

Assignment 3: DATA ANALYSIS

R AND POWERBI
NOORULAIN FAHAD

Table of Contents

INTRODUCTION	2
TASK – 1 : DATA ANALYSIS IN R	2
1. IMPORT DATASET	2
2. VIEWING THE DATA FRAME	2
3. INSTALLING TIDYVERSE	3
4. CHECK DATA TYPES.....	3
5. CHECKING DIMESIONS OF DATA FRAME	4
6. CHECK FOR MISSING VALUES.....	4
7. REMOVE MISSING VALUES.....	5
8. CHECK DUPLICATES.....	7
9. ROUND-OFF VALUES.....	7
10. CHECK FOR OUTLIERS ON BOX PLOT	8
11. SUMMARY STATISTICS (UNIVARIATE ANALYSIS)	11
12. BIVARIATE ANALYSIS.....	12
a. SCATTER PLOT.....	12
b. BAR CHART.....	13
13. EXPORT CLEANED DATA	14
TASK – 2 : DASHBOARD IN POWERBI	14
1. IMPORT DATA.....	14
2. BUILD VISUALS.....	18
3. BUILD DASHBOARD	24
REFLECTIVE	25

INTRODUCTION

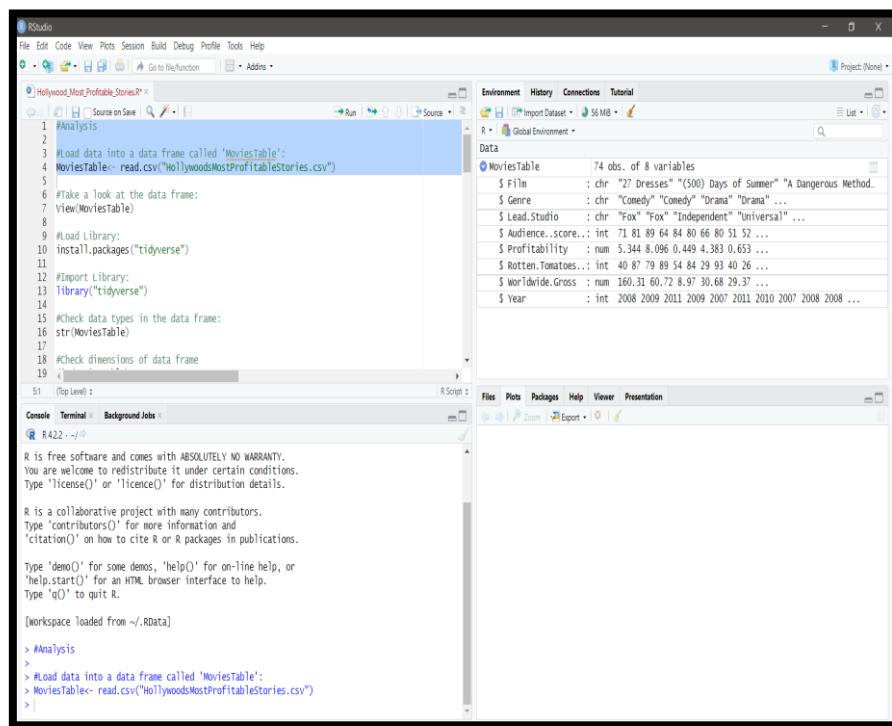
The aim of this project is to analyse the performance of Hollywood movies. The dataset used for this project is 'Hollywood's Most Profitable Stories' which contains the Title, Genre, Lead Studio, Profitability and Rotten Tomatoes Ratings for movies released between 2007 and 2012 in CSV format. The source of this dataset is given below :

<https://public.tableau.com/app/sample-data/HollywoodsMostProfitableStories.csv>

TASK – 1 : DATA ANALYSIS IN R

To perform data analysis in R, the workflow consists of importing the data, removing missing values, removing duplicates, removing outliers, feature engineering (creating new column using existing data), summary statistics, EDA(Exploratory Data Analysis) and export the cleaned data.

1. **IMPORT DATASET** To import the dataset in R studio, "read.csv()" command is used. It reads the file and stores it into a data frame. In current scenario, the dataset is imported and stored into a data frame called "MoviesTable".



The screenshot shows the RStudio interface. The code editor pane displays the following R script:

```
#Analysis
#Load data into a data frame called 'MoviesTable'
MoviesTable<- read.csv("HollywoodsMostProfitableStories.csv")
#Take a look at the data frame:
View(MoviesTable)
#Load Library:
install.packages("tidyverse")
#Import Library:
library("tidyverse")
#Check data types in the data frame:
str(MoviesTable)
#Check dimensions of data frame

```

The environment pane shows the following data frame structure:

	MoviesTable	74 obs. of 8 variables
\$ Film	: chr	"27 Dresses" "(500) Days of Summer" "A Dangerous Method..."
\$ Genre	: chr	"Comedy" "Comedy" "Drama" "Drama" ...
\$ Lead.Studio	: chr	"Fox" "Fox" "Independent" "Universal" ...
\$ Audience..score.:	: int	71 81 89 64 84 80 66 80 51 52 ...
\$ Profitability	: num	5.344 8.096 0.449 4.383 0.653 ...
\$ Rotten.Tomatoes.:	: int	40 87 76 89 54 81 29 93 40 26 ...
\$ Worldwide.Gross	: num	160.31 60.72 8.97 30.68 29.37 ...
\$ Year	: int	2008 2009 2011 2009 2007 2013 2010 2007 2008 2008 ...

2. **VIEWING THE DATA FRAME** To view the data frame, "View(MoviesTable)" command is used. To run a command, you can either press "alt + Enter" or simply click the "Run" on top pane of the editor. The View command generated the data frame with column

headers Film, Genre, Lead Studio, Audience score, Profitability, Rotten Tomatoes, Worldwide Gross and Year as shown below.

	Film	Genre	Lead.Studio	Audience.score..	Profitability	Rotten.Tomatoes..	Worldwide.Gross
1	27 Dresses	Comedy	Fox	71	5.3436218	40	
2	(500) Days of Summer	Comedy	Fox	81	8.0960000	87	
3	A Dangerous Method	Drama	Independent	89	0.4486447	79	
4	A Serious Man	Drama	Universal	64	4.3828571	89	
5	Across the Universe	Romance	Independent	84	0.6526032	54	
6	Beginners	Comedy	Independent	80	4.4716750	84	
7	Dear John	Drama	Sony	66	4.5988000	29	
8	Enchanted	Comedy	Disney	80	4.0057371	93	
9	Fireproof	Drama	Independent	51	66.9340000	40	
10	Four Christmases	Comedy	Warner Bros.	52	2.0228250	26	
11	Ghosts of Girlfriends Past	Comedy	Warner Bros.	47	2.0444000	27	
12	Gnomeo and Juliet	Animation	Disney	52	5.3879722	56	
13	Going the Distance	Comedy	Warner Bros.	56	1.3140625	53	
14	Good Luck Chuck	Comedy	Lionsgate	61	2.3676851	3	
15	He's Just Not That Into You	Comedy	Warner Bros.	60	7.1536000	42	
16	High School Musical 3: Senior Year	Comedy	Disney	76	22.9131365	65	
17	I Love You Phillip Morris	Comedy	Independent	57	1.3400000	71	
18	It's Complicated	Comedy	Universal	63	2.6423529	56	
19	Jane Eyre	Romance	Universal	77	NA	85	
20	Just Wright	Comedy	Fox	58	1.7974167	45	
21	Killers	Action	Lionsgate	45	1.2453333	11	
22	Knocked Up	Comedy	Universal	83	6.6364018	91	
23	Leap Year	Comedy	Universal	49	1.7152632	21	
24	Letters to Juliet	Comedy	Summit	62	2.3993333	40	
25	License to Wed	Comedy	Warner Bros.	55	1.9802064	8	

3. INSTALLING TIDYVERSE R has integrated packages which are designed to work together to make common data science operations more user-friendly. The "Tidyverse" suite of integrated packages have functions used for data wrangling, tidying, reading/writing, parsing, and visualizing. The commands "install.pacakages(tidyverse)" and "library(tidyverse)" are used to install this package in R studio.

4. CHECK DATA TYPES Checking the data type of all the fields in the dataset is a foundation step in performing data analysis. It aligns with the principles of data validation and integrity. To perform this in R, "str(MoviesTable)" command is used.

```

13 library("tidyverse")
14
15 #Check data types in the data frame:
16 str(MoviesTable)
17
18 #Check dimensions of data frame
19 dim(MoviesTable)
20
21 #Check for missing values
22 colSums(is.na(MoviesTable))
23
24 #Remove missing values in rows
25 movies <- na.omit(MoviesTable)
26
27 #Check for missing values
28 colSums(is.na(movies))
29
30 #Check dimensions of data frame
31
31 <-- Conflicts
32 dplyr::filter() masks stats::filter()
33 dplyr::lag() masks stats::lag()
34 > #Check data types in the data frame:
35 > str(MoviesTable)
36 'data.frame': 74 obs. of  8 variables:
37   $ Film      : chr "27 Dresses" "(500) Days of Summer" "A Dangerous Method" "A Serious
38   Man" ...
39   $ Genre     : chr "Comedy" "Comedy" "Drama" "Drama" ...
40   $ Lead.Studio: chr "Fox" "Fox" "Independent" "Universal" ...
41   $ Audience.score.: int 78 81 80 64 84 80 66 80 51 52 ...
42   $ Profitability: num 5.344 8.096 0.449 4.383 0.653 ...
43   $ Rotten.Tomatoes.: int 40 87 79 89 54 84 29 93 40 26 ...
44   $ Worldwide.Gross: num 160.31 60.72 8.97 30.68 29.37 ...
45   $ Year       : int 2008 2009 2011 2009 2007 2011 2010 2007 2008 2008 ...
46
47 Conflicts:
48 - dplyr::filter() masks stats::filter()
49 - dplyr::lag() masks stats::lag()
50 > #Check data types in the data frame:
51 > str(MoviesTable)
52 'data.frame': 74 obs. of  8 variables:
53   $ Film      : chr "27 Dresses" "(500) Days of Summer" "A Dangerous Method" "A Serious
54   Man" ...
55   $ Genre     : chr "Comedy" "Comedy" "Drama" "Drama" ...
56   $ Lead.Studio: chr "Fox" "Fox" "Independent" "Universal" ...
57   $ Audience.score.: int 78 81 80 64 84 80 66 80 51 52 ...
58   $ Profitability: num 5.344 8.096 0.449 4.383 0.653 ...
59   $ Rotten.Tomatoes.: int 40 87 79 89 54 84 29 93 40 26 ...
60   $ Worldwide.Gross: num 160.31 60.72 8.97 30.68 29.37 ...
61   $ Year       : int 2008 2009 2011 2009 2007 2011 2010 2007 2008 2008 ...
62
63 tidyverse_conflicts()
64
```

- 5. CHECKING DIMENSIONS OF DATA FRAME** To check the dimensions of the data frame "MoviesTable", the command "dim(MoviesTable)" is used. It shows that our data frame has 74 rows and 8 columns.

```

library("tidyverse")
#Check data types in the data frame:
str(MoviesTable)
#Check dimensions of data frame
dim(MoviesTable)

#Check for missing values
colSums(is.na(MoviesTable))
#Remove missing values in rows
movies <- na.omit(MoviesTable)
#Check for missing values
colSums(is.na(movies))
#Check dimensions of data frame
dim(MoviesTable)

```

```

'data.frame': 74 obs. of 8 variables:
$ Film           : chr "27 Dresses" "(500) Days of Summer" "A Dangerous Method" "A Serious Man" ...
$ Genre          : chr "Comedy" "Comedy" "Drama" "Drama" ...
$ Lead.Studio    : chr "Fox" "Fox" "Independent" "Universal" ...
$ Audience..score.: int 71 81 89 64 84 80 66 80 51 52 ...
$ Profitability  : num 5.344 8.096 0.449 4.383 0.653 ...
$ Rotten.Tomatoes.: int 40 87 79 89 54 84 29 93 40 26 ...
$ Worldwide.Gross : num 160.31 60.72 8.97 30.68 29.37 ...
$ Year           : int 2008 2009 2011 2009 2007 2011 2010 2007 2008 2008 ...
> #Check dimensions of data frame
> dim(MoviesTable)
[1] 74 8

```

- 6. CHECK FOR MISSING VALUES** The next step is to check for missing values in the dataset before performing the analysis. To check for missing values, "colSums(is.na(MoviesTable))" command is used. After running this command, you get the number of missing values in each column of the table. There are 1, 3 and 1 missing values in columns Audience score, Profitability and Rotten Tomatoes columns respectively, as shown below.

```

library("tidyverse")
#Check data types in the data frame:
str(MoviesTable)
#Check dimensions of data frame
dim(MoviesTable)

#Check for missing values
colSums(is.na(MoviesTable))

#Remove missing values in rows
movies <- na.omit(MoviesTable)
#Check for missing values
colSums(is.na(movies))
#Check dimensions of data frame
dim(MoviesTable)

