

5: Part 1 - Data Visualization Basics

Environmental Data Analytics | John Fay and Luana Lima | Developed by Kateri Salk

Spring 2022

Objectives

1. Perform simple data visualizations in the R package `ggplot`
2. Develop skills to adjust aesthetics and layers in graphs
3. Apply a decision tree framework for appropriate graphing methods

Opening discussion

Effective data visualization depends on purposeful choices about graph types. The ideal graph type depends on the type of data and the message the visualizer desires to communicate. The best visualizations are clear and simple. A good resource for data visualization is Data to Viz, which includes both a decision tree for visualization types and explanation pages for each type of data, including links to R resources to create them. Take a few minutes to explore this website.

Set Up

```
getwd()

## [1] "/Users/ataliefischer/Desktop/EDA/Environmental_Data_Analytics_2022/Lessons"

library(tidyverse)
#install.packages("ggridges")
library(ggridges)

PeterPaul.chem.nutrients <-
  read.csv("../Data/Processed/NTL-LTER_Lake_Chemistry_Nutrients_PeterPaul_Processed.csv", stringsAsFactors = FALSE)
PeterPaul.chem.nutrients.gathered <-
  read.csv("../Data/Processed/NTL-LTER_Lake_Nutrients_PeterPaulGathered_Processed.csv", stringsAsFactors = FALSE)
EPAair <- read.csv("../Data/Processed/EPAair_03_PM25_NC2021_Processed.csv", stringsAsFactors = TRUE)

EPAair$Date <- as.Date(EPAair$Date, format = "%Y-%m-%d")
PeterPaul.chem.nutrients$sampldate <- as.Date(PeterPaul.chem.nutrients$sampldate, format = "%Y-%m-%d")
PeterPaul.chem.nutrients.gathered$sampldate <- as.Date(PeterPaul.chem.nutrients.gathered$sampldate, format = "%Y-%m-%d")
```

ggplot

`ggplot`, called from the package `ggplot2`, is a graphing and image generation tool in R. This package is part of `tidyverse`. While base R has graphing capabilities, `ggplot` has the capacity for a wider range and more sophisticated options for graphing. `ggplot` has only a few rules:

- The first line of `ggplot` code always starts with `ggplot()`
- A data frame must be specified within the `ggplot()` function. Additional datasets can be specified in subsequent layers.

- Aesthetics must be specified, most commonly x and y variables but including others. Aesthetics can be specified in the `ggplot()` function or in subsequent layers.
- Additional layers must be specified to fill the plot.

Geoms

Here are some commonly used layers for plotting in ggplot:

- `geom_bar`
- `geom_histogram`
- `geom_freqpoly`
- `geom_boxplot`
- `geom_violin`
- `geom_dotplot`
- `geom_density_ridges`
- `geom_point`
- `geom_errorbar`
- `geom_smooth`
- `geom_line`
- `geom_area`
- `geom_abline` (plus `geom_hline` and `geom_vline`)
- `geom_text`

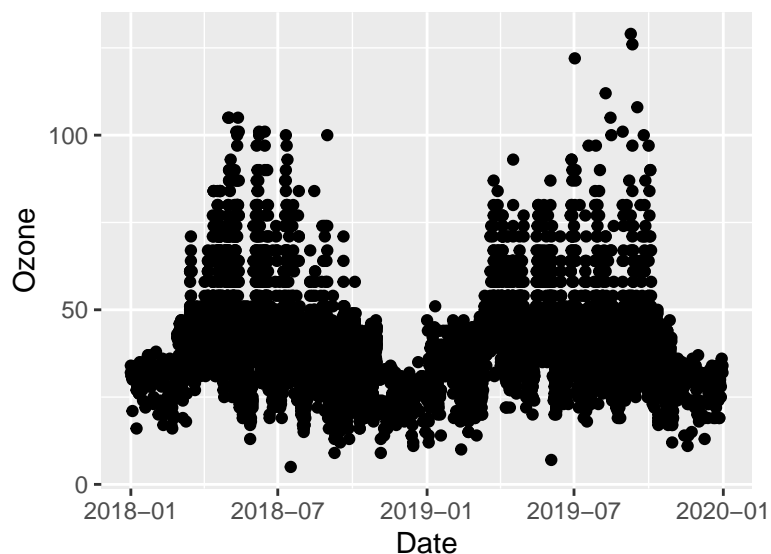
Aesthetics

Here are some commonly used aesthetic types that can be manipulated in ggplot:

