

# Stat-295 HW1

Sabahattin Alp Kocabaş, Başak Kabaloğlu, Efe Örencik

2024-04-20

## Question 1

Codes for Linux

- **To work on a file.**

```
#cd C:/Users/kocab/Desktop/Stat_295_HW1
```

- **(i) To read the data.**

```
#wget --content-disposition https://raw.githubusercontent.com/dhavalpotdar/College-Score-card-Data-Analysis/master/MERGED_2017_2018_cleaned.csv
```

- **(i) To change the name of data set.**

```
#mv MERGED_2017_2018_cleaned.csv college_score.csv
```

- **(i) To print the lines 10 through 60.**

```
#head -n 60 college_score.csv | tail -n 51
```

- **(ii) To create sub-sample of data set with appropriate conditions.**

```
#grep "Public" college_score.csv | grep ",Montgomery," > subsample.csv
```

- **(iii) To obtain the frequencies of each cities.**

```
#cut -d ',' -f2 college_score.csv | sort | uniq -c
```

## Outputs for 1

- **(i) The lines 10 through 60.**

- (iii) Head of obtained frequencies of each cities.

## Question 2

- ```
library(dplyr)
library(tidyverse)
```

- ```
chocolate <- read.csv("https://raw.githubusercontent.com/rfordatascience/tidytuesday/master/data/2022/2022-01-18/chocolate.csv")
head(chocolate)
```

```
##      ref company_manufacturer company_location review_date country_of_bean_origin
## 1 2454                5150          U.S.A.      2019          Tanzania
## 2 2458                5150          U.S.A.      2019      Dominican Republic
## 3 2454                5150          U.S.A.      2019          Madagascar
## 4 2542                5150          U.S.A.      2021             Fiji
## 5 2546                5150          U.S.A.      2021          Venezuela
## 6 2546                5150          U.S.A.      2021             Uganda
##      specific_bean_origin_or_bar_name cocoa_percent ingredients
## 1      Kokoa Kamili, batch 1          76% 3- B,S,C
## 2      Zorzal, batch 1          76% 3- B,S,C
## 3      Bejofo Estate, batch 1          76% 3- B,S,C
## 4      Matasawalevu, batch 1          68% 3- B,S,C
## 5      Sur del Lago, batch 1          72% 3- B,S,C
## 6      Semuliki Forest, batch 1          80% 3- B,S,C
##      most_memorable_characteristics rating
## 1      rich cocoa, fatty, bready 3.25
## 2      cocoa, vegetal, savory 3.50
## 3      cocoa, blackberry, full body 3.75
## 4      chewy, off, rubbery 3.00
## 5      fatty, earthy, moss, nutty,chalky 3.00
## 6      mildly bitter, basic cocoa, fatty 3.25
```

- (i) Examining the structure and comments.

```
str(chocolate)
```

```
## 'data.frame':    2530 obs. of  10 variables:
## $ ref              : int  2454 2458 2454 2542 2546 2546 2542 797 797
1011 ...
## $ company_manufacturer : chr  "5150" "5150" "5150" "5150" ...
## $ company_location    : chr  "U.S.A." "U.S.A." "U.S.A." "U.S.A." ...
## $ review_date         : int  2019 2019 2019 2021 2021 2021 2021 2012 201
2 2013 ...
## $ country_of_bean_origin : chr  "Tanzania" "Dominican Republic" "Madagasca
r" "Fiji" ...
## $ specific_bean_origin_or_bar_name: chr  "Kokoa Kamili, batch 1" "Zorzal, batch 1"
"Bejofo Estate, batch 1" "Matasawalevu, batch 1" ...
## $ cocoa_percent      : chr  "76%" "76%" "76%" "68%" ...
## $ ingredients        : chr  "3- B,S,C" "3- B,S,C" "3- B,S,C" "3- B,S,C"
...
## $ most_memorable_characteristics : chr  "rich cocoa, fatty, bready" "cocoa, vegeta
l, savory" "cocoa, blackberry, full body" "chewy, off, rubbery" ...
## $ rating              : num  3.25 3.5 3.75 3 3 3.25 3.5 3.5 3.75 2.75
...
```

The data set likely contains information about chocolate products, including details like manufacturer, origin, ingredients, characteristics, and ratings.

