# Assignment 7: Time Series Analysis

## Emily Wood

## OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics on time series analysis.

## Directions

1. Change "Student Name" on line 3 (above) with your name.
2. Work through the steps, **creating code and output** that fulfill each instruction.
3. Be sure to **answer the questions** in this assignment document.
4. When you have completed the assignment, **Knit** the text and code into a single PDF file.
5. After Knitting, submit the completed exercise (PDF file) to the dropbox in Sakai. Add your last name into the file name (e.g., "Fay_A07_TimeSeries.Rmd") prior to submission.

The completed exercise is due on Tuesday, March 16 at 11:59 pm.

## Set up

1. Set up your session:

- Check your working directory
- Load the tidyverse, lubridate, zoo, and trend packages
- Set your ggplot theme

2. Import the ten datasets from the Ozone_TimeSeries folder in the Raw data folder. These contain ozone concentrations at Garinger High School in North Carolina from 2010-2019 (the EPA air database only allows downloads for one year at a time). Import these either individually or in bulk and then combine them into a single dataframe named `GaringerOzone` of 3589 observation and 20 variables.

```
# 1
getwd()
```

```
## [1] "/home/guest/EDA_2022/EDA-Fall2022"
```

```
library(tidyverse)
library(lubridate)
# install.packages('trend')
library(trend)
# install.packages('zoo')
library(zoo)
# install.packages('Kendall')
library(Kendall)
# install.packages('tseries')
library(tseries)
library(ggplot2)
library("formatR")
library(agricolae)
```

```
mynewtheme <- theme_grey(base_size = 12) + theme(axis.text = element_text(color = "Dark Green"),
    legend.position = "top")

theme_set(mynewtheme)

Ozone_2010 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2010_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2011 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2011_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2012 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2012_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2013 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2013_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2014 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2014_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2015 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2015_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2016 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2016_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2017 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2017_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2018 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2018_raw.csv",
    stringsAsFactors = TRUE)
Ozone_2019 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2019_raw.csv",
    stringsAsFactors = TRUE)

GaringerOzone <- rbind(Ozone_2010, Ozone_2011, Ozone_2012, Ozone_2013, Ozone_2014,
    Ozone_2015, Ozone_2016, Ozone_2017, Ozone_2018, Ozone_2019)
```

### Wrangle

3. Set your date column as a date class.

4. Wrangle your dataset so that it only contains the columns Date, Daily.Max.8.hour.Ozone.Concentration, and DAILY_AQI_VALUE.

5. Notice there are a few days in each year that are missing ozone concentrations. We want to generate a daily dataset, so we will need to fill in any missing days with NA. Create a new data frame that contains a sequence of dates from 2010-01-01 to 2019-12-31 (hint: `as.data.frame(seq())`). Call this new data frame Days. Rename the column name in Days to "Date".

6. Use a `left_join` to combine the data frames. Specify the correct order of data frames within this function so that the final dimensions are 3652 rows and 3 columns. Call your combined data frame GaringerOzone.

```
# 3

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

class(GaringerOzone$Date)

## [1] "Date"
# 4

GaringerOzone_new <- GaringerOzone %>%
```

```
    select(Date, Daily.Max.8.hour.Ozone.Concentration, DAILY_AQI_VALUE)

# 5

Days <- as.data.frame(seq(as.Date("2010/01/01"), as.Date("2019/12/31"), by = "days"))
colnames(Days) <- c("Date")

# 6

GaringerOzone <- left_join(Days, GaringerOzone_new, by = "Date")
```

