

Assignment 7: Time Series Analysis

Langston ALEXander

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/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Assignments"
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)
```

```
## Warning: package 'zoo' was built under R version 4.1.3
```

```
##
```

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

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      as.Date, as.Date.numeric
```

```
library(trend)
```

```
## Warning: package 'trend' was built under R version 4.1.3
```

```
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
```

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

```
OzoneFiles = list.files(  
  path = "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
OzoneFiles
```

```
## [1] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
## [2] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
## [3] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
## [4] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
## [5] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
## [6] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
## [7] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
```

```
## [8] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
## [9] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E
## [10] "C:/Users/lwa8/Documents/R/ENV872/Environmental_Data_Analytics_2022/Data/Raw/Ozone_TimeSeries/E

GaringerOzone <- OzoneFiles %>%
  ldply(read.csv)

GaringerOzone$Date <- as.factor(GaringerOzone$Date)
```

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")

# 4

GaringerOzoneSelect <- GaringerOzone %>%
  select(Date, Daily.Max.8.hour.Ozone.Concentration, DAILY_AQI_VALUE)

GaringerOzoneSelect$Date <- as.character(
  GaringerOzoneSelect$Date)

# 5

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

names(Days) <- "Date"

Days$Date <- as.character(Days$Date)

# 6

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

GaringerOzone$Date <- as.Date(GaringerOzone$Date)
```

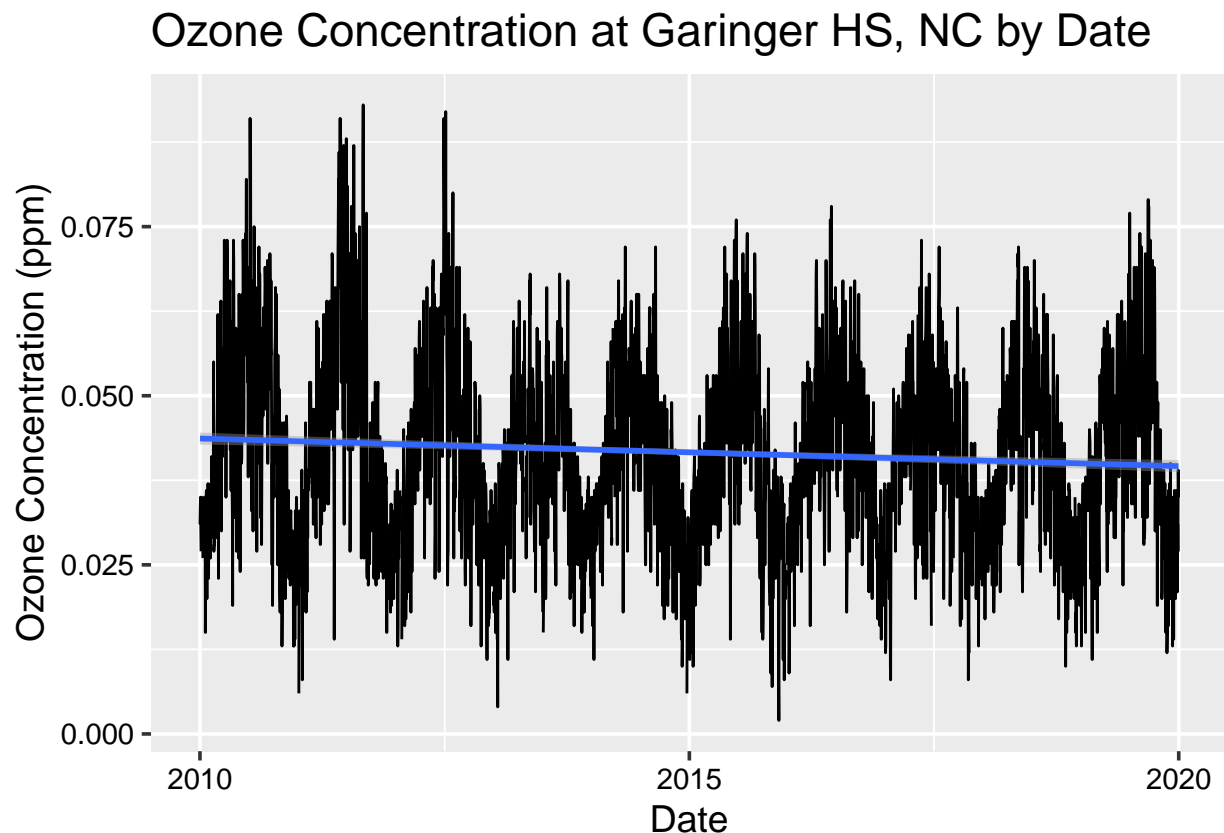
Visualize

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")+
  labs(
    title = "Ozone Concentration at Garinger HS, NC by Date ",
    y = "Ozone Concentration (ppm)"
  )

## `geom_smooth()` using formula 'y ~ x'
## Warning: Removed 63 rows containing non-finite values (stat_smooth).
```



Answer: Visually, the line plot above suggests that ozone levels at Garinger Highschool have slightly declined on average in the past 10 years. We see this trend mirrored in the smoothed linear trend line of the data.

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 <- mutate(
  GaringerOzone, Daily.Max.8.hour.Ozone.Concentration = zoo::na.approx(Daily.Max.8.hour.Ozone.Concentra
```

Answer: For the most part, the short-term trends throughout the data are increasing or decreasing linearly. The missing data point should fit into this short-term linear trend. Neither a piecewise constant nor a spline interpolation would replace missing data along a linear trend. A piecewise constant would replace the missing data by making them equal to their nearest neighbor, while the spline interpolation would use quadratic function to interpolate.

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

#group by year, then month, create new column of monthly average, THERE SHOULD BE NO DAILY DATA

```
GaringerOzone.monthly <- GaringerOzone %>%
  mutate(month = month(Date)) %>%
  mutate(year = year(Date)) %>%
  group_by(year, month) %>%
  dplyr::summarize(Ozone_Avg = mean(Daily.Max.8.hour.Ozone.Concentration)) %>%
  as.data.frame()
```

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

```
GaringerOzone.monthly$Date <- seq(
  from = as.Date("2010-01-01"),
  to = as.Date("2019-12-01"),
  by = "1 month")
```

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
)

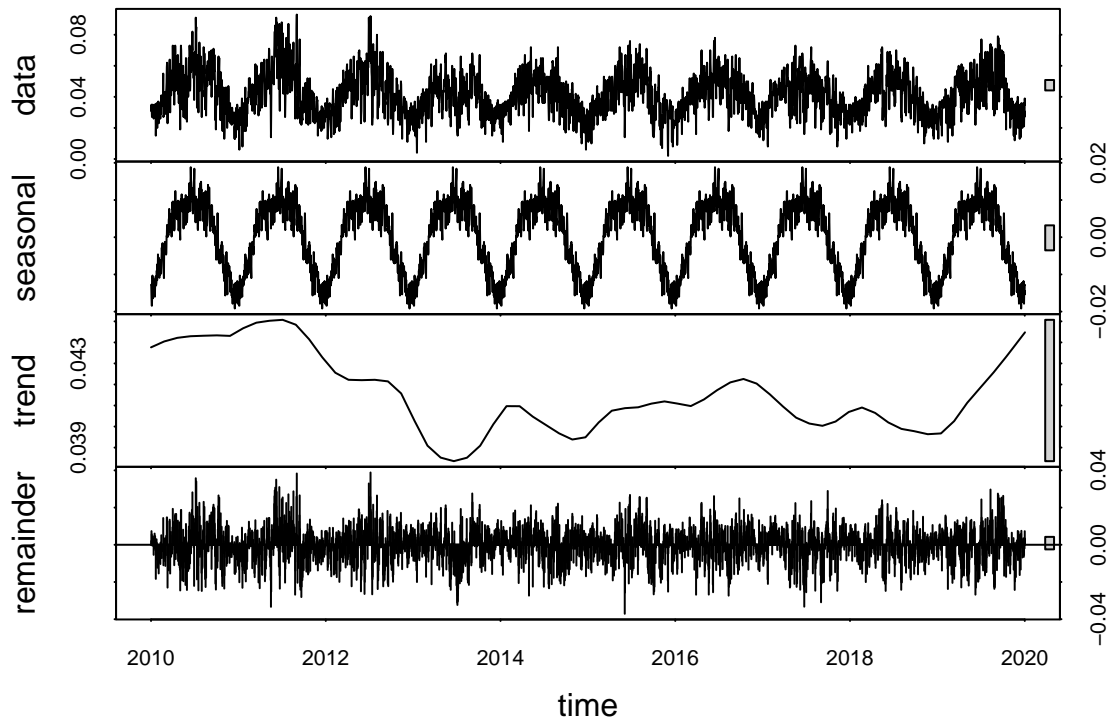
GaringerOzone.monthly.ts <- ts(
  GaringerOzone.monthly$Ozone_Avg,
  start(2010,1),
  frequency = 12
)
```

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

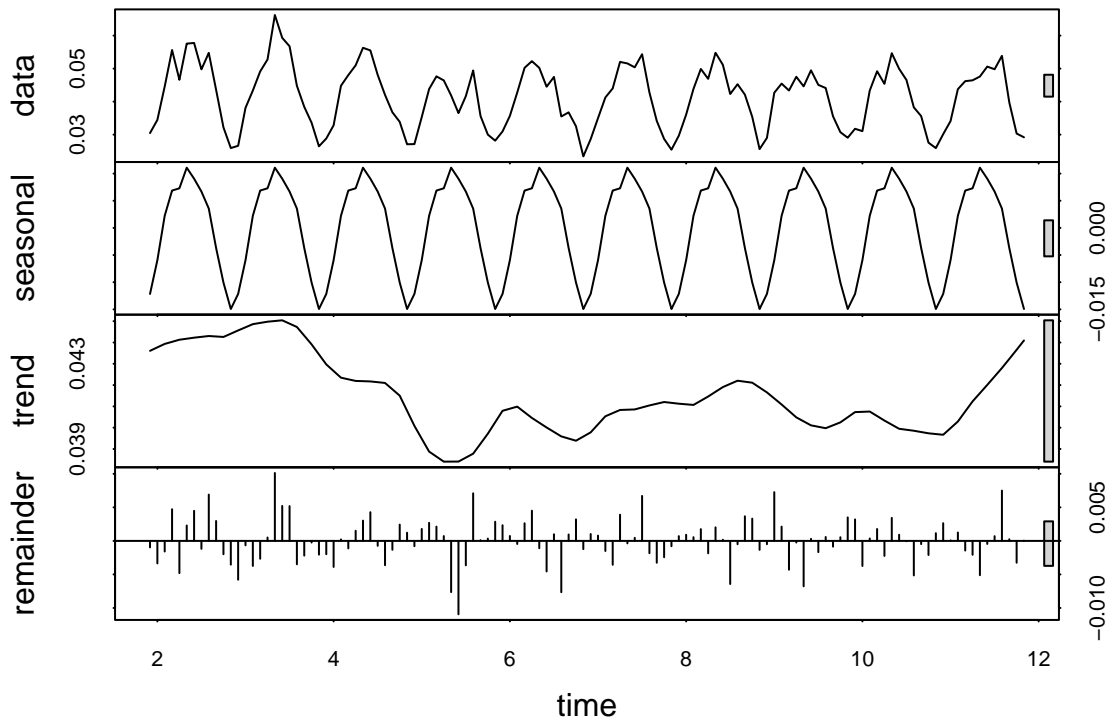
#11

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

```
plot(GaringerOzone.daily.decomp)
```



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

```
GaringerOzone.trend1 <- trend::smk.test(
  GaringerOzone.monthly.ts)
summary(GaringerOzone.trend1)
```

```
##
## Seasonal Mann-Kendall trend test (Hirsch-Slack test)
##
## data: GaringerOzone.monthly.ts
## alternative hypothesis: two.sided
##
## Statistics for individual seasons
##
## H0
##
```

	S	varS	tau	z	Pr(> z)
## Season 1:	S = 0	-1	125	-0.022	0.000
## Season 2:	S = 0	-4	124	-0.090	-0.269
## Season 3:	S = 0	-17	125	-0.378	-1.431
## Season 4:	S = 0	-15	125	-0.333	-1.252
## Season 5:	S = 0	-17	125	-0.378	-1.431
## Season 6:	S = 0	-11	125	-0.244	-0.894
## Season 7:	S = 0	-7	125	-0.156	-0.537
## Season 8:	S = 0	-5	125	-0.111	-0.358

```
Pr(>|z|)
1.00000
0.78762
0.15241
0.21050
0.15241
0.37109
0.59151
0.72051
```

```
## Season 9:   S = 0  -13  125 -0.289 -1.073  0.28313
## Season 10:  S = 0  -13  125 -0.289 -1.073  0.28313
## Season 11:  S = 0   11  125  0.244  0.894  0.37109
## Season 12:  S = 0   15  125  0.333  1.252  0.21050
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

GaringerOzone.trend2 <- Kendall::SeasonalMannKendall(
  GaringerOzone.monthly.ts)
summary(GaringerOzone.trend2)

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

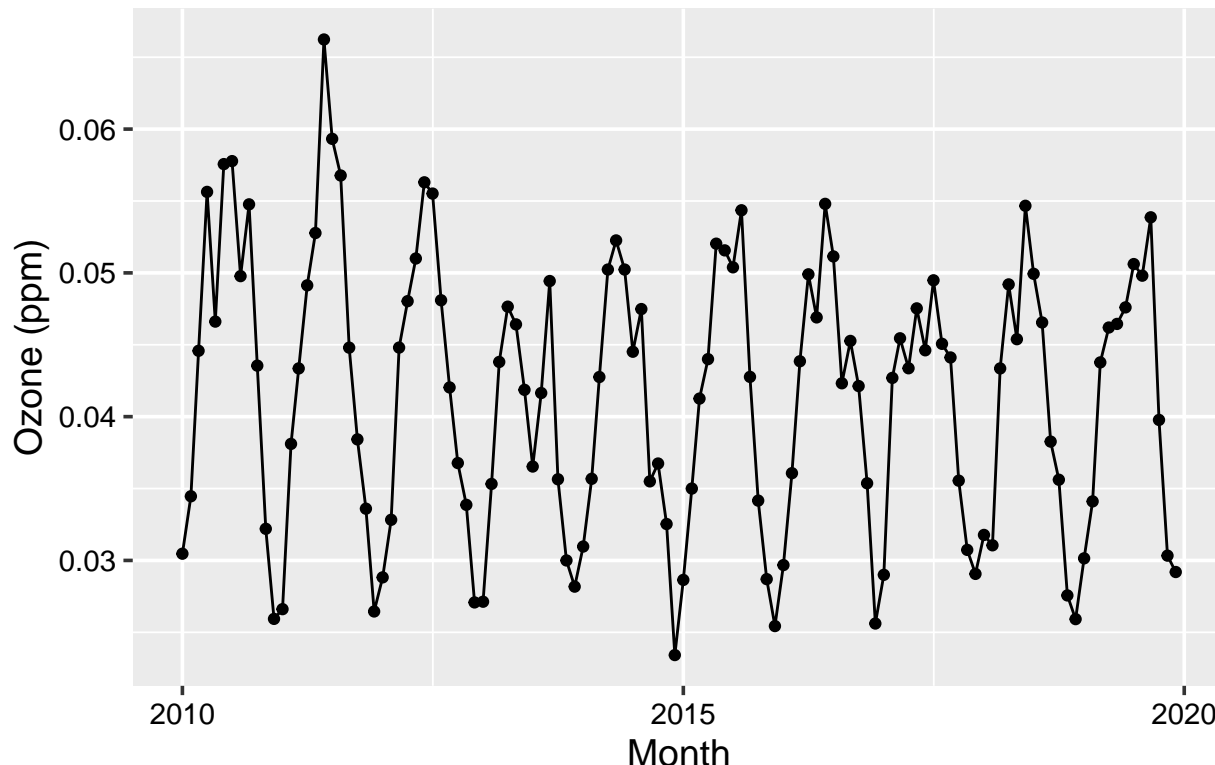
Answer: We can see from the decomposition graphs above that there is a regular seasonal trend over time. The only trend analysis that can be used with seasonal data in the seasonal Mann-Kendall test.

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(x = Date, y = Ozone_Avg)) +
  geom_point() +
  geom_line()+
  labs(
    title = "Ozone Concentration at Garinger HS, NC by Month",
    x = "Month", y = "Ozone (ppm)")
```


Ozone Concentration at Garinger HS, NC by Month



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: In response to the study question, the ozone concentration has changed between 2010-2020 at Garinger HS. According to the Seasonal Mann-Kendall test we are statistically confident that we can reject the null hypothesis that the ozone level has not changed in the 2010s. It is important to note though that we replaced some data points using a linear model. (Seasonal Mann-Kendall results : Score = -77 , Var(Score) = 1499, denominator = 539.4972, tau = -0.143, 2-sided pvalue = 0.046724)

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

```
Garinger.monthly.Components <- as.data.frame(
  GaringerOzone.monthly.decomp$time.series[,2:3])
```

#16

```
GaringerOzone.nonSeasonal.trend <- Kendall::MannKendall(
  Garinger.monthly.Components$trend)
```

```
summary(GaringerOzone.nonSeasonal.trend)
```

```
## Score = -1922 , Var(Score) = 194366.7  
## denominator = 7140  
## tau = -0.269, 2-sided pvalue =1.3168e-05
```

Answer: The p-value from the non-seasonal Mann Kendall test was much smaller than the seasonal Mann Kendall test we ran earlier.