

# Report for R script

Monika Aulich-Olszewska

<sup>1</sup>Data used for this report comes from Kaggle

monika.aulich.olszewska@gmail.com

## 1. Introduction

This file is a report with proper plots, information and tables from Superstore\_data.Rmd file.

This report is made only for training purposes.

The script was made to show some of my skills in R and statistics.

Right now my goal is to find an interesting job as Data Analyst, Data Scientist and related.

## 2. Description

Here will be listed code, plots and some resume. Whole script with commentary is on the repository.

Repository can be found here: <https://github.com/senpai-ao/training>

Data that was used can be found here: <https://www.kaggle.com/datasets/vivek468/superstore-dataset-final>

```
1 library(ggplot2)
2
3 csv_file <- read.csv("Sample-Superstore.csv", header = TRUE)
4 df <- data.frame(csv_file)
5
6 #to check if data is loaded one can use head() function
7 head(df)
```

Next step is to choose interesting data and what to do with them. At first I want to see what product is the most bought

```
1 ggplot(df, aes(x=Product.ID, y=Row.ID)) +
2 geom_bar(stat = "identity", color="lightpink") +
3 labs(x="Product_ID", y="Number_of_products",
4       title="Bar_plot_with_numbers_of_sold_products"))
5
6 #to check if data is loaded one can use head() function
7 head(df)
```

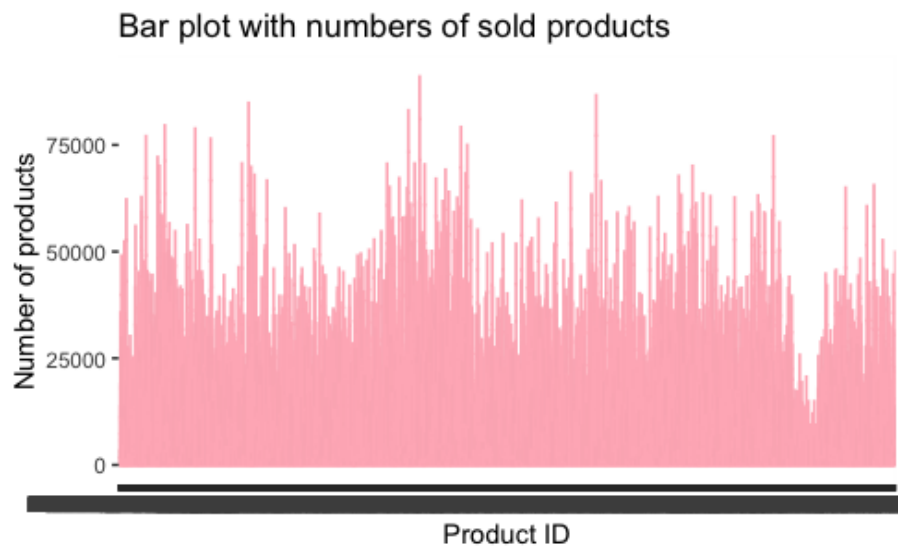


Figura 1. Plot 1

As we can see, this plot won't tell us much. There is too many different data in one category. Let's try wider category - Segment.

```
1 ggplot(df, aes(x=Segment, y=Row.ID)) +
2 geom_bar(stat = "identity", color="lightblue") +
3 labs(x="Cathegory_name", y="Number_of
4     sold_products_in_segment",
5     title="Bar_plot_with_numbers_of_sold_products")
```

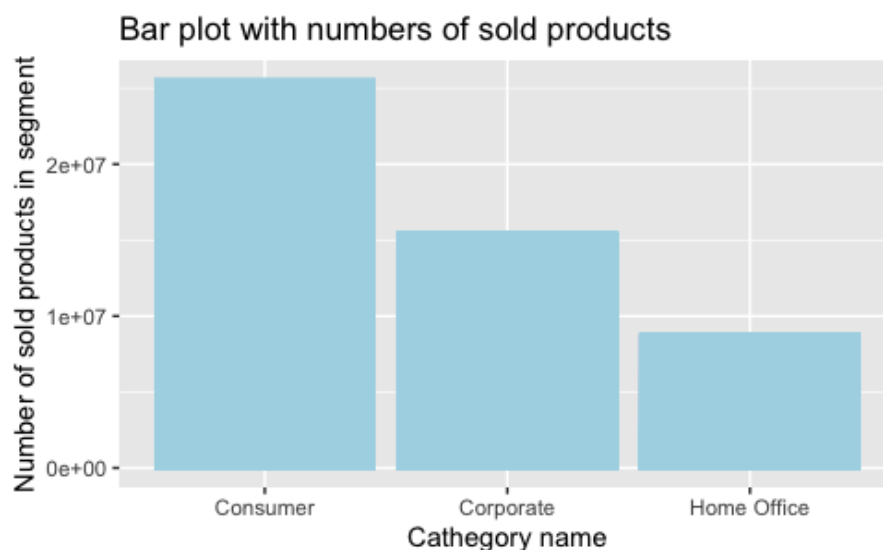


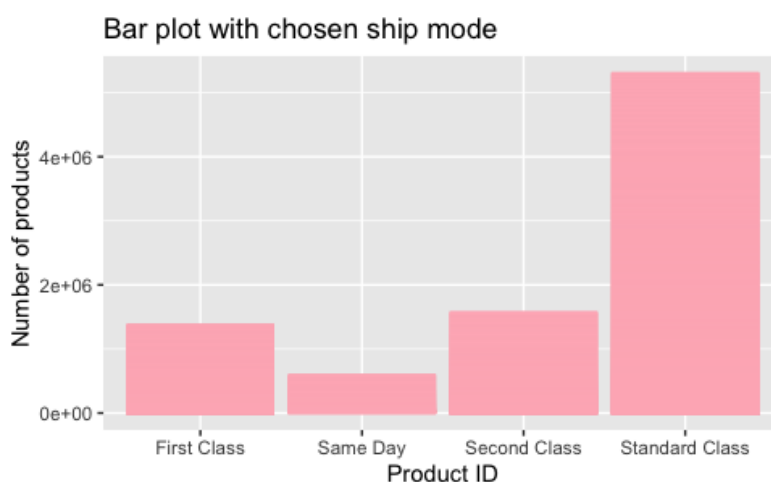
Figura 2. Plot 2

So right now one can see, that the biggest group in Superstore is Consumer group. Let's see what info we can get from Home Office data. First step here is to get data from Home office group only.

```

1 df_ho <- df[df$Segment=="Home_Office",]
2 head(df_ho)
3 print(paste("Number_of_sold_products_for_Home_Office
4     Group_is",nrow(df_ho)))
5
6 > "Number_of_sold_products_for_Home_Office_Group_is_1783"
7
8 ggplot(df_ho, aes(x=Ship.Mode, y=Row.ID)) +
9 geom_bar(stat = "identity", color="lightpink") +
10 labs(x="Product_ID",y="Number_of_products",
11       title="Bar_plot_with_chosen_ship_mode")

```



**Figura 3. Plot 3**

As we can see the most used ship mode is Standard Class. I would like to see if there is a correlation between used ship mode and product category. I used easy function cor().

```

1 df_corr <- data.frame(df_ho$Ship.Mode,df_ho$Category)
2 colnames(df_corr) <- c("Ship.Mode","Category")
3
4 #only numeric variables can be used to compute correlation
5 df_corr$Ship_num <- as.numeric(as.factor(df_corr$Ship.Mode))
6 df_corr$Categ_num <- as.numeric(as.factor(df_corr$Category))
7
8 cor_ship_cat <- cor(df_corr$Ship_num,df_corr$Categ_num)
9
10 print(paste("Correlation_is", round(cor_ship_cat,3)))

```

```
>"Correlation is -0.016"
```

The correlation is very close to 0, so one can say that there is no correlation between used

ship mode and product category. What interesting is for me, I would also like to see, how many orders were in what year in chosen group. At first I need to change data format.

```
1 df_ho$newdate <- strptime(as.character(df_ho$Order.Date),
2   "%m/%d/%Y")
3 df_ho$Date <- format(df_ho$newdate, "%Y-%m-%d")
4 splitted_data <- split(df_ho,
5   format(as.Date(df_ho$Date, format="%Y-%m-%d"),
6     format = "%Y"))
7 df_dates <- data.frame()
8 for(i in 1:length(splitted_data)){
9   data <- splitted_data[[i]]
10  year <- format(as.Date(data$Date[1],
11    format="%Y-%m-%d"), "%Y")
12  print(paste("In_year", year[1], "there_was",
13    nrow(data), "solds"))
14  line <- c(year, nrow(data))
15  df_dates <- rbind(df_dates, line)}
16 colnames(df_dates) <- c("Year", "Solds")
17 ggplot(data=df_dates, aes(x=Year, y=Solds, group=1)) +
18 geom_line(color="pink")+ geom_point() +
19 labs(x="Years", y="Number_of_sold_products",
20   title="Plot_with_number_of_sold_products_per_year")
```

```
>[1] "In year 2014 there was 312 solds"
```

```
>[1] "In year 2015 there was 341 solds"
```

```
>[1] "In year 2016 there was 466 solds"
```

```
>[1] "In year 2017 there was 664 solds"
```

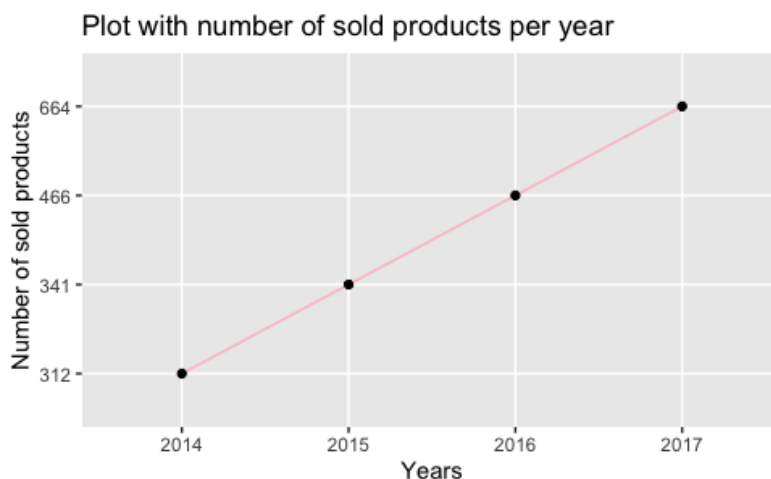


Figura 4. Plot 4