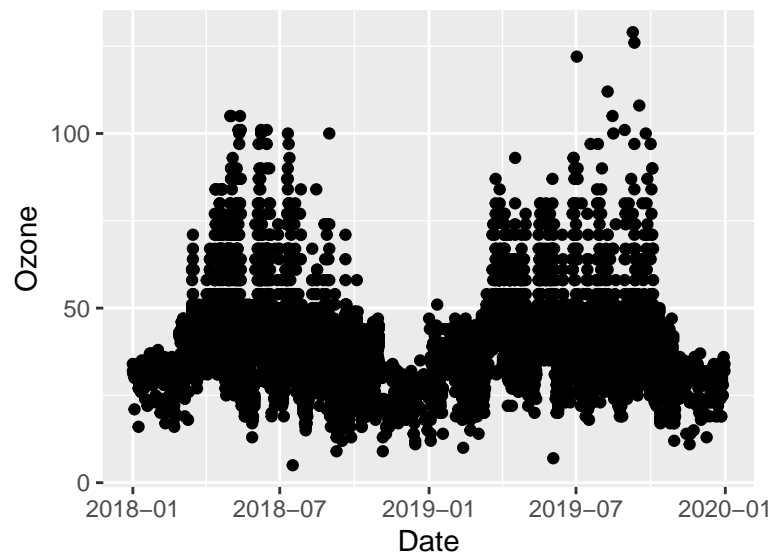
- `color`
- `fill`
- `shape`
- `size`
- `transparency`

Plotting continuous variables over time: Scatterplot and Line Plot

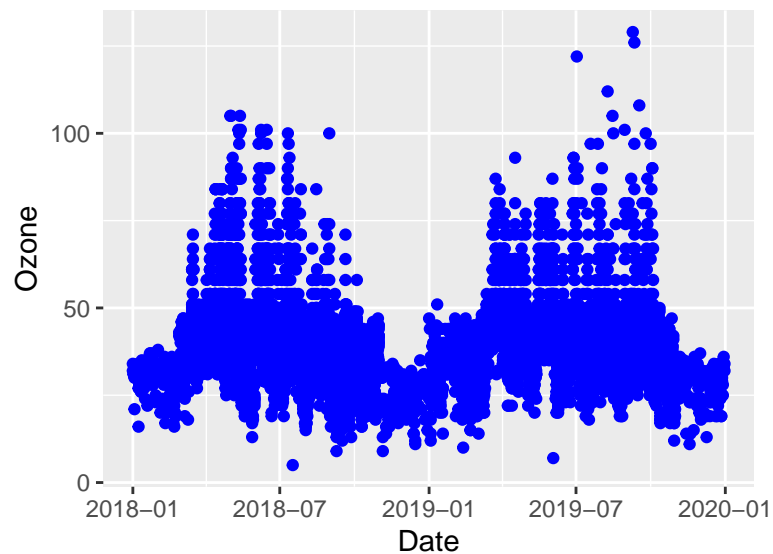
```
# Scatterplot
ggplot(EPAair, aes(x = Date, y = Ozone)) +
  geom_point()
```



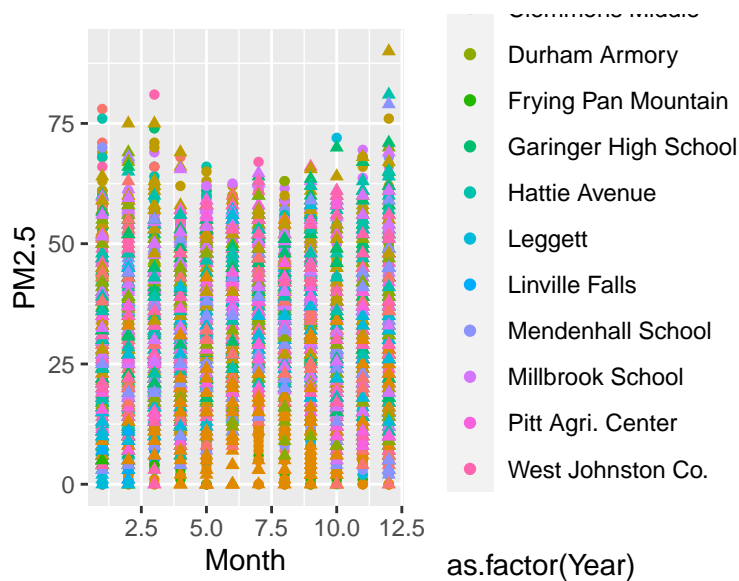
```
O3plot <- ggplot(EPAair) +
  geom_point(aes(x = Date, y = Ozone))
print(O3plot)
```



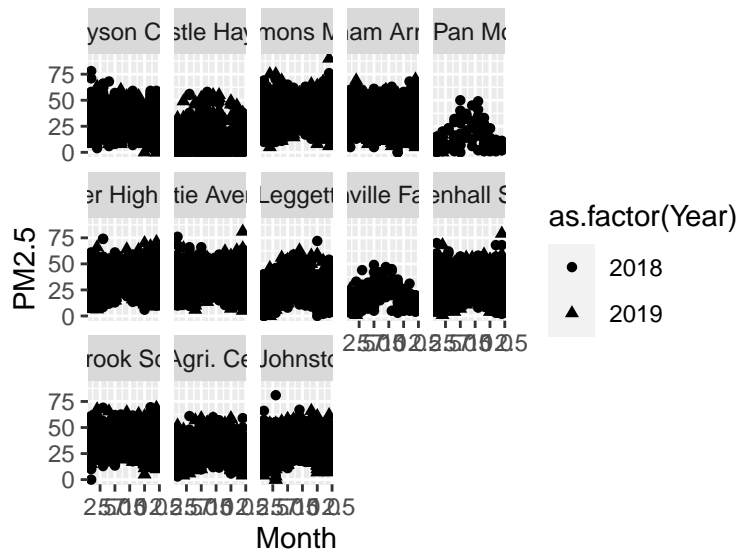
```
# Fix this code
O3plot2 <- ggplot(EPAair) +
  #geom_point(aes(x = Date, y = Ozone, color = "blue"))
  geom_point(aes(x = Date, y = Ozone), color = "blue")
print(O3plot2)
```



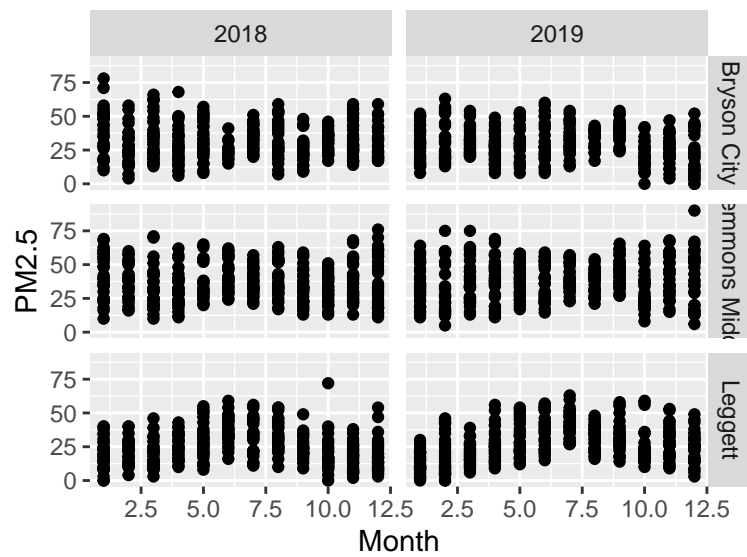
```
# Add additional variables
# How could you automatically assign a marker color to a variable?
PMplot <-
  ggplot(EPAair, aes(x = Month, y = PM2.5, shape = as.factor(Year), color = Site.Name)) +
  geom_point()
print(PMplot)
```



