

Introduction

When you're making a new game, it's important to keep potential sales in mind, especially if the goal is to make money. This issue applies to any game developer looking to make a profit out of a game. This analysis is meant to show developers what kind of game they should make. It would be in their best interest to create a game with a specific genre or on a specific platform if it is shown that a certain one is shown to have historically impacted global sales. For an example, if Platform and Genre were found to be good predictors, a game developer would know exactly which platform contributes the most to global sales and what kind of game generates the most sales. If Critic score was found to be a good predictor, a game developer could research the games whose critic score is high and get a general sense of how critics rate games. They would be able to use those features and implement them into their own games to maximize their sales. There would be similar reasoning for the remaining possible predictors.

I aim to find the best set of predictors excluding sales from the different regions for global sales using a statistical model. After that, I plan to compare the accuracy of predictions from the statistical model to that of a model built from machine learning.

Dataset

The dataset that I will use is located on Kaggle at: <https://www.kaggle.com/rush4ratio/video-game-sales-with-ratings>

It combines a web scrape of VGChartz data with a web scrape from Metacritic. It acknowledges that there are null values because Metacritic only includes a subset of the platforms.

It includes the following columns:

Name – The title of the game

Platform – What platform the game is on. I.e. PC, PS2

Year_of_Release – The year the game was released

Genre – Genre of the game (Only one is listed)

Publisher – Publisher of the game

NA_Sales – Sales in North America (in millions)

EU_Sales – Sales in Europe (in millions)

JP_Sales – Sales in Japan (in millions)

Other_Sales – Sales in the rest of the world (in millions)

Global_Sales – Total worldwide sales (in millions)

Critic_score – Aggregate score compiled by Metacritic staff

Critic_count – Number of critics used in Critic_Score

User_score – Score by Metacritic's subscribers

User_count – Number of users in the User_score

Developer – Party responsible for developing the game

Rating – The ESRB ratings

I do not know of any other datasets I can use. If I were to find one, it would need to be games listed in the same range of data from the same year. If it is data used in any other year, it would not accurately be matched with the global sales found in the data set. For an example, the global sales could be 20 million higher in the 2018 data for 'Overwatch' compared to the data in 2016 (if this data point exists in the dataset). So, this makes finding a suitable dataset to combine with difficult. If there are games found outside what is available in the dataset, there would be null values in the other columns.

If another dataset is used, it would most likely be better to webscrape it from scratch.

Data Cleaning

The first cleaning step I performed was removing all null values. I checked beforehand that none of the null values came from irrelevant columns (regional sales). My question requires complete cases and for the sake of simplicity, imputation is not used. The data set is large enough for imputation to not be impactful.

I noticed some of the columns were of the wrong data type. User_Score was an object instead of a float64 value. I changed it to a float64 data type. I also changed Platform, Genre, Rating, and Publisher as a category instead of an object.

Under Platform, DC (Dreamcast) has a fairly low sample size. I combine it with the WiiU category and renamed it 'Other'.

There are many categories under the Publisher column. I found the first Publisher with less than 30 observations and then took everything below that and combined it into the 'Other' category.

The Developer column was done the same since there many distinct categories. I took the top 50 categories with the remaining ones collapsed in the 'Other' category.

Next, I created a new column which shows if the Publisher is the same value as the Developer.

I notice the Rating column has only 1 count for RP, K-A, and AO. The game tagged as AO was GTA: San Andreas. There was a controversy about a scene in that game which made it AO several years after release. Originally, it was tagged as M, so I collapse it with the M category. K-A is kids to adults, so that

belongs in the E category. RP means Rating Pending. I collapse it with the T category because it is the most populated one.

I examine how many rows have Critic_Count or User_Count less than 10. These are the number of votes that contribute to Critic_Score and User_Score respectively. There are 1954 rows. I don't remove any because there are so many.

I change the year_of_release column to years_since_release. The data was collected at the end of 2016, so I assume that any game released in 2016 has been out for a year, 2015 for 2 years, etc.

Finally, I created a new column to determine if a game is good or not. The definition I used is if User_Score \geq 8, the game is good.

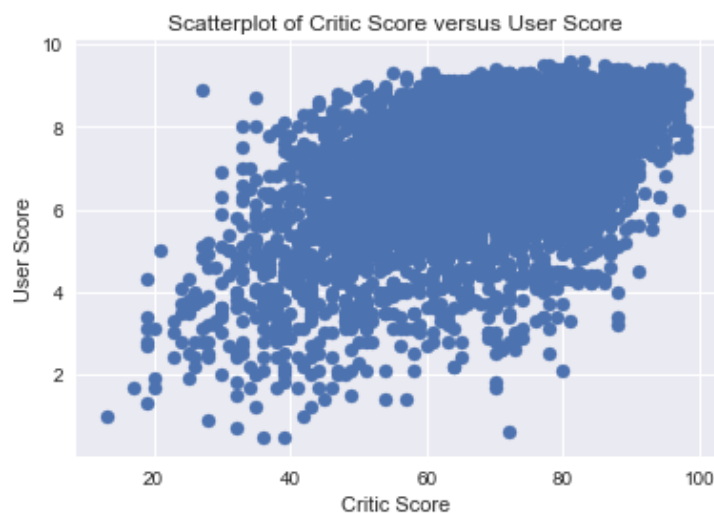
There are outliers, however, they are valid, so I don't remove them. For an example, one of the outliers is the highest Global_Sales contributor, Wii Sports. It has over 2 times the value of the 2nd highest Global_Sales. The outliers can be examined further than the analysis to see what effect they have if necessary.

Exploratory Data Analysis

There are 16719 rows of data in the original dataset. In the cleaned data, there are 6825 complete cases. The games with the highest global Sales is Wii Sports at 82.53 million.

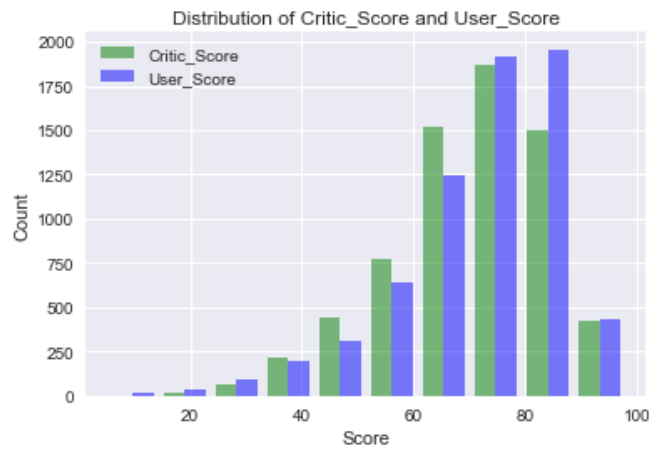
In the EDA, I will go through a few visualizations from comparing a variable with global sales and what it tells us. I will also use the same visualizations to look at the full dataset to see if anything changes if we had more observations in our cleaned dataset.

How does Critic Score compare with User Score?



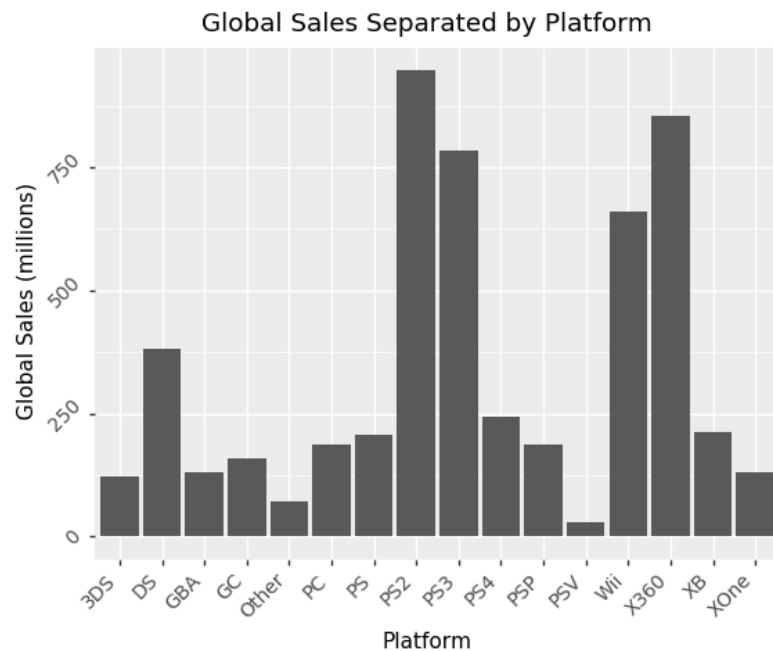
We see that there is a positive correlation between these two variables. The pearson correlation coefficient is 0.58 showing a moderate-strong positive correlation. This means that when critics rate

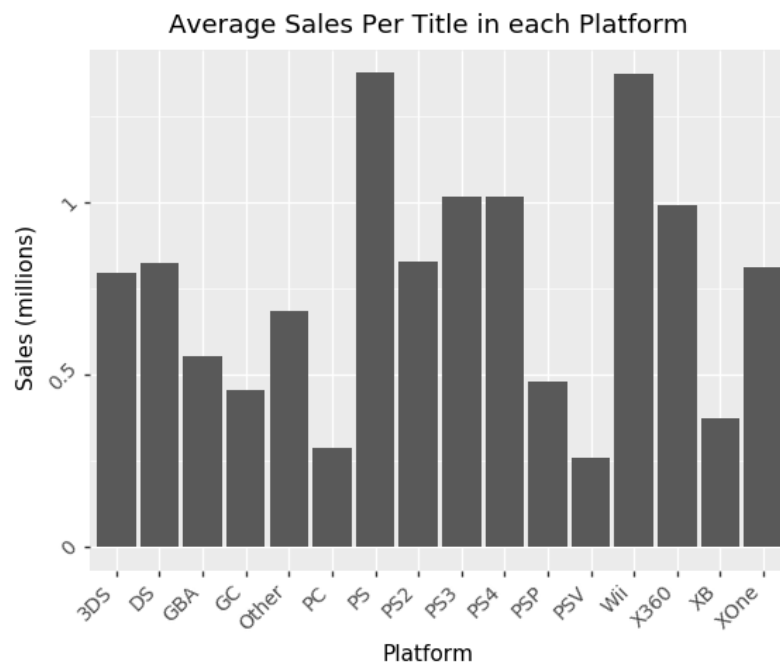
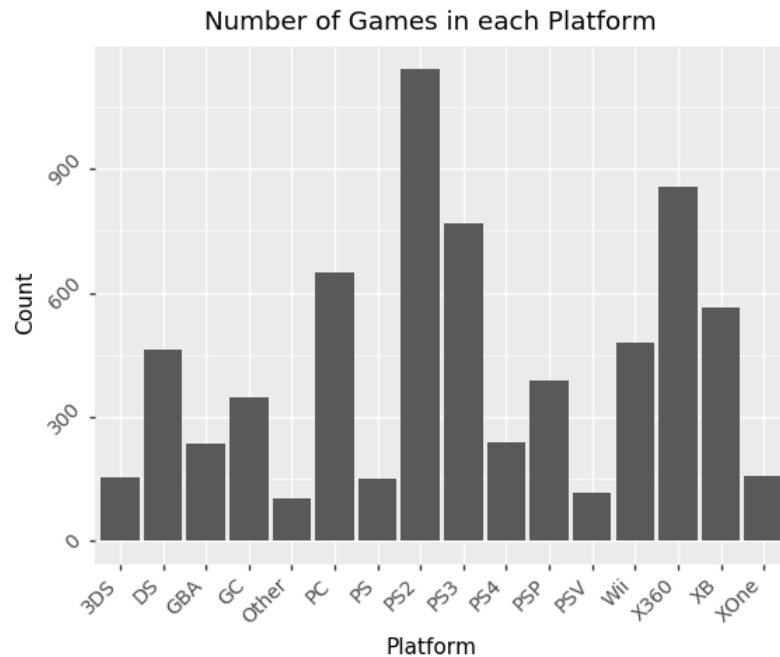
games low, users tend to also rate games low. When critics rate games high, users tend to also rate games high.



From this we see that the distribution of scores are fairly similar between the Critic group and the User group. We see a left skewed distribution for both groups. This shows that both users and critics give higher ratings more often than lower ratings.

How does each platform contribute to global sales?





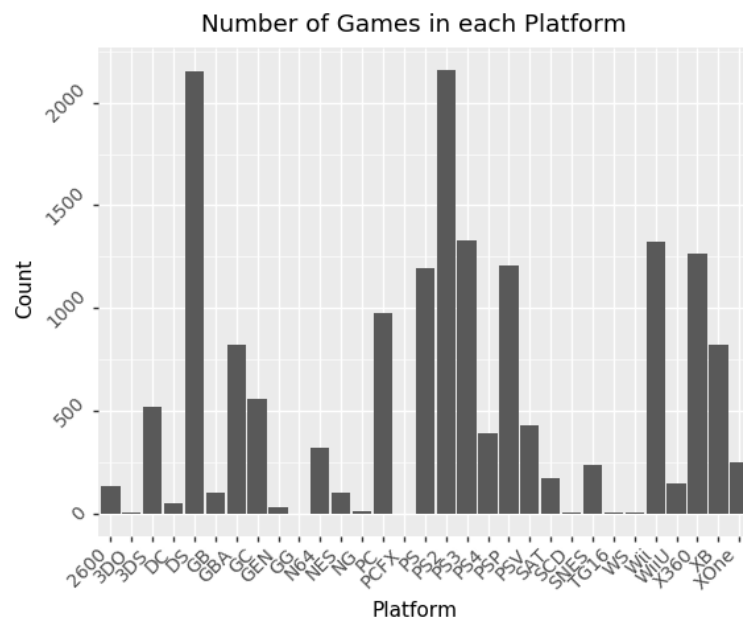
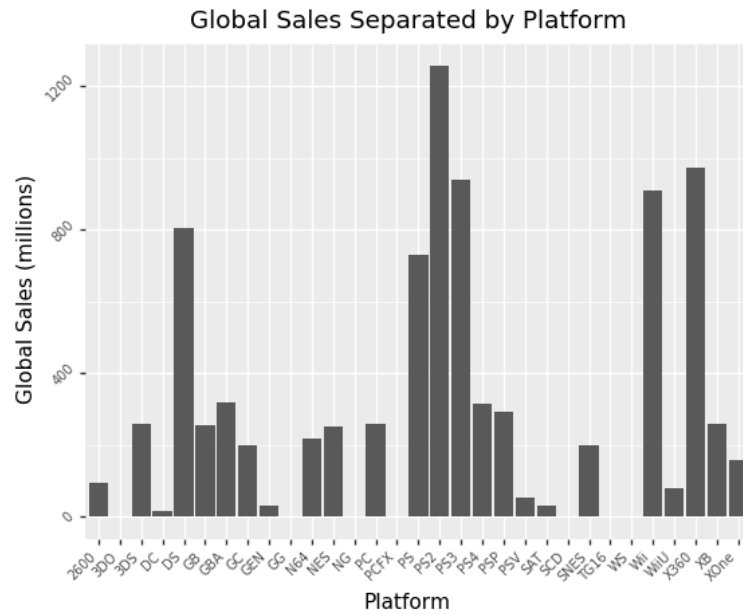
We see the lowest contributor to Global Sales is the PSV (PS Vita) and the highest platform contributor is the PS2 (Playstation 2). So despite that the Wii Sports game is the highest individual game contributor to Global Sales, PS2 games made more money than the Wii overall. This is most likely attributed to the fact that the PS2 has more titles included in the sum. The average Wii title has made more on average compared to the PS2. Part of this reason may be due to the Wii Sports game and that Wii does not have many other games in the data.

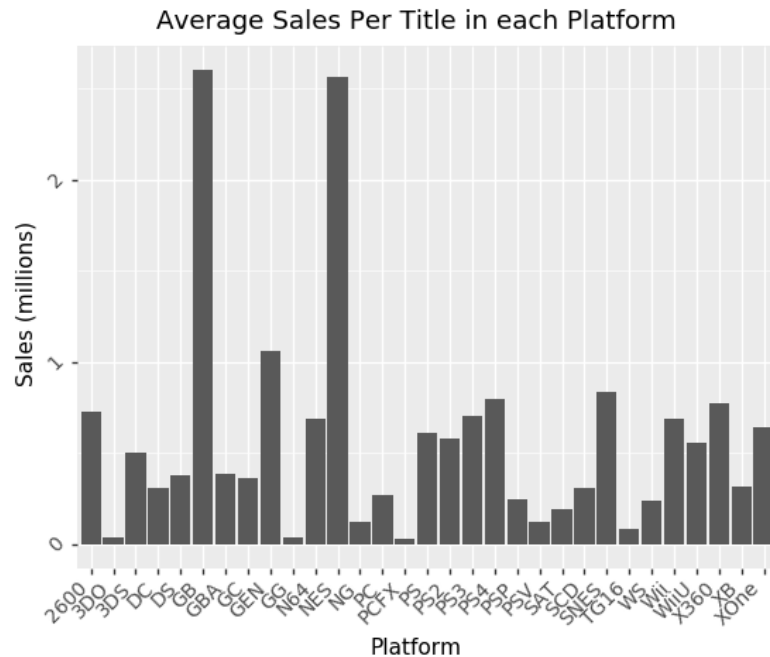
It is also interesting to see that the PS3 and PS4 hasn't surpassed the PS2 in sales yet, which may be due to that not enough years has passed since the release of those consoles. We see something similar happen with the Xbox 360 to the Xbox One. The Xbox 360 came out in 2005 while the Xbox One came out in 2013. The PS2 came out in 2000 and the PS3 came out in 2006. We see a six year difference for PS2 and PS3 and a 8 year difference for Xbox 360 and Xbox One. We also note that the average PS2 title has been less successful in terms of sales when compared to the average PS3 or PS4 title. Since this is the case, when there is a larger sample of PS3 and PS4 games in the data, we would see to it that PS3 and PS4 will eventually surpass the PS2 in global sales. Since PS4's count is still not as high as PS3's count, new titles added to the PS4 group will influence the average title's worth a lot more. A really successful title or several unsuccessful titles can easily influence it. Since the PS3 titles have over 3 times as many games and the average sales is about the same, we'd most likely see greater growth in the global sales for PS3 games. On the other hand, the average Xbox 360 title has been more successful in terms of sales when compared to the average Xbox One title. Assuming that the sample in my data is representative of the population, the Xbox One has not been as successful in terms of sales compared to its predecessor in both total sales and average sales per title. With this trend, the Xbox One would need to have a much higher count compared to the Xbox 360 in order to surpass it in global sales.

With this current dataset we can't really investigate why there is a large gap in Global Sales between these platforms other than the lack of counts in some platforms, but we can speculate.

One possible reason for lower Global Sales in the successor systems is technology advancement. This would also apply to average sales per title in each platform. What I mean by that is that it is possible to play games without buying them. Also known as pirating. There are ways around systems to jailbreak them and be able to download the game online and burn them to a disc and then run them on the system. This would lower the sales for games and overall lower global sales. We notice that Xbox 360 and PS3 came out roughly around the same time. Xbox 360 came out about a year later. These sales are greater in PS3 than Xbox 360. They also see similar decreases in Global Sales in their successor systems. Although this is more than likely attributed to a low count especially as seen in the PS3-PS4 because they have similar average sales per title. Xbox One however, may just be not as successful. Pirating may also be a reason why the PC average sales and global sales are both fairly low. It is much easier to pirate on the PC than other platforms because you don't need to do any extra steps other than download a file. That may be why we see PC being 4th highest count in the data but have fairly low global sales. This is attributed to a low average sales per game.

Another reason is that there is missing data on some of the newer titles, hence not giving these other platforms enough observations to include more sales as seen with the PS3-PS4 case discussed earlier. In the cleaning stage, roughly 10000 observations were removed. So it could be that these are all observations that belong in the newer generation systems which we can verify.





From this we see that the relative distribution of the PS2, Xbox 360, and the newer gen systems are the same for global sales. We see increases in Global Sales for systems such as the PS1, NES, DS which are all much older. It makes sense because users and critics are less likely to have played and rated much older games. So despite that the data of global sales were available, users and critics have not rated those games. If that is the case, those null values would be removed, which is our current dataframe.

In the average sales plot, see that with a higher count for PS4 games, the PS4 on average has been more successful in terms of sales per title compared to PS3 and the previous generations. Since we would unlikely continue to see an increase in older generation sales, the PS4 would eventually overtake the previous generations for global sales since it has the highest average per game. This can change if the newer counts of PS4 games aren't as successful, so it may be too early to tell, especially because it is not much higher than the previous generations. For Xbox 360 comparing to Xbox One, we still see the same trend. Xbox 360 is more successful compared to the Xbox One in both global sales and average sales per title.

The PC average game value actually decreased in the full data set which is not too surprising. Especially because I think smaller, not as popular games are more likely to be pirated due to not wanting to spend money on a game someone may not like. It also may be more convenient. A large number of these games are single player games. There are a limited amount of multiplayer games which require players to buy the game. These games may have already been included in the cleaned data set so that including more PC games lowered the average game value due to a greater ratio of single player games compared to popular multiplayer games.

In the very old generation platforms, we see two interesting points. We note that the gameboy and the original nintendo have extremely high average sales per game values. These are the one of the few

earliest systems and are regarded as classics. There was also little technology available at the time so pirating games was virtually impossible. We see the effect here on the older games have much higher sales per title than newer systems. The Super Nintendo may have also been fairly successful as well, but our data just may not include it due to the nature of the website the data is collected on.

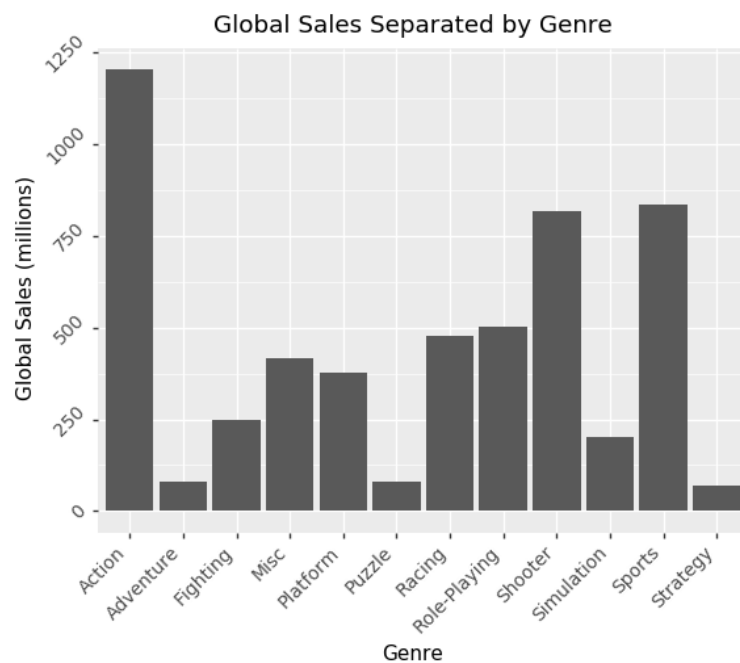
One other point sticks out from the full data set which is the number of DS games. It has fairly high global sales, but low average game sales. Each game individually didn't perform well in sales most likely due to pirating and jailbreaking, but it made it up with the number of DS games.

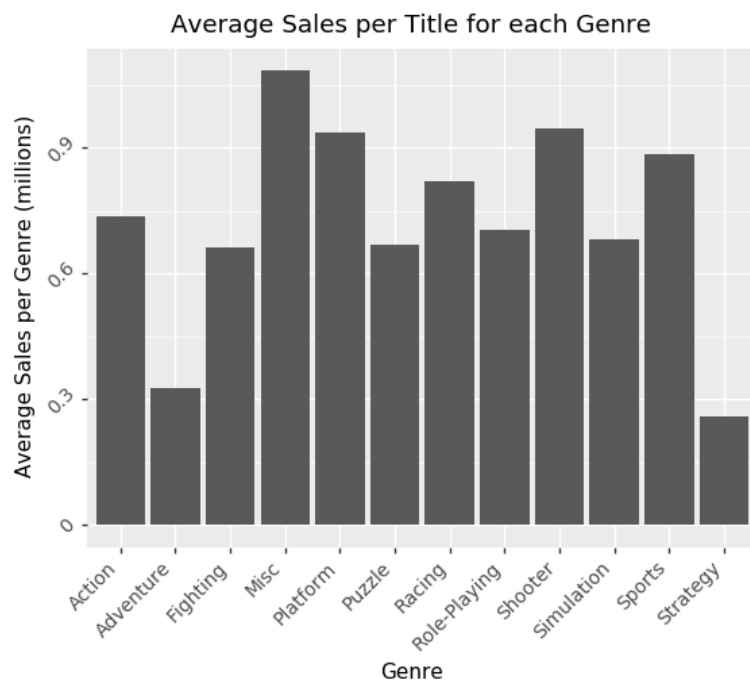
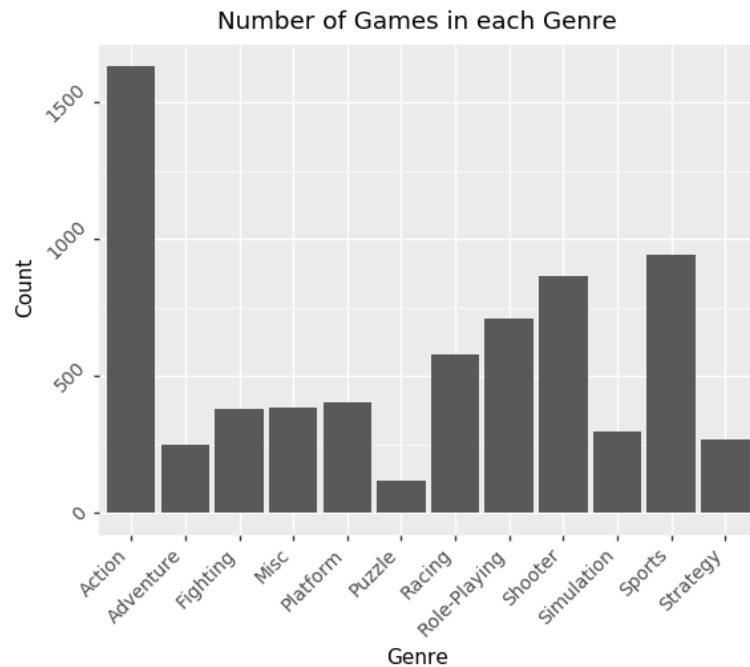
For the platforms with low count that were pooled into the 'Other' category in the cleaned data, the included platforms may have been the more popular titles on that particular platform. This may lead to a high average sales value for the 'Other' category. A closer inspection on what games were included could be insightful.

Since GEN(Genesis) has an average sales of 1 million with 29 values, we will examine this platform. The other low count platforms have low contribution.

After examination, I discover that the top 5 contributors to the Genesis system are all very well known titles. I.e. Sonic the Hedgehog 2, Sonic the Hedgehog, Mortal Kombat, Streets of Rage, and Mortal Kombat 2. So it is the case that a few well known titles are bringing up the 'Other' platform in our cleaned dataset average global sales than overall global sales. This is more than likely the case for other platforms as well with high global sales as seen in the case with Wii Sports.

How does each genre contribute to global sales?





From this we see that the Action genre is the highest contributor to Global Sales. But this doesn't necessarily mean that Action by itself is the most popular. Games aren't necessarily only one genre. It could be that Action paired with another genre contributes the most sales, but in this data set only one genre is listed for each game. For an example, Action may be frequently paired with Shooter. It is hard to imagine a shooter game that is not "action"-based.

It looks like the Action genre may just be an average of all the other categories assuming that most games are "action"-based. We can check this by averaging all the other genres.

I pooled together all the other genres except Action and pooled them together. I found the average to be 0.7284. The value for Action is 0.738135.

From this we see that is pretty much the case. The value for average sales for Action games is 0.738135 while our computed average sales for every genre except Action combined is 0.7248. This shows that it is very likely that the Action genre is including at least one of another genre.

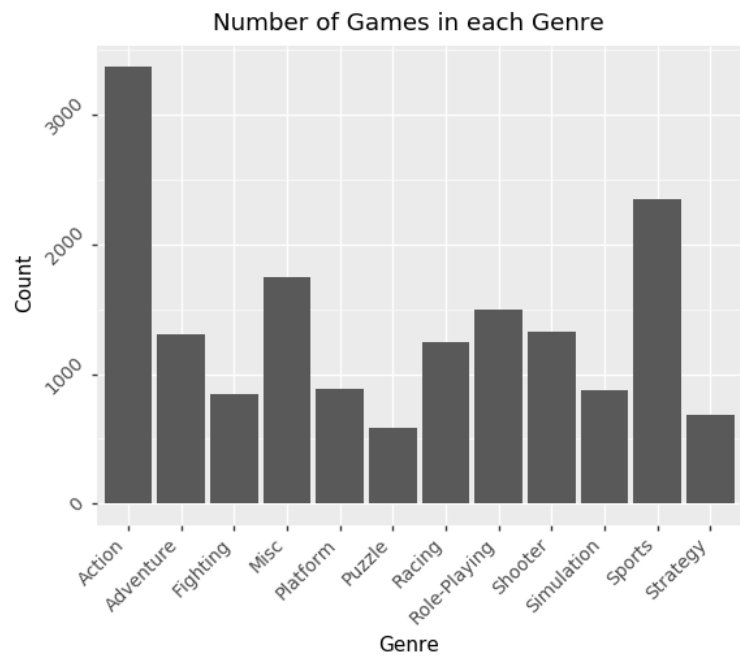
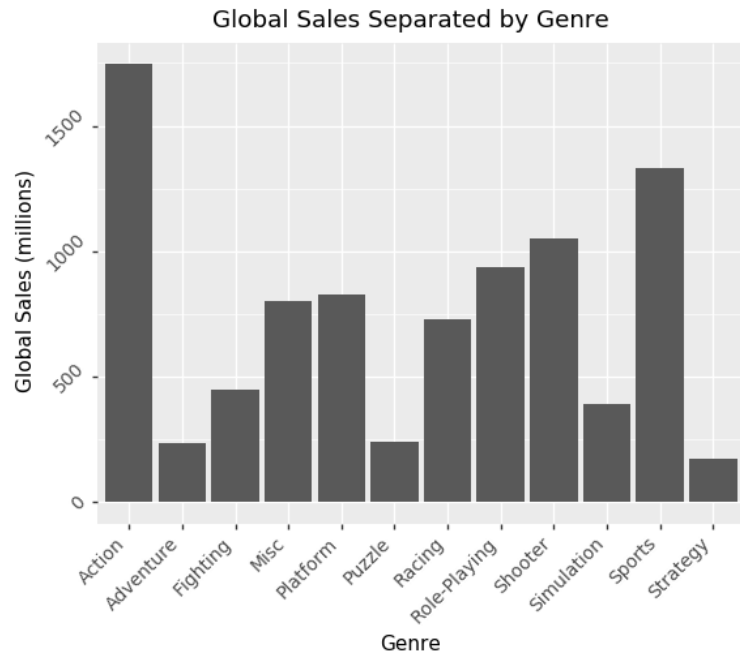
It is the same for other genres. To analyze the data better we would need to have better data, such as all of the genres associated for each game rather than a single one.

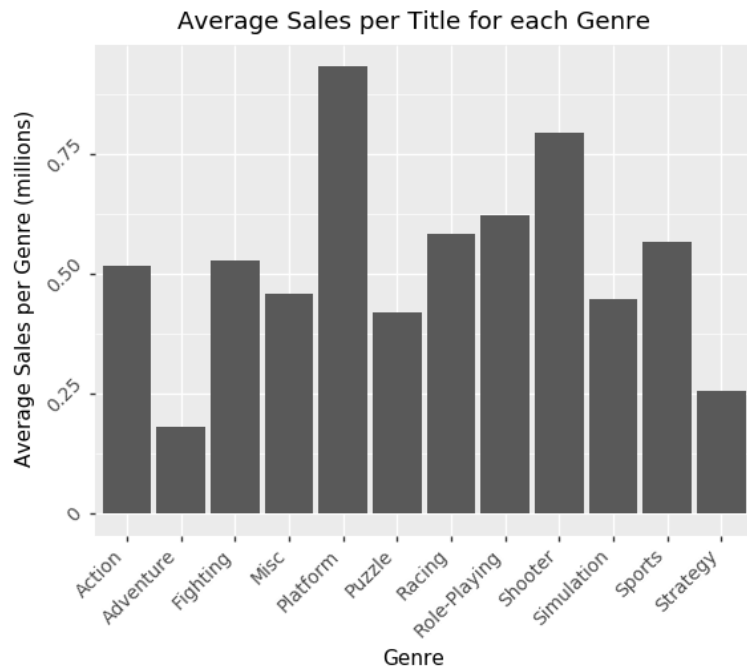
We see that the Misc genre on average has more successful games in terms of sales. The Misc genre is somewhat vague. We can examine what kind of games are under the Misc genre and its highest contributor by examining the data.

I discover that most of these titles are music games, 'Party' games which comprise of mini games, and brain academy. One of the reasons why this genre in particular has a high average sales value is because there aren't as many of these types of games around so more people want to buy them. Mario party is the first of its type of game and there hasn't been many other similar games which may cause the sales of these games to increase. Music or rhythm based games aren't as saturated as other times of games either.

I also notice that the majority of these titles come from the Wii, which we saw earlier have very high average game sales attributed to mostly a few titles with high global sales. This also comes to the fact that the Wii caters to a wider audience in general compared to other systems, so it ends in more sales.

We can look at the full data set to see if anything changes with more data.





The global sales looks relatively the same with an increase in Sports and Role-Playing games. Since many null values were in the Critics Count and User Count category, it could be that these Genres are not popular on the website in the population of users and critics that use that website.

The distribution of average sales remains mostly the same but there are two notable changes. The misc genre went down and the Platform genre went up for average sales. The number of misc games increased dramatically about 5-6 times more games were included in the full data. So my earlier discussion about the presence of party games being more popular and the misc genre itself unsaturated with games was incorrect. It still could be true that these games are the more popular games and just have a large number of unsuccessful games. We would need to look at a large sample of the unsuccessful games to know for sure which isn't the focus of this analysis.

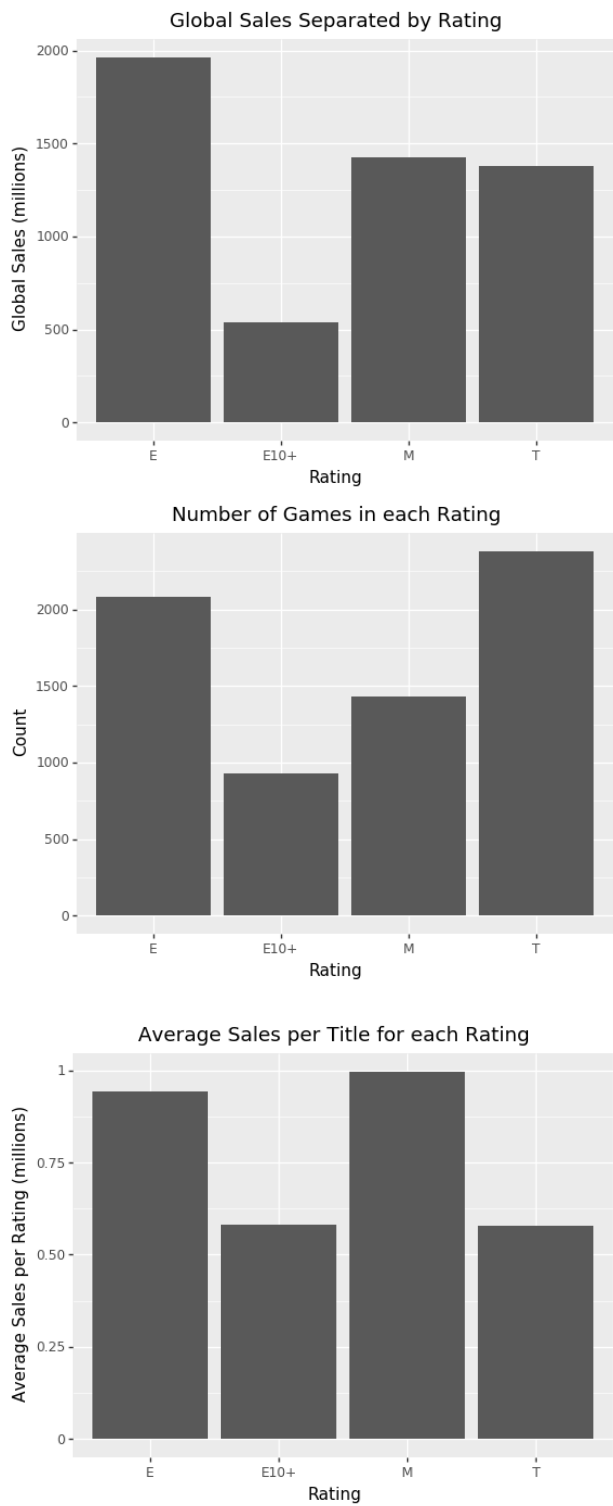
We can examine further what titles exactly are pulling the Platform genre upward. Since we observed earlier that the average game sales values for Nintendo and Game boy it is more than likely that these games are pulling the Platform genre up as well.

After examining the data, I discover that several NES titles are on the top 20 global sales list for the Platform genre. I also see a few titles from Wii and some of the other smaller count data such as the N64 and the SNES. This is also expected earlier as we saw with the Genesis data. The more popular games were included with the smaller count rather than the unpopular ones which also may be due to that these are old games. Not many people may know about the older games and don't bother listing them or have data on them.

I also see that most of these games are the Mario games. Since my analysis is not focused on subsetting the data further, I will not examine this data further. It could be insightful to subset the Platforming data

more so that it is separated into groups such as "Mario", "Donkey Kong", etc to understand if its the Platform genre itself having an effect or a specific franchise.

How does each rating contribute to global sales?



From this we see that the E rating contributes the most to Global Sales. This is attributed to a fairly high number of E rated games and having a high average sales value. The most saturated Rating is the "T" for Teens rating, but since they have a low average sales value it resulted in a lower global sales compared to the E rating. Since games rated for everyone are accessible to everyone, it has a larger population buying this type of game which results in higher global sales. T and M rated games have about the same contribution to global sales despite having very different average sales values.

The highest average sales value belongs to the M rating. This could be due to that the high contributing Shooter games are rated M. Since we know that Shooter games have a fairly high average sales value it comes to no surprise that this would have an effect on the M games. These games involve more violence and gore which would make them inappropriate for some audiences. There are also less M rated games compared to E rated games. This is either due to a lack of data for M rated games or that there are just more E rated games in general which would require additional data or research.

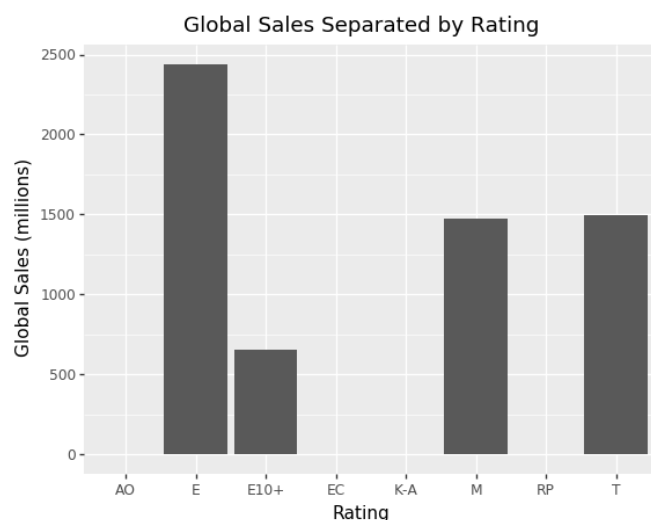
We examine some of the top contributors of the M rating to what titles are pulling it upward.

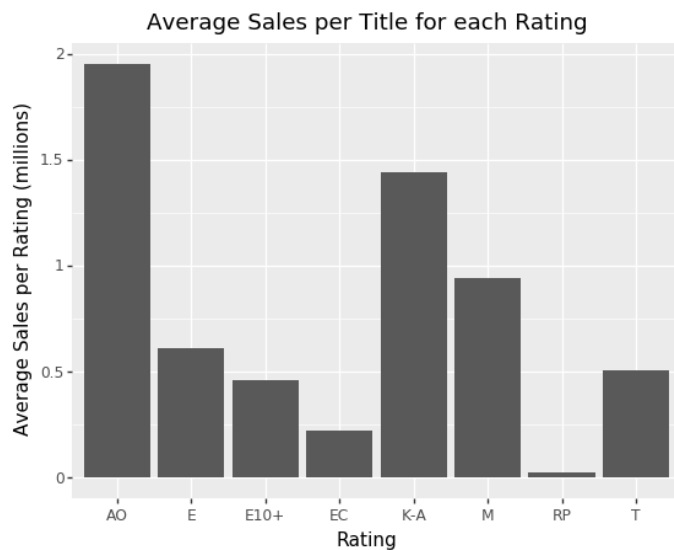
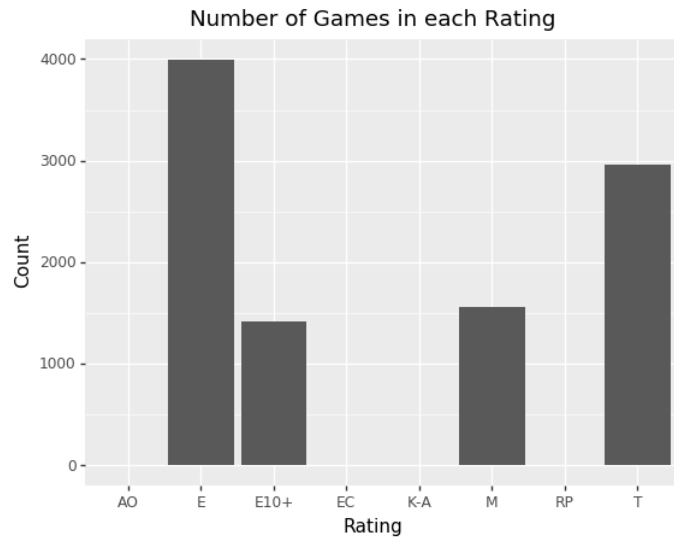
After examination, I discover that my assumptions were correct. The top 20 contributors to the M rating are all of the Shooting genre. Since we found that Action was sharing its genre with Shooting, it is no surprise to find a few games tagged as Action here such as the Grand Theft Auto series.

The E rated games are most likely being pulled up by Misc, Platforming, and Sports games including Wii sports since those don't involve as much violence which we can check.

I examine the data and do find that the top 20 includes those genres. The top contributors are also Wii so it lines up what we found earlier about platforms.

We can check the full data next.

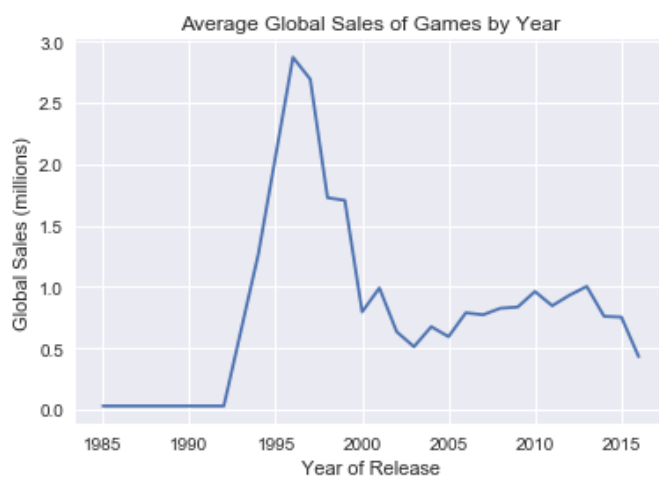
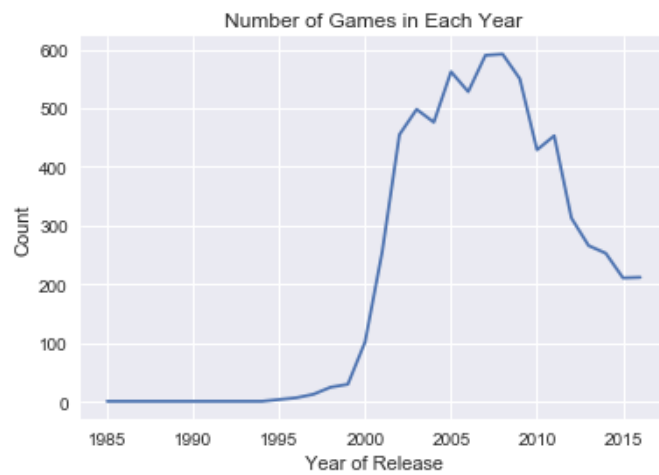
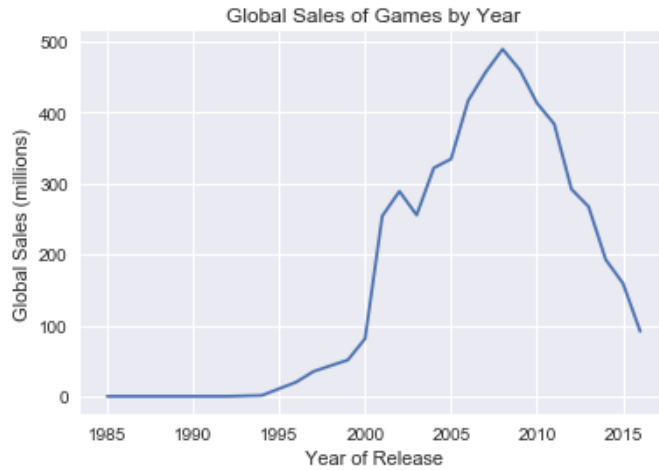




We see the same distribution for global sales and average sales in the full data. Despite having 10000 less observations, we have the same shape. This could mean that this is the "true" distribution of global sales by rating, meaning that there is a sufficient sample size in each rating we have in our clean data. So our earlier possibility of either having less M rated games compared to E rated games or just not having enough data leans more on the side of just being less M rated games in general. This means that E rated games in general will have a higher overall global sales, but M rated games tend to have more popular games because of. Assuming that high individual contributors to global sales means that it's popular.

Since there is only 1 title in AO, EC, and RP, these values are inflated, especially AO. As we know this value is GTA San Andreas which we re-categorized as M. The EC category has null values so it was dropped and wasn't in the clean data.

Does years since release correlate to higher global sales?

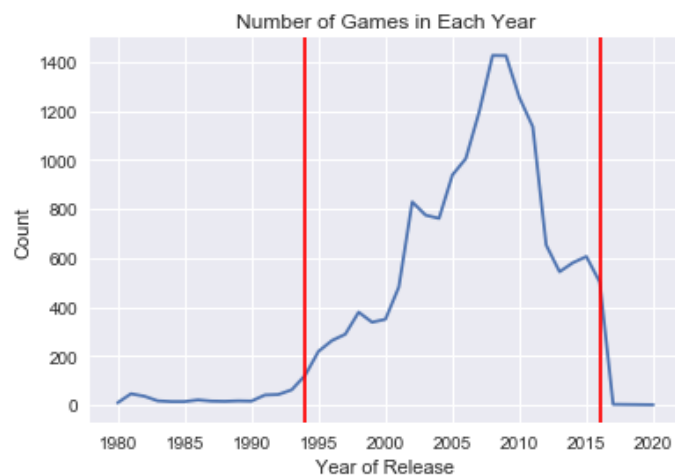


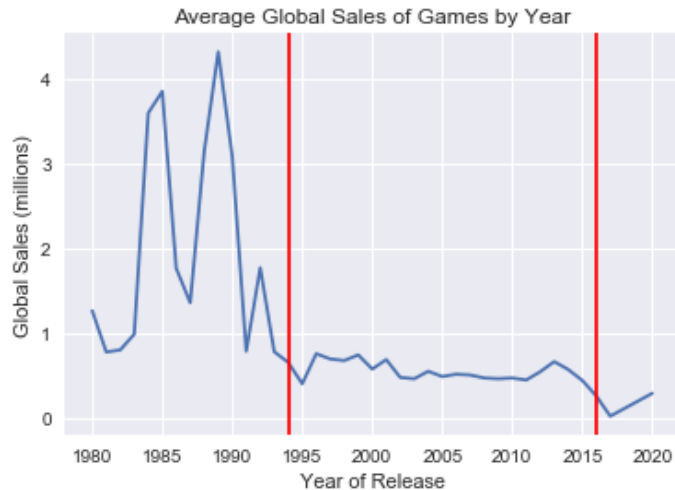
Based on total global sales alone, we see that the games released around 2007 and 2008 has seen the highest overall sales. This is also about the time where the count of games were the highest. So this

peak could just be attributed to a larger abundance of data for these years. We see the average global sales of games peak around 1996. It is likely that these are just those popular games for the older systems such as the NES, SNES, etc inflating the average values upward. We see under the number of games released in 1996 is very low.

There doesn't seem to be any linear relationship between years since release and global sales overall or the average. If we compare the data that have enough counts from year 2000 and onward, we see a non linear relationship in overall global sales and average global sales of games seem to go down as the years progress. This means that our years_since_release variable would have a slight positive correlation. This could be attributed to the technology/pirating explanation from earlier because it is easier to pirate games and play them illegally nowadays.

It may be more insightful to look at the full data.





The vertical red lines indicate the years of data that we have that have enough data points. Looking at the full range of data, we see very inconsistent values for the average global sales before 1994. These are definitely due to those NES and GB games. For global sales, we see a steady positive increase in global sales up until 2007-2008 which is the same in our clean data set. This may indicate that games released around 8-9 years ago see the highest peak in global sales. This is also when we had the highest count of data so that may not be accurate. But it could be that a year after this data was collected, more data on games released in 2009 were made available which make it have as high as count or higher than 2007 or 2008. Without more data. If we compare our past data, it has always been the case that global sales for games were greater than the previous year up until 2007 or 2008 in general.

In the average global sales plot, we see a similar decrease in global sales as we did in our clean data that I speculate is due to the technology/pirating effect.

In general, we can't be sure that the low count or average game sales observed above are caused by unpopularity or just a lack of titles in those genres/platforms etc on the website in which the data was collected from except in the case of the Ratings because we found that the distributions of the full data set versus the clean dataset were relatively the same.

There are also a few outliers that are mistakes in the data. We have 3 2017 games and 1 2020 game in 2016 data which has to be an error, unless it counts preorders. But I find it hard to believe that preordering a game 4 years in advance is possible. We can check exactly what 4 games these are to see if there is an error in the data. Since our clean data doesn't have these points it doesn't matter for the analysis but will check nonetheless.

After comparing the names and cross checking with quick google searches, I found that they were incorrect values. Imagine: Makeup Artist was tagged as 2020 but it was released in 2009. While these data points don't matter in the analysis, the presence of these errors could hurt the analysis or our previously found results if these errors are prevalent in the data set. There's no way to know if there are more errors without cross checking every single title manually. But if there are more errors like this and

a lot of them for other information such as year of release, global sales, critic score, etc then we have more problems other than underrepresented categories.

Are there correlations between our possible quantitative predictors?

I examine a correlation matrix between all the quantitative variables. Other than our previously found moderate-strong correlation between User score and Critic score, there are no other correlations between critic score, user score, critic count, user count, and years since release.