

# Assignment 7: Time Series Analysis

Amanda Booth

## 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

#Check working directory

getwd()

## [1] "/Users/amandabooth/Documents/GitHub/Environmental_Data_Analytics_2022/Assignments"

#Load packages

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   1.0.8
## v tidyr   1.2.0      v stringr 1.4.0
## v readr   2.1.2      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)
#install.packages("Kendall")
library(Kendall)

#Build theme
mytheme <- theme_linedraw(base_size = 10) +
  theme(axis.text = element_text(color = "rosybrown"),
        legend.position = "bottom")

#Set as default theme
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

#Import datasets

Garinger2010 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2010_raw.csv"
, stringsAsFactors = TRUE)

Garinger2011 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2011_raw.csv"
, stringsAsFactors = TRUE)

Garinger2012 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2012_raw.csv"
, stringsAsFactors = TRUE)
```

```

Garinger2013 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2013_raw.csv"
, stringsAsFactors = TRUE)

Garinger2014 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2014_raw.csv"
, stringsAsFactors = TRUE)

Garinger2015 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2015_raw.csv"
, stringsAsFactors = TRUE)

Garinger2016 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2016_raw.csv"
, stringsAsFactors = TRUE)

Garinger2017 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2017_raw.csv"
, stringsAsFactors = TRUE)

Garinger2018 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2018_raw.csv"
, stringsAsFactors = TRUE)

Garinger2019 <- read.csv("../Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2019_raw.csv"
, stringsAsFactors = TRUE)

#Combine datasets

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

```

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

# 5

```

```

#New dataframe

Days <- as.data.frame(seq(as.Date("2010/1/1"), as.Date("2019/12/31"), "day"))

#Rename column

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

# 6

GaringerOzone <- left_join(Days, GaringerOzone)

## 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

Plot1 <- ggplot(GaringerOzone, aes(x = Date, y = Daily.Max.8.hour.Ozone.Concentration)) +
  geom_point() +
  labs(y = "Daily ozone concentration") +
  geom_smooth(method = "lm")

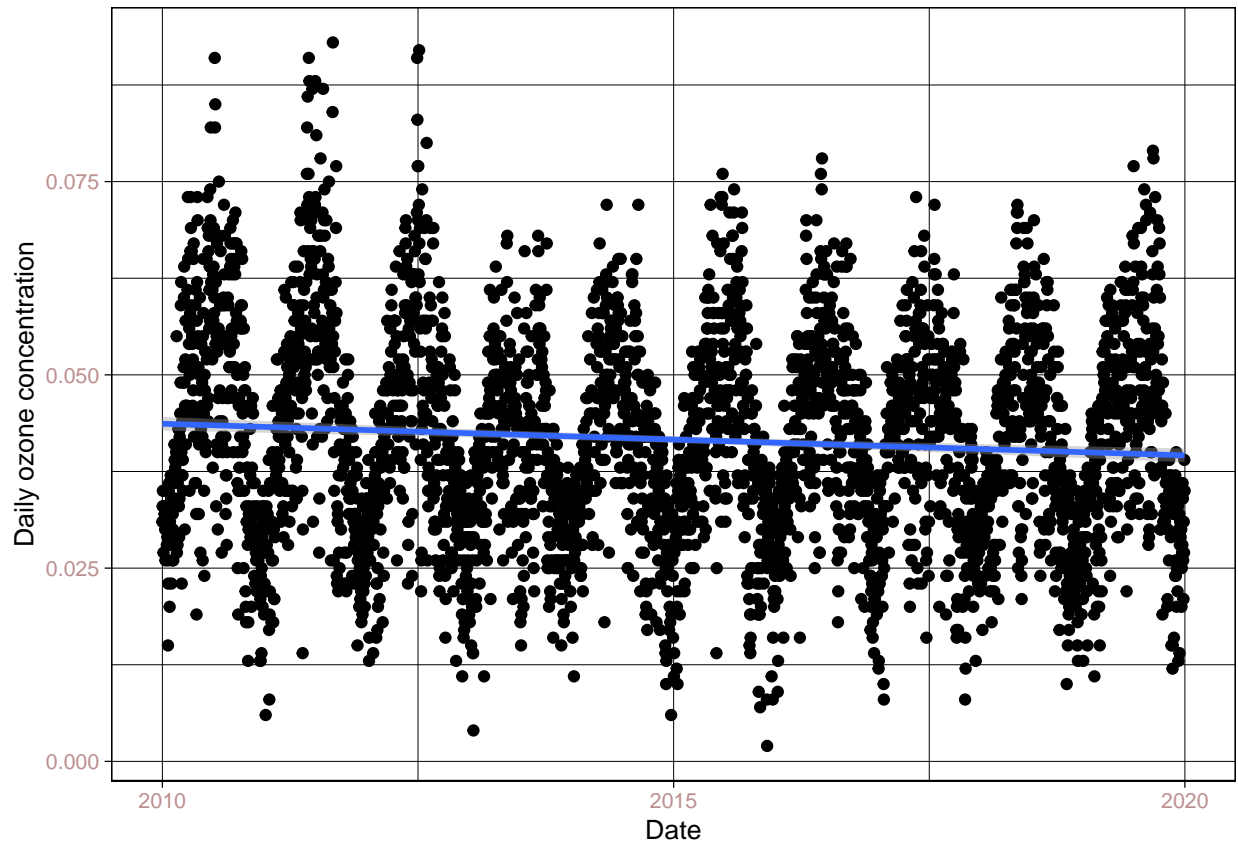
print(Plot1)

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

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

## Warning: Removed 63 rows containing missing values (geom_point).

