

Assignment 3: Data Exploration

Jiaqi Li

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.

```
#Check working directory
getwd()

## [1] "/Users/ljq/Desktop/Blue Devils/Data Analysis/ENV872_02/Assignments"

setwd("/Users/ljq/Desktop/Blue Devils/Data Analysis/ENV872_02")
# Load package
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.2.1 --
## v ggplot2 3.1.0      v purrr  0.2.5
## v tibble  1.4.2      v dplyr  0.7.7
## 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()

#Import Data
lake <- read.csv("../Data/Raw/NTL-LTER_Lake_ChemistryPhysics_Raw.csv")
```

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: Three salient pieces of information in the README file are database information, data content information, and naming conventions and file formats. Database information describes where, how, and when the data was collected. Data content information contains definitions and collection methods of each variable from the NTL-LTER database website. Naming conventions and file formats contains the rules of naming files.

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 lakenname, sampleddate, depth, and temperature
5. summary of lakenname, depth, and temperature

```
# 1 Dimension of the dataset
dim(lake)
```

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

```
# 2 Class of the dataset
class(lake)
```

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

```
# 3 First 8 rows of the dataset
head(lake, 8)
```

```
##   lakeid lakenname 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 of the variables
class(lake$lakenname)
```

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

```
class(lake$sampdate)
```

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

```
class(lake$depth)
```

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

```
class(lake$temperature_C)
```

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

```
# 5 Summary of the variables
```

```
summary(lake$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(lake$depth)
```

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

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

```
lake$sampdate <- as.Date(lake$sampdate, format = "%m/%d/%y")
```

```
class(lake$sampdate)
```

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

```
head(lake$sampdate, 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: No, the reason that NAs should not be removed is that NAs have meanings that some of the data was not collected or not available in this dataset. It may cause inaccuracy in the data analysis if NAs are removed.

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

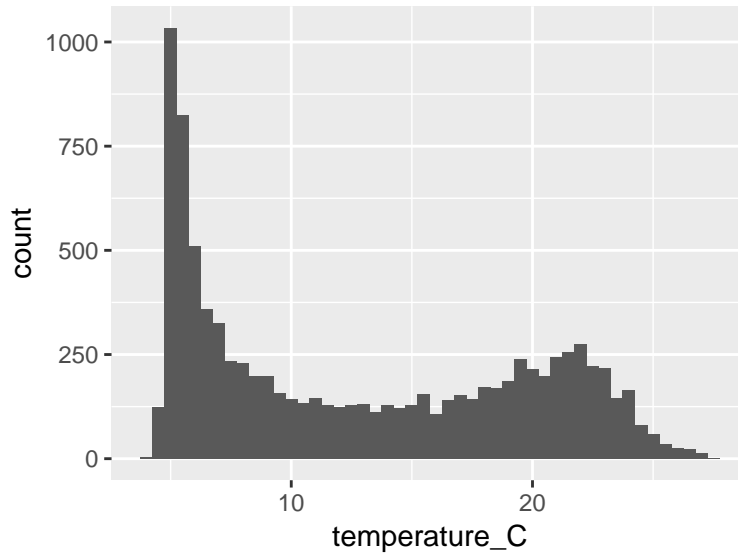


Figure 1: Temperature at Paul Lake

6. Boxplot of temperature based on depth, with depth divided into 0.25 m increments
7. Scatterplot of temperature by depth

1 Bar chart of temperature counts for each lake

```
lakenames <- unique(lake$lakename)
```

```
lakenames <- c("Paul Lake", "Peter Lake", "Tuesday Lake", "East Long Lake", "West Long Lake", "Central Lake")
```

```
ggplot(lake[lake$lakename == c('Paul Lake'),], aes(x = temperature_C)) +  
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use  
## `geom_histogram()` instead.
```

```
## Warning: Removed 1072 rows containing non-finite values (stat_bin).
```

```
ggplot(lake[lake$lakename == c('Peter Lake'),], aes(x = temperature_C)) +  
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use  
## `geom_histogram()` instead.
```

```
## Warning: Removed 1099 rows containing non-finite values (stat_bin).
```

```
ggplot(lake[lake$lakename == c('Tuesday Lake'),], aes(x = temperature_C)) +  
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use  
## `geom_histogram()` instead.
```

```
## Warning: Removed 604 rows containing non-finite values (stat_bin).
```

```
ggplot(lake[lake$lakename == c('East Long Lake'),], aes(x = temperature_C)) +  
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use  
## `geom_histogram()` instead.
```

```
## Warning: Removed 355 rows containing non-finite values (stat_bin).
```

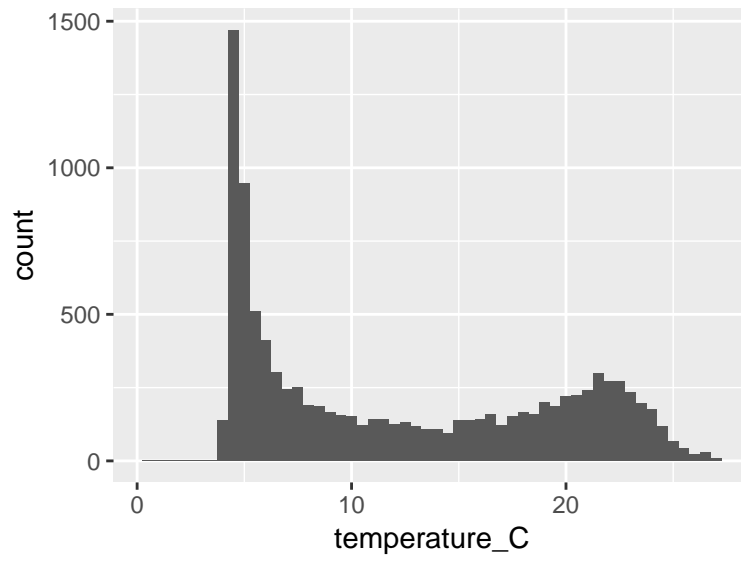


Figure 2: Temperature at Peter Lake

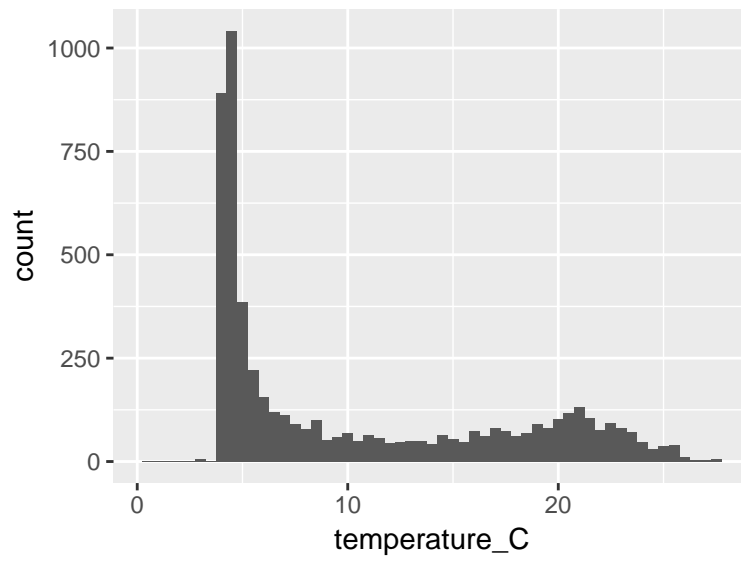


Figure 3: Temperature at Tuesday Lake

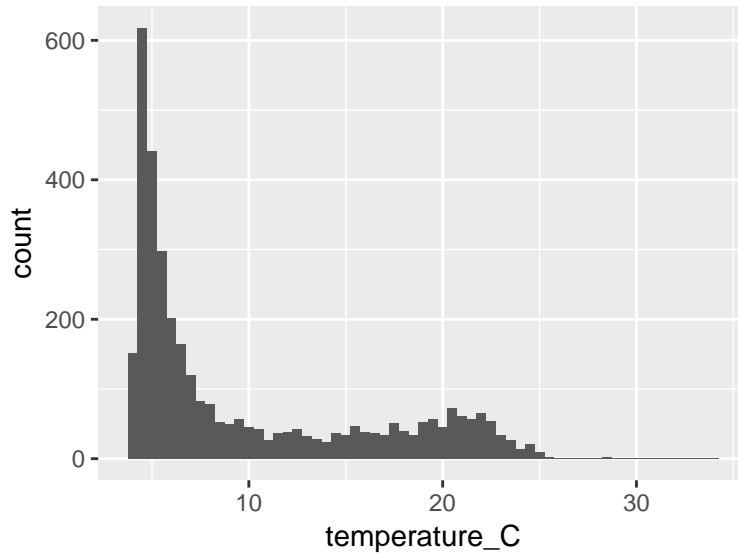


Figure 4: Temperature at East Long Lake

```
ggplot(lake[lake$lakename == c('West Long Lake'),], aes(x = temperature_C)) +
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use
## `geom_histogram()` instead.
```

```
## Warning: Removed 383 rows containing non-finite values (stat_bin).
```

```
ggplot(lake[lake$lakename == c('Central Long Lake'),], aes(x = temperature_C)) +
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use
## `geom_histogram()` instead.
```

```
## Warning: Removed 96 rows containing non-finite values (stat_bin).
```

```
ggplot(lake[lake$lakename == c('Hummingbird Lake'),], aes(x = temperature_C)) +
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use
## `geom_histogram()` instead.
```

```
## Warning: Removed 52 rows containing non-finite values (stat_bin).
```

```
ggplot(lake[lake$lakename == c('Crampton Lake'),], aes(x = temperature_C)) +
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use
## `geom_histogram()` instead.
```

```
## Warning: Removed 126 rows containing non-finite values (stat_bin).
```

```
ggplot(lake[lake$lakename == c('Ward Lake'),], aes(x = temperature_C)) +
  geom_bar(binwidth=0.5)
```

```
## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use
## `geom_histogram()` instead.
```

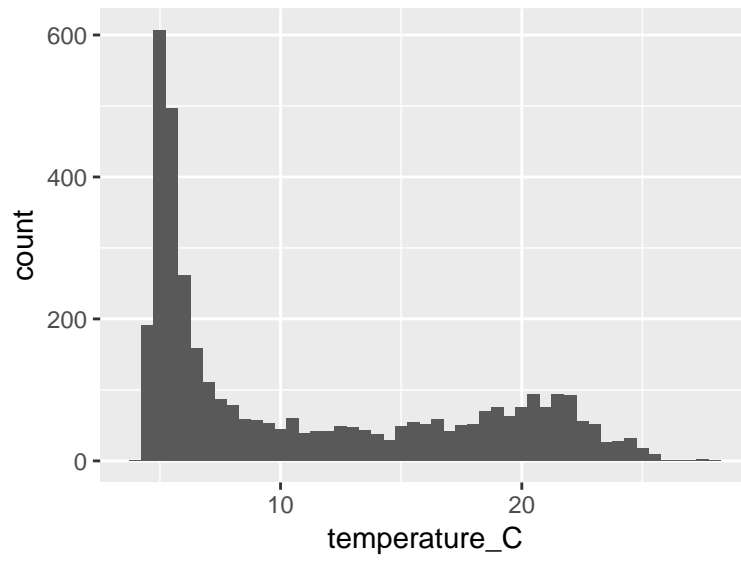


Figure 5: Temperature at West Long Lake

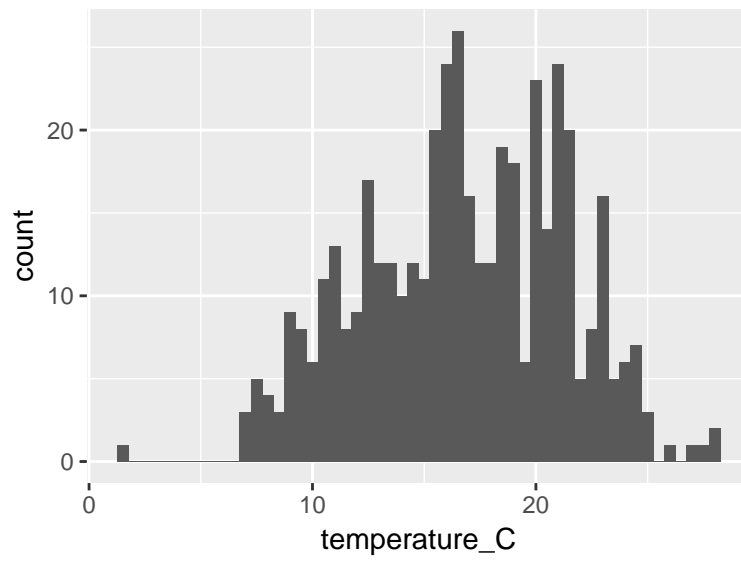


Figure 6: Temperature at Central Long Lake

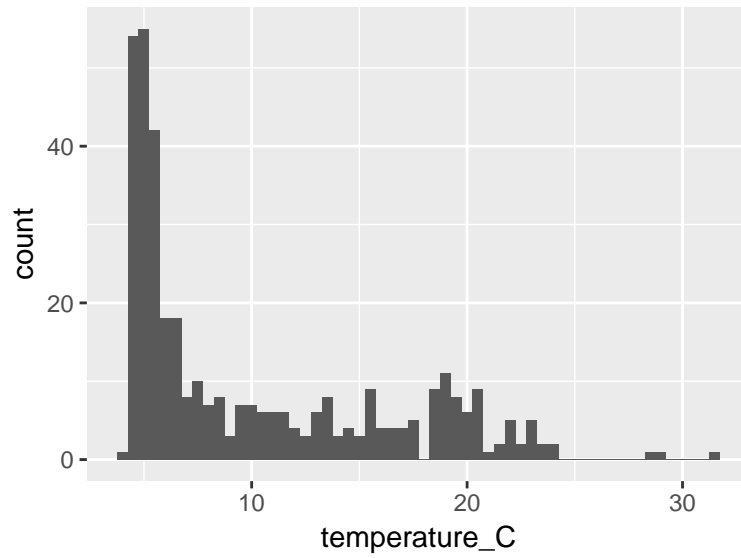


Figure 7: Temperature at Hummingbird Lake

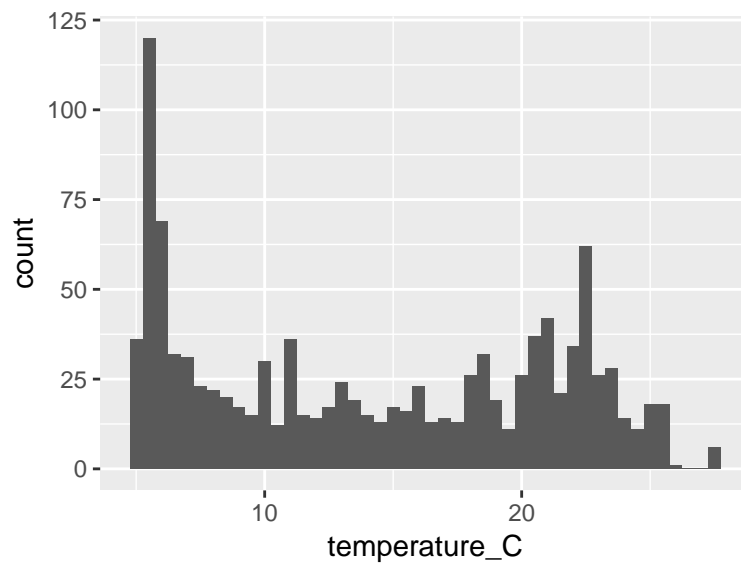


Figure 8: Temperature at Crampton Lake

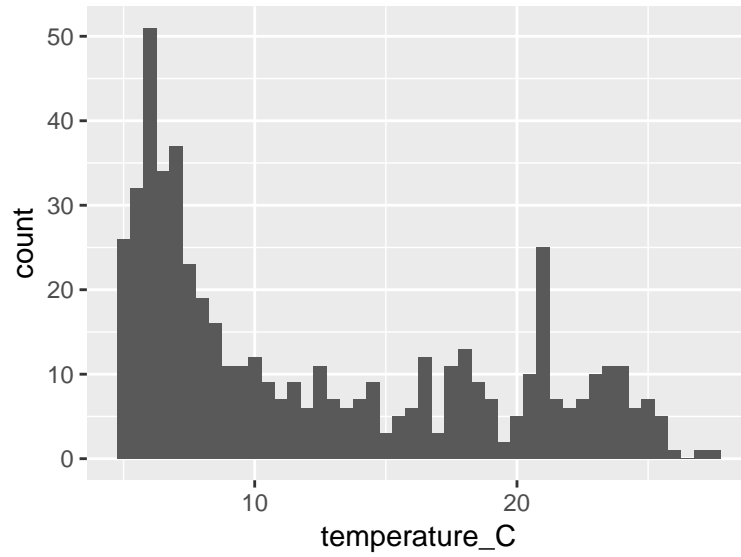


Figure 9: Temperature at Ward Lake

```
## Warning: Removed 71 rows containing non-finite values (stat_bin).
```

```
# 2
# 3
# 4
# 5
# 6
# 7
```

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:

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

ANSWER 1:

ANSWER 2:

ANSWER 3: