

Assignment 5: Data Visualization

Xin Zhang

OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics (ENV872L) on data wrangling.

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_A04_DataWrangling.pdf”) prior to submission.

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

Set up your session

1. Set up your session. Upload the NTL-LTER processed data files for chemistry/physics for Peter and Paul Lakes (tidy and gathered), the USGS stream gauge dataset, and the EPA Ecotox dataset for Neonicotinoids.
2. Make sure R is reading dates as date format, not something else (hint: remember that dates were an issue for the USGS gauge data).

```
#1
```

```
getwd()
```

```
## [1] "C:/Users/Xin Zhang/Desktop/EDA"
```

```
library(tidyverse)
```

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

```
## v ggplot2 3.1.0      v purrr    0.3.0  
## v tibble   2.0.1      v dplyr    0.7.8  
## v tidyverse 0.8.2      v stringr  1.3.1  
## v readr    1.3.1      v forcats  0.3.0
```

```
## -- Conflicts ----- tidyverse_conflicts
```

```
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()   masks stats::lag()
```

```
library(viridis)
```

```
## Loading required package: viridisLite
```

```
library(RColorBrewer)
```

```
library(colormap)
```

```
PeterPaul.chem.phys <- read.csv("./Data/Processed/NTL-LTER_Lake_Chemistry_Nutrients_PeterPaul_Processed.csv")
```

```

PeterPaul.nutrients.gathered <- read.csv("./Data/Processed/NTL-LTER_Lake_Nutrients_PeterPaulGathered_Pro")
Stream.guage <- read.csv("./Data/Raw/USGS_Site02085000_Flow_Raw.csv")
Ecotox <- read.csv("./Data/Raw/ECOTOX_Neonicotinoids_Mortality_raw.csv")

#2
#check classes for PeterPaul Lakes data and USGS data
class(PeterPaul.chem.phys$sampleddate)

## [1] "factor"
class(PeterPaul.nutrients.gathered$sampleddate)

## [1] "factor"
class(Stream.guage$datetime)

## [1] "factor"

#reformat Peter and Pual Lakes date data
PeterPaul.chem.phys$sampleddate <- as.Date(PeterPaul.chem.phys$sampleddate, format = "%m/%d/%y")
PeterPaul.nutrients.gathered$sampleddate <- as.Date(PeterPaul.nutrients.gathered$sampleddate, format = "%m/%d/%y")
#reformat USGS date data
Stream.guage$datetime <- as.Date(Stream.guage$datetime, format = "%m/%d/%y")
Stream.guage$datetime <- format(Stream.guage$datetime, format = "%Y%m%d")
create.early.dates <- (function(d) {
  paste0(ifelse(d > 181231, "19", "20"), d)
})
Stream.guage$datetime <- create.early.dates(Stream.guage$datetime)
Stream.guage$datetime <- as.Date(Stream.guage$datetime, format = "%Y%m%d")

```

Define your theme

3. Build a theme and set it as your default theme.

```

#3
mytheme <- theme_classic(base_size = 14) +
  theme(axis.text = element_text(color = "black"),
        legend.position = "top")

```

Create graphs

For numbers 4-7, create graphs that follow best practices for data visualization. To make your graphs “pretty,” ensure your theme, color palettes, axes, and legends are edited to your liking.

Hint: a good way to build graphs is to make them ugly first and then create more code to make them pretty.

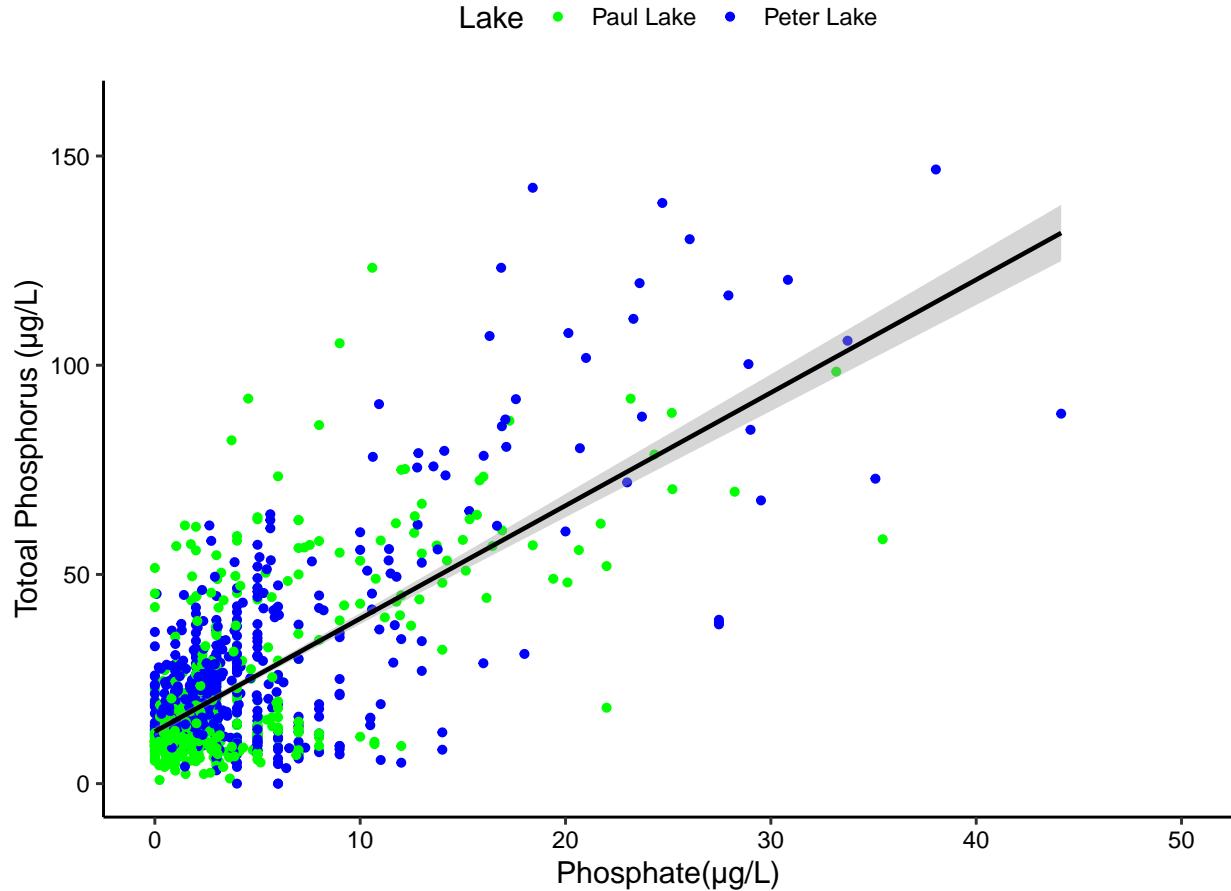
4. [NTL-LTER] Plot total phosphorus by phosphate, with separate aesthetics for Peter and Paul lakes.
Add a line of best fit and color it black.

```

#4
TotalP <- ggplot(PeterPaul.chem.phys,aes(x = po4, y = tp_ug, color = lakename)) +
  xlim(0,50)+ ylim(0,160)+ 
  geom_point()+
  geom_smooth(method = lm, color = "black")+
  labs(x ="Phosphate(\u003BCg/L)", y="Total Phosphorus (\u003BCg/L)", color ="Lake")+
  scale_color_manual(values=c("green", "blue"))+
  mytheme
print(TotalP)

```

```
## Warning: Removed 22312 rows containing non-finite values (stat_smooth).
## Warning: Removed 22312 rows containing missing values (geom_point).
```



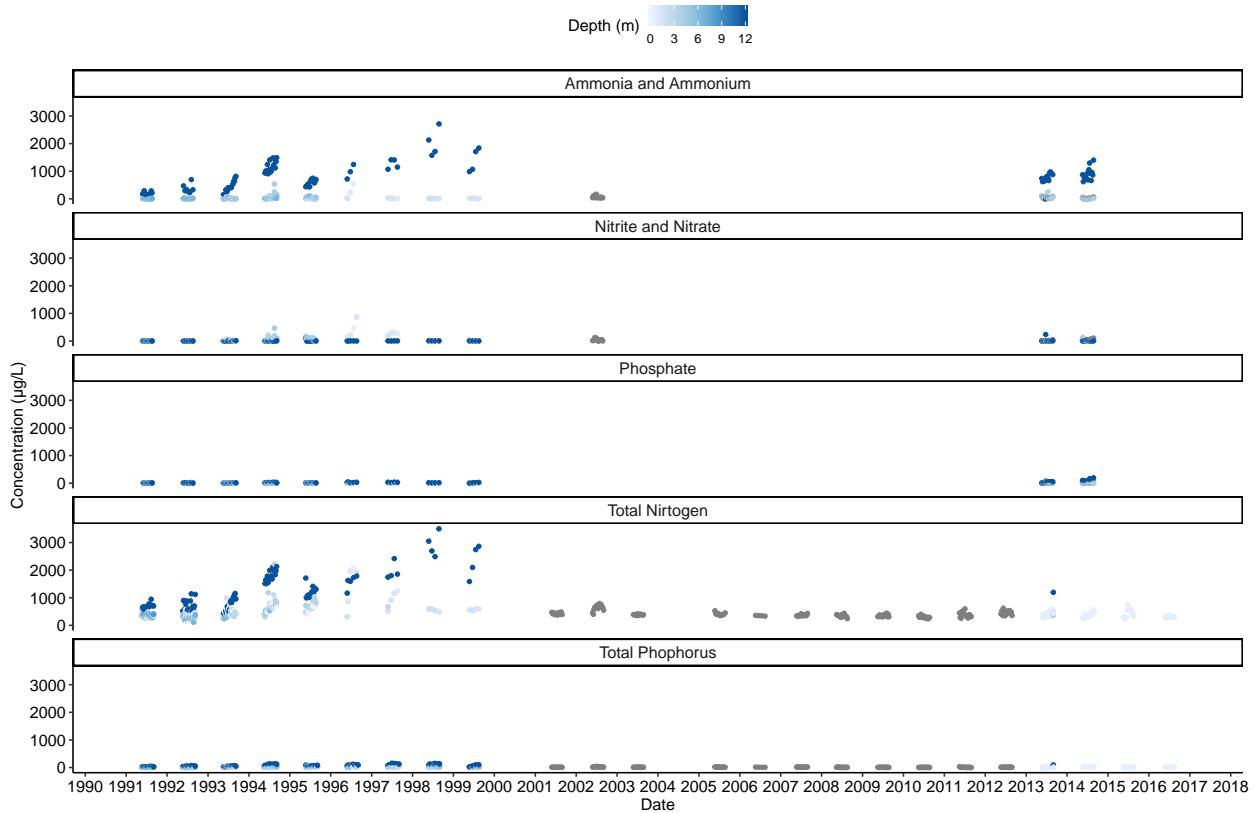
- [NTL-LTER] Plot nutrients by date for Peter Lake, with separate colors for each depth. Facet your graph by the nutrient type.

```
#5
```

```
peterlake <- PeterPaul.nutrients.gathered[PeterPaul.nutrients.gathered$lakename == "Peter Lake",]
levels(peterlake$nutrient)

## [1] "nh34"   "no23"   "po4"    "tn_ug"   "tp_ug"
levels(peterlake$nutrient) <- c("Ammonia and Ammonium", "Nitrite and Nitrate", "Phosphate", "Total Nitrogen", "Total Phosphorus")
Nutrients <- ggplot(peterlake, aes(x = sampledate, y = concentration, color = depth)) +
  geom_point() +
  facet_wrap(vars(nutrient), nrow=5) +
  labs(x = "Date", y = "Concentration (\u003bcg/L)", color = "Depth (m)") +
  scale_x_date(limits = as.Date(c("1991-01-01", "2016-12-31")),
               date_breaks = "1 year", date_labels = "%Y") +
  scale_color_distiller(palette = "Blues", direction = 1) +
  mytheme +
  theme(axis.text=element_text(size=14), strip.text = element_text(size = 14))
```

```
print(Nutrients)
```



#What's the names for the facet labels

6. [USGS gauge] Plot discharge by date. Create two plots, one with the points connected with geom_line and one with the points connected with geom_smooth (hint: do not use method = "lm"). Place these graphs on the same plot (hint: ggarrange or something similar) ?

```
#6
discharge1 <-
ggplot(Stream.guage,aes(x = datetime, y = X165986_00060_00001)) +
  geom_point()+
  geom_smooth()+
  labs(x ="Date", y= expression(paste("Mean Discharge (ft"^-3,"/s)")))+
```

mytheme

```
discharge2 <-
ggplot(Stream.guage,aes(x = datetime, y = X165986_00060_00001)) +
  geom_point()+
  geom_line()+
  labs(x ="Date", y= expression(paste("Mean Discharge (ft"^-3,"/s)")))+
```

mytheme

```
library(ggpubr)
```

```
## Loading required package: magrittr
##
## Attaching package: 'magrittr'
## The following object is masked from 'package:purrr':
```

```

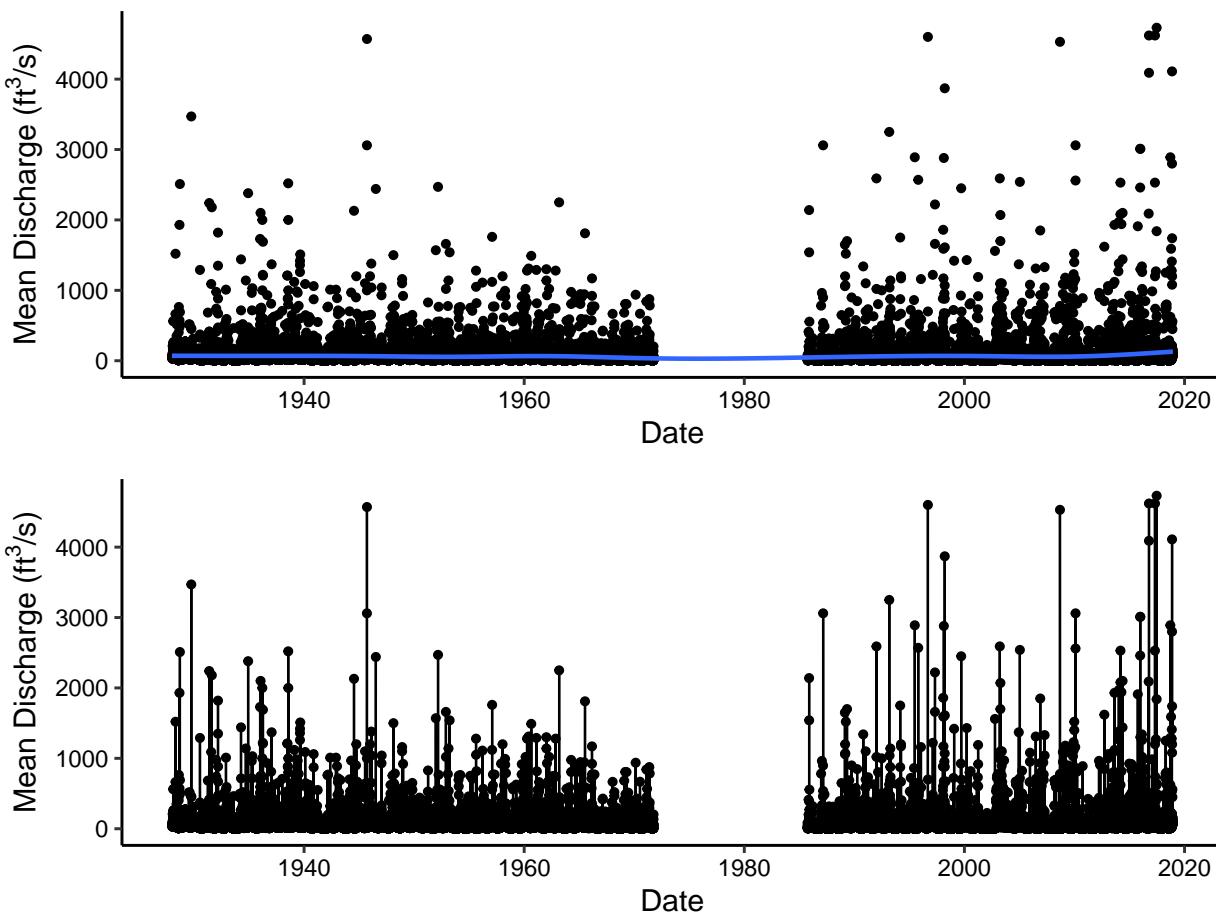
## 
##     set_names
## The following object is masked from 'package:tidyverse':
## 
##     extract
ggarrange(dischARGE1,dischARGE2, nrow=2)

## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
## Warning: Removed 5113 rows containing non-finite values (stat_smooth).
## Warning: Removed 5113 rows containing missing values (geom_point).

## Warning: Removed 5113 rows containing missing values (geom_point).

## Warning: Removed 5113 rows containing missing values (geom_point).

```



```

dischARGE3 <-
ggplot(Stream.guage,aes(x = datetime, y = X165986_00060_00001)) +
  geom_point()+
  geom_smooth()+
  labs(x ="Date", y= expression(paste("Mean Discharge (ft"^-3,"/s"))))+ 
  scale_x_date(limits = as.Date(c("1990-01-01", "2018-12-31")),
  date_breaks = "5 years", date_labels = "%Y") +
  ylim(0,500)+
  mytheme

```

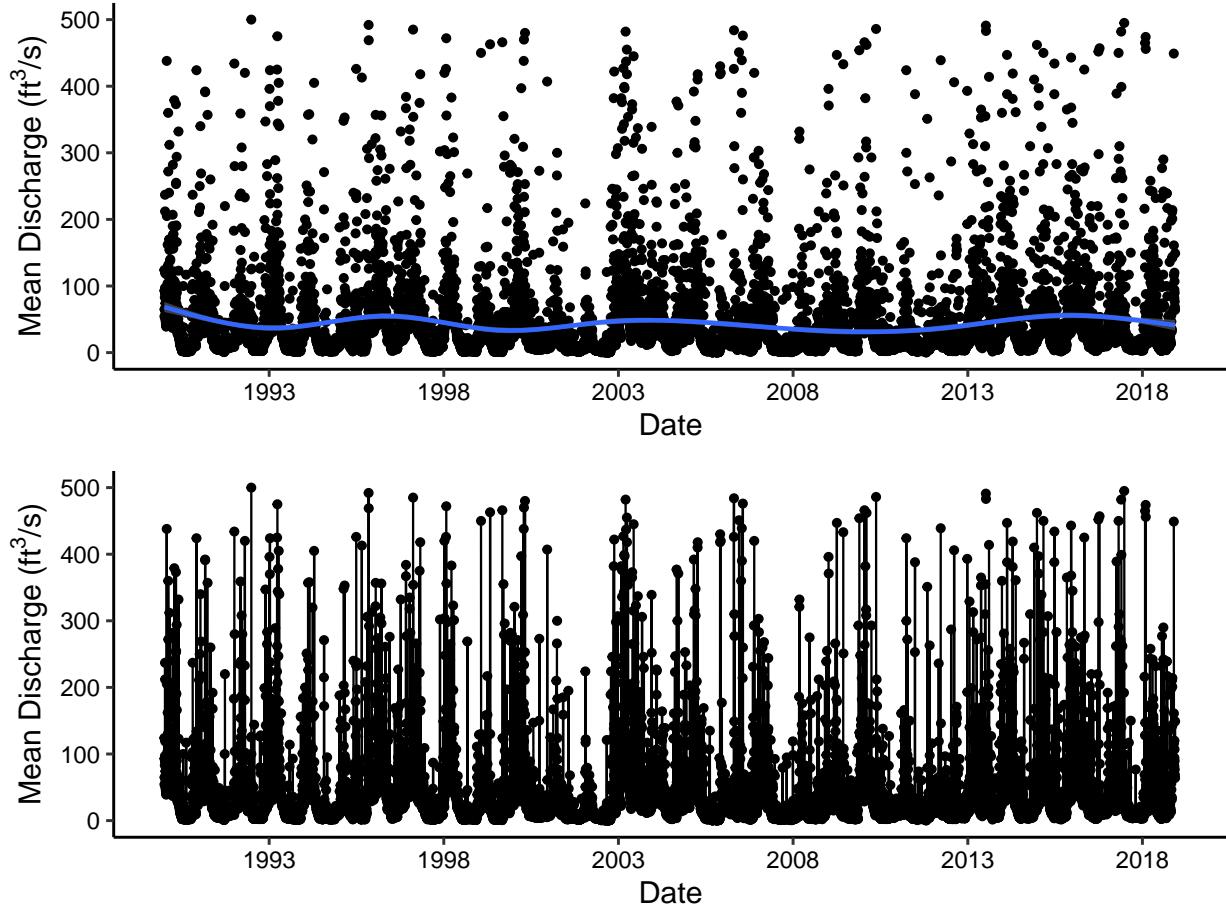
```

discharge4 <-
ggplot(Stream.guage,aes(x = datetime, y = X165986_00060_00001)) +
  geom_point()+
  geom_line()+
  labs(x ="Date", y= expression(paste("Mean Discharge (ft"^-3,"/s)")))+
  scale_x_date(limits = as.Date(c("1990-01-01", "2018-12-31")),
date_breaks = "5 years", date_labels = "%Y") +
  ylim(0,500)+ 
  mytheme
ggarrange(discharge3,discharge4, nrow=2)

## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
## Warning: Removed 22918 rows containing non-finite values (stat_smooth).
## Warning: Removed 22918 rows containing missing values (geom_point).

## Warning: Removed 22918 rows containing missing values (geom_point).
## Warning: Removed 22646 rows containing missing values (geom_path).

```



Question: How do these two types of lines affect your interpretation of the data?

Answer: `geom_smooth` can create a line that shows the trend of the discharge data, but `geom_line` can only connect all the points. The former can help see the data trend, which is nearly a linear line that didn't change a lot over time, but the latter can just let us know there is fluctuation of discharge data over time.

7. [ECOTOX Neonicotinoids] Plot the concentration, divided by chemical name. Choose a geom that accurately portrays the distribution of data points.

```
#7
subEcotox <- Ecotox[Ecotox$Conc..Units..Std. == "AI mg/L",]
con <- ggplot(subEcotox) +
  geom_freqpoly(aes(x = Conc..Mean..Std., color= Chemical.Name), bins=20) +
  labs(x ="Concentration (mg/L)", y="Count", color = "Chemical") +
  xlim(0,20) +
  scale_color_manual(values = c(1:8)) +
  mytheme
print(con)
```

Warning: Removed 265 rows containing non-finite values (stat_bin).

Warning: Removed 14 rows containing missing values (geom_path).

