

Video Game Analysis: Trends and Insights

Daniel Ethridge¹ Atharva Patil², Deep Shukhla³, Sujith Battu⁴

University of Colorado, Boulder

Author Emails

¹Corresponding author: daet2304@colorado.edu

²atpa5127@colorado.edu

³desh3965@colorado.edu

⁴suba2075@colorado.edu

Abstract. This project examines patterns and trends within the video game sector, utilizing information from diverse sources to gain insights on game popularity, sales statistics, player demographics, reviews analysis and genre success. The study's goal is to pinpoint crucial factors that impact the success of video games and to represent these patterns through graphs and charts. This report offers a thorough summary of the video game market's current condition by gathering and organizing pertinent data, facilitating well-informed decision-making for stakeholders.

INTRODUCTION

The video game industry has quickly transformed from a specialized market to a worldwide force, propelled by technological progress and a variety of player demographics. This project examines patterns in the video game industry, with a specific emphasis on sales figures, genre trends, and player involvement.

Our goal is to discover important insights for developers and marketers by analyzing data from different sources like sales numbers and critical reviews. What types of genres are the most prevalent in the market? What is the relationship between game ratings and sales? By carefully gathering, scrubbing, and presenting data, this document offers a thorough examination of the gaming industry, emphasizing the correlation between consumer actions and market patterns. It is essential to grasp these dynamics in order to successfully navigate the future of the industry.

DATA COLLECTION/PREPARATION

This section discusses the procedure and measures we choose to collect data and the sources of our database. This section also details the issues and errors while preparing the data for visualization process.

Data Sources

1. SteamDB: Sales and user reviews for PC games.
2. API Access: The data utilized in this work comes directly from Steam via an API. There are three main API calls utilized to collect the data:
 - <https://api.steampowered.com/ISteamApps/GetAppList/v2/>. This is used to acquire all Steam app IDs.
 - <https://store.steampowered.com/api/appdetails?appids=<APPID>>, where APPID is the steam application ID number. This is used to collect details about a single Steam application.
 - https://store.steampowered.com/appreviews/<APPID>?json=1&num_per_page=100&cursor=<CUR>&filter=recent&purchase_type=all, where CUR is a string representing the next page of reviews if there are multiple pages.

Fig.1 is the image of the json file of the collected data.



FIGURE 1. Image of json file of data

Data Schema

Here is the schema of the final database we collected to make sure o the structure of the dataset.

Fig.2 Below shows the clear schema of the database.

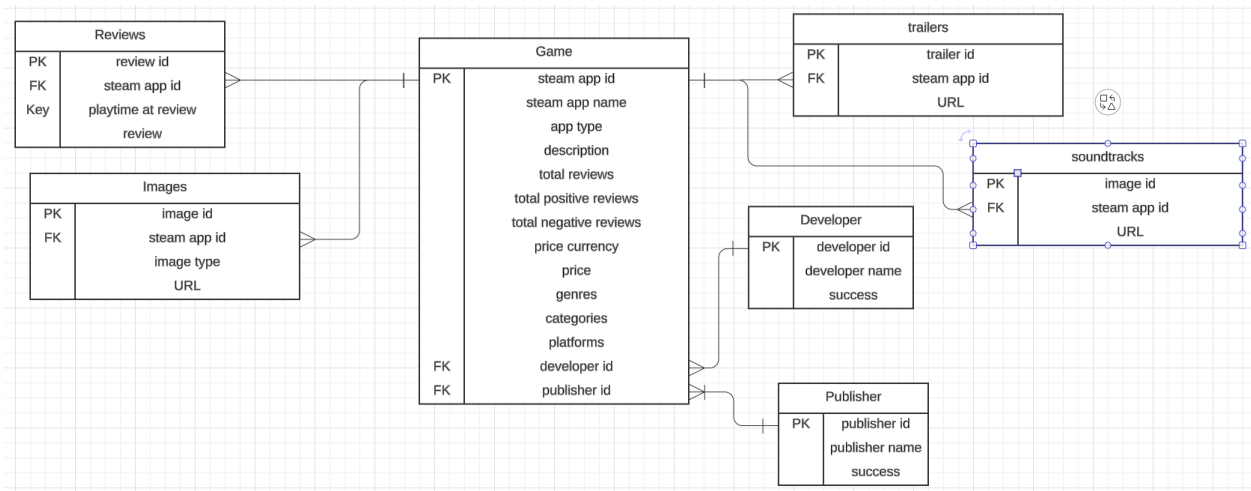


FIGURE 2. Data Schema

Fig.3 shows the exact data file we collected.

Tools Used

- Pandas: Used for data cleaning and manipulation in Python.
- NumPy is used for performing numerical computations and managing missing data.

Data Visualization

Plots used:

1. Bar Charts
2. Line Graphs
3. Heatmaps
4. Pie Charts
5. Bubble graph
6. Scatter plot

Visualizations:

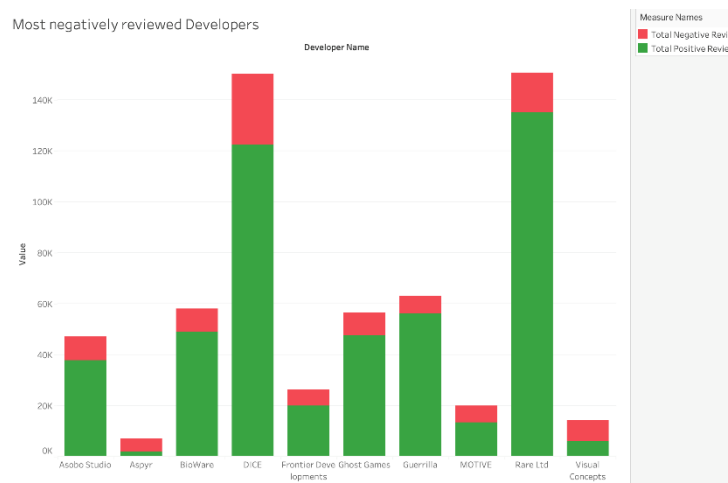


FIGURE 4. Negatively reviewed reviews 10 developers

The [Fig.4](#) below show initial insight into the data. Below each is a caption describing them. The x axis is the top 10 most negatively review game developers and y axis is the total positive and negative reviews. We can see that rare ltd and DICE are the two most negatively reviewed game developers and that most developers have almost nearly equal positive reviews as well.

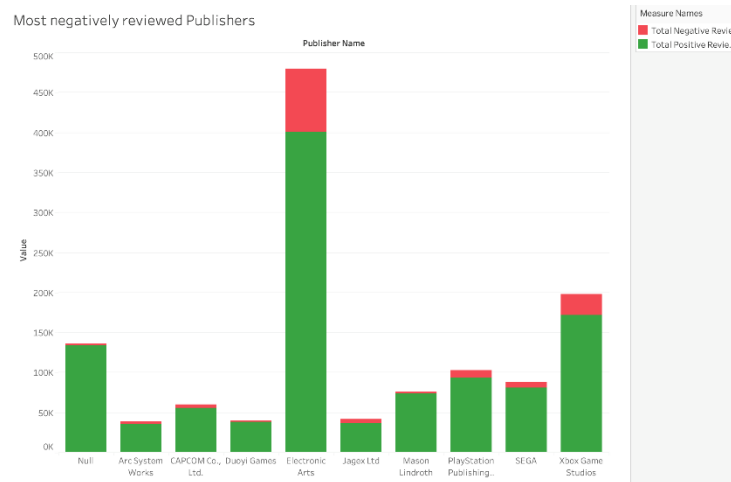


FIGURE 5. Most negatively reviewed 10 publishers

This graph in Fig.5 displays the top 10 most negatively reviewed game publishers and the number of positive (green) reviews and negative (red) reviews they get. We can observe that electronic arts has the most negative reviews as well, and that most of these publishers have comparable number of positive reviews to the negative reviews. Implying that people tend to have mixed opinions on game publisher games.

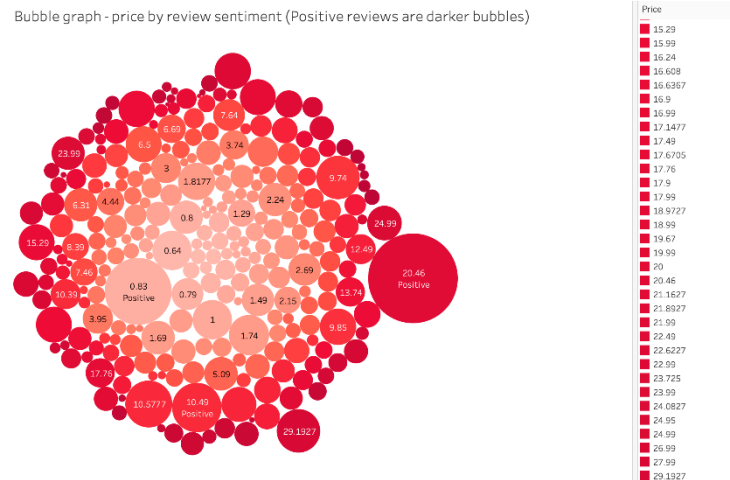


FIGURE 6. price vs reviews

The bubble plot Fig.6 shows the price vs review sentiment (which is the ratio of positive and negative reviews). The bubble size denotes the value of review sentiment (larger bubble implies positive review). While the colour denotes the price, more red implies more pricey game. We can observe that more pricey games tend to have more positive review sentiment, and vice versa. Which can be confirmed as less pricey games tend to be made by cheaper indie game developers who cannot make the best of games.

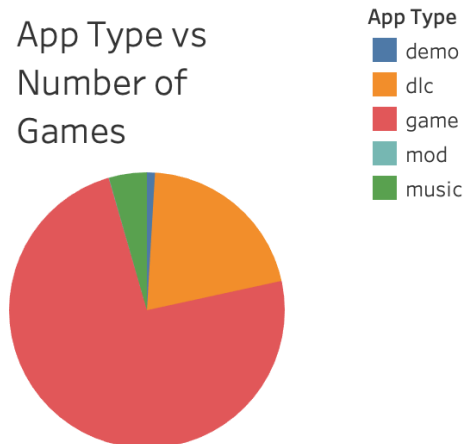


FIGURE 7. Game types

The pie chart [Fig.7](#) shows the app types like game, demo, dlc, etc vs the total number of games. We can observe that the game app type is most common followed by dlc and music.

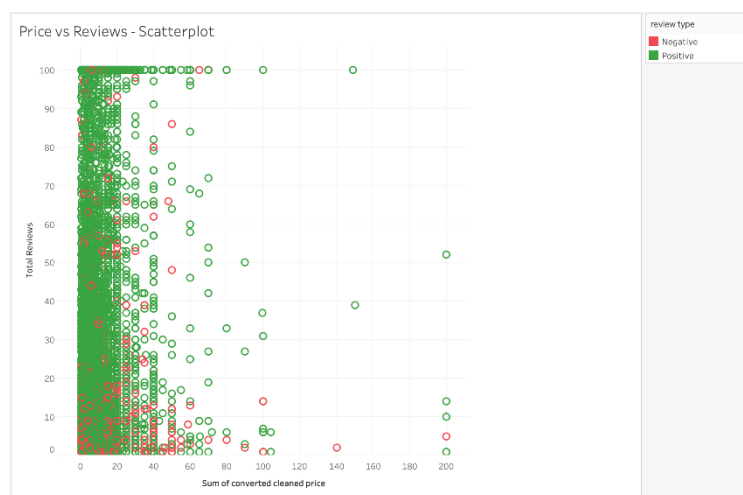


FIGURE 8. Price vs Reviews

The graph in [Fig.8](#) total number of reviews for each price range of each game. Red are net negatively reviewed games and green and net positively reviewed games. We have removed outliers to make the graph more digestible. We can see that most games are positively reviewed, only a few are negatively reviewed. This mean the data is more skewed towards positive reviews. We can also observe that the data is more concentrated towards the left end, meaning games tend to be less pricey. But on the y axis the total reviews are quite equally spread apart, meaning there is no visible trend with games and number of reviews.

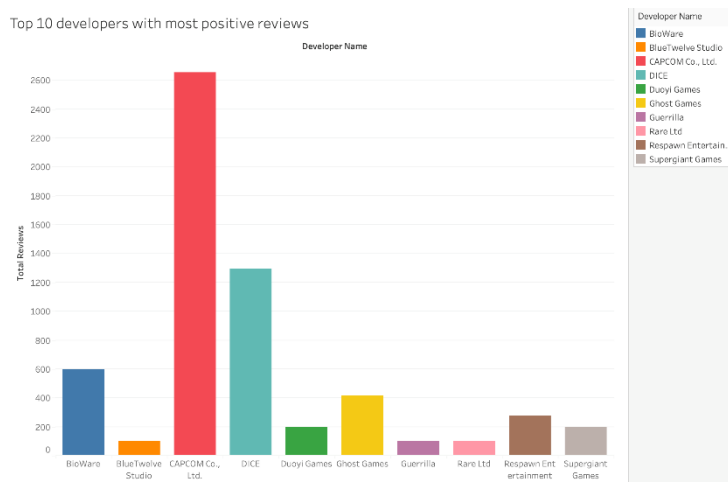


FIGURE 9. Most positively reviewed 10 developers

In above graph Fig.9 X axis is the top 10 game developers based on positive user reviews. Y axis is the total reviews received by each developer. We can see that capcom co has the highest reviews and is followed by dice , but there is a huge difference between the two. We can see that BioWare is the most positively reviewed developer.

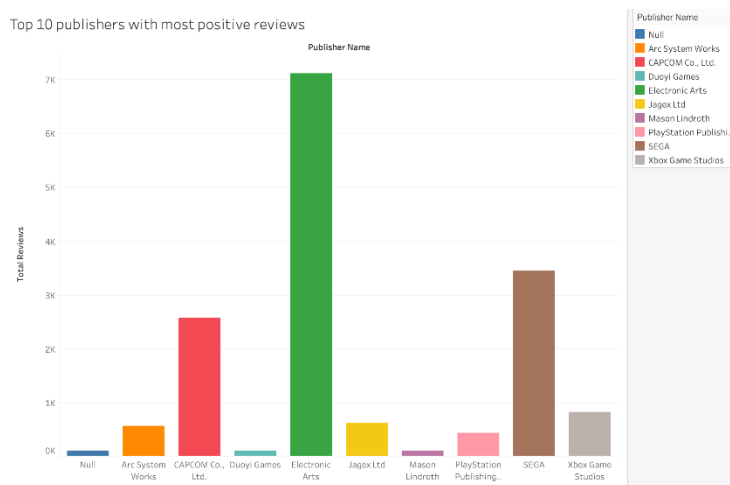


FIGURE 10. Most positively reviewed 10 developers

Above plot Fig.10 has X axis as the names of the top 10 publisher names, which were found using the most positive reviews. Y axis has the total number of reviews. This graph shows the top 10 publisher names, and that electronic arts has the highest number of reviews, and arc system works has the highest number of positive reviews.

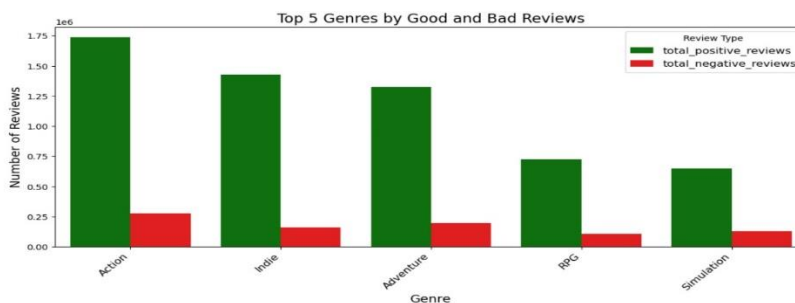


FIGURE 11. Genre analysis (Good vs Bad reviews)

Fig.11 X axis has the names of the top 5 game genres, which were found using the most positive reviews. Y axis has the total number of reviews. This graph shows the top 5 genre names, and that action is the most liked game genre and these have very little negative reviews.

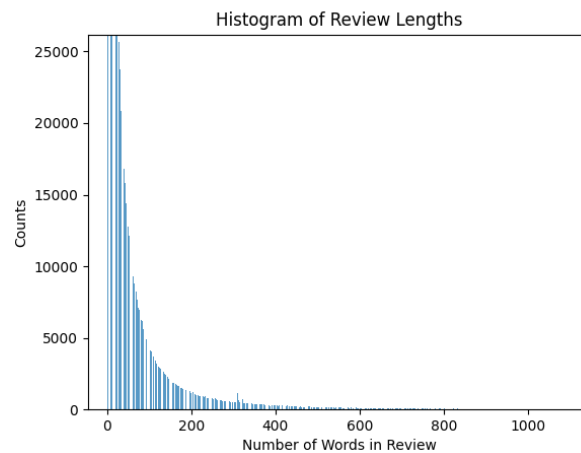


FIGURE 12. Number of words in reviews

The x axis in Fig.12 shows the number of words in a review and the y axis shows how many reviews had that many words. Perhaps unsurprisingly, words with fewer words are more common while reviews with more words are less common.

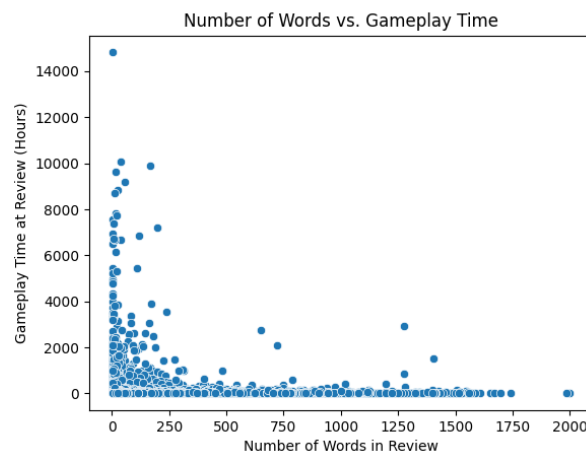


FIGURE 13. Number of words vs Gametime

The Fig.13 shows number of words in a review plotted against how long someone had played the game before leaving a review. Intuitively, one may believe that longer, more indepth reviews would come from people who have played the game longer. However, this is not the case. Longer reviews typically come from people with a shorter amount of playtime. This could present interesting ideas for analysis later on.

Variables in the dataset are, reviews, price, developer, publisher, genre, currency, platforms, images, number of reviews, positive reviews, negative reviews, brightness, contrast, saturation and hue.

1. CLIP MODEL: CLIP is a multimodal model, proposed by OpenAI, that aligns textual and image embeddings into a joint space. It was trained to predict which caption fits best for an image among candidates, allowing it to learn associations of visual and textual data. Trained on diverse image-text pairs, CLIP can execute tasks such as image classification, retrieval, and zero-shot learning without fine-tuning specifically for individual tasks. Using a dual encoder architecture, it processes images through a vision model-e.g., ResNet or ViT-and processes text using a transformer to produce embeddings to be compared. Such flexibility and generalization make CLIP a very powerful model for different visual-linguistic applications.

WHY WE USED IT:

We utilized the CLIP model in our project to examine and find associations between game images and their textual descriptions. The ability of this model to create a common embedding space for text and images allowed us to efficiently perform tasks like image classification and retrieval, helping uncover patterns, understand context, and improve insight into the relationships between visuals and game descriptions.

DATA CLEANING:

How data looked before cleaning:

[illegible]

After cleaning:

steam_app_id	cleaned_description
22135	br game party skills unique team idle gems cla...
22714	br risk system deluxe content original light e...
14642	stage healthy box vr fitness game people willi...
21109	lord zedd character power rangers battle grid...
1045	chess game games remix world rules new play sh...
20667	br roman romans evil britons level welcome sea...
2714	unique steam extras strong malice platformer a...
16973	borg instructors mountain faulty apprentice dl...
5312	game heroes players battle deck zodiac abiliti...
10495	maze steam extras shadow reborn ego shape expe...
9939	electronic components power steam extras robot...
15236	format
2533	steam extras person night fang mysterious new ...
14438	chemistry steam extras molecules simulation re...
7047	dlc valentine xi churchill
14618	br treatment dungeon health way tests obstacle...
13771	people horror game city person thrill feelings...
24033	coffee run game upgrade graphics ui gameplay c...
4059	love note steam extras notebook japanese roman...
7346	steam game mini keys games points characters k...
23391	br artefacts player mode maldrin journey game ...
20063	supergiant god like steam extras strong underw...
15538	br steve human silent hill game average day in...
20807	picture bylo quick horror game experience back...
10436	bear ears cave lantern footsteps hunted large ...
18571	br characters battlers rpg maker hero bust pac...
18092	steam extras game explosions floor weapons str...
7819	grappling players physics hook movement active...
20363	steam extras level action strategy battlefield...

This would serve to clean game descriptions of noise like irrelevant words and characters, hence making the text data concise and meaningful. The preprocessing step helps in improving the quality of the textual embeddings that are generated by the CLIP model, therefore improving its ability to relate game descriptions with their corresponding images accurately. Cleaned data reduces computational complexity and improves overall model performance.

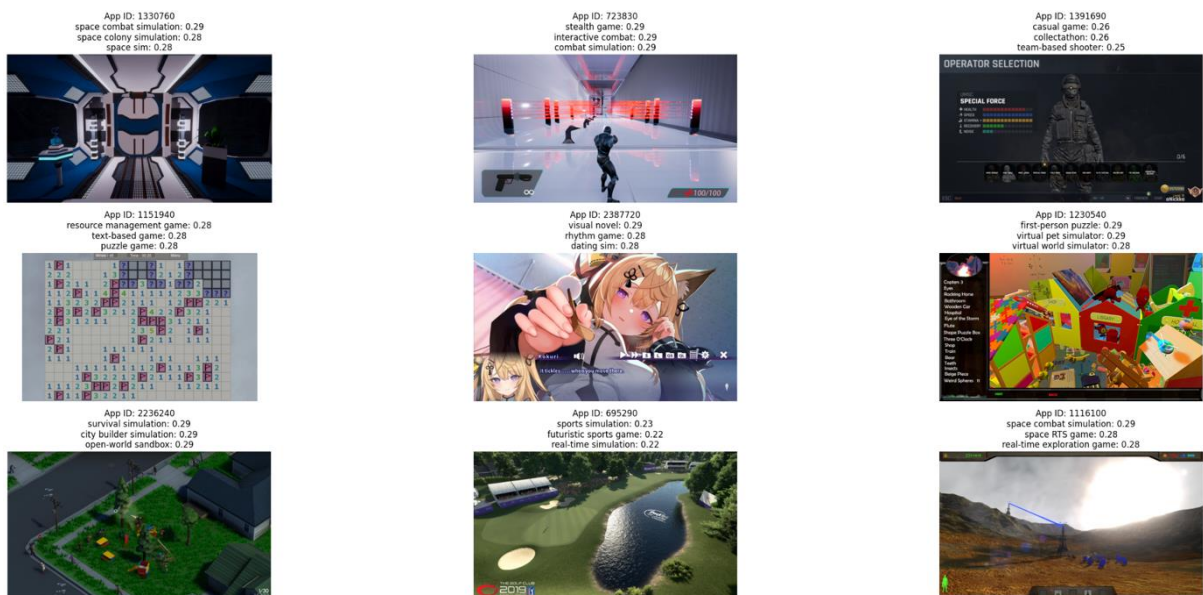
Cleaning was also done for image dataset, such as removing 404 error URLs and only filtering screenshots for more consistent training.

MODEL EVALUATION (using similarity scores):

```
Similarity scores:
[[0.1638 0.2218 0.1312 ... 0.169 0.1782 0.1259]
 [0.1913 0.2017 0.1605 ... 0.1886 0.2017 0.1562]
 [0.22 0.232 0.2188 ... 0.2515 0.23 0.208 ]
 ...
 [0.2034 0.251 0.2646 ... 0.2292 0.2242 0.2021]
 [0.2034 0.251 0.2646 ... 0.2292 0.2242 0.2021]
 [0.2034 0.251 0.2646 ... 0.2292 0.2242 0.2021]]
```

The similarity scores indicate how well text embeddings (e.g., game descriptions) align with image embeddings in the CLIP model. Each value represents the similarity between a text-image pair where higher scores indicate stronger associations. Rows correspond to text embeddings, columns correspond to image embeddings, and together this allows for an evaluation of CLIP's effectiveness at multimodal alignment. The scores are cosine similarity, between 0 and 1. Since we are finding custom descriptions for completely new images, even a score above 0.2 is good.

OUTPUT:



The CLIP model is very effective in describing what each game image is. In a matrix of 10 images, it was able to correctly describe each image. Thus proving CLIP is an effective image-to-text model. The examples of results

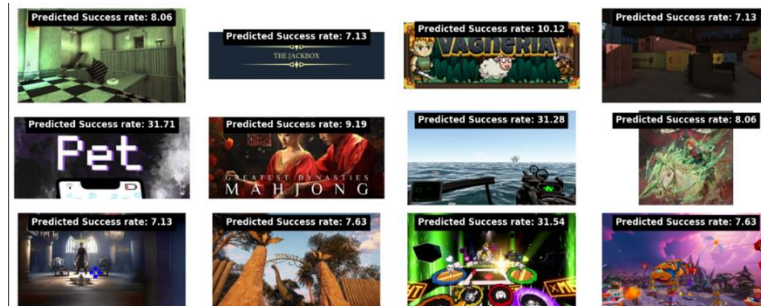
To train the gradient boosting model, we had to create a new dataset from the image URLs, which contained the app ids and image values such as brightness, contrast, saturation and hue. We made this so that we can find patterns in underlying aspects of the images and game success.

MODEL EVALUATION (using Mean Square Error):

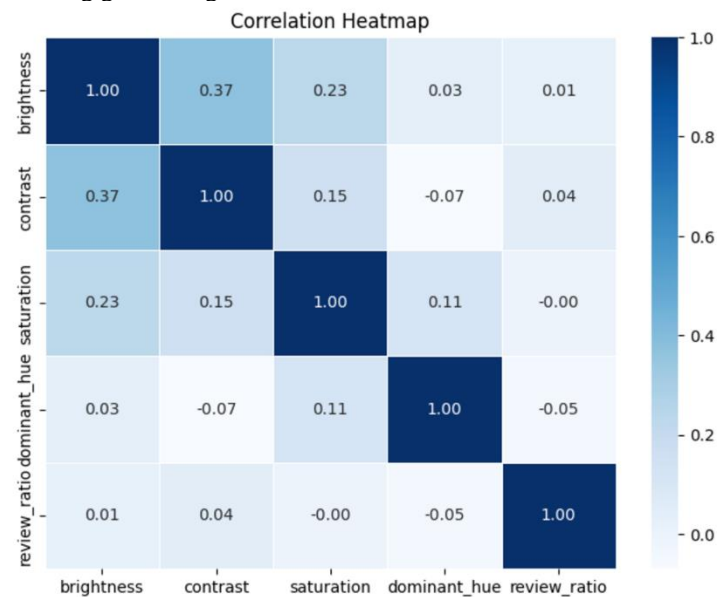
Gradient Boosting Mean Squared Error: 130.68720097122053

We use Mean Squared Error to calculate an average of the squared difference between predicted and actual values. The lower the MSE, the better the accuracy of the model. At a value of 130, the model performance is good, showing that the predictions are relatively near to the actual values; thus, effective learning takes place.

OUTPUT:



As we can see, the model predicts the success rate, which is the ratio of number of positive reviews to number of negative reviews. It rates higher success rate to games with more interesting game image screenshots and lesser success rates to dull or confusing game image screenshots.



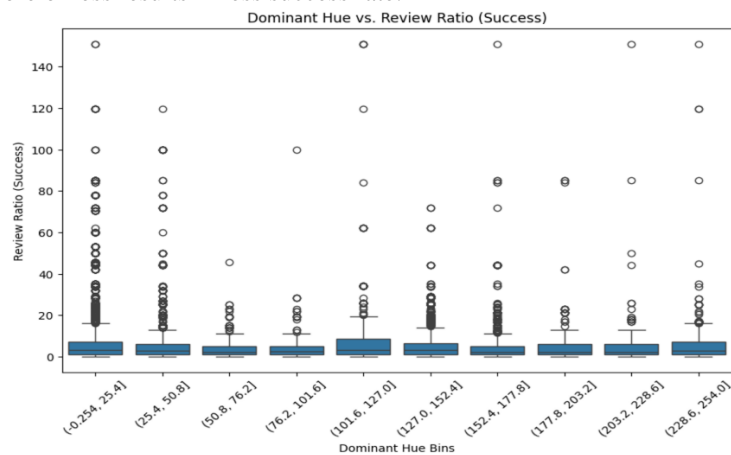
This is a correlation heatmap, showing correlation between various features, as we can see most features are unrelated except contrast and brightness which have some correlation.



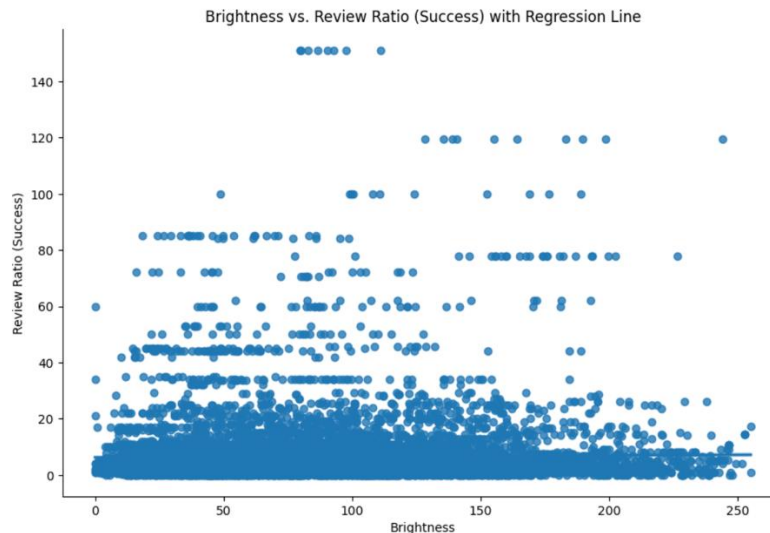
This scatterplot shows the relationship between review ratio (success rate) and saturation of images. As we can see the more saturation implies less review ratio, implying more saturated images might be worse for game success.



This scatterplot is the relation between contrast and success rate. As we can see a contrast between 50 and 80 is the sweet spot. Anything more or less results in less success rate.



This graph is the boxplot relationship between hue and success rate. It shows that all colors have similar success rate and there is no correlation between color/Hue and success of a game. It means that consumers don't prefer one color or the other. Hues between 101-127 (green-yellow) have slightly higher success rate, but nothing can be said with certainty.



This graph shows the relationship between brightness of game images and success rate of those games. As we can see, there is no clear relationship between brightness and game success.

WHAT WE LEARNED:

From this model, we extracted a few important points about how far image features such as brightness, contrast, and saturation will reflect success rates of games through reviews. Thus, the model predicted an optimal range of 50 to 80 for contrast for high success rates, while very low or high values of contrast will negatively affect the review score. Interestingly, saturation is negatively correlated with the success rate, which might indicate that highly saturated images could turn players off. On the other hand, hue did not correlate significantly with success, which would indicate that the consumer's preference is not strongly connected with specific colors. Also, no obvious influence of brightness could be found on the success rate, meaning it's not a determining factor. The model performs well with an MSE of 130, showing that it learns the relationship between the image features and game success quite effectively. This shows the importance of specific image attributes such as contrast in predicting game success.

3. ASSOCIATION RULE MINING:

WHY WE USED IT:

We applied association rule mining to the reviews of the games. We looked at Steam Apps that were classified as a "game" (in other words, we didn't look at DLC or other non-game Steam Apps), and we looked at differences among the top 10 genres. We filtered game reviews by the ratio of positive reviews to negative reviews of a game. If a game had more positive reviews than negative, it was considered a positive and vice versa. We acknowledge that this is not the best method, especially considering Steam reviews are tagged as either positive or negative. A bug in code caused that review tag to not be collected. Grouping reviews by how positively overall the game they are associated with is the best approximation.

ARM requires transaction data. To create the transaction data, we first extracted the necessary reviews according to liked and disliked games and genre. We removed any line break characters, removed capital letters, punctuation, numbers, non-latin characters, duplicate words, and a set of words (such as "game", "ever", "it", "etc") that provide no meaning. Afterwards, we saved the cleaned words as basket data. The two images below show messy data (top) and clean data (right).

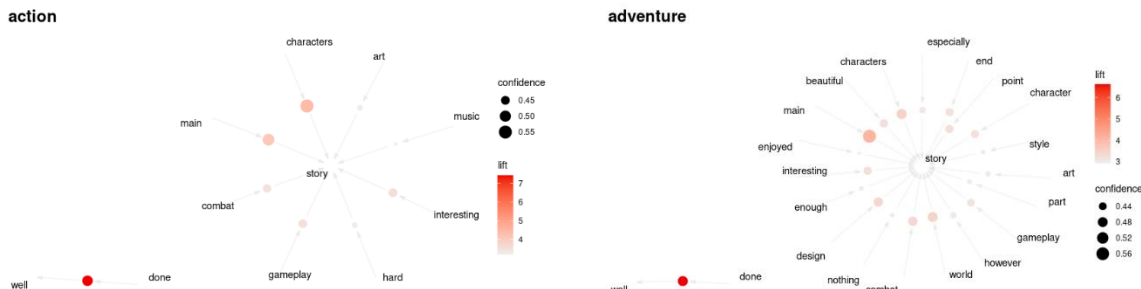
Before:

	steam_app_id	[projective_at_review]
560029	1405180	563 don't waste your money!!!!!!!!!!!!!!
560626	1405180	765H needs LOTS of work and I think this games dead in the water. Don't Buy IT!!!!!!!!! DO NOT BUY
560627	1405180	154A big waste of money they not even working on it should be removed by steam
560631	1405180	This game has got to be the worst game i've ever played on Steam. The game itself is very buggy, but the worse part of it is that I paid \$29.99 for this. Anyway, with their recent update to adding multiplayer, they've changed the map, and it's both this game was good after 5 min can do anything
560636	1405180	190there is very little control in this game. Potential is there, but they are way too slow getting things done. I like playing as a wolf, but there, again, nothing much to do. I wasn't able to rescue animals and then the "real" guy stopped letting me eat.
560638	1405180	447
560643	1405180	142/At this point game still has too many bugs. sometimes you can cash items, sometimes you cant. I have never ask for a refund, but had to on this one.
560647	1405180	20Needs a lot of work but it has potential. New hunter imagine that you gain UTTERALLY decided to do that RECENT KNOW, not bringing any of your belongings when you except your pants and t-shirt, and you asked an REDUCER prior to intended you, try you or two some random Canadian!
<hr/>		
That's the game.		
<hr/>		
At first, your dream has come true - nature around you look realistic and incredibly pretty, albeit a bit repetitive since you will soon discover that the forest you're being placