

# IMDB Video Games Data Analysis

Qingyue Tian

2025-06-28

## Introduction and Background Information

This study aims to investigate the relationship between game ratings and sales, as well as the types, publishers, and platforms of games, between 2000 and 2020. For the ratings, I chose to use data from the large rating website IMDb because its ratings have a relatively large user base, the ratings are relatively reliable, and the dataset from its website can be found on Kaggle. I will present my research using various visual effects so that the audience can understand the content of my research most intuitively. Specifically, I will divide my research content into three sections. The first section examines the relationship between game ratings and sales. Is it true that the more people play the game, the higher the rating? This is what I care about. Can a good game in the traditional sense get the return it deserves in business? In the second section, we will delve deeper and add data such as the game's age, publisher, and platform to explore the characteristics of highly rated games in this era. In the third section, we will focus on one platform, the Wii platform, and examine the sales and ratings differences between first-party games and other third-party games on the platform. We will also investigate the impact of video game compatibility with the platform on their ratings.

## About Our Data

My primary dataset is a dataset from Kaggle, which captures data from the IMDB rating website. Its credibility can be guaranteed by the credibility of the IMDB website. It is a dataset from three years ago, which can fully cover the data up to 2020. The extra two years will also make the data in 2020 more referenceable after time has passed. It includes the following attributes: name, url, year, certificate, rating, votes, plot, Action, Adventure, Comedy, Crime, Family, Fantasy, Mystery, Sci-Fi and Thriller. Among them, we will choose the name, year, rating, and votes columns to use. We will not use the genre column, which ranges from Action to Thriller, as the genre criterion. This is because, in the True and False criteria of these genre columns, as long as the game touches the edge of this category, it will be classified as True, which causes almost all games in the Action category to be classified as True, making the data

inconvenient to analyze. We only focus on works with more than 500 ratings.

My secondary datasets also come from Kaggle and Wikipedia. My first secondary dataset comes from Kaggle, which was uploaded two years ago by ULRIK THYGE PEDERSEN, a Kaggle-certified database expert. The upload time of this data is not significantly different from that of the previous data, so the two datasets can be compared and supplemented. The primary purpose of this dataset is to augment the first dataset with additional columns, including a Platform column, a Genre column, a Publisher column, and four sales columns. For the merger of dataset two and dataset one, since there are sufficient datasets and it is cumbersome to check and standardize the different names of each game, we will only select games with the same name in both datasets for analysis.

My second secondary data set comes from Wikipedia, which lists Wii games released. We primarily focus on its developers' column, from which we select Nintendo's first-party developed games and non-Nintendo developed games to compare and determine how the compatibility of games and platforms affects their ratings and sales.

## **Exploratory Data Analysis**

First, we import the three data sets into the file and organize them, including selecting the required columns, filtering out rows whose data do not meet the requirements, and merging the data into a new dataset to facilitate our next step of visualizing the data.

## **Data Visualizations**

Our first research topic is to consider the relationship between game sales and game ratings. Under this topic, we not only explore the simple relationship between total game sales and game ratings but also consider that IMDB is an American rating website, so sales in North America may be more critical to the rating than sales in other regions. We will use visual data to answer whether this hypothesis is true.

## Rating vs Sales

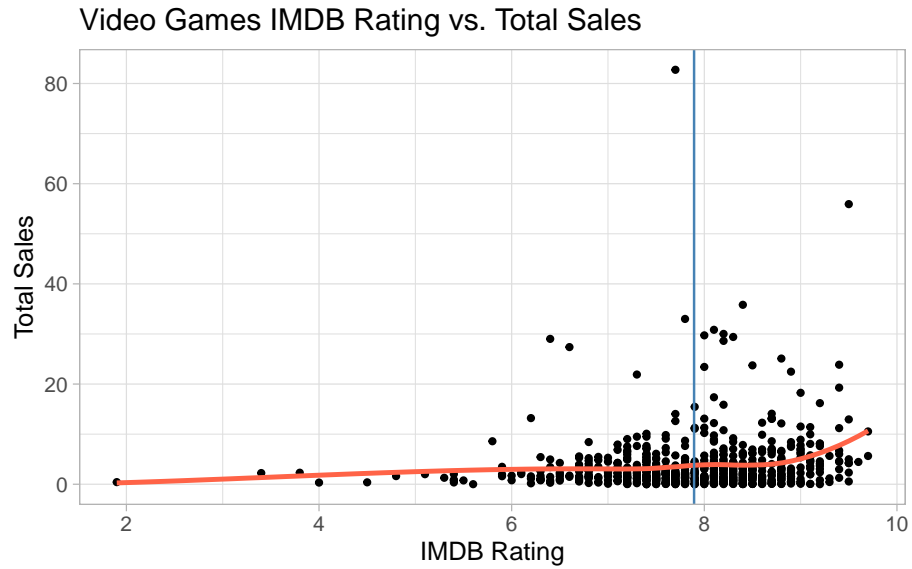


Figure 1

From the plot in Figure 1, we can see that in the score range of 6 to 10, games with significantly high sales (sales exceeding 10 million) appear to be evenly distributed, with a slight concentration at a score of 8. The blue vertical line represents the average of all scores, as shown in Table 1, with a value of 7.9 points. The red line in the figure represents the trend line of this scatter plot. It can be seen from the trend chart that when the score exceeds 8.7 points, sales exhibit a significant upward trend compared to before. This shows that the score is indeed proportional to sales to a certain extent, and the sales of games with a score higher than a specific value (8.7 points) will also increase significantly. The good games considered by the public are indeed reflected in their commercial value.

Let's take a look at the visualization with North American data added.

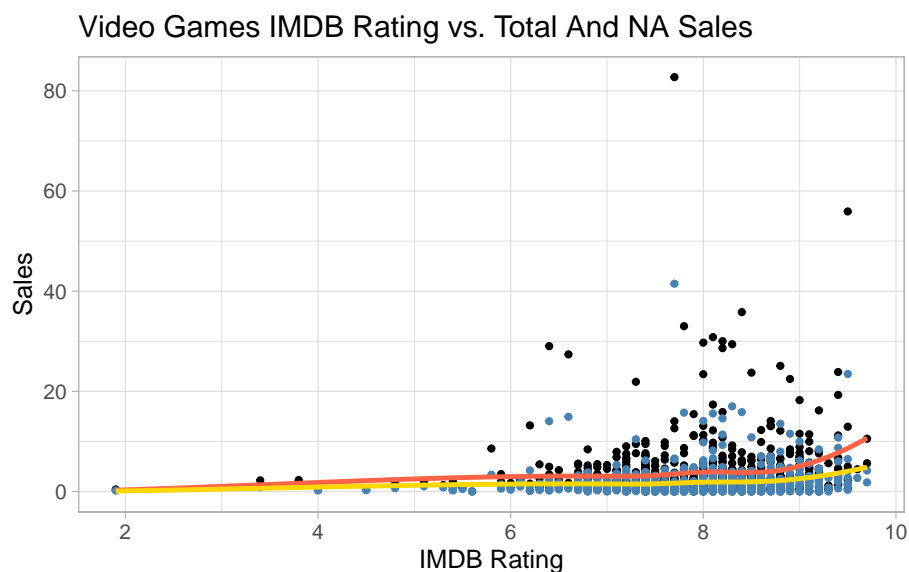


Figure 2

From the plot Figure 2, we can see that, compared to the total sales data, the North American sales data appears to be perfectly cut in half. This number seems to be inconsistent with the demographics. The total population of the United States and Canada is about 335 million, while the total population of Japan and Europe is about 865 million. This indicates that video game culture is more prevalent in North America than in other regions. In terms of the trend line, the general trend of the North American sales trend line is not significantly different from the total sales; however, the trend after the games with a score higher than 8.7 is relatively flat compared to the total sales. This indicates that games with higher scores will stimulate the purchasing desire of game enthusiasts worldwide.

It seems like North American sales are not more important than other countries in terms of ratings.

Our second research topic is whether the era, publisher, platform, and genre affect a game's rating data. In this section, we will explore whether a specific company is more successful with games of a particular genre.

## Rating vs Published Year

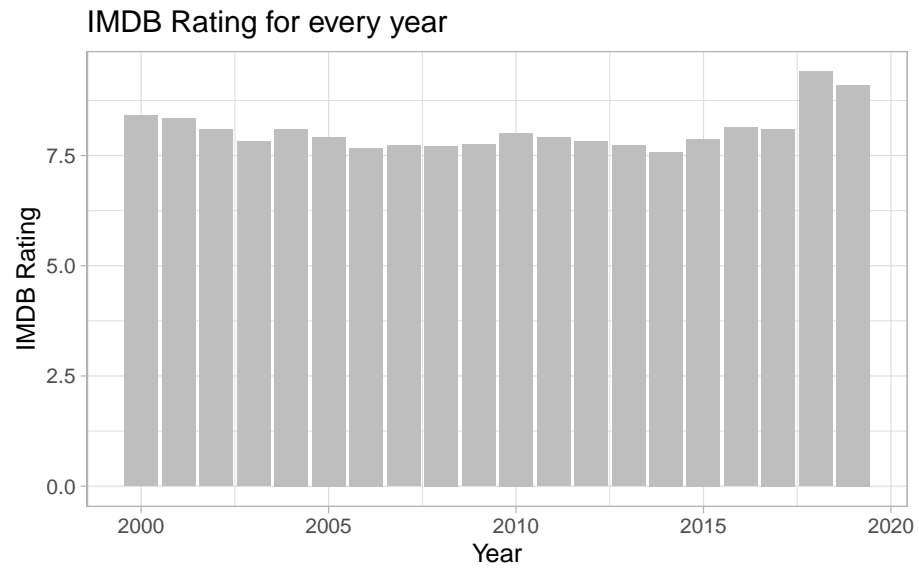


Figure 3

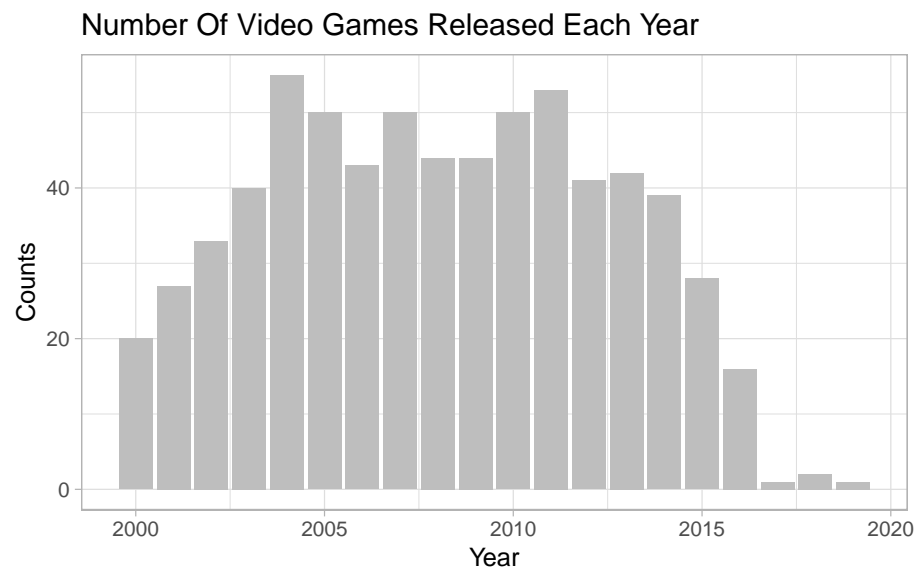


Figure 4

According to Figure 3, the ratings do not correlate with the year. In the second half of the figure, the difference between the data after 2006 and before is slightly larger because the amount of data has decreased significantly since 2017, making the average value look less average. I initially expected that in the release year of new game consoles, such as 2006 when the Wii was released, 2001 when the Xbox was released, and 2016 when the PlayStation 4 was released, the overall game ratings would have an upward trend because manufacturers would launch some big-budget games in conjunction with the consoles. However, the results in Figure 4 show that although there are specific significant differences in the number of game releases depending on the year, the release of a game console has no significant impact on the game's annual ratings and sales.

Since there are no significant fluctuations in game ratings in the years when game consoles are released, can we infer that game consoles (platforms) will not have a substantial impact on game ratings, and further, that the game manufacturers that manufacture game consoles (platforms) will not have a significant effect on game ratings?

### Rating vs Platform

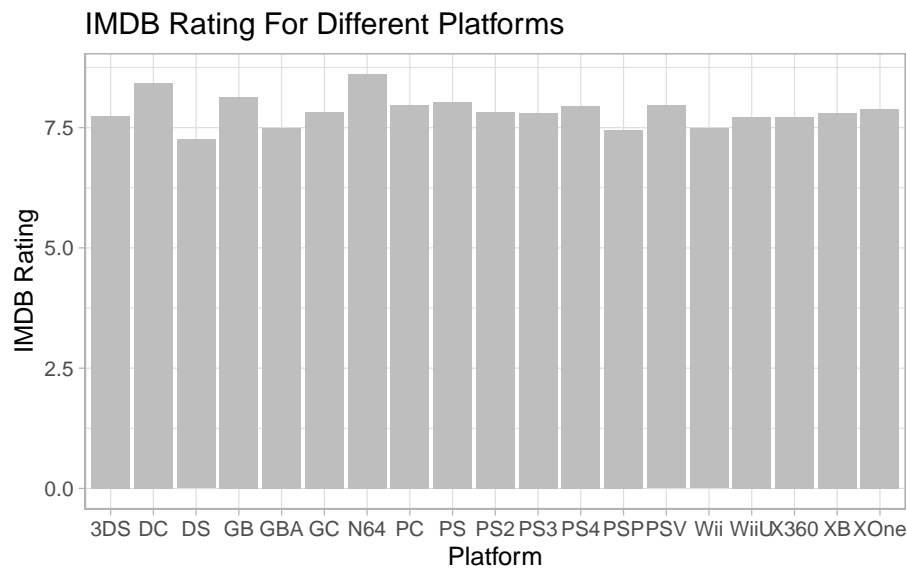


Figure 5

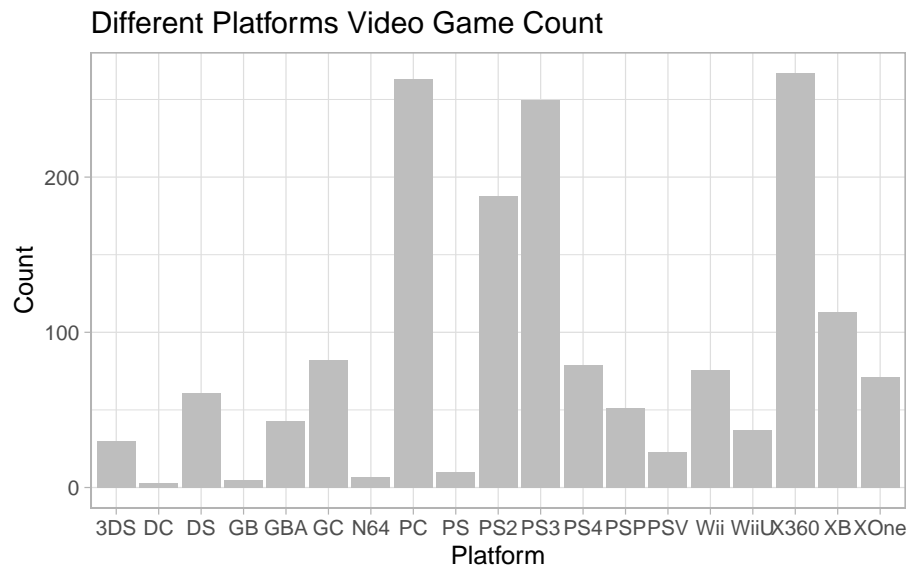


Figure 6

According to Figure 6 and Figure 5, although the number of games released on different platforms varies significantly, the difference in their average ratings is not substantial. The ratings of relatively popular platforms, such as PC, PS3, and X360, have not increased significantly compared to other ratings, but are closer to the average. Although games on different platforms have different characteristics, the image quality of PS4 is much better than that of PS2; the Wii, which must be connected to a TV, is significantly heavier than the 3DS, which is lightweight and can be played in the hand; however, the characteristics of lightness and image quality do not appear to impact the ratings of the game by players significantly.

### Rating vs Publisher

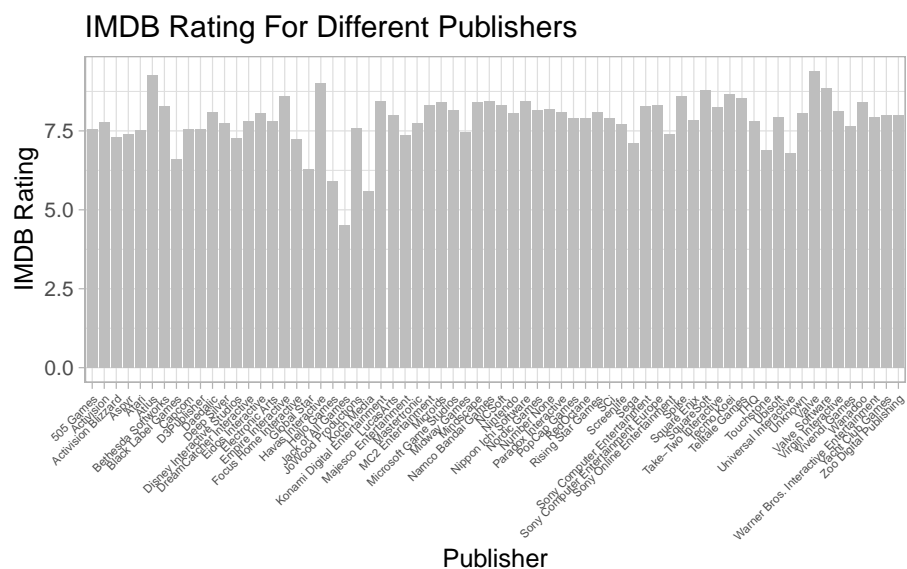


Figure 7

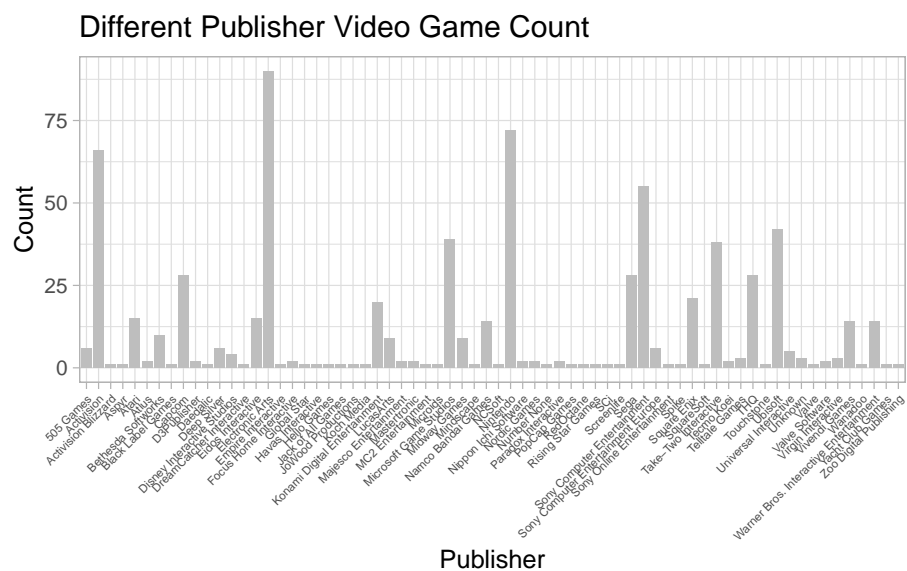


Figure 8



Table 1: The four publishers with the most games released

Publisher	count
Electronic Arts	90
Nintendo	72
Activision	66
Sony Computer Entertainment	55

From the above Figure 7 and Figure 8, we can see that the differences in scores between publishers are much greater than those between platforms. This is reasonable because it requires significantly more R&D investment to develop a platform than a game. From Figure 8, we can see that several large game manufacturers, such as Nintendo, Electronic Arts, Sony Computer Entertainment, and Ubisoft, all have scores of at least 7.5, indicating that the quality of games released by these manufacturers is generally high.

Now, let's turn our attention to the most basic attribute of the game: genre. Which games are sold the most in the market, and which games are likely to attract high ratings?

### Rating vs Genre

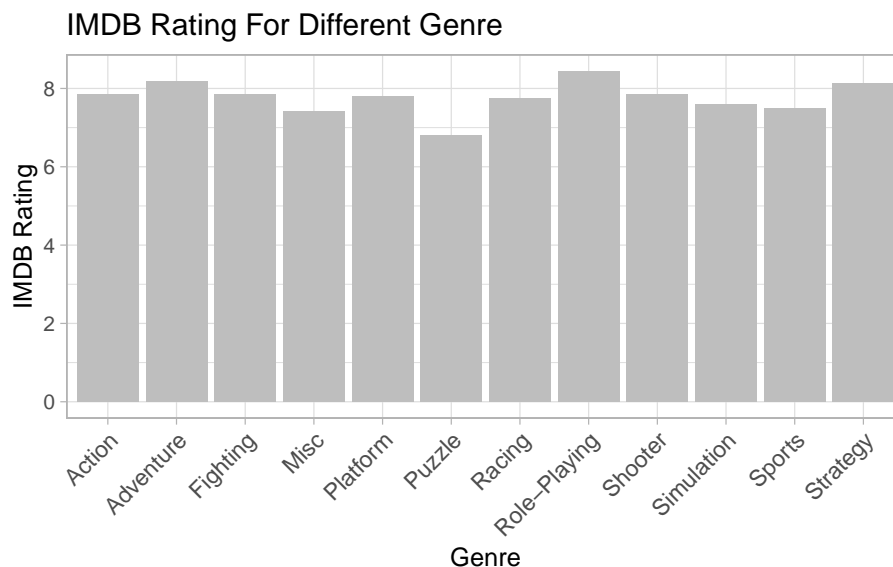


Figure 9

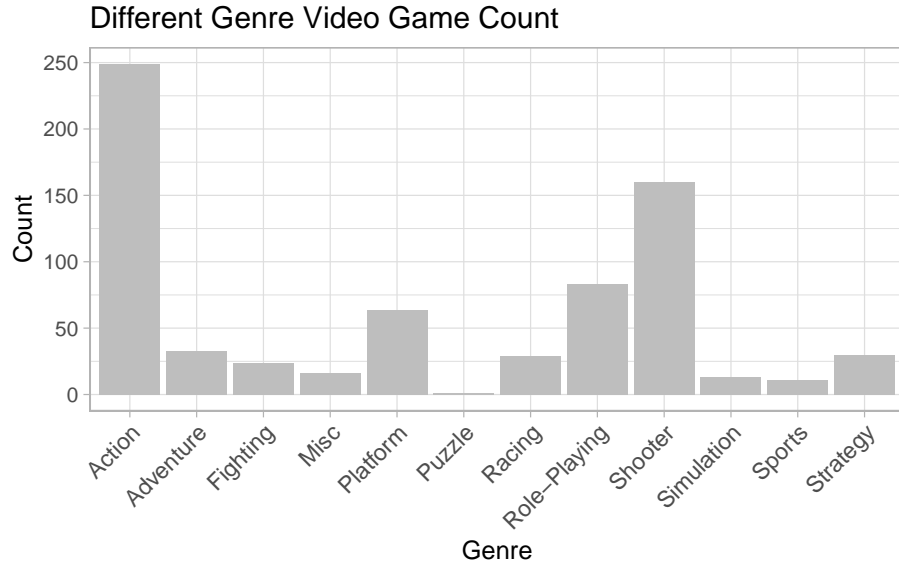


Figure 10

As can be seen from Figure 9 and Figure 10, the number of action games is significantly higher than that of other games. Design games are second, followed by role-playing games. In terms of ratings, role-playing games are the highest, followed by adventure games and strategy games. It is worth noting that puzzle games tend to have significantly lower sales and ratings compared to other types. As a game player, I can understand this phenomenon. Shooting games and action games have several “classic modes”, such as horizontal action games and first-person shooter games, which are relatively easy for game manufacturers to reproduce. Role-playing games require a complete storyline, which makes their production somewhat challenging, but they are relatively easy to attract attention. After all, a successful character will make many groups that were not initially interested in it aware of its existence. For example, not many people in the world watch anime, but almost everyone is familiar with the existence of Doraemon.

As for adventure and strategy games, adventure games typically require more resources to create the stage of the adventure world during the production process, so there is a certain threshold for their release. This threshold also ensures their level to a certain extent. Strategy games focus on gameplay, and it is usually challenging to conceive new gameplay. However, good strategy games can be played almost infinitely, unlike adventure games or role-playing games with fixed clearance processes and time, so their scores will also be relatively high. Regarding puzzle games, their primary target audience is younger children, so they may not be as attractive to specific groups, especially those who

can consume and enter the website to score.

## Nintendo

As shown in Figure 8, Nintendo is the second-largest game manufacturer in the data, releasing the second-most games between 2000 and 2020. At the same time, it also develops game consoles part-time. For example, Wii is a console designed by Nintendo. In Figure 6, you can see that our data includes more than 70 games. Next, we will narrow our horizons and focus on Nintendo’s Wii consoles. By comparing Nintendo’s first-party games with those of third-party developers, they will further explore the relationship between game consoles and ratings.

Pie chart of first-party and third-party game types

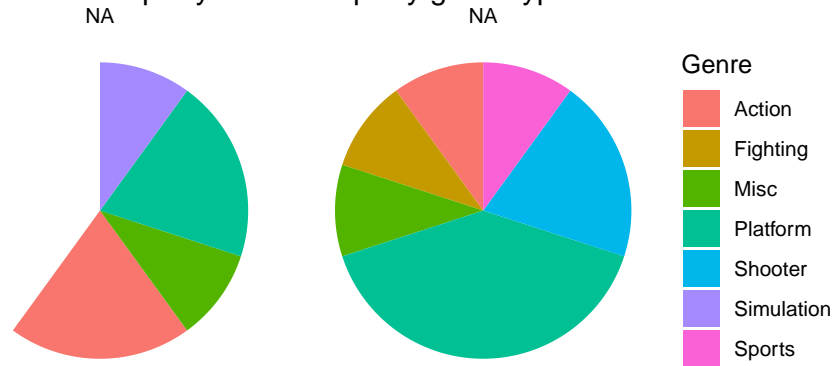


Figure 11

Table 2: A summary of the genres of first-party games and third-party games

Genre	Third-Party	First-Party
Action	1	2
Fighting	1	0
Misc	1	1
Platform	4	2
Shooter	2	0
Sports	1	0
Simulation	0	1

From Figure 11, we can see the types of first-party games and third-party games on the Wii. Third-party games tend to be more platform games and shooting games, while Nintendo's first-party games tend to be more platform and action games. From this, we can see that the Wii-type console is particularly well-suited for platform games, for example, the popular Super Monkey Ball series.

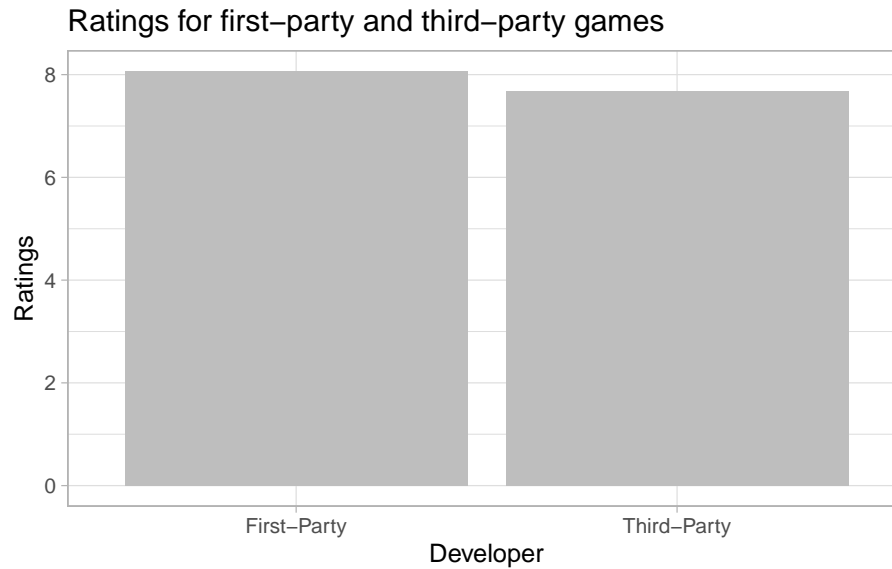


Figure 12

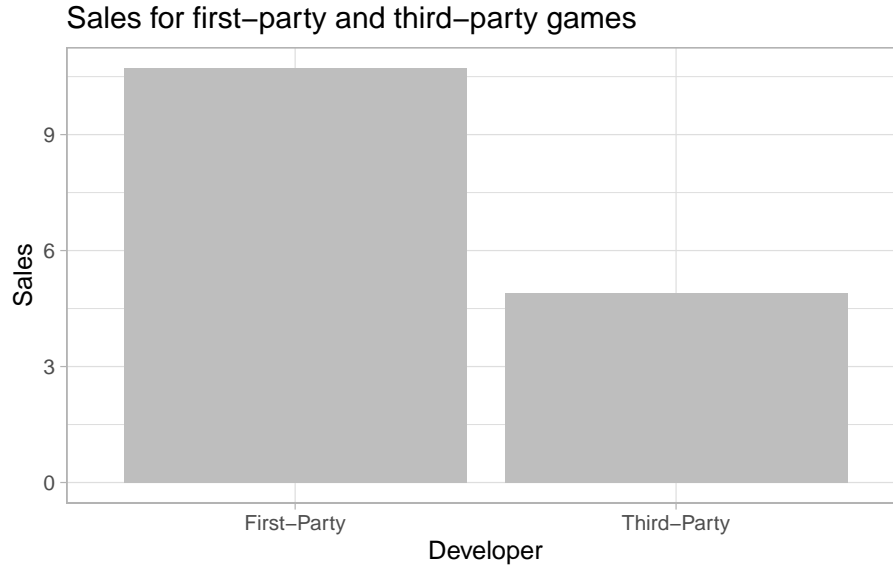


Figure 13

From Figure 12 and Figure 13, we can see that the average score of first-party games on the Nintendo Wii console is not significantly higher than that of third-party games. Still, the average sales of first-party games on the Nintendo console are considerably higher than those of third-party games. This shows that people are still more willing to buy first-party games that are compatible with the console, and their quality is guaranteed.

## Conclusion

The analysis of video games from 2000 to 2020 helps us gain a better understanding of the game market during that period. From the game sales records and the ratings on the IMDb website, we can conclude that, although it is not immediately apparent, the game rating is proportional to sales, and the sales of games with a score of 9 or above exhibit a clear upward trend. The year of the game does not significantly impact the game rating, which suggests that whether it is an old, pixel-style game or a new 3A masterpiece, there will be a group of gamers who appreciate this aspect. The playing method, such as handhelds like the 3DS, consoles like the PS4, and PC, has little impact on the game rating. We found that game companies tend to develop shooting games and action games, while games with higher ratings are often role-playing games and adventure games. We found that games released by large game companies typically come with a high-quality guarantee. Finally, by comparing the data of first-party games and third-party games for Nintendo Wii, we found that first-party games are similar to third-party games in terms of ratings. Still, they

are far better than third-party games in terms of sales, which suggests that a certain segment of gamers is more inclined to buy first-party games.

## Reference

- Wikimedia Foundation. (2025, May 22). *List of best-selling Wii Video Games*.  
Wikipedia. [https://en.wikipedia.org/wiki/List\\_of\\_best-selling\\_Wii\\_video\\_games](https://en.wikipedia.org/wiki/List_of_best-selling_Wii_video_games)
- Feroze, Z. (2025, March 21). *Video games sale*. Kaggle.  
<https://www.kaggle.com/datasets/zahidmughal2343/video-games-sale>
- Talay, M. A. (2022, September 2). *IMDB video games*. Kaggle.  
<https://www.kaggle.com/datasets/muhammadadiltalay/imdb-video-games>

## Code Appendix

```
1 # Load necessary packages
2 # The code in this pdf follows Google's R Style Guide
3 library(ggplot2)
4 library(dplyr)
5 library(tidyverse)
6 library(tidyr)
7 library(rvest)
8 library(readr)
9 library(knitr)
10
11 #Importing the Data
12 IMDBDataSetRaw <- read.csv(
13   "/Users/maming/Downloads/imdb-videogames.csv"
14 )
15 SalesDataSetRaw <- read.csv(
16   "/Users/maming/Downloads/video games sales.csv"
17 )
18 WebsitePage <-
19   "https://en.wikipedia.org/wiki/List_of_best-selling_Wii_video_games"
20 TableList <- WebsitePage %>%
21   read_html() %>%
22   html_nodes(css = "table") %>%
23   html_table(fill = TRUE)
24 WiiDataSetRaw <- TableList[[3]]
25
26 #cleaning the three raw data
27 IMDBCleaned <- IMDBDataSetRaw %>%
28   filter(year >= 2000 & year <= 2020) %>%
29   select(name, year, rating, votes) %>%
30   mutate(votes = gsub(",", "", votes)) %>%
31   mutate(votes = as.integer(votes)) %>%
32   filter(votes >= 500) %>%
33   filter(!is.na(rating) & rating != "")
34
35 SalesCleaned <- SalesDataSetRaw %>%
36   filter(Year >= 2000 & Year <= 2020) %>%
37   filter(!is.na(Platform) & Platform != "") %>%
38   filter(!is.na(Genre) & Genre != "") %>%
39   filter(!is.na(Publisher) & Publisher != "") %>%
40   filter(!is.na(NA_Sales) & NA_Sales != "") %>%
41   filter(!is.na(EU_Sales) & EU_Sales != "") %>%
42   filter(!is.na(JP_Sales) & JP_Sales != "") %>%
43   filter(!is.na(Other_Sales) & Other_Sales != "") %>%
```



```

44   filter(!is.na(Global_Sales) & Global_Sales != "") %>%
45   select(-Rank)
46
47   WiiCleaned <- WiiDataSetRaw %>%
48     mutate(Game = gsub("†", "", Game)) %>%
49     mutate(Game = gsub("\\[e\\]", "", Game)) %>%
50     rename(releaseDate = "Release date[a]",
51            Developer = "Developer(s)", Name = Game) %>%
52     mutate(releaseDate = format(as.Date(
53       releaseDate, format = "%B %d, %Y"), "%Y")) %>%
54     mutate(Nintendo = grepl("Nintendo EAD", Developer)) %>%
55     select(Name, releaseDate, Nintendo)
56
57   #Merge data frames into new data frames for easy visualization
58   SalesMerge <- SalesCleaned %>%
59     select(Name, Global_Sales, JP_Sales,
60            NA_Sales, EU_Sales, Other_Sales) %>%
61     group_by(Name) %>%
62     summarise(
63       TotalSales = sum(Global_Sales, na.rm = TRUE),
64       JPSale = sum(JP_Sales, na.rm = TRUE),
65       NASales = sum(NA_Sales, na.rm = TRUE),
66       EUSales = sum(EU_Sales, na.rm = TRUE),
67       OtherSales = sum(Other_Sales, na.rm = TRUE),
68       .groups = "drop"
69     )
70   IMDBMerge <- IMDBCleaned %>%
71     group_by(name) %>%
72     slice_max(votes, with_ties = FALSE) %>%
73     ungroup()
74
75   # Data frames for rating vs sales
76   RatingSales <- inner_join(
77     IMDBMerge, SalesMerge, by = c("name" = "Name"))
78   AverageRating <- RatingSales %>%
79     summarize(avgRating = mean(rating))
80
81   # Data frames for rating vs year
82   RatingYear <- RatingSales %>%
83     select(year, rating) %>%
84     group_by(year) %>%
85     summarize(RatY = mean(rating))
86   AmontVG <- RatingSales %>%
87     select(year, rating) %>%
88     group_by(year) %>%

```

```

89     summarise(count = n())
90
91 # Data frames for rating vs platform
92 SalesMergePlatform <- SalesCleaned %>%
93   select(Name, Platform, Genre, Publisher, Global_Sales)
94 RatingPlatformMerge <- inner_join(
95   IMDBMerge, SalesMergePlatform, by = c("name" = "Name"))
96 RatingPlatform <- RatingPlatformMerge %>%
97   select(Platform, rating) %>%
98   group_by(Platform) %>%
99   summarize(RatP = mean(rating))
100 RatingPlatformCount <- RatingPlatformMerge %>%
101   select(Platform, rating) %>%
102   group_by(Platform) %>%
103   summarise(count = n())
104
105 # Data frames for rating vs publishers
106 SalesMergePG <- SalesCleaned %>%
107   select(Name, Genre, Publisher) %>%
108   distinct()
109 RatingPGMerge <- inner_join(
110   IMDBMerge, SalesMergePG, by = c("name" = "Name"))
111 RatingPublisher <- RatingPGMerge %>%
112   select(Publisher, rating) %>%
113   group_by(Publisher) %>%
114   summarize(RatPub = mean(rating))
115 RatingPublisherCount <- RatingPGMerge %>%
116   select(Publisher, rating) %>%
117   group_by(Publisher) %>%
118   summarise(count = n())
119
120 # Data frames for rating vs genre
121 RatingGenre <- RatingPGMerge %>%
122   select(Genre, rating) %>%
123   group_by(Genre) %>%
124   summarize(RatG = mean(rating))
125 RatingGenreCount <- RatingPGMerge %>%
126   select(Genre, rating) %>%
127   group_by(Genre) %>%
128   summarise(count = n())
129
130 #Nintendo Wii data frames
131 SalesWii <- SalesCleaned %>%
132   filter(Platform == "Wii") %>%
133   select(Name, Genre, Global_Sales)

```

```

134 RatingWii <- inner_join(
135   IMDBMerge, SalesWii, by = c("name" = "Name"))
136 RatingWiiN <- inner_join(
137   RatingWii, WiiCleaned, by = c("name" = "Name"))
138 RatingWiiNG <- RatingWiiN %>%
139   group_by(Nintendo, Genre) %>%
140   summarise(count = n(), .groups = "drop") %>%
141   mutate(Nintendo = factor(Nintendo,
142                             levels = c("TRUE", "FALSE"),
143                             labels = c("First-Party", "Third-Party")))
144 RatingWiiNG_wide <- RatingWiiNG %>%
145   pivot_wider(
146     names_from = Nintendo,
147     values_from = count,
148     values_fill = 0
149   )
150 RatingWiiNR <- RatingWiiN %>%
151   group_by(Nintendo) %>%
152   summarize(RatN = mean(rating), .groups = "drop") %>%
153   mutate(Nintendo = factor(Nintendo,
154                             levels = c("TRUE", "FALSE"),
155                             labels = c("First-Party", "Third-Party")))
156 RatingWiiNS <- RatingWiiN %>%
157   group_by(Nintendo) %>%
158   summarize(SalN = mean(Global_Sales), .groups = "drop") %>%
159   mutate(Nintendo = factor(Nintendo,
160                             levels = c("TRUE", "FALSE"),
161                             labels = c("First-Party", "Third-Party")))
162
163 # Visualization of video game total sales and IMDB ratings
164 ggplot()+
165   geom_point(data = RatingSales,
166             mapping = aes(x = rating, y = TotalSales), size = 1)+
167   geom_smooth(data = RatingSales,
168             mapping = aes(x = rating, y = TotalSales),
169             method = "loess", color = "tomato", se = FALSE)+
170   geom_vline(data = AverageRating,
171             aes(xintercept = avgRating), color = "steelblue")+
172   labs(
173     x = "IMDB Rating",
174     y = "Total Sales",
175     title = "Video Games IMDB Rating vs. Total Sales"
176   )+
177   theme_light()
178

```

```

179 # Visualization of video game north american sales and IMDB ratings
180 ggplot( data = RatingSales, )+
181   geom_point(mapping = aes(x = rating, y = TotalSales),
182             size = 1)+
183   geom_point(mapping = aes(x = rating, y = NASales),
184             size = 1, color = "steelblue")+
185   geom_smooth(data = RatingSales,
186             mapping = aes(x = rating, y = TotalSales),
187             method = "loess", color = "tomato", se = FALSE)+
188   geom_smooth(data = RatingSales,
189             mapping = aes(x = rating, y = NASales),
190             method = "loess", color = "gold", se = FALSE)+
191   labs(
192     x = "IMDB Rating",
193     y = "Sales",
194     title = "Video Games IMDB Rating vs. Total And NA Sales"
195   )+
196   theme_light()
197
198 # Visualization of video games published year and IMDB ratings
199 ggplot(RatingYear, aes(x = year, y = RatY)) +
200   geom_col(fill = "grey") +
201   theme_minimal() +
202   labs(title = "IMDB Rating for every year",
203        x = "Year",
204        y = "IMDB Rating") +
205   theme_light()
206
207 # Visualization of number of video games released by year
208 ggplot(AmontVG, aes(x = year, y = count)) +
209   geom_col(fill = "grey") +
210   theme_minimal() +
211   labs(title = "Number Of Video Games Released Each Year",
212        x = "Year",
213        y = "Counts") +
214   theme_light()
215
216 # Visualization of video games on different plantform and their IMDB ratings
217 ggplot(RatingPlatform, aes(x = Platform, y = RatP)) +
218   geom_col(fill = "grey") +
219   theme_minimal() +
220   labs(title = "IMDB Rating For Different Platforms",
221        x = "Platform",
222        y = "IMDB Rating") +
223   theme_light()

```

```

224
225 # Visualization of number of video games on different platform
226 ggplot(RatingPlatformCount, aes(x = Platform, y = count)) +
227   geom_col(fill = "grey") +
228   theme_minimal() +
229   labs(title = "Different Platforms Video Game Count",
230        x = "Platform",
231        y = "Count") +
232   theme_light()
233
234 # Visualization of video games on different publishers and their IMDB ratings
235 ggplot(RatingPublisher, aes(x = Publisher, y = RatPub)) +
236   geom_col(fill = "grey") +
237   theme_minimal() +
238   labs(title = "IMDB Rating For Different Publishers",
239        x = "Publisher",
240        y = "IMDB Rating") +
241   theme_light() +
242   theme(axis.text.x = element_text(angle = 45, hjust = 1,
243                                     size = 5))
244
245 # Visualization of number of video games on different publisher
246 ggplot(RatingPublisherCount, aes(x = Publisher, y = count)) +
247   geom_col(fill = "grey") +
248   theme_minimal() +
249   labs(title = "Different Publisher Video Game Count",
250        x = "Publisher",
251        y = "Count") +
252   theme_light() +
253   theme(axis.text.x = element_text(angle = 45, hjust = 1,
254                                     size = 5))
255
256 # Table of number of video games on different publisher
257 RatingPublisherCount %>%
258   arrange(desc(count)) %>%
259   slice_head(n = 4) %>%
260   kable(caption = "The four publishers with the most games released")
261
262 # Visualization of video games on different genre and their IMDB ratings
263 ggplot(RatingGenre, aes(x = Genre, y = RatG)) +
264   geom_col(fill = "grey") +
265   theme_minimal() +
266   labs(title = "IMDB Rating For Different Genre",
267        x = "Genre",
268        y = "IMDB Rating") +

```

```

269 theme_light() +
270 theme(axis.text.x = element_text(angle = 45, hjust = 1,
271                                   size = 10))
272
273 # Visualization of number of video games on different genre
274 ggplot(RatingGenreCount, aes(x = Genre, y = count)) +
275   geom_col(fill = "grey") +
276   theme_minimal() +
277   labs(title = "Different Genre Video Game Count",
278        x = "Genre",
279        y = "Count") +
280   theme_light() +
281   theme(axis.text.x = element_text(angle = 45, hjust = 1,
282                                     size = 10))
283
284 # Visualization of video games on Nintendo Wii Genre by developer
285 ggplot(RatingWiiNG, aes(x = "", y = count, fill = Genre)) +
286   geom_bar(stat = "identity", width = 1) +
287   coord_polar(theta = "y") +
288   facet_wrap(~ Nintendo, labeller = as_labeller(c(
289     "TRUE" = "First-Party",
290     "FALSE" = "Third-Party")))) +
291   theme_void() +
292   labs(title = "Pie chart of first-party and third-party game types")
293
294 # A table of the genres of first-party games and third-party games
295 kable(RatingWiiNG_wide,
296       caption = "A summary of the genres of first-party games and third-party games")
297
298 # Visualization of video games on Nintendo Wii Rating by developer
299 ggplot(RatingWiiNR, aes(x = Nintendo, y = RatN)) +
300   geom_col(fill = "grey") +
301   theme_minimal() +
302   labs(title = "Ratings for first-party and third-party games",
303        x = "Developer",
304        y = "Ratings") +
305   theme_light()
306
307 # Visualization of video games on Nintendo Wii sales by developer
308 ggplot(RatingWiiNS, aes(x = Nintendo, y = SalN)) +
309   geom_col(fill = "grey") +
310   theme_minimal() +
311   labs(title = "Sales for first-party and third-party games",
312        x = "Developer",
313        y = "Sales") +

```

314

```
theme_light()
```