## Visualize

7. Create a line plot depicting ozone concentrations over time. In this case, we will plot actual concentrations in ppm, not AQI values. Format your axes accordingly. Add a smoothed line showing any linear trend of your data. Does your plot suggest a trend in ozone concentration over time?
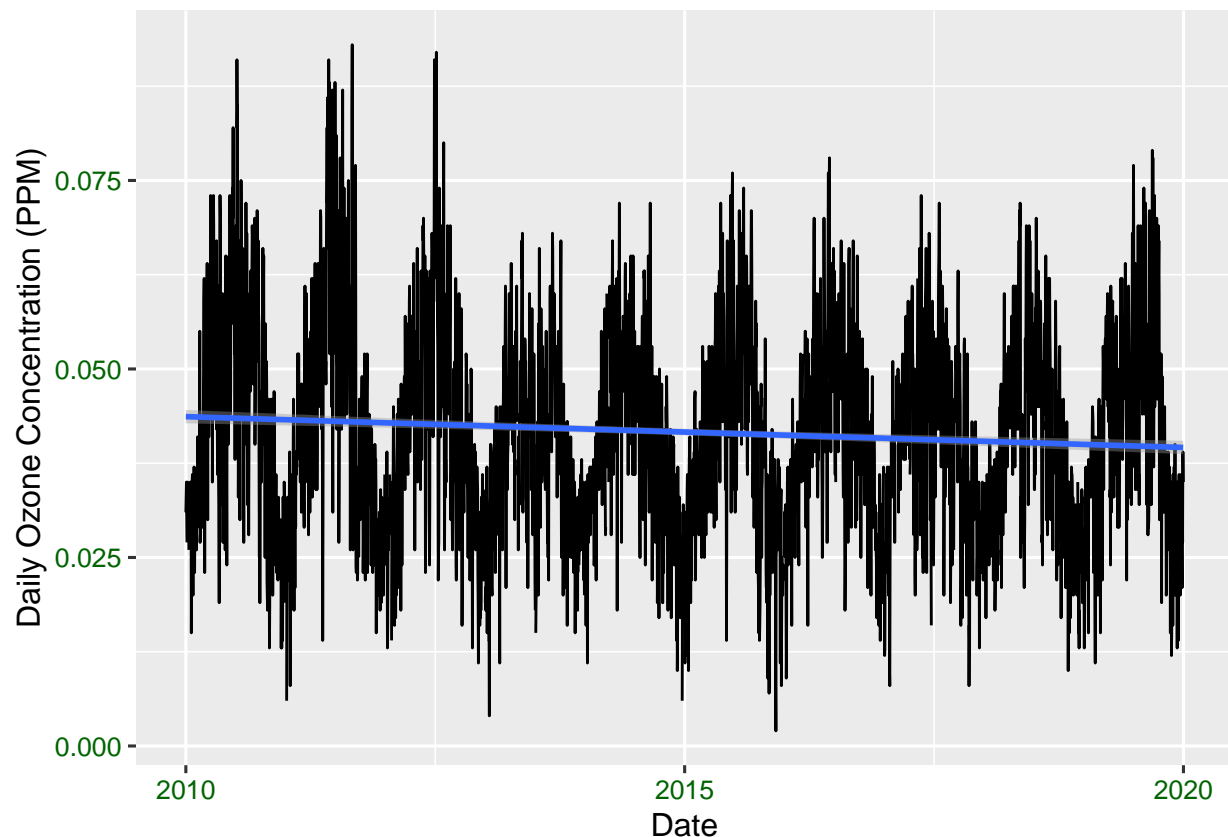
```
# 7

ggplot(GaringerOzone, aes(x = Date, y = Daily.Max.8.hour.Ozone.Concentration)) +
    geom_line() + geom_smooth(method = "lm") + ylab("Daily Ozone Concentration (PPM)")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

```
## Warning: Removed 63 rows containing non-finite values (stat_smooth).
```



Answer: My plot suggests a negative trend over time.

## Time Series Analysis

Study question: Have ozone concentrations changed over the 2010s at this station?

8. Use a linear interpolation to fill in missing daily data for ozone concentration. Why didn't we use a piecewise constant or spline interpolation?

```
# 8

Garinger_Ozone_clean <- GaringerOzone %>%
    mutate(Daily.Max.8.hour.Ozone.Concentration = zoo::na.approx(Daily.Max.8.hour.Ozone.Concentration))
```

Answer: We use the linear interpolation as it connects the gap between missing data points using a connect-the-dot method. This is the right choice as our data because the spline interpolation method uses an quadratic formula which wouldn't make sense for our trends and the piecewise constant interpolation uses the nearest neighbor approach which would not make sense as our NA's are isolated.

9. Create a new data frame called `GaringerOzone.monthly` that contains aggregated data: mean ozone concentrations for each month. In your pipe, you will need to first add columns for year and month to form the groupings. In a separate line of code, create a new Date column with each month-year combination being set as the first day of the month (this is for graphing purposes only)

```
# 9

GaringerOzone.monthly <- Garinger_Ozone_clean %>%
    mutate(month = month(Date), year = year(Date)) %>%
    mutate(Month_year = my(paste0(month, "-", year))) %>%
    select(Month_year, Daily.Max.8.hour.Ozone.Concentration) %>%
    group_by(Month_year) %>%
    summarise(Mean_PPM = mean(Daily.Max.8.hour.Ozone.Concentration, n = n()))
```

10. Generate two time series objects. Name the first `GaringerOzone.daily.ts` and base it on the dataframe of daily observations. Name the second `GaringerOzone.monthly.ts` and base it on the monthly average ozone values. Be sure that each specifies the correct start and end dates and the frequency of the time series.

