

Assignment 8: Time Series Analysis

Felipe Raby Amadori

OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics (ENV872L) on time series analysis.

Directions

1. Change “Student Name” on line 3 (above) with your name.
2. Use the lesson as a guide. It contains code that can be modified to complete the assignment.
3. Work through the steps, **creating code and output** that fulfill each instruction.
4. Be sure to **answer the questions** in this assignment document. Space for your answers is provided in this document and is indicated by the “>” character. If you need a second paragraph be sure to start the first line with “>”. You should notice that the answer is highlighted in green by RStudio.
5. When you have completed the assignment, **Knit** the text and code into a single PDF file. You will need to have the correct software installed to do this (see Software Installation Guide) Press the **Knit** button in the RStudio scripting panel. This will save the PDF output in your Assignments folder.
6. After Knitting, please submit the completed exercise (PDF file) to the dropbox in Sakai. Please add your last name into the file name (e.g., “Salk_A08_TimeSeries.pdf”) prior to submission.

The completed exercise is due on Tuesday, 19 March, 2019 before class begins.

Brainstorm a project topic

1. Spend 15 minutes brainstorming ideas for a project topic, and look for a dataset if you are choosing your own rather than using a class dataset. Remember your topic choices are due by the end of March, and you should post your choice ASAP to the forum on Sakai.

Question: Did you do this?

ANSWER: Yes

Set up your session

2. Set up your session. Upload the EPA air quality raw dataset for PM2.5 in 2018, and the processed NTL-LTER dataset for nutrients in Peter and Paul lakes. Build a ggplot theme and set it as your default theme. Make sure date variables are set to a date format.

```
knitr::opts_chunk$set(echo = TRUE, message = FALSE, warning = FALSE, eval=TRUE)
getwd()
```

```
## [1] "C:/Users/Felipe/OneDrive - Duke University/1. DUKE/1. Ramos 2 Semestre/EOS-872 Env. Data Analyt.
```

```
library(FSA)
```

```
## ## FSA v0.8.22. See citation('FSA') if used in publication.
```

```
## ## Run fishR() for related website and fishR('IFAR') for related book.
```

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.2.1 --
```

```

## v ggplot2 3.0.0      v purrr  0.2.5
## v tibble  1.4.2      v dplyr  0.7.6
## v tidyr   0.8.1      v stringr 1.3.1
## v readr   1.1.1      v forcats 0.3.0

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

library(RColorBrewer)
library(ggpubr)

## Loading required package: magrittr

##
## Attaching package: 'magrittr'

## The following object is masked from 'package:purrr':
##
##   set_names

## The following object is masked from 'package:tidyr':
##
##   extract

library(viridis)

## Loading required package: viridisLite

library(colormap)
library(lubridate)

##
## Attaching package: 'lubridate'

## The following object is masked from 'package:base':
##
##   date

library(nlme)

##
## Attaching package: 'nlme'

## The following object is masked from 'package:dplyr':
##
##   collapse

library(lsmeans)

## Loading required package: emmeans

## The 'lsmeans' package is now basically a front end for 'emmeans'.
## Users are encouraged to switch the rest of the way.
## See help('transition') for more information, including how to
## convert old 'lsmeans' objects and scripts to work with 'emmeans'.

library(multcompView)
library(trend)

felipe_theme <- theme_light(base_size = 12) +

```

```

    theme(axis.text = element_text(color = "grey8"),
          legend.position = "right", plot.title = element_text(hjust = 0.5))
theme_set(felipe_theme)

EPA_raw_AQ_PM25_2018 <- read.csv("./Data/Raw/EPAair_PM25_NC2018_raw.csv")
NTL_TER_raw_ChemistryNutr_PeterPaul <-
  read.csv("./Data/Processed/NTL-LTER_Lake_Chemistry_Nutrients_PeterPaul_Processed.csv")

class(EPA_raw_AQ_PM25_2018$Date)

## [1] "factor"

class(NTL_TER_raw_ChemistryNutr_PeterPaul$sampldate)

## [1] "factor"

EPA_raw_AQ_PM25_2018$Date <- as.Date(EPA_raw_AQ_PM25_2018$Date, format = "%m/%d/%y")
NTL_TER_raw_ChemistryNutr_PeterPaul$sampldate <-
  as.Date(NTL_TER_raw_ChemistryNutr_PeterPaul$sampldate, format = "%Y-%m-%d")

```

Run a hierarchical (mixed-effects) model

Research question: Do PM2.5 concentrations have a significant trend in 2018?

3. Run a repeated measures ANOVA, with PM2.5 concentrations as the response, Date as a fixed effect, and Site.Name as a random effect. This will allow us to extrapolate PM2.5 concentrations across North Carolina.

3a. Illustrate PM2.5 concentrations by date. Do not split aesthetics by site.

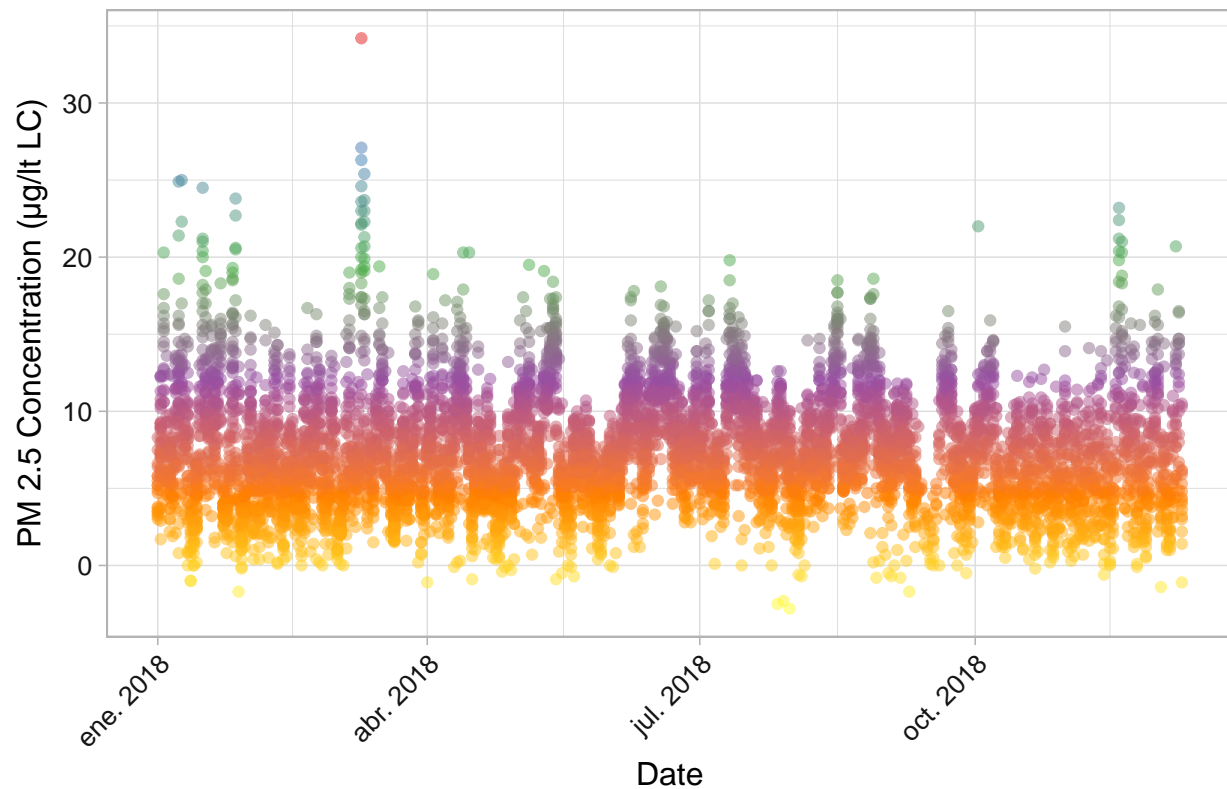
```

#3a.

ggplot(EPA_raw_AQ_PM25_2018,
       aes(x = Date, y = Daily.Mean.PM2.5.Concentration, color = Daily.Mean.PM2.5.Concentration)) +
  geom_point(alpha = 0.5) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
  scale_color_distiller(palette = "Set1") +
  xlab(expression("Date")) +
  ylab(expression(paste("PM 2.5 Concentration (\U003BCg/lt LC)"))) +
  theme(legend.position = "none") +
  ggtitle("2018 Daily PM2.5 concentration")

```

2018 Daily PM2.5 concentration



3b. Insert the following line of code into your R chunk. This will eliminate duplicate measurements on single dates for each site. `PM2.5 = PM2.5[order(PM2.5[, 'Date'], -PM2.5[, 'Site.ID']),]` `PM2.5 = PM2.5[!duplicated(PM2.5$Date),]`

3c. Determine the temporal autocorrelation in your model.

3d. Run a mixed effects model.

