == 'Surry'|County_Info$County =='Stok
                        ifelse(County_Info$County == 'Scotland'|County_Info$Cour

```

Data exploration of the County_Info dataframe.

```
dim(County_Info)
```

```
## [1] 100 3
```

```
str(County_Info)
```

```

## 'data.frame': 100 obs. of 3 variables:
## $ County : chr "Alamance" "Alexander" "Alleghany" "Anson" ...
## $ Population: num 157844 37159 10935 25531 26833 ...
## $ Zone : chr "Piedmont" "Piedmont" "Mountains" "Piedmont" ...

```

```
colnames(County_Info)
```

```
## [1] "County" "Population" "Zone"
```

```
summary(County_Info)
```

```

##      County      Population      Zone
## Length:100      Min.   : 4090      Length:100
## Class :character 1st Qu.: 25001      Class :character
## Mode  :character Median : 55311      Mode  :character
##                      Mean  : 100526
##                      3rd Qu.: 114764
##                      Max.   :1034290

```

```
unique(County_Info$Zone)
```

```
## [1] "Piedmont" "Mountains" "Coastal"
```

Adding the County Information to the PM2.5_PM10_sf_utm dataframe.

```
PM2.5_PM10_Info_sf_utm <- PM2.5_PM10_sf_utm %>%  
left_join(County_Info, by = c("COUNTY"="County"))
```

3.3 US Census Bureau US counties shapefile

Reading the USA county shapefile, sub-setting for NC.

```
counties_sf<- st_read('./Data/Spatial/cb_2017_us_county_20m.shp') %>%  
filter(STATEFP == 37) #Filter for just NC Counties
```

```
## Reading layer `cb_2017_us_county_20m' from data source `C:\Users\Felipe\OneDrive - Du  
## Simple feature collection with 3220 features and 9 fields  
## geometry type: MULTIPOLYGON  
## dimension: XY  
## bbox: xmin: -179.1743 ymin: 17.91377 xmax: 179.7739 ymax: 71.35256  
## epsg (SRID): 4269  
## proj4string: +proj=longlat +datum=NAD83 +no_defs
```

```
#CRS
```

```
st_crs(counties_sf) #crs=4269 = NAD83.
```

```
## Coordinate Reference System:
```

```
## EPSG: 4269
```

```
## proj4string: "+proj=longlat +datum=NAD83 +no_defs"
```

Converting the counties_sf to UTM Zone 17

```
#Convert all to UTM Zone 17 (crs = 26917)
```

```
counties_sf_utm <- st_transform(counties_sf, c=26917)
```

```
#Adding the Zone Info
```

```
counties_sf_utm <- counties_sf_utm %>%
```

```
left_join(County_Info, by = c("NAME"="County"))
```

Data exploration of the County_Info dataframe.

```
dim(counties_sf_utm)
```

```
## [1] 100 12
```

```
str(counties_sf_utm)
```

```
## Classes 'sf' and 'data.frame': 100 obs. of 12 variables:
## $ STATEFP : Factor w/ 52 levels "01","02","04",...: 34 34 34 34 34 34 34 34 34 34 .
## $ COUNTYFP : Factor w/ 325 levels "001","003","005",...: 116 71 94 54 62 125 84 102
## $ COUNTYNS : Factor w/ 3220 levels "00023901","00025441",...: 1672 1649 1659 1639 16
## $ AFFGEOID : Factor w/ 3220 levels "0500000US01001",...: 1981 1944 1964 1931 1937 19
## $ GEOID : Factor w/ 3220 levels "01001","01003",...: 1981 1944 1964 1931 1937 198
## $ NAME : chr "Vance" "Lenoir" "Pitt" "Guilford" ...
## $ LSAD : Factor w/ 9 levels "00","03","04",...: 5 5 5 5 5 5 5 5 5 5 ...
## $ ALAND : num 6.54e+08 1.03e+09 1.69e+09 1.67e+09 1.01e+09 ...
## $ AWATER : num 42187365 5900300 8248766 30723331 3981006 ...
## $ Population: num 44420 57934 176484 517197 52571 ...
## $ Zone : chr "Piedmont" "Coastal" "Coastal" "Piedmont" ...
## $ geometry :sfc_MULTIPOLYGON of length 100; first list element: List of 1
## ..$ :List of 1
## .. ..$ : num [1:10, 1:2] 724066 727615 739570 744478 741926 ...
## ..- attr(*, "class")= chr "XY" "MULTIPOLYGON" "sfg"
## - attr(*, "sf_column")= chr "geometry"
## - attr(*, "agr")= Factor w/ 3 levels "constant","aggregate",...: NA NA NA NA NA NA NA
## ..- attr(*, "names")= chr "STATEFP" "COUNTYFP" "COUNTYNS" "AFFGEOID" ...
```

```
colnames(counties_sf_utm)
```

```
## [1] "STATEFP" "COUNTYFP" "COUNTYNS" "AFFGEOID" "GEOID"
## [6] "NAME" "LSAD" "ALAND" "AWATER" "Population"
## [11] "Zone" "geometry"
```

```
summary(counties_sf_utm)
```

```
## STATEFP COUNTYFP COUNTYNS AFFGEOID GEOID
## 37 :100 001 : 1 01008531: 1 0500000US37001: 1 37001 : 1
## 01 : 0 003 : 1 01008532: 1 0500000US37003: 1 37003 : 1
## 02 : 0 005 : 1 01008533: 1 0500000US37005: 1 37005 : 1
## 04 : 0 007 : 1 01008534: 1 0500000US37007: 1 37007 : 1
## 05 : 0 009 : 1 01008535: 1 0500000US37009: 1 37009 : 1
## 06 : 0 011 : 1 01008536: 1 0500000US37011: 1 37011 : 1
## (Other): 0 (Other):94 (Other) :94 (Other) :94 (Other):94
## NAME LSAD ALAND AWATER
## Length:100 06 :100 Min. :4.472e+08 Min. :4.453e+05
## Class :character 00 : 0 1st Qu.:8.936e+08 1st Qu.:7.287e+06
## Mode :character 03 : 0 Median :1.192e+09 Median :1.595e+07
## 04 : 0 Mean :1.259e+09 Mean :1.347e+08
## 05 : 0 3rd Qu.:1.501e+09 3rd Qu.:3.831e+07
## 12 : 0 Max. :2.453e+09 Max. :3.001e+09
## (Other): 0
## Population Zone geometry
## Min. : 4090 Length:100 MULTIPOLYGON :100
```

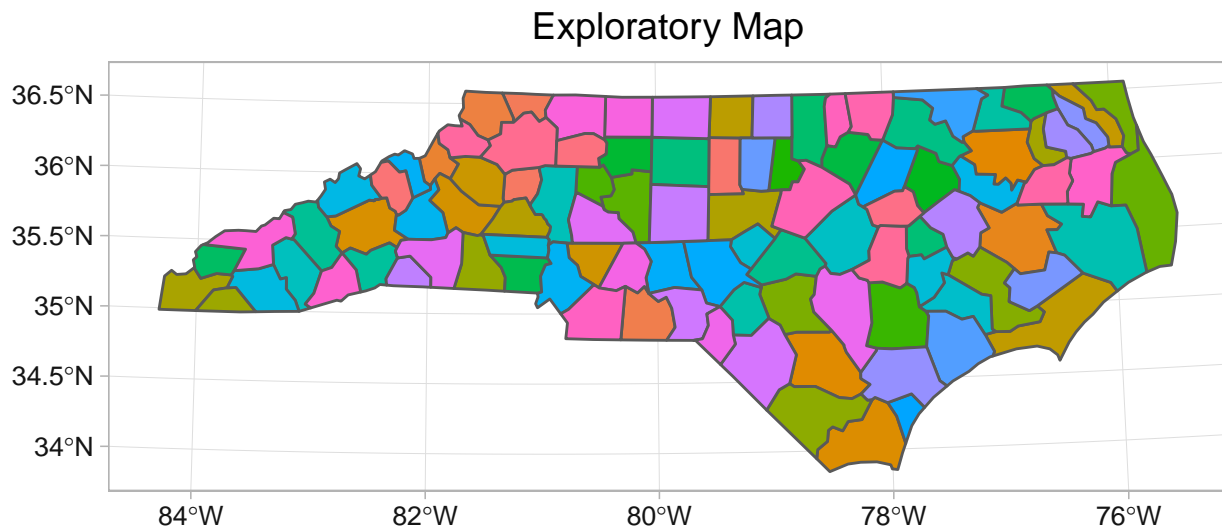



Figure 6: Counties exploratory map.

```
## 1st Qu.: 25001   Class :character   epsg:26917   : 0
## Median : 55311   Mode  :character   +proj=utm ...: 0
## Mean   : 100526
## 3rd Qu.: 114764
## Max.   :1034290
##
```

Visual data exploration of the `counties_sf_utm` dataframe in Figure 6.

3.4 NOAA Average Temperature Dataset

Reading the 2018 North Carolina Air Temperature data.

```
#Read the 2018 Air Temperature data
NOAA_DTAVG_NC2018_raw <- read.csv("../Data/Raw/NOAA_TAVG_NC2018_raw.csv")
```

Data exploration of the `NOAA_DTAVG_NC2018_raw` dataframe.

```
dim(NOAA_DTAVG_NC2018_raw)
```

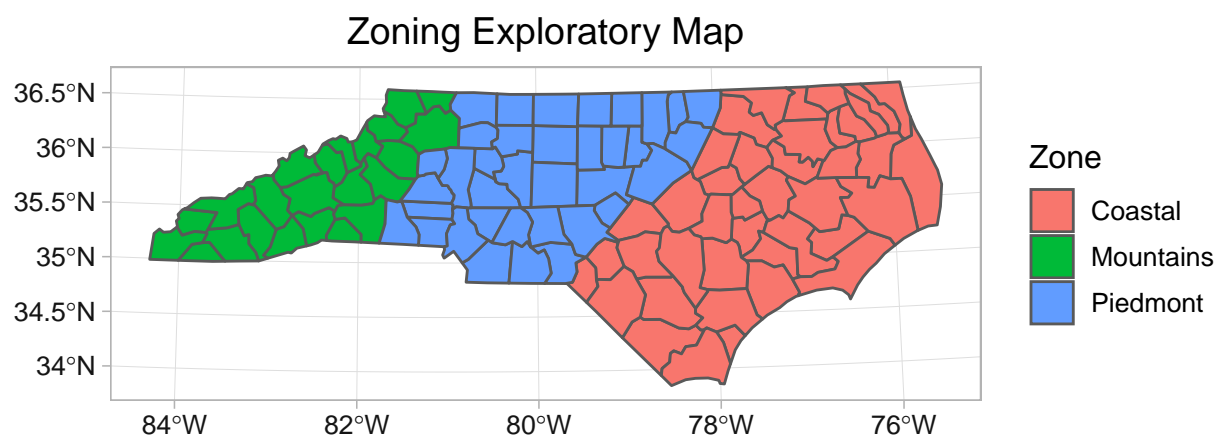


Figure 7: NC Zoning exploratory map.

```
## [1] 283423      7
```

```
str(NOAA_DTAVG_NC2018_raw)
```

```
## 'data.frame':    283423 obs. of  7 variables:
## $ STATION : Factor w/ 1066 levels "US1NCAG0001",...: 217 217 217 217 217 217 217 217
## $ NAME    : Factor w/ 1063 levels "ABERDEEN 0.7 NW, NC US",...: 716 716 716 716 716
## $ LATITUDE: num  34.8 34.8 34.8 34.8 34.8 ...
## $ LONGITUDE: num -76.9 -76.9 -76.9 -76.9 -76.9 ...
## $ ELEVATION: num  8.5 8.5 8.5 8.5 8.5 8.5 8.5 8.5 8.5 8.5 ...
## $ DATE     : Factor w/ 365 levels "01/01/2018","01/02/2018",...: 37 133 145 205 265 2
## $ TAVG     : int  NA NA NA NA NA NA NA NA NA NA ...
```

```
colnames(NOAA_DTAVG_NC2018_raw)
```

```
## [1] "STATION" "NAME" "LATITUDE" "LONGITUDE" "ELEVATION" "DATE"
## [7] "TAVG"
```

```
summary(NOAA_DTAVG_NC2018_raw)
```

```
##          STATION          NAME          LATITUDE
## US1NCBC0005: 365 SPARTA 3.5 SSW, NC US : 545 Min. :33.88
## US1NCBC0041: 365 HILLSBOROUGH 5.6 NNW, NC US: 502 1st Qu.:35.16
## US1NCBK0004: 365 ADVANCE 0.2 ESE, NC US : 365 Median :35.56
## US1NCCH0004: 365 ALBEMARLE, NC US : 365 Mean :35.49
## US1NCCS0002: 365 ARDEN 1.6 ENE, NC US : 365 3rd Qu.:35.90
## US1NCCY0003: 365 ASHEBORO 1.3 SSE, NC US : 365 Max. :36.56
## (Other) :281233 (Other) :280916
## LONGITUDE ELEVATION DATE TAVG
## Min. : -84.30 Min. : 0.0 16/04/2018: 887 Min. : 9.00
## 1st Qu.: -81.67 1st Qu.: 29.3 17/05/2018: 874 1st Qu.:47.00
## Median : -79.16 Median : 150.9 20/03/2018: 871 Median :63.00
## Mean : -79.70 Mean : 279.9 12/06/2018: 870 Mean :60.19
## 3rd Qu.: -78.01 3rd Qu.: 389.5 30/05/2018: 869 3rd Qu.:75.00
## Max. : -75.46 Max. :1902.0 01/08/2018: 867 Max. :87.00
## (Other) :278185 NA's :269572
```

Data wrangling of the NOAA_DTAVG_NC2018_raw dataframe.

```
#Remove stations without Temperature information
```

```
NOAA_DTAVG_NC2018_Complete <- na.omit(NOAA_DTAVG_NC2018_raw)
```

```
#Convert the dataset to a spatially enabled "sf" data frame
```

```
NOAA_DTAVG_NC2018_sf <- st_as_sf(NOAA_DTAVG_NC2018_Complete, coords = c('LONGITUDE', 'LATITUDE'))
```

```
#Convert all to UTM Zone 17 (crs = 26917)
```

```
NOAA_DTAVG_NC2018_sf_utm <- st_transform(NOAA_DTAVG_NC2018_sf, crs=26917)
```

```
#Formatting dates
```

```
NOAA_DTAVG_NC2018_sf_utm$DATE <- as.Date(NOAA_DTAVG_NC2018_sf_utm$DATE, format = "%d/%m/%Y")
```

The 2018 Air Temperature data does not have County information, so the location is used with the counties_sf_utm dataframe to locate the county of each station.

```
#Adding the county and zone information to the Temperature dataframe
```

```
#Index of the matching feature
```

```
county_index <- st_nearest_feature(NOAA_DTAVG_NC2018_sf_utm, counties_sf_utm)
```

```
#Eliminates geo info
```

```
aux1 <- st_set_geometry(counties_sf_utm[county_index,"NAME"], value=NULL)
```

```
#adds the columns
```

```
NOAA_DTAVG_NC2018_sf_utm$COUNTY <- aux1$NAME
```

```
#Reordering
```

```
NOAA_DTAVG_NC2018_sf_utm <- NOAA_DTAVG_NC2018_sf_utm[,c(1,2,3,4,5,7,6)]
```

Visual data exploration of the 2018 North Carolina Air Temperature data in Figure 8, Figure 9, Figure 10, and Figure 11..

Next, the temperature of the nearest Temperature Station is added to each PM2.5 Station in the PM2.5_PM10_Info_sf_utm dataframe.

```
#Create a Data frame with only the PM2.5 station info
```

```
PM2.5_Stations <- PM2.5_PM10_Info_sf_utm %>%
```

```
  select(Site.ID, geometry) %>%
```

```
  subset(!duplicated(Site.ID))
```

```
#Distances bewteen the PM2.5 stations and the Temperature Stations
```

```
Nearest <- st_nearest_feature(PM2.5_Stations, NOAA_DTAVG_NC2018_sf_utm)
```

```
a <- length(unique(PM2.5_Stations$Site.ID))
```

```
NOAA_DTAVG_NC2018_sf_utm$NAME <- as.character(NOAA_DTAVG_NC2018_sf_utm$NAME)
```

```
#Assingning the nearest Temperature Station to each PM2.5 station.
```

```
for (i in 1:a){
```

```
  PM2.5_Stations$Temp_Est[i] <- NOAA_DTAVG_NC2018_sf_utm$NAME[Nearest[i]]
```

```
}
```

```
#Drop the geo data
```

```
aux2 <- st_set_geometry(PM2.5_Stations, value=NULL)
```

NC 2018 Temperature Exploratory Map, Mean Annual Temperature

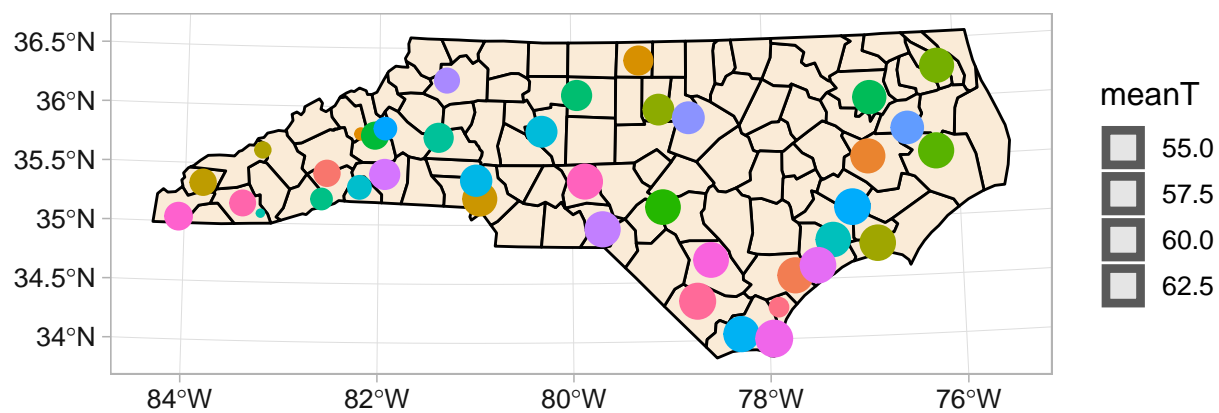


Figure 8: Mean Annual Temperature exploratory map.

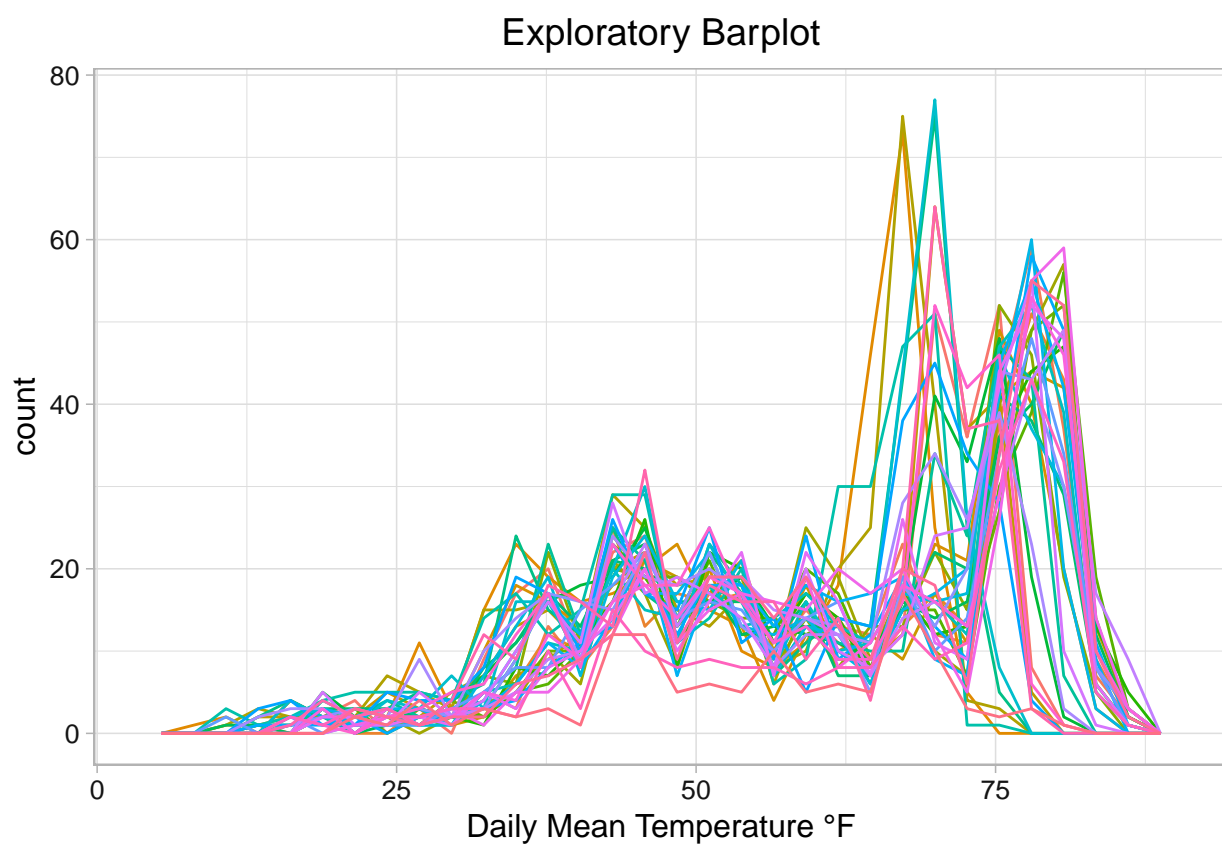


Figure 9: Daily Mean Temperature NC 2018 frequency polygon.

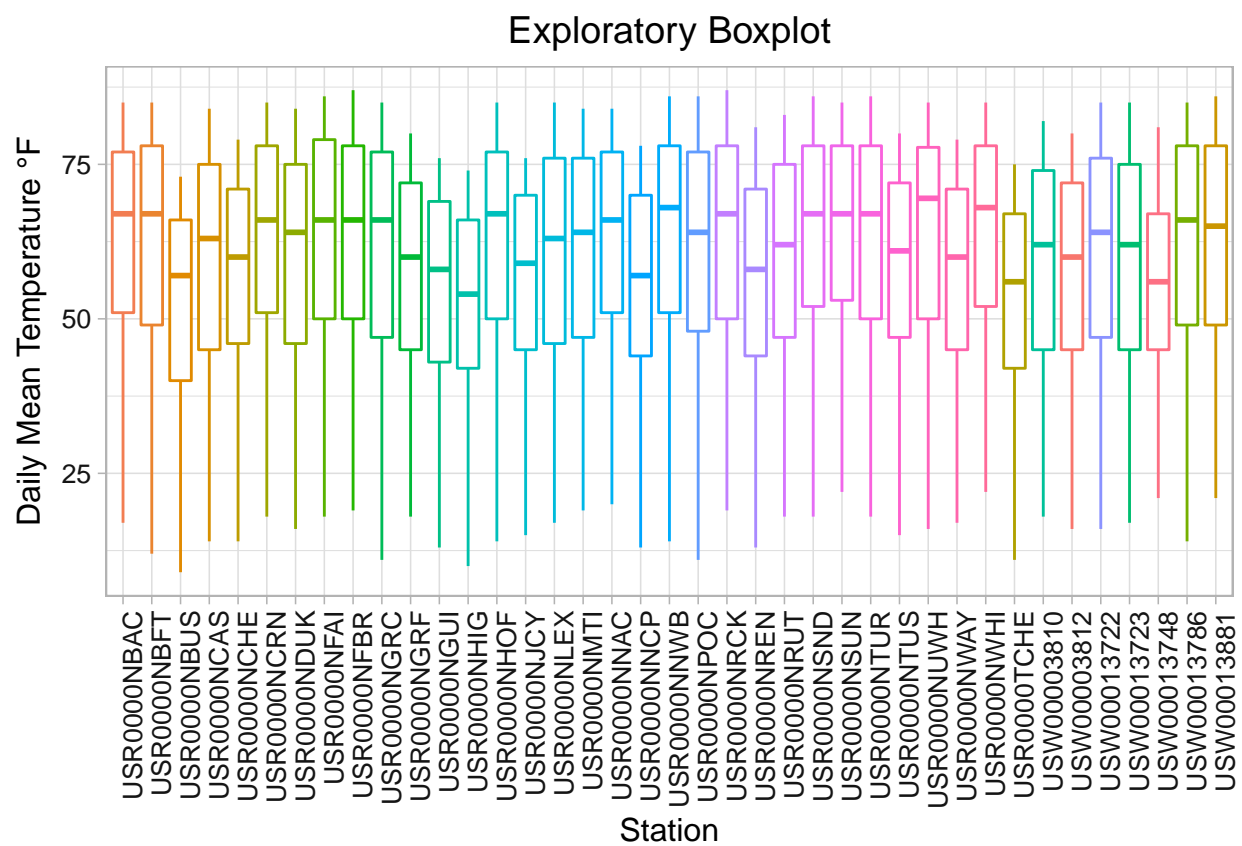


Figure 10: Daily Mean Temperature NC 2018 boxplot.

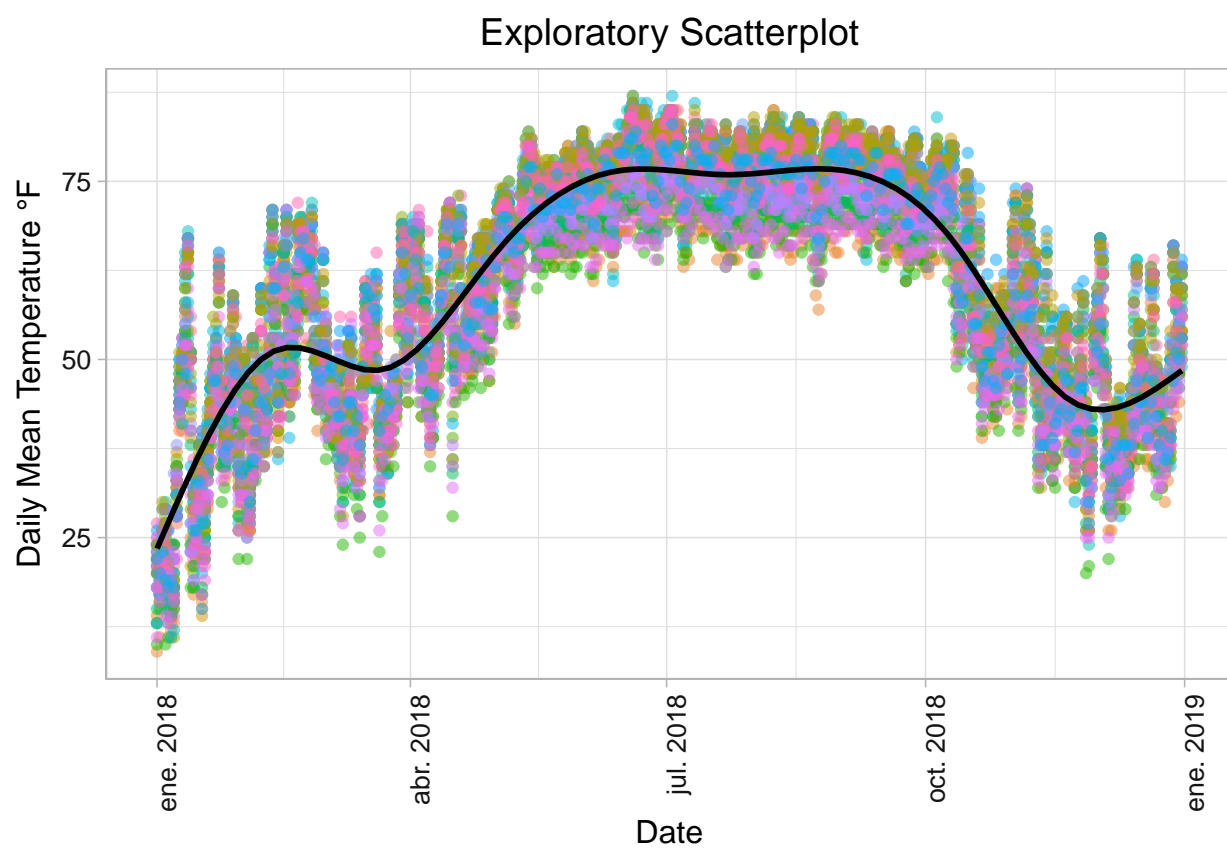


Figure 11: Daily Mean Temperature NC 2018 scatterplot.


```

#Left_join the data
PM2.5_PM10_Temp_sf_utm <- PM2.5_PM10_Info_sf_utm %>%
left_join(aux2)

#Assingning the Temperature of the nearest Temperature Station to each PM2.5 station.

#Drops the geo data
aux3 <- st_set_geometry(NOAA_DTAVG_NC2018_sf_utm, value=NULL)

#Left_join the data
PM2.5_PM10_Temp_sf_utm <- PM2.5_PM10_Temp_sf_utm %>%
left_join(aux3, by = c("Temp_Est"="NAME", "Date"="DATE", "COUNTY")) %>%
select(Date,Site.ID,COUNTY,Population,Zone,PM2.5,PM10,TAVG,geometry)

```

3.5 EPA combustion points for electricity generation in the US Dataset

Reading the Electricity Generation via Combustion data.

```
EPA_US_CombEmissions <- st_read("./Data/Raw/EPA_ElecGenComb_US_raw.kml")
```

```

## Reading layer `Electricity Generation via Combustion' from data source `C:\Users\Feli
## Simple feature collection with 2042 features and 2 fields
## geometry type:  POINT
## dimension:      XYZ
## bbox:           xmin: -176.6593 ymin: 19.63283 xmax: -67.00325 ymax: 71.29221
## epsg (SRID):    4326
## proj4string:     +proj=longlat +datum=WGS84 +no_defs

```

Wrangling the data

```
st_crs(EPA_US_CombEmissions) #crs=4326 = WGS 84
```

```

## Coordinate Reference System:
##   EPSG: 4326
##   proj4string: "+proj=longlat +datum=WGS84 +no_defs"

```

```
#Convert all to UTM Zone 17 (crs = 26917)
```

```
EPA_US_CombEmissions_utm <- st_transform(EPA_US_CombEmissions, c=26917)
```

```
#Clip the EPA_US_CombEmissions data set by the NC State boundary dataset
```

```
#First create a State_sf file
```

```
#Aggregate the data using group_by and summarize, just as you would a non-spatial data
```