```



Answer: The plot suggests a very slight decrease in ozone concentration 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

#Remove NAs

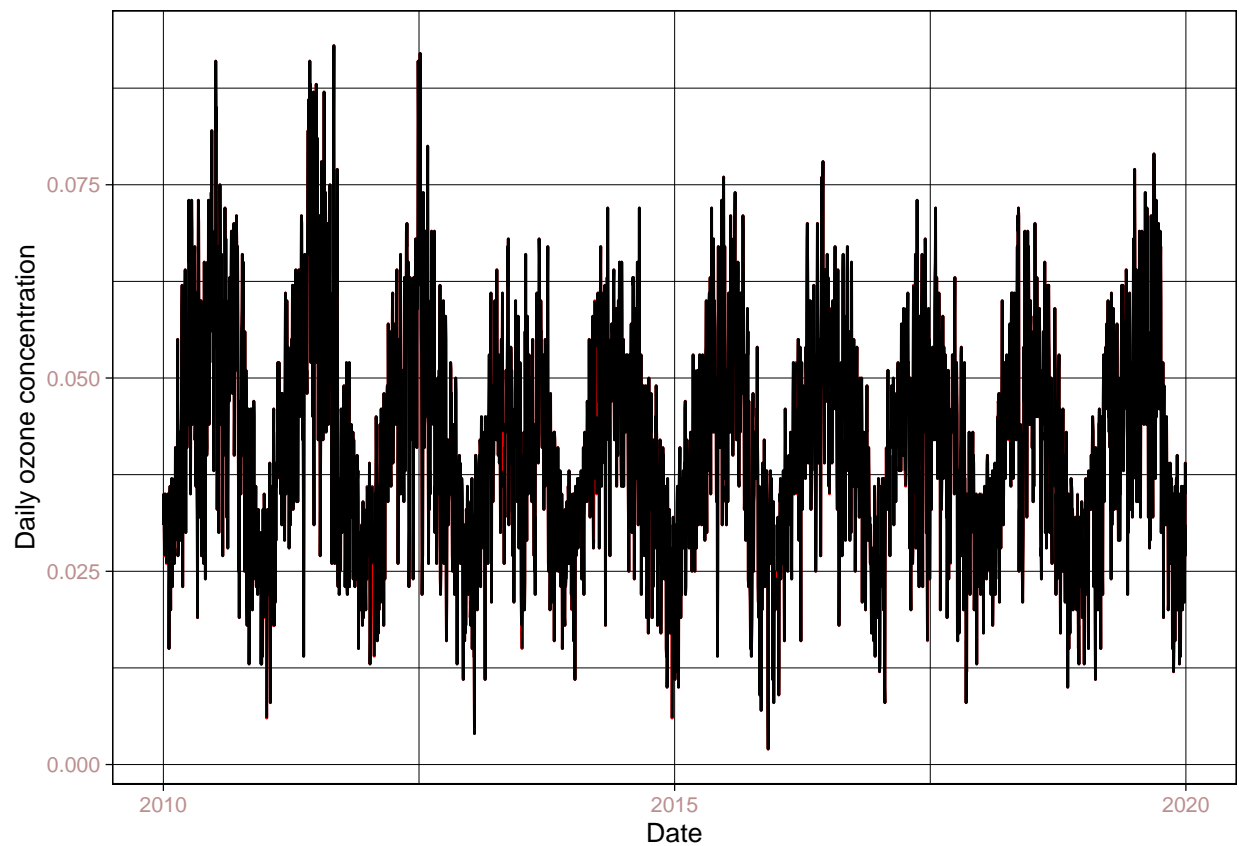
GaringerOzone_clean <-
  GaringerOzone %>%
  mutate(Daily.Max.8.hour.Ozone.Concentration.clean =
    zoo::na.approx(Daily.Max.8.hour.Ozone.Concentration))

summary(GaringerOzone_clean$Daily.Max.8.hour.Ozone.Concentration.clean)

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

```
#Visualize interpolation
```

```
Plot2 <- ggplot(GaringerOzone_clean) +  
  geom_line(aes(x = Date, y = Daily.Max.8.hour.Ozone.Concentration_clean),  
            color = "red") +  
  geom_line(aes(x = Date, y = Daily.Max.8.hour.Ozone.Concentration),  
            color = "black") +  
  ylab("Daily ozone concentration")  
  
print(Plot2)
```



Answer: A piecewise constant would fill in a date's missing data with that of the closest date with data, which might interfere with identifying trends over time. A spline interpolation assumes that the trend is quadratic, which we have no reason to believe.

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(Date_month = my(paste0(Month,"-",Year))) %>%
filter(!is.na(Daily.Max.8.hour.Ozone.Concentration)) %>%
group_by(Date_month) %>%
summarise(Mean = mean(Daily.Max.8.hour.Ozone.Concentration))
```

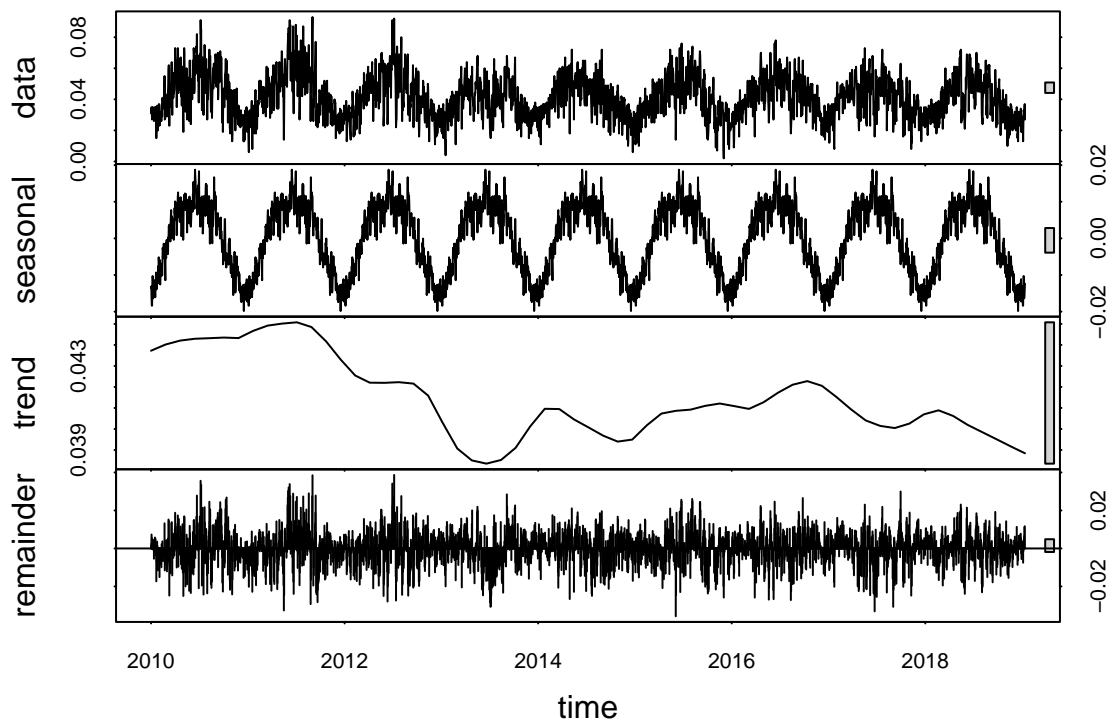
10. Generate two time series objects. Name the first `GaringerOzone.daily.ts` and base it on the dataframe of daily observations. Name the second `GaringerOzone.monthly.ts` and base it on the monthly average ozone values. Be sure that each specifies the correct start and end dates and the frequency of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

```
#10
GaringerOzone.daily.ts <- ts(GaringerOzone_clean$Daily.Max.8.hour.Ozone.Concentration.clean,
                             start=c(2010,01),
                             end= c(2019,12),
                             frequency=365)

GaringerOzone.monthly.ts <- ts(GaringerOzone.monthly$Mean,
                               start=c(2010,01),
                               end= c(2019,12),
                               frequency=12)
```

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

```
#11
#Daily decomposed
Daily_decomposed <- stl(GaringerOzone.daily.ts, s.window = "periodic")
plot(Daily_decomposed)
```

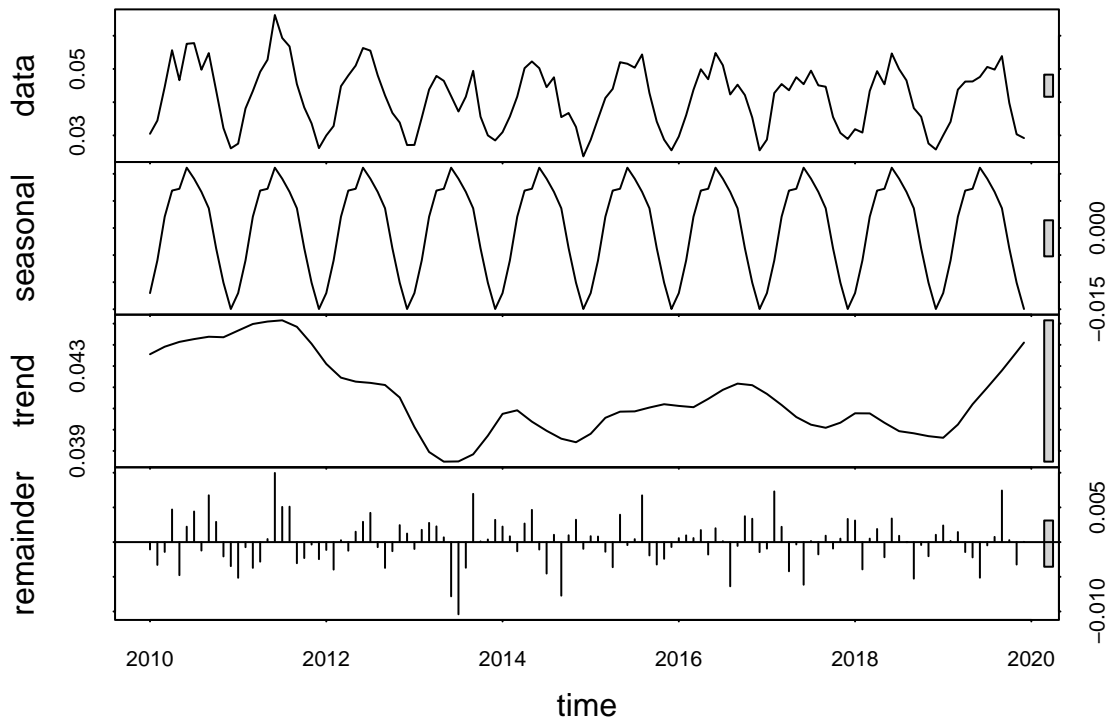


```
#Monthly decomposed
```

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

```
plot(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

#Run test

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

#View results

summary(Monthly_trend)

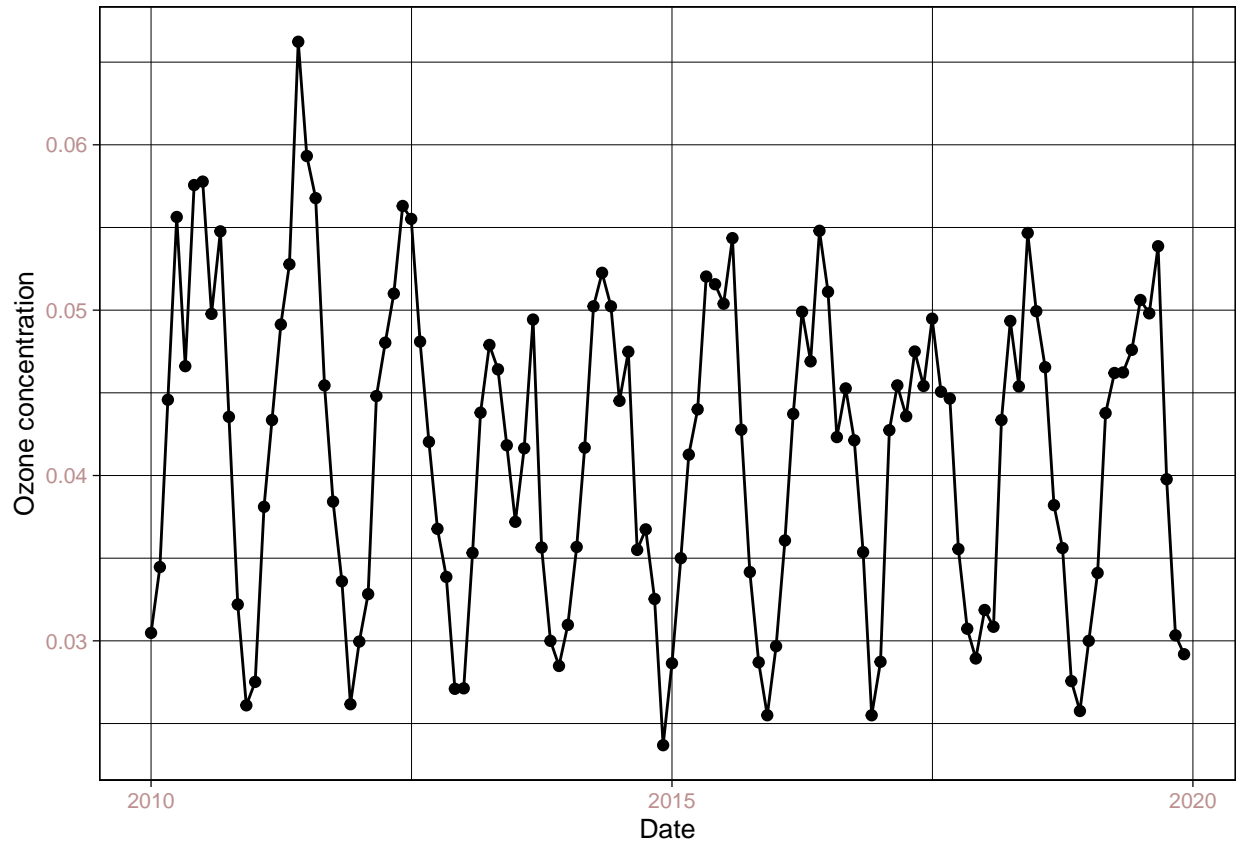
## Score = -88 , Var(Score) = 1498
## denominator = 538.9944
## tau = -0.163, 2-sided pvalue =0.022986
```

Answer: We must account for seasonal changes.

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

```
Plot3 <- ggplot(GaringerOzone.monthly, aes(x = Date_month, y = Mean)) +  
  geom_point() +  
  geom_line() +  
  labs(x = "Date", y = "Ozone concentration")  
  
print(Plot3)
```



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: Assessing the research question, “Have ozone concentrations changed over the 2010s at this station?” we can reject the null hypothesis at a 95% confidence interval and conclude that ozone concentrations changed at a statistically significant rate over the 2010s ( $\tau = -0.163$ , 2-sided  $p\text{-value} = 0.022986$ ). We can also determine that the ozone concentration trends downward.

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

*#Create data frame of components*

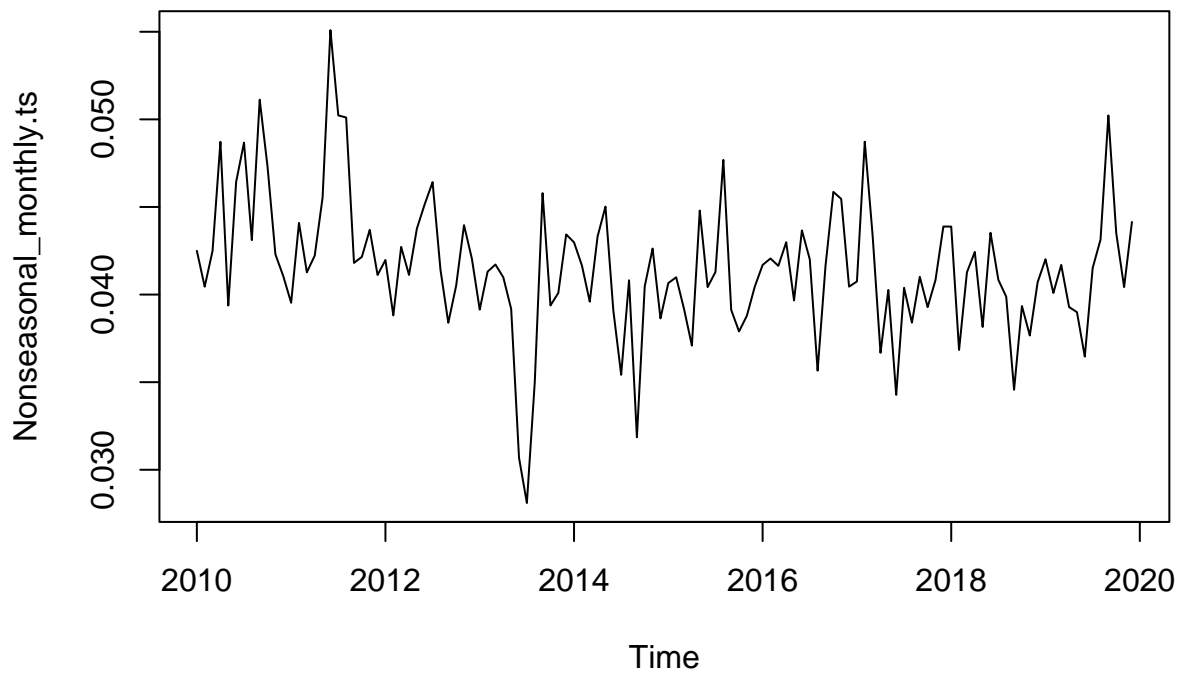
```
Ozone_Components <- as.data.frame(Monthly_decomposed$time.series[,1:3])
```

*#Add observations to data frame*

```
Ozone_Components <- Ozone_Components %>%  
  mutate(Ozone_Components,  
    Observed = GaringerOzone.monthly$Mean,  
    Date = GaringerOzone.monthly$Date_month)
```

*#Subtract seasonal component*

```
Nonseasonal_monthly.ts <- GaringerOzone.monthly.ts - Ozone_Components$seasonal  
plot(Nonseasonal_monthly.ts)
```



#16

```
Monthly_trend2 <- MannKendall(Nonseasonal_monthly.ts)
```

```
summary(Monthly_trend2)
```

```
## Score = -1278 , Var(Score) = 194364.7
## denominator = 7139
## tau = -0.179, 2-sided pvalue =0.0037728
```

Answer: With a p-value of 0.0037728, the non-seasonal Ozone monthly series also rejects the null hypothesis, but can do so at a higher confidence level of 99%. The ozone concentration still trends downward over time, but does so at a somewhat steeper rate of -0.179 per year.