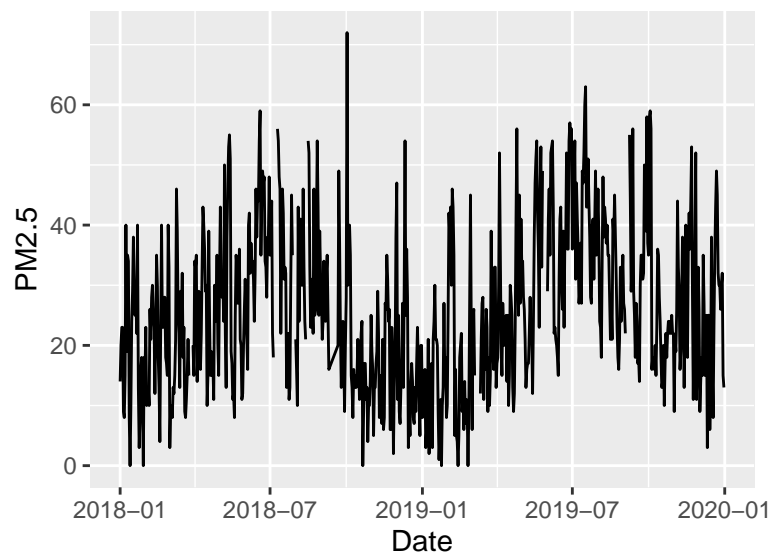
```
# Separate plot with facets: facet_wrap --> creates one plot for each specified variable
PMplot.faceted <-
  ggplot(EPAair, aes(x = Month, y = PM2.5, shape = as.factor(Year))) +
  geom_point() +
  facet_wrap(vars(Site.Name), nrow = 3)
print(PMplot.faceted)
```



```
# Filter dataset within plot building and facet by multiple variables: subset --> only use few sites with
PMplot.faceted2 <-
  ggplot(subset(EPAair, Site.Name == "Clemmons Middle" | Site.Name == "Leggett" |
    Site.Name == "Bryson City"),
    aes(x = Month, y = PM2.5)) +
  geom_point() +
  facet_grid(Site.Name ~ Year)
print(PMplot.faceted2)
```

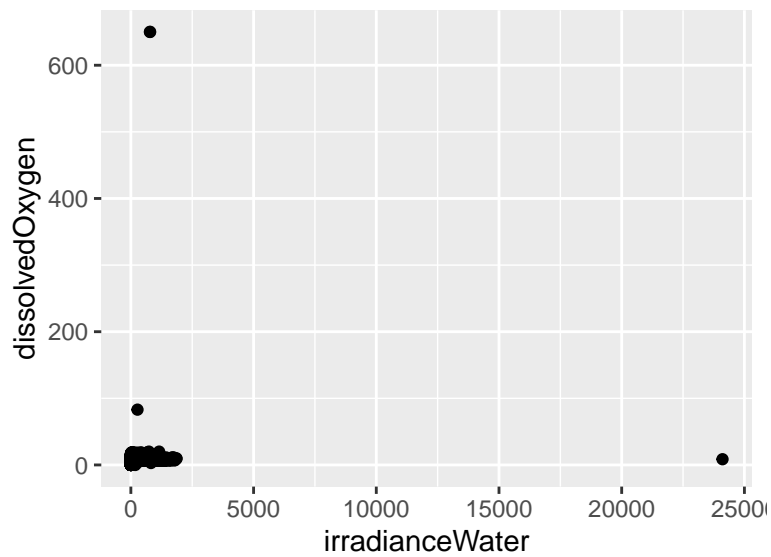


```
# Plot true time series with geom_line
PMplot.line <-
  ggplot(subset(EPAair, Site.Name == "Leggett"),
    aes(x = Date, y = PM2.5)) +
    geom_line()
print(PMplot.line)
```

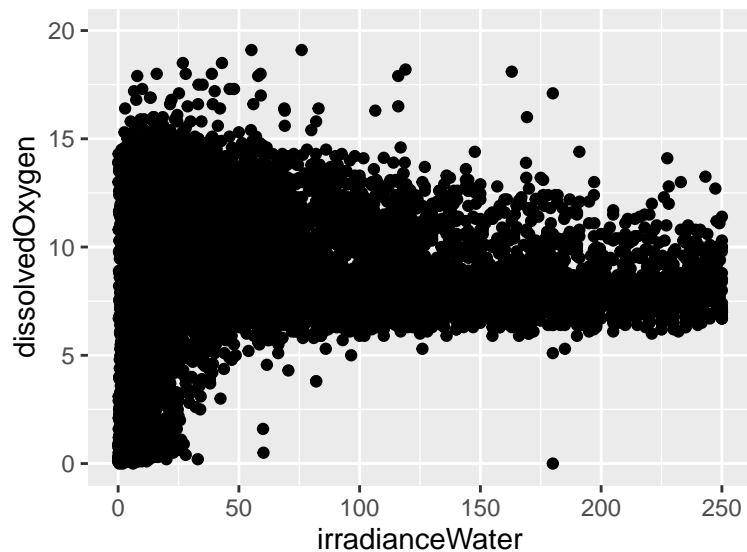


Plotting the relationship between two continuous variables: Scatterplot

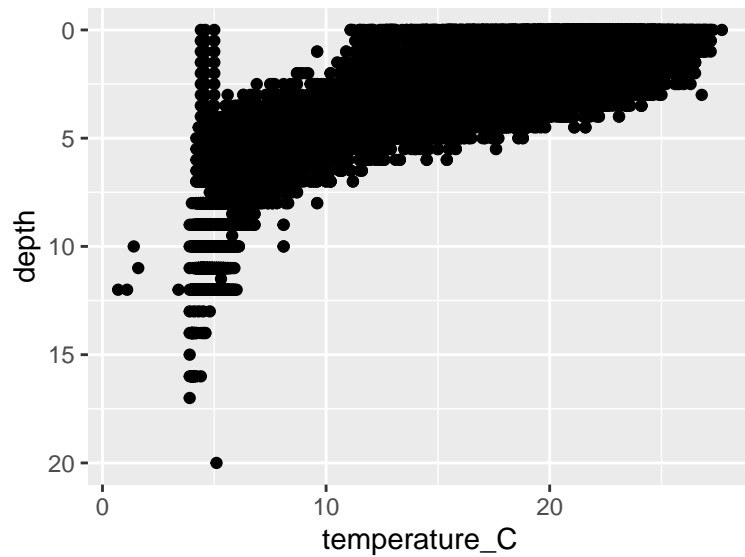
```
# Scatterplot
lightvsDO <-
  ggplot(PeterPaul.chem.nutrients, aes(x = irradianceWater, y = dissolvedOxygen)) +
  geom_point()
print(lightvsDO)
```



```
# Adjust axes
lightvsDOfixed <-
  ggplot(PeterPaul.chem.nutrients, aes(x = irradianceWater, y = dissolvedOxygen)) +
  geom_point() +
  xlim(0, 250) +
  ylim(0, 20)
print(lightvsDOfixed)
```

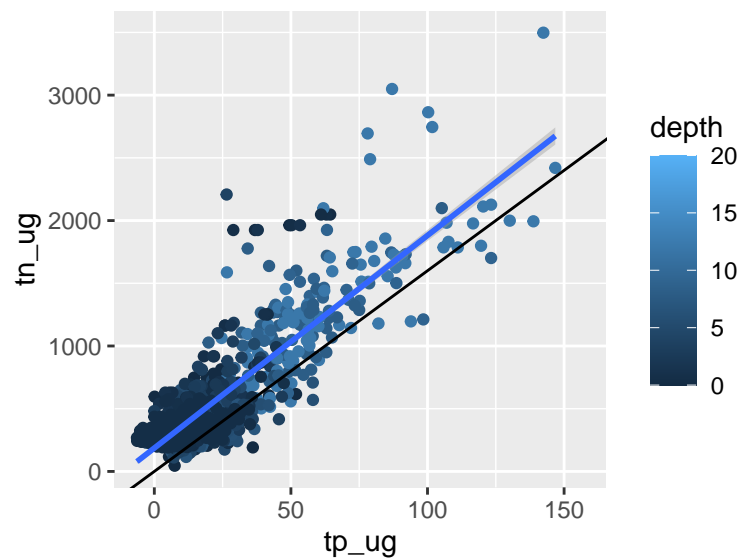


```
# Depth in the fields of limnology and oceanography is on a reverse scale
tempvsdepth <-
  ggplot(PeterPaul.chem.nutrients, aes(x = temperature_C, y = depth)) +
  #ggplot(PeterPaul.chem.nutrients, aes(x = temperature_C, y = depth, color = daynum)) +
  geom_point() +
  scale_y_reverse()
print(tempvsdepth)
```



```
NvsP <-
  ggplot(PeterPaul.chem.nutrients, aes(x = tp_ug, y = tn_ug, color = depth)) +
  geom_point() +
  geom_smooth(method = lm) +
  geom_abline(aes(slope = 16, intercept = 0))
print(NvsP)
```

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



Plotting continuous vs. categorical variables

A traditional way to display summary statistics of continuous variables is a bar plot with error bars. Let's explore why this might not be the most effective way to display this type of data. Navigate to the Caveats page on Data to Viz (<https://www.data-to-viz.com/caveats.html>) and find the page that explores barplots and error bars.

What might be more effective ways to display the information? Navigate to the boxplots page in the Caveats section to explore further.

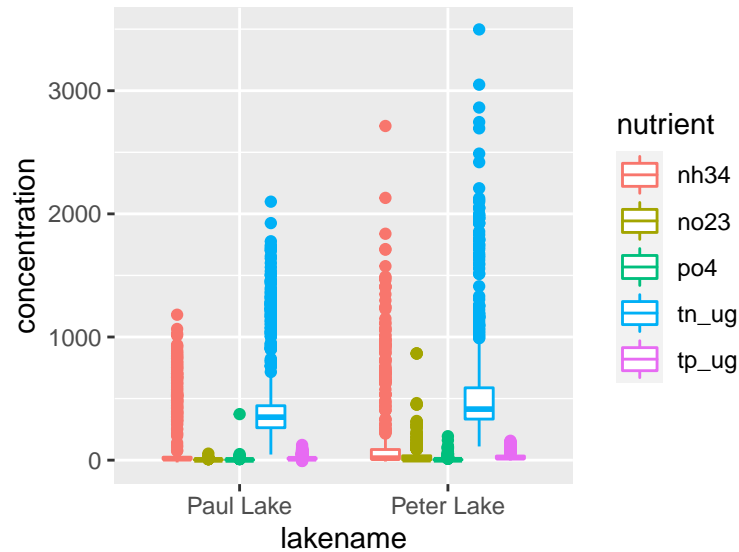
```
# Box and whiskers plot
```

```
Nutrientplot3 <-
```

```
  ggplot(PeterPaul.chem.nutrients.gathered, aes(x = lakename, y = concentration)) +
```

```
  geom_boxplot(aes(color = nutrient)) # Why didn't we use "fill"? color --> contour. When we use fill f
```

```
print(Nutrientplot3)
```



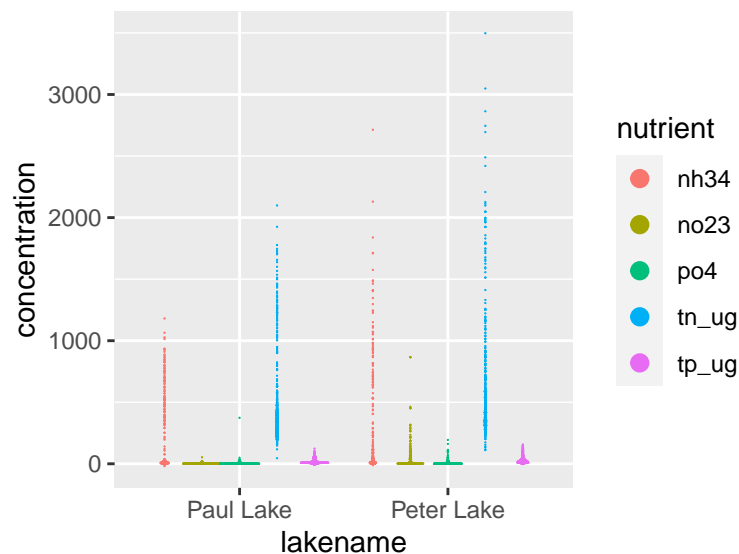
```
# Dot plot: binaxis --> which direction to stack. dodge --> pull same value points to the side without
```

```
Nutrientplot4 <-
```

```
  ggplot(PeterPaul.chem.nutrients.gathered, aes(x = lakename, y = concentration)) +
```

```
  geom_dotplot(aes(color = nutrient, fill = nutrient), binaxis = "y", binwidth = 1,  
              stackdir = "center", position = "dodge", dotsize = 2) #
```

```
print(Nutrientplot4)
```



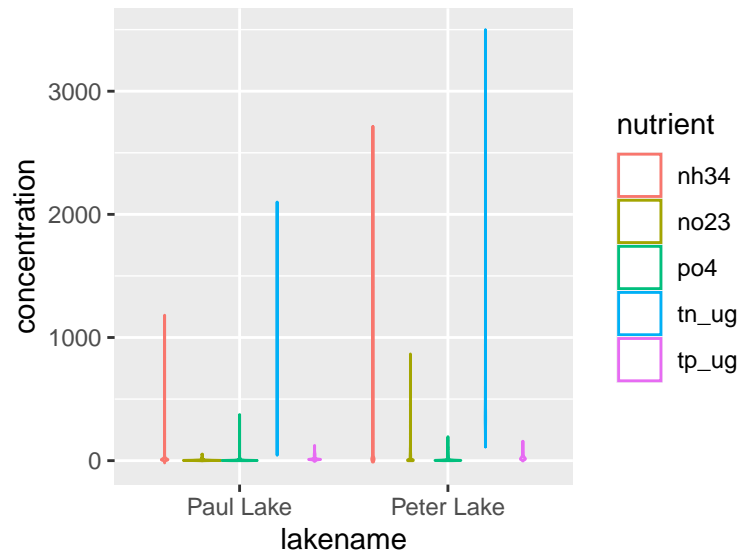
```
# Violin plot: like a box plot but with density as well
```

```
Nutrientplot5 <-
```

```
  ggplot(PeterPaul.chem.nutrients.gathered, aes(x = lakename, y = concentration)) +
```

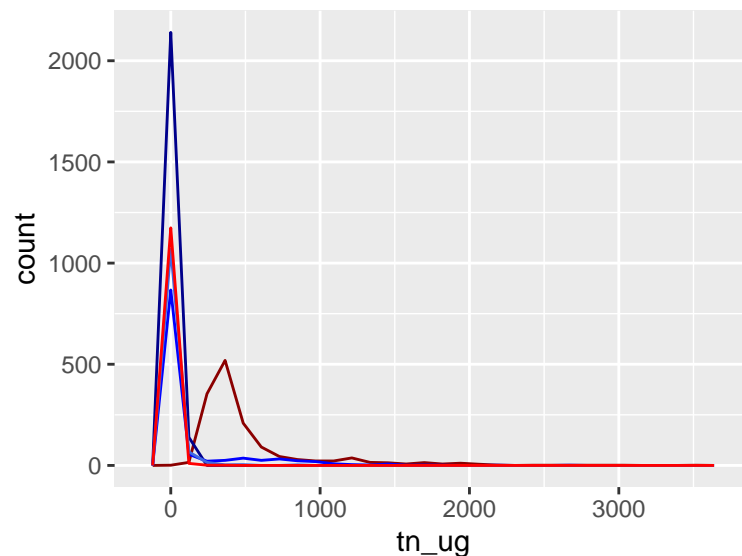
```
  geom_violin(aes(color = nutrient)) #
```

```
print(Nutrientplot5)
```

```
# Frequency polygons
# Using a tidy dataset --> no legend
Nutrientplot6 <-
  ggplot(PeterPaul.chem.nutrients) +
    geom_freqpoly(aes(x = tn_ug), color = "darkred") +
    geom_freqpoly(aes(x = tp_ug), color = "darkblue") +
    geom_freqpoly(aes(x = nh34), color = "blue") +
    geom_freqpoly(aes(x = no23), color = "royalblue") +
    geom_freqpoly(aes(x = po4), color = "red")
print(Nutrientplot6)
```

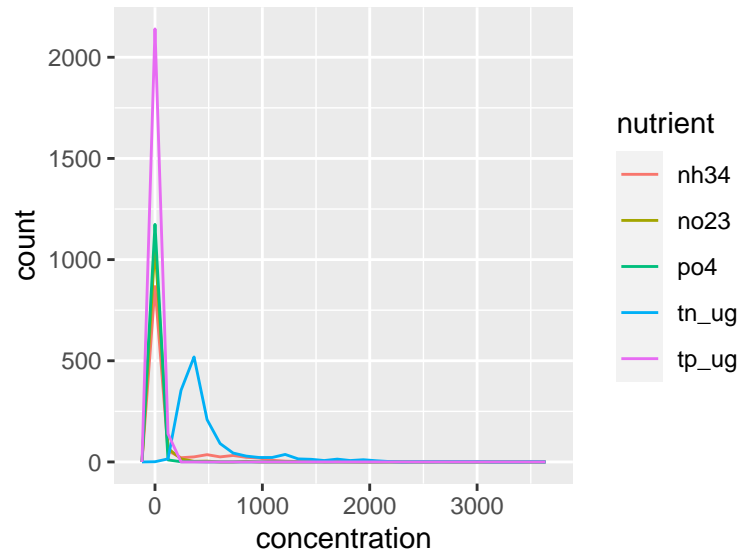
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
# Using a gathered dataset --> legend included
Nutrientplot7 <-
```

```
ggplot(PeterPaul.chem.nutrients.gathered) +
  geom_freqpoly(aes(x = concentration, color = nutrient))
print(Nutrientplot7)
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



*# Frequency polygons have the risk of becoming spaghetti plots.
See <<https://www.data-to-viz.com/caveat/spaghetti.html>> for more info.*

```
# Ridgeline plot --> density distribution over time
Nutrientplot6 <-
  ggplot(PeterPaul.chem.nutrients.gathered, aes(y = nutrient, x = concentration)) +
  geom_density_ridges(aes(fill = lakename), alpha = 0.5) #
print(Nutrientplot6)
```

Picking joint bandwidth of 10.9

