# Assignment 8: Time Series Analysis

## Vignesh Arunkumar

### Fall 2024

## OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics on generalized linear models.

## Directions

1. Rename this file <FirstLast>_A08_TimeSeries.Rmd (replacing <FirstLast> with your first and last name).
2. Change "Student Name" on line 3 (above) with your name.
3. Work through the steps, **creating code and output** that fulfill each instruction.
4. Be sure to **answer the questions** in this assignment document.
5. When you have completed the assignment, **Knit** the text and code into a single PDF file.

## Set up

1. Set up your session:

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

```
getwd()
```

```
## [1] "/home/guest/EDE_Fall2024"
```

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----------------------- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr    1.5.1
## v ggplot2    3.5.1      v tibble     3.2.1
## v lubridate  1.9.3      v tidyr      1.3.1
## v purrr      1.0.2
## -- Conflicts ------------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```r
library(lubridate)
library(zoo)


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

library(trend)
library(Kendall)
theme_update(legend.position="left",
             plot.background = element_rect(fill ="Forest green"))
```

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.

```r
#1
O3_2010 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2010_raw.csv")
O3_2011 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2011_raw.csv")
O3_2012 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2012_raw.csv")
O3_2013 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2013_raw.csv")
O3_2014 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2014_raw.csv")
O3_2015 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2015_raw.csv")
O3_2016 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2016_raw.csv")
O3_2017 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2017_raw.csv")
O3_2018 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2018_raw.csv")
O3_2019 <- read.csv("~/EDE_Fall2024/Data/Raw/Ozone_TimeSeries/EPAair_O3_GaringerNC2019_raw.csv")

GaringerOzone <- rbind(O3_2010,O3_2011,O3_2012,O3_2013,O3_2014,O3_2015,
                       O3_2016,O3_2017,O3_2018,O3_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")

# 4
GaringerOzone.wrangled <-
  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
NewGaringerOzone <- left_join(Days, GaringerOzone.wrangled, 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
Garinger_plot1 <- ggplot(NewGaringerOzone, aes(x = NewGaringerOzone$Date,
                y = NewGaringerOzone$Daily.Max.8.hour.Ozone.Concentration))+
  geom_line()+
  geom_smooth()+
  ylab("Max 8 hour Ozone Concentration" )+
  xlab("Time")
Garinger_plot1
```

```
## Warning: Use of 'NewGaringerOzone$Date' is discouraged.
## i Use 'Date' instead.


## Warning: Use of 'NewGaringerOzone$Daily.Max.8.hour.Ozone.Concentration' is discouraged.
## i Use 'Daily.Max.8.hour.Ozone.Concentration' instead.


## Warning: Use of 'NewGaringerOzone$Date' is discouraged.
## i Use 'Date' instead.


## Warning: Use of 'NewGaringerOzone$Daily.Max.8.hour.Ozone.Concentration' is discouraged.
## i Use 'Daily.Max.8.hour.Ozone.Concentration' instead.


