

Segment 8 Extra Credit

Britney Brown

2/8/2022

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

```
original <- read.csv("TV-Sales.csv", stringsAsFactors=FALSE)
df <- original

# change to numeric
change_to_int <- colnames(df)[-1]
for(i in 1:length(change_to_int)){
  df[,change_to_int[i]] <- as.numeric(df[,change_to_int[i]])}

# converting to datetime object
df[['Date']] <- as.POSIXct(df[['Date']],format = "%Y-%m-%d")

head(df)
```

```
##           Date S1 S2 S3 S4 S5 S6 S7 S8 S9 S10
## 1 2017-01-01 76 79 57 50 30 77 62 20 40 42
## 2 2017-01-02 64 91 42 59 46 64 61 25 35 51
## 3 2017-01-03 53 66 60 70 82 28 NA 34 49 69
## 4 2017-01-04 65 72 43 53 39 75 67 35 43 49
## 5 2017-01-05 58 69 43 71 42 64 52 28 29 43
## 6 2017-01-06 55 74 66 72 73 44 47 46 65 56
```

4.1 Which store had the highest mean sale in 2017?

```
df_dates <- df %>% dplyr::mutate(year = lubridate::year(Date),
                                month = lubridate::month(Date),
                                day = lubridate::day(Date))

df_2017 <- df_dates[which(df_dates$year == 2017),-c(1,12,13,14)]

means_2017 <- as.data.frame(colMeans(df_2017, na.rm = TRUE))
means_2017 <- cbind(store = rownames(means_2017), means_2017)
colnames(means_2017)[2] <- "means"
means_2017[which(means_2017$means == max(means_2017$means)),1]
```

```
## [1] "S2"
```

4.2 Which day showed the highest variance in sales across different stores?

```
library(matrixStats)
dailyvar <- rowVars(as.matrix(df[, -1]), na.rm = TRUE)
maxvar <- max(dailyvar)
df[which(dailyvar == maxvar), 'Date']
```

```
## [1] "2017-11-07 PST"
```

4.3 Which year showed the highest median sale for the store S5?

```
df$Year <- format(df$Date, format="%Y")
highest_mean <- df %>% select(S5, Year) %>% group_by(Year) %>%
  summarise(median = median(S5, na.rm = TRUE))

as.numeric(highest_mean[which(highest_mean$median == max(highest_mean$median)), 'Year'])
```

```
## [1] 2019
```

4.4 Which store recorded the highest number of sales for the largest number of days?

```
df2 <- original
df2[is.na(df2)] <- 0 # replace NA with 0
tbl <- table(colnames(df2[, -1])[max.col(df2[, -1], ties.method="last")])
rownames(as.data.frame(tbl[order(-tbl)][1]))
```

```
## [1] "S2"
```

4.5 Which store ranks 5th in the cumulative number of units sold over the 3-year interval?

```
stores <- as.data.frame(sapply(original[, -1], as.numeric))
cumSales <- colSums(stores, na.rm = TRUE)
rownames(as.data.frame(cumSales[order(-cumSales)][5]))
```

```
## [1] "S7"
```

4.6 Your program should create a file named `repaired.csv` in the directory which contains the same data as `TV-Sales.csv`, but with “N/A” values replaced with the median sale of that store, over the entire 3-year interval. Retain the header row found in `TV-Sales.csv`.

```
medians <- as.numeric(apply(stores, 2, median, na.rm = TRUE))

for(i in 1:10){
  stores[,i][is.na(stores[,i])] <- medians[i]
}

new_data <- cbind(Date = original[,1], stores)
write.csv(new_data, "repaired.csv", row.names = FALSE)
```

4.7 After imputing the missing values, plot the cumulative sale of each store as a function of time. Your plot should have a legend, a grid along the y-axis with a cell size of 10, labeled y-ticks at each multiple of 10, and x-ticks corresponding to the first days of each quarter (3 month period). Save the plot as `hw1/q4/plot.png`.