

Assignment 3: Data Exploration

Laurie F Muzzy

OVERVIEW

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

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_A02_DataExploration.pdf”) prior to submission.

The completed exercise is due on Thursday, 31 January, 2019 before class begins.

1) Set up your R session

Check your working directory, load necessary packages (tidyverse), and upload the North Temperate Lakes long term monitoring dataset for the light, temperature, and oxygen data for three lakes (file name: NTL-LTER_Lake_ChemistryPhysics_Raw.csv). Type your code into the R chunk below.

```
getwd()

## [1] "/Users/laurie/Desktop/Envtl_Data_Analytics/MuzzyGitFile/Assignments"

library(tidyverse)

## Loading tidyverse: ggplot2
## Loading tidyverse: tibble
## Loading tidyverse: tidyr
## Loading tidyverse: readr
## Loading tidyverse: purrr
## Loading tidyverse: dplyr

## Warning: package 'ggplot2' was built under R version 3.4.4
## Warning: package 'tibble' was built under R version 3.4.3
## Warning: package 'dplyr' was built under R version 3.4.4
## Conflicts with tidy packages -----
## filter(): dplyr, stats
## lag():    dplyr, stats

lter.chemphys <- read.csv("../Data/Raw/NTL-LTER_Lake_ChemistryPhysics_Raw.csv") #the data is in the Raw
```

2) Learn about your system

Read about your dataset in the NTL-LTER README file. What are three salient pieces of information you gained from reading this file?

ANSWER: Data is from 1984-2016; the Physical and Chemical Limnology is measured at the deepest location in each of the several lakes; the chemical data is sometimes pool mixed layer sample, sometimes at 3 locations, and sometimes in vertical profiles.

3) Obtain basic summaries of your data

Write R commands to display the following information:

1. dimensions of the dataset
2. class of the dataset
3. first 8 rows of the dataset
4. class of the variables lakename, sampleddate, depth, and temperature
5. summary of lakename, depth, and temperature

```
# 1 dimensions of the dataset
dim(lter.chemphys)
```

```
## [1] 38614    11
```

```
# 2
class(lter.chemphys)
```

```
## [1] "data.frame"
```

```
# 3
head(lter.chemphys, 8)
```

```
##   lakeid  lakename year4 daynum sampleddate depth temperature_C
## 1      L Paul Lake 1984   148    5/27/84  0.00           14.5
## 2      L Paul Lake 1984   148    5/27/84  0.25             NA
## 3      L Paul Lake 1984   148    5/27/84  0.50             NA
## 4      L Paul Lake 1984   148    5/27/84  0.75             NA
## 5      L Paul Lake 1984   148    5/27/84  1.00           14.5
## 6      L Paul Lake 1984   148    5/27/84  1.50             NA
## 7      L Paul Lake 1984   148    5/27/84  2.00           14.2
## 8      L Paul Lake 1984   148    5/27/84  3.00           11.0
##   dissolvedOxygen irradianceWater irradianceDeck comments
## 1              9.5             1750             1620    <NA>
## 2              NA             1550             1620    <NA>
## 3              NA             1150             1620    <NA>
## 4              NA              975             1620    <NA>
## 5              8.8              870             1620    <NA>
## 6              NA              610             1620    <NA>
## 7              8.6              420             1620    <NA>
## 8             11.5              220             1620    <NA>
```

```
# 4
class(lter.chemphys$lakename)
```

```
## [1] "factor"
```

```
class(lter.chemphys$sampledate)
```

```
## [1] "factor"
class(lter.chemphys$depth)

## [1] "numeric"
class(lter.chemphys$temperature_C)

## [1] "numeric"
# 5
summary(lter.chemphys$lakename)

## Central Long Lake      Crampton Lake      East Long Lake      Hummingbird Lake
##           539           1234           3905           430
##      Paul Lake      Peter Lake      Tuesday Lake      Ward Lake
##      10325      11288      6107      598
## West Long Lake
##      4188
```

```
summary(lter.chemphys$depth)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00   1.50   4.00   4.39   6.50   20.00
```

```
summary(lter.chemphys$temperature_C)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##      0.30   5.30   9.30   11.81   18.70   34.10   3858
```

Change sampleddate to class = date. After doing this, write an R command to display that the class of sampleddate is indeed date. Write another R command to show the first 10 rows of the date column.

#sampledate is in class factor, so we have to convert its class to date and make it year/month/day

#lter.chemphys\$sampledate <- class(lter.chemphys\$sampledate) This changed all the actual dates to the w
#lter.chemphys\$sampledate <- as.Date(lter.chemphys\$sampledate, "%y/%m/%d") #This changed all the sample
#lter.chemphys\$sampledate <- format(as.Date(lter.chemphys\$sampledate, "%y%m%d")) #This ALSO changed all
#lter.chemphys\$sampledate <- factor("%m/%d/%y"), as.Date(lter.chemphys\$sampledate, format = "%Y/%m/%d")

```
lter.chemphys$sampledate <- as.Date(lter.chemphys$sampledate, format = "%m/%d/%y") #telling it to conver
```

```
## Warning in strptime(x, format, tz = "GMT"): unknown timezone 'default/'
## America/New_York'
```

```
class(lter.chemphys$sampledate)
```

```
## [1] "Date"
```

```
head(lter.chemphys$sampledate, 10)
```

```
## [1] "1984-05-27" "1984-05-27" "1984-05-27" "1984-05-27" "1984-05-27"
## [6] "1984-05-27" "1984-05-27" "1984-05-27" "1984-05-27" "1984-05-27"
```

Question: Do you want to remove NAs from this dataset? Why or why not?

ANSWER: Yes, probably, depending on how we were displaying the data. We don't want it to skew our displays in any way. but when I used na.omit, all but 149 entries (out of over 38,000) remained.

#complete.cases(lter.chemphys) #this returned a list of all observations as FALSE ...?
lter.chemphys.complete <- na.omit(lter.chemphys) #but now there's only 149 entries

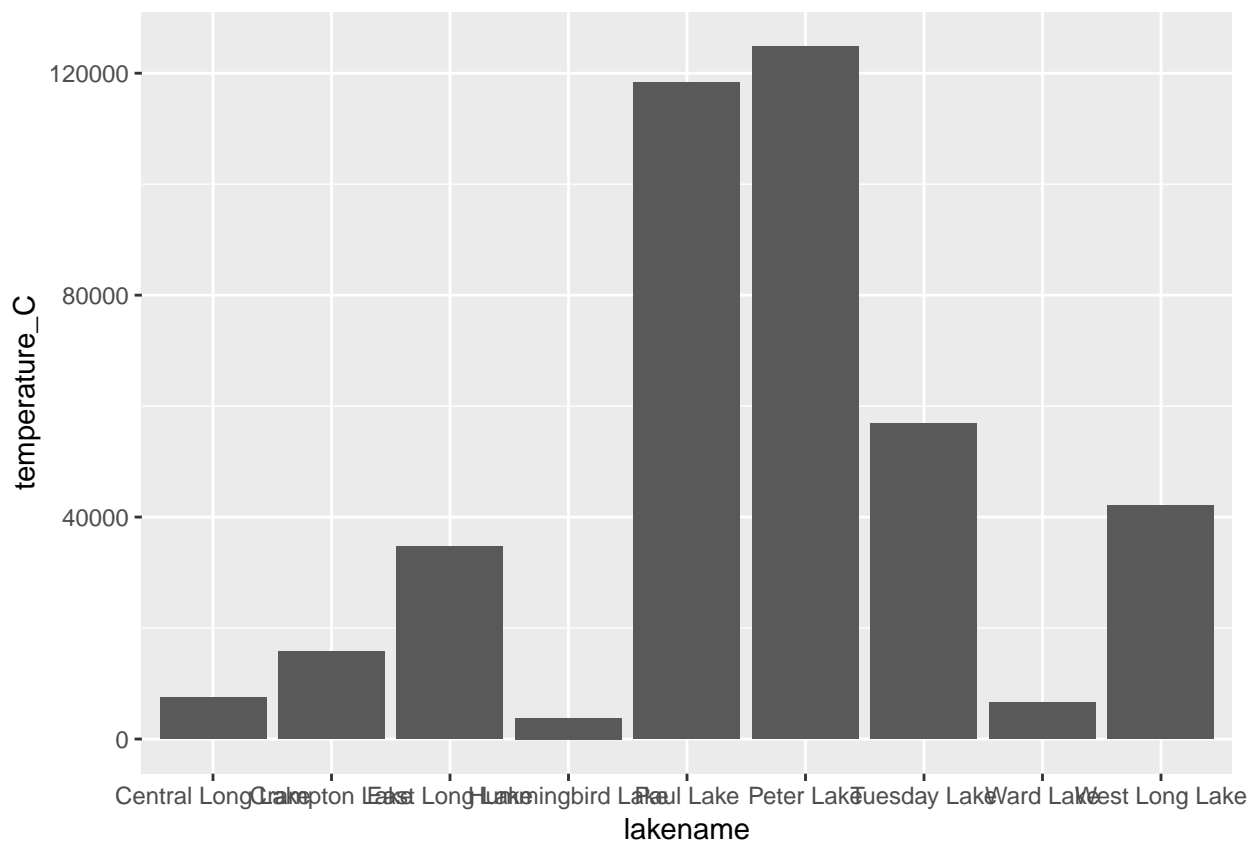
4) Explore your data graphically

Write R commands to display graphs depicting:

1. Bar chart of temperature counts for each lake
2. Histogram of count distributions of temperature (all temp measurements together)
3. Change histogram from 2 to have a different number or width of bins
4. Frequency polygon of temperature for each lake. Choose different colors for each lake.
5. Boxplot of temperature for each lake
6. Boxplot of temperature based on depth, with depth divided into 0.25 m increments
7. Scatterplot of temperature by depth

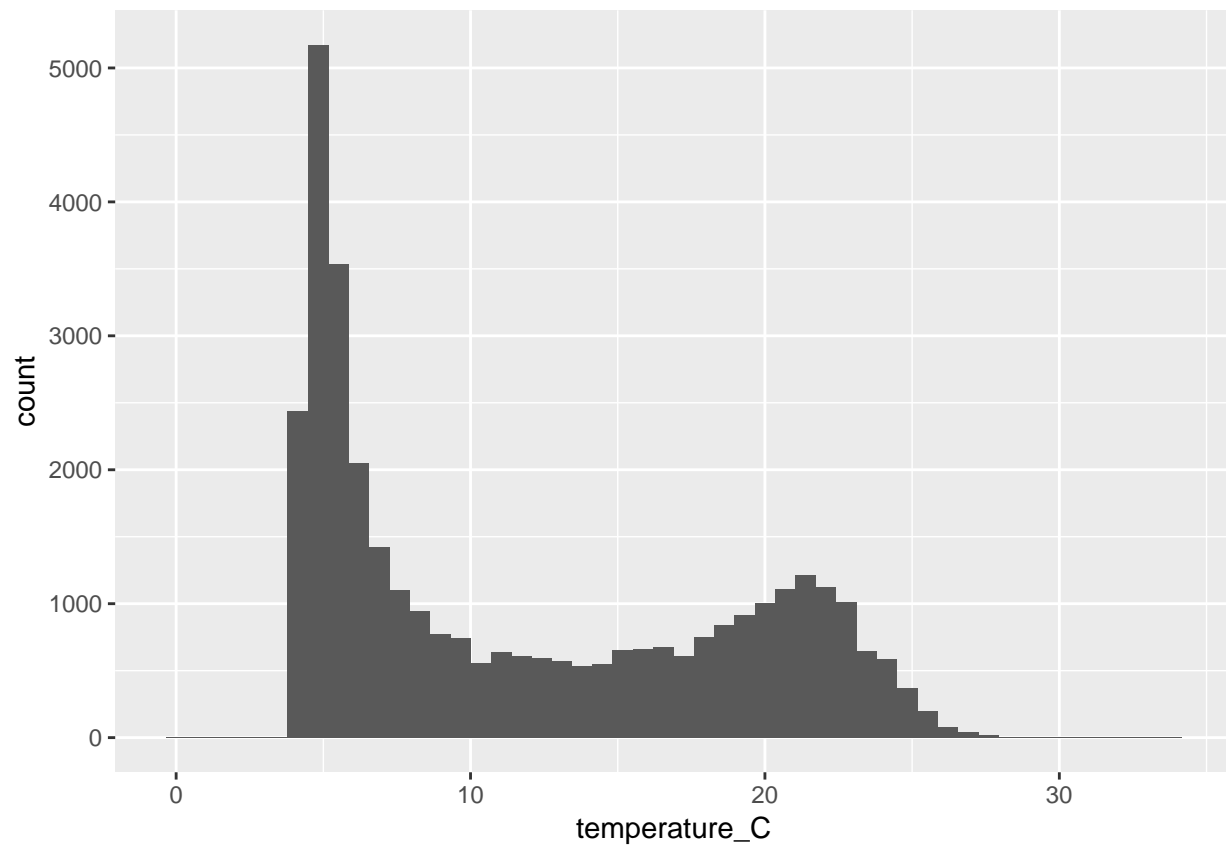
```
# 1
ggplot(lter.chemphys, aes(x = lakename, y = temperature_C)) +
  geom_col() #can't figure out how to specify according to each lakename
```

```
## Warning: Removed 3858 rows containing missing values (position_stack).
```



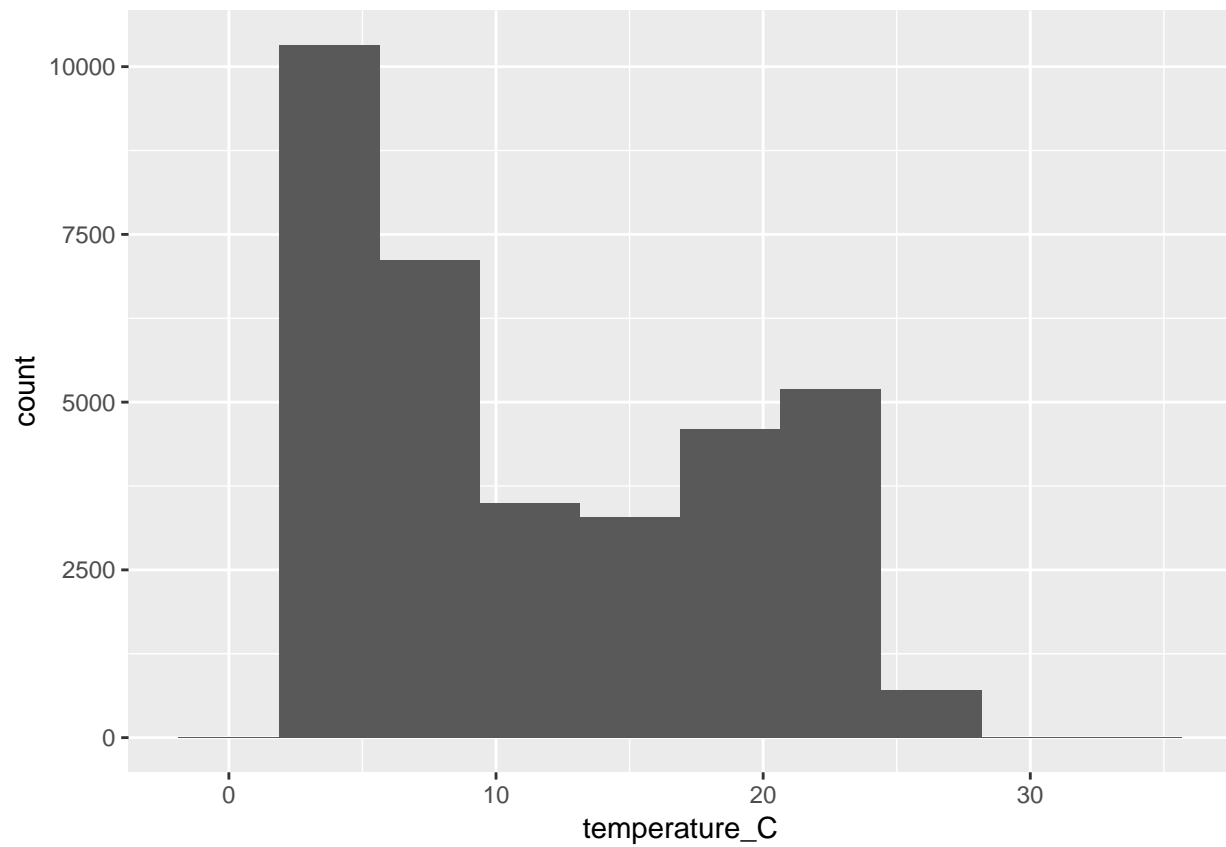
```
# 2
ggplot(lter.chemphys) +
  geom_histogram(aes(x = temperature_C), bins = 50)
```

```
## Warning: Removed 3858 rows containing non-finite values (stat_bin).
```



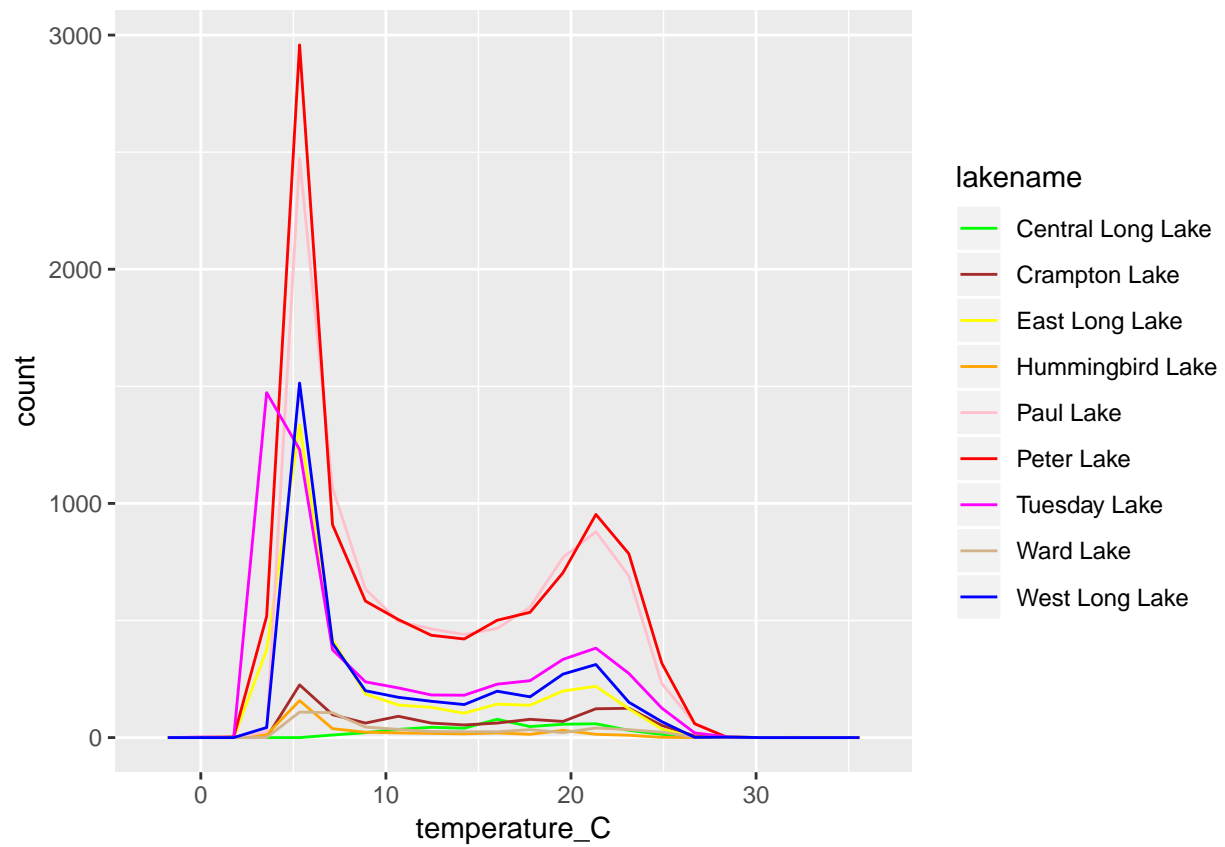
```
# 3  
ggplot(lter.chemphys) +  
  geom_histogram(aes(x = temperature_C), bins = 10)
```

```
## Warning: Removed 3858 rows containing non-finite values (stat_bin).
```



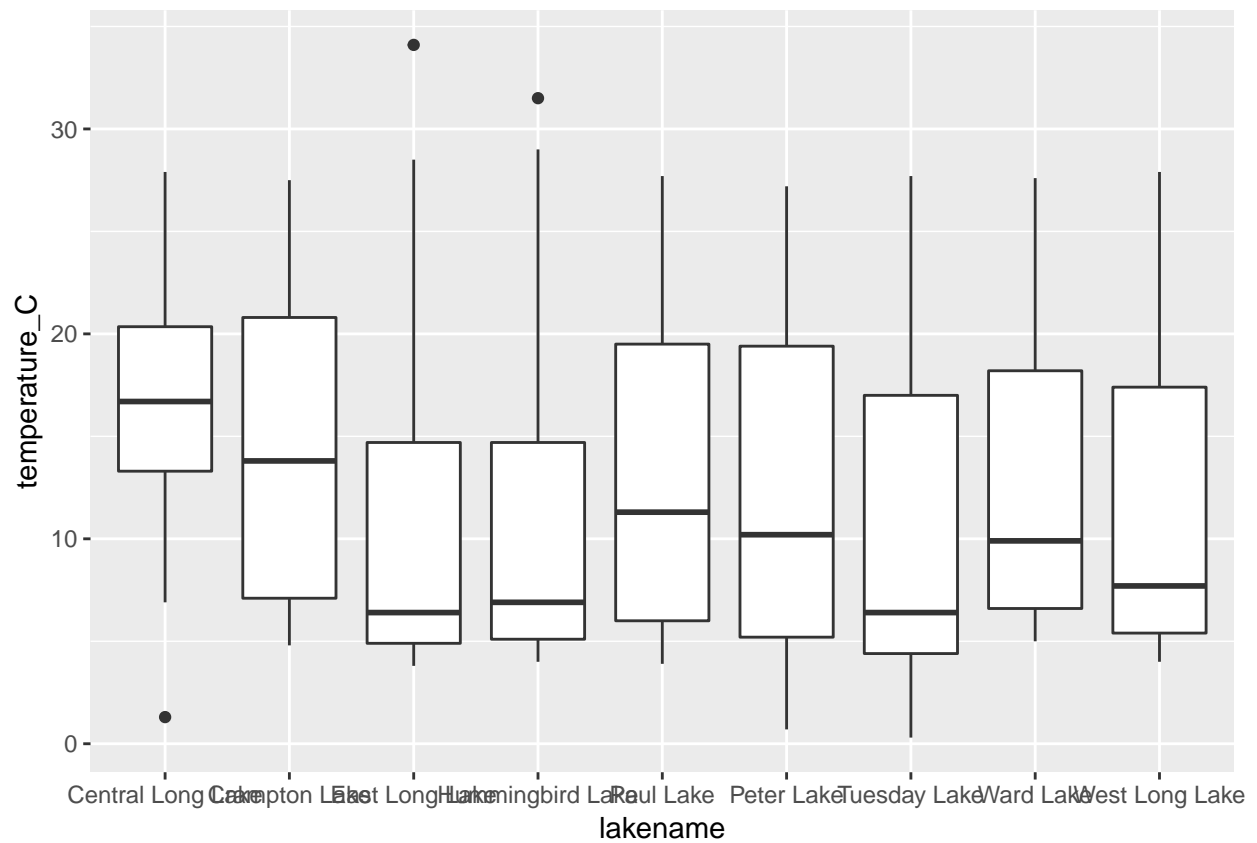
```
# 4
ggplot(lter.chemphys) +
  geom_freqpoly(aes(x = temperature_C, color = lakename), bins = 20) +
  scale_color_manual(values = c("green", "brown", "yellow", "orange", "pink", "red", "magenta", "tan", "b

## Warning: Removed 3858 rows containing non-finite values (stat_bin).
```



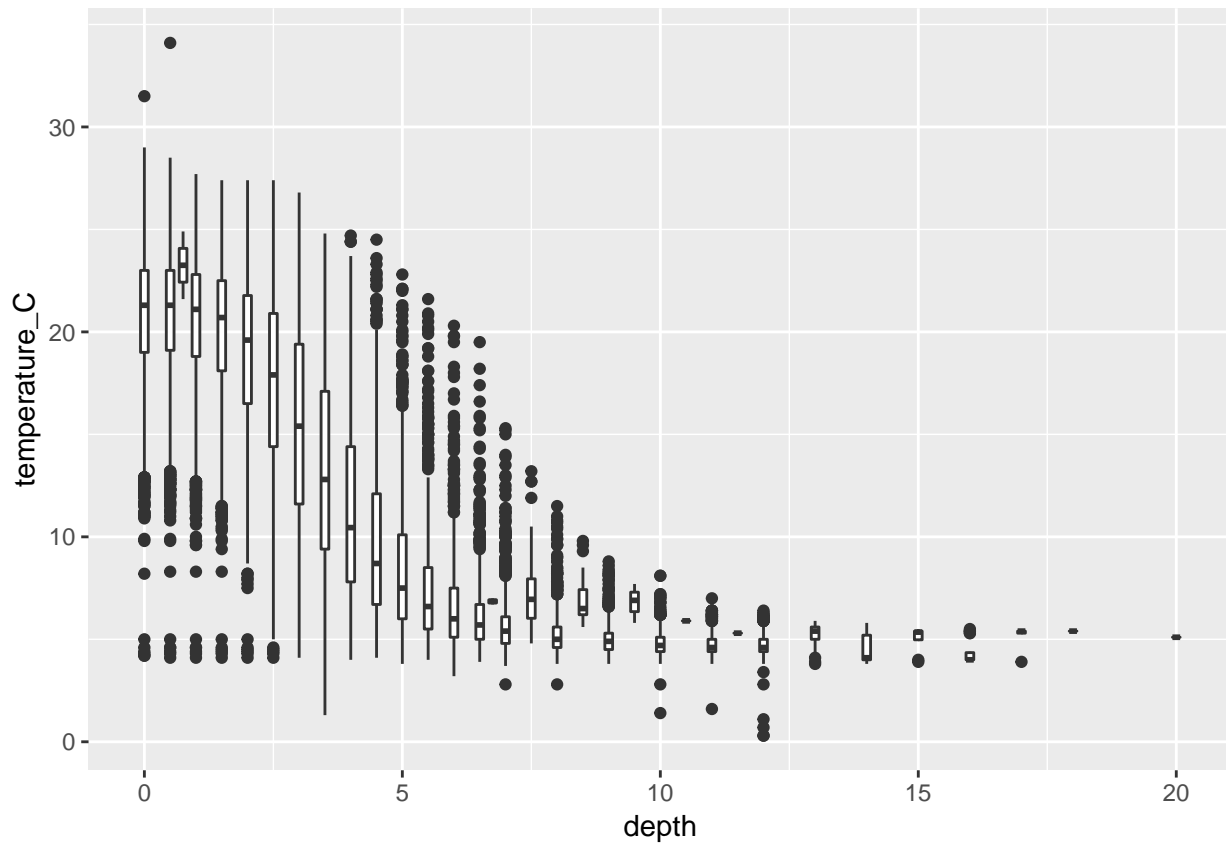
```
# 5
ggplot(lter.chemphys) +
  geom_boxplot(aes(x = lakename, y = temperature_C))
```

```
## Warning: Removed 3858 rows containing non-finite values (stat_boxplot).
```



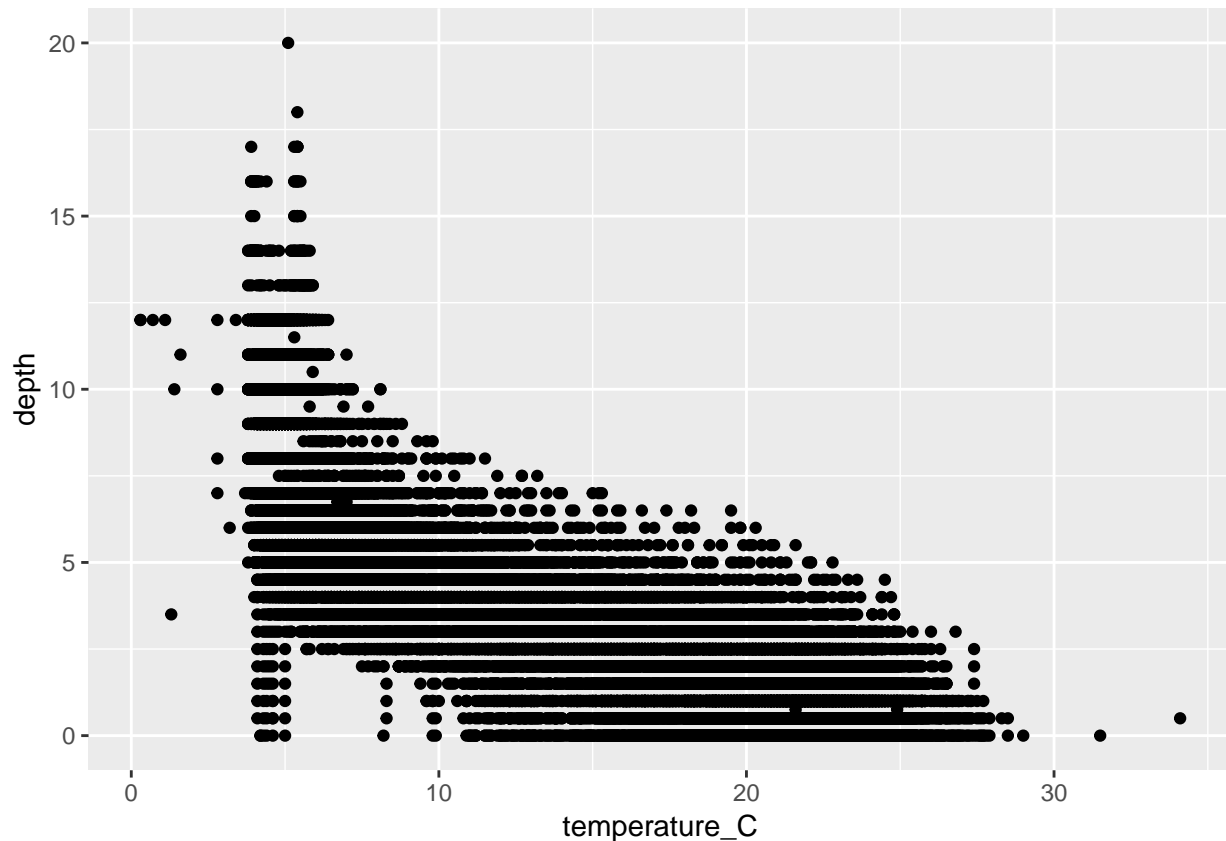
```
# 6
ggplot(lter.chemphys) +
  geom_boxplot(aes(x = depth, y = temperature_C, group = cut_width(depth, 0.25)))
```

```
## Warning: Removed 3858 rows containing non-finite values (stat_boxplot).
```

```
# 7  
ggplot(lter.chemphys) +  
  geom_point(aes(x = temperature_C, y = depth))
```

```
## Warning: Removed 3858 rows containing missing values (geom_point).
```



5) Form questions for further data analysis

What did you find out about your data from the basic summaries and graphs you made? Describe in 4-6 sentences.

ANSWER: 3,858 rows contain missing values. Paul Lake and Peter Lake have the highest counts for observations of temperature. The most common temperature is around 5 degreesC. East Long Lake had the highest recorded temperature. The coldest temp was at a depth of about 12m. The scatterplot seemed easier to read the relationship btwn depth and temp.

What are 3 further questions you might ask as you move forward with analysis of this dataset?

ANSWER 1: Is there a relationship btwn temperature and sample date?

ANSWER 2: Can we explore dissolved oxygen and how it might influence irradiance values?

ANSWER 3: Can we pull out the areas of the lakes that have the highest temperatures and try to find trends ?