

Assignment 7: Time Series Analysis

Ben Egan

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/benja/OneDrive/Documents/R/win-library/Environmental_Data_Analytics_2022"

library(tidyverse)

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

## v ggplot2 3.3.5      v purrr  0.3.4
## v tibble  3.1.2      v dplyr  1.0.7
## v tidyr   1.1.3      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)  
  
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
```

```
Garinger2010 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2010_raw.csv", stringsAsFactors = FALSE)  
Garinger2011 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2011_raw.csv", stringsAsFactors = FALSE)  
Garinger2012 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2012_raw.csv", stringsAsFactors = FALSE)  
Garinger2013 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2013_raw.csv", stringsAsFactors = FALSE)  
Garinger2014 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2014_raw.csv", stringsAsFactors = FALSE)  
Garinger2015 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2015_raw.csv", stringsAsFactors = FALSE)  
Garinger2016 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2016_raw.csv", stringsAsFactors = FALSE)  
Garinger2017 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2017_raw.csv", stringsAsFactors = FALSE)  
Garinger2018 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2018_raw.csv", stringsAsFactors = FALSE)  
Garinger2019 <- read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2019_raw.csv", stringsAsFactors = FALSE)
```

```
GaringerOzone <- rbind(Garinger2010,Garinger2011,Garinger2012,Garinger2013,Garinger2014,Garinger2015,Garinger2016,Garinger2017,Garinger2018,Garinger2019)  
  
summary(GaringerOzone)
```

```
##           Date           Source      Site.ID           POC  
## 01/01/2010:    1    AQS      :3588    Min.    :371190041    Min.    :1
```

