



KOLEJ PROFESIONAL MARA BERANANG

DIPLOMA IN COMPUTER SCEINCE

COURSE NAME	: DATA ANALYSIS & VISUALIZATION
COURSE CODE	: CSC 2823
SESSION	: 3 2024/2025
TYPE OF ASSESSMENT	: FINAL PROJECT
DURATION	: 03/02/2025 – 21/02/2025

CLO 3 Prepare data visualization for effective presentation using computer software tool.

INSTRUCTION TO CANDIDATES:

1. Students are required to answer ALL questions.
2. Students need to submit the assignment report in hardcopy.

Personal Details	
Name	
I/D Number	
Class	<input type="checkbox"/> DCS 5A <input type="checkbox"/> DCS 5B <input type="checkbox"/> DCS 5C <input type="checkbox"/> DCS 5D
Lecturer	<input type="checkbox"/> DR.ZALINA AYOB <input type="checkbox"/> MAWARWIDURI BT AB HALIK

Task	Marks
1	
2	
3	
4	
5	
Total	/ 50

PROJECT SCENARIO

Assume you are a data analyst working for a company, government agency, or nonprofit organization. Your team has been tasked with providing actionable insights using data to address specific challenges or make better decisions. An example of dataset relevant to a specific domain as shown in Figure 1 below. You are required to identify and analyze the selected dataset, create visualizations, and present your findings clearly, informative, and useful for stakeholders.

Domain	Example of dataset
Healthcare	<ul style="list-style-type: none">▪ Patient health records (age, gender, diagnosis, treatments, recovery rates)▪ Disease outbreaks (geographical spread, case counts, recovery rates)▪ Hospital performance (admissions, bed availability, patient satisfaction)
Retail & E-commerce	<ul style="list-style-type: none">▪ Product sales (date, category, region, revenue, units sold)▪ Customer behaviours (age group, purchase frequency, cart abandonment rate)▪ Inventory data (stock levels, restocking frequency).
Education	<ul style="list-style-type: none">▪ Student performance (grades, attendance, participation)▪ Enrolments trends (age, gender, region, program type)▪ Online learning metrics (time spent on platform, quiz scores, course completion).
Transportation	<ul style="list-style-type: none">▪ Public transport usage (daily ridership, routes, timings)▪ Traffic data (vehicle count, congestion levels, accident locations)▪ Flight or train delays (departure time, delay length, reason for delay).
Sports	<ul style="list-style-type: none">▪ Player statistics (points, assists, performance over seasons)▪ Match results (team, score, location, audience size)▪ Fan engagement (ticket sales, merchandise, social media activity).
Environmental science	<ul style="list-style-type: none">▪ Climate data (temperature, precipitation, CO2 levels, time)▪ Air quality (pollutants, AQI, region, time)▪ Wildlife populations (species, count, region, time).
Energy and utilities	<ul style="list-style-type: none">▪ Energy consumption (time, location, source type, usage)▪ Renewable energy output (solar, wind, hydro generation by region)▪ Utility service complaints (frequency, type, resolution time).
Public safety and law enforcement	<ul style="list-style-type: none">▪ Crime statistics (type, location, time, resolution status).▪ Emergency response times (type of emergency, location).▪ Traffic violations (type, location, time of day).
Gaming industry	<ul style="list-style-type: none">▪ Player activity (time spent, levels completed, purchases).▪ Game performance metrics (bugs, crashes, frame rate).▪ Monetization (in-app purchases, subscriptions, ad revenue).
Media and entertainment	<ul style="list-style-type: none">▪ Streaming platform data (viewership, time spent, genres watched).▪ Box office sales (ticket sales, revenue, regions).▪ Social media activity (hashtags, trends, engagement).

Urban Planning and smart city	<ul style="list-style-type: none"> ▪ City infrastructure usage (traffic, utilities, public transport). ▪ Population density (by region, time of day). ▪ Pollution levels (air, water, noise by region).
-------------------------------	--

Figure 1

The following steps outline the processes and the tasks that you need to complete for the project:

Task 1: Identification of dataset and analysis background

- Identify domain/key areas based on Figure 1. You may use suggested dataset or other relevant examples.

Task 2: Gather Relevant Data

- Gather data to support the generation of effective visualization.

Task 3: Clean and Preprocess the Data

- Raw data is cleaned and pre-processed to ensure it can be used effectively.

Task 4: Analyze the data to extract meaningful insights and trends.

- Analyze the data to extract meaningful insights and trends using:
 - ☐ Univariate Data Analysis
 - ☐ Bivariate Data analysis

Task 5: Create visually appealing and informative data visualizations to convey the identified insights

- Utilize charts, graphs, maps, and other visualization techniques to make the information easily understandable. The visualization should use different types of variables and suitable presentation approaches.

Assessment Rubrics:

No	Attribute	Task	(1 mark)	(2 Marks)	(3 Marks)	(4 Marks)	Marks Weighted	Marks Obtained
1	Gather	Identification of dataset and analysis background					1	
		Identify domain area and the analysis background	Able to provide at least ONE (1) of the following: <ul style="list-style-type: none"> <input type="checkbox"/> Provide objective of the analysis correctly <input type="checkbox"/> Able to identify target audience appropriately 	Able to provide All of the following: <ul style="list-style-type: none"> <input type="checkbox"/> Provide objectives of the analysis correctly. <input type="checkbox"/> Able to identify target audience appropriately 	Able to provide All of the following: <ul style="list-style-type: none"> <input type="checkbox"/> Precisely explain the objective of the analysis correctly. <input type="checkbox"/> Able to identify target audiences appropriately 	Able to provide All of the following: <ul style="list-style-type: none"> <input type="checkbox"/> Precisely explain the objective of the analysis correctly. <input type="checkbox"/> establish the scope of the analysis, outlining what will be included and excluded <input type="checkbox"/> Able to identify target audience appropriately and explain the benefit gained from the analysis 		
2		Gather Relevant Data					1	
		Gather data to support the generation of effective visualization.	Able to provide ALL the following: <ul style="list-style-type: none"> <input type="checkbox"/> Able to justify appropriately the reason for your data selection. <input type="checkbox"/> Able to select suitable dataset to fulfill the objective(s) 	Able to provide ALL the following: <ul style="list-style-type: none"> <input type="checkbox"/> Able to justify appropriately the reason for your data selection. <input type="checkbox"/> Able to select suitable dataset to fulfill the objective(s) 				

			<input type="checkbox"/> Provide evidence for the dataset from various sources (website/pre-installed RStudio etc)	<input type="checkbox"/> Provide evidence for the dataset from various sources (website/pre-installed RStudio etc) <input type="checkbox"/> Able to specify suitable variable(s) from the dataset for the analysis				
3		Clean and Preprocess the Data						
		Clean and Preprocess the Data	<input type="checkbox"/> Provide a screen shot of with insufficient steps to treat any missing values or quantities of zero in the dataset	<input type="checkbox"/> Provide a screen shot of the correct steps to treat any missing values or quantities of zero in the dataset	<input type="checkbox"/> Provide a screen shot of the correct steps to treat any missing values or quantities of zero in the dataset. <input type="checkbox"/> Provide a minimal explanation of the steps.	<input type="checkbox"/> Provide a shot of the correct steps to treat any missing values or quantities of zero in the dataset. <input type="checkbox"/> Provide a precise and clear explanation of each step in the screen shots	1	
4	Reproduce and Process Information	Analyze the data to extract meaningful insights and trends.						
		Univariate Data analysis	<input type="checkbox"/> Able to provide one category of univariate analysis <input type="checkbox"/> Able to display one visual on category with a unsuitable selection of a plot or graph for any category	<input type="checkbox"/> Able to provide one category of univariate analysis <input type="checkbox"/> Able to display one visual on category with a suitable selection of a plot or graph for any category	<input type="checkbox"/> Able to provide more than one category of univariate analysis <input type="checkbox"/> Able to display visual on categories with a suitable selection of a plot or graph for each category	<input type="checkbox"/> Able to provide more than one category of univariate analysis <input type="checkbox"/> Able to display visual on categories with a suitable selection of a plot or graph for each category <input type="checkbox"/> Provide a clear and concise	2	

						finding explanation based on the analysis conducted for each category.		
		Bivariate Data Analysis	<ul style="list-style-type: none"> ○ Able to provide one category of Bivariate analysis ○ Able to display one visual on category with a suitable selection of a plot or graph for any category 	<ul style="list-style-type: none"> ○ Able to provide one category of Bivariate analysis ○ Able to display one visual on category with a suitable selection of a plot or graph for any category 	<ul style="list-style-type: none"> ○ Able to provide more than one category of Bivariate analysis ○ Able to display visual on categories with a suitable selection of a plot or graph for each category 	<ul style="list-style-type: none"> ○ Able to provide more than one category of univariate analysis ○ Able to display visual on categories with a suitable selection of a plot or graph for each category ○ Provide a clear and concise finding explanation based on the analysis conducted for each category 	2	
5		Create visually appealing and informative data visualizations to convey the identified insights						
		Use suitable chart /graph type	<ul style="list-style-type: none"> □ Provide suitable types of presentations for the targeted audience. □ Use the ggplot2 library to create the graphic. □ The created visualisation is simple and only caters to ONE (1) of 	<ul style="list-style-type: none"> □ Provide suitable types of presentations for the targeted audience. □ Use the ggplot2 library to create the graphic. □ The created visualisation is simple and only caters to 	<ul style="list-style-type: none"> ○ Provide suitable types of presentations for the targeted audience. ○ Use the ggplot2 library to create the graphic. ○ The created visualisation is simple and only caters to THREE (3) of the following: <ul style="list-style-type: none"> ○ Colors and contrast 	<ul style="list-style-type: none"> ○ Provide suitable types of presentations for the targeted audience. ○ Use the ggplot2 library to create the graphic. ○ The created visualisation is simple and only cater to ALL of the following: <ul style="list-style-type: none"> ○ Colors and contrast 	2	

	Convey		<p>the following:</p> <ul style="list-style-type: none"> ○ Colors and contrast ○ size ○ theme ○ scale 	<p>TWO (2) of the following:</p> <ul style="list-style-type: none"> ○ Colors and contrast ○ size ○ theme ○ scale 	<ul style="list-style-type: none"> ○ size ○ theme ○ scale 	<ul style="list-style-type: none"> ○ size ○ theme ○ scale 		
		Customize graph in creative ways	<p><input type="checkbox"/> The visualisation produced covers only ONE(1) of the aspects of below:</p> <ul style="list-style-type: none"> ○ Good quality with clear and accurate ○ Have aesthetic (colors,label & formatting) appeal with visually pleasing <p><input type="checkbox"/> easy for the audience to understand</p> <p><input type="checkbox"/> relevance and effective to communicate the intended message</p>	<p><input type="checkbox"/> The visualisation produced covers only TWO(1) of the aspects below:</p> <ul style="list-style-type: none"> ○ Good quality with clear and accurate ○ Have aesthetic (colors,label & formatting) appeal with visually pleasing ○ easy for the audience to understand ○ relevance and effective to communicate the intended message 	<p><input type="checkbox"/> The visualisation produced covers only THREE (3) of the aspects of below:</p> <ul style="list-style-type: none"> ○ Good quality with clear and accurate ○ Have aesthetic (colors, label & formatting) appeal with visually pleasing ○ easy for the audience to understand ○ relevance and effective to communicate the intended message 	<p><input type="checkbox"/> The visualisation produced covers ALL of the aspects of below:</p> <ul style="list-style-type: none"> ○ Good quality with clear and accurate ○ Have aesthetic (colors, label & formatting) appeal with visually pleasing ○ easy for the audience ○ relevance and effective to communicate the intended message 	2	

		Summarize the findings	<input type="checkbox"/> Able to highlight some of the important findings/result. <input type="checkbox"/> No interpretation of the result provided.	<input type="checkbox"/> Able to highlight some of the important findings. <input type="checkbox"/> Poor interpretation of results provided.	<input type="checkbox"/> Able to highlight all the important findings. <input type="checkbox"/> The interpretation of results is briefly explained but inappropriate, based on the analysis of the data.	<input type="checkbox"/> Able to highlight most of the important findings. <input type="checkbox"/> The interpretation of results is briefly explained and appropriate, based on the analysis of the data.	2	
Total Marks Earned								/50
Total Percentage (40%)								

Table of Contents

PROJECT SCENARIO	2
1.0 Introduction.....	2
2.0 Identification of dataset.....	2
2.1 Objectives.....	3
2.2 Scope.....	3
2.3 Target Audience.....	3
3.0 Dataset.....	3
3.1 Justification of selected dataset.....	3
3.2 Evidence of selected dataset.....	4
3.3 Suitable dataset variables.....	4
4.0 Data Cleaning.....	4
5.0 Data Analysis.....	5
5.1 Univariate Data Analysis (Categorical).....	6
5.2 Univariate Data Analysis (Numerical).....	7
5.3 Bivariate Data Analysis (Categorical vs Numerical).....	8
5.4 Bivariate Data Analysis (Numerical vs. Numerical).....	9
6.0 Data Visualization.....	10
6.1 Average IMDb score.....	10
6.2 Number of movies by genre.....	10
6.3 Average box office earnings per genre.....	11
6.4 Budget vs box office earnings.....	11
6.5 Summarize the Findings.....	12
7.0 References.....	12
8.0 Appendix.....	13

1.0 Introduction

In today's digital era, movies remain one of the most popular forms of entertainment across all generations. The film industry has evolved with changing audience preferences, shifting from traditional cinemas to streaming platforms. By analyzing a dataset from Kaggle.com, I explored key aspects of the movie industry, such as average ratings, popular genres over time, and the financial impact of production budgets on revenue. This analysis helps determine trends in audience preferences, the success factors behind high-grossing films, and whether investing in movies is a profitable venture

2.0 Identification of dataset

a. Objective

The objective of this analysis is to extract meaningful insights from a dataset containing movie information from 1930 to 2016. This analysis explores the average movie ratings over time and identifies the most popular genres during different periods. Additionally, it examines the revenue generated based on genre and investigates whether a high budget is necessary to achieve high revenue.

Scope

Included (state all variables that use to achieve objective)

Genre - Category of the film

Budget - Money spent on making the movie and its publicity

Box Office(revenue) - Money from ticket sales

IMDb score - Score out of 10 that is calculated from the votes of registered IMDb users on the movie

Excluded

Movie - Name of the film

Director - Film director that controls the making of the movie and supervises the actors and technical crew

Running time - The length of time in minutes of the movie

Actor 1, Actor 2 and Actor 3 - Three different actors that participate in the movie

Actors Box Office % - Percentage that reflects how many times the actors managed to at least double the budget in their other films movies

Director Box Office % - Percentage that reflects how many times the director managed to at least double the budget in their other films movies

Oscars and Golden Globes nominations - Amount of nominations that the movie had in the Oscars and the Golden Globes

Oscars and Golden Globes awards - Amount of awards that the movie had in the Oscars and the Golden Globes

Release year - Year when the movie was first released

b. Target audience and benefit of the analysis

- Movie productions and producers: better decisions making leads to avoid financial loss and what drives the movie to success.
- Investor: to get high return on investment. By understanding the revenue and budget pattern this can lead to better decision making.
- Movie critics and influences: play a crucial role to shape public by analyzing the movie, is it worth the hype like overrated or underrated or not?

3.0 Dataset

a. Justification using dataset

The reason this dataset is chosen is because it meets the requirements to do univariate and bivariate analysis, which includes 2 types of variables which is categorical and numerical data types. Another thing is, the datasets provide budget and box office(revenue) perfect to analyze movies profitability. Also, it provides IMDb scores for the public to criticize the movie. Lastly, genre to know what movie production should focus on.

b. Evidence of the data

<https://www.kaggle.com/delfinaoliva/movies>

Appendix1

c. Suitable dataset used (from included)

Genre - Category of the film

Budget - Money spent on making the movie and its publicity

Box Office(revenue) - Money from ticket sales

IMDb score - Score out of 10 that is calculated from the votes of registered IMDb users on the movie

4.0 Data cleaning

```
> df_status(data)
```

	variable	q_zeros	p_zeros	q_na	p_na	q_inf	p_inf	type	unique
1	Movie	0	0.00	0	0.00	0	0	character	3907
2	Director	0	0.00	0	0.00	0	0	character	1760
3	Running.time	0	0.00	0	0.00	0	0	integer	155
4	Actor.1	0	0.00	0	0.00	0	0	character	1591
5	Actor.2	0	0.00	0	0.00	0	0	character	2367
6	Actor.3	0	0.00	0	0.00	0	0	character	2782
7	Genre	0	0.00	0	0.00	0	0	character	14
8	Budget	0	0.00	0	0.00	0	0	integer	374
9	Box.Office	0	0.00	0	0.00	0	0	numeric	994
10	Actors.Box.Office..	321	8.08	0	0.00	0	0	numeric	234
11	Director.Box.Office..	854	21.49	0	0.00	0	0	numeric	39
12	Earnings	68	1.71	0	0.00	0	0	numeric	1244
13	Oscar.and.Golden.Globes.nominations	2839	71.44	0	0.00	0	0	integer	22
14	Oscar.and.Golden.Globes.awards	3504	88.17	3	0.08	0	0	integer	14
15	Release.year	0	0.00	0	0.00	0	0	integer	86
16	IMDb.score	0	0.00	0	0.00	0	0	numeric	77

The image above shows that this data set has two variables that have percentage zero more than 60 percent. Variable 60 percent can lead to bias decisions.

```
12 # 2. Data cleaning remove percentage zero value
13 my_data = df_status(data)
14 arrange(my_data, -p_zeros) %>% select(variable, q_zeros, p_zeros)
15 vars_to_remove <- filter(my_data, p_zeros > 60) %>% .$variable
16 vars_to_remove
17 movies = select(data, -one_of(vars_to_remove))
18 df_status(movies)
```

- In line 13, assign a new variable to be use in line 14
- In line 14, arrange data in descending order by percentage zero. The output should be displayed like the image below:

```
> arrange(my_data, -p_zeros) %>% select(variable, q_zeros, p_zeros)
```

	variable	q_zeros	p_zeros
1	Oscar.and.Golden.Globes.awards	3504	88.17
2	Oscar.and.Golden.Globes.nominations	2839	71.44
3	Director.Box.Office..	854	21.49
4	Actors.Box.Office..	321	8.08
5	Earnings	68	1.71
6	Movie	0	0.00
7	Director	0	0.00
8	Running.time	0	0.00
9	Actor.1	0	0.00
10	Actor.2	0	0.00
11	Actor.3	0	0.00
12	Genre	0	0.00
13	Budget	0	0.00
14	Box.Office	0	0.00
15	Release.year	0	0.00
16	IMDb.score	0	0.00

- In line 15, is the code to remove variables that have percentage zero more than 60%
- Line 16 is to execute the code in line 15. The output should be displayed like the image below

```
> vars_to_remove <- filter(my_data, p_zeros > 60) %>% .$variable
> vars_to_remove
[1] "Oscar.and.Golden.Globes.nominations" "Oscar.and.Golden.Globes.awards"
```

- In line 17, assign new variable to keeping all columns except the ones present in "vars_to_remove".
- Line 18 is to see variable. The output should be displayed like the image below:

```
> movies = select(data, -one_of(vars_to_remove))
> df_status(movies)
```

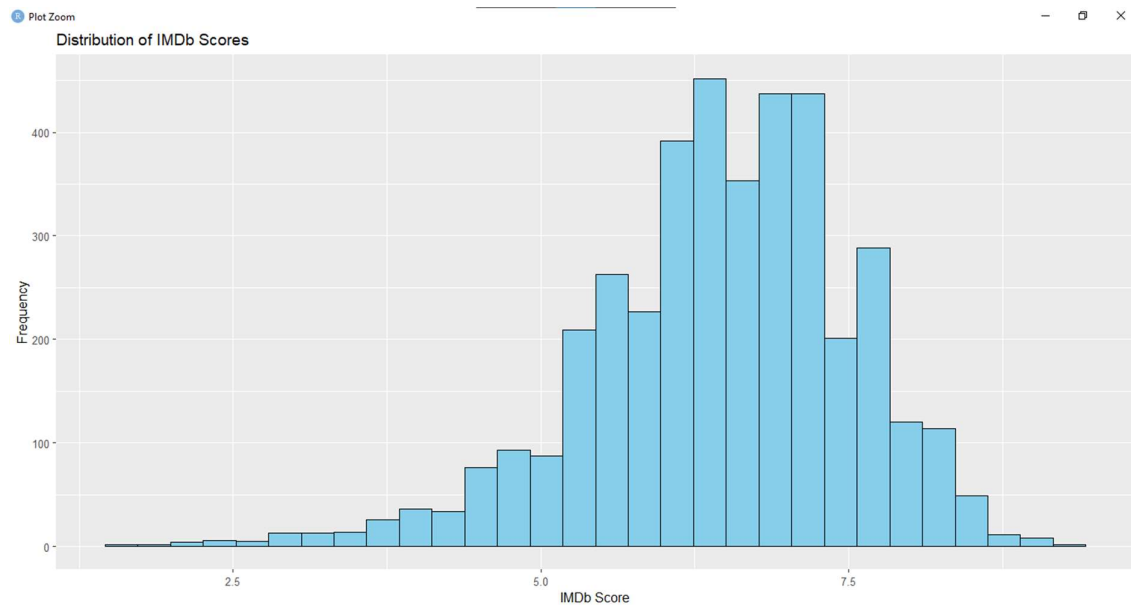
	variable	q_zeros	p_zeros	q_na	p_na	q_inf	p_inf	type	unique
1	Movie	0	0.00	0	0	0	0	character	3907
2	Director	0	0.00	0	0	0	0	character	1760
3	Running.time	0	0.00	0	0	0	0	integer	155
4	Actor.1	0	0.00	0	0	0	0	character	1591
5	Actor.2	0	0.00	0	0	0	0	character	2367
6	Actor.3	0	0.00	0	0	0	0	character	2782
7	Genre	0	0.00	0	0	0	0	character	14
8	Budget	0	0.00	0	0	0	0	integer	374
9	Box.Office	0	0.00	0	0	0	0	numeric	994
10	Actors.Box.Office..	321	8.08	0	0	0	0	numeric	234
11	Director.Box.Office..	854	21.49	0	0	0	0	numeric	39
12	Earnings	68	1.71	0	0	0	0	numeric	1244
13	Release.year	0	0.00	0	0	0	0	integer	86
14	IMDb.score	0	0.00	0	0	0	0	numeric	77

5.0 EDA Data analysis

Univariate analysis

Numerical analysis: average IMDb score

```
25 ggplot(movies, aes(x = IMDb.score)) +  
26   geom_histogram(fill = "skyblue", color = "black") +  
27   labs(title="Distribution of IMDb Scores", x="IMDb Score", y="Frequency" )  
28  
29 summary(movies)
```



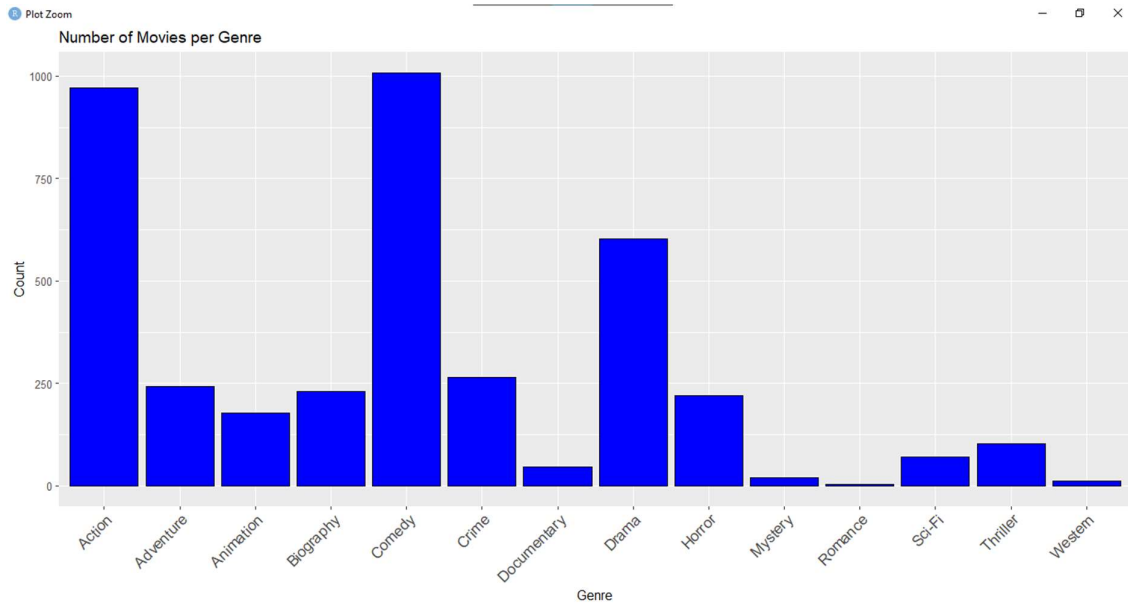
by using summary () function:

- The lowest rating is 1.6/10
- Highest rating is 9.3/10
- And the average rating is 6.48/10
- In conclusion, most of the public audience give moderate rating more than high and low.

```
IMDb.score  
Min.   :1.600  
1st Qu.:5.900  
Median :6.600  
Mean   :6.468  
3rd Qu.:7.200  
Max.   :9.300
```

Categorical analysis: Number of movies by genre

```
31 # categorical analysis: Number of movies by genre
32 ggplot(movies, aes(x = Genre)) +
33   geom_bar(fill = "blue", color = "black") +
34   theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 12)) +
35   labs(title = "Number of Movies per Genre", x="Genre", y="Count")
```



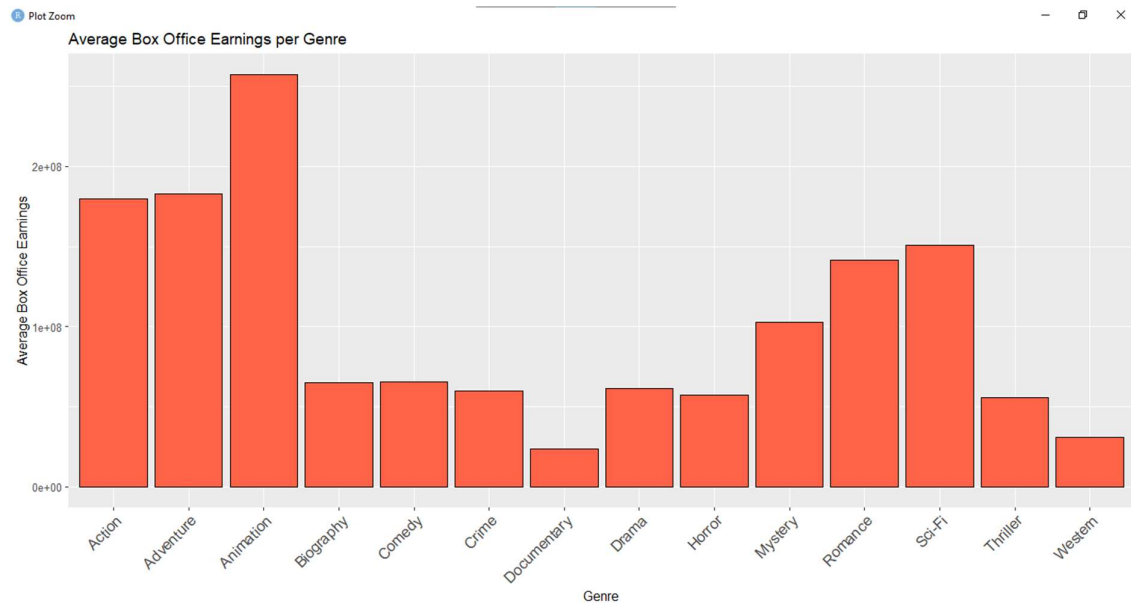
Based on the bar image:

- Comedy genre is the most produced movies
- The second goes to action
- The third is Drama
- And the least produce goes to romance
- In conclusion, movies production tend to produce more comedy and action genre.

Bivariate analysis

categorical analysis vs numerical analysis: average box office earnings per genre

```
38 # categorical vs numerical: average box office per genre
39 movies %>%
40   group_by(Genre) %>%
41   summarise(Average_Box_Office = mean(Box.Office)) %>%
42   ggplot(aes(x = Genre, y = Average_Box_Office)) +
43   geom_bar(stat = "identity", fill = "tomato", color = "black") +
44   theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 12)) +
45   labs(title="Average Box Office Earnings per Genre", x="Genre",
46        y="Average Box Office Earnings")
```



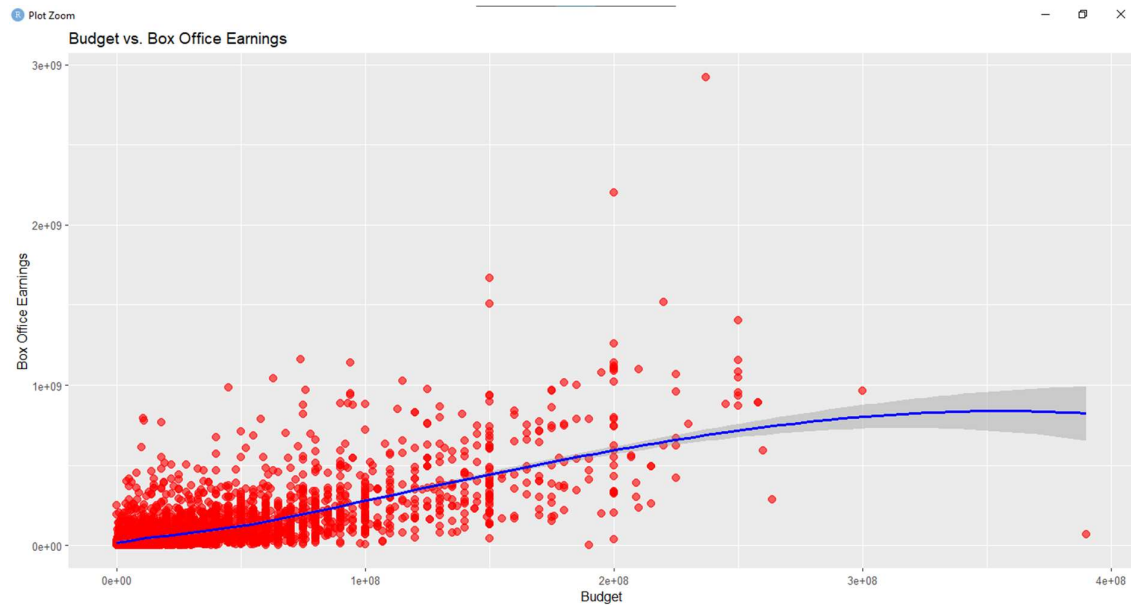
Based on the above image:

- Animation performs better than any other genre
- Meanwhile average box office earnings go to adventure and action
- Least, box office earnings(revenue) is documentary
- From summary() function:

```
Box.Office
Min.   :5.000e+04
1st Qu.:1.200e+07
Median :4.300e+07
Mean   :1.087e+08
3rd Qu.:1.250e+08
Max.   :2.923e+09
```


numerical analysis vs numerical analysis

```
48 # numerical vs numerical:
49 ## Numerical vs. Numerical: Relationship between Budget and Box Office Earnings
50 ggplot(movies, aes(x = Budget, y = Box.Office)) +
51   geom_point(color = "red", alpha = 0.6, size = 3) +
52   geom_smooth(method = "loess", color = "blue", size = 1.2) +
53   ggtitle("Budget vs. Box Office Earnings") +
54   xlab("Budget") +
55   ylab("Box Office Earnings")
```



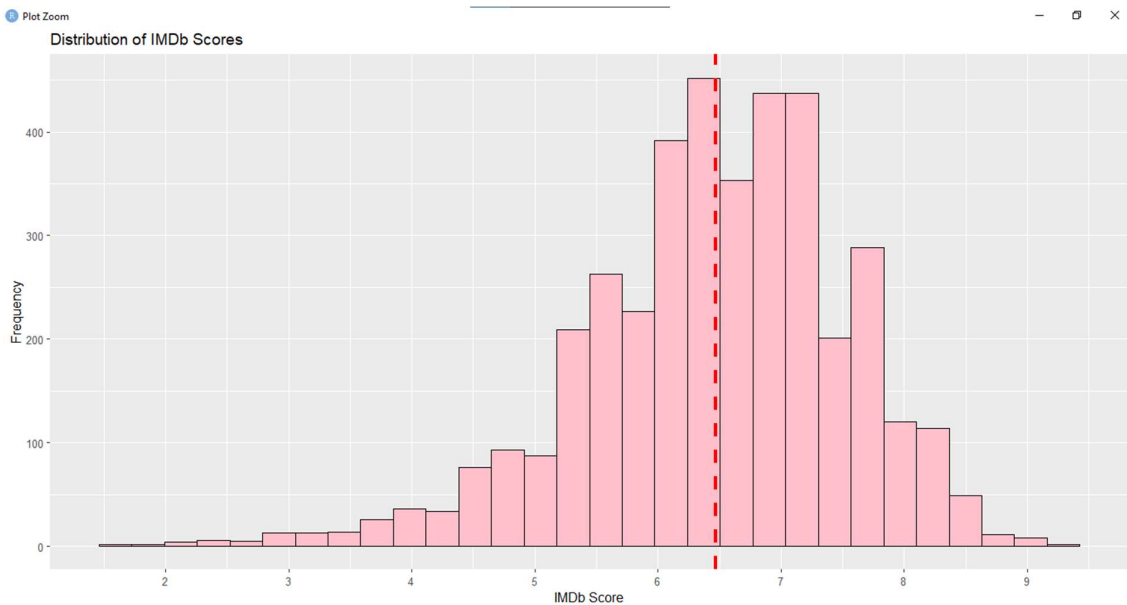
Based on the image above:

- The scatter plot shows a positive correlation, higher budget led to higher box office earnings(revenue)
- But it is not guaranteed that higher budgets lead to higher box office earnings(revenue) because some movies with higher budgets perform poorly
- From summary function ()

Budget		Box.Office	
Min.	: 1100	Min.	: 5.000e+04
1st Qu.	: 9000000	1st Qu.	: 1.200e+07
Median	: 22000000	Median	: 4.300e+07
Mean	: 36906392	Mean	: 1.087e+08
3rd Qu.	: 50000000	3rd Qu.	: 1.250e+08
Max.	: 390000000	Max.	: 2.923e+09

6.0 Data visualization

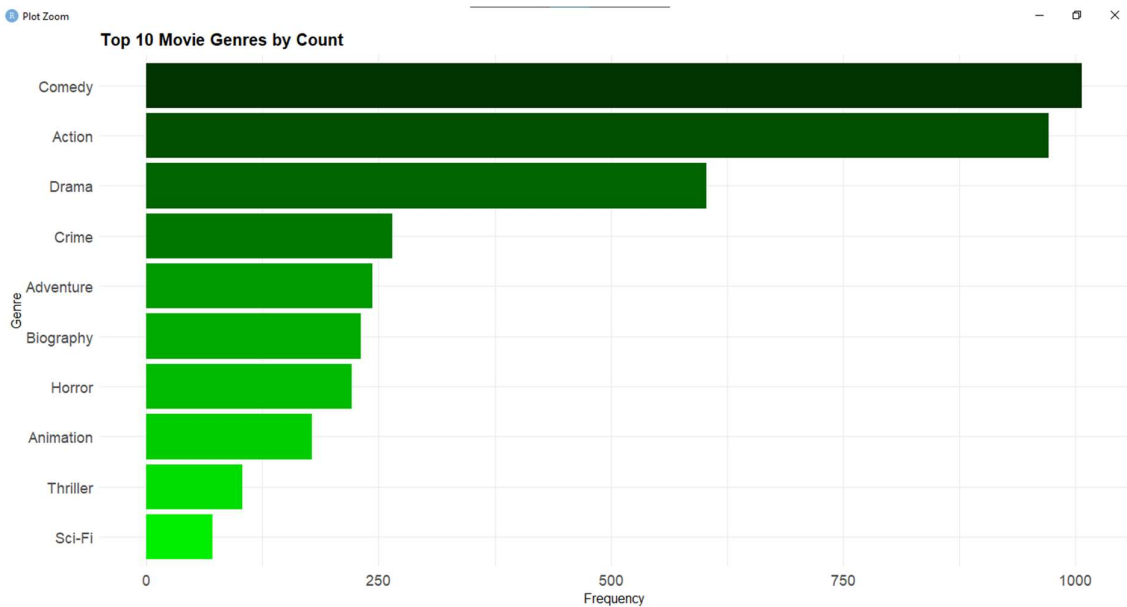
Average IMDb score



Most movies receive mid-range scores rather than being considered masterpieces or complete failures.

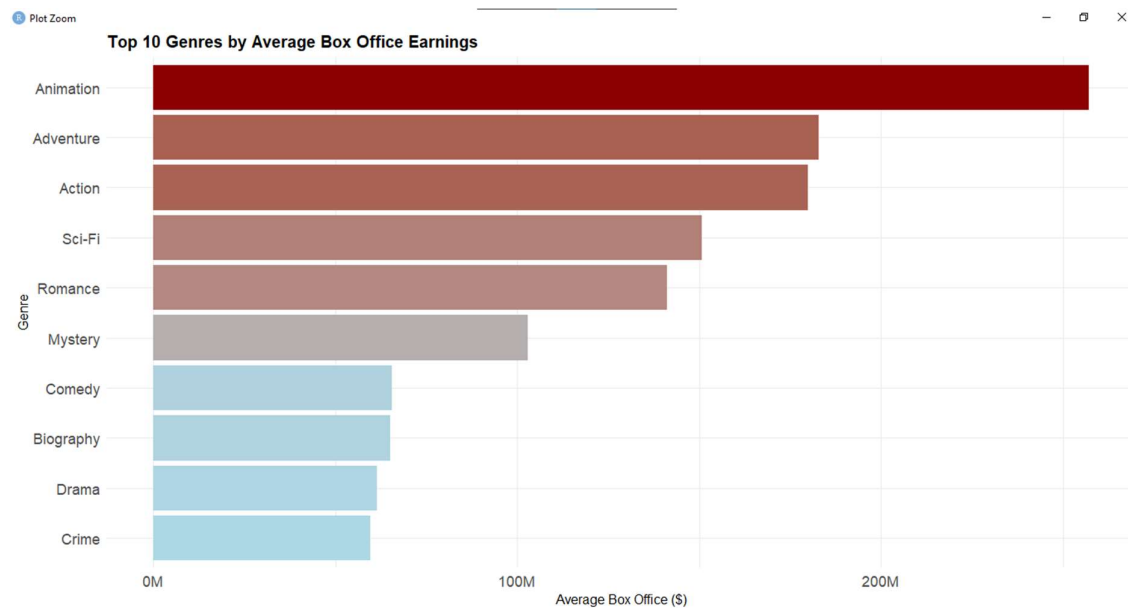
The data suggests that getting an IMDb score above 8 is rare, but there are also fewer movies rated below 3.

Number of movies by genre



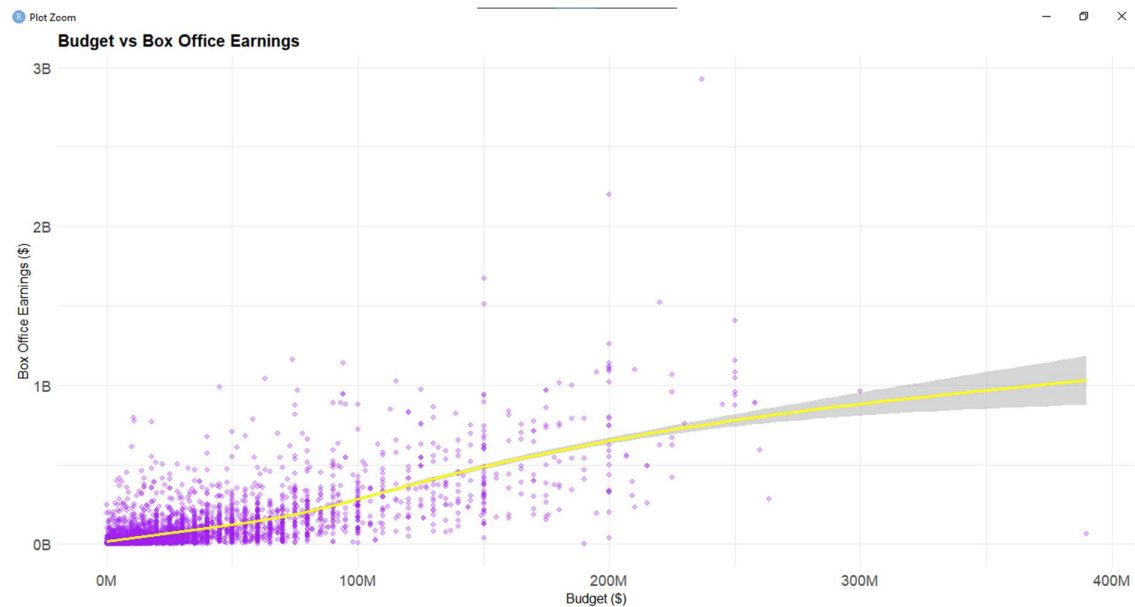
The top three genres (Comedy, Action, and Drama) are significantly more common than other genres, indicating a strong preference for these types of films in movie production.

Average box office earnings per genre



Animation has the highest average box office earnings, suggesting that animated movies are highly profitable. This could be due to their broad appeal to both children and adults, strong franchise potential, and high merchandising revenue.

Budget vs box office earnings



The yellow trend line shows a general upward trend, indicating that movies with higher budgets tend to generate higher box office earnings.

However, the relationship is not perfectly linear, meaning a higher budget does not always guarantee proportionally higher earnings.

6.5 Summarize findings

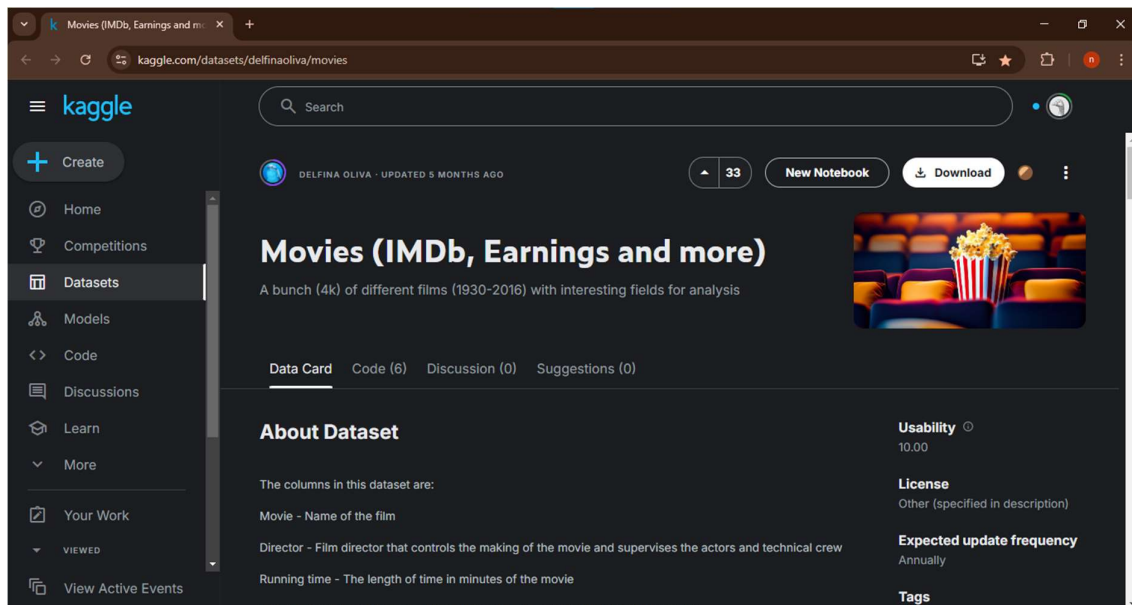
Most movies receive mid-range scores, with very few rated above 8 or below 3.

Comedy, Action, and Drama are the most common genres, indicating a strong production preference. Animation has the highest average earnings, likely due to its broad appeal, franchise potential, and merchandising opportunities. Higher budgets generally lead to higher box office earnings, but the relationship is not strictly linear, meaning a large budget does not always guarantee success.

7.0 References

1. GeeksforGeeks. (2024, October 3). *Display an axis value in millions in ggplot using R*. GeeksforGeeks. https://www.geeksforgeeks.org/display-an-axis-value-in-millions-in-ggplot-using-r/?ref=header_outind
2. GeeksforGeeks. (2023, December 20). *Histogram in R using ggplot2*. GeeksforGeeks. <https://www.geeksforgeeks.org/histogram-in-r-using-ggplot2/>
3. Pedersen, H. W. D. N. a. T. L. (n.d.). *ggplot2: Elegant Graphics for Data Analysis (3e) - 11 Colour scales and legends*. <https://ggplot2-book.org/scales-colour>
4. Mph, S. P. S. I., & Mph, S. P. S. I. (2024, April 24). *A practical guide to selecting top N values by group in R | R-Bloggers*. R-bloggers. https://www.r-bloggers.com/2024/04/a-practical-guide-to-selecting-top-n-values-by-group-in-r/#google_vignette

8.0 Appendix



Appendix 1

Display an axis value in million: X +

← → ↻

geeksforgeeks.org/display-an-axis-value-in-millions-in-ggplot-using-r/?ref=header_outind

☆ 📄 📑

Sign In

Courses ▾

Tutorials ▾

DSA ▾

Data Science ▾

Web Tech ▾

geeksforgeeks

🔍 🔄 📄 📑

Sign In

Data Visualization

Statistics in R

Machine Learning in R

Data Science in R

Packages in R

Data Types

String

Array

Vector

Lists

Matrices

Oops in R

Next Article:

Display Only Integer Values on ggplot2 Axis in R →

0 Ads with

GeeksforGeeks

Upgrade Now

Display an axis value in millions in ggplot using R

Last Updated : 03 Oct, 2024

🔗 💬 ✎ ⋮

When working with large numerical data in R using ggplot2, axis values can sometimes become cumbersome and hard to read, especially when the numbers are in the millions or billions. Displaying these values in a more readable format, such as in millions (e.g., 1,000,000 as 1M), enhances the clarity and presentation of your plot. This article will guide you through the steps of formatting your axis values to display in millions using ggplot2 in [R Programming Language](#).

Why Format Axis Values in Millions?

Playing large numbers in full form can make your plots difficult to interpret. For example:

- 1,000,000 can be more intuitively represented as 1M.
- Reduces clutter and ensures better readability.
- Provides a more professional and polished appearance.

Setting Up the Environment

Appendix 2

Histogram in R using ggplot2 X +

← → ↻

geeksforgeeks.org/histogram-in-r-using-ggplot2/

☆ 📄 📑

Sign In

Courses ▾

Tutorials ▾

DSA ▾

Data Science ▾

Web Tech ▾

geeksforgeeks

🔍 🔄 📄 📑

Sign In

Data Visualization

Statistics in R

Machine Learning in R

Data Science in R

Packages in R

Data Types

String

Array

Vector

Lists

Matrices

Oops in R

Next Article:

Box plot in R using ggplot2 →

0 Ads with

GeeksforGeeks

Upgrade Now

Histogram in R using ggplot2

Last Updated : 20 Dec, 2023

🔗 💬 ✎ ⋮

ggplot2 is an R Package that is dedicated to Data visualization. ggplot2 Package Improve the quality and the beauty (aesthetics) of the graph. By Using ggplot2 we can make almost every kind of graph In [RStudio](#).

What is Histogram?

A [histogram](#) is an approximate representation of the distribution of numerical data. In a histogram, each bar groups numbers into ranges. Taller bars show that more data falls in that range. A histogram displays the shape and spread of continuous sample data.

Basic ggplot2 Histogram in R

Histograms roughly give us an idea about the probability distribution of a given variable by depicting the frequencies of observations occurring in certain ranges of values. Histograms are used to show distributions of a given variable while bar charts are used to compare variables. Histograms plot quantitative data with ranges of the data grouped into intervals while bar charts plot categorical data.

Appendix 3