```
#3b.
PM2.5 = EPA_raw_AQ_PM25_2018[order(EPA_raw_AQ_PM25_2018[, 'Date'], -EPA_raw_AQ_PM25_2018[, 'Site.ID']),]
PM2.5 = PM2.5[!duplicated(PM2.5$Date),]

#3c.
PM252018.auto <- lme(data = PM2.5, Daily.Mean.PM2.5.Concentration ~ Date,
                     random = ~1|Site.Name)
PM252018.auto

## Linear mixed-effects model fit by REML
##   Data: PM2.5
##   Log-restricted-likelihood: -928.6076
##   Fixed: Daily.Mean.PM2.5.Concentration ~ Date
##   (Intercept)          Date
## 90.465022634 -0.004727976
##
## Random effects:
##   Formula: ~1 | Site.Name
##   (Intercept) Residual
## StdDev:      1.650184 3.559209
```

```
##
## Number of Observations: 343
## Number of Groups: 3
```

```
ACF(PM252018.auto)
```

```
##      lag      ACF
## 1      0 1.000000000
## 2      1 0.513829909
## 3      2 0.194512680
## 4      3 0.117925187
## 5      4 0.126462863
## 6      5 0.100699787
## 7      6 0.058215891
## 8      7 -0.053090104
## 9      8 0.017671857
## 10     9 0.012177847
## 11    10 -0.003699721
## 12    11 -0.020305291
## 13    12 -0.044621086
## 14    13 -0.055602646
## 15    14 -0.065787345
## 16    15 -0.123987593
## 17    16 -0.055414056
## 18    17 0.002911218
## 19    18 0.025133456
## 20    19 -0.015306468
## 21    20 -0.143472007
## 22    21 -0.155495492
## 23    22 -0.060369985
## 24    23 0.003954231
## 25    24 0.042295682
## 26    25 0.001320007
```

```
#3d.
PM252018.mixed <- lme(data = PM2.5, Daily.Mean.PM2.5.Concentration ~ Date,
                      random = ~1|Site.Name,
                      correlation = corAR1(form = ~ Date|Site.Name, value = 0.513829909),
                      method = "REML")
summary(PM252018.mixed)
```

```
## Linear mixed-effects model fit by REML
## Data: PM2.5
##      AIC      BIC    logLik
## 1756.622 1775.781 -873.311
##
## Random effects:
## Formula: ~1 | Site.Name
##      (Intercept) Residual
## StdDev: 0.001025094 3.597269
##
## Correlation Structure: ARMA(1,0)
## Formula: ~Date | Site.Name
## Parameter estimate(s):
##      Phi1
## 0.5384349
```

```
## Fixed effects: Daily.Mean.PM2.5.Concentration ~ Date
##           Value Std.Error   DF   t-value p-value
## (Intercept) 83.14801  60.63585 339   1.371268  0.1712
## Date       -0.00426   0.00342 339  -1.244145  0.2143
## Correlation:
##      (Intr)
## Date -1
##
## Standardized Within-Group Residuals:
##      Min      Q1      Med      Q3      Max
## -2.3220745 -0.6187194 -0.1116751  0.6164257  3.4192603
##
## Number of Observations: 343
## Number of Groups: 3
```

Is there a significant increasing or decreasing trend in PM2.5 concentrations in 2018?

ANSWER: There is not a significant decrease trend in PM2.5 concentrations in 2018 (Linear mixed-effects model, p-value = 0.2143 > 0.05).

3e. Run a fixed effects model with Date as the only explanatory variable. Then test whether the mixed effects model is a better fit than the fixed effect model.

```
PM252018.fixed <- gls(data = PM2.5, Daily.Mean.PM2.5.Concentration ~ Date, method = "REML")
summary(PM252018.fixed)
```

```
## Generalized least squares fit by REML
## Model: Daily.Mean.PM2.5.Concentration ~ Date
## Data: PM2.5
##      AIC      BIC    logLik
## 1865.202 1876.698 -929.6011
##
## Coefficients:
##           Value Std.Error   t-value p-value
## (Intercept) 98.57796  34.60285  2.848840  0.0047
## Date       -0.00513   0.00195 -2.624999  0.0091
##
## Correlation:
##      (Intr)
## Date -1
##
## Standardized residuals:
##      Min      Q1      Med      Q3      Max
## -2.3531000 -0.6348100 -0.1153454  0.6383004  3.4063068
##
## Residual standard error: 3.584321
## Degrees of freedom: 343 total; 341 residual
```

```
anova(PM252018.mixed, PM252018.fixed)
```

```
##           Model df      AIC      BIC    logLik  Test  L.Ratio
## PM252018.mixed   1  5 1756.622 1775.781 -873.3110
## PM252018.fixed   2  3 1865.202 1876.698 -929.6011 1 vs 2 112.5802
##           p-value
## PM252018.mixed
## PM252018.fixed  <.0001
```

Which model is better?

ANSWER: The mixed effects model has a lower AIC value (1756.62 vs 1865.2); hence, it is a better fit than the fixed effect model.

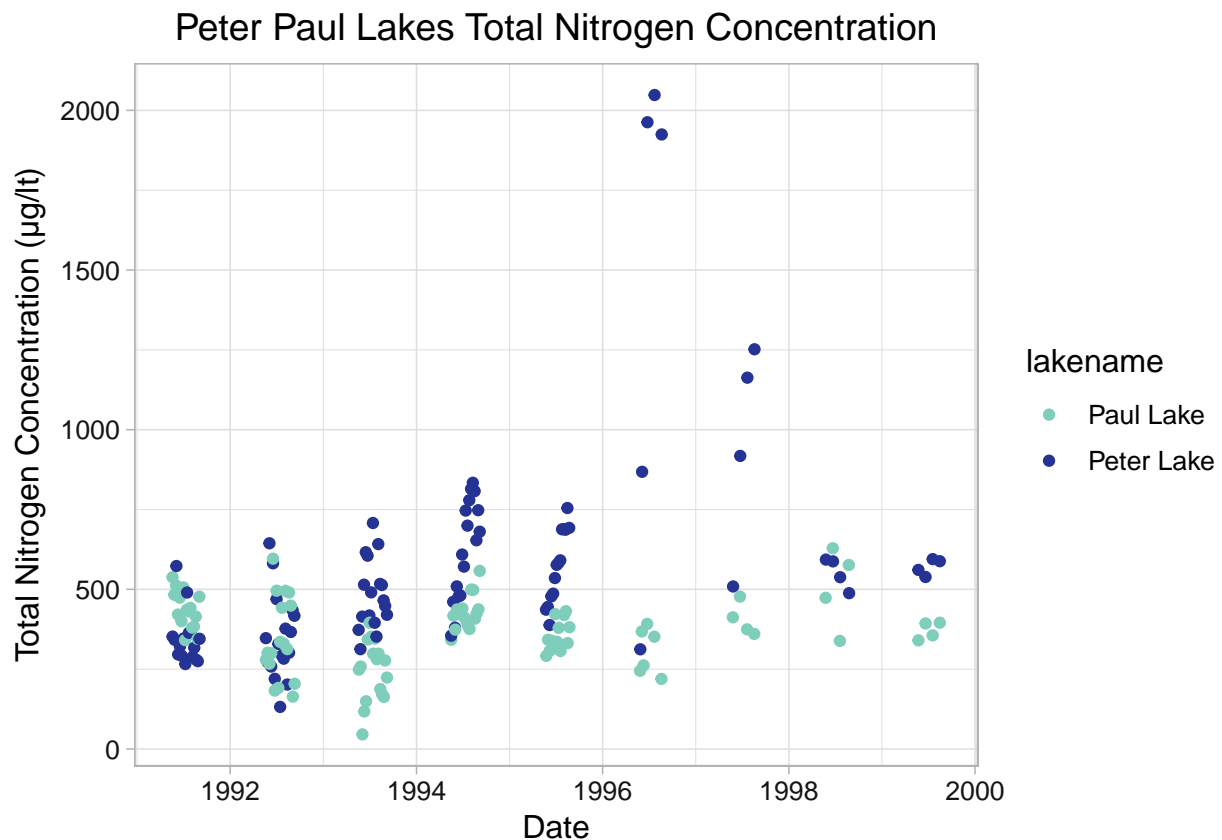
Run a Mann-Kendall test

Research question: Is there a trend in total N surface concentrations in Peter and Paul lakes?

4. Duplicate the Mann-Kendall test we ran for total P in class, this time with total N for both lakes. Make sure to run a test for changepoints in the datasets (and run a second one if a second change point is likely).

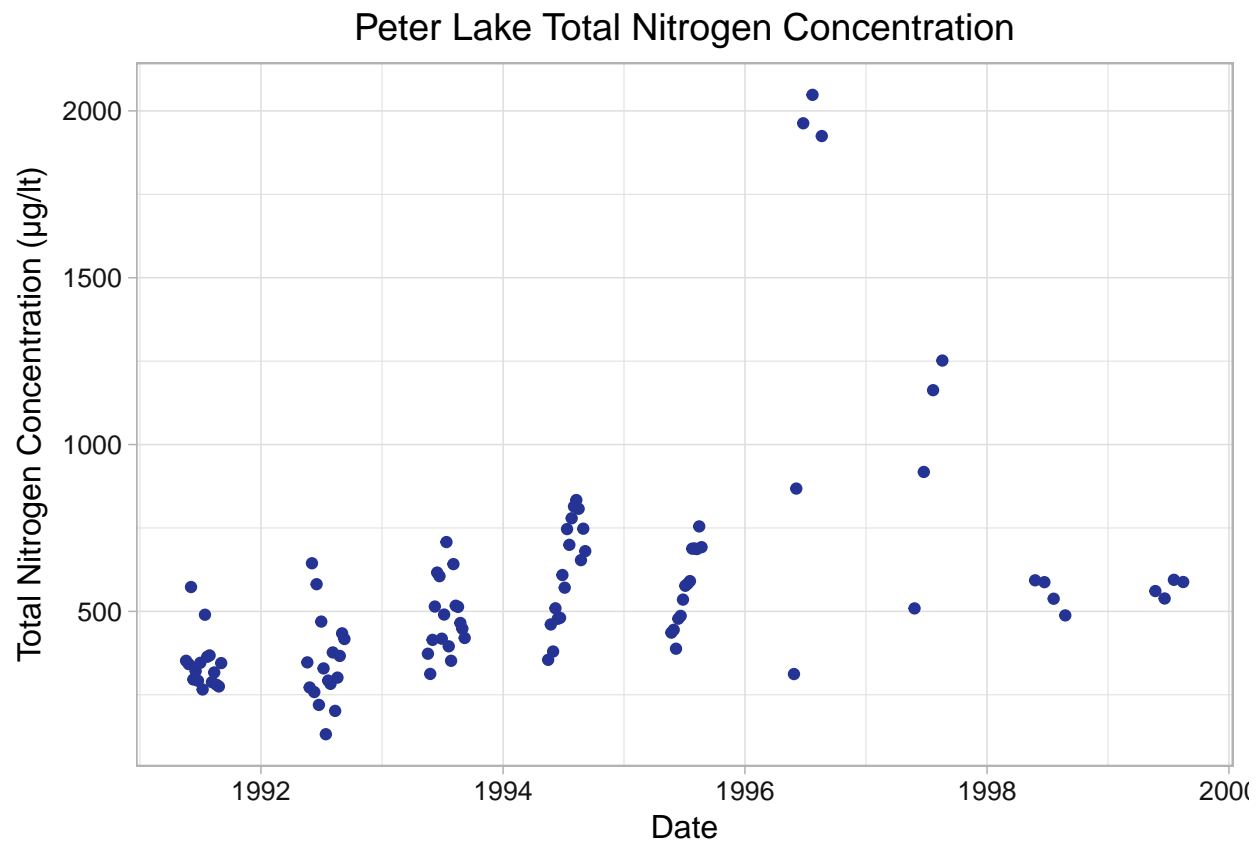
```
# Wrangle our dataset
PeterPaul.tN.surface <-
  NTL_TER_raw_ChemistryNutr_PeterPaul %>%
  select(lakename, sampleddate, depth, tn_ug) %>%
  filter(depth == 0) %>%
  filter(!is.na(tn_ug))

ggplot(PeterPaul.tN.surface, aes(x = sampleddate, y = tn_ug, color = lakename)) +
  geom_point() +
  xlab(expression("Date")) +
  ylab(expression(paste("Total Nitrogen Concentration (\u003BCg/lt)"))) +
  ggtitle("Peter Paul Lakes Total Nitrogen Concentration") +
  scale_color_manual(values = c("#7fcdbb", "#253494"))
```



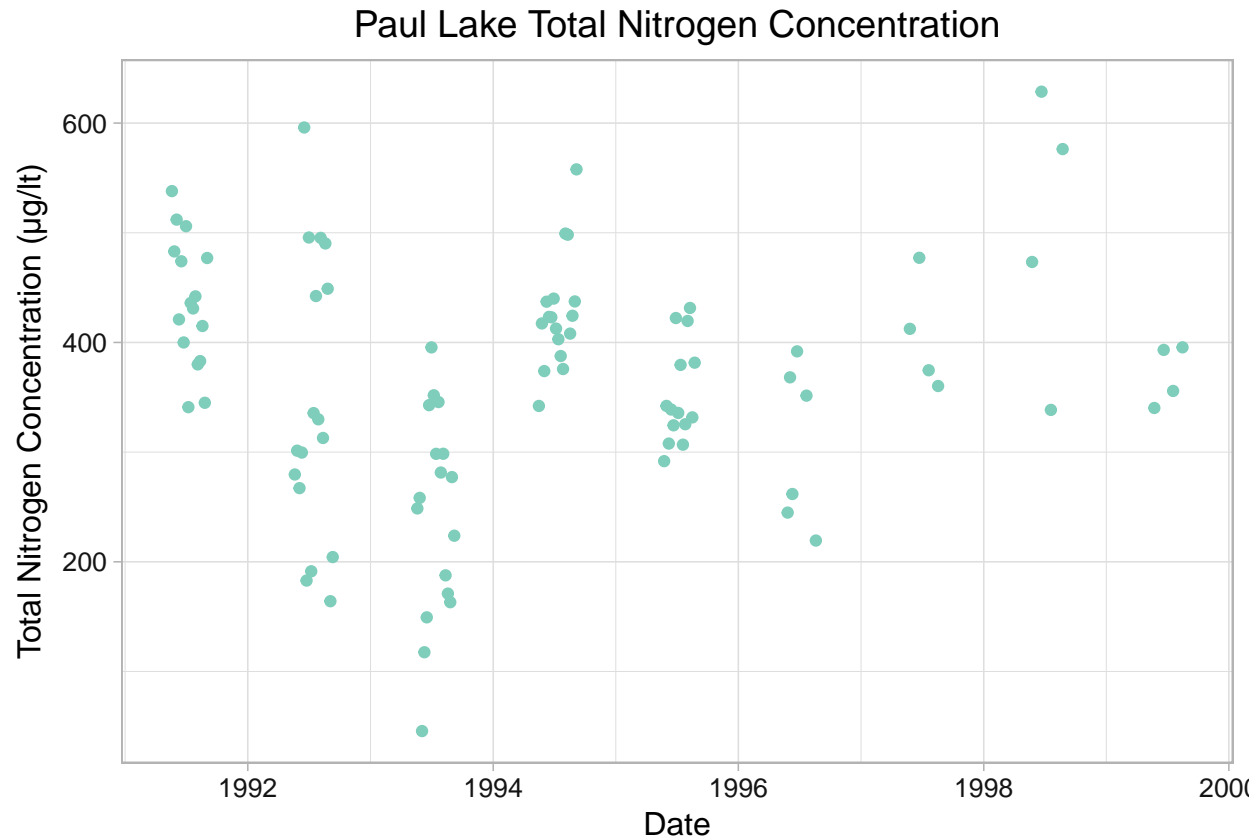
```
# Split dataset by lake
Peter.tN.surface <- filter(PeterPaul.tN.surface, lakename == "Peter Lake")

ggplot(Peter.tN.surface, aes(x = sampleddate, y = tn_ug)) +
  geom_point(color = c("#253494")) +
  xlab(expression("Date")) +
  ylab(expression(paste("Total Nitrogen Concentration (\u003BCg/l\u003C)"))) +
  ggtitle("Peter Lake Total Nitrogen Concentration")
```



```
Paul.tN.surface <- filter(PeterPaul.tN.surface, lakename == "Paul Lake")

ggplot(Paul.tN.surface, aes(x = sampleddate, y = tn_ug)) +
  geom_point(color = c("#7fcdbb")) +
  xlab(expression("Date")) +
  ylab(expression(paste("Total Nitrogen Concentration (\u003BCg/l\u003C)"))) +
  ggtitle("Paul Lake Total Nitrogen Concentration")
```

```
# Run a Mann-Kendall test for Peter Lake
```

```
mk.test(Peter.tN.surface$tn_ug)
```

```
##
## Mann-Kendall trend test
##
## data: Peter.tN.surface$tn_ug
## z = 7.2927, n = 98, p-value = 3.039e-13
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##      S      varS      tau
## 2.377000e+03 1.061503e+05 5.001052e-01
```

```
# There is a significant trend in total N surface concentrations in Peter
# lake (Mann-Kendall trend test, z = 7.2927, n = 98, p-value = 3.039e-13 < 0.05).
```

```
# Run a Pettitt's Test for changepoints in the datasets.
```

```
pettitt.test(Peter.tN.surface$tn_ug)
```

```
##
## Pettitt's test for single change-point detection
##
## data: Peter.tN.surface$tn_ug
## U* = 1884, p-value = 3.744e-10
## alternative hypothesis: two.sided
## sample estimates:
```

```

## probable change point at time K
##                                     36

# There is a significant changepoint in total N surface concentrations in Peter lake.
# Probable change point at time K 36 = 06-02-1993
 #(Pettitt's test,  $U^* = 1884$ ,  $p\text{-value} = 3.744e-10 < 0.05$ ).

# Run separate Mann-Kendall for each changepoint segment

mk.test(Peter.tN.surface$tn_ug[1:35])

##
## Mann-Kendall trend test
##
## data: Peter.tN.surface$tn_ug[1:35]
## z = -0.22722, n = 35, p-value = 0.8203
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S          varS          tau
## -17.00000000 4958.33333333 -0.02857143

# There is not a significant trend in total N surface concentrations in Peter
# lake between 05-20-1991 and 05-26-1993
# (Mann-Kendall trend test,  $z = -0.22722$ ,  $n = 35$ ,  $p\text{-value} = 0.8203 > 0.05$ ).

mk.test(Peter.tN.surface$tn_ug[36:98])

##
## Mann-Kendall trend test
##
## data: Peter.tN.surface$tn_ug[36:98]
## z = 3.1909, n = 63, p-value = 0.001418
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S          varS          tau
## 5.390000e+02 2.842700e+04 2.759857e-01

# There is not a significant trend in total N surface concentrations in Peter
# lake between 06-02-1993 and 08-16-1999
# (Mann-Kendall trend test,  $z = 3.1909$ ,  $n = 63$ ,  $p\text{-value} = 0.001418 < 0.05$ ).

# Is there a second change point?
pettitt.test(Peter.tN.surface$tn_ug[1:35])

##
## Pettitt's test for single change-point detection
##
## data: Peter.tN.surface$tn_ug[1:35]
##  $U^* = 72$ ,  $p\text{-value} = 0.9879$ 
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                                     30

# There is not a significant changepoint in total N surface concentrations in Peter lake
# between 05-20-1991 and 05-26-1993. (Pettitt's test,  $U^* = 72$ ,  $p\text{-value} = 0.9879$ ).

```

```

pettitt.test(Peter.tN.surface$tn_ug[36:98])

##
## Pettitt's test for single change-point detection
##
## data: Peter.tN.surface$tn_ug[36:98]
## U* = 560, p-value = 0.001213
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                                21
# There is a significant changepoint in total N surface concentrations in Peter
# lake between 06-02-1993 and 08-16-1999. Probable change point at time
# K 21 = 06-22-1994 (Pettitt's test, U* = 560, p-value = 0.001213 < 0.05).

# Run another Mann-Kendall for the second change point
mk.test(Peter.tN.surface$tn_ug[36:55])

##
## Mann-Kendall trend test
##
## data: Peter.tN.surface$tn_ug[36:55]
## z = -1.2004, n = 20, p-value = 0.23
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##      S  varS  tau
## -38.0 950.0 -0.2
# There is not a significant trend in total N surface concentrations in Peter
# lake between 06-02-1993 and 06-15-1994
# (Mann-Kendall trend test, z = -1.2004, n = 20, p-value = 0.23 > 0.05).

mk.test(Peter.tN.surface$tn_ug[56:98])

##
## Mann-Kendall trend test
##
## data: Peter.tN.surface$tn_ug[56:98]
## z = 0.48141, n = 43, p-value = 0.6302
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##      S      varS      tau
## 4.700000e+01 9.130333e+03 5.204873e-02
# There is not a significant trend in total N surface concentrations in Peter
# lake between 06-22-1994 and 08-16-1999
# (Mann-Kendall trend test, z = 0.48141, n = 43, p-value = 0.6302 > 0.05).

# Is there a third change point?
pettitt.test(Peter.tN.surface$tn_ug[36:55])

##
## Pettitt's test for single change-point detection
##

```

```

## data: Peter.tN.surface$tn_ug[36:55]
## U* = 42, p-value = 0.5673
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                                     12

# There is not a significant changepoint in total N surface concentrations in Peter lake
# between 06-02-1993 and 06-15-1994. (Pettitt's test, U* = 42, p-value = 0.5673).

pettitt.test(Peter.tN.surface$tn_ug[56:98])

##
## Pettitt's test for single change-point detection
##
## data: Peter.tN.surface$tn_ug[56:98]
## U* = 128, p-value = 0.5974
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                                     21

# There is not a significant changepoint in total N surface concentrations in Peter lake
# between 06-22-1994 and 08-16-1999. (Pettitt's test, U* = 42, p-value = 0.5673).

# In Peter lake data visually it can be seen a decreasing trend after point 84 (1996-06-25),
# so we check for that.

mk.test(Peter.tN.surface$tn_ug[56:84])

##
## Mann-Kendall trend test
##
## data: Peter.tN.surface$tn_ug[56:84]
## z = 0.69405, n = 29, p-value = 0.4877
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S          varS          tau
## 3.800000e+01 2.842000e+03 9.359606e-02

# There is not a significant trend in total N surface concentrations in Peter lake between 06-22-1994
# and 06-25-1996 (Mann-Kendall trend test, z = 0.69405, n = 29, p-value = 0.4877 > 0.05).

mk.test(Peter.tN.surface$tn_ug[84:98])

##
## Mann-Kendall trend test
##
## data: Peter.tN.surface$tn_ug[84:98]
## z = -2.2764, n = 15, p-value = 0.02282
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S          varS          tau
## -47.000000 408.333333 -0.447619

# There is a significant trend in total N surface concentrations in Peter lake between
# 06-25-1996 and 08-16-1999 (Mann-Kendall trend test, z = -2.2764, n = 15, p-value = 0.02282 < 0.05).

```

```

# Run a Mann-Kendall test for Paul Lake

mk.test(Paul.tN.surface$tn_ug)

##
## Mann-Kendall trend test
##
## data: Paul.tN.surface$tn_ug
## z = -0.1572, n = 99, p-value = 0.8751
## alternative hypothesis: true S is not equal to 0
## sample estimates:
##          S          varS          tau
## -5.300000e+01  1.094170e+05 -1.092558e-02

# There is not a significant trend in total N surface concentrations in Paul
# lake (Mann-Kendall trend test, z = -0.1572, n = 99, p-value = 0.8751 > 0.05).

# Run a Pettitt's Test for changepoints in the datasets.

pettitt.test(Paul.tN.surface$tn_ug)

##
## Pettitt's test for single change-point detection
##
## data: Paul.tN.surface$tn_ug
## U* = 704, p-value = 0.09624
## alternative hypothesis: two.sided
## sample estimates:
## probable change point at time K
##                                16

# There is not a significant changepoint in total N surface concentrations in Paul lake.
# (Pettitt's test, U* = 704, p-value = 0.09624 > 0.05).

```

What are the results of this test?

ANSWER: In Peter Lake, there is a significant increasing trend in total N surface concentrations between 05-20-1991 and 08-16-1999 (Mann-Kendall trend test, $z = 7.2927$, $n = 98$, $p\text{-value} = 3.039e-13 < 0.05$). There is also a significant changepoint in total N surface concentrations in Peter lake at time K 36 = 06-02-1993 (Pettitt's test, $U = 1884$, $p\text{-value} = 3.744e-10 < 0.05$), and a second one at time K 21 = 06-22-1994 (Pettitt's test, $U = 560$, $p\text{-value} = 0.001213 < 0.05$).

There is also a significant decreasing trend in total N surface concentrations between 06-25-1996 and 08-16-1999 (Mann-Kendall trend test, $z = -2.2764$, $n = 15$, $p\text{-value} = 0.02282 < 0.05$), which was identified visually.

In Paul Lake, there is not a significant trend in total N surface concentrations (Mann-Kendall trend test, $z = -0.1572$, $n = 99$, $p\text{-value} = 0.8751 > 0.05$).

5. Generate a graph that illustrates the TN concentrations over time, coloring by lake and adding vertical line(s) representing changepoint(s).

```

ggplot(PeterPaul.tN.surface, aes(x = sampleddate, y = tn_ug, color = lakename)) +
  geom_point() +
  xlab(expression("Date")) +
  ylab(expression(paste("Total Nitrogen Concentration (\u003BCg/lt)"))) +
  ggtitle("Peter Paul Lakes Total Nitrogen Concentration") +

```

```
scale_color_manual(values = c("#7fcdbb", "#253494")) +  
geom_vline(xintercept=as.Date("1993/06/02"), color= "#253494", lty=2) +  
geom_vline(xintercept=as.Date("1994/06/22"), color= "#253494", lty=2)
```

