## R Studio

**Step 1:** Initial Exploratory Analysis



- Import the data into R by using the data frame function
- Install the package **'tidyverse'** (it helps to transform and better present data)
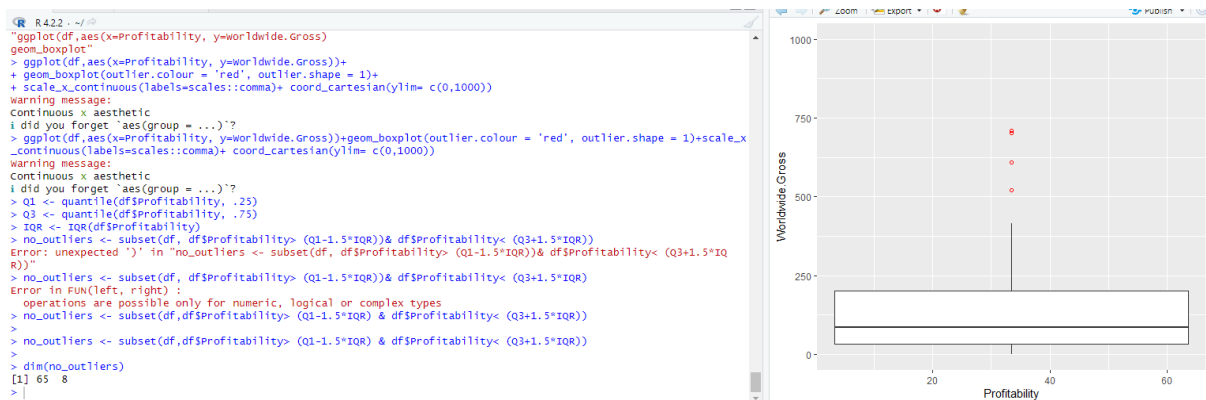


- Import library **'tidyverse'**
- Check the data type of each variable

**Step 2**: Clean Data

```
> colSums(is.na(df))
            Film              Genre      Lead.Studio Audience..score..    Profitability Rotten.Tomatoes..
               0                  0                0                 1                3                 1
   worldwide.Gross               Year
               0                  0
> df <- na.omit(df)
> colSums(is.na(df))
            Film              Genre      Lead.Studio Audience..score..    Profitability Rotten.Tomatoes..
               0                  0                0                 0                0                 0
   worldwide.Gross               Year
               0                  0
> dim(df[duplicated(df$Film,)])[1]
[1] 70
> df$Profitability <- round(df$Profitability, digit=2)
> dim (df)
[1] 70  8
> df$Worldwide.Gross <- round(df$worldwide.Gross, digit=2)
> dim (df)
[1] 70  8
>
```

- Using **'colSums(is.na)'** to count the NA value in each variable inside the data frame
- **'Na.omit'** is used to drop the missing value (NA)
- Using **'duplicated'** to check for duplicate
- **'Round'** is used to round off values
- **'Dim'** is used to get the dimensions of the data frame

**Step 2.1** Outlier Removal



- **Other outliers** are problematic and should be removed because they represent measurement errors, data entry or processing errors, or poor sampling
- The boxplot is shown on the right-hand side
- To remove outliers in 'Profitability', we first need to calculate the value of Q1(25%), Q3(75%) and IQR (Q3-Q1). Then, find
    1. Upper boundary (Anything above Q3 + 1.5 x IQR is an outlier)
    2. Lower boundary (Anything below Q1 - 1.5 x IQR is an outlier)

The value that is out of this range will be removed to increase the accuracy

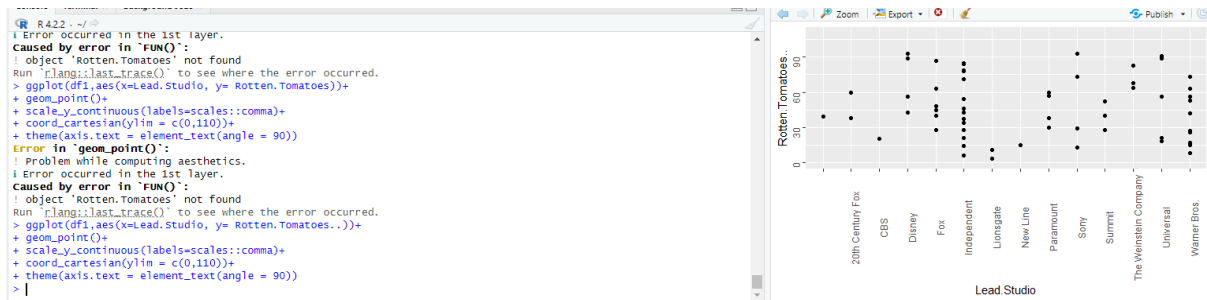The syntax of 'no_outliers' is getting the data in the range of the upper boundary and lower boundary

Dimension of 'no_outliers' data **65, 8**

```
> Q1 <- quantile(no_outliers$worldwide.Gross.25)
> Q3 <- quantile(no_outliers$worldwide.Gross, .75)
> Q1 <- quantile(no_outliers$worldwide.Gross, .25)
> IQR <- IQR(no_outliers$worldwide.Gross)
> df1 <- subset(no_outliers, no_outliers$worldwide.Gross>(Q1-1.5*IQR) & no_outliers$worldwide.Gross<(Q3+1.5*IQR))
> dim(df1)
[1] 61  8
```
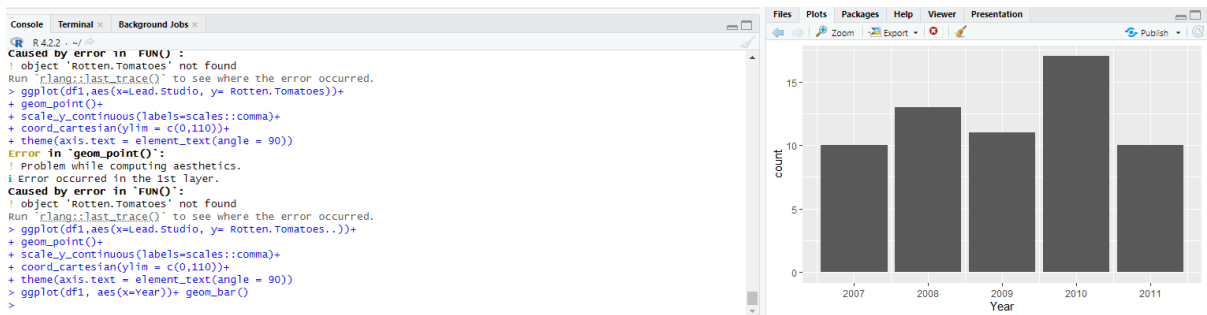
- Use the 'no_outliers' data to continue to remove outer outliers in 'Worldwide. Gross'

- The data frame dimension has now been reduced to **61, 8**
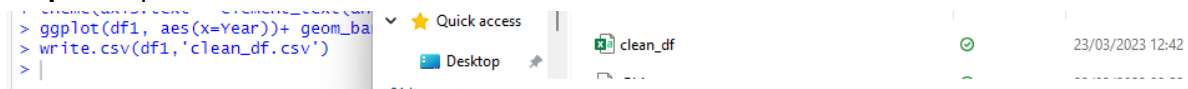
**Step 3**: Exploratory Data Analysis



- Scatter plot of the **df1**, showing the rotten tomatoes rating for every movie per studio
- According to rotten tomatoes, 'Independent' produced the highest number of movies and it also has a few movies rate above 60%. Whereas, overall 'Lionsgate' produce movies with the lowest rating
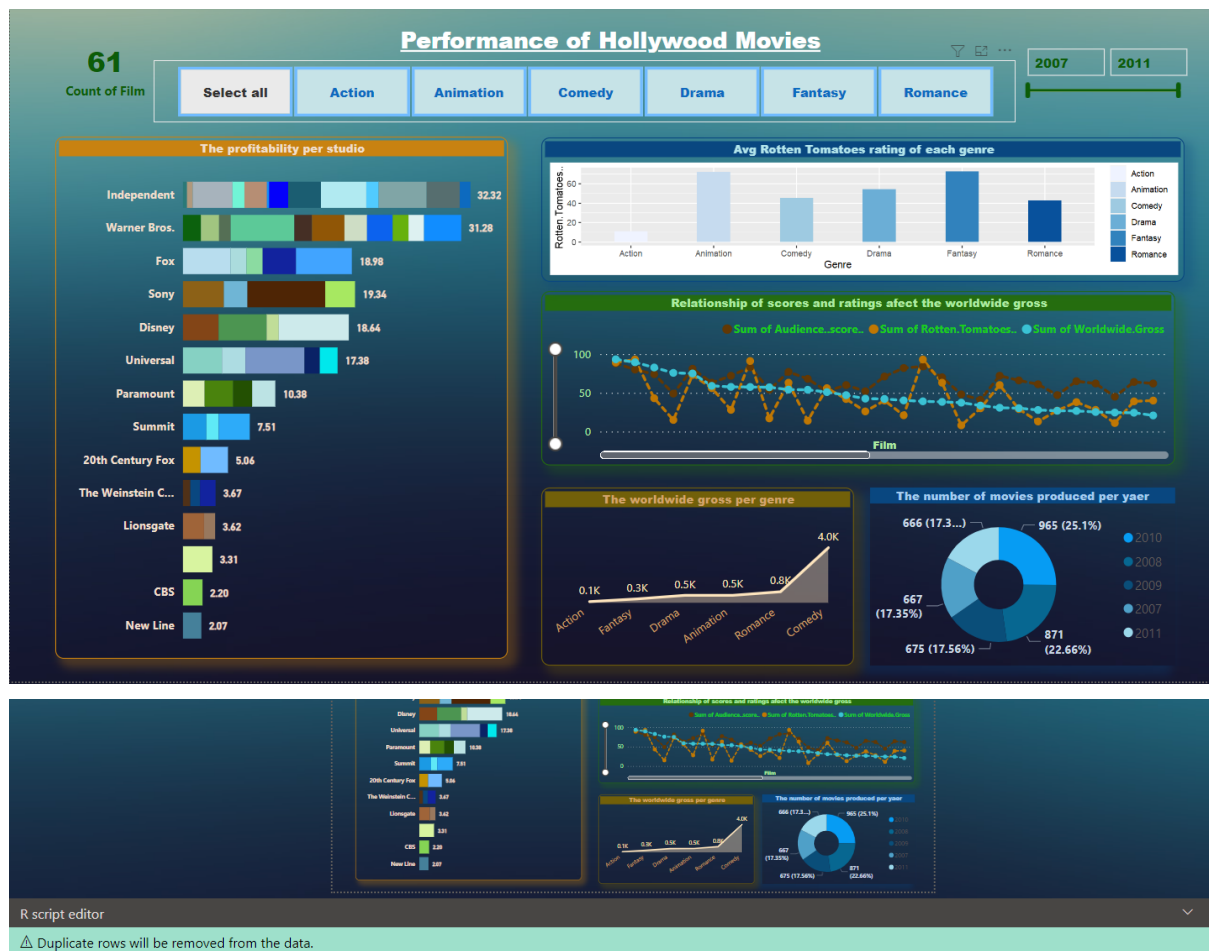


- Bar chart of the **df1**, count the year
- From this graph, we can tell that 2010 produce the most movies and a sharp decline in 2011

**Step 4**: Export Data



- **Write.csv** is used to export data

## Power BI



In this dashboard, I used a range of visualisation and embedded R script into the Power BI to show the performance of Hollywood Movies. Also, I mainly use the colour blue, brown and green to meet the client's criteria

Type of visualisations used:
- Use **'card'** to display the number of films in the dataset
- **'Slicer'** is used to show the 6 types of genres and years (2007-2011)
- **'Stacked bar chart'** shows the profitability per studio. Different types of films show in a different colour on each bar
- Using the **'ggplot'** in R script to create a bar chart for finding the average rotten tomatoes rating of each genre. Gradient blue is used to show the type of genre
- **'Line chart'** is used to see how can movies ratings and scores affect the worldwide gross
- **'Area chart'** suggest the worldwide gross per genre
- **'Donut chart'** shows the number of movies produced per year