```
# 10

GaringerOzone.daily.ts <- ts(Garinger_Ozone_clean$Daily.Max.8.hour.Ozone.Concentration,
    start = c(2010), frequency = 365)


GaringerOzone.monthly.ts <- ts(GaringerOzone.monthly$Mean_PPM, start = c(2010, 1),
    frequency = 12)

# Adding the end function cut off data so I chose to leave it out
```
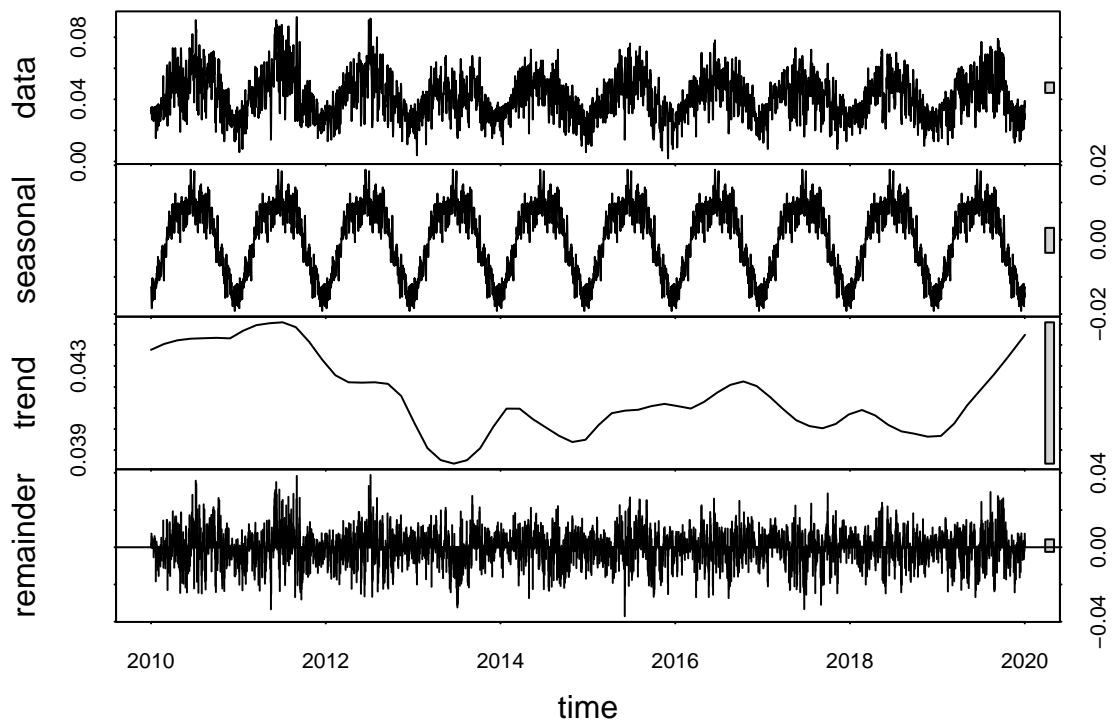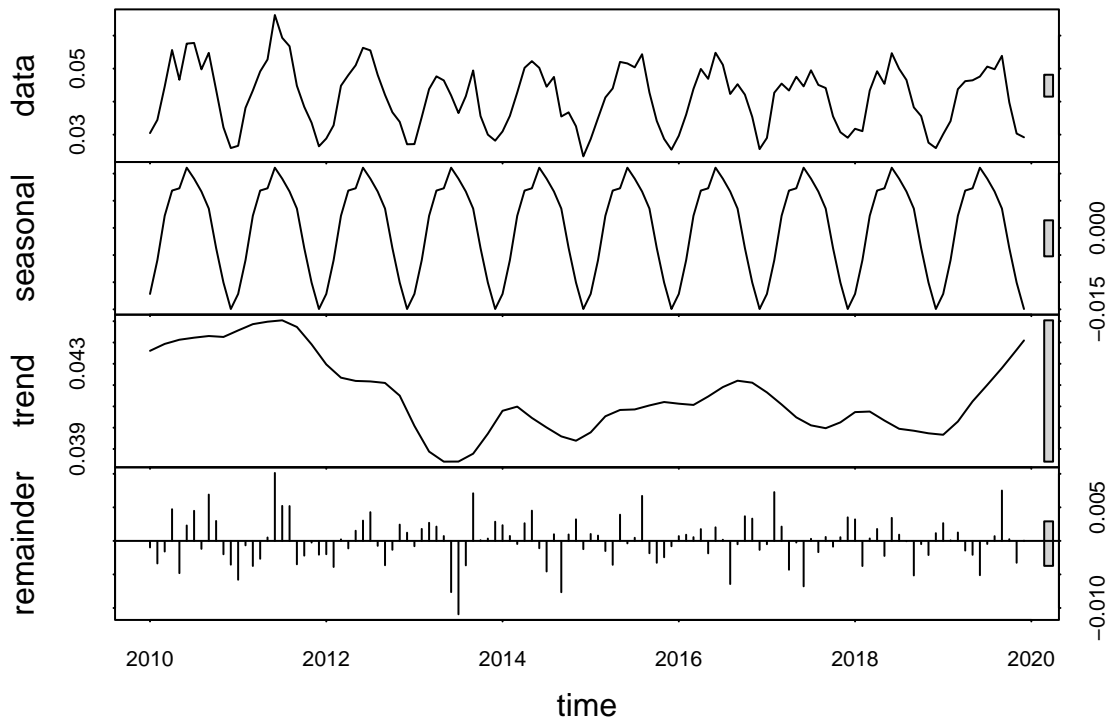
11. Decompose the daily and the monthly time series objects and plot the components using the `plot()` function.

```
# 11

GaringerOzone.daily.decomposed <- stl(GaringerOzone.daily.ts, s.window = "periodic")
plot(GaringerOzone.daily.decomposed)
```

```
GaringerOzone.monthly.decomposed <- stl(GaringerOzone.monthly.ts, s.window = "periodic")
plot(GaringerOzone.monthly.decomposed)
```

12. Run a monotonic trend analysis for the monthly Ozone series. In this case the seasonal Mann-Kendall is most appropriate; why is this?

```
# 12

Garinger_Monthly_trend <- Kendall::SeasonalMannKendall(GaringerOzone.monthly.ts)

Garinger_Monthly_trend
```

```
## tau = -0.143, 2-sided pvalue =0.046724
```

```
summary(Garinger_Monthly_trend)
```

```
## Score =  -77 , Var(Score) = 1499
## denominator =  539.4972
## tau = -0.143, 2-sided pvalue =0.046724
```

Answer: We are using the seasonal Mann-Kendall because in our plot we can see that seasonally does explain well the variability of ozone PPM over time.
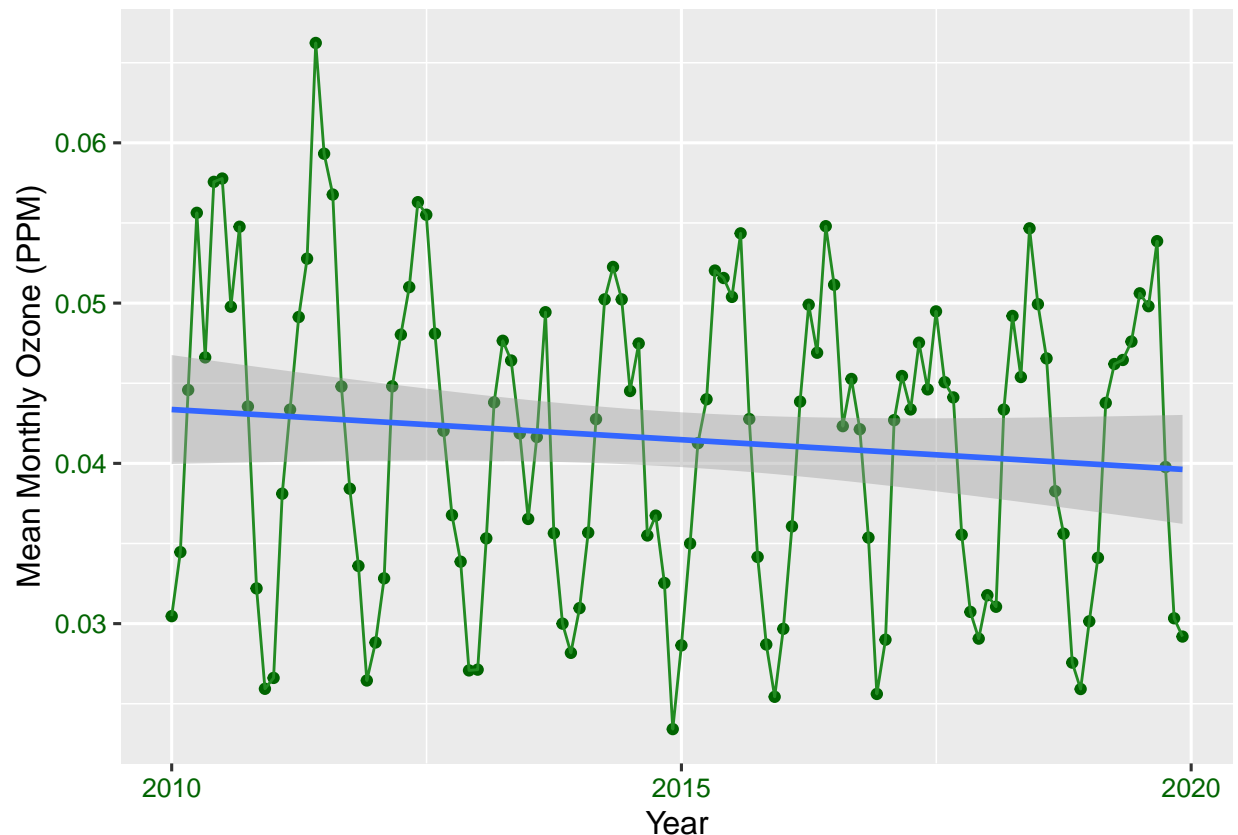
13. Create a plot depicting mean monthly ozone concentrations over time, with both a geom_point and a geom_line layer. Edit your axis labels accordingly.

```
# 13

Ozone_monthly_plot <- ggplot(GaringerOzone.monthly, aes(x = Month_year, y = Mean_PPM)) +
    geom_point(color = "dark green") + geom_line(color = "forest green") + geom_smooth(method = "lm") +
    ylab("Mean Monthly Ozone (PPM)") + xlab("Year")
```

```
Ozone_monthly_plot
```

```
## `geom_smooth()` using formula 'y ~ x'
```



14. To accompany your graph, summarize your results in context of the research question. Include output from the statistical test in parentheses at the end of your sentence. Feel free to use multiple sentences in your interpretation.

