

# Assignment 7: Time Series Analysis

Suad Muradov

## 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

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     vforcats 0.5.1

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

library(dplyr)
library(lubridate)

##
## Attaching package: 'lubridate'

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

```

library(ggplot2)
library(trend)
library(zoo)

##
## Attaching package: 'zoo'

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

library(Kendall)
library(tseries)

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

own.theme<-theme_classic(base_size = 18) +
  theme(axis.text = element_text(color = "black"),
        legend.position = "right", legend.title = element_text(size = 14))
theme_set(own.theme)

```

- 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

EPA.Gar.2010 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2010_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2011 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2011_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2012 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2012_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2013 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2013_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2014 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2014_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2015 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2015_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2016 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2016_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2017 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2017_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2018 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2018_raw.csv",stringsAsFactors=TRUE)
EPA.Gar.2019 <-read.csv("./Data/Raw/Ozone_TimeSeries/EPAair_03_GaringerNC2019_raw.csv",stringsAsFactors=TRUE)

GaringerOzone<-rbind(EPA.Gar.2010, EPA.Gar.2011, EPA.Gar.2012, EPA.Gar.2013, EPA.Gar.2014, EPA.Gar.2015)

```

## Wrangle

- Set your date column as a date class.
- Wrangle your dataset so that it only contains the columns Date, Daily.Max.8.hour.Ozone.Concentration, and DAILY\_AQI\_VALUE.
- 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”.
- 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 Creating a pipe to select 3 columns

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

# 6

GaringerOzone<-left_join(Days, GaringerOzone.processed)

## Joining, by = "Date"
summary(GaringerOzone)

##           Date      Daily.Max.8.hour.Ozone.Concentration DAILY_AQI_VALUE
##   Min.   :2010-01-01   Min.   :0.00200               Min.   : 2.00
##   1st Qu.:2012-07-01   1st Qu.:0.03200               1st Qu.: 30.00
##   Median :2014-12-31   Median :0.04100               Median : 38.00
##   Mean    :2014-12-31   Mean    :0.04163               Mean    : 41.57
##   3rd Qu.:2017-07-01   3rd Qu.:0.05100               3rd Qu.: 47.00
##   Max.    :2019-12-31   Max.    :0.09300               Max.    :169.00
##   NA's    :63          NA's    :63                  NA's    :63
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?

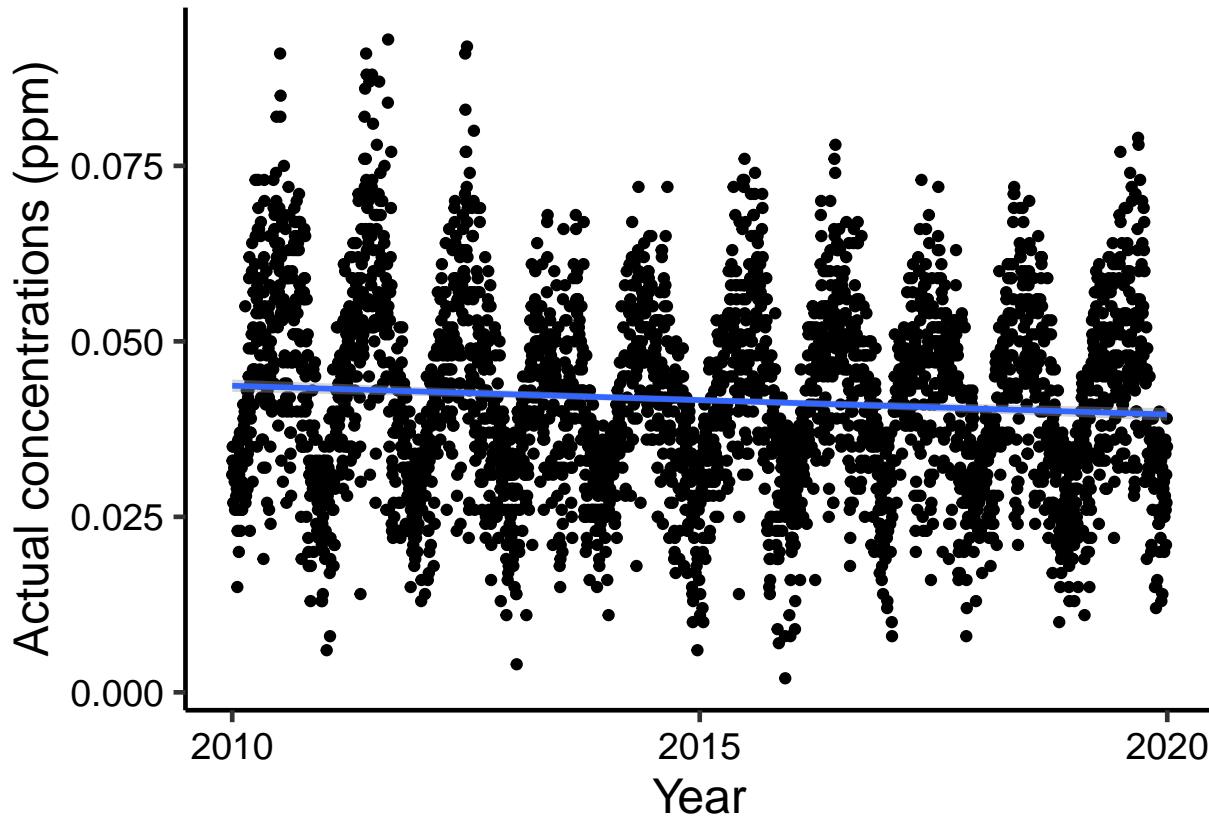
```

#7

GaringerOzone.ppmbydate<-
  ggplot(GaringerOzone, aes(x=Date, y=Daily.Max.8.hour.Ozone.Concentration)) +
  ylab("Actual concentrations (ppm)")+
  xlab("Year")+
  geom_point()+
  geom_smooth(method=lm)
print(GaringerOzone.ppmbydate)

## `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 we have sketched illustrates a sinusoid-shaped pattern that means ozone concentrations are changing at a predictable rate within the certain range.

## 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
```

```
head(GaringerOzone)
```

```
##           Date Daily.Max.8.hour.Ozone.Concentration DAILY_AQI_VALUE
## 1 2010-01-01                      0.031                  29
## 2 2010-01-02                      0.033                  31
## 3 2010-01-03                      0.035                  32
## 4 2010-01-04                      0.031                  29
## 5 2010-01-05                      0.027                  25
## 6 2010-01-06                      NA                    NA
```

```
summary(GaringerOzone)
```

```
##      Date          Daily.Max.8.hour.Ozone.Concentration DAILY_AQI_VALUE
## Min.   :2010-01-01  Min.   :0.00200                  Min.   : 2.00
## 1st Qu.:2012-07-01  1st Qu.:0.03200                  1st Qu.:30.00
## Median :2014-12-31  Median :0.04100                  Median :38.00
## Mean   :2014-12-31  Mean    :0.04163                  Mean   :41.57
## 3rd Qu.:2017-07-01  3rd Qu.:0.05100                  3rd Qu.:47.00
```

```

##   Max.    :2019-12-31    Max.    :0.09300      Max.    :169.00
##                NA's     :63                  NA's     :63
GaringerOzone.clean<-  

  GaringerOzone%>%
  mutate(Daily.Max.8.hour.Ozone.Concentration.clean=zoo::na.approx(Daily.Max.8.hour.Ozone.Concentration))
  mutate(DAILY_AQI_VALUE.clean=zoo::na.approx(DAILY_AQI_VALUE))%>%
  select(Date,Daily.Max.8.hour.Ozone.Concentration.clean,DAILY_AQI_VALUE.clean)

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

```

Answer: We use linear interpolation because it is easier to use and the dataset is tiny and linear which gives more accurate results. Certainly, we would not use linear interpolation with the large dataset which has a complex function (many variables).

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.clean%>%
  mutate(Month=month(Date))%>%
  mutate(Year=year(Date))%>%
  mutate(Day=my(paste0(Month, " - ", Year)))%>%
  dplyr::group_by(Date, Month, Year)%>%
  dplyr::summarise(mean_03=mean(Daily.Max.8.hour.Ozone.Concentration.clean))%>%
  select(mean_03, Date, Month, Year)

```

```

## `summarise()` has grouped output by 'Date', 'Month'. 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

```

f_year<- year(first(GaringerOzone.clean$date))
l_year<- year(last(GaringerOzone.clean$date))
f_month <- month(first(GaringerOzone.clean$date))
l_month<- month(last(GaringerOzone.clean$date))
f_day <- day(first(GaringerOzone.clean$date))
l_day<- day(last(GaringerOzone.clean$date))

GaringerOzone.monthly.ts <- ts(GaringerOzone.clean$Daily.Max.8.hour.Ozone.Concentration.clean,
                                start=c(f_year,f_month), end = c(l_year, l_month),
                                frequency=12)
GaringerOzone.daily.ts <- ts(GaringerOzone.clean$Daily.Max.8.hour.Ozone.Concentration.clean,
                               start=c(f_month,f_day), end = c(l_year, l_month),
                               frequency=365)

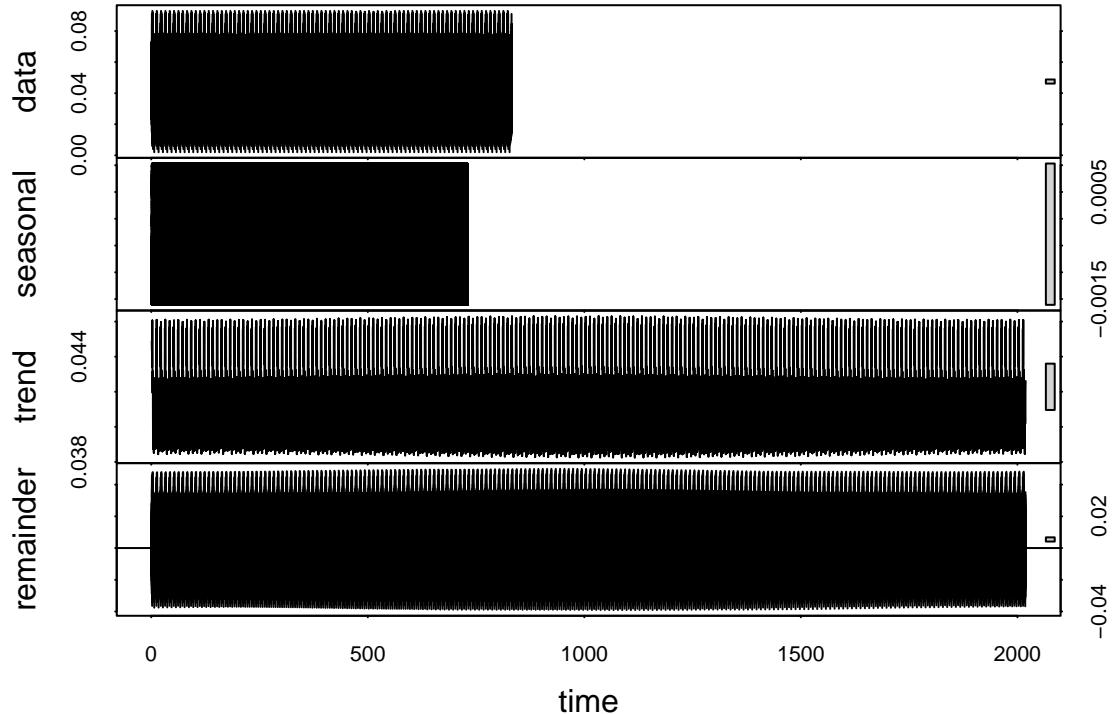
```

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

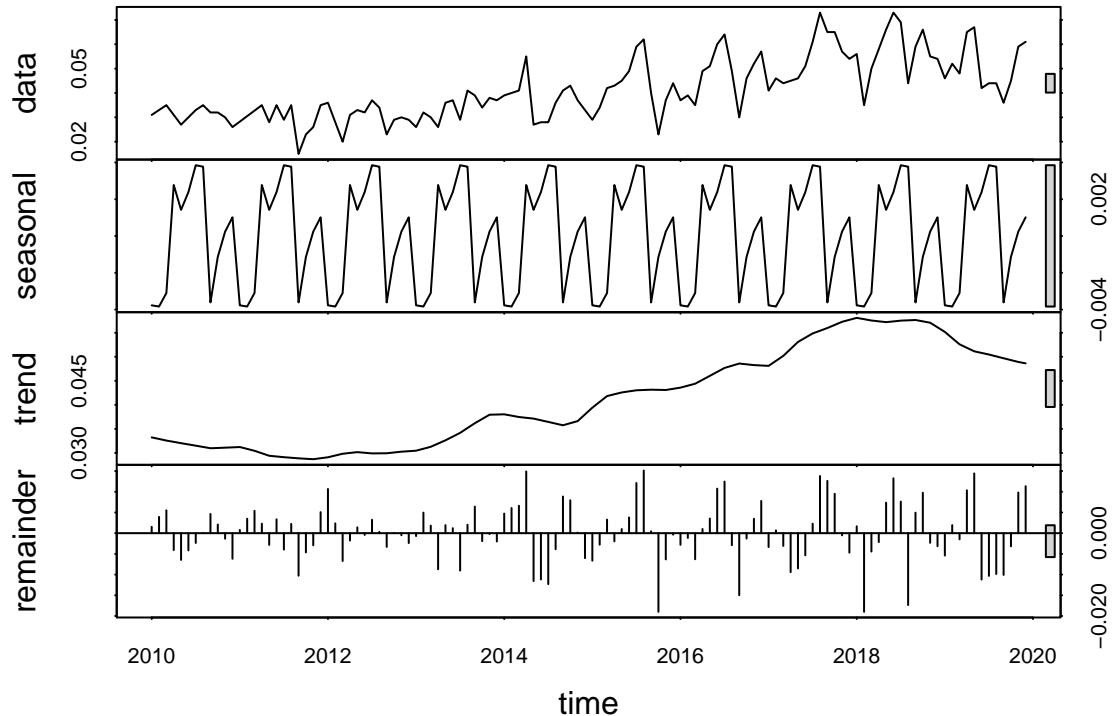
function.

```
#11
```

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



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



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.ts.trend1<-Kendall::SeasonalMannKendall(GaringerOzone.monthly.ts)
summary(GaringerOzone.monthly.ts.trend1)
```

```
## Score = 328 , Var(Score) = 1490
## denominator = 534.9604
## tau = 0.613, 2-sided pvalue <= 2.22e-16
```

```
GaringerOzone.monthly.ts.trend2<-trend::smk.test (GaringerOzone.monthly.ts)
summary(GaringerOzone.monthly.ts.trend2)
```

```
##
## 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  27  125 0.600 2.326 0.0200447 *
## Season 2: S = 0  27  125 0.600 2.326 0.0200447 *
## Season 3: S = 0  28  124 0.629 2.425 0.0153222 *
## Season 4: S = 0  28  124 0.629 2.425 0.0153222 *
## Season 5: S = 0  36  124 0.809 3.143 0.0016717 **
## Season 6: S = 0  25  125 0.556 2.147 0.0318231 *
## Season 7: S = 0  20  124 0.449 1.706 0.0879615 .
## Season 8: S = 0  23  123 0.523 1.984 0.0472923 *
## Season 9: S = 0  19  125 0.422 1.610 0.1074046
## Season 10: S = 0 26  124 0.584 2.245 0.0247639 *
## Season 11: S = 0 35  123 0.796 3.066 0.0021718 **
## Season 12: S = 0 34  124 0.764 2.963 0.0030417 **
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Answer: In our case we cannot exactly expect seasonality from month to month regarding the ozone distributions, because the trends occur in various leanings as we change from one season to another. Seasonal Mann Kendall analysis is therefore helpful to solve this problem.

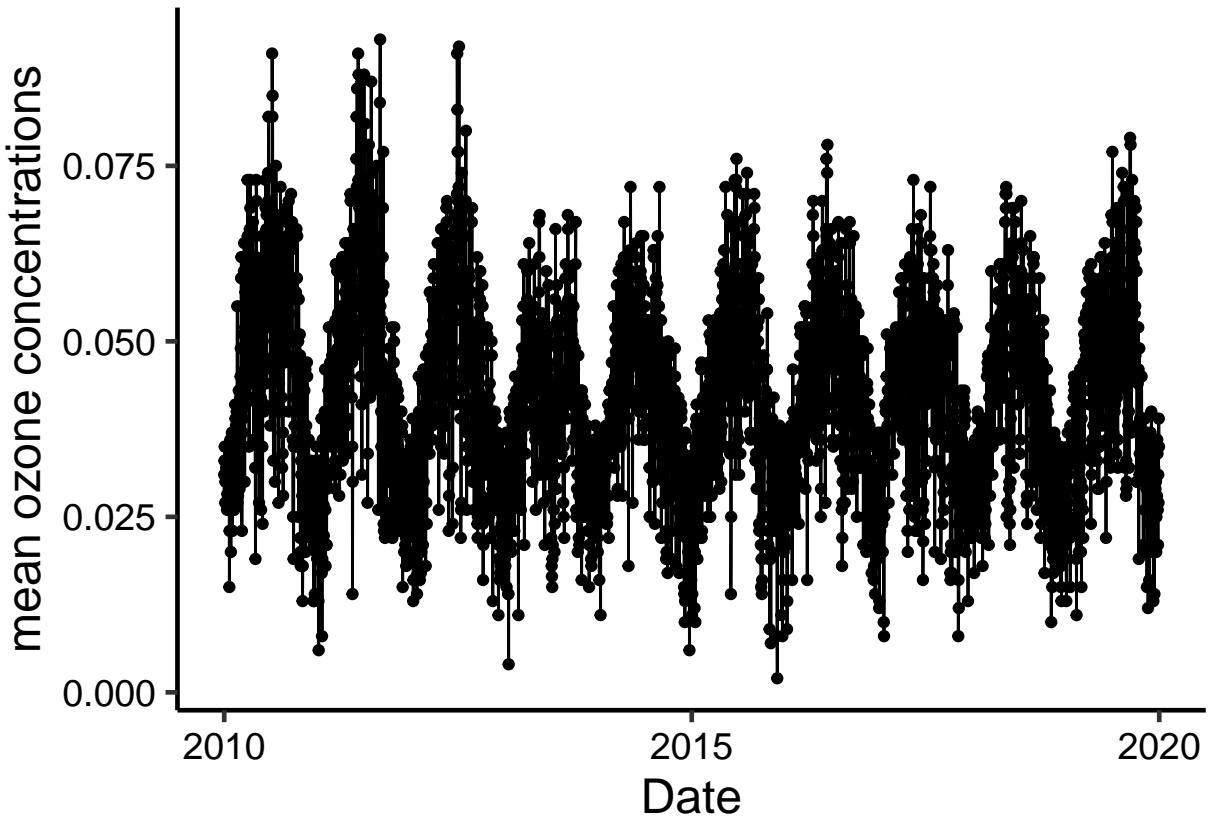
13. Create a plot depicting mean monthly ozone concentrations over time, with both a geom\_point and a geom\_line layer. Edit your axis labels accordingly.

```
# 13
```

```
is.data.frame(GaringerOzone.monthly)
```

```
## [1] TRUE
```

```
GaringerOzone.monthly.plot <-
ggplot(GaringerOzone.monthly, aes(x=Date, y=mean_03)) +
ylab("mean ozone concentrations")+
geom_point()+
geom_line()
print(GaringerOzone.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: Our question was whether ozone concentrations changed over the 2010s, and based on our p-values across individual seasons, we can find that some months have performed statistically significant results that differed from the other years' values: fx, for May (p-value of 0.0016717), November (p-value of 0.0021718) and December (p-value of 0.0030417) our findings indicated that the across years our ozone concentrations have differed.

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
```

```
#16
```

Answer: