# Assignment 7: Time Series Analysis

## Lambert Ngenzi

## 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 Monday, March 14 at 7:00 pm.

### Set up

1. Set up your session:

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

```
#1
getwd()
```

```
## [1] "C:/Users/ln113/Documents/EDA_R/Environmental_Data_Analytics_2022"
```

```
library(tidyverse)
```

```
## -- Attaching packages --------------------------------------- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr   0.3.4
## v tibble  3.1.6      v dplyr   1.0.7
## v tidyr   1.1.4      v stringr 1.4.0
## v readr   2.1.1      v forcats 0.5.1
```

```
## -- Conflicts ------------------------------------------ tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
library(lubridate)
```

```
##
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':
##
##     date, intersect, setdiff, union
```

```
library(zoo)
```

```
##
## Attaching package: 'zoo'

## The following objects are masked from 'package:base':
##
##     as.Date, as.Date.numeric
```

```
library(trend)
library(plyr)
```

```
## ------------------------------------------------------------------------------

## You have loaded plyr after dplyr - this is likely to cause problems.
## If you need functions from both plyr and dplyr, please load plyr first, then dplyr:
## library(plyr); library(dplyr)

## ------------------------------------------------------------------------------

##
## Attaching package: 'plyr'

## The following objects are masked from 'package:dplyr':
##
##     arrange, count, desc, failwith, id, mutate, rename, summarise,
##     summarize

## The following object is masked from 'package:purrr':
##
##     compact
```

```
library(Kendall)
library(tseries)
```

```
## Registered S3 method overwritten by 'quantmod':
##   method             from
##   as.zoo.data.frame zoo
```

```
mytheme <- theme_classic(base_size = 14) +
  theme(axis.text = element_text(color = "black"),
        legend.position = "top")
theme_set(mytheme)
```

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.

```
#2
Ozone_TimeSeries = list.files("./Data/Raw/Ozone_TimeSeries/", pattern = "*csv"
                              , full.names = TRUE)

GaringerOzone <- Ozone_TimeSeries %>%
  ldply(read.csv)
```

## 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
Wrangled_GO <- GaringerOzone %>%
  select(Date, Daily.Max.8.hour.Ozone.Concentration, DAILY_AQI_VALUE)
min(Wrangled_GO$Date)
```

```
## [1] "2010-01-01"
```

```
max(Wrangled_GO$Date)
```

```
## [1] "2019-12-31"
```

```
# 5

Days <- as.data.frame(seq(as.Date("2010-01-01"), as.Date("2019-12-31"),
                          by = 1))


names(Days)[1] <- 'Date'

# 6
GaringerOzone <- left_join(Days, Wrangled_GO)
```

```
## Joining, by = "Date"
```

```
dim(GaringerOzone)
```
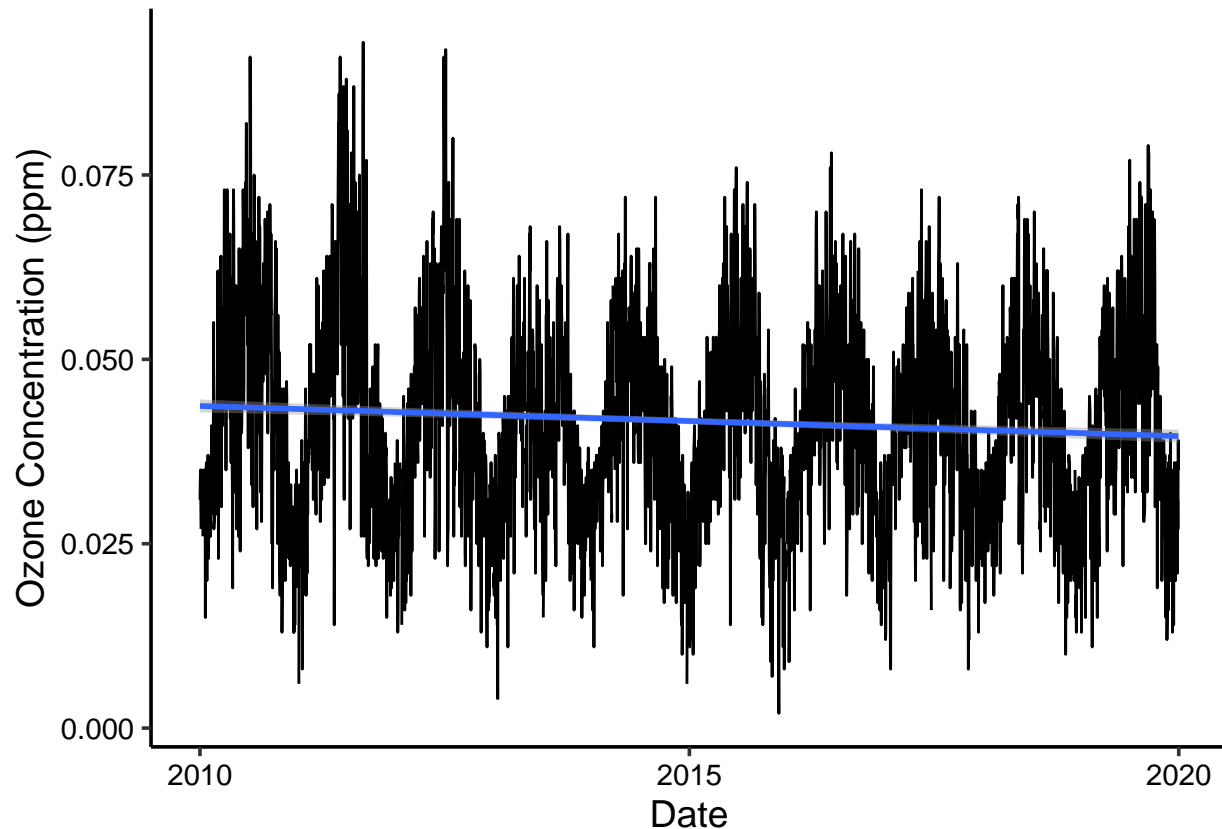
```
## [1] 3652    3
```

## 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(Date, Daily.Max.8.hour.Ozone.Concentration)) +
  geom_line() +
  geom_smooth(method = "lm") +
  ylab("Ozone Concentration (ppm)")
```

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

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



Answer: There is a slight negative trend over time of Ozone Concentration (ppm)

## 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
summary(GaringerOzone$Daily.Max.8.hour.Ozone.Concentration)

##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
## 0.00200 0.03200 0.04100 0.04163 0.05100 0.09300      63

GaringerOzone$Daily.Max.8.hour.Ozone.Concentration <-
  zoo::na.approx(GaringerOzone$Daily.Max.8.hour.Ozone.Concentration)
summary(GaringerOzone$Daily.Max.8.hour.Ozone.Concentration)

##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.00200 0.03200 0.04100 0.04151 0.05100 0.09300
```

> Answer: This is due to the fact that the trend has mostly linear trend, it makes sense that we will use a linear interpolation to fill the missing daily data while thr piecewise constant there is not jump, spline interpolation

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

4

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 <- GaringerOzone %>%
  mutate(Year = year(Date)) %>%
 mutate(Month = month(Date)) %>%
  group_by(Year, Month) %>%
 dplyr::summarise(OzoneMeanCon = mean(Daily.Max.8.hour.Ozone.Concentration))
```

```
## `summarise()` has grouped output by 'Year'. You can override using the `.groups` argument.
```

```
GaringerOzone.monthly$MY <- as.yearmon(paste(GaringerOzone.monthly$Year,
                                             GaringerOzone.monthly$Month), "%Y %m")
```

```
GaringerOzone.monthly$MY <- as.Date(GaringerOzone.monthly$MY, format = "%m-%Y")
```

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(GaringerOzone$Daily.Max.8.hour.Ozone.Concentration,
                             start = c(2010,1), frequency = 365)

GaringerOzone.monthly.ts <- ts(GaringerOzone.monthly$OzoneMeanCon,
                               start = c(2010,1), frequency = 12)
```
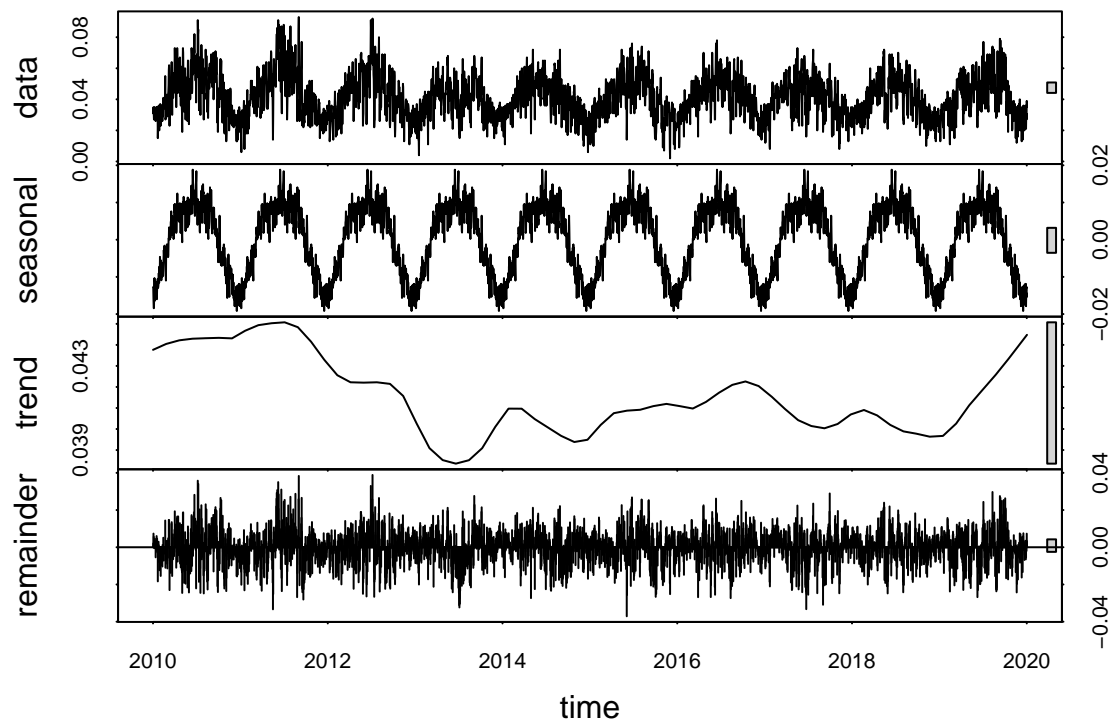
11. Decompose the daily and the monthly time series objects and plot the components using the `plot()` function.

```
#11
GaringerOzone.daily.ts.Dec <- stl(GaringerOzone.daily.ts, s.window = "periodic")

GaringerOzone.monthly.ts.Dec <- stl(GaringerOzone.monthly.ts,
                                    s.window = "periodic")

plot(GaringerOzone.daily.ts.Dec)
```

```
plot(GaringerOzone.monthly.ts.Dec)
```

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

```
Garinger.monthly.ts.mk <- Kendall::SeasonalMannKendall(GaringerOzone.monthly.ts)
summary(Garinger.monthly.ts.mk)
```

```
## Score =  -77 , Var(Score) = 1499
## denominator =  539.4972
## tau = -0.143, 2-sided pvalue =0.046724
```
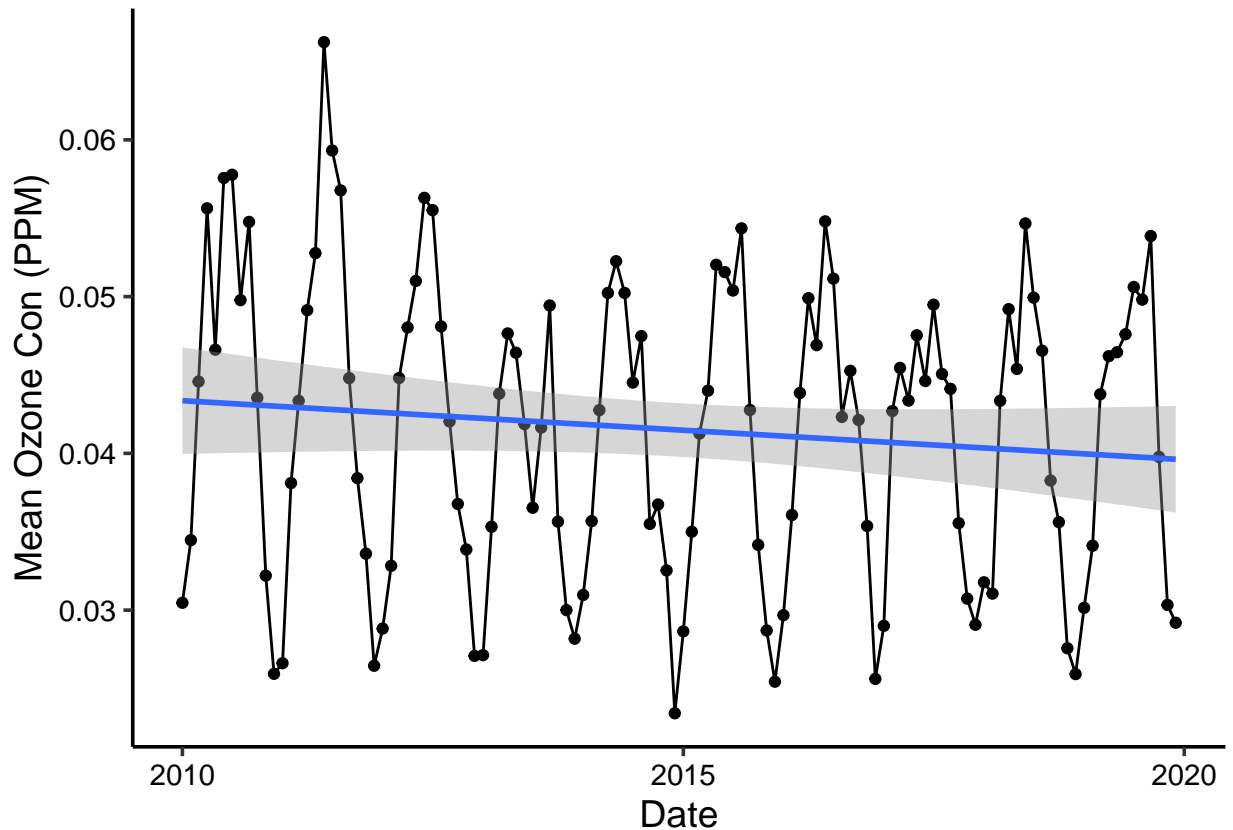
Answer: Due to the seasonality trend

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

```
ggplot(GaringerOzone.monthly, aes(MY, OzoneMeanCon)) +
  geom_point() +
  geom_line() +
  geom_smooth( method = "lm") +
  xlab("Date") +
  ylab("Mean Ozone Con (PPM)")
```

```
## `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: Overall, the Ozone mean concenreation over the entire data from 2010 to 2020 has not changed much, despite a slight a negative monotonic trend. Although it shows a strong seasonal trend. With the seasonal mean kendall, our 2 sided p value is equal 0.0467, therefore it is statiscally significant.

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
view(GaringerOzone.monthly.ts)

GaringerOzone.Subs <- as.data.frame(GaringerOzone.monthly.ts.Dec$time.series[,2:3])

GaringerOzone.Subs <- mutate(GaringerOzone.Subs,
                             obersved = GaringerOzone.monthly$OzoneMeanCon,
Date = GaringerOzone.monthly$MY)



#16
```

```
GaringerOzone.Subs.ts  <- ts(GaringerOzone.Subs$obersved, start = c(2010,1),
                             frequency = 12)

GMK <- MannKendall(GaringerOzone.Subs.ts)
summary(GMK)
```

```
## Score =  -424 , Var(Score) = 194364.7
## denominator =  7139
## tau = -0.0594, 2-sided pvalue =0.33732
```

```
GT <- mk.test(GaringerOzone.Subs.ts)
summary(GT)
```

```
##              Length Class  Mode
## data.name    1      -none- character
## p.value      1      -none- numeric
## statistic    1      -none- numeric
## null.value   1      -none- numeric
## parameter    1      -none- numeric
## estimates    3      -none- numeric
## alternative  1      -none- character
## method       1      -none- character
## pvalg        1      -none- numeric
```

Answer: We will reject the null hypotesis, there is a stronger tendency of decrease trend. The difference is much greater (score of -424 and p-value equals 0.337 smaller than the Seasonal Mann Kendal)