```

```

'data.frame': 74 obs. of 8 variables:
$ Film           : chr "27 Dresses" "(500) Days of Summer" "A Dangerous Method" "A Serious Man" ...
$ Genre          : chr "Comedy" "Comedy" "Drama" "Drama" ...
$ Lead.Studio    : chr "Fox" "Fox" "Independent" "Universal" ...
$ Audience..score.: int 71 81 89 64 84 80 66 80 51 52 ...
$ Profitability  : num 5.344 8.096 0.449 4.383 0.653 ...
$ Rotten.Tomatoes.: int 40 87 79 89 54 84 29 93 40 26 ...
$ Worldwide.Gross : num 160.31 60.72 8.97 30.68 29.37 ...
$ Year           : int 2008 2009 2011 2009 2007 2011 2010 2007 2008 2008 ...
> #Check dimensions of data frame
> dim(MoviesTable)
[1] 74 8
> #Check for missing values
> colSums(is.na(MoviesTable))
Film          Genre        Lead.Studio Audience..score.. Profitability
0             0            0                 1                  3
Rotten.Tomatoes.. Worldwide.Gross Year
1             0            0            0

```

7. REMOVE MISSING VALUES The next step is to remove the missing values in the rows of the data frame `MoviesTable` and assign the name "Movies" to the new data frame. The command "`na.omit()`" is used.

```

File Edit Code View Plots Session Build Debug Profile Tools Help
File Edit Code View Plots Session Build Debug Profile Tools Help
Hollywood_Most_Profitable_Stories.Rx MoviesTable x
Source On Save Run Source
13 library("tidyverse")
14
15 #Check data types in the data frame:
16 str(MoviesTable)
17
18 #Check dimensions of data frame
19 dim(MoviesTable)
20
21 #Check for missing values
22 colSums(is.na(MoviesTable))
23
24 #Remove missing values in rows of MoviesTable and assign it a new data frame called 'movies'
25 movies <- na.omit(MoviesTable)
26
27 #Check for missing values
28 colSums(is.na(movies))
29
30 #Check dimensions of data frame
31 dim(movies)
32
33 #Check for missing values again in the new data frame "Movies" and there
34 #are no missing values as shown below.
35
36 #Check for duplicates
37
38 dim(movies[duplicated(movies$Film),]) [1]
39
(R 4.2.2 - ~/)
  Film   Genre   Lead.Studio Audience..score.. Profitability
  0       0       0           1               3
Rotten.Tomatoes.. Worldwide.Gross      Year
  1       0       0           0               0
> #Remove missing values in rows of MoviesTable and assign it a new data frame called 'movies'
> movies <- na.omit(MoviesTable)

```

After this, we check for missing values again in the new data frame "Movies" and there are no missing values as shown below.

```

File Edit Code View Plots Session Build Debug Profile Tools Help
File Edit Code View Plots Session Build Debug Profile Tools Help
Hollywood_Most_Profitable_Stories.Rx MoviesTable x
Source On Save Run Source
21 #Check for missing values
22 colSums(is.na(MoviesTable))
23
24 #Remove missing values in rows of MoviesTable and assign it a new data frame called 'movies'
25 movies <- na.omit(MoviesTable)
26
27 #Check for missing values
28 colSums(is.na(movies))
29
30 #Check dimensions of new data frame 'movies':
31 dim(movies)
32
33 #Take a look at the data frame 'movies':
34 View(movies)
35
36 #Check for duplicates
37
38 dim(movies[duplicated(movies$Film),]) [1]
39
(R 4.2.2 - ~/)
  Worldwide.Gross      Year
  0                   2008 2009 2011 2009 2007 2011 2010 2007 2008 ...
> #Check dimensions of data frame
> dim(MoviesTable)
[1] 74 8
> #Check for missing values
> colSums(is.na(MoviesTable))
  Film   Genre   Lead.Studio Audience..score.. Profitability
  0       0       0           1               3
Rotten.Tomatoes.. Worldwide.Gross      Year
  1       0       0           0               0
> #Remove missing values in rows of MoviesTable and assign it a new data frame called 'movies'
> movies <- na.omit(MoviesTable)
> #Check for missing values
> colSums(is.na(movies))
  Film   Genre   Lead.Studio Audience..score.. Profitability
  0       0       0           0               0
Rotten.Tomatoes.. Worldwide.Gross      Year
  0       0       0           0               0
>

```

The dimensions of the new data frame "movies" are now 70 rows and 8 columns as shown below.

A screenshot of the RStudio interface. In the top right, the 'Environment' tab is selected. Under the 'Data' section, there are two entries: 'movies' (70 obs. of 8 variables) and 'MoviesTable' (74 obs. of 8 variables). The main workspace shows a partial R script:

```

data frame called 'movies'
8 ...
Profitability
3
ame called 'movies'
Profitability
0

```

The new data frame movies with no missing values is shown as follows.

A screenshot of the RStudio interface showing the 'movies' data frame. The table has 70 rows and 8 columns, with the following headers: Film, Genre, Lead.Studio, Audience..score., Profitability, Rotten.Tomatoes., Worldwide.Gross, and Year. The table lists various movies from 1 to 26, along with their respective details. The bottom status bar indicates 'Showing 1 to 26 of 70 entries, 8 total columns'.

	Film	Genre	Lead.Studio	Audience..score..	Profitability	Rotten.Tomatoes..	Worldwide.Gross	Year
1	27 Dresses	Comedy	Fox	71	5.3436218	40	160.308654	
2	(500) Days of Summer	Comedy	Fox	81	8.0960000	87	60.720000	
3	A Dangerous Method	Drama	Independent	89	0.4486447	79	8.972895	
4	A Serious Man	Drama	Universal	64	4.3828571	89	30.680000	
5	Across the Universe	Romance	Independent	84	0.6526032	54	29.367143	
6	Beginners	Comedy	Independent	80	4.4718750	84	14.310000	
7	Dear John	Drama	Sony	66	4.5988000	29	114.970000	
8	Enchanted	Comedy	Disney	80	4.0057371	93	340.487652	
9	Fireproof	Drama	Independent	51	66.9340000	40	33.467000	
10	Four Christmases	Comedy	Warner Bros.	52	2.0229250	26	161.834000	
11	Ghosts of Girlfriends Past	Comedy	Warner Bros.	47	2.0444000	27	102.220000	
12	Gnomeo and Juliet	Animation	Disney	52	5.3879722	56	193.967000	
13	Going the Distance	Comedy	Warner Bros.	56	1.3140625	53	42.050000	
14	Good Luck Chuck	Comedy	Lionsgate	61	2.3676851	3	59.192128	
15	He's Just Not That Into You	Comedy	Warner Bros.	60	7.1536000	42	178.840000	
16	High School Musical 3: Senior Year	Comedy	Disney	76	22.9131365	65	252.044501	
17	I Love You Phillip Morris	Comedy	Independent	57	1.3400000	71	20.100000	
18	It's Complicated	Comedy	Universal	63	2.6423529	56	224.600000	
20	Just Wright	Comedy	Fox	58	1.7974167	45	21.569000	
21	Killers	Action	Lionsgate	45	1.2453333	11	93.400000	
22	Knocked Up	Comedy	Universal	83	6.6364018	91	219.001261	
23	Leap Year	Comedy	Universal	49	1.7152632	21	32.590000	
24	Letters to Juliet	Comedy	Summit	62	2.6393333	40	79.180000	
25	License to Wed	Comedy	Warner Bros.	55	1.9802064	8	69.307224	
26	Life as We Know It	Comedy	Independent	62	2.5305263	28	96.160000	

8. CHECK DUPLICATES The next step is to check for duplicate records in the data frame. The command used is shown below and you can see there are no duplicate records in the dataset.

```

File Edit Code View Plots Session Build Debug Profile Tools Help
Hollywood_Most_Profitable_Stories.R movies MoviesTable
Source | Run | Source
32
33 #Take a look at the data frame 'movies':
34 view(movies)
35
36 #Check for duplicates
37
38 dim(movies[duplicated(movies$Film),])[1]
39
40 #round off values to 2 places
41
42 movies$Profitability <- round(movies$Profitability ,digit=2)
43
44 movies$Worldwide.Gross <- round(movies$Worldwide.Gross ,digit=2)
45
46
47 head(movies)
48
49 view(movies)
50

```

R Script

```

R422 - ~/Desktop/Hollywood_Most_Profitable_Stories.R
R422 - ~/Desktop/Hollywood_Most_Profitable_Stories.R
1 0 0 0 0 0
> #Remove missing values in rows of MoviesTable and assign it a new data frame called 'movies'
> movies <- na.omit(MoviesTable)
> #Check for missing values
> colSums(is.na(movies))
  Film   Genre Lead.Studio Audience..score.. Profitability 
    0      0       0          0            0 
Rotten.Tomatoes.. Worldwide.Gross Year
  0      0       0
> #Check dimensions of new data frame 'movies':
> dim(movies)
[1] 70 8
> #Take a look at the data frame 'movies':
> view(movies)
> #Check for duplicates
>
> dim(movies[duplicated(movies$Film),])[1]
[1] 0
>

```

Environment History Connections Tutorial

Data

- movies 70 obs. of 8 variables
- MoviesTable 74 obs. of 8 variables

Files Plots Packages Help Viewer Presentation

9. ROUND-OFF VALUES The next step is to round off the values of variable Profitability and Worldwide Gross to two decimal places to make the numerical data consistent.

```

File Edit Code View Plots Session Build Debug Profile Tools Help
Hollywood_Most_Profitable_Stories.R movies MoviesTable
Source | Run | Source
32
33 #Take a look at the data frame 'movies':
34 view(movies)
35
36 #Check for duplicates
37
38 dim(movies[duplicated(movies$Film),])[1]
39
40 #round off Profitability values to 2 places
41
42 movies$Profitability <- round(movies$Profitability ,digit=2)
43
44 #round off Worldwide Gross values to 2 places
45
46 movies$Worldwide.Gross <- round(movies$Worldwide.Gross ,digit=2)
47
48
49 head(movies)
50

```

R Script

```

R422 - ~/Desktop/Hollywood_Most_Profitable_Stories.R
R422 - ~/Desktop/Hollywood_Most_Profitable_Stories.R
1 0 0 0 0 0
> #Check dimensions of new data frame 'movies':
> dim(movies)
[1] 70 8
> #Take a look at the data frame 'movies':
> view(movies)
> #Check for duplicates
>
> dim(movies[duplicated(movies$Film),])[1]
[1] 0
> #round off Profitability values to 2 places
>
> movies$Profitability <- round(movies$Profitability ,digit=2)
> #round off Worldwide Gross values to 2 places
>
> movies$Worldwide.Gross <- round(movies$Worldwide.Gross ,digit=2)
>

```

Environment History Connections Tutorial

Data

- movies 70 obs. of 8 variables
 - \$ Film : chr "27 Dresses" "(500) Days"
 - \$ Genre : chr "Comedy" "Comedy" "Drama"
 - \$ Lead.Studio : chr "Fox" "Fox" "Independent"
 - \$ Audience..score.. : int 71 81 89 64 84 80 66 80
 - \$ Profitability : num 5.34 8.1 0.45 4.38 0.65
 - \$ Rotten.Tomatoes. : int 40 87 79 89 54 84 29 93
 - \$ Worldwide.Gross : num 160.31 60.72 8.97 30.68
 - \$ Year : int 2008 2009 2011 2009 2009 2000 2001 2010
 - attr(*, "na.action")= 'omit' Named int [1:4] 19
 - ..- attr(*, "names")= chr [1:4] "19" "42" "51" "
- MoviesTable 74 obs. of 8 variables

Files Plots Packages Help Viewer Presentation

You can now see that the Profitability and Worldwide gross columns both have values with two decimal places.

The screenshot shows the RStudio interface. In the top menu, 'File', 'Edit', 'Code', 'View', 'Plots', 'Session', 'Build', 'Debug', 'Profile', 'Tools', and 'Help' are visible. Below the menu, there are tabs for 'Hollywood_Most_Profitable_Stories.R', 'movies', and 'MoviesTable'. The 'Environment' tab is selected in the top right. The 'Data' pane shows the 'movies' dataset with 70 observations and 8 variables. The 'Console' tab at the bottom has some initial R code and output.

10. CHECK FOR OUTLIERS ON BOX PLOT To check for outliers (a data point that lie away from the other points of the dataset), we first need to make a box plot using the following command. A box plot in R, also known as box and whisker plot, is a graphical representation which allows you to summarize the main characteristics of the data (position, dispersion, skewness etc) and identify the presence of outliers.

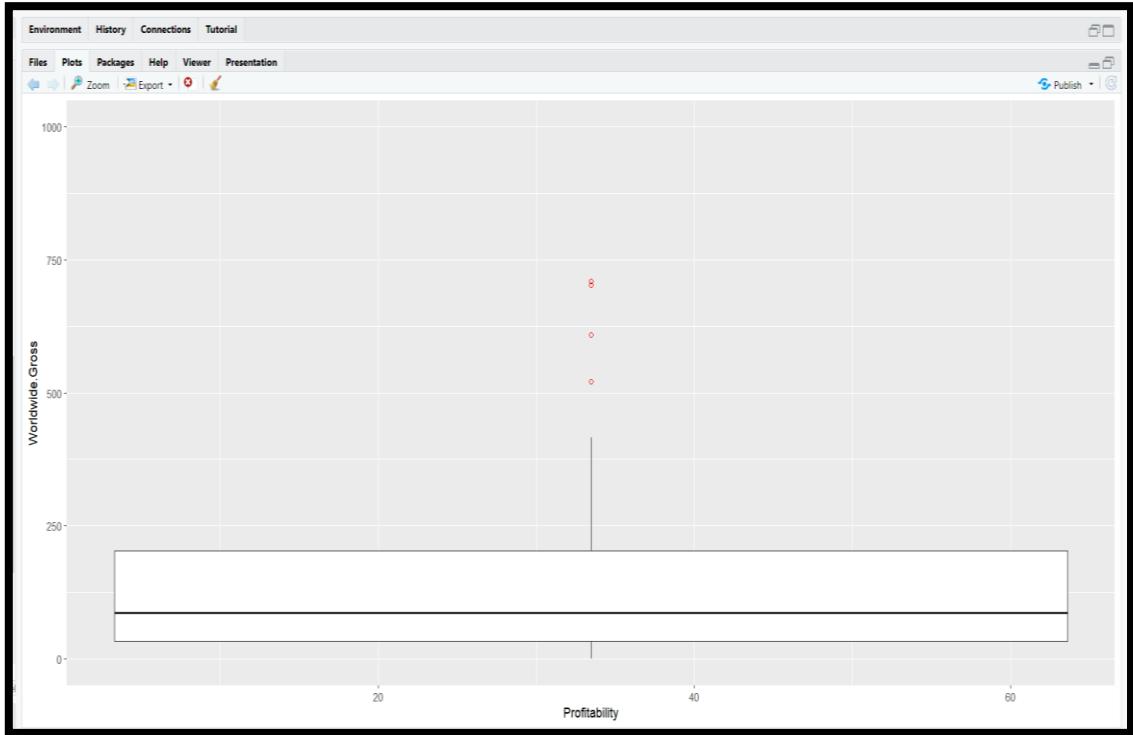
The screenshot shows the RStudio interface with the 'R Script' tab active. The code in the console includes:

```

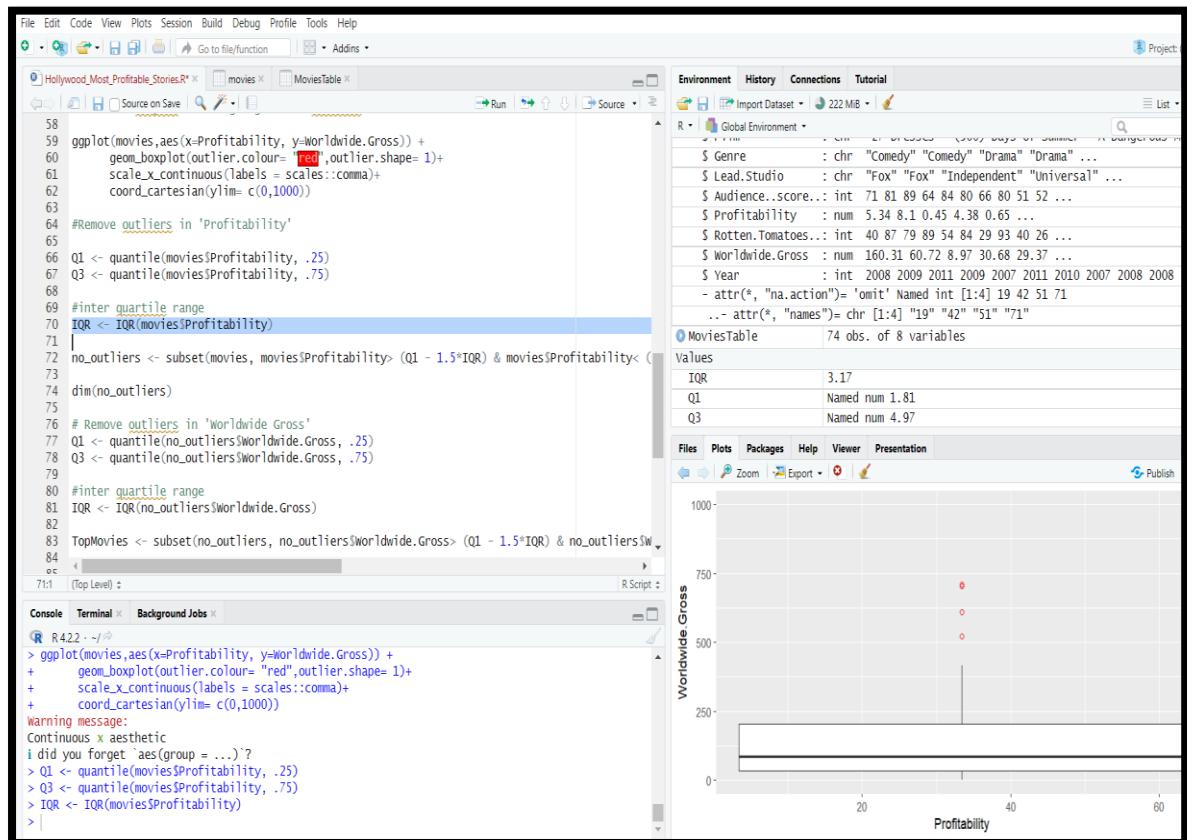
52 library("ggplot2")
53 
54 #Check for outliers on box plot
55 
56 #Create a boxplot that highlights the outliers
57 
58 ggplot(movies,aes(x=Profitability, y=Worldwide.Gross)) +
59   geom_boxplot(outlier.colour = "red",outlier.shape = 1) +
60   scale_x_continuous(labels = scales::comma) +
61   coord_cartesian(ylim= c(0,1000))
62 
```

The 'Plots' pane displays a box plot of Worldwide Gross vs Profitability. The x-axis is Profitability (ranging from 0 to 60+) and the y-axis is Worldwide Gross (ranging from 0 to 1000). Three red outlier points are visible above the upper whisker.

The following images shows the box plot for the movies data frame.



The next step is to remove the outliers in the "Profitability" variable and assign the subset of movies data frame with no outliers to a new data frame called "no_outliers".



The dimensions of new data frame "no_outliers" after the removal of outliers from the profitability variable are now 65 rows and 8 columns.

```

76 #Check dimension of no_outliers data frame
77 dim(no_outliers)
78
79 # Remove outliers in 'Worldwide.Gross'
80 Q1 <- quantile(no_outliers$Worldwide.Gross, .25)
81 Q3 <- quantile(no_outliers$Worldwide.Gross, .75)
82
83 #inter quartile range
84 IQR <- IQR(no_outliers$Worldwide.Gross)
85
86 TopMovies <- subset(no_outliers, no_outliers$Worldwide.Gross > (Q1 - 1.5*IQR) & no_outliers$W
87
88 dim(TopMovies)
89
90 #Summary Statistics / Univariate Analysis:
91
92
93

```

(Top Level) R Script

Console Terminal × Background Jobs ×

R 4.2.2 · ~/

```

i did you forget `aes(group = ...)`?
> Q1 <- quantile(movies$profitability, .25)
> Q3 <- quantile(movies$profitability, .75)
> IQR <- IQR(movies$Profitability)
> no_outliers <- subset(movies, movies$Profitability > (Q1 - 1.5*IQR) &
+           movies$Profitability < (Q3 + 1.5*IQR))
> #Check dimension of no_outliers data frame
>
> dim(no_outliers)
[1] 65 8
> 
```

Repeat the same procedure to remove the outliers from the Worldwide Gross variable and assign the subset of data frame no_outliers to a new data frame "TopMovies".

```

80 # Remove outliers in 'Worldwide.Gross'
81 Q1 <- quantile(no_outliers$Worldwide.Gross, .25)
82 Q3 <- quantile(no_outliers$Worldwide.Gross, .75)
83
84 #inter quartile range
85 IQR <- IQR(no_outliers$Worldwide.Gross)
86
87 #Remove outliers in the 'no_outliers' data frame and assign the new subset to 'TopMovies'
88
89 TopMovies <- subset(no_outliers, no_outliers$Worldwide.Gross > (Q1 - 1.5*IQR) &
90                   no_outliers$Worldwide.Gross < (Q3 + 1.5*IQR))
91
92 #Check dimension of TopMovies data frame
93
94 dim(TopMovies)
95
96 #Summary Statistics / Univariate Analysis:
97 summary(TopMovies)
98
99 #Bivariate analysis:
100
101 #1- Scatter plot
102 
```

(Top Level) R Script

Console Terminal × Background Jobs ×

R 4.2.2 · ~/

```

> # Remove outliers in 'Worldwide.Gross'
> Q1 <- quantile(no_outliers$Worldwide.Gross, .25)
> Q3 <- quantile(no_outliers$Worldwide.Gross, .75)
>
> #inter quartile range
> IQR <- IQR(no_outliers$Worldwide.Gross)
> #Remove outliers in the 'no_outliers' data frame and assign the new subset to 'TopMovies'
>
> TopMovies <- subset(no_outliers, no_outliers$Worldwide.Gross > (Q1 - 1.5*IQR) &
+                     no_outliers$Worldwide.Gross < (Q3 + 1.5*IQR))
> 
```

Properties	
\$ Film	: chr "27 Dresses"
\$ Genre	: chr "Comedy"
\$ Lead.Studio	: chr "Fox" "Fox"
\$ Audience..score..	: int 71 81 89 0
\$ Profitability	: num 5.34 8.1 0
\$ Rotten.Tomatoes..	: int 40 87 79 8
\$ Worldwide.Gross	: num 160.31 60.
\$ Year	: int 2008 2009

Values

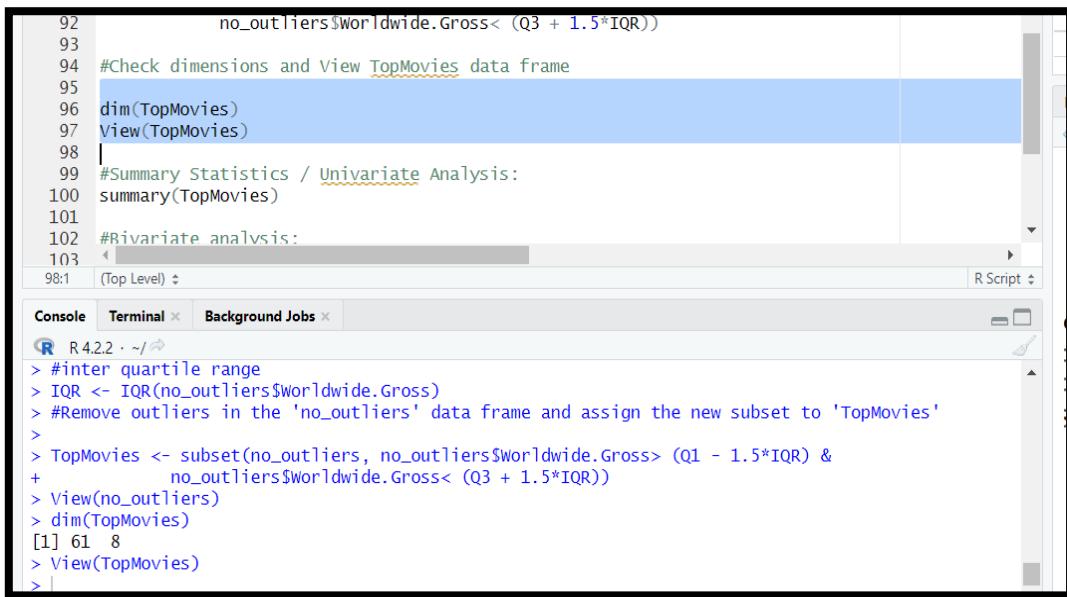
IQR	146.25
Q1	Named num 32.6
Q3	Named num 179

Files Plots Packages Help Viewer Presentation

Zoom Export

Worldwide.Gross

The dimensions of new data frame "TopMovies" after the removal of outliers from the Worldwide Gross variable are now 61 rows and 8 columns.



```

92     no_outliers$Worldwide.Gross < (Q3 + 1.5*IQR))
93
94 #Check dimensions and View TopMovies data frame
95
96 dim(TopMovies)
97 View(TopMovies)
98
99 #Summary Statistics / Univariate Analysis:
100 summary(TopMovies)
101
102 #Bivariate analysis:
103
103:1 (Top Level) R Script

```

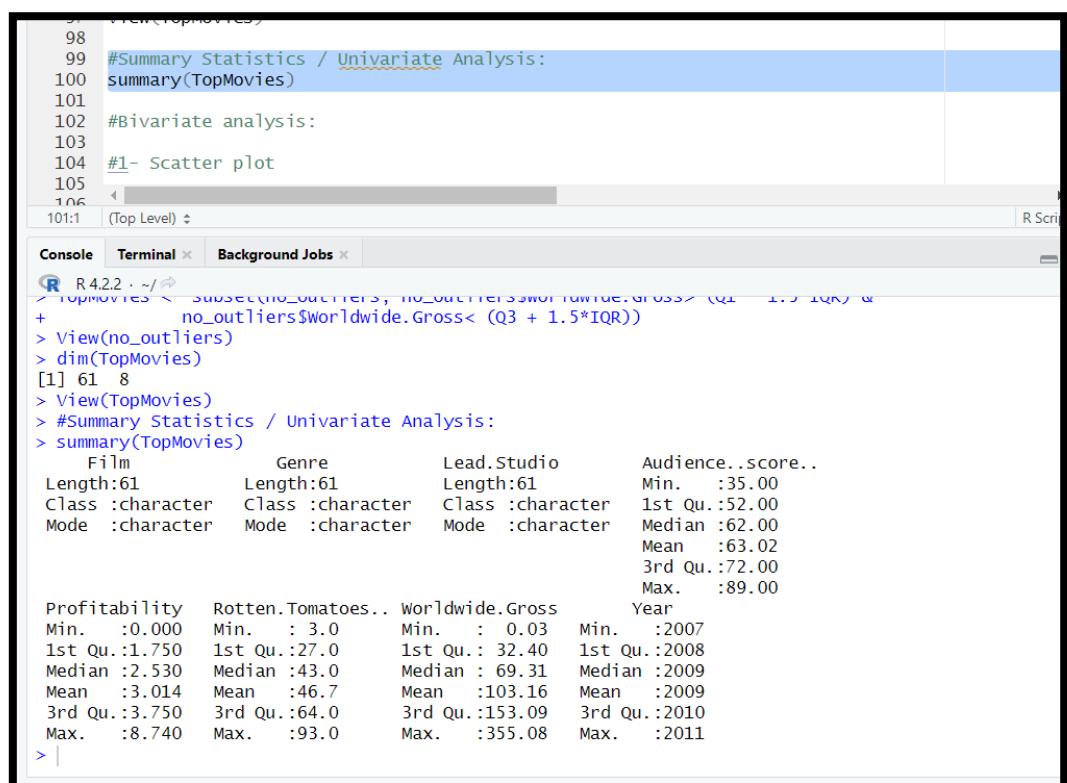
Console Terminal Background Jobs

```

R 4.2.2 · ~/ ...
> #inter quartile range
> IQR <- IQR(no_outliers$Worldwide.Gross)
> #Remove outliers in the 'no_outliers' data frame and assign the new subset to 'TopMovies'
>
> TopMovies <- subset(no_outliers, no_outliers$Worldwide.Gross > (Q1 - 1.5*IQR) &
+   no_outliers$Worldwide.Gross < (Q3 + 1.5*IQR))
> View(no_outliers)
> dim(TopMovies)
[1] 61 8
> View(TopMovies)
> |

```

11. SUMMARY STATISTICS (UNIVARIATE ANALYSIS) R provides a wide range of functions for obtaining summary statistics or single variable analysis. One method of obtaining descriptive statistics of a data set is to use the `summary()` function. The summary statistics for each of the 8 Variables in the dataset are shown below.



```

97 View(TopMovies)
98
99 #Summary Statistics / Univariate Analysis:
100 summary(TopMovies)
101
102 #Bivariate analysis:
103
104 #1- Scatter plot
105
106
106:1 (Top Level) R Script

```

Console Terminal Background Jobs

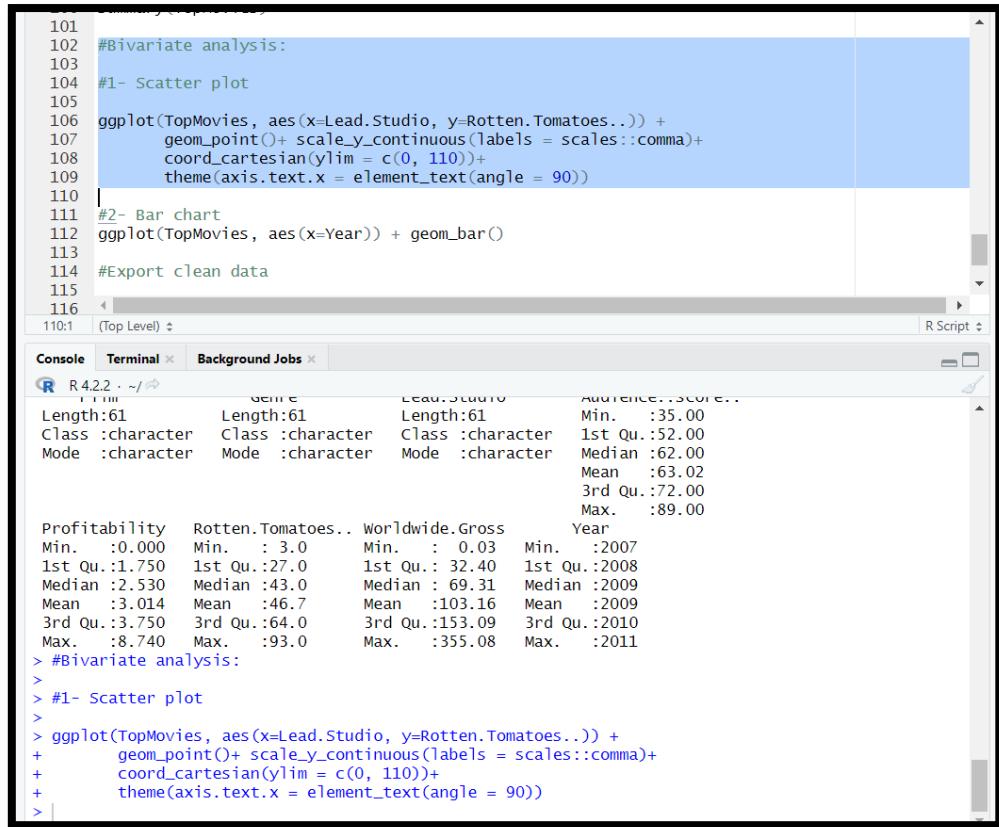
```

R 4.2.2 · ~/ ...
> TopMovies <- subset(no_outliers, no_outliers$Worldwide.Gross > (Q1 - 1.5*IQR) &
+   no_outliers$Worldwide.Gross < (Q3 + 1.5*IQR))
> View(no_outliers)
> dim(TopMovies)
[1] 61 8
> View(TopMovies)
> #Summary Statistics / Univariate Analysis:
> summary(TopMovies)
   Film          Genre          Lead.Studio      Audience..score..
Length:61    Length:61    Length:61      Min. :35.00
Class :character Class :character Class :character  1st Qu.:52.00
Mode  :character Mode  :character Mode  :character  Median :62.00
                                         Mean  :63.02
                                         3rd Qu.:72.00
                                         Max. :89.00
Profitability  Rotten.Tomatoes.. Worldwide.Gross      Year
Min.   :0.000  Min.   : 3.0  Min.   : 0.03  Min.   :2007
1st Qu.:1.750 1st Qu.:27.0 1st Qu.: 32.40 1st Qu.:2008
Median :2.530  Median :43.0  Median : 69.31  Median :2009
Mean   :3.014  Mean   :46.7  Mean   :103.16  Mean   :2009
3rd Qu.:3.750 3rd Qu.:64.0 3rd Qu.:153.09 3rd Qu.:2010
Max.   :8.740  Max.   :93.0  Max.   :355.08  Max.   :2011
> |

```

12. BIVARIATE ANALYSIS The term bivariate refers to the analysis of two variables. It is one of the simplest forms of statistical analysis. It is generally used to find out if there is a relationship between two sets of values or two variables. In present case, we have used scatter plot and bar chart to do the bivariate analysis.

- a. **SCATTER PLOT** A scatter plot is used to display the relationship between two variables Lead Studio and Rotten Tomatoes Ratings.



The screenshot shows the RStudio environment. The top pane displays R code for performing bivariate analysis on movie data. The bottom pane shows the R console output, which includes summary statistics for various variables like genre, lead studio, audience score, and year, followed by the ggplot command for the scatter plot.

```

101
102 #Bivariate analysis:
103
104 #1- Scatter plot
105
106 ggplot(TopMovies, aes(x=Lead.Studio, y=Rotten.Tomatoes..)) +
107   geom_point() + scale_y_continuous(labels = scales::comma) +
108   coord_cartesian(ylim = c(0, 110)) +
109   theme(axis.text.x = element_text(angle = 90))
110
111 #2- Bar chart
112 ggplot(TopMovies, aes(x=Year)) + geom_bar()
113
114 #Export clean data
115
116
110:1 (Top Level) ◊ R Script ◊

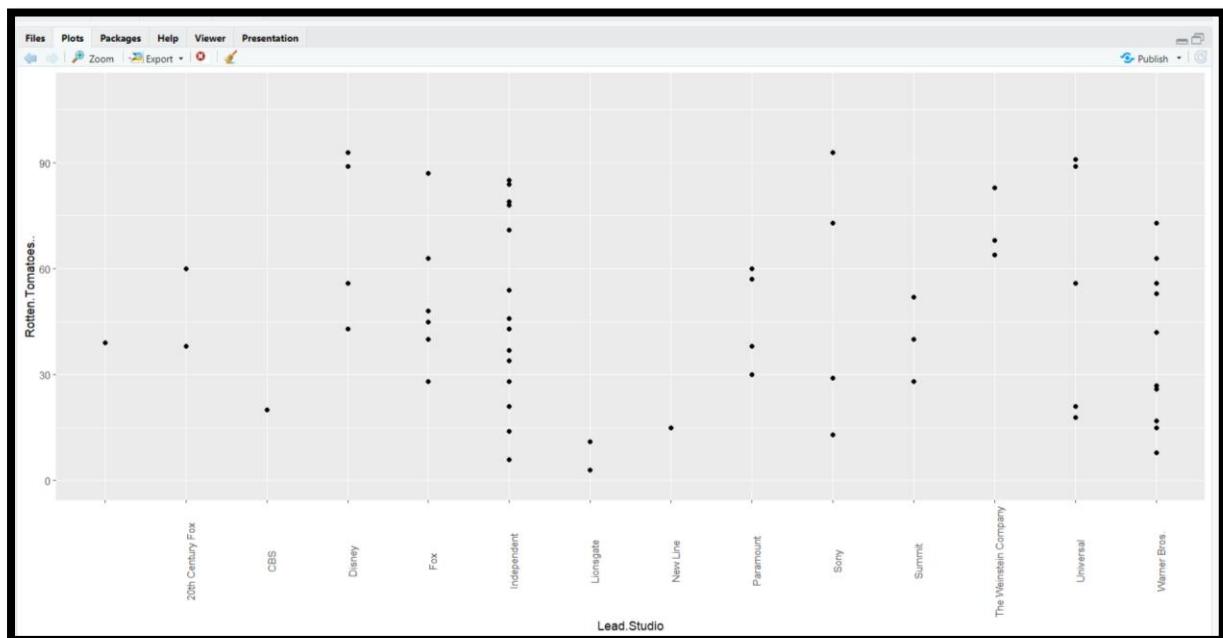
```

```

R 4.2.2 : ~/ ◊
Length:61      Length:61      Length:61      Audience..Score...
Class :character Class :character Class :character 1st Qu.:52.00
Mode  :character Mode  :character Mode  :character Median :62.00
                                         Mean  :63.02
                                         3rd Qu.:72.00
                                         Max.  :89.00
Profitability    Rotten.Tomatoes.. Worldwide.Gross      Year
Min.   :0.000   Min.   : 3.0   Min.   : 0.03   Min.   :2007
1st Qu.:1.750   1st Qu.:27.0   1st Qu.: 32.40  1st Qu.:2008
Median :2.530   Median :43.0   Median : 69.31  Median :2009
Mean   :3.014   Mean   :46.7   Mean   :103.16  Mean   :2009
3rd Qu.:3.750   3rd Qu.:64.0   3rd Qu.:153.09  3rd Qu.:2010
Max.   :8.740   Max.   :93.0   Max.   :355.08  Max.   :2011
> #Bivariate analysis:
>
> #1- Scatter plot
>
> ggplot(TopMovies, aes(x=Lead.Studio, y=Rotten.Tomatoes..)) +
+   geom_point() + scale_y_continuous(labels = scales::comma) +
+   coord_cartesian(ylim = c(0, 110)) +
+   theme(axis.text.x = element_text(angle = 90))
>

```

Following is the image of the scatter plot between Rotten Tomatoes Ratings and Lead Studios.



- b. **BAR CHART** A bar chart is used to display the relationship between two variables Year and Count (Number of Films produced each year).

```

105 ggplot(TopMovies, aes(x=Lead.Studio, y=Rotten.Tomatoes..)) +
106   geom_point() + scale_y_continuous(labels = scales::comma) +
107   coord_cartesian(ylim = c(0, 110)) +
108   theme(axis.text.x = element_text(angle = 90))
109
110 #2- Bar chart
111 ggplot(TopMovies, aes(x=Year)) + geom_bar()
112
113
114 #Export clean data
115
116 write.csv(TopMovies, "Hollywood_Most_Profitable_Stories_Cleaned.csv")
117
118 (Top Level) R Script

```

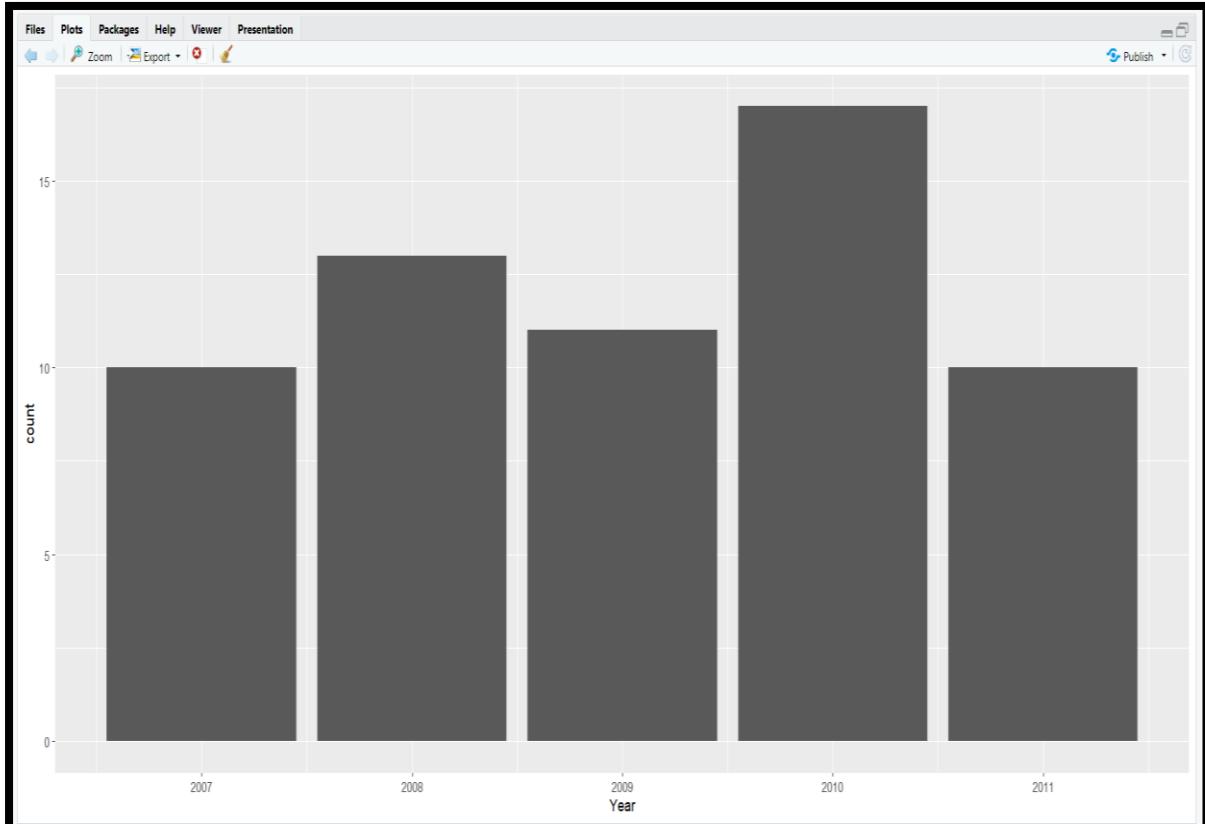
Console Terminal < Background Jobs <

```

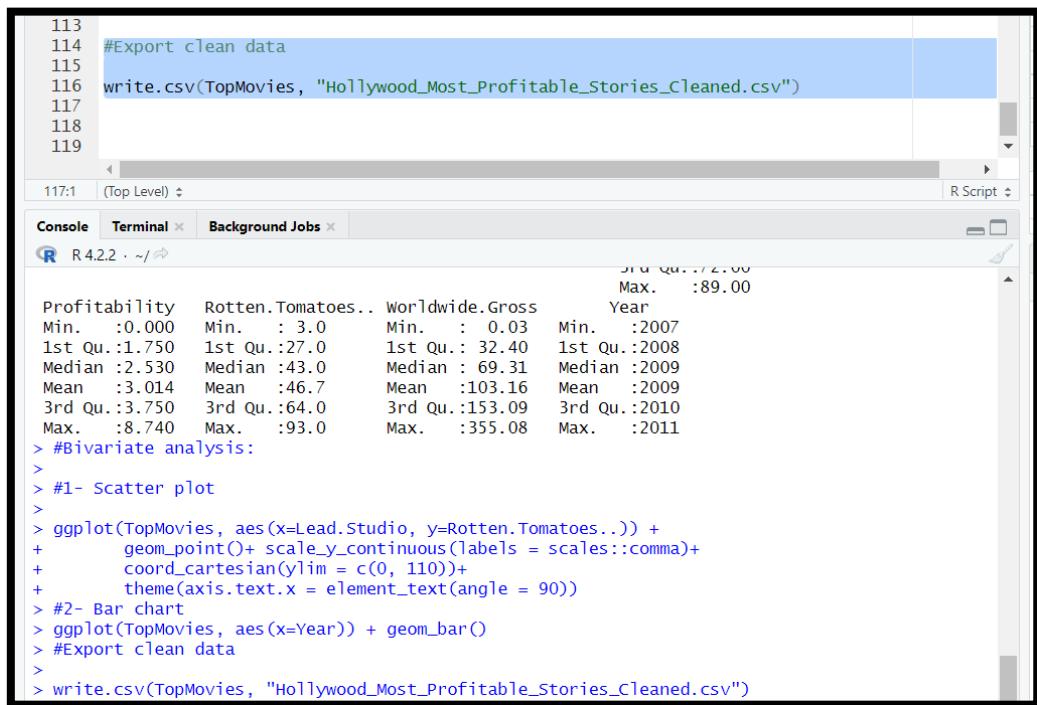
R 4.2.2 · ~/Desktop
Class : character   Class : character   Class : character   1st Qu.:52.00
Mode  :character   Mode  :character   Mode  :character   Median :62.00
                                         Mean   :63.02
                                         3rd Qu.:72.00
                                         Max.   :89.00
Profitability    Rotten.Tomatoes.. Worldwide.Gross      Year
Min.   :0.000   Min.   : 3.0   Min.   : 0.03   Min.   :2007
1st Qu.:1.750   1st Qu.:27.0   1st Qu.: 32.40   1st Qu.:2008
Median :2.530   Median :43.0   Median : 69.31   Median :2009
Mean   :3.014   Mean   :46.7   Mean   :103.16   Mean   :2009
3rd Qu.:3.750   3rd Qu.:64.0   3rd Qu.:153.09   3rd Qu.:2010
Max.   :8.740   Max.   :93.0   Max.   :355.08   Max.   :2011
> #Bivariate analysis:
>
> #1- Scatter plot
>
> ggplot(TopMovies, aes(x=Lead.Studio, y=Rotten.Tomatoes..)) +
+   geom_point() + scale_y_continuous(labels = scales::comma) +
+   coord_cartesian(ylim = c(0, 110)) +
+   theme(axis.text.x = element_text(angle = 90))
> #2- Bar chart
> ggplot(TopMovies, aes(x=Year)) + geom_bar()
>

```

Following is the image of the scatter plot between Rotten Tomatoes Ratings and Lead Studios.



13. EXPORT CLEANED DATA Finally, the cleaned data is exported to csv file named "Hollywood_Most_Profitable_Stories_Cleaned.csv" and then it can be used in any platform like Tableau or Power BI for data visualisation and story-telling.



The screenshot shows an RStudio interface. The code editor pane contains the following R script:

```
113
114 #Export clean data
115
116 write.csv(TopMovies, "Hollywood_Most_Profitable_Stories_Cleaned.csv")
117
118
119
```

The console pane shows the output of the script, including summary statistics for 'Profitability' and 'Year' columns:

```
Profitability   Rotten.Tomatoes.. Worldwide.Gross      Year
Min.    :0.000   Min.   : 3.0   Min.   : 0.03   Min.   :2007
1st Qu.:1.750   1st Qu.:27.0   1st Qu.: 32.40   1st Qu.:2008
Median  :2.530   Median  :43.0   Median  : 69.31   Median  :2009
Mean    :3.014   Mean    :46.7   Mean    :103.16   Mean    :2009
3rd Qu.:3.750   3rd Qu.:64.0   3rd Qu.:153.09   3rd Qu.:2010
Max.    :8.740   Max.   :93.0   Max.   :355.08   Max.   :2011
```

The script also includes comments for bivariate analysis, scatter plots, and bar charts, followed by the final command to export the data:

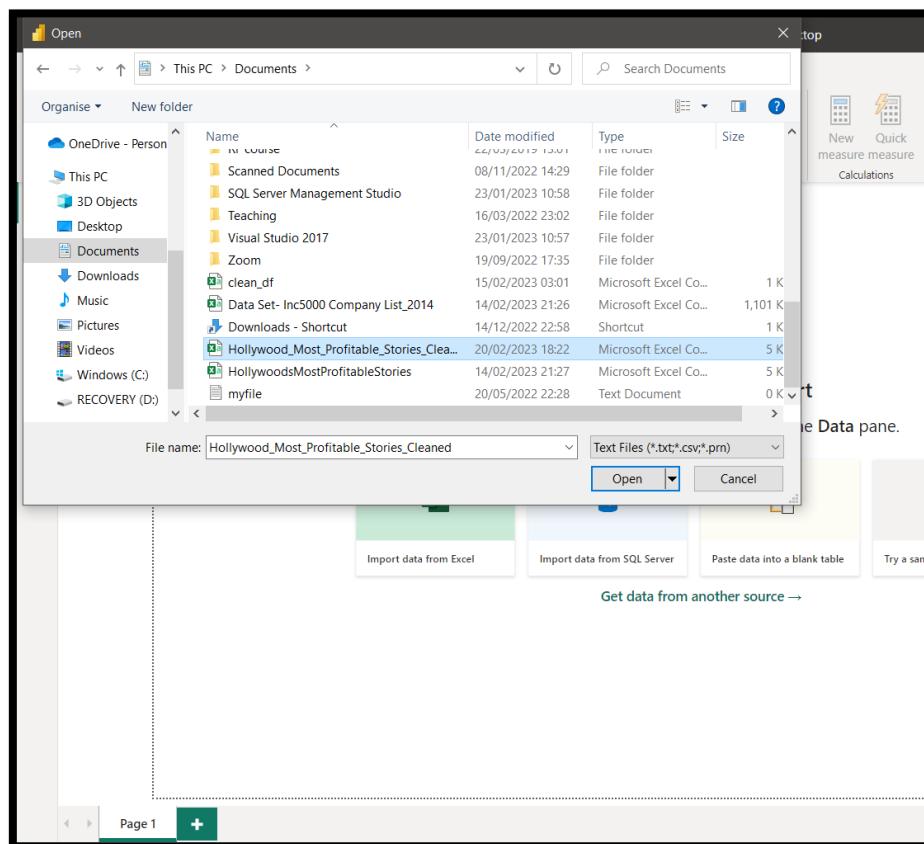
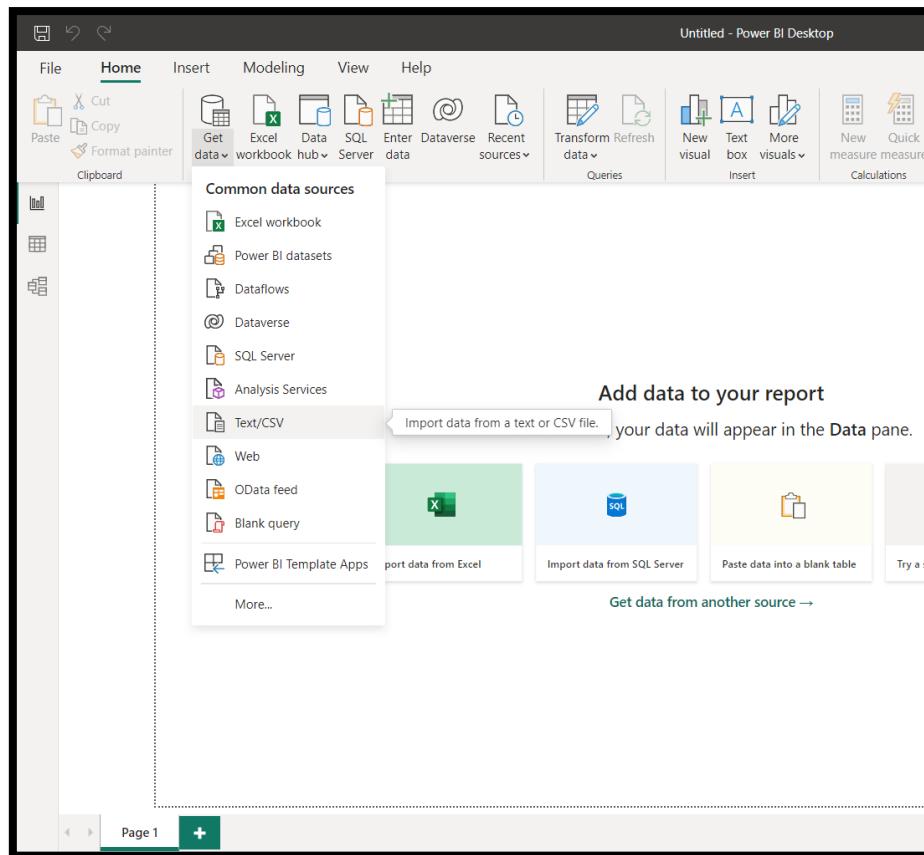
```
> #Bivariate analysis:
>
> #1- Scatter plot
>
> ggplot(TopMovies, aes(x=Lead.Studio, y=Rotten.Tomatoes..)) +
+     geom_point() + scale_y_continuous(labels = scales::comma) +
+     coord_cartesian(ylim = c(0, 110)) +
+     theme(axis.text.x = element_text(angle = 90))
> #2- Bar chart
> ggplot(TopMovies, aes(x=Year)) + geom_bar()
> #Export clean data
>
> write.csv(TopMovies, "Hollywood_Most_Profitable_Stories_Cleaned.csv")
```

TASK – 2 : DASHBOARD IN POWERBI

In this task, "Hollywood_Most_Profitable_Stories_Cleaned.csv" dataset is visualised in a dashboard using PowerBI. Following are the charts required by the client, but additional charts can be used for effective story telling:

- The average Rotten Tomatoes Ratings of each Genre
- The number of Movies produced per Year
- The Audience Scores for each Movie
- The Profitability per Studio
- The Worldwide Gross per Genre

1. IMPORT DATA Open PowerBI desktop and click on "Get Data" in "Text/CSV" header under the "Common Data Sources" section. Select the required file and click Open. This will import the "Hollywood_Most_Profitable_Stories_Cleaned.csv" dataset in PowerBI.



After opening the file, click Load and it will open the csv file in a table format as shown below.

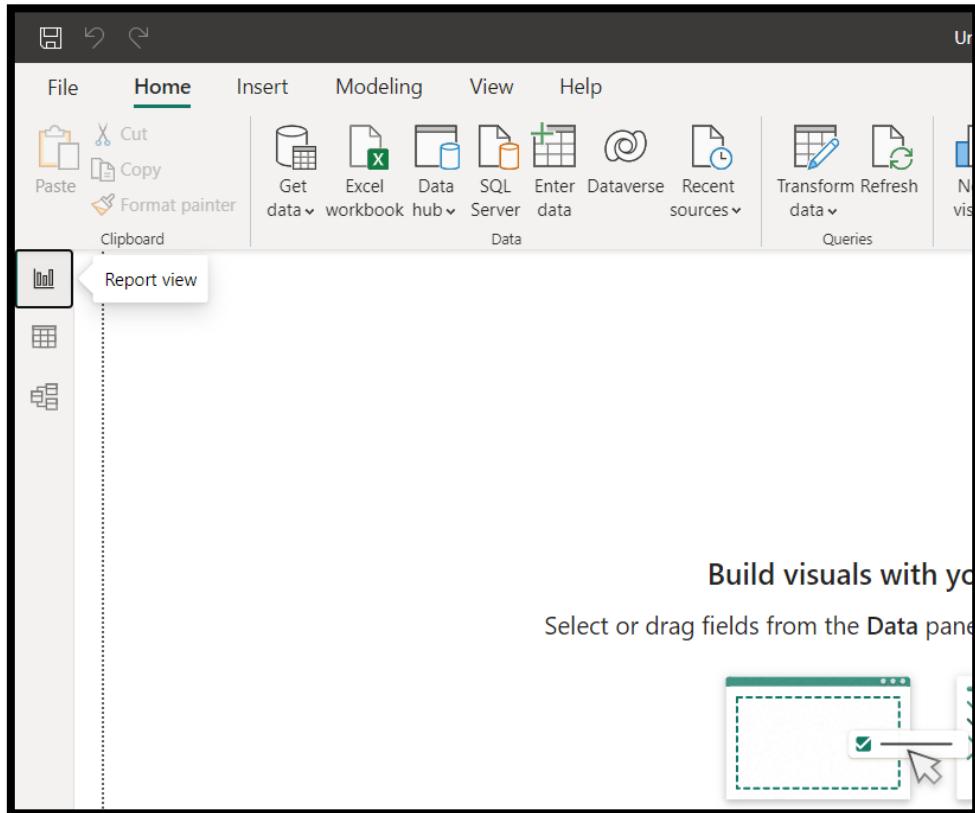
The data in the preview has been truncated due to size limits.

Extract Table Using Examples

Load Transform Data Cancel

The dataset is imported and ready for building the dashboard. You can check the header names on the right hand side of the PowerBI screen.

You can also switch between the report view and data view on the left hand side of the screen as shown below.



Untitled - Power BI Desktop

File Home Help Table tools

Name Hollywood_Most_P... Structure

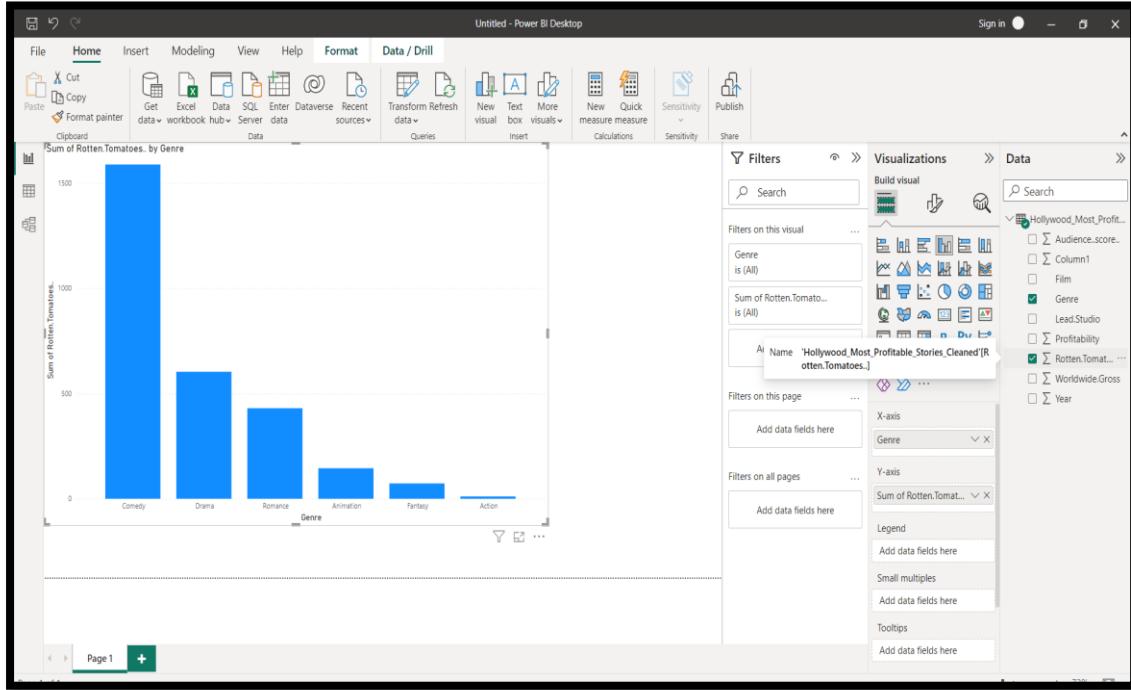
Mark as date table Calendars Manage relationships Relationships New measure Quick New measure column New table Calculations

Data view Film Genre Lead.Studio Audience..score.. Profitability Rotten.Tomatoes.. Worldwide.Gross Year

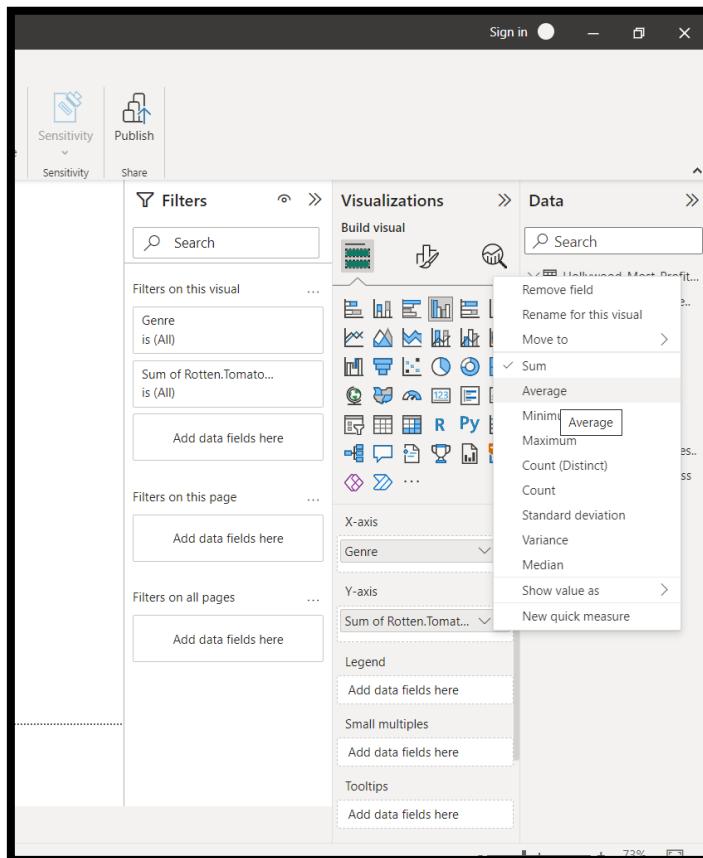
	Film	Genre	Lead.Studio	Audience..score..	Profitability	Rotten.Tomatoes..	Worldwide.Gross	Year
1	27 Dresses	Comedy	Fox	71	5.34	40	160.31	2008
2	(500) Days of Summer	Comedy	Fox	81	8.1	87	60.72	2009
3	A Dangerous Method	Drama	Independent	89	0.45	79	8.97	2011
4	A Serious Man	Drama	Universal	64	4.38	89	30.68	2009
5	Across the Universe	Romance	Independent	84	0.65	54	29.37	2007
6	Beginners	Comedy	Independent	80	4.47	84	14.31	2011
7	Dear John	Drama	Sony	66	4.6	29	114.97	2010
8	Enchanted	Comedy	Disney	80	4.01	93	340.49	2007
10	Four Christmases	Comedy	Warner Bros.	52	2.02	26	161.83	2008
11	Ghosts of Girlfriends Past	Comedy	Warner Bros.	47	2.04	27	102.22	2009
12	Gnomeo and Juliet	Animation	Disney	52	5.39	56	193.97	2011
13	Going the Distance	Comedy	Warner Bros.	56	1.31	53	42.05	2010
14	Good Luck Chuck	Comedy	Lionsgate	61	2.37	3	59.19	2007
15	He's Just Not That Into You	Comedy	Warner Bros.	60	7.15	42	178.84	2009
17	I Love You Phillip Morris	Comedy	Independent	57	1.34	71	20.1	2010
18	It's Complicated	Comedy	Universal	63	2.64	56	224.6	2009
20	Just Wright	Comedy	Fox	58	1.8	45	21.57	2010
21	Killers	Action	Lionsgate	45	1.25	11	93.4	2010
22	Knocked Up	Comedy	Universal	83	6.64	91	219	2007
23	Leap Year	Comedy	Universal	49	1.72	21	32.59	2010
24	Letters to Juliet	Comedy	Summit	62	2.64	40	79.18	2010
25	License to Wed	Comedy	Warner Bros.	55	1.98	8	69.31	2007
26	Life as We Know It	Comedy	Independent	62	2.53	28	96.16	2010
27	Love & Other Drugs	Comedy	Fox	55	1.82	48	54.53	2010
28	Love Happens	Drama	Universal	40	2	18	36.08	2009
29	Made of Honor	Comedy	Sony	61	2.65	13	105.96	2008
31	Marley and Me	Comedy	Fox	77	3.75	63	206.07	2008
32	Midnight in Paris	Romance	Sony	84	8.74	93	148.66	2011

Table: Hollywood_Most_Profitable_Stories_Cleaned (61 rows)

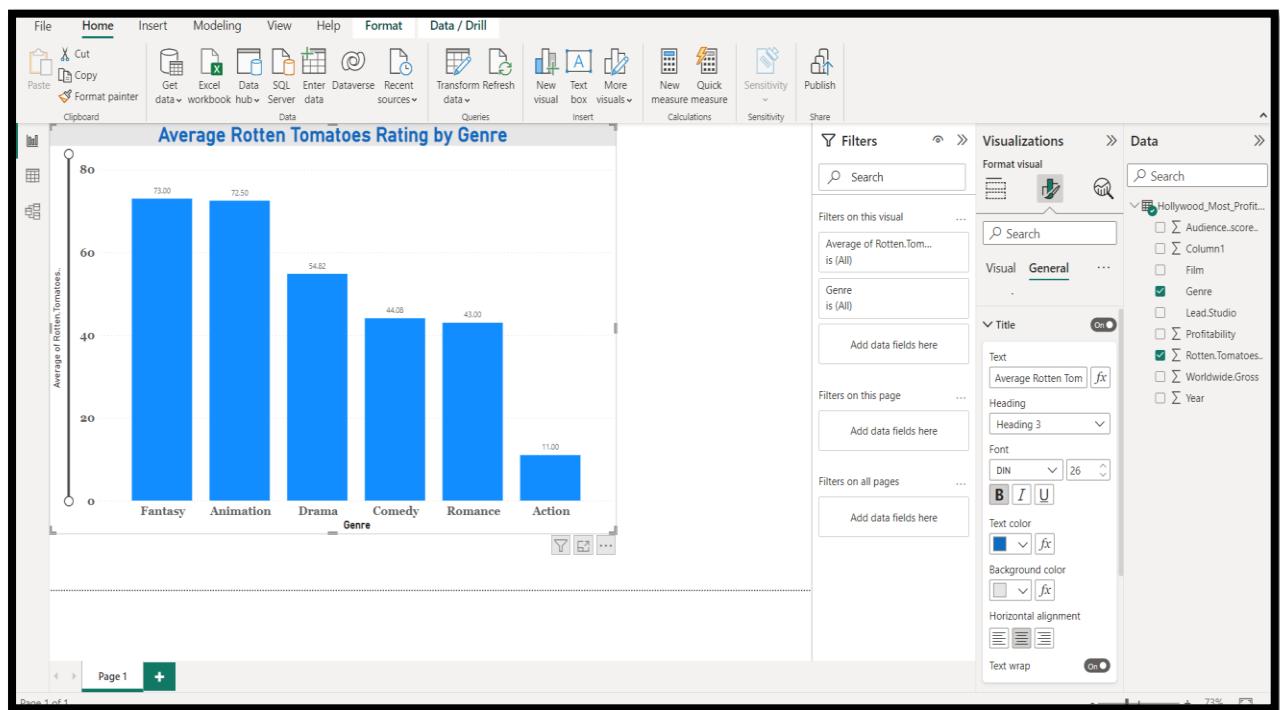
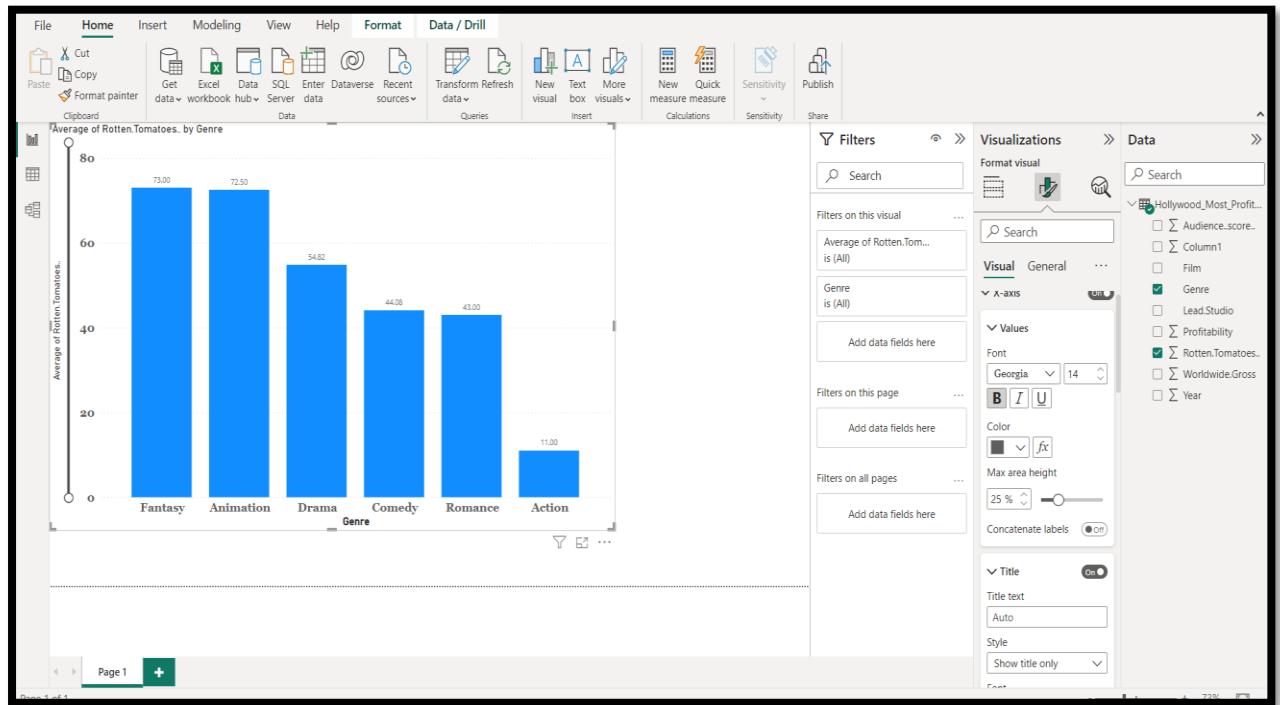
2. **BUILD VISUALS** To start building , visualisations, you can select or drag and drop the required variables to data field placeholder and then select the visual type.



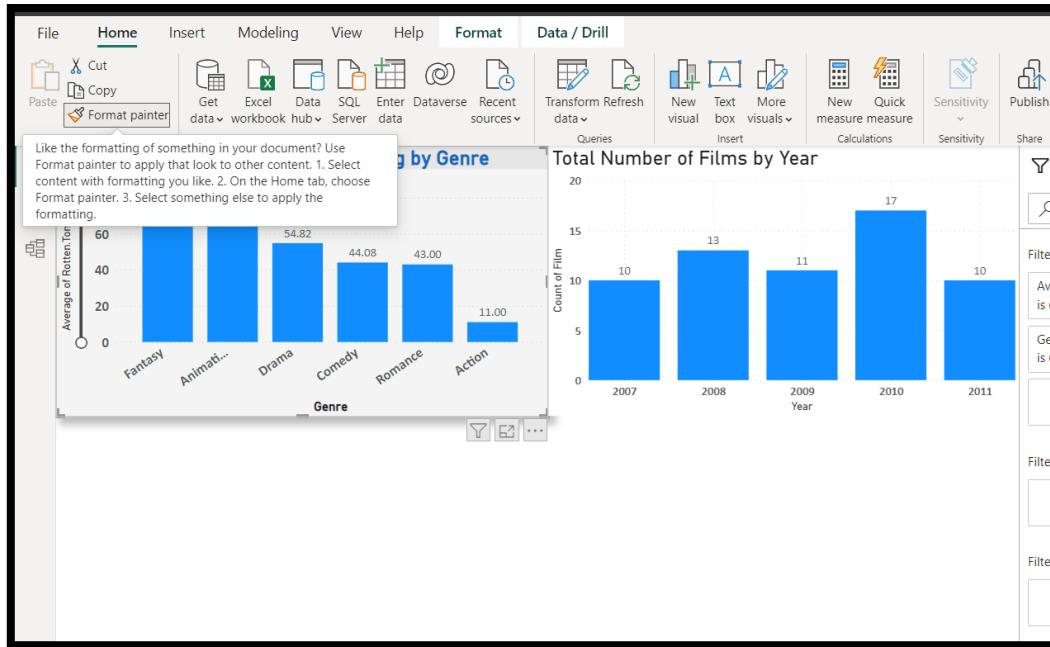
After building the chart, you can change the value to be shown as sum, average etc. For the first chart, which requires Average Rotten Tomatoes Ratings by Genre. We will display the variable as an average.



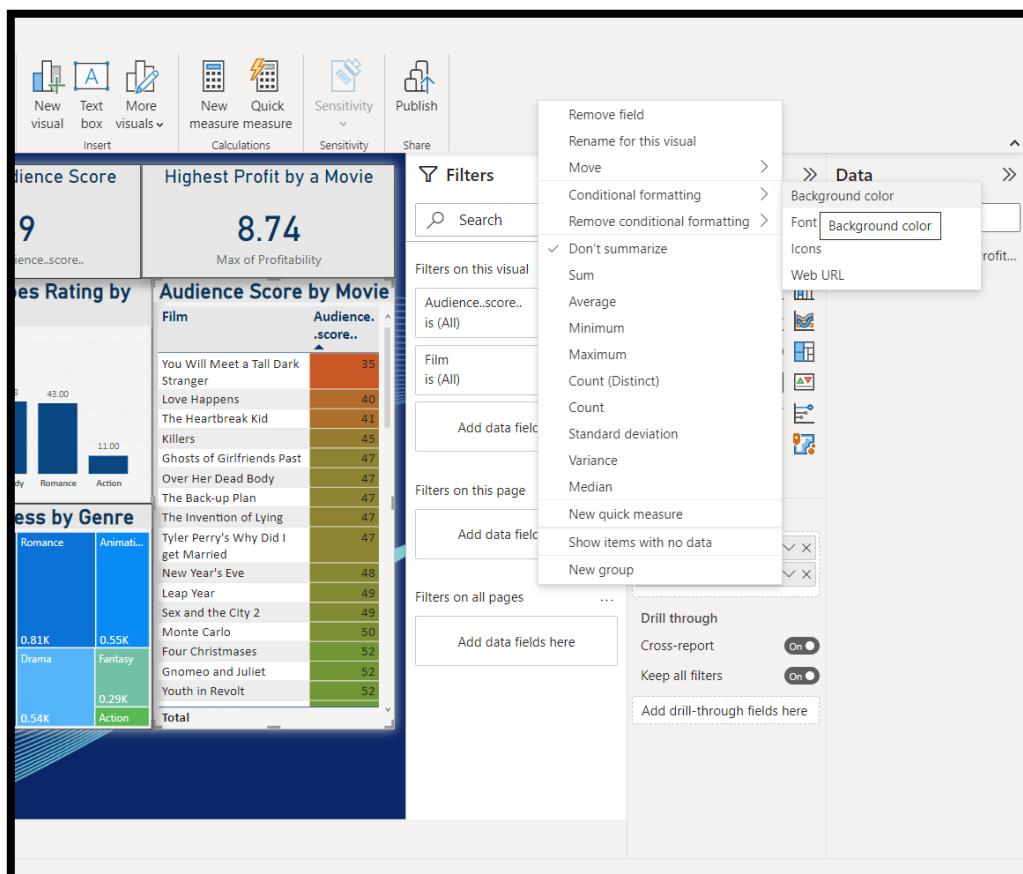
You can also add title, data labels, fonts and colour attributes to your chart to make it more visually appealing and more understandable for the stakeholders.



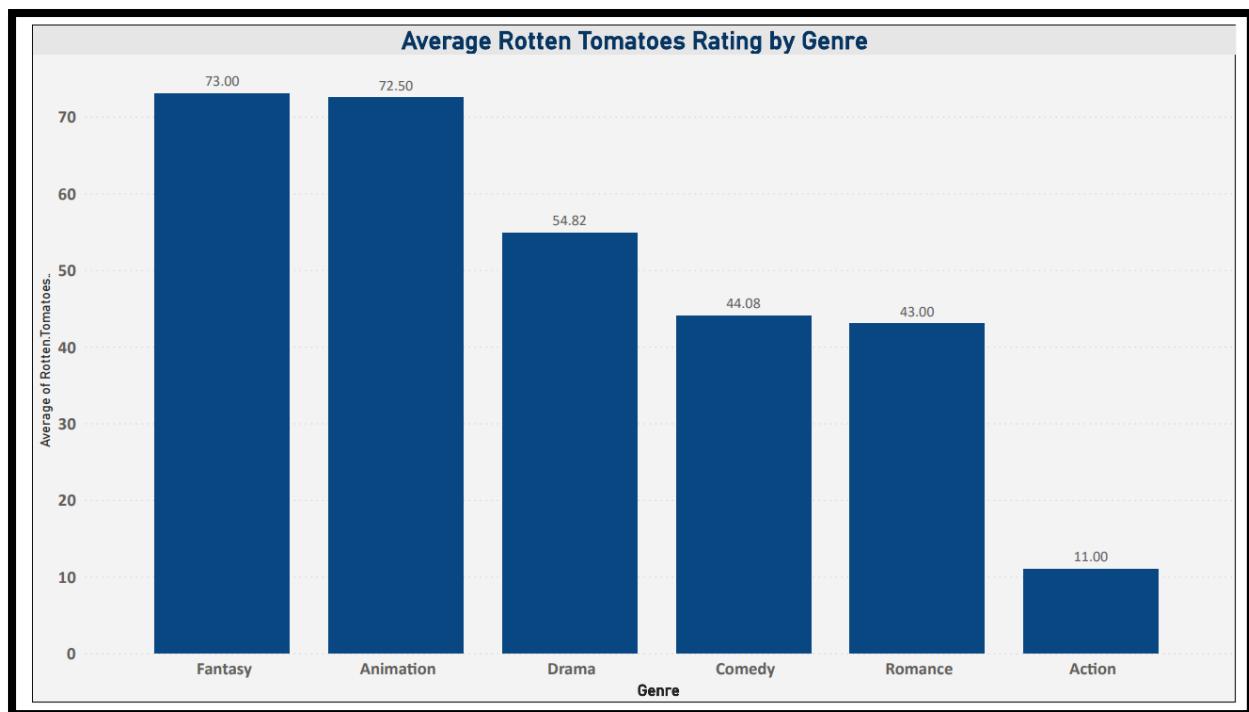
Once you have created one visualisation, you can simply select it and click format painter and paste it the rest of the visualisations to keep formatting consistent throughout your dashboard. It also makes the building process seamless and efficient.



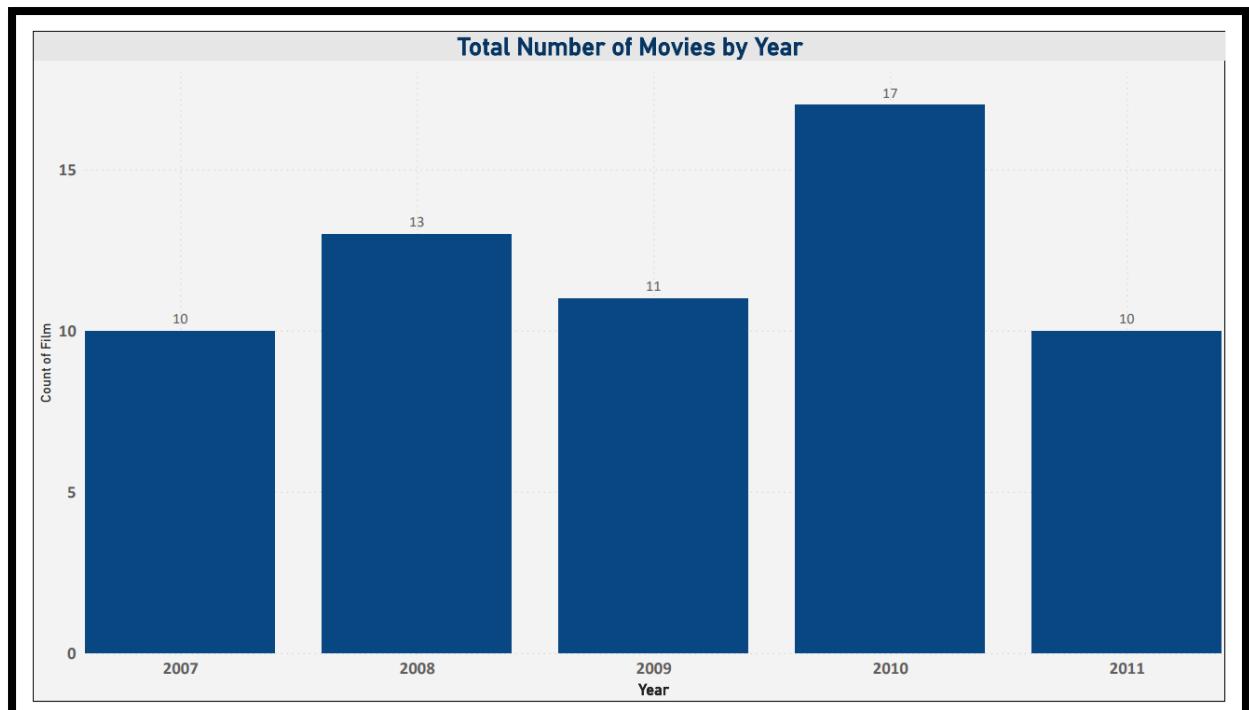
You can add features like conditional formatting to the table sorting the values from low to high or vice versa and assigning a colour palette to it.



To bar chart displaying the Average Rotten Tomatoes Ratings by Genre is shown below. The Genre Fantasy has the highest rating with Animation as a close second while Action Genre has the lowest Rotte Tomatoes Ratings by user.



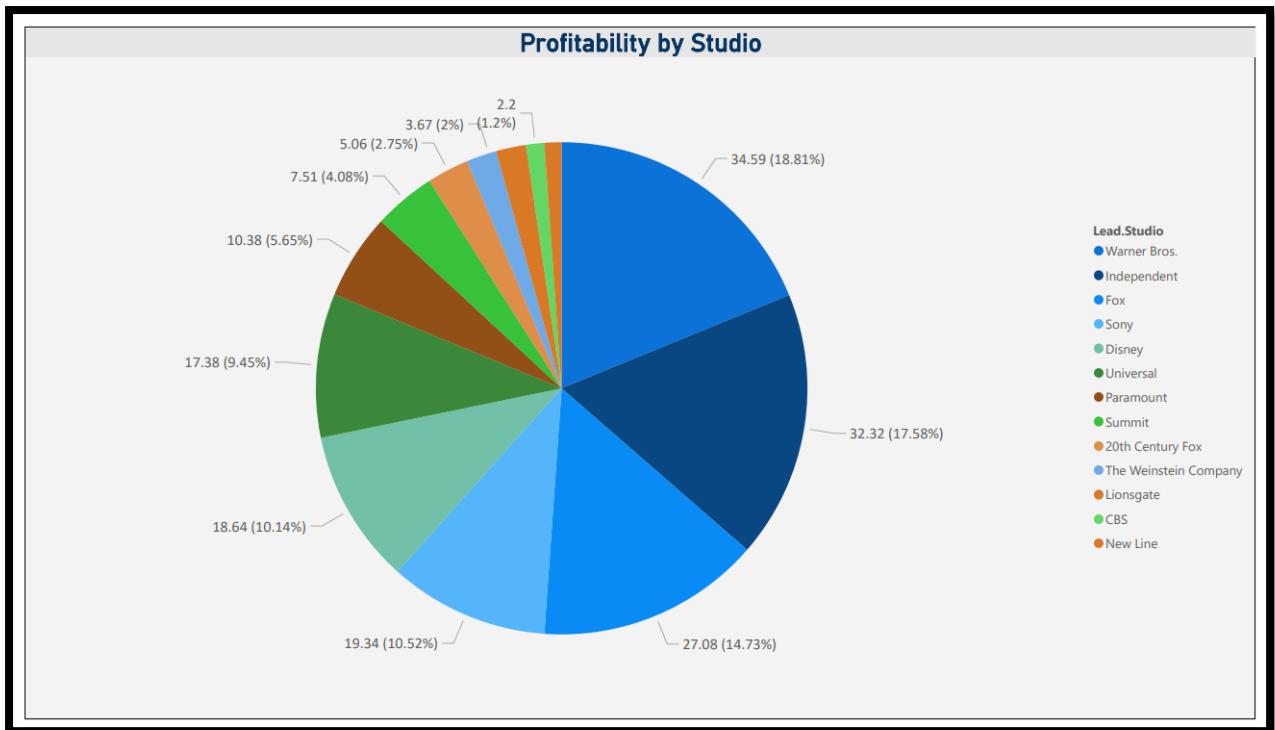
To bar chart displaying the Total Number of Movies produced each year is shown below. The highest number of movies (17) were produced in year 2010 while 2007 and 2010 has the same movie count of 10.



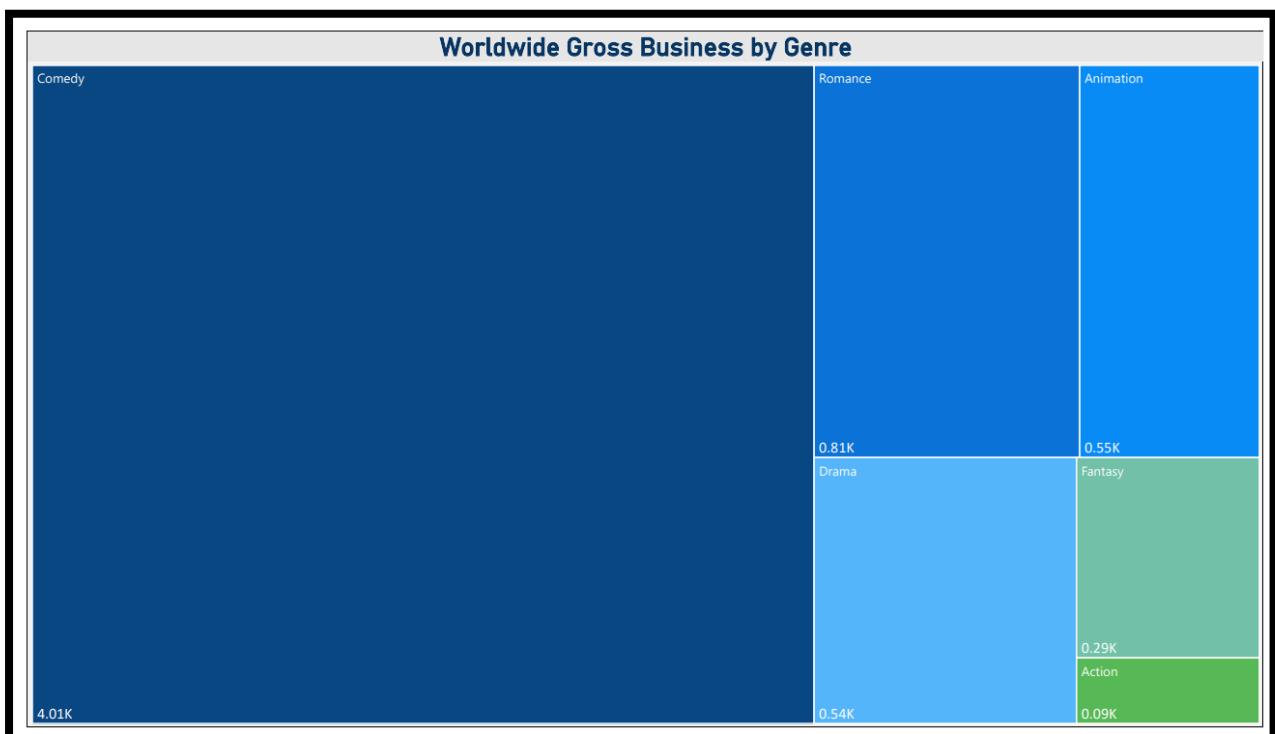
The table displaying Audience Score for each film is shown below, displaying the values in ascending order and with the conditional formatting in colour palette (green, blue and brown) as required by the client .

Film	Audience. .score..
Zack and Miri Make a Porno	70
Youth in Revolt	52
You Will Meet a Tall Dark Stranger	35
What Happens in Vegas	72
Water For Elephants	72
Waiting For Forever	53
Valentine's Day	54
Tyler Perry's Why Did I get Married	47
The Ugly Truth	68
The Time Traveler's Wife	65
The Proposal	74
The Invention of Lying	47
The Heartbreak Kid	41
The Duchess	68
The Curious Case of Benjamin Button	81
The Back-up Plan	47
Tangled	88
She's Out of My League	60
Sex and the City 2	49
Remember Me	70
Total	61

The pie chart displaying the Profitability by Studio is shown below. Warner Bros. is the studio with highest Profitability while New line studio has the lowest Profitability.



The tree map displaying the Worldwide Gross Business by Genre is shown below. The Genre Comedy has the highest worldwide gross business while action is the lowest.



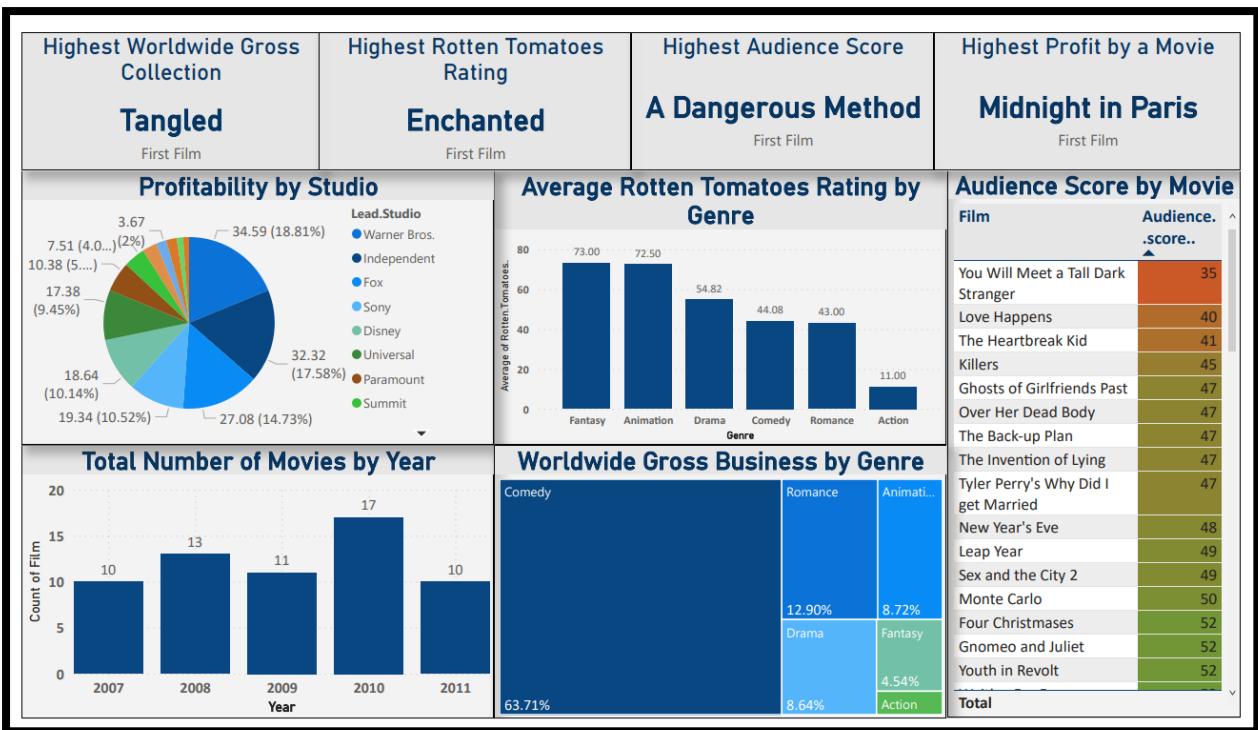
To display cards with highest values, you can use the filter and select Top N values as shown below.

The screenshot shows the Power BI Desktop interface with the title "Hollywood Most Profitable Movies - Power BI Desktop". The ribbon at the top has a "Drill" tab selected, followed by "Queries", "New visual", "Text box", "More visuals", "Insert", "New measure", "Quick measure", "Calculations", "Sensitivity", and "Publish". On the right side, there's a "Filters" pane and a "Visualizations" pane. The "Filters" pane shows a search bar and a section for "Filters on this visual" with a dropdown for "First Film" set to "(All)". It also includes a "Filter type" dropdown set to "Top N", a "Show items" dropdown set to "Top 1", and a "By value" field containing "Max of Worldwide.Gross". The "Visualizations" pane shows a list of chart types like Line, Area, Bar, etc., with "Line" currently selected. A card visual titled "Highest Worldwide Gross Collection" is displayed, showing "Tangled" as the First Film. The main canvas area features a blue background with abstract wavy lines.

3. BUILD DASHBOARD

Once you have completed building all your individual charts, you can re-arrange them in an appropriate and visually appealing order to build your dashboard. For this project, client asked for a colour-palette consisting of blue, green and brown colours. The dashboard shown below has been created following the guidelines and principles outlined by the client.

The final dashboard for the dataset is shown below. It gives insight to the Hollywood most Profitable Movies between the year 2007-2012.



REFLECTIVE

This project has helped me to utilise the knowledge of programming language R for data cleaning and analysis and Microsoft Power BI to visualise the data for effective storytelling. I found it highly interesting to learn to write the script in R as I enjoy coding. I particularly enjoyed using PowerBI to visualise the data and see the story behind the data set come to life. It provided me insights into dataset which didn't make much sense to begin with.