- (ii) Converting all the characters into factors.

```
chocolate <- chocolate %>%
  mutate_if(is.character, factor)
```

- (iii) Obtaining some statistics with respect to different company locations. Printing the 10 observations. And comments.

```
chocolate %>%
  group_by(company_location) %>%
  summarise(mean_rating = mean(rating),
            sd_rating = sd(rating),
            median_rating = median(rating),
            range_rating = diff(range(rating))) %>%
  head(10)
```

```
## # A tibble: 10 × 5
##   company_location mean_rating sd_rating median_rating range_rating
##   <fct>            <dbl>     <dbl>         <dbl>         <dbl>
## 1 Amsterdam        3.31      0.264         3.25          0.75
## 2 Argentina        3.31      0.349         3.5           1
## 3 Australia        3.36      0.409         3.5           1.5
## 4 Austria          3.26      0.325         3.25           1
## 5 Belgium          3.10      0.661         3             3
## 6 Bolivia          3.25      0.707         3.25           1
## 7 Brazil           3.28      0.356         3.25           1.5
## 8 Canada           3.30      0.416         3.25           2
## 9 Chile            3.75      0           3.75           0
## 10 Colombia        3.20      0.425         3.25           1.75
```

If we consider 10 rows of statistics, mean ratings vary 3.1 to 3.75 , standart deviatons is also vary and there are differences among locations. Median ratings are between 3.00 to 3.75 for this 10 examination. And there are range ratings that varies from location to location between 0 to 3.

- (iv) Finding the chocolates that its review date is equal to 2020 and country of bean origin is equal to Colombia.

```
chocolate %>%
  filter(review_date == 2020 & country_of_bean_origin == "Colombia")
```

```
##   ref company_manufacturer company_location review_date country_of_bean_origin
## 1 2466      Crow and Moss      U.S.A.      2020      Colombia
## 2 2534      El Buen      U.S.A.      2020      Colombia
## 3 2482      Finnia      Canada      2020      Colombia
## 4 2478      Kin + Pod      Canada      2020      Colombia
## 5 2498      Odyssey      U.S.A.      2020      Colombia
##   specific_bean_origin_or_bar_name cocoa_percent ingredients
## 1      Aruaca, batch 39      70%      2- B,S
## 2      Tumaco      70%      2- B,S
## 3      Chigorodo, batch 001      70%      3- B,S,C
## 4      Tumaco, batch 113      70%      3- B,S,C
## 5      Arhuaca      70%      2- B,S
##   most_memorable_characteristics rating
## 1      nutty, citrus      3.25
## 2      cocoa, spice, alcohol, dirty      3.00
## 3      nutty, melon, vinegar      2.75
## 4      poor finish, cocoa, grassy      3.25
## 5      walnut, tobacco, grassy      3.25
```

- (v) Taking the mean of chocolate rating and cocoa percent according to the company location.

```
chocolate %>%
  group_by(company_location) %>%
  summarise(mean_rating = mean(rating),
            mean_cocoa_percent = mean(as.numeric(factor(cocoa_percent))))
```

```
## # A tibble: 67 × 3
##   company_location mean_rating mean_cocoa_percent
##   <fct>             <dbl>             <dbl>
## 1 Amsterdam         3.31                 3
## 2 Argentina         3.31                1.67
## 3 Australia         3.36                5.89
## 4 Austria           3.26                6.47
## 5 Belgium           3.10                7.95
## 6 Bolivia           3.25                 1.5
## 7 Brazil            3.28                 6
## 8 Canada            3.30                9.84
## 9 Chile             3.75                 1
## 10 Colombia         3.20                6.90
## # i 57 more rows
```