```
state_sf_utm <- st_union(counties_sf_utm)
```

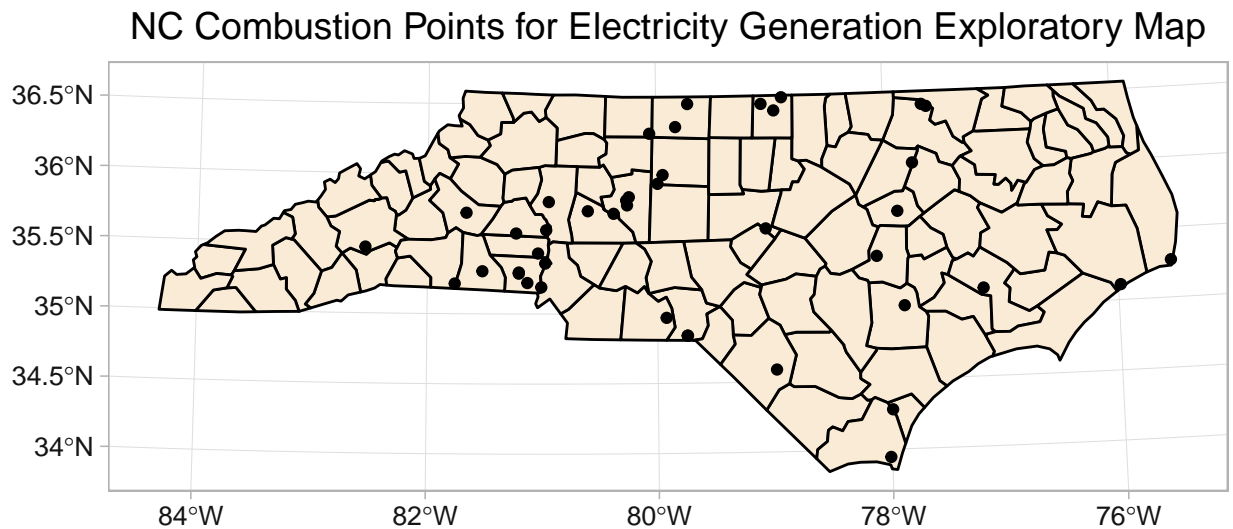


Figure 12: Combustion points for electricity generation in the North Carolina.

```
#Eliminate the emission points outside NC
EPA_NC_CombEmissions_utm <- st_intersection(EPA_US_CombEmissions_utm,state_sf_utm)
```

Visual data exploration of the EPA combustion points for electricity generation in the North Carolina in Figure 8, Figure 9, Figure 12, and Figure 11..

Now the distance between PM2.5 stations and Electricity Generation via Combustion points is determined and added to the PM2.5_PM10_Temp_sf_utm dataframe.

```
#Distances between PM2.5 stations and Electricity Generation via Combustion points
Distances <- st_distance(PM2.5_Stations, EPA_NC_CombEmissions_utm)

a <- length(unique(PM2.5_Stations$Site.ID))

#Determining the minimum distance of each PM2.5 station to a combustion point in meters
for (i in 1:a){
  PM2.5_Stations$Emiss_Dist[i] <- min(Distances[i,])
}
```

```
#Filling the PM2.5_PM10_Temp_sf_utm file with the distances
```

```
#Drops the geo data
```

```
aux4 <- PM2.5_Stations %>%  
  st_set_geometry(value=NULL) %>%  
  select(Site.ID,Emiss_Dist)
```

```
#Left_join the data
```

```
PM2.5_Full_utm <- PM2.5_PM10_Temp_sf_utm %>%  
left_join(aux4, by = c("Site.ID")) %>%  
  select(Date,Site.ID,COUNTY,Population,Zone,PM2.5,PM10,TAVG,Emiss_Dist,geometry)
```

Finally, using the `elevatr` package, elevation information is added to the PM 2.5 stations in the `PM2.5_Full_utm` dataframe, creating the `PM2.5_Full_Elev_utm`, which is saved in the Project folder `./Data/Processed`.

Elevations for the PM2.5 Stations

```
prj_dd <- "+proj=utm +zone=17 +ellps=GRS80 +towgs84=0,0,0,0,0,0,0 +units=m +no_defs"  
PM2.5_Full_Elev_utm <- get_elev_point(PM2.5_Full_utm, prj = prj_dd, src = "epqs")
```

```
st_write(PM2.5_Full_Elev_utm,  
  "./Data/Processed/PM2.5_Full_Elev_utm.shp", driver = "ESRI Shapefile")
```

```
PM2.5_Full_Elev_utm <- st_read('./Data/Processed/PM2.5_Full_Elev_utm.shp')
```

```
## Reading layer `PM2.5_Full_Elev_utm' from data source `C:\Users\Felipe\OneDrive - Duke  
## Simple feature collection with 7460 features and 11 fields  
## geometry type:  POINT  
## dimension:      XY  
## bbox:           xmin: 278314.3 ymin: 3807066 xmax: 935107.5 ymax: 3996703  
## epsg (SRID):    NA  
## proj4string:     +proj=utm +zone=17 +ellps=GRS80 +units=m +no_defs
```

3.6 Additional previsualization of the data

4 Analysis

In 2012, the United States Environmental Protection Agency (USEPA) established two complementary primary regulatory standards for PM_{2.5}. The first is based on a yearly average value and is set at 12 micrograms per cubic meter, ug/m³,

First statistical test should look at the standard in each station.

FRA: This model structure should look familiar, with a typical linear model structure and dataframe defined. The addition here is that we have defined Week as a random variable. Essentially, we are interested not in the specific effects of each week but in the variability among weeks, so we have defined it as a random effect (essentially coming from a larger distribution of seasonal variability).

5 Summary and Conclusions