## 'geom_smooth()' using method = 'gam' and formula = 'y ~ s(x, bs = "cs")'


## Warning: Removed 63 rows containing non-finite outside the scale range
## ('stat_smooth()').
```
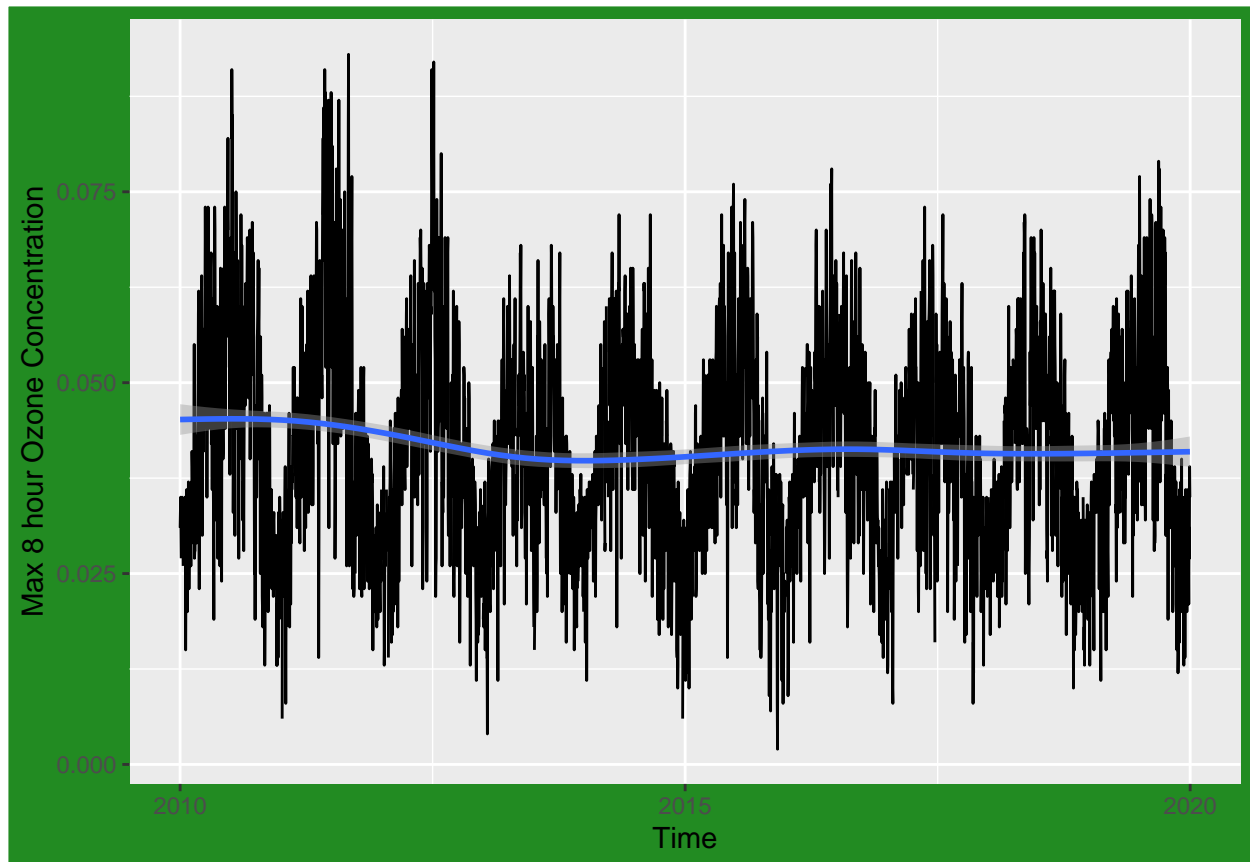
Answer: The plot suggest that the trend in ozone concentration overtime decreases slighlty overtime.

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

Answer: There is not local region where you can use piecewise constants that is because the data is cyclical. Spline interpolation doesnt work the best because our trend is not quadratic.

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)

4

```
#9
GaringerOzone.monthly <- Garinger_fill %>%
  mutate(year=year(Garinger_fill$Date)) %>%
  mutate(month=month(Garinger_fill$Date)) %>%
  group_by(year, month) %>%
  summarize(meanppm=mean(Daily.Max.8.hour.Ozone.Concentration))%>%
      mutate(Date=as.Date(paste(year, month, "01", sep="-"), "%Y-%m-%d"))
```

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

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.
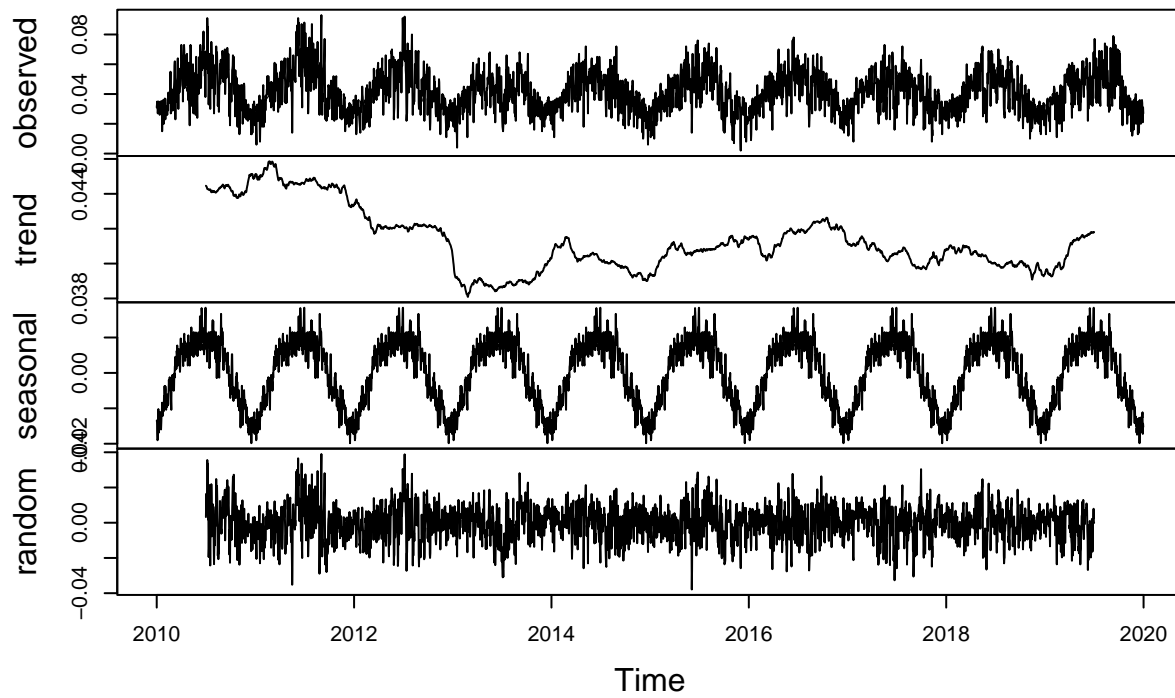
```
#10

GaringerOzone.daily.ts <- ts(Garinger_fill$Daily.Max.8.hour.Ozone.Concentration,
                    start = c(2010,1), end = c(2019, 365),frequency = 365)

GaringerOzone.monthly.ts <- ts(GaringerOzone.monthly$meanppm, start = c(2010),
                        end = c(2019),frequency = 12)
```

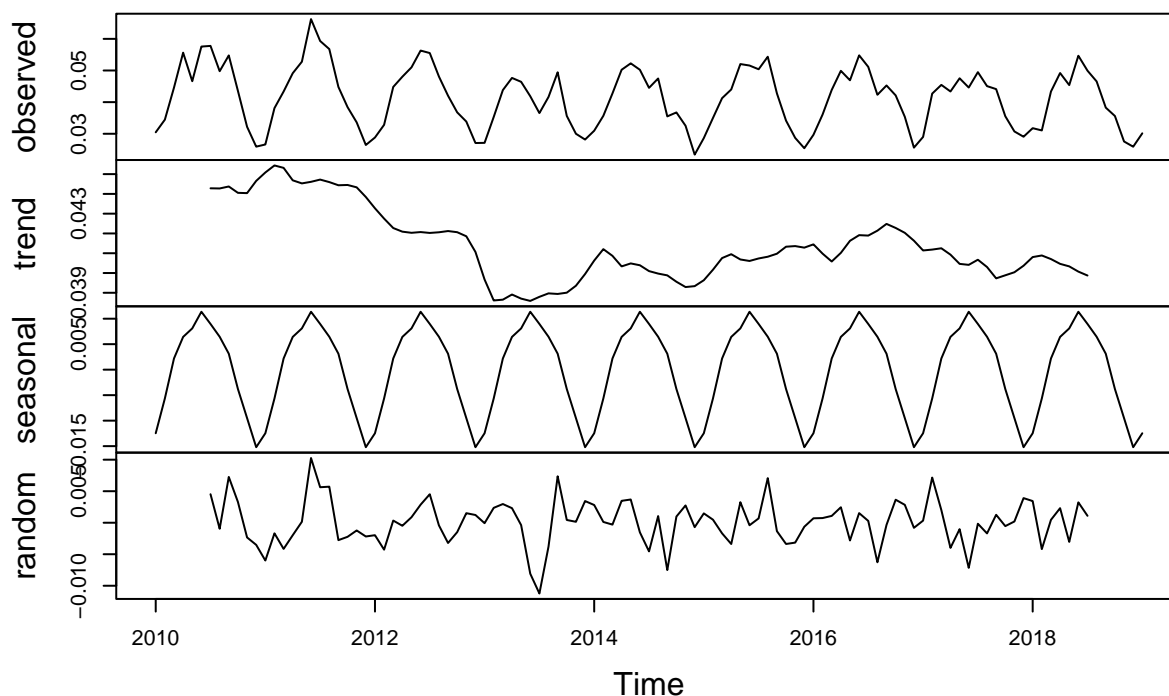11. Decompose the daily and the monthly time series objects and plot the components using the `plot()` function.

