



UTP
UNIVERSITI TEKNOLOGI PETRONAS

TEB2164

INTRODUCTION TO DATA SCIENCE

FINAL PROJECT

“ANALYSING PERFORMANCE OF VIDEO GAME SALES BASED ON REGION”

Prepared for:

Ts Ahmad Izuddin B Zainal Abidin

Dr. Said Jadid Abdulkadir

Prepared by:

NAME	ID
ENGKU MUHAMMAD ZAHIMI BIN CHE ENGKU IBRAHIM	16003930
IQBAL NURIL ANWAR BIN TORMIZI	17004961
KHAIRUL ILHAM	17004942
AHMAD ALIF BIN MOHAMAD AFIFI	17004106
NATHANEEL NAAVIN AISU	16004681
TIA FARISHA BINTI SHAHARUL MIZA	17002588

TABLE OF CONTENTS

NO.	CONTENT	PAGE
1.0	EXECUTIVE SUMMARY	2
2.0	PROBLEM DESCRIPTION 2.1 Business Goal 2.2 Technical Goal	3 - 4
3.0	DATA DESCRIPTION 3.1 Data Quantity 3.2 Data Quality	5
4.0	DATA PREPARATION 4.1 Selecting Data 4.2 Cleaning Data 4.3 Visualization	6
5.0	SOLUTION	7 - 8
6.0	CONCLUSION	9
7.0	APPENDIX	10 - 20

1.0 EXECUTIVE SUMMARY

Topic :Analysing Performance of Video Games Based On Region

Category :Sales

Sources :

- Original Owners : Maksim Klimentyev
- Donor of database : Kaggle.com

Problem Description : Annual sales of video games from various genres, developers and even various years based on each country. In this research, our task with the analysis is to examine the success of the year-to-year sales of video games.1.54mb/1.01mb

Data Description : The annual performance of each developer around the country is obtained from Kaggle, while the year begins from 1985 to 2020. The size of the raw database is 1.54 MB, which includes 16 row numbers and 16720 column numbers.

Data Preparation : As there are so many redundant quantities of data in the raw data collection, we continue our analysis by using data preparation measures such as data cleaning, removing and replacing null values for standardized purposes. The technique for preparing the data is CRISP-DM.

Solution and Finding : The elements that can be obtained from the data set include Global Sales, Game Produced, and Distribution of Release Year by Year. In the final section, the identification of the Top values in the data set and game count by Developer is created according to the results of "Platforms with most dataset games."

2.0 PROBLEM DESCRIPTION

2.1 Business Goal

The Data Science Concepts are designed to incorporate mathematics, programming, and business analysis to generate consistent trends within the data itself. Knowledge is power in business, and data is the fuel that generates the power. It is extremely necessary to be able to harness the power of this data through data science. Through drawing numbers and statistics from **data analysis**, an organization will create statistical models to simulate a variety of choices. As a result, companies can learn which path to take to the best possible outcome of how to better market their games on the basis of derived data analysis.

As an example, valuable trends can help businesses accomplish their revenues appropriately. Many raw data can be obtained from the sales of video games produced annually by the developers. This can be used to adapt or customize programs and goods to different markets, the producer can concentrate more on what kind of games to develop and distribute according to the area they want to increase their sales in. It means that organizations may tailor games and other products to cater to specific audiences, for example, by knowing which genres have been most played in one area, will help the company develop new promotions or deals for audiences that may not have been available before.

2.2 Technical Goal

The main aim of this project is to evaluate and determine the annual revenue of each developer for each given year. Nevertheless, in the process of analyzing these results, there have been few difficulties in achieving the objectives of the project.

Below is a list of problems that we encountered :-

- **Datasets are incomplete (null values)**

Unfilled data tables (null values) causes errors when visualising said databases. The graphs would become inconsistent and drastically askew in various phases. This can be seen when null values are kept and the graph spikes in random phases or that the average is unrealistically lower than predicted.

- **Datasets that are irrelevant.**

Raw databases had many categories of data, some of these categories are unrelated to our project goals. Including said data categories would be redundant as they do not provide the necessary relations that we aim for in our goals, hence the data must be removed to streamline our data analysis.

- **Wrong Data Visualisation Method.**

Many methods of visualization can misrepresent data, as can be clearly seen in pie chart visualization where the partition is a misrepresentation of facts and the consequence is that an incorrect conclusion can be drawn from it. It can happen if it unintentionally obstructs the correct understanding of the data or is mistakenly due to a lack of familiarity with the graphics program, a misinterpretation of the data, or because the data can not be properly conveyed.

3.0 DATA DESCRIPTION

3.1 Data Quantity

The format of the data is Comma Separated Values (CSV) format.

The method used to capture the data is Web Scraping. Web scraping is a technique used to extract a large amount of the data from websites where data is exported to be extracted and saved to a local file or to a database in a table/spreadsheet. This data is scraped from Metacritic.com, a website where critics leave ratings on Movies and games.

The size of the database is 1.54MB which is the number of rows is 17 and the number of columns is 16720.

3.2 Data Quality

This data is about the performance of games sales performance in different regions around the world. This data shows the figure of year released, sales in Europe, Japan, North America and Global region, the platform developed, the publisher, the developer and critics score. This data is mostly relevant to conclude key sales performances for further analysis. Our main purpose for using this data is to identify the performance of games based on regions and platforms. The data types are present in numbering and text. Data will be computed and visualized to give a better picture of the performance and where each developer/publisher can focus on.

4.0 DATA PREPARATION

The data that was collected could not be used directly as there are flaws that must be corrected. In order for the dataset to be usable for our project, one that provides accurate and consistent data, certain steps must be taken to clean and prepare the database. Selection and cleansing of data was done using R.studio. We then saved the new dataset and used Power BI for our visualisation as we found it easier and simpler to use. Mistakes in coding would affect our overall analysis, hence to minimize that, we opted for Power Bi.

4.1 Selecting Data

The database that we used consists of several categories. However, not all categories are relevant to our goal. Because our main focus is on the global sales of video games, we decided that ratings and critic score is irrelevant due to those being of a qualitative nature and does not relate with the quantity aspect. Hence, we removed those two categories in order to streamline our database and provide a more consistent dataset. Removing those two categories also reduced the size of our database for further efficiency.

4.2 Cleaning Data

Cleaning data is necessary to provide an accurate data analysis. Some of the datasets had null values. If no cleansing was done, the analysis would provide inconsistent and flawed results as certain values would manipulate the results and cause inaccurate analysis. The dataset had 6815 rows of data that was incomplete. From the original 16720 data, 9905 data was left after we removed those with null values. Even though a large number of data was removed, this cleansed dataset will provide a more consistent relationship when processed visually.

4.3 Visualisation

Data in its table form is hard to analyse and draw relations from. Because of this, we need to convert this into something whereby patterns can be concluded from. With the available data, we picked two categories and produced graphs that showed the patterns between them. We took the Game genre and compared it to Country sales to find out which games were more popular regionally. We did this with several pairs of categories to provide as much relevant information from the data that we had.

5.0 SOLUTION

In this sales analysis report, it shows an overview of performance of video game sales based on region. It shows the different trends happening in the production of video games over the year until today. As we can conclude that in early 1985 until 2000, there was very low production of game consoles. But the starting peak of rising in the production of console games comes in the year of 2008 where many people were introduced to console games. As for today, not so many developers produce console games because people nowadays are more likely to use social media and do other activities compared to playing games.

Next, the sales for countries like Europe (EU), North America (NA), and Japan (JP) can be shown in this report. As we can see in each country, Nintendo's developer has the highest sales in EU, NA and even JP and then followed by games developed from EU and NA itself. To conclude, from Japanese sales by developer, it is hard for developers from EU and NA to infiltrate the market in Japan as the people in that regional area are more interested in playing games from the developer of its country. For EU and NA, maybe the developer can focus more on selling their console games in other parts of the world except Japan.

Besides, in this report, it also shows the sales for each country by genre. It helps developers to decide which type of genre they want to develop for that particular region. For example, EU and NA have the most sales for the action genre. But for JP, the role-playing genre has the highest sales. By doing this, it will help to reduce the cost of developing unused console games by using target segmentation.

Sales reports help to identify possible new business markets where results may be enhanced. So, this report's aim is to make it crystal clear to the sales team as to what is important and what they should be focusing on.

As stated, everyday sales is different and there are many things that are simply uncontrollable. Now we will concentrate on some of the problems or concerns that you can monitor by looking at the report on a regular basis.

Identified problem:

• What can we understand from this data set?

- 1.The changing trend in using technology and video games.
2. Different preferred developers game based on region.
- 3.Different preferred genre based on region.

6.0 CONCLUSION

In conclusion, video game sales can help keep the above parameters in mind and a well-planned game launch could be successful. Knowing the trends in genre, regions preferences, popularity of games and presently used platforms can help in predicting how successful the game could be. As for the data preparation, we defined the data using a specific method which is CRISP- DM. Here, we gather the data, discover and assess data, cleanse and validate the data, transform and enrich data and finally restore the data. In short, we do think that if we concentrate on certain issues or difficulties we can increase the performance of video game sales.

7.0 APPENDIX

PROJECT SOURCE CODE: <https://github.com/engkuzahimi/introds-jan20>

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
	Name	Platform	Year of Release	Genre	Publisher	NA Sales	EU Sales	JP Sales	Other Sales	Global Sales	Critic Score	Critic Count	User Score	User Count	Developer	Rating		
1	Wii Sports	Wii	2006	Sports	Nintendo	41.36	28.96	3.77	8.45	82.53	76	51	8	322	Nintendo	E		
2	Super Mario	NES	1985	Platform	Nintendo	29.08	3.58	6.81	0.77	40.24								
3	Mario Kart	Wii	2008	Racing	Nintendo	15.68	12.76	3.79	3.29	35.52	82	73	8.3	709	Nintendo	E		
4	Wii Sports	Wii	2009	Sports	Nintendo	15.61	10.93	3.28	2.95	32.77	80	73	8	192	Nintendo	E		
5	Pokemon	GB	1996	Role-Playing	Nintendo	11.27	8.89	10.22	1	31.37								
6	Tetris	GB	1989	Puzzle	Nintendo	23.2	2.26	4.22	0.58	30.26								
7	New Super	DS	2006	Platform	Nintendo	11.28	9.14	6.5	2.88	29.8	89	65	8.5	431	Nintendo	E		
8	Wii Play	Wii	2006	Misc	Nintendo	13.96	9.18	2.93	2.84	28.92	58	41	6.6	129	Nintendo	E		
9	New Super	Wii	2009	Platform	Nintendo	14.44	6.94	4.7	2.24	28.32	87	80	8.4	594	Nintendo	E		
10	Duck Hunt	NES	1984	Shooter	Nintendo	26.93	0.63	0.28	0.47	28.31								
11	Nintendogs	DS	2005	Simulation	Nintendo	9.05	10.95	1.93	2.74	24.67								
12	Mario Kart	DS	2005	Racing	Nintendo	9.71	7.47	4.13	1.9	23.21	91	64	8.6	464	Nintendo	E		
13	Pokemon	GB	1999	Role-Playing	Nintendo	9	6.18	7.2	0.71	23.1								
14	Wii Fit	Wii	2007	Sports	Nintendo	8.92	8.03	3.6	2.15	22.7	80	63	7.7	146	Nintendo	E		
15	Kinect	Adi X360	2010	Misc	Microsoft	15	4.89	0.24	1.69	21.81	61	45	6.3	106	Good Scien	E		
16	Wii Fit	Plu Wii	2009	Sports	Nintendo	9.01	8.49	2.53	1.77	21.79	80	33	7.4	52	Nintendo	E		
17	Grand Theft	PS3	2013	Action	Take-Two	7.02	9.09	0.98	3.96	21.04	97	50	8.2	3994	Rockstar Nc	M		
18	Grand Theft	PS2	2004	Action	Take-Two	9.43	0.4	0.41	10.57	20.81	95	80	9	1588	Rockstar Nc	M		
19	Super Mario	SNES	1990	Platform	Nintendo	12.78	3.75	3.54	0.55	20.61								
20	Brain Age	DS	2005	Misc	Nintendo	4.74	9.2	4.16	2.04	20.15	77	58	7.9	50	Nintendo	E		
21	Pokemon	DS	2006	Role-Playing	Nintendo	6.38	4.46	6.04	1.36	18.25								
22	Super Mario	GB	1989	Platform	Nintendo	10.83	2.71	4.18	0.42	18.14								
23																		

Figure 1.0 Original Data Before Cleansing

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
	Name	Platform	Year of Release	Genre	Publisher	NA Sales	EU Sales	JP Sales	Other Sales	Global Sales	Critic Score	Critic Count	User Score	User Count	Developer	Rating		
178	Assassin's	X360	2012	Action	Ubisoft	3.13	1.69	0.03	0.44	5.29	84	61	6.7	1196	Ubisoft	M		
179	Final Fant	PS2	2003	Role-Playing	Electronic	1.92	1.08	2.11	0.17	5.29	85	45	6.6	400	SquareSoft	T		
180	Donkey K	N64	1999	Platform	Nintendo	3.33	0.79	1.09	0.06	5.27								
181	Call of Du	XOne	2014	Shooter	Activision	3.22	1.55	0.01	0.48	5.27	81	53	5.4	898	Sledgeham	M		
182	Minecraft	PS3	2014	Misc	Sony Com	2.03	2.37	0	0.87	5.26								
183	Assassin's	X360	2012	Action	Ubisoft	3.11	1.55	0.08	0.51	5.21								
184	Tomb Raic	PS	2003	Adventure	Eidos Inte	2.3	2.46	0.2	0.28	5.24								
185	Madden N	PS2	2003	Sports	Electronic	4.26	0.26	0.01	0.71	5.23								
186	Tomodach	3DS	2013	Simulation	Nintendo	0.97	2.11	1.9	0.24	5.23								
187	New Super	WiiU	2012	Platform	Nintendo	2.3	1.34	1.27	0.32	5.22								
188	Dragon Q	PS2	2004	Role-Playing	Square En	0.65	0.75	3.61	0.2	5.21								
189	Super Mar	GBA	2003	Platform	Nintendo	2.93	1.25	0.83	0.2	5.1								
190	Professor	DS	2007	Puzzle	Nintendo	1.21	2.43	1.03	0.52	5.19								
191	Super Mar	GB	1994	Platform	Nintendo	2.49	0.98	1.57	0.15	5.19								
192	FIFA Socce	X360	2012	Action	Electronic	1.09	3.47	0.03	0.57	5.16	90	48	6.1	403	Electronic A	E		
193	Donkey K	SNES	1995	Platform	Nintendo	2.1	0.74	2.2	0.11	5.15								
194	Diablo III	PC	2012	Role-Playing	Activision	2.44	2.16	0	0.54	5.14	88	86	4	9629	Blizzard Ent	M		
195	Medal of	IP	2003	Shooter	Electronic	1.98	2.23	0.13	0.8	5.13	68	30	7.6	100	EA LA	T		
196	Kirby's Dr	GB	1992	Platform	Nintendo	2.71	0.61	1.7	0.11	5.13								
197	Microsoft	PC	1996	Simulation	Microsoft	3.22	1.69	0	0.2	5.12								
198	Guitar Her	PS2	2006	Misc	RedOctan	3.81	0.63	0	0.68	5.12	92	69	8.5	112	Harmonix A	T		
199	Fable III	X360	2010	Role-Playing	Microsoft	3.59	1.08	0.05	0.38	5.1	80	88	6.5	604	Lionhead S	M		
200	Mario & S	DS	2008	Sports	Sega	1.63	2.45	0.44	0.57	5.09								

Figure 2.0 Data Problem

```
server <- function(input, output) {
  # Reading the dataset into data and then omit rows with empty cells
  data <- na.omit(read.csv('s.csv', stringsAsFactors = FALSE, na.strings=c("N/A")))

  #
  # # Removing rows with any empty cells inside it
  data <- data[!(is.na(data$Name) | data$Name==""), ]
  data <- data[!(is.na(data$Platform) | data$Platform==""), ]
  data <- data[!(is.na(data$Year) | data$Year==""), ]
  data <- data[!(is.na(data$Genre) | data$Genre==""), ]
  data <- data[!(is.na(data$Publisher) | data$Publisher==""), ]
  data <- data[!(is.na(data$NA_Sales) | data$NA_Sales==""), ]
  data <- data[!(is.na(data$EU_Sales) | data$EU_Sales=="") , ]
  data <- data[!(is.na(data$JP_Sales) | data$JP_Sales==""), ]
  data <- data[!(is.na(data$Other_Sales) | data$Other_Sales==""), ]
  data <- data[!(is.na(data$Global_Sales) | data$Global_Sales==""), ]
  data <- data[!(is.na(data$Developer) | data$Developer==""), ]

  #reexport cleansed data
  write.csv(data, "C:\\Users\\Jim\\Desktop\\CleanedData.csv", row.names = TRUE)
}
```

Figure 3.0 Cleansing Code

	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales	Developer
1	Wii Sports	Wii	2006	Sports	Nintendo	41.36	28.96	3.77	8.45	82.53	Nintendo
3	Mario Kart Wii	Wii	2008	Racing	Nintendo	15.68	12.76	3.79	3.29	35.52	Nintendo
4	Wii Sports Resort	Wii	2009	Sports	Nintendo	15.61	10.93	3.28	2.95	32.77	Nintendo
7	New Super Mario Bros.	DS	2006	Platform	Nintendo	11.28	9.14	6.50	2.88	29.80	Nintendo
8	Wii Play	Wii	2006	Misc	Nintendo	13.96	9.18	2.93	2.84	28.92	Nintendo
9	New Super Mario Bros. Wii	Wii	2009	Platform	Nintendo	14.44	6.94	4.70	2.24	28.32	Nintendo
12	Mario Kart DS	DS	2005	Racing	Nintendo	9.71	7.47	4.13	1.90	23.21	Nintendo
14	Wii Fit	Wii	2007	Sports	Nintendo	8.92	8.03	3.60	2.15	22.70	Nintendo
15	Kinect Adventures!	X360	2010	Misc	Microsoft Game Studios	15.00	4.89	0.24	1.69	21.81	Good Science Studio

Figure 4.0 Data After Cleaning

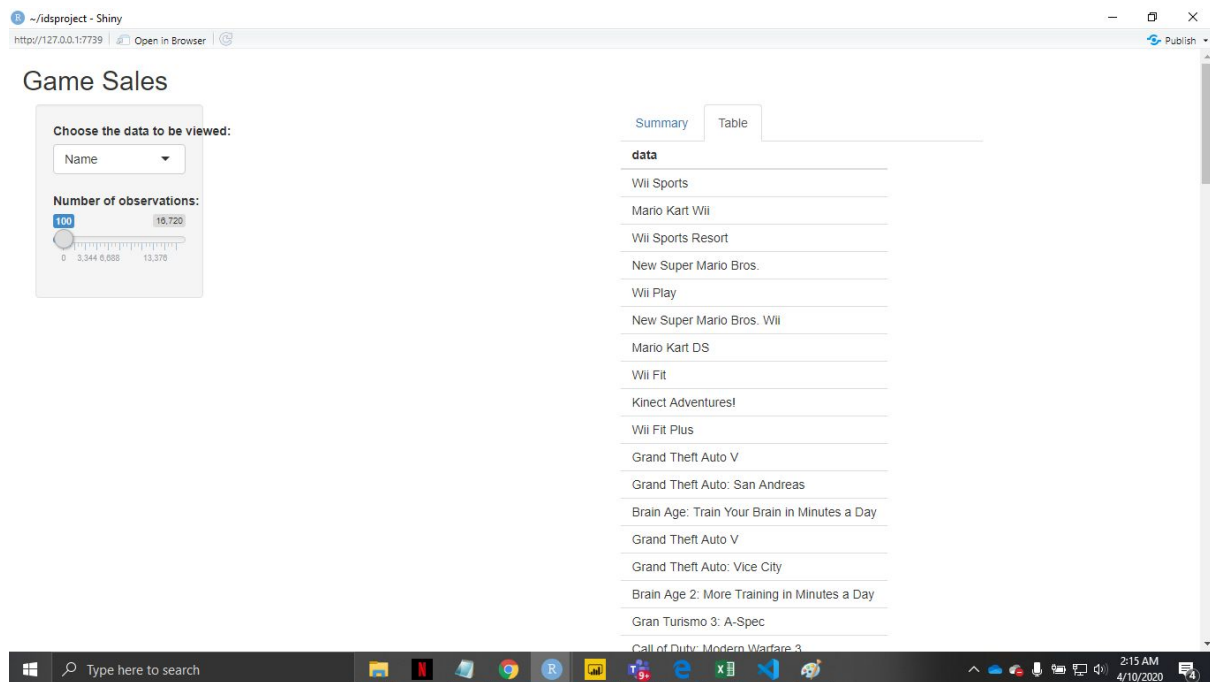


Figure 5.0 Data Presentation Using Shiny - Table View

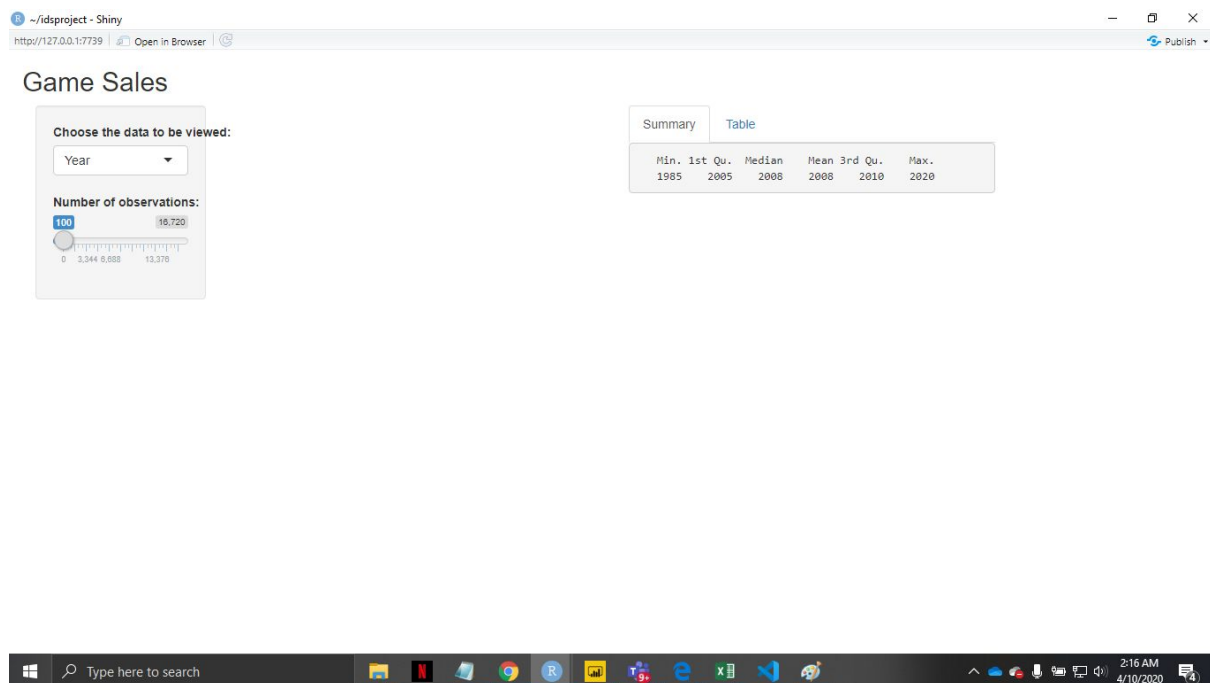


Figure 6.0 Data Presentation Using Shiny - Summary

```

1 ui <- fluidPage(
2   titlePanel('Game Sales'),
3   splitLayout(
4     sidebarPanel(
5       selectInput(
6         inputId = 'dataset',
7         label = 'Choose the data to be viewed:',
8         choices = c('Name', 'Platform', 'Year', 'Genre', 'Publisher', 'SalesNorthAmerica', 'SalesEurope', 'SalesJapan', 'SalesOther', 'SalesGlobal'),
9       ),
10      sliderInput(
11        inputId = 'obs',
12        label = 'Number of observations: ',
13        min = 0,
14        max = 16720,
15        value = 100
16      )
17    ),
18   mainPanel(
19     tabsetPanel(type = "tabs",
20       #tabPanel("Plot", plotOutput("plot")),
21       tabPanel("Summary", verbatimTextOutput("summary")),
22       tabPanel("Table", tableOutput("view"))
23     )
24   )
25 )
26 )
27 )

```

Figure 7.0 Shiny Code - UI

```

1 # Getting the input from user on which column of table to display
2 dataInput <- reactive({
3   switch(input$dataset,
4     'Name' = data$Name,
5     'Platform' = data$Platform,
6     'Year' = data$Year,
7     'Genre' = data$Genre,
8     'Publisher' = data$Publisher,
9     'SalesNorthAmerica' = data$NA_Sales,
10    'SalesEurope' = data$EU_Sales,
11    'SalesJapan' = data$JP_Sales,
12    'SalesOther' = data$Other_Sales,
13    'SalesGlobal' = data$Global_Sales,
14    'Developer' = data$Developer
15  )
16 })
17
18 output$summary <- renderPrint({
19   dataColumn <- dataInput()
20   summary(dataColumn)
21 })
22
23 output$view <- renderTable({
24   dataColumn <- dataInput()
25   head(dataColumn, n = input$obs)
26 })
27

```

Figure 8.0 Shiny Code - Server

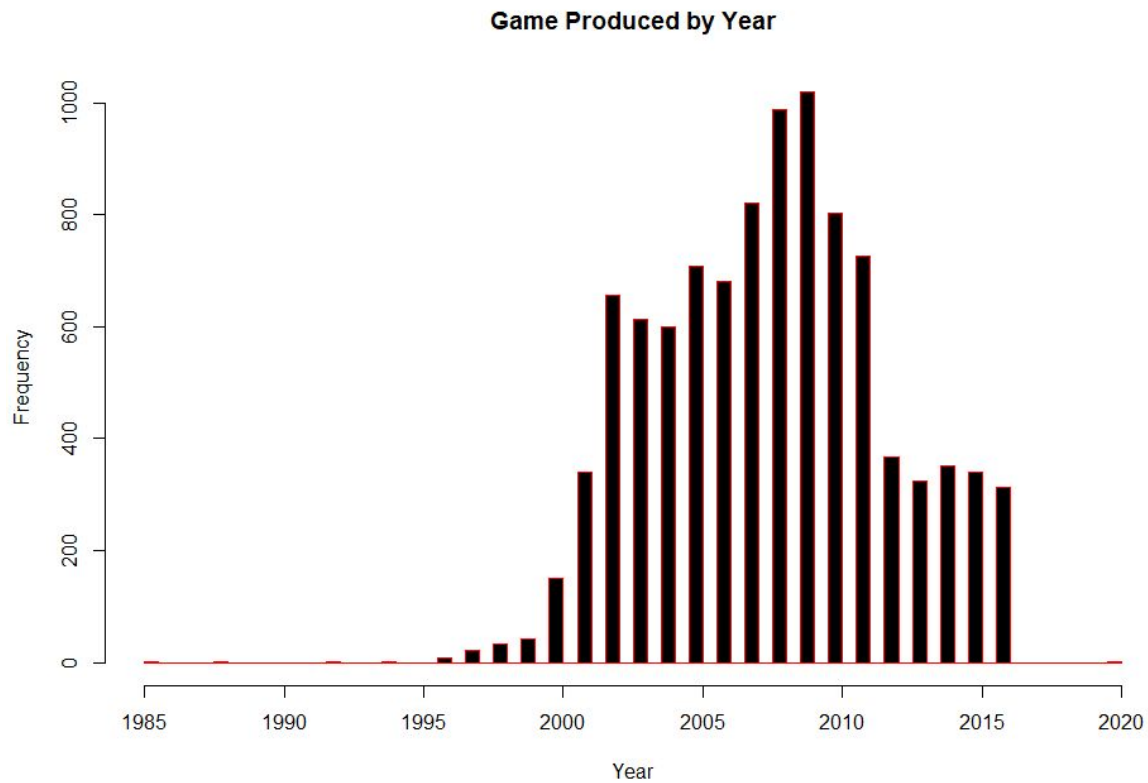


Figure 9.0 Graph made using R - Game Produced by Year

```
1 hist(data$Year,  
2     main="Game Produced by Year",  
3     xlab="Year",  
4     border="red",  
5     col="black",  
6     xlim=c(1985,2020),  
7     breaks=55)  
8
```

Figure 10.0 Graph Code

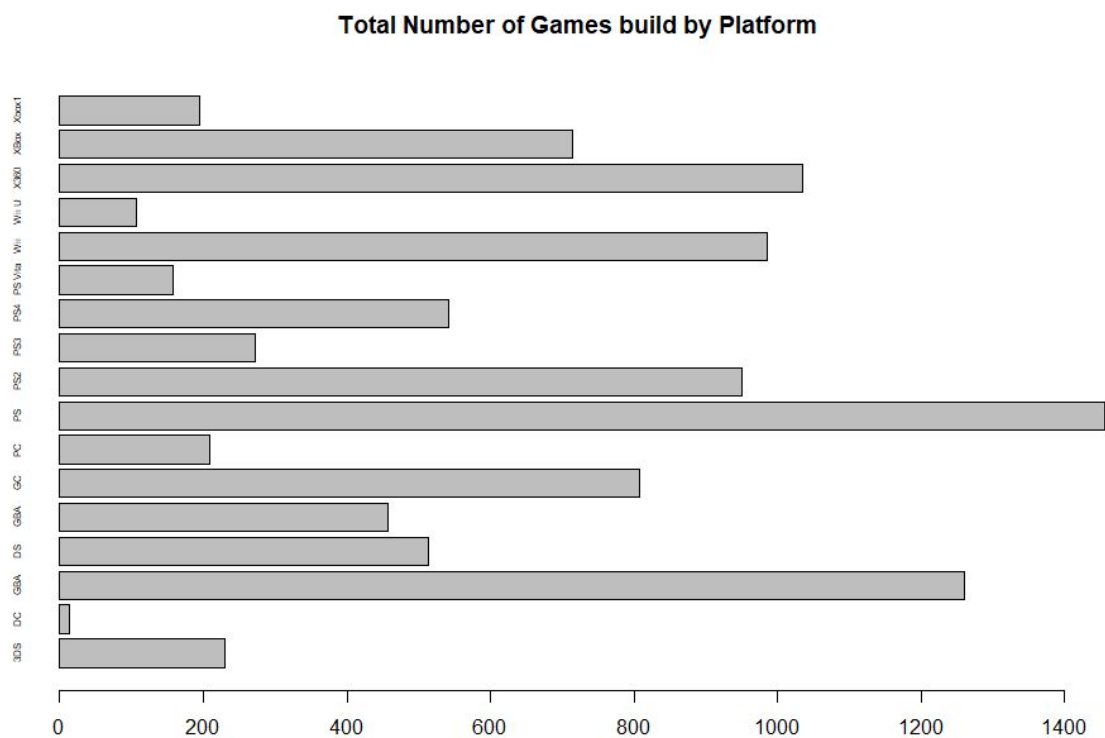
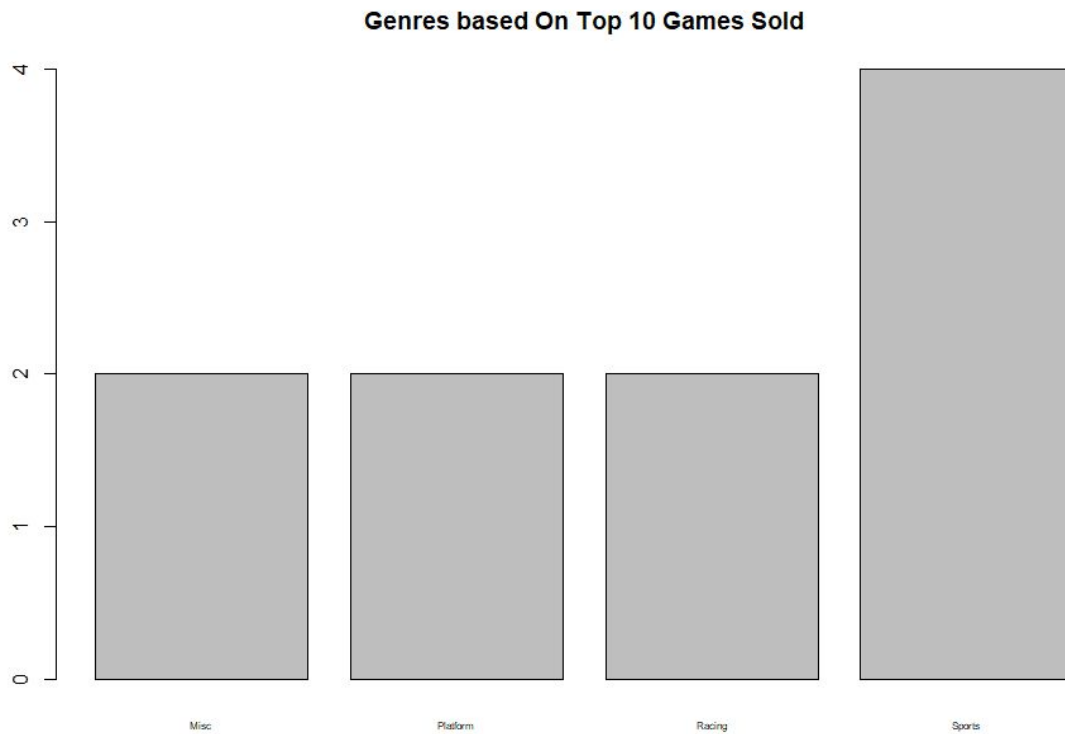


Figure 11.0 Graph made using R - Genres based on Top 10 Games & Total Number of Games by Platform

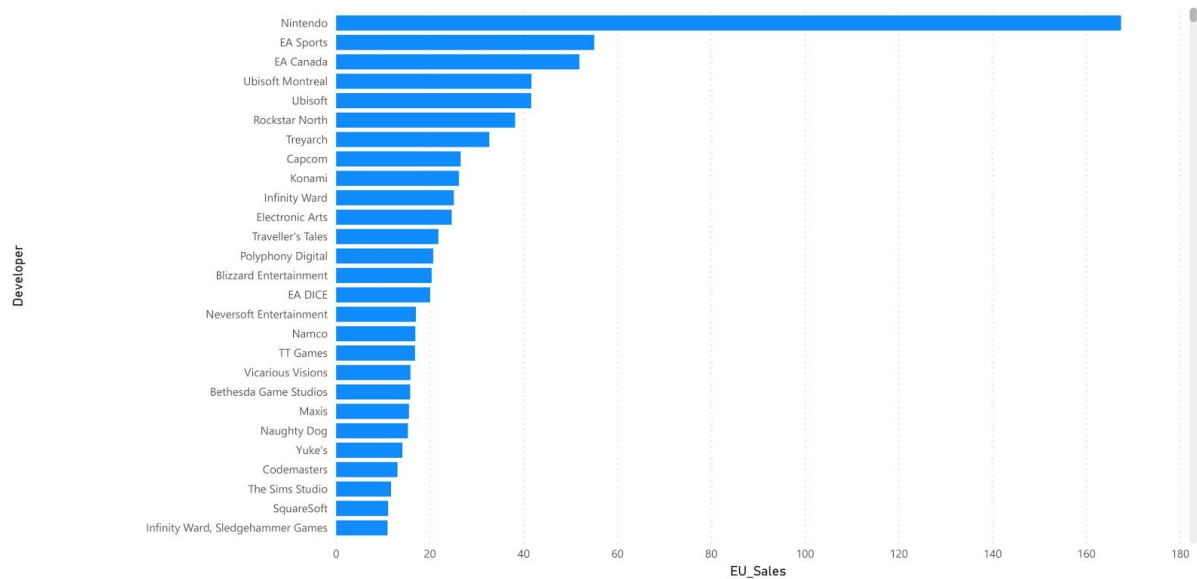

```

1 top10 <- head(data, 10)
2 barplot(table(head(data$Genre,10)),cex.names=0.5, main="Genres based On Top 10 Games Sold",
3         names.arg=c("Misc","Platform", "Racing", "Sports"))
4
5 table(data$Platform)
6 barplot(table(data$Platform),cex.names=0.5, main="Total Number of Games build by Platform", horiz=TRUE,
7         mes.arg=c("3DS", "DC", "GBA", "DS", "GBA", "GC", "PC", "PS", "PS2", "PS3", "PS4", "PS Vita", "Wii", "Wii U", "X360", "XBox", "Xbox1"))

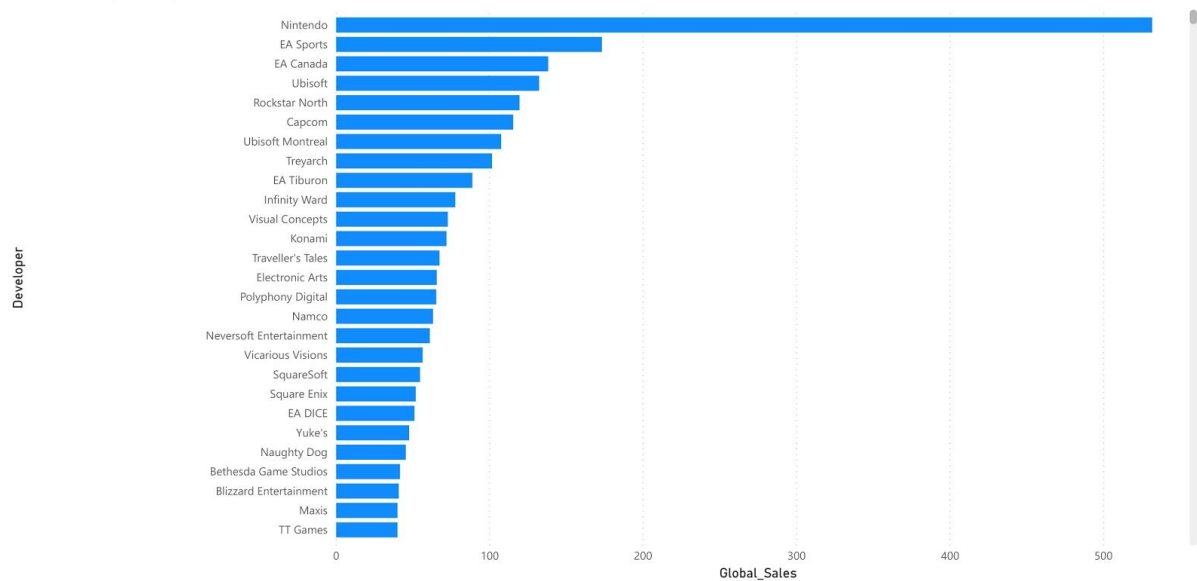
```

Figure 12.0 Graph Code

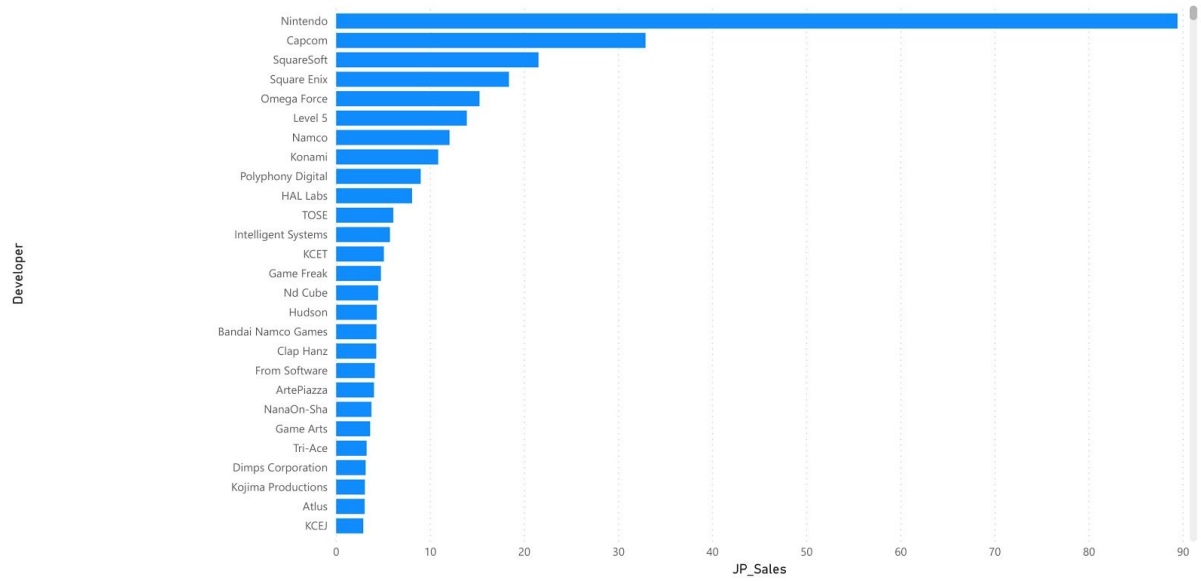
EU_Sales by Developer



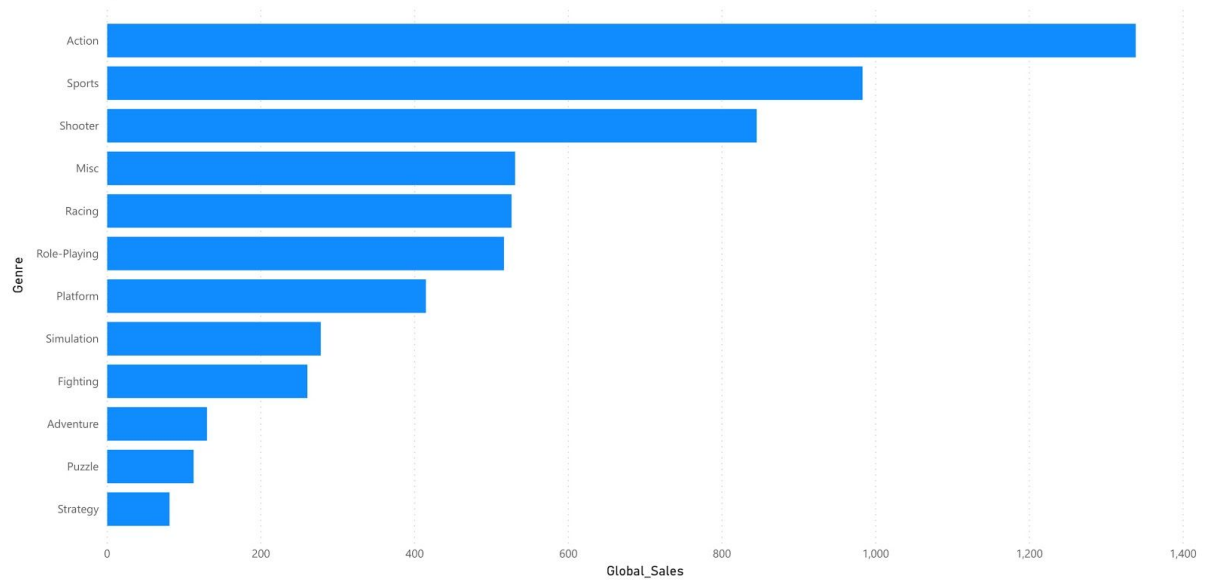
Global_Sales by Developer



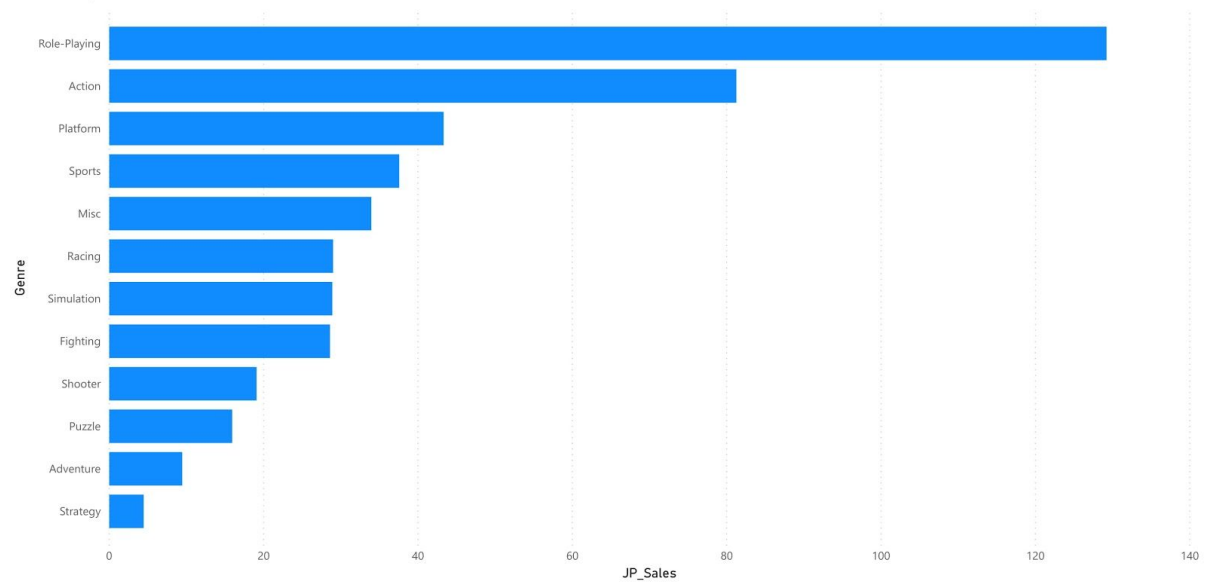
JP_Sales by Developer



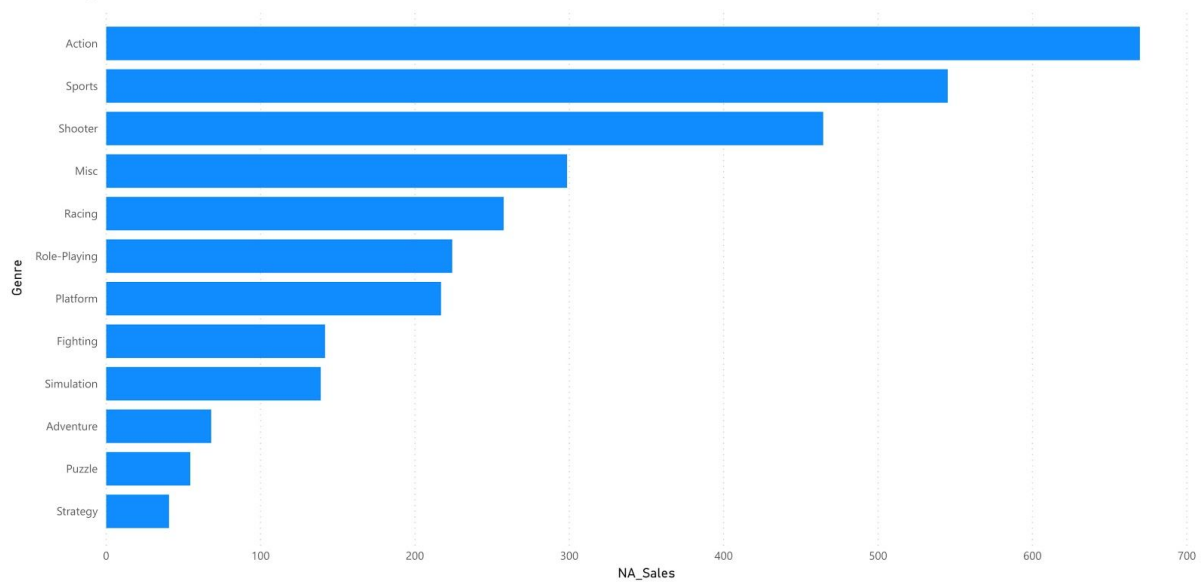
Global_Sales by Genre



JP_Sales by Genre



NA_Sales by Genre



Preferred Genre based on Regions

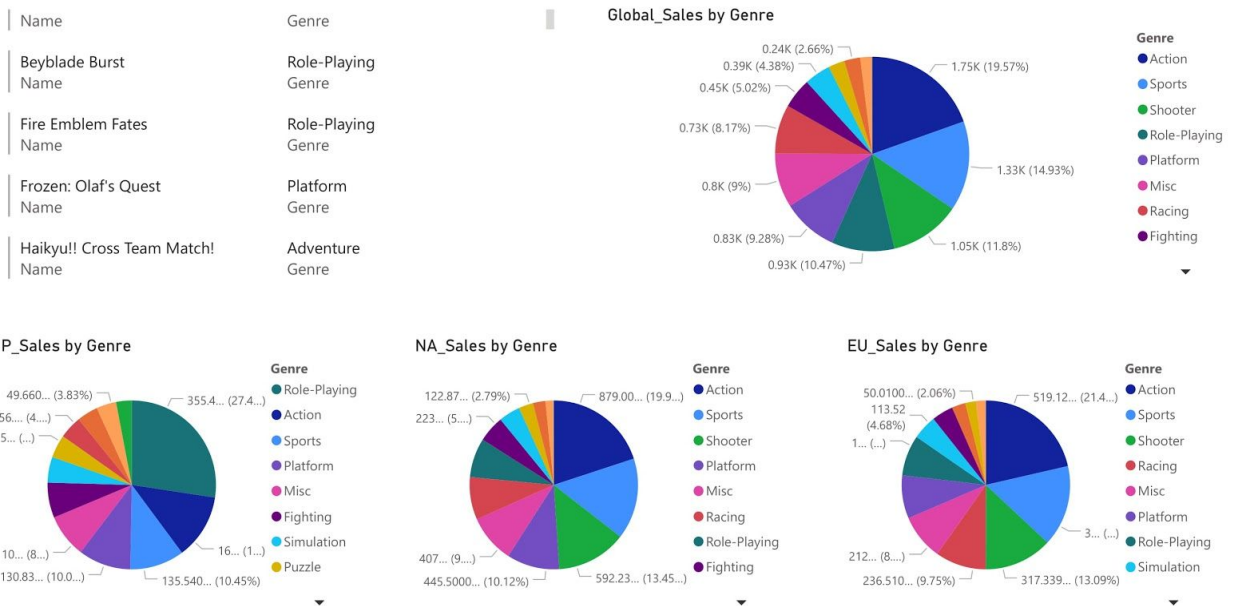


Figure 13.0 Additional Graph made using Power BI