

Assignment 3: Data Exploration

Jake Greif

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.

```
# Working directory
getwd()
```

```
## [1] "/Users/jakegreif/Environmental_Data_Analytics/Assignments"
```

```
# Packages
library(tidyverse)
```

```
## -- Attaching packages -----
## v ggplot2 3.1.0      v purrr  0.2.5
## v tibble  1.4.2      v dplyr  0.7.8
## v tidyr   0.8.1      v stringr 1.3.1
## v readr   1.1.1      v forcats 0.3.0

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

```
# Dataset
```

```
N.Temp.Lake.data <- read.csv("/Users/jakegreif/Environmental_Data_Analytics/Data/Raw/NTL-LTER_Lake_ChemistryPhysics_Raw.csv")
```

```
# View dataset
```

```
View(N.Temp.Lake.data)
```

```
## Warning in system2("/usr/bin/otool", c("-L", shQuote(DSO)), stdout = TRUE):
## running command '/usr/bin/otool' -L '/Library/Frameworks/R.framework/
## Resources/modules/R_de.so' had status 1
```

```
# Define dataset type; get familiar with dataset layout
class(N.Temp.Lake.data)
```

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

```
colnames(N.Temp.Lake.data)
```

```
## [1] "lakeid"          "lakename"         "year4"
## [4] "daynum"          "sampledate"       "depth"
## [7] "temperature_C"   "dissolvedOxygen"  "irradianceWater"
## [10] "irradianceDeck"  "comments"
```

```
class(N.Temp.Lake.data$sampledate)
```

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

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: The three most important pieces of information gathered from the README file are the context of our data (where/when/how it was collected), what type of data was collected, and the naming conventions/file formats.

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, sampledate, depth, and temperature
5. summary of lakename, depth, and temperature

```
# 1
dim(N.Temp.Lake.data)
```

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

```
# 2
class(N.Temp.Lake.data)
```

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

```
# 3
head(N.Temp.Lake.data, 8)
```

```
##   lakeid lakename year4 daynum sampledate 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(N.Temp.Lake.data$lakename)
```

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

```
class(N.Temp.Lake.data$sampldate)
```

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

```
class(N.Temp.Lake.data$depth)
```

```
## [1] "numeric"
```

```
class(N.Temp.Lake.data$temperature_C)
```

```
## [1] "numeric"
```

```
# 5
```

```
summary(N.Temp.Lake.data$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(N.Temp.Lake.data$depth)
```

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

```
summary(N.Temp.Lake.data$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 sampldate to class = date. After doing this, write an R command to display that the class of sampldate is indeed date. Write another R command to show the first 10 rows of the date column.

```
N.Temp.Lake.data$sampldate <- as.Date(N.Temp.Lake.data$sampldate, format = "%m/%d/%y")
```

```
class(N.Temp.Lake.data$sampldate)
```

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

```
head(N.Temp.Lake.data$sampldate, 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: I want to remove the NAs from this dataset because I'd like to quantitatively analyze the data. If I do not remove the NAs, I won't be able to run some statistical models, such as a timeseries.

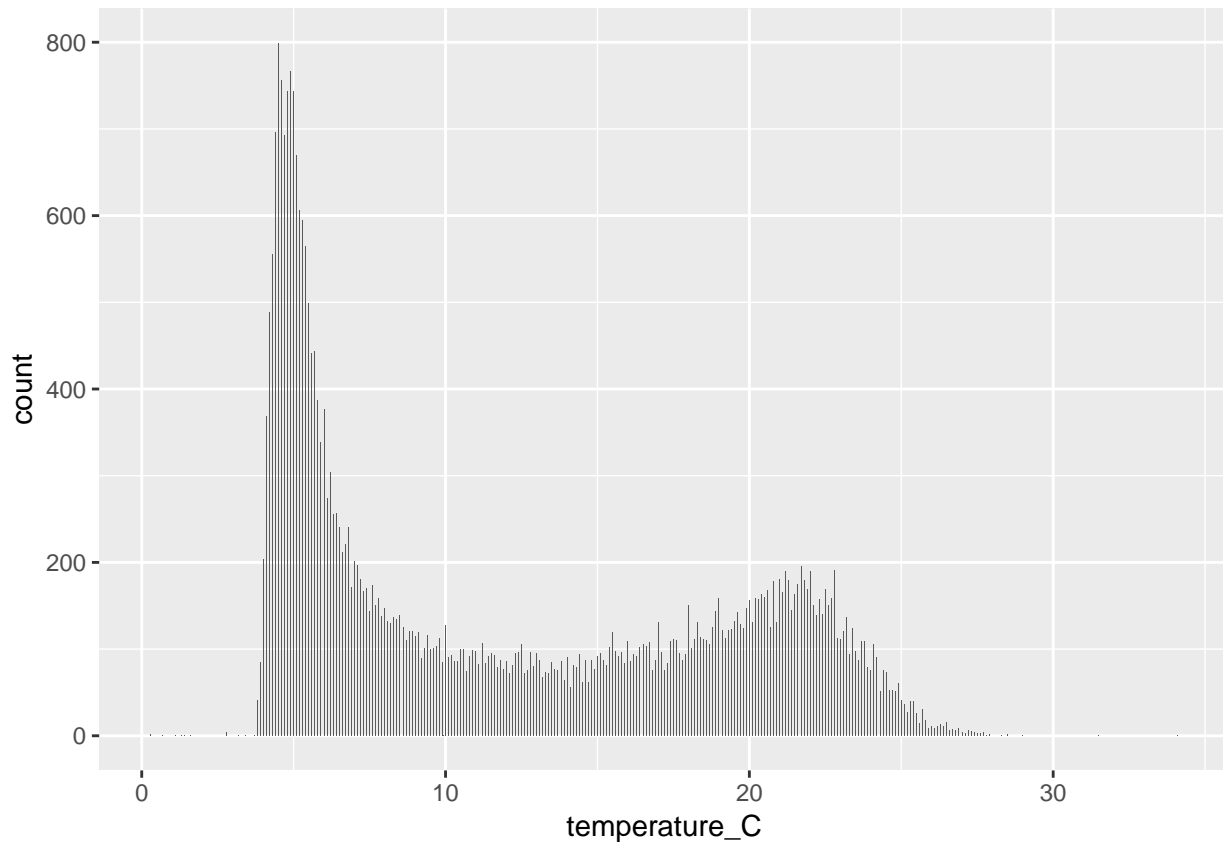
4) Explore your data graphically

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(N.Temp.Lake.data, aes(x = temperature_C)) +
  geom_bar()
```

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

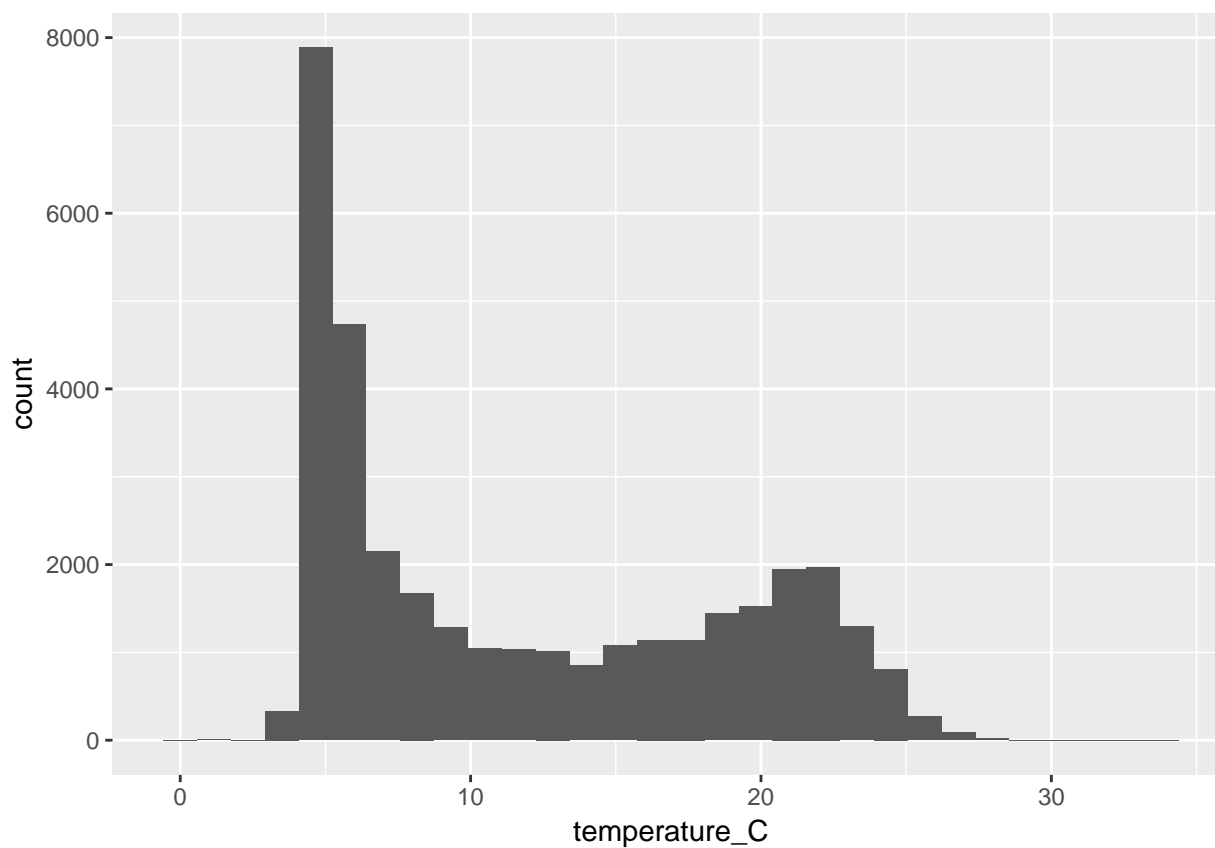


```
# 2
```

```
ggplot(N.Temp.Lake.data) +  
  geom_histogram(aes(x = temperature_C))
```

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

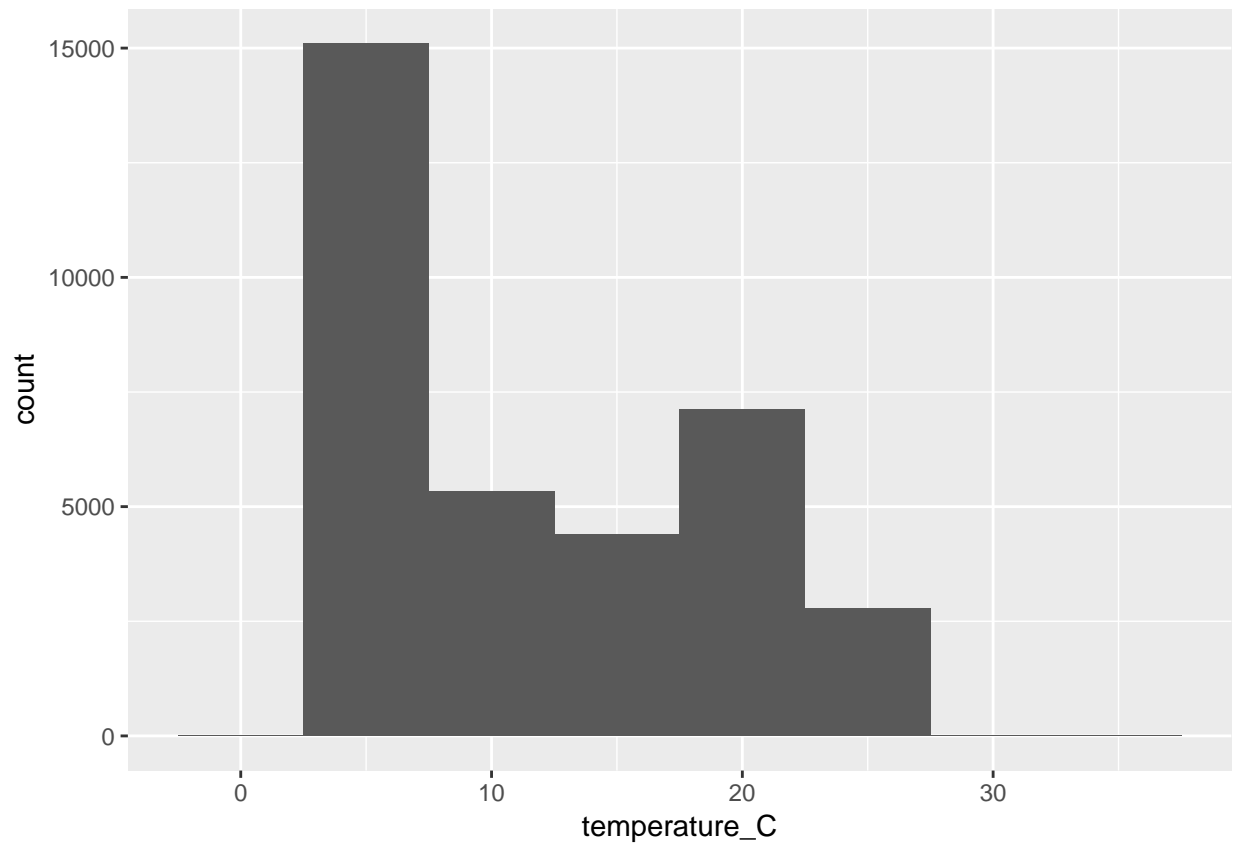
```
## Warning: Removed 3858 rows containing non-finite values (stat_bin).
```



```
# 3
```

```
ggplot(N.Temp.Lake.data) +  
  geom_histogram(aes(x = temperature_C), binwidth = 5)
```

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



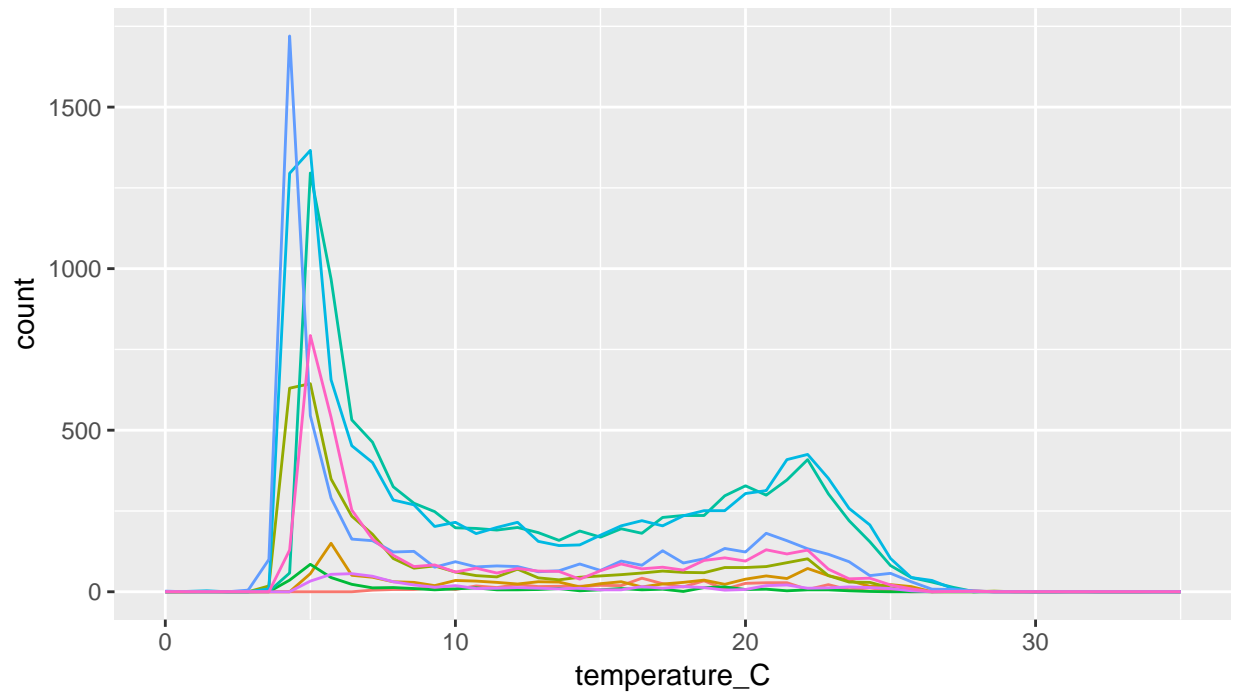
```
# 4
ggplot(N.Temp.Lake.data) +
  geom_freqpoly(aes(x = temperature_C, color = lakename), bins = 50) +
  scale_x_continuous(limits = c(0, 35)) +
  theme(legend.position = "top")
```

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

```
## Warning: Removed 18 rows containing missing values (geom_path).
```

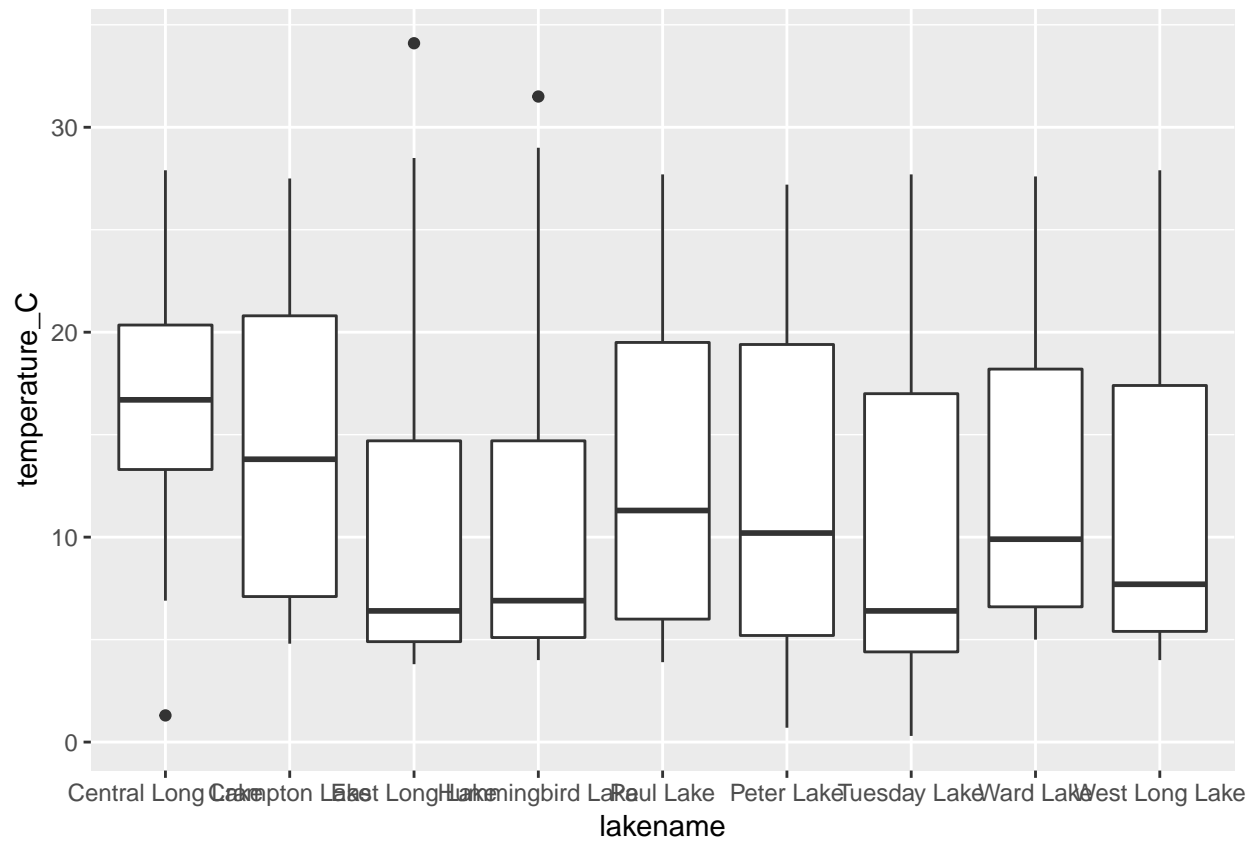
lakename

| | | | | |
|-------------------|------------------|------------|--------------|--------|
| Central Long Lake | East Long Lake | Paul Lake | Tuesday Lake | West L |
| Crampton Lake | Hummingbird Lake | Peter Lake | Ward Lake | |



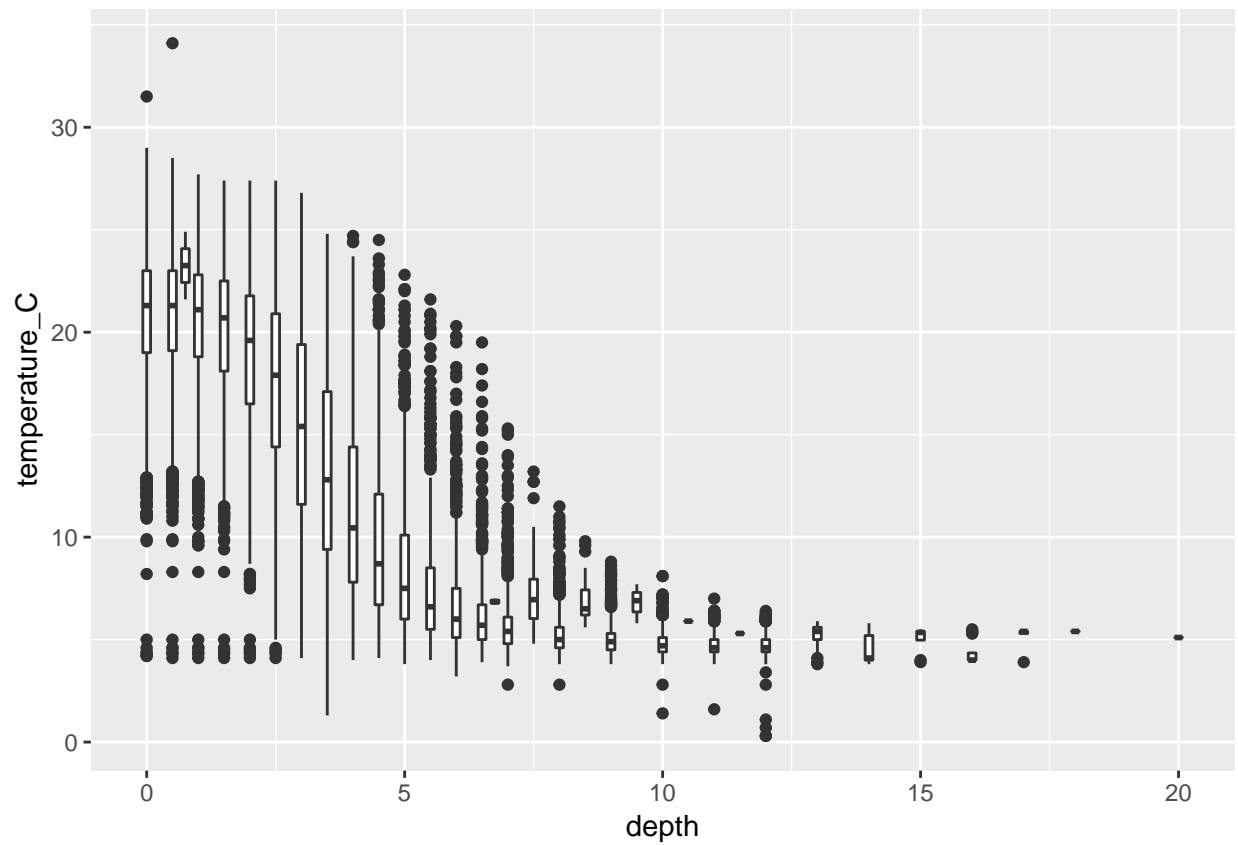
```
# 5
ggplot(N.Temp.Lake.data) +
  geom_boxplot(aes(x = lakename, y = temperature_C))
```

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



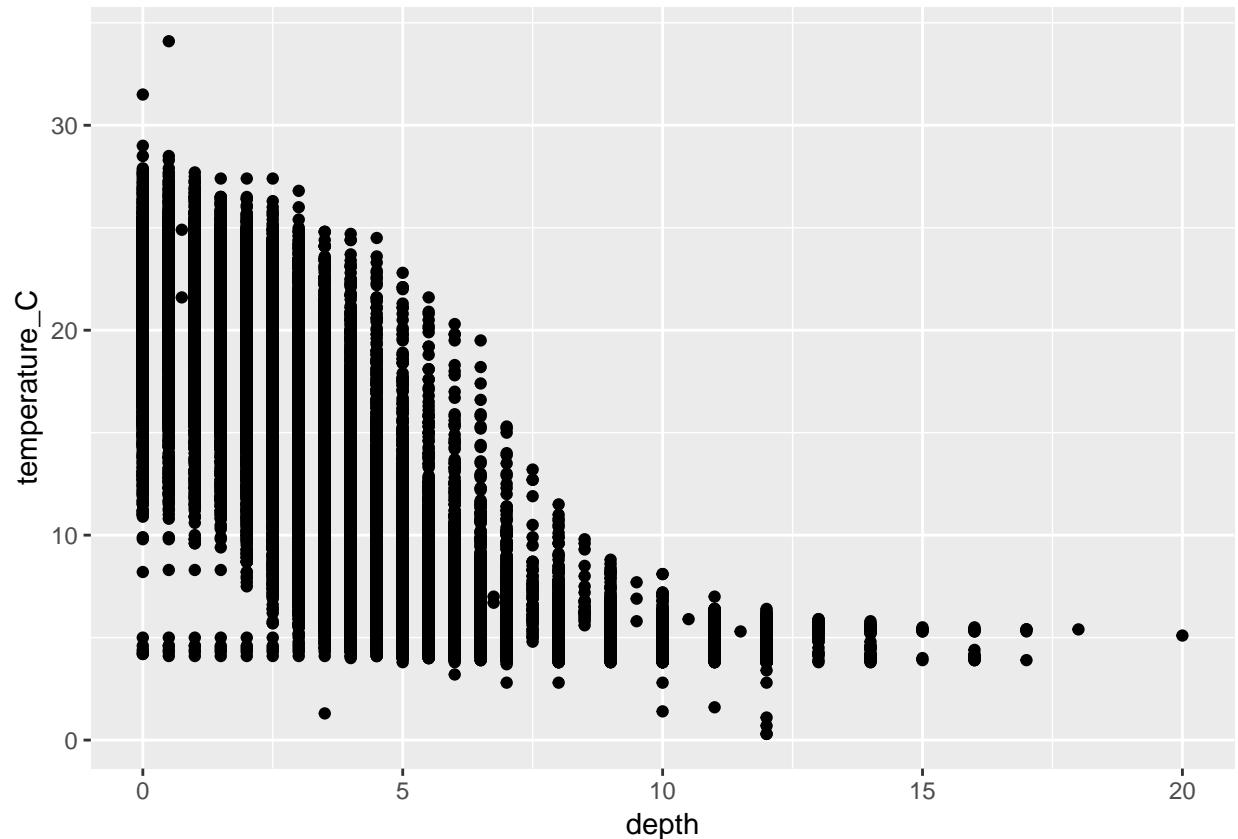
```
# 6
ggplot(N.Temp.Lake.data) +
  geom_boxplot(aes(x = depth, y = temperature_C, group = cut_width(depth, 0.25)))
```

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

```
# 7
ggplot(N.Temp.Lake.data) +
  geom_point(aes(y = temperature_C, x = 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: I learned that temperatures both decreases and become less variable as depth increases. I also learned that the lakes being monitored are generally cold, and have a proportional distribution of temperature readings, which is a good indication that they all reside in the same climate and are likely similar in size.

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

ANSWER 1: Does dissolved oxygen change with depth?

ANSWER 2: How does the temperature at a depth of X meters change temporally in each lake?

ANSWER 3: How does the dissolved oxygen content in each lake compare?