```

## 01/02/2010: 1 AirNow: 1 1st Qu.:371190041 1st Qu.:1
## 01/03/2010: 1 Median :371190041 Median :1
## 01/04/2010: 1 Mean :371190041 Mean :1
## 01/05/2010: 1 3rd Qu.:371190041 3rd Qu.:1
## 01/07/2010: 1 Max. :371190041 Max. :1
## (Other) :3583
## Daily.Max.8.hour.Ozone.Concentration UNITS DAILY_AQI_VALUE
## Min. :0.00200 ppm:3589 Min. : 2.00
## 1st Qu.:0.03200 1st Qu.: 30.00
## Median :0.04100 Median : 38.00
## Mean :0.04163 Mean : 41.57
## 3rd Qu.:0.05100 3rd Qu.: 47.00
## Max. :0.09300 Max. :169.00
##
## Site.Name DAILY_OBS_COUNT PERCENT_COMPLETE
## Garinger High School:3589 Min. : 6.00 Min. : 35.0
## 1st Qu.:17.00 1st Qu.:100.0
## Median :17.00 Median :100.0
## Mean :16.97 Mean : 99.8
## 3rd Qu.:17.00 3rd Qu.:100.0
## Max. :19.00 Max. :100.0
##
## AQS_PARAMETER_CODE AQS_PARAMETER_DESC CBSA_CODE
## Min. :44201 Ozone:3589 Min. :16740
## 1st Qu.:44201 1st Qu.:16740
## Median :44201 Median :16740
## Mean :44201 Mean :16740
## 3rd Qu.:44201 3rd Qu.:16740
## Max. :44201 Max. :16740
##
## CBSA_NAME STATE_CODE STATE
## Charlotte-Concord-Gastonia, NC-SC:3589 Min. :37 North Carolina:3589
## 1st Qu.:37
## Median :37
## Mean :37
## 3rd Qu.:37
## Max. :37
##
## COUNTY_CODE COUNTY SITE_LATITUDE SITE_LONGITUDE
## Min. :119 Mecklenburg:3589 Min. :35.24 Min. : -80.79
## 1st Qu.:119 1st Qu.:35.24 1st Qu.: -80.79
## Median :119 Median :35.24 Median : -80.79
## Mean :119 Mean :35.24 Mean : -80.79
## 3rd Qu.:119 3rd Qu.:35.24 3rd Qu.: -80.79
## Max. :119 Max. :35.24 Max. : -80.79
##

```

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_Processed <- 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"))

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

# 6
GaringerOzone <- left_join(x = Days, y = GaringerOzone_Processed)

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

Visualize

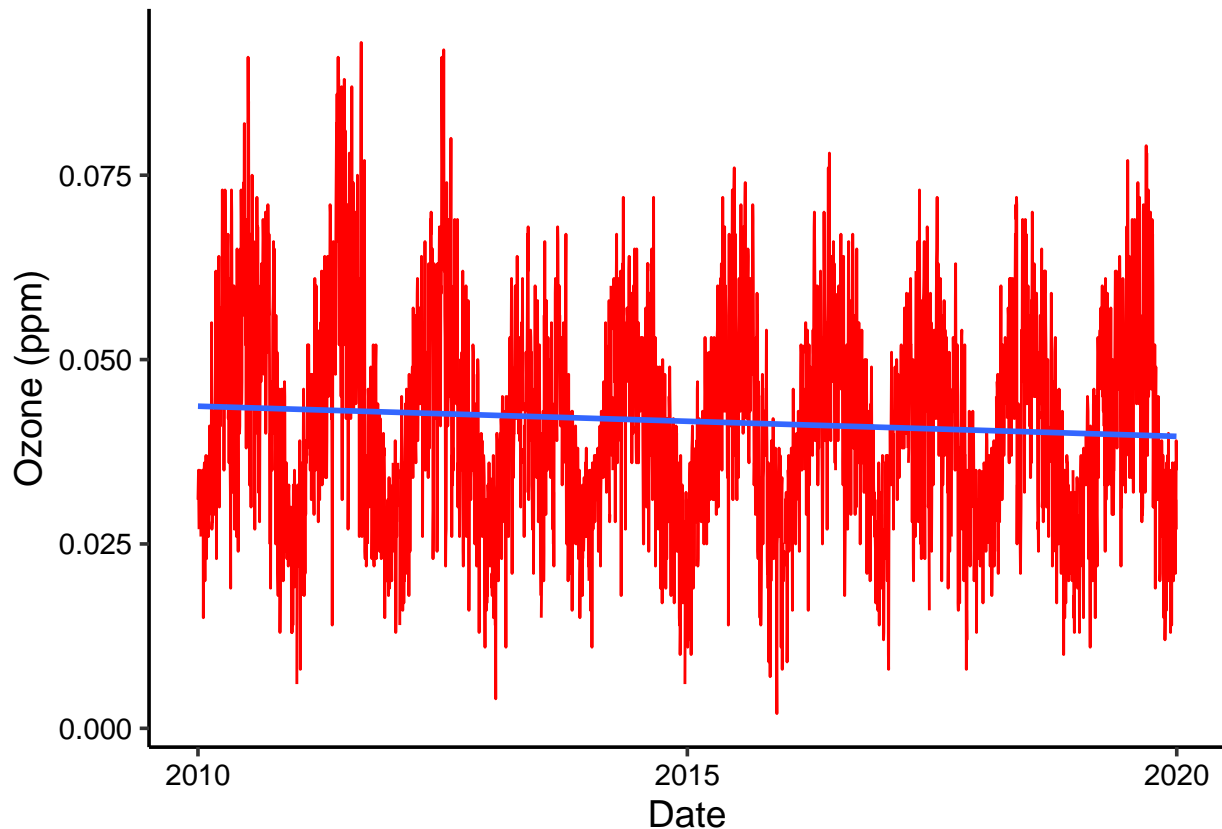
7. Create a line plot depicting ozone concentrations over time. In this case, we will plot actual concentrations in ppm, not AQI values. Format your axes accordingly. Add a smoothed line showing any linear trend of your data. Does your plot suggest a trend in ozone concentration over time?

```
#7
GaringerOzone_plot <-
  ggplot(GaringerOzone, aes(x = Date, y = Daily.Max.8.hour.Ozone.Concentration )) +
  geom_line(color = "red") +
  geom_smooth(method = lm, se = FALSE) +
  labs(x = "Date", y = "Ozone (ppm)")

print(GaringerOzone_plot)

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

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



Answer: There appears to be a slight downward trend of ozone concentration with 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

```
GaringerOzone$Daily.Max.8.hour.Ozone.Concentration =  
zoo::na.approx(GaringerOzone$Daily.Max.8.hour.Ozone.Concentration)
```

Answer: The Data has a regular seasonality therefore it can be concluded that the missing data is either rising or falling between two data points before or after it. The best assumption is to pick a value between this two data points using linear interpolation as opposed to picking the same number using piecewise. Linear is the best option as the change in data is closer to a line than a quadratic curve therefore spline should also not be used.

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 <- GaringerOzone %>%
  mutate(Month = month(Date), Year = year(Date)) %>%
  mutate(Month_Year = my(paste0(Month, "-", Year))) %>%
  dplyr::group_by(Month_Year) %>%
  dplyr::summarise(MeanOzone =
    mean(Daily.Max.8.hour.Ozone.Concentration))
```

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

```
f_month <- month(first(GaringerOzone$Date))
f_year <- year(first(GaringerOzone$Date))

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

fm_month <- month(first(GaringerOzone.monthly$Month_Year))
fm_year <- year(first(GaringerOzone.monthly$Month_Year))

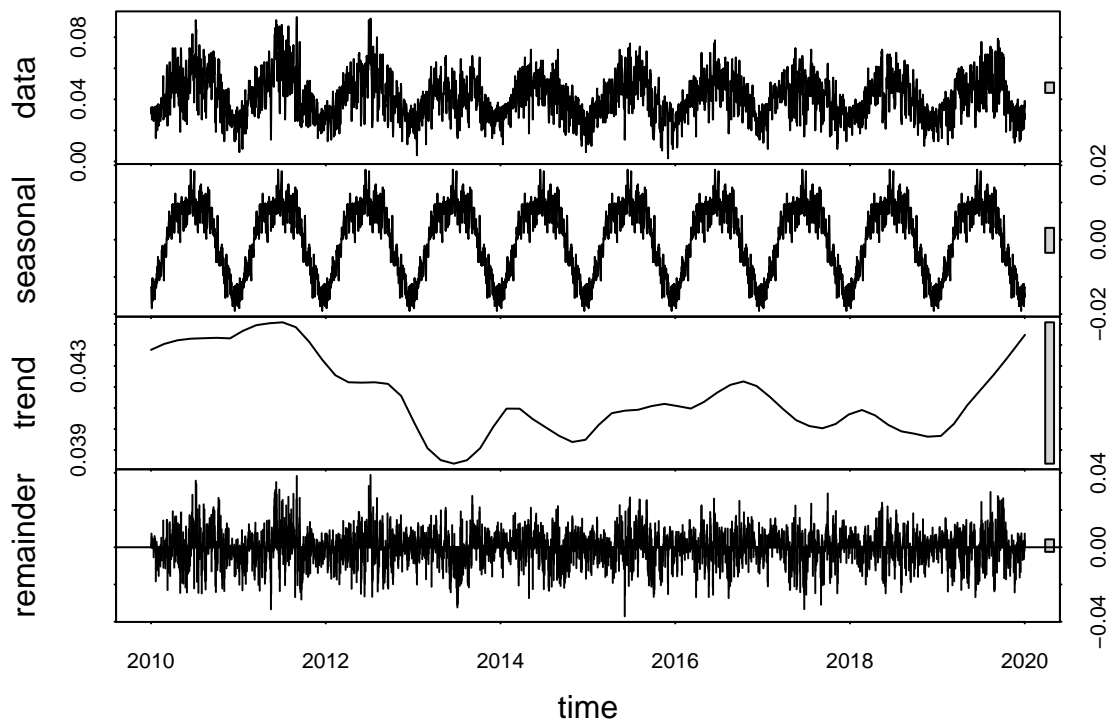
GaringerOzone.monthly.ts <-
  ts(GaringerOzone.monthly$MeanOzone,
     start = c(fm_year,fm_month), frequency = 12)
```

11. Decompose the daily and the monthly time series objects and plot the components using the `plot()` function.

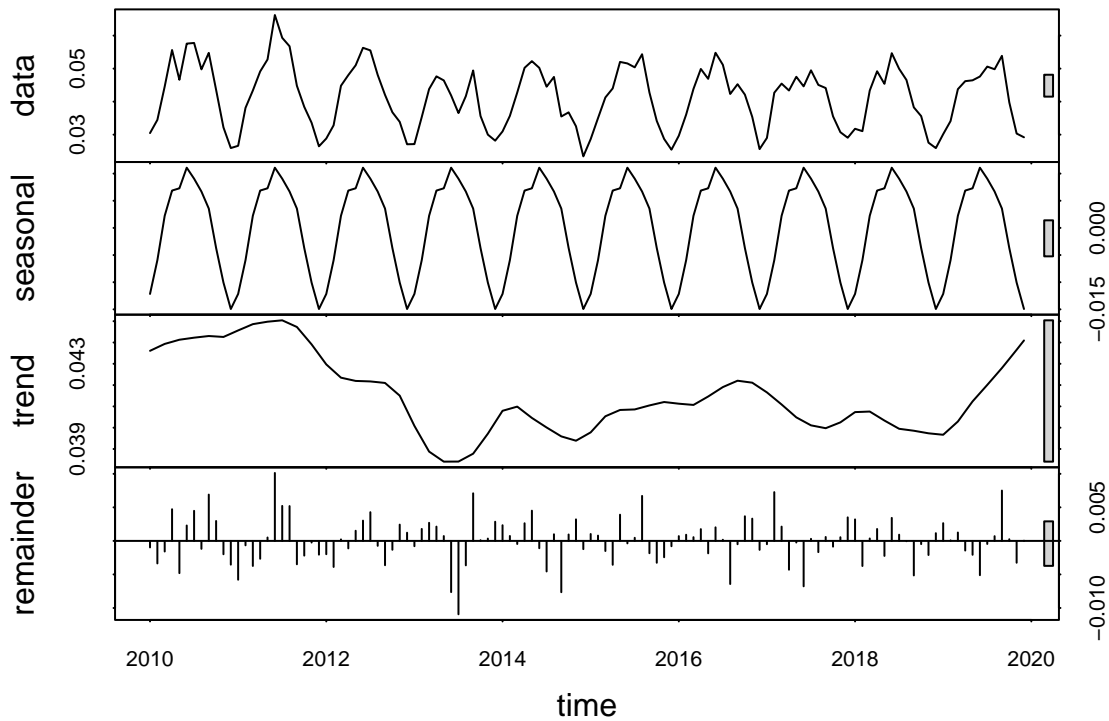
#11

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

plot(daily_data_decomp)
```



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



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

Ozone_data_trendm

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

summary(Ozone_data_trendm)

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

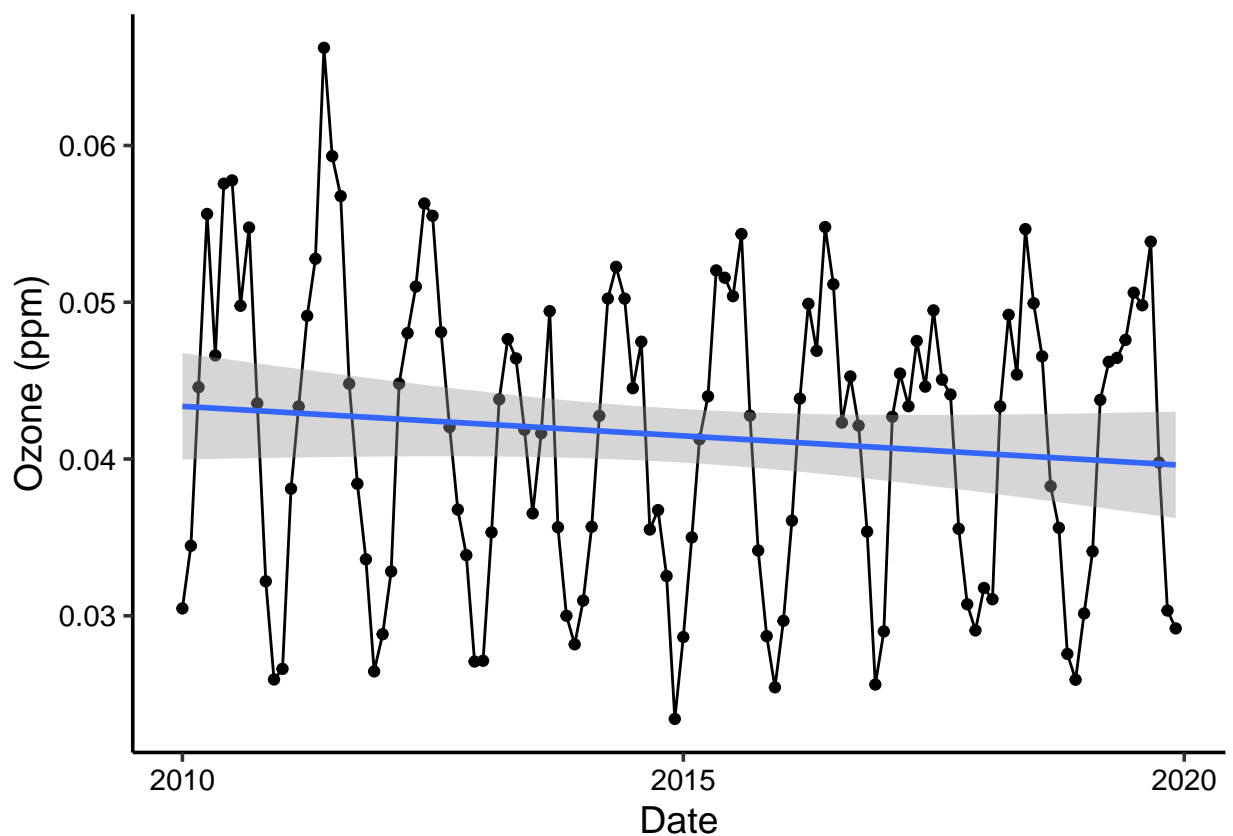
Answer: Because the data is seasonal, the Seasonal mann kendall allows us to interpret the data while comparing multiple seasons.

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_data_plot <-  
ggplot(GaringerOzone.monthly, aes(x = Month_Year, y = MeanOzone)) +  
  geom_point() +  
  geom_line() +  
  ylab("Mean Ozone Concentration") +  
  geom_smooth( method = lm ) +  
  labs( x = "Date", y = "Ozone (ppm)" )  
  
print(Ozone_data_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: Visually, there is a downward trend. the P value is less than .05 at .047. Over time, Ozone concentration does decrease and it is statistically 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

```
GaringerOzone.monthly_Components <-  
  as.data.frame(monthly_data_decomp$time.series)  
  
GaringerOzone.monthly_Components <-  
  mutate(GaringerOzone.monthly_Components,  
    Ozone = GaringerOzone.monthly$MeanOzone,  
    Date = GaringerOzone.monthly$Month_Year) %>%  
  mutate(NonSeason = Ozone-seasonal)
```

#16

```
Ozone_data_trendm2 <-  
  Kendall::MannKendall(GaringerOzone.monthly_Components$NonSeason)  
  
summary(Ozone_data_trendm2)
```

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

Answer: The results are more significant removing the seasonality. using the seasonal Mannk-
endall the p value was .046 while the Mann Kendall with seasonality removed returned a Pvalue
of 0.0075.