```
#11
GaringerOzone.daily.decomp <- decompose(GaringerOzone.daily.ts)
plot(GaringerOzone.daily.decomp)
```

**Decomposition of additive time series**



```
GaringerOzone.monthly.decomp <- decompose(GaringerOzone.monthly.ts)
plot(GaringerOzone.monthly.decomp)
```

## Decomposition of additive time series



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.monthly.trend <- SeasonalMannKendall(GaringerOzone.monthly.ts)
GaringerOzone.monthly.trend
```

```
## tau = -0.182, 2-sided pvalue =0.017617
```

```
summary(GaringerOzone.monthly.trend)
```

```
## Score =  -80 , Var(Score) = 1136
## denominator =  440.4965
## tau = -0.182, 2-sided pvalue =0.017617
```
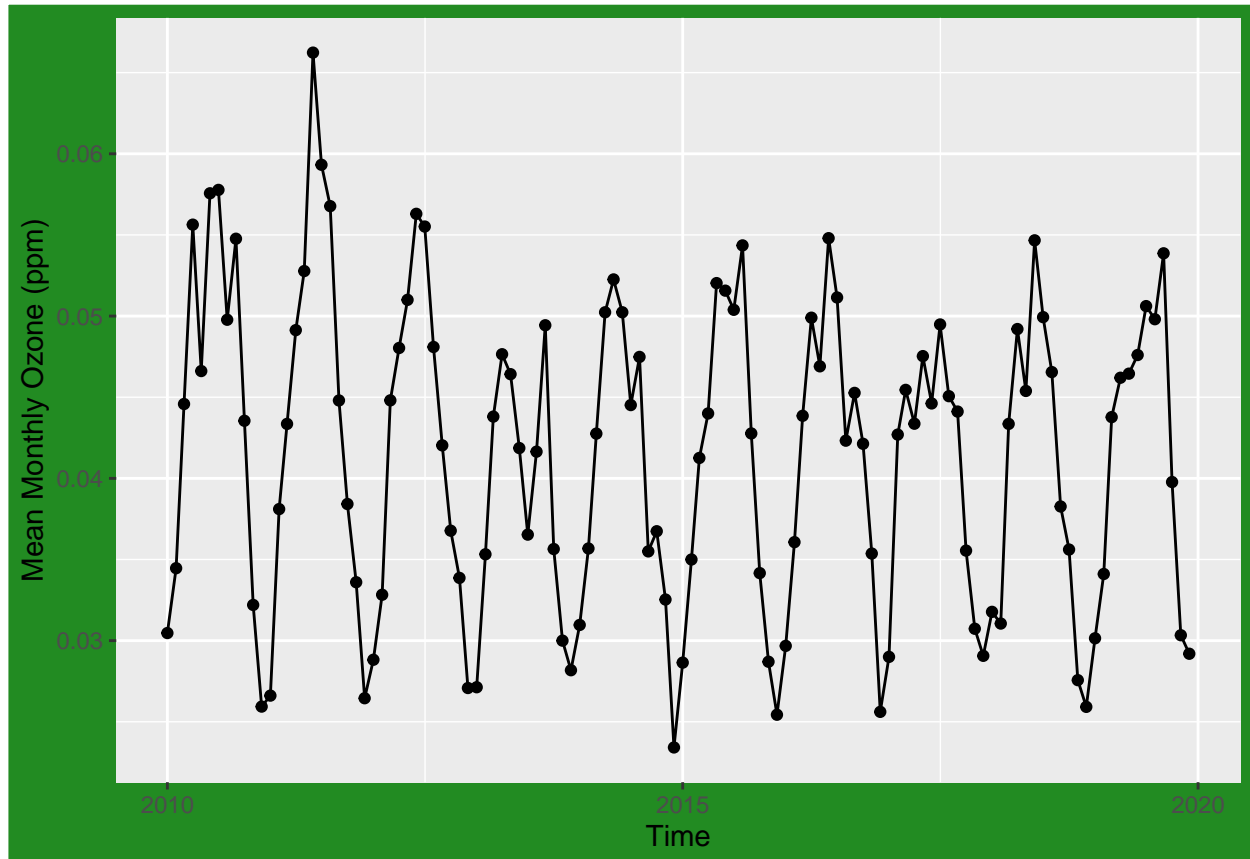
Answer: It accounts for our seasonal data in time series where data is collected at irregular intervals.

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.

7

```
# 13

GaringerOzone.mean.monthly.plot <- ggplot(GaringerOzone.monthly, aes(x= GaringerOzone.monthly$Date, y =
  geom_point()+
  geom_line()+
  ylab("Mean Monthly Ozone (ppm)")+
  xlab("Time")

GaringerOzone.mean.monthly.plot
```



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.
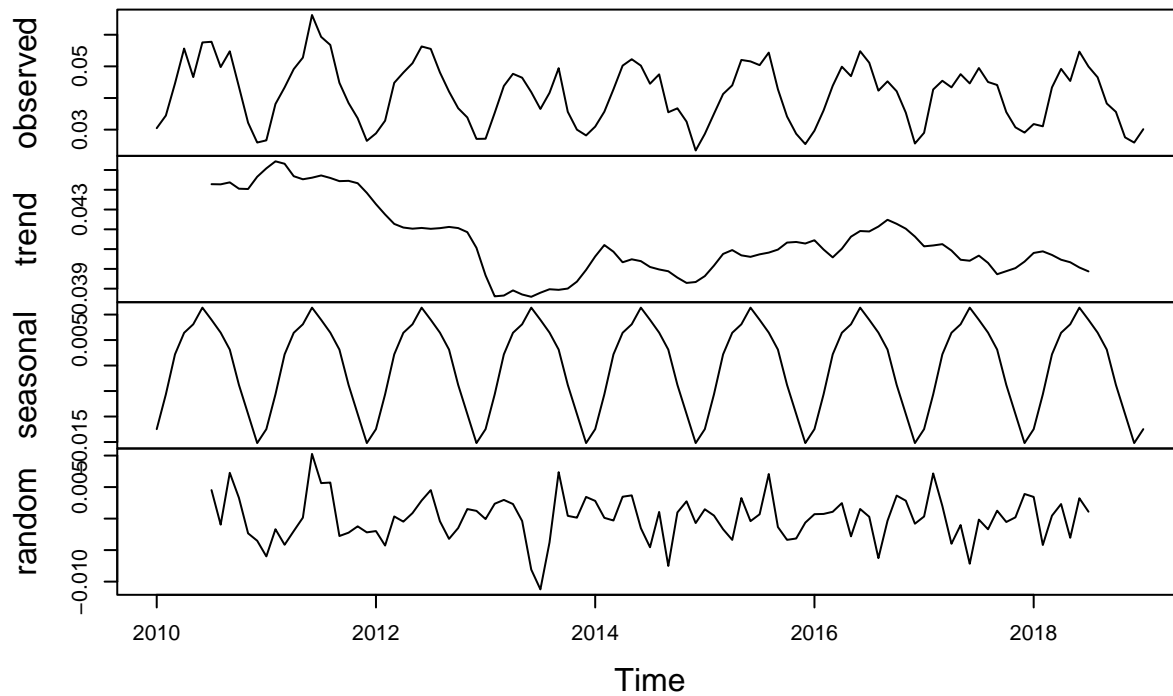
    Answer: We have a significant pvalue of .0176 so we reject the null hypothesis. This means our mean monthly ozone data does not change in the 2010s. We are also given a tau value of -.182 which shows a negative correlation.

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.
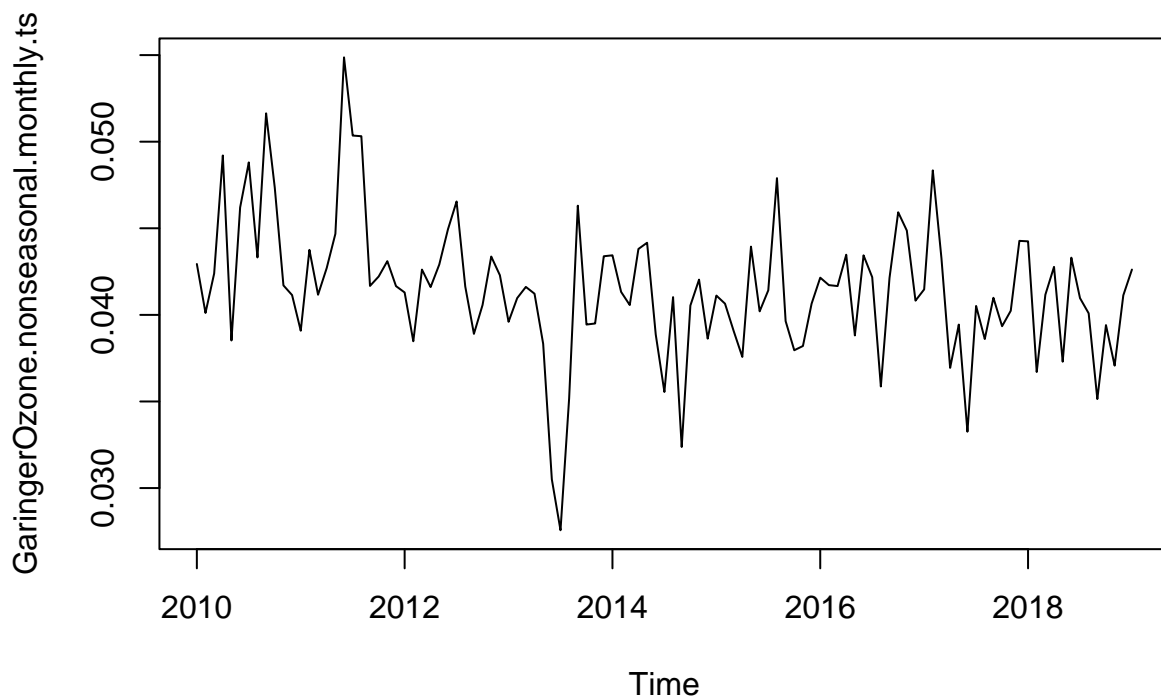
```
#15
GaringerOzone.seasonal.monthly.ts <-
  decompose(GaringerOzone.monthly.ts)$seasonal
plot(decompose(GaringerOzone.monthly.ts))
```

## Decomposition of additive time series



```
GaringerOzone.nonseasonal.monthly.ts<-
  (GaringerOzone.monthly.ts - GaringerOzone.seasonal.monthly.ts)
plot(GaringerOzone.nonseasonal.monthly.ts)
```

```
#16
GaringerOzone.monthly.trend2 <- MannKendall(GaringerOzone.nonseasonal.monthly.ts)
GaringerOzone.monthly.trend2
```

```
## tau = -0.206, 2-sided pvalue =0.0015052
```

```
summary(GaringerOzone.monthly.trend2)
```

```
## Score =  -1213 , Var(Score) = 145841
## denominator =  5885.5
## tau = -0.206, 2-sided pvalue =0.0015052
```

Answer: Both Seasonal and nonseasonal resulted in significant p-values. The nonseasonal p-value is .00150 which is drastically more significant then the P-value of seasonal. So this shows removing seasonality gives us stronger reasoning to reject the null hypothesis. Our Tau value is also lower (more negative) so this should mean our trend is even more negativily correlated now that we removed seasonal data.