   Answer: Although there might appear to be a slight decrease in ozone (ppm) over time the rounded p-value tells us this decrease is not significant (p-value = 0.046724). The p-value can be rounded to .05 and therefore not statistically significant so I rejected the null hypothesis that there is no change in Ozone (PPM) over time.

15. Subtract the seasonal component from the `GaringerOzone.monthly.ts`. Hint: Look at how we extracted the series components for the EnoDischarge on the lesson Rmd file.

16. Run the Mann Kendall test on the non-seasonal Ozone monthly series. Compare the results with the ones obtained with the Seasonal Mann Kendall on the complete series.

```
# 15

GaringerOzone.monthly_Components <- GaringerOzone.monthly.ts - GaringerOzone.monthly.decomposed$time.se
    1]

# 16

GaringerOzone.monthly_Components_trend <- Kendall::MannKendall(GaringerOzone.monthly_Components)
```

```
GaringerOzone.monthly_Components_trend
```

```
## tau = -0.165, 2-sided pvalue =0.0075402
```

```
summary(GaringerOzone.monthly_Components_trend)
```

```
## Score =  -1179 , Var(Score) = 194365.7
## denominator =  7139.5
## tau = -0.165, 2-sided pvalue =0.0075402
```

Answer: The results of the Man Kendall test without the influence of seasonality tell us that the decrease in Ozone (pmm) over time is statistically significant (p-value =0.0075402). This p-value is lower than .05 unlike our p-value produced in the previous test, therefore we reject the null hypothesis that there is no decrease in Ozone (ppm) over time.