- (vi) Selecting company manufacturer, company location and country of bean origin shortly. Printing the first 10 rows of the data frame.

```
chocolate %>%
  select(starts_with("c")) %>%
  head(10)
```

```
##   company_manufacturer company_location country_of_bean_origin cocoa_percent
## 1                    5150           U.S.A.             Tanzania          76%
## 2                    5150           U.S.A. Dominican Republic          76%
## 3                    5150           U.S.A.             Madagascar          76%
## 4                    5150           U.S.A.                Fiji          68%
## 5                    5150           U.S.A.             Venezuela          72%
## 6                    5150           U.S.A.                Uganda          80%
## 7                    5150           U.S.A.                India          68%
## 8                   A. Morin           France             Bolivia          70%
## 9                   A. Morin           France                Peru          63%
## 10                  A. Morin           France             Panama          70%
```

- (vii) Filtering that company location in Switzerland whose rating between 3.25 and 3.5. Taking the five observations.

```
chocolate %>%
  filter(company_location == "Switzerland" & rating >= 3.25 & rating <= 3.5) %>%
  head(5)
```

```
##   ref company_manufacturer company_location review_date country_of_bean_origin
## 1 508 Beschle (Felchlin)      Switzerland      2010      Venezuela
## 2 508 Beschle (Felchlin)      Switzerland      2010      Venezuela
## 3 508 Beschle (Felchlin)      Switzerland      2010      Indonesia
## 4 508 Beschle (Felchlin)      Switzerland      2010      Venezuela
## 5 636 Beschle (Felchlin)      Switzerland      2011      Venezuela
##       specific_bean_origin_or_bar_name cocoa_percent ingredients
## 1   Carenero S., Barlovento, Grand Cru      70%    3- B,S,C
## 2 Porcelana, Premier Cru, Quizas No. 1      74%    3- B,S,C
## 3                               Java, Grand Cru      64%    3- B,S,C
## 4   Ocumare, Premier Cru, Quizas No. 2      72%    3- B,S,C
## 5 Indigena Amazonia, Grand Cru, Quizas      72%    3- B,S,C
##       most_memorable_characteristics rating
## 1      creamy, macadamia, pepper    3.25
## 2 nutty, light toffee, mild musty    3.25
## 3      ham-like, smokey, banana    3.50
## 4      dark cocoa, spicy pepper    3.50
## 5      creamy, banana, rich    3.50
```

- (viii) Mean of the rating column for each company locations that ordered by descending.

```
chocolate %>%
  group_by(company_location) %>%
  summarise(mean_rating = mean(rating)) %>%
  arrange(desc(mean_rating))
```

```
## # A tibble: 67 × 2
##   company_location mean_rating
##   <fct>           <dbl>
## 1 Chile           3.75
## 2 U.A.E.          3.4
## 3 Poland          3.38
## 4 Vietnam         3.36
## 5 Australia       3.36
## 6 Guatemala       3.35
## 7 Denmark         3.34
## 8 Norway          3.33
## 9 Switzerland    3.32
## 10 Amsterdam      3.31
## # i 57 more rows
```

- (ix) Counting the observations are assigned Bonnat for each country of bean origin.

```
chocolate %>%
  filter(company_manufacturer == "Bonnat") %>%
  count(country_of_bean_origin, sort = TRUE)
```

```
##      country_of_bean_origin n
## 1                      Blend 4
## 2                      Brazil 4
## 3                      Peru 4
## 4                      Venezuela 4
## 5                      Mexico 3
## 6                      Madagascar 2
## 7                      Cuba 1
## 8      Dominican Republic 1
## 9                      Ecuador 1
## 10                     Gabon 1
## 11                     Haiti 1
## 12                     Ivory Coast 1
## 13                     Jamaica 1
## 14                     Nicaragua 1
## 15                     Sri Lanka 1
```

- **(x) Creating a new column called Rating Percentage, which is percentage version of the rating column & rating the chocolates.**

```
chocolate <- chocolate %>%
  mutate(Rating_Percentage = rating / 4 * 100,
         Class = case_when(
           Rating_Percentage < 25 ~ "Low",
           Rating_Percentage < 50 ~ "Medium",
           Rating_Percentage <= 87.5 ~ "Tasty",
           Rating_Percentage >87.5 ~ "Excellent"
         ))
```

## Question 3

- **Calling necessary libraries.**

```
library(ggplot2)
library(dplyr)
library(tidyverse)
```

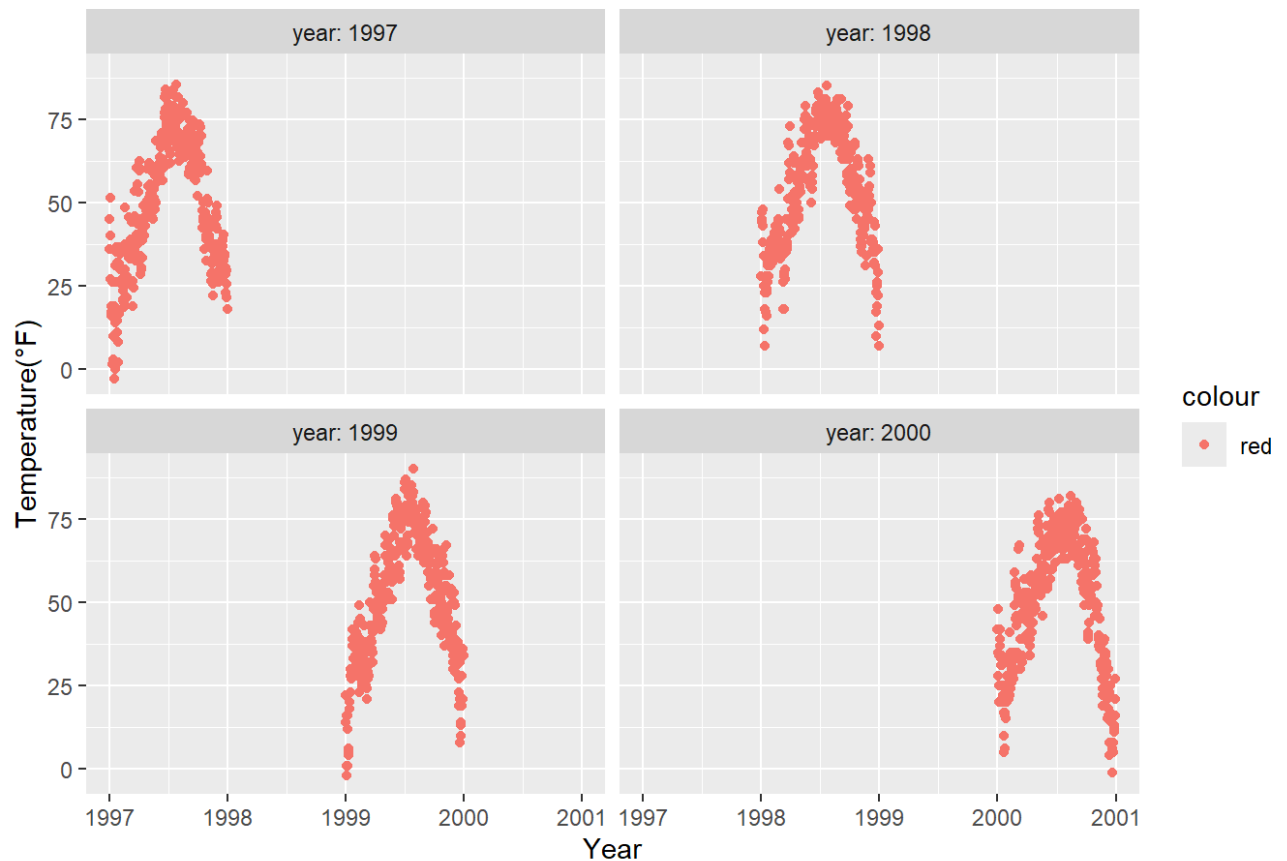
- **(i) Reading data and named it as nmmaps.**

```
nmmaps<-read.csv("https://www.cedricscherer.com/data/chicago-nmmaps-custom.csv")
```

- **(i)Examining the relationship between date and temp based on the year.**

```
ggplot(nmmaps,aes(as.Date(date),temp,color="red"))+
  geom_point()+
  labs(x = "Year", y = "Temperature(°F)",
       title = "Relationship between Date & Temperature based on Year")+
  facet_wrap(~year, labeller = label_both)
```

## Relationship between Date & Temperature based on Year



- **Interpretation for plot 1.**

Although there are no major changes in general, temperatures in the last months of the year have started to decrease over the years. The lowest temperature was reached in early 1997, and the highest temperature was reached in mid-1999. In general, the temperature first increased and then decreased within a year.

- **(ii) Factoring season variables.**

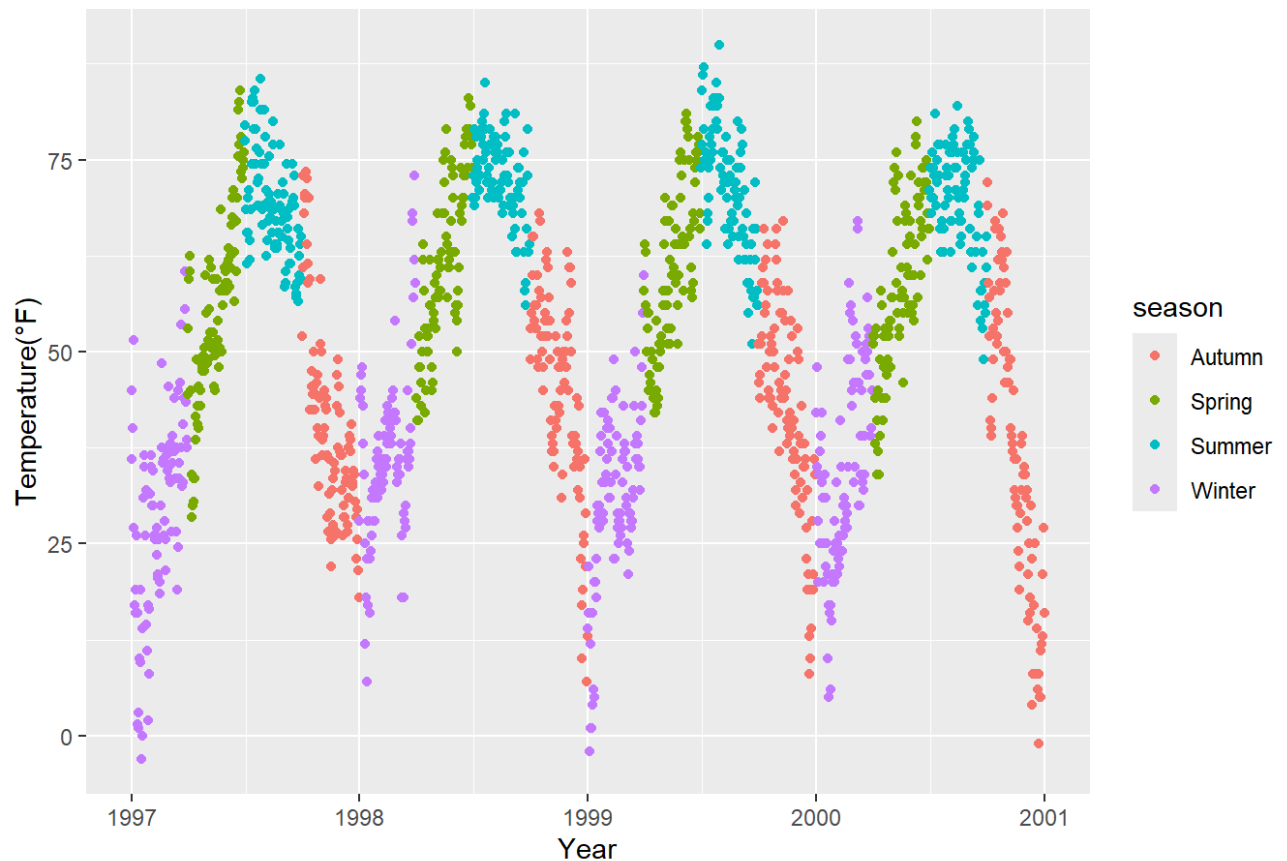
```
nmmaps$season<-factor(nmmaps$season,labels = c("Autumn", "Spring", "Summer", "Winter"))
```

- **(ii) Examining the relationship between date, temp and season.**

```
ggplot(nmmaps,aes(as.Date(date),temp,color=season))+
  geom_point()+
  labs(x = "Year", y = "Temperature(°F)", title = "Relationship between Date, Temperature, and Season" )
```



Relationship between Date, Temperature, and Season

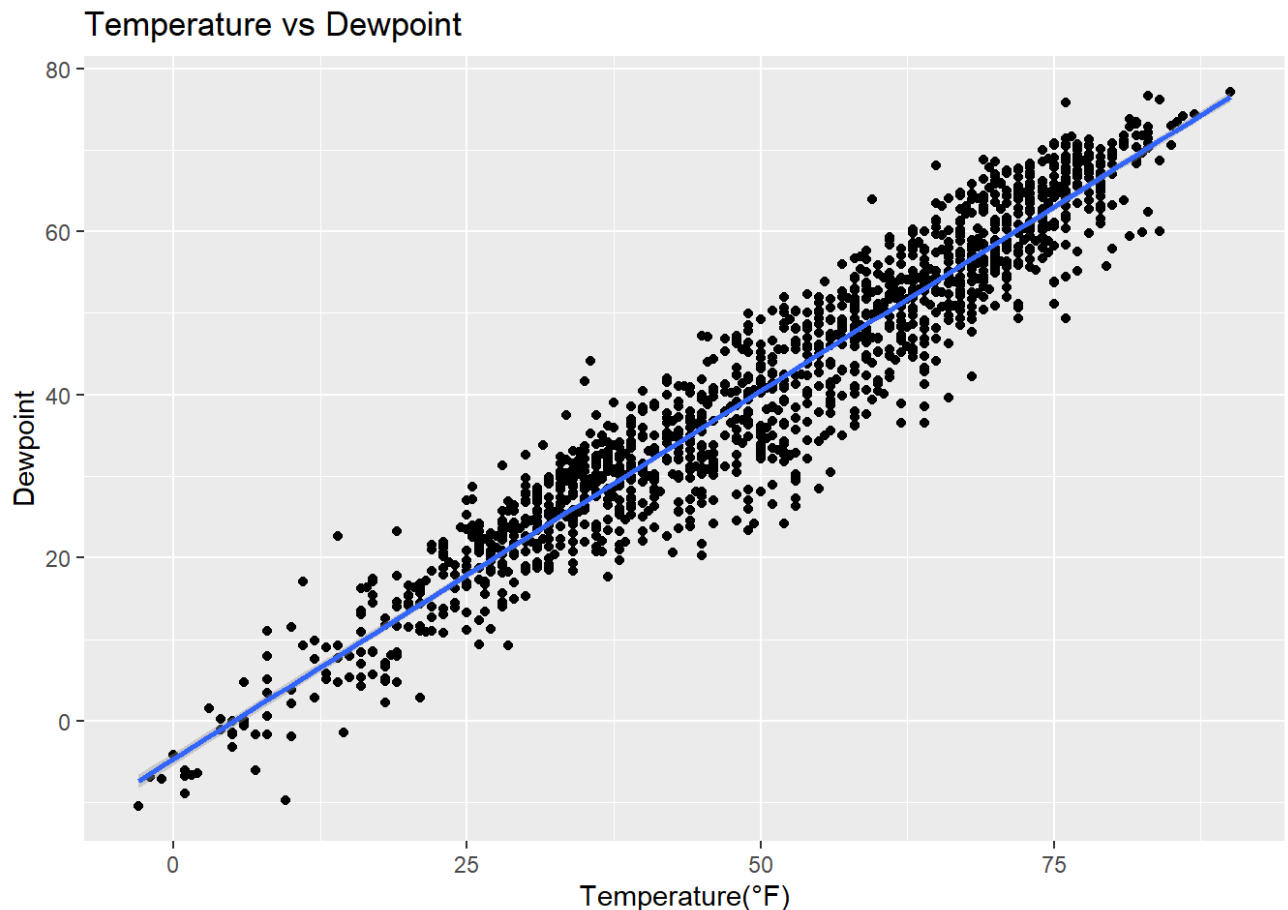


- **Interpretation for plot 2.**

The temperature change generally followed a regular pattern. Temperature change was more vary in winter and autumn. As we see in the first graph, the temperature reached its lowest value at the beginning of 1997 and its peak in the middle of 1999.

- **(iii) The relationship between temperature and dewpoint, and the correlation.**

```
ggplot(nmmaps,aes(temp,dewpoint))+
  geom_point()+
  labs(x = "Temperature(°F)",
       y = "Dewpoint",
       title = "Temperature vs Dewpoint") +
  geom_smooth(method="lm")
```



- **Interpretation for plot 3.**

In this graph we observe the relationship between temperature and dewpoint. There is a strong linear relationship between temperature and dewpoint with small deviations. Simply we can conclude that the dewpoint increases as the temperature increases.

## Question 4

- In this part, we tried to fetch earthquake data from the Terremoti website, analyzed the data, categorized the earthquakes into different classes based on their magnitudes, and visualized it on a Leaflet map.
- **Calling necessary libraries.**

```
libraries <- c("ggplot2","sf","rworldmap","tidyverse","magrittr",
              "leaflet", "dplyr", "rvest", "xml2","rvest",
              "maps","mapdata","RgoogleMaps","lubridate","rnaturalearth","dplyr","rnaturalearthdata",
              "RColorBrewer","httr")
lapply(libraries, require, character.only = TRUE)
```

- **Scraping data from web.**

```
url <- "https://terremoti.ingv.it/en/events?starttime=2024-04-12%2B00%253A00%253A00&endtime=2024-04-19%2B23%253A59%253A59&last_nd=7&minmag=2&maxmag=10&mindepth=-10&maxdepth=1000&minlat=-90&maxlat=90&minlon=-180&maxlon=180&minversion=100&limit=30&orderby=mag-asc&lat=0&lon=0&maxradiuskm=-1&wheretype=area&box_search=Mondo&timezone=UTC&page=2"
res <- GET(url)
html_con <- content(res, "text")
```

- **Reading the html content.**

```
html_data <- read_html(html_con)
```

```
tables <- html_data %>%  
  html_nodes("table") %>%  
  html_table()
```

```
earthquake <- as.data.frame(tables)  
str(earthquake)
```

```
## 'data.frame':    30 obs. of  6 variables:  
## $ Origin.time..UTC.: chr  "2024-04-13 14:59:00" "2024-04-16 10:17:55" "2024-04-16 0  
2:11:32" "2024-04-14 04:20:06" ...  
## $ Magnitude       : chr  "ML 2.3" "Md 2.3" "Md 2.3" "ML 2.4" ...  
## $ Region          : chr  "3 km SW Arpino (FR)" "Campi Flegrei" "Campi Flegrei" "6 k  
m SE Mistretta (ME)" ...  
## $ Depth           : int   9 3 3 4 126 15 106 4 1 11 ...  
## $ Latitude        : num   41.6 40.8 40.8 37.9 38.5 ...  
## $ Longitude       : num   13.6 14.1 14.1 14.4 15.6 ...
```

- **Removing non-numeric characters from the Magnitude column and converting it to numeric.**

```
earthquake$Magnitude <- as.numeric(gsub("[^0-9.]", "", earthquake$Magnitude))
```

- **Adding a new column Class based on earthquake magnitudes, categorizing them as “Minor”, “Light”, and “Major”.**

```
earthquake <- earthquake %>%  
  mutate(Class = ifelse(Magnitude < 3, "Minor",  
                        ifelse(Magnitude <= 3.7, "Light", "Major")))
```

- **Defining colors for different earthquake classes.**

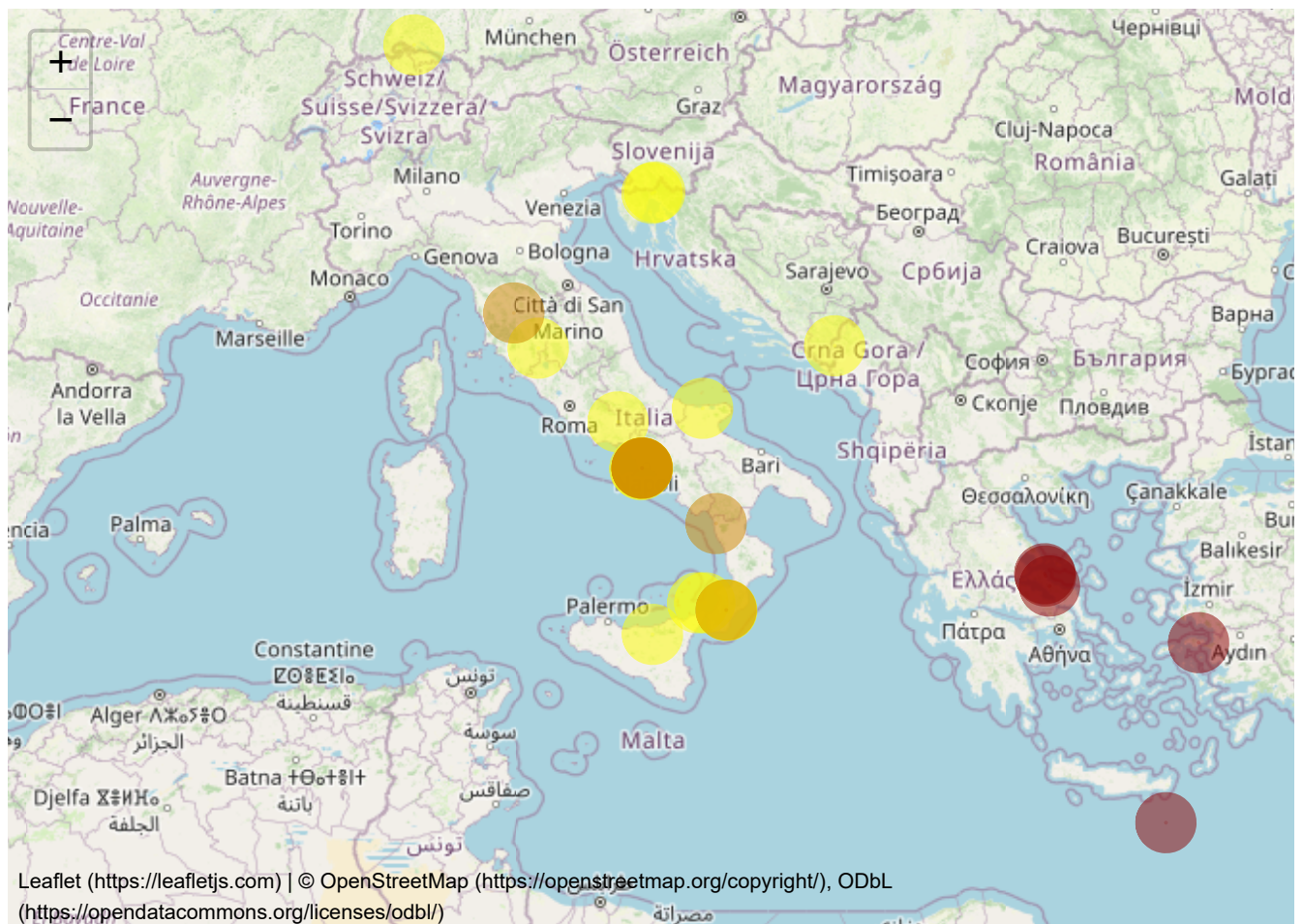
```
colors <- c("yellow", "orange3", "darkred")  
color_vector <- colorFactor(colors, levels = c("Minor", "Light", "Major"))
```

- **Creating Leaflet map. (Please click the circles to open popups.)**

```

leaflet() %>%
  addTiles() %>%
  addCircles(data = earthquake,
    ~Longitude, ~Latitude,
    weight = 30,
    radius = 100,
    popup = paste0(
      "<b>Date: </b>", earthquake$Origin.time..UTC.,
      "<br>",
      "<b>Place: </b>", earthquake$Region,
      "<br>",
      "<b>Depth in km: </b>", earthquake$Depth,
      "<br>",
      "<b>Magnitude: </b>", earthquake$Magnitude),
    label = ~Region,
    color = ~color_vector(Class)) %>%
  setView(lng = median(earthquake$Longitude),
    lat = median(earthquake$Latitude),
    zoom = 5)

```



## End Of The Work.

Note: Dear teacher, even though we used resources on the internet when we got stuck while answering the questions, our goal for each question was to learn, and it was a very educational assignment for us. Thank you, have a good day.