

# Assignment 7: Time Series Analysis

Ricky Prophete

## 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] "/Users/ricpro/Documents/Duke MPP Coursework/Spring 2022/Environmental Data Analysis 872/Environm

#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  1.0.7
## 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)

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

library(trend)
#set theme
new_theme <- theme_bw(base_size = 13) +
  theme(axis.text = element_text(color = "black"),
        legend.position = "bottom")
theme_set(new_theme)
```

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
GaringerOzone.Files = list.files(path = "../Data/Raw/Ozone_TimeSeries/", pattern="*.csv",
                                full.names = TRUE)

#check files
GaringerOzone.Files

## [1] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2010_raw.csv"
## [2] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2011_raw.csv"
## [3] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2012_raw.csv"
## [4] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2013_raw.csv"
## [5] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2014_raw.csv"
## [6] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2015_raw.csv"
## [7] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2016_raw.csv"
## [8] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2017_raw.csv"
## [9] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2018_raw.csv"
## [10] "../Data/Raw/Ozone_TimeSeries//EPAair_03_GaringerNC2019_raw.csv"

#combine files
library(plyr)
GaringerOzone <- GaringerOzone.Files %>% ldply(read.csv)
#check
dim(GaringerOzone)

## [1] 3589    20
```

## 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.Columns <- 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"))

#change name
colnames(Days) <- "Date"

# 6
GaringerOzone <- left_join(Days, GaringerOzone.Columns, by = "Date")
#check
dim(GaringerOzone)

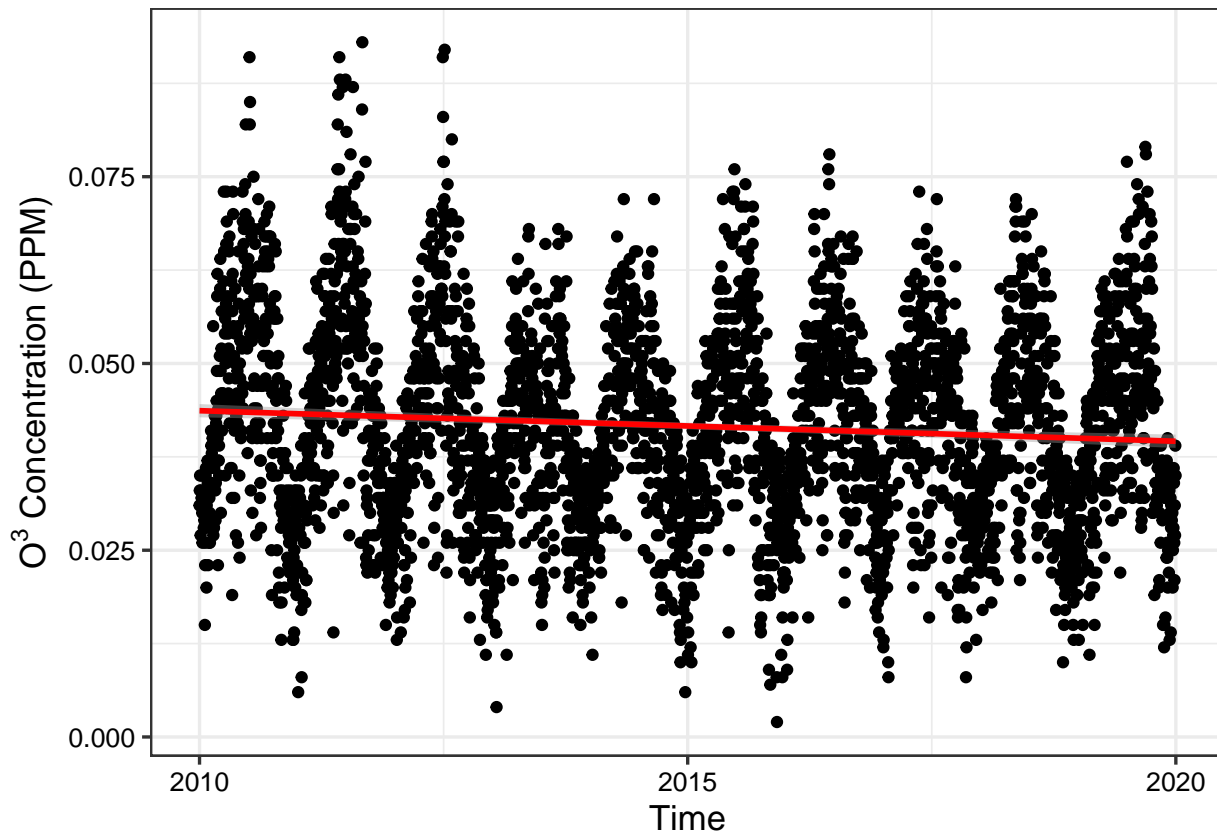
## [1] 3652    3
```

## 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
Plot1 <- ggplot(GaringerOzone, aes(x = Date, y = Daily.Max.8.hour.Ozone.Concentration))+
  geom_point()+
  labs(y = expression("O"~3*" Concentration (PPM)"), x = "Time")+
  geom_smooth(method = "lm", color = "red")
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: There seems to be a slight declining trend 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
GaringerOzone.Clean <- GaringerOzone %>%
  mutate(Daily.Max.8.hour.Ozone.Concentration = zoo::na.approx(Daily.Max.8.hour.Ozone.Concentration))
#check
summary(GaringerOzone$Daily.Max.8.hour.Ozone.Concentration)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.      NA's
## 0.00200 0.03200 0.04100 0.04163 0.05100 0.09300      63

summary(GaringerOzone.Clean$Daily.Max.8.hour.Ozone.Concentration)

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

Answer: A linear interpolation makes sense here as the existing data is generated from observations taken on specific dates and follows a fairly predictable linear pattern between sequential observations. The missing data maps to dates in between dates with recorded observations. A piecewise constant would assume that a missing measurement is equal to an adjacent observation, while a Spline interpolation assumes that the data takes a quadratic form - neither of these assumptions is borne out by the existing data.

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(Date = my(paste0(Month, "-", Year))) %>%
  dplyr::group_by(Date, Year, Month) %>%
  dplyr::summarise(Mean_Ozone = mean(Daily.Max.8.hour.Ozone.Concentration)) %>%
  select(Date, Mean_Ozone)
```

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

## Adding missing grouping variables: ``Year``

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.

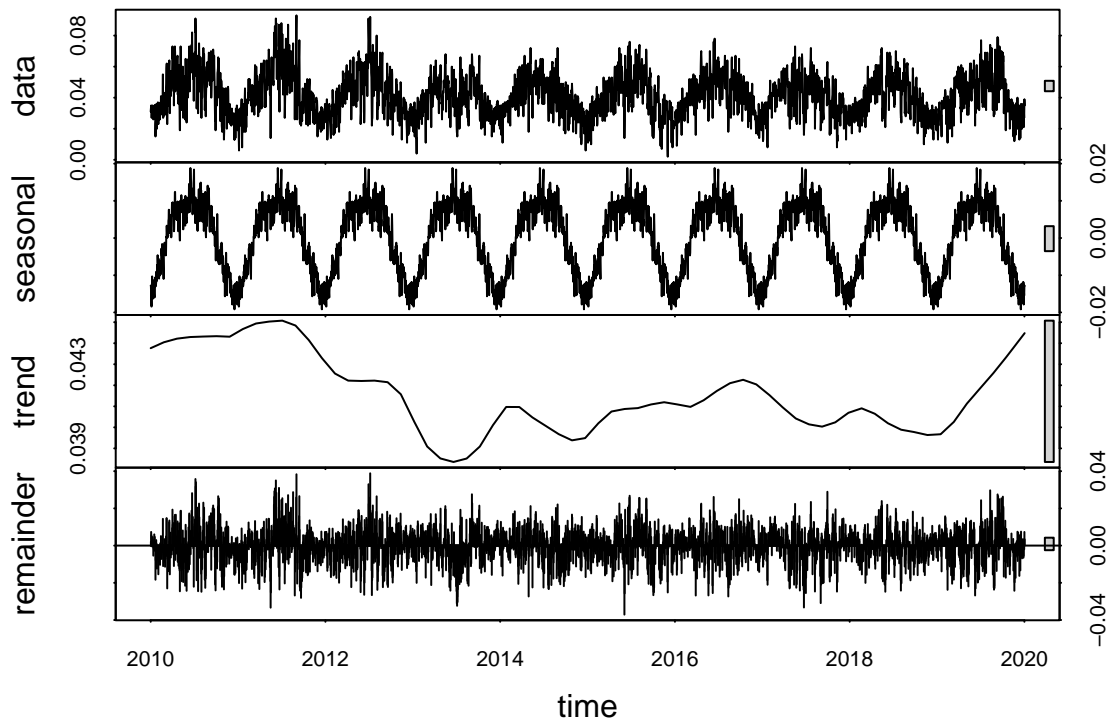
#10

```
GaringerOzone.daily.ts <- ts(GaringerOzone.Clean$Daily.Max.8.hour.Ozone.Concentration, start = c(2010,1),
                             end = c(2014,12), frequency = 365)
GaringerOzone.monthly.ts <- ts(GaringerOzone.monthly$Mean_Ozone, start = c(2010,1), frequency = 12)
```

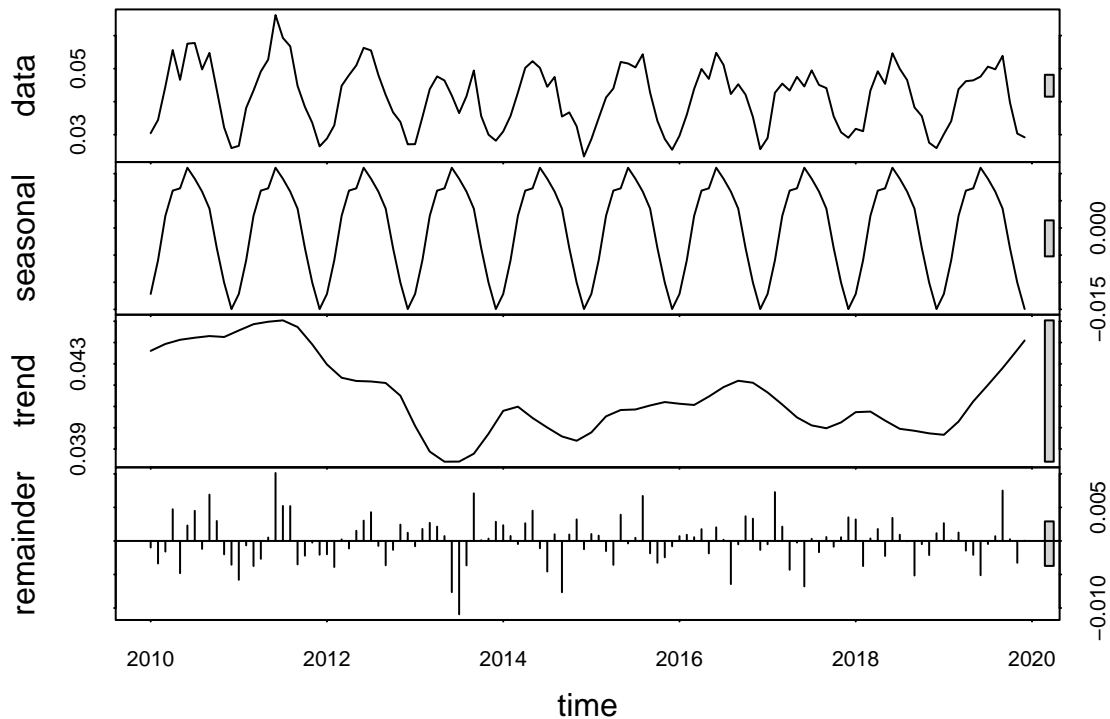
11. Decompose the daily and the monthly time series objects and plot the components using the `plot()` function.

#11

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



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

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

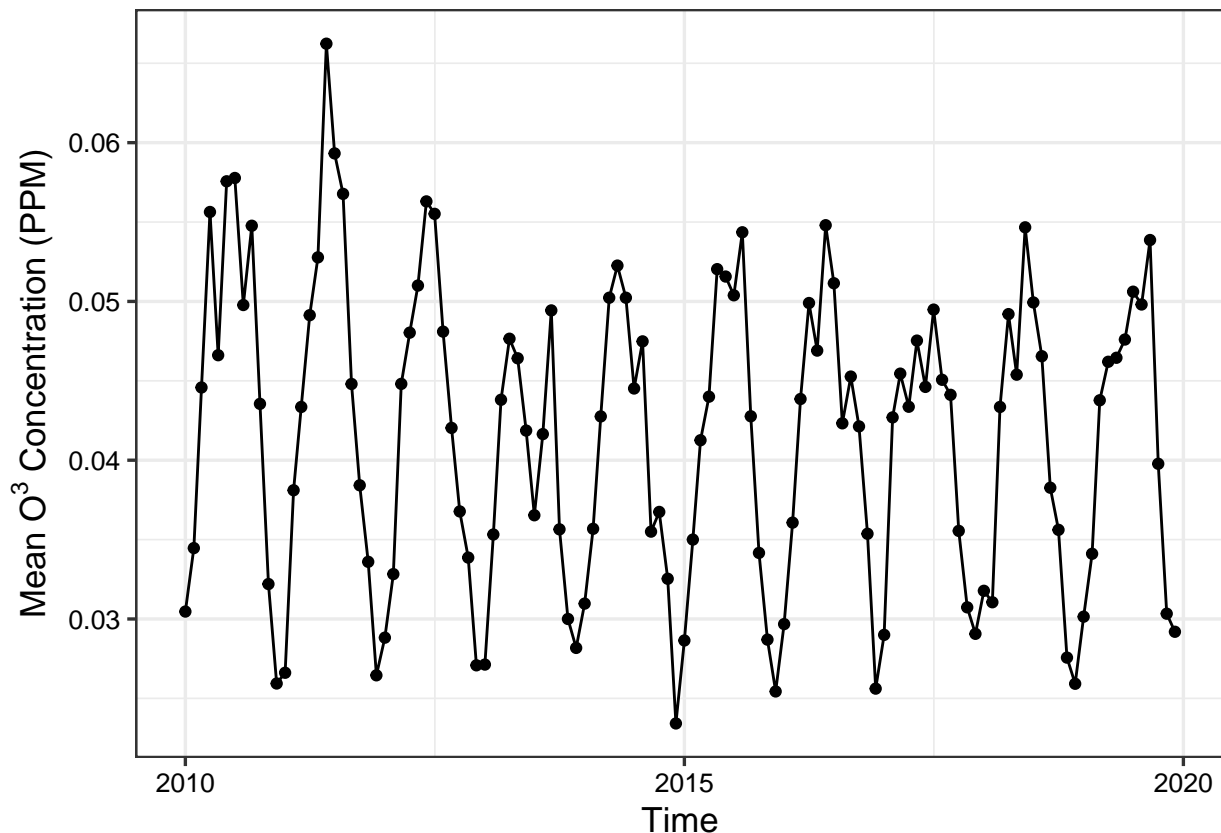
```
GaringerOzone.monthly.trend
```

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

Answer: This is most appropriate as the Ozone dataset exhibits seasonality

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
Plot2 <- ggplot(GaringerOzone.monthly, aes(x = Date, y = Mean_Ozone))+
  geom_point()+
  geom_line()+
  labs(x = "Time", y = expression ("Mean O"3*" Concentration (PPM)"))
print(Plot2)
```



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: There is a declining trend in Ozone concentration over the 2010s. This is supported by our ability to reject the null given the output (p-value <.05) of the Seasonal Mann Kendall test (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
O3.components <- as.data.frame(GaringerOzone.monthly.decomposed$time.series[,1:3])

O3.components <- O3.components %>%
  mutate(O3.components,
         Observed = GaringerOzone.monthly$Mean_Ozone,
         Date = GaringerOzone.monthly$Date)

#remove seasonal component
removed.monthly <- GaringerOzone.monthly.ts - O3.components$seasonal

#16
GraingerOzone.monthly.trend2 <- Kendall::MannKendall(removed.monthly)
summary(GraingerOzone.monthly.trend2)

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

```

Answer: These results point in the same direction, with the Mann Kendall test indicating a trend. The p-value of .0075402 allows us to reject the null hypothesis. I would note that this is a more confident statistical result than was produced by the Seasonal Mann Kendall test.