

## Buoy Project: Analysis of Climate Change

Jung Hwa Yeom

MA615: Data Science in R

September 25, 2020

## 1. Introduction

To begin with, we need to find whether there is evidence of global warming in the weather buoy. We read 20 years of data NDBC Station 44013 in NOAA National Data Buoy Center. This project requires us to process from acquiring data to presenting results. Additionally, there are tasks to make substitutions for NA data and decide an appropriate sampling frequency. To determine sampling frequency, we reduce the number of observations in the dataset. The dataset consists of year, month, date, and time with other various variables. Our group focuses on monthly average air temperature to observe global warming trends. We compare time versus air temperature from 1988 to 2018. If there is a positive slope in the plots, we can conclude global warming has progressed. Furthermore, as the slope is steeper, it indicates global warming drastically proceeds.

## 2. Explanation of R codes

```
1 library(tidyverse)
2 library(ggplot2)
3 library(rstanarm)
4 library(ggpubr)
5 library(dplyr)
6 library(tidyr)
7 # library(stringr)
8
9
```

At the beginning of the codes, we show all the packages we use in this project. We used tidyverse, ggplot2, rstanarm, ggpubr, dplyr, and tidyr.

```

10
11 ### make URLs
12
13 url1 <- "https://www.ndbc.noaa.gov/view_text_file.php?filename=44013h"
14 url2 <- ".txt.gz&dir=data/historical/stdmet/"
15
16 years <- c(1998:2018)
17
18 urls <- str_c(url1, years, url2, sep = "")
19
20 filenames <- str_c("buoy", years, sep = "")
21

```

Before reading data, we need to bring data from NDBC Station 44013 in NOAA National Data Buoy Center. Both “url1” and “url2” enable reading data from NDBC Station 44013. We use the data from 1988 to 2018, so we extracted those data by setting “years.” The “urls” shows an url in the formats of url1 and url2. If the data is from 1988, then the filename is “buoy1988.” We include buoy and years in the filenames to easily indicate the data year.

```

22 ### Read the data from the website
23
24 N <- length(urls)
25
26 for (i in 1:N){
27   if (i <= 7) {
28     suppressMessages( ### This stops the annoying messages on your screen. Do this last.
29       assign(filenames[i], read_table(urls[i], col_names = c("YYYY", "MM", "DD", "hh",
30         "WDIR", "WSPD", "GST", "WVHT", "DPD", "APD",
31         "MWD", "PRES", "ATMP", "WTMP", "DEWP", "VIS", "TIDE"), skip = 1)
32     )}
33   else {
34     suppressMessages( ### This stops the annoying messages on your screen. Do this last.
35       assign(filenames[i], read_table(urls[i], col_names = c("YYYY", "MM", "DD", "hh", "mm",
36         "WDIR", "WSPD", "GST", "WVHT", "DPD", "APD",
37         "MWD", "PRES", "ATMP", "WTMP", "DEWP", "VIS", "TIDE"), skip = 2))
38     )}
39 }
40
41 file <- get(filenames[i])
42 colnames(file)[1] <- "YYYY"
43
44 # put '19' in front of 98
45 if (i == 1) {
46   file[i] = file[i] + 1900
47 }
48 }
49
50 if (i <= 3) {
51   file[, 'TIDE'] = 0
52 }
53
54 if (i <= 7) {
55   file[, 'mm'] = 0
56 }
57
58 if(i == 1){
59   MR <- file
60 }
61
62 else{
63   MR <- rbind.data.frame(MR, file)
64 }
65 }
66 }

```

This part is reading the data from the website. First of all, we define “N” as the length of urls. We used a for loop to read tables and to put ‘19’ in front of years. To further illustrate, in the tables, years before 2000 show only two last digits of the year. For example, if the year is 1988, then only ‘88’ is shown in the table. After processing the for loop, then all the tables will store in “MR” by using `rbind.data.frame` function.

```
67 MR <- filter(MR, MR$ATMP < 50)
68 MR$MM <- as.numeric(MR$MM)
69 two_pm <- MR[MR$hh %in% c("14"), ]
70
```

One of the tasks in this project is to make substitutions for NA data. For instance, there is data that air temperature is 999 celsius degrees, which is impossible and meaningless to use. Therefore, we remove the data which has over 50 Celsius degrees in the air temperature (ATMP) column. In addition, `as.numeric` function allows us to make “MR\$MM” a numeric variable. We select 2 pm because 2 pm is the hottest time in a day and we analyze the air temperature at 2 pm. “two\_pm” chooses “14” which means 2 pm in the “hh” columns.

```
71 ~ for (i in 1998:2018){
72   year <- two_pm[two_pm$YYYY %in% c(i),]
73   for (j in 1:12){
74     month <- year[year$MM %in% c(j),]
75     if (j == 1){
76       ave = mean(month$ATMP)
77     } else {
78       ave <- rbind(ave, mean(month$ATMP))
79     }
80   }
81   if (i == 1998) {
82     data_frame <- rbind.data.frame("year" = i, "ave_temp" <- ave)
83   } else {
84     data_frame <- cbind.data.frame(data_frame, rbind.data.frame("year" = i, "ave_temp" <- ave))
85   }
86 }
87 final.df <- as.data.frame(t(data_frame))
88
89 colnames(final.df) = c("year", "Jan", "Feb", "Mar", "Apr", "May", "Jun", "Jul", "Aug", "Sep", "Oct", "Nov", "Dec")
90
```

In this part, we use another for loop to calculate the average of each month. The for loop has similar properties with the first our for loop. We take an average of the month in 20 years by writing a mean function and store the average by using `rbind` function. Basically, we are making a new data frame that only has the average of the month at 2 pm from 1988 to 2018. “final.df” is our final data frame before we start to plot a graph. Lastly, we set up new

column names and the column names are years and months to better indicate which month of the average air temperature is.

```

91 # do the plot from here
92
93 jj<-ggplot(data=final.df, aes(x=year, y=Jan)) +
94   geom_point() +
95   labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
96 ff<- ggplot(data=final.df, aes(x=year, y=Feb)) +
97   geom_point()+
98   labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
99 mm <-ggplot(data=final.df, aes(x=year, y=Mar)) +
100  geom_point() +
101  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
102 aa<-ggplot(data=final.df, aes(x=year, y=Apr)) +
103  geom_point() +
104  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
105 mayy<-ggplot(data=final.df, aes(x=year, y=May)) +
106  geom_point() +
107  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
108 junn<-ggplot(data=final.df, aes(x=year, y=Jun)) +
109  geom_point() +
110  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
111 jull<-ggplot(data=final.df, aes(x=year, y=Jul)) +
112  geom_point() +
113  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
114 augg<-ggplot(data=final.df, aes(x=year, y=Aug)) +
115  geom_point() +
116  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
117 ss<-ggplot(data=final.df, aes(x=year, y=Sep)) +
118  geom_point() +
119  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
120 oo<-ggplot(data=final.df, aes(x=year, y=Oct)) +
121  geom_point() +
122  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
123 novv<-ggplot(data=final.df, aes(x=year, y=Nov)) +
124  geom_point() +
125  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)
126 dd<-ggplot(data=final.df, aes(x=year, y=Dec)) +
127  geom_point() +
128  labs(x="Year", y="ATMP", title= "") + geom_smooth(method = lm)

```

We make 12 plots because there are 12 months in a year; thus, it is effective to compare each month. “jj”, “ff”, “mm”, etc. mean each month. We use ggplot, geom\_plot, and geom\_smooth to plot the graph. The x-axis is “Year” and the y-axis is “ATMP” to compare time and air temperature. Furthermore, we do not remove shades in the plot because we want to observe confidence intervals.

```

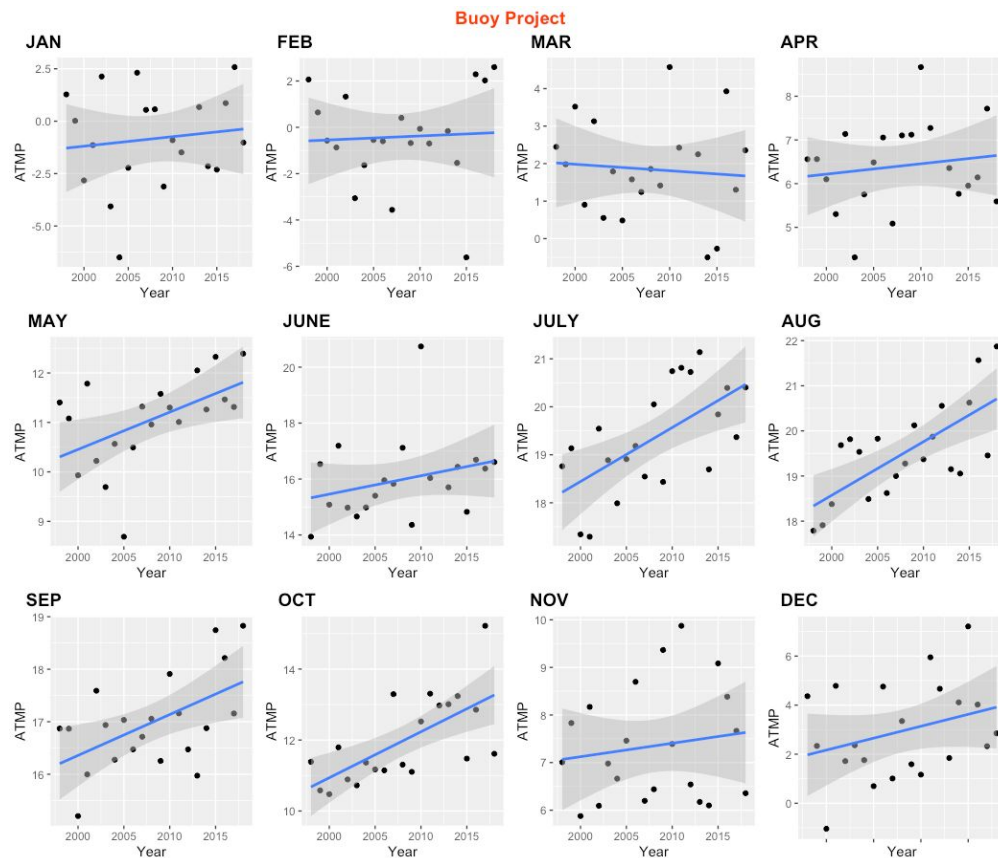
129 figure <- ggarrange(jj, ff, mm, aa, mayy, junn, jull, augg, ss, oo, novv, dd + rremove("x.text"),
130   labels = c("JAN", "FEB", "MAR", "APR", "MAY", "JUNE", "JULY", "AUG", "SEP", "OCT", "NOV", "DEC"),
131   ncol = 4, nrow = 3)
132 annotate_figure(figure,
133   top = text_grob("Buoy Project", color = "red", face = "bold", size = 14))

```

Finally, we need to put 12 graphs on one page to view the graphs at once. To combine all graphs on one page, we use ggarange function and label each month on the top left corner.

Annotate\_figure function allows us to format the graph such as labeling a title, assigning colors and font size.

### 3. Interpretation



The above graphs are the average air temperature of each month from 1998 to 2018. Grey shades around the lines are confidence intervals. If a graph has a wide range of the confidence interval, it means the global warming trend is not strong enough compared to a narrow one. The summer season has steeper positive slopes and narrow confidence intervals, so we can interpret air temperature is rising as time goes, showing global warming is rapidly progressing. To be specific, July and August are the months when temperatures had risen the most in 20 years. On the other hand, the winter season does not show a clear positive slope or narrow confidence intervals. However, it does not mean that global warming does not happen

in the winter season because not only air temperature is the factor that determines the level of global warming. Global warming leads to extreme weather both hotter summer and colder winter. To sum up, we can conclude global warming progressed from 1988 to 2018 and expect global warming will become more serious in the future.

## Bibliography

```
> citation(package = "tidyverse")
```

Wickham et al., (2019). Welcome to the tidyverse. Journal of Open Source Software, 4(43), 1686, <https://doi.org/10.21105/joss.01686>

```
> citation(package = "ggplot2")
```

To cite ggplot2 in publications, please use:

H. Wickham. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York, 2016.

```
> citation(package= "rstanarm")
```

To cite rstanarm in publications please use the first citation entry. If you were using the 'stan\_jm' modelling function then, where possible, please consider including the second citation entry as well.

Goodrich B, Gabry J, Ali I & Brilleman S. (2020). rstanarm: Bayesian applied regression modeling via Stan. R package version 2.21.1 <https://mc-stan.org/rstanarm>.

Brilleman SL, Crowther MJ, Moreno-Betancur M, Bueros Novik J & Wolfe R. Joint longitudinal and time-to-event models via Stan. StanCon 2018. 10-12 Jan 2018. Pacific Grove, CA, USA. [https://github.com/stan-dev/stancon\\_talks/](https://github.com/stan-dev/stancon_talks/)

```
> citation(package= "ggpubr")
```

To cite package 'ggpubr' in publications use:

Alboukadel Kassambara (2020). ggpubr: 'ggplot2' Based Publication Ready Plots. R package version 0.4.0. <https://CRAN.R-project.org/package=ggpubr>

```
> citation(package = "dplyr")
```

To cite package 'dplyr' in publications use:

Hadley Wickham, Romain François, Lionel Henry and Kirill Müller (2020). dplyr: A Grammar of Data Manipulation. R package version 1.0.2. <https://CRAN.R-project.org/package=dplyr>

```
> citation(package = "tidyr")
```

To cite package 'tidyr' in publications use:

Hadley Wickham (2020). tidyr: Tidy Messy Data. R package version 1.1.2. <https://CRAN.R-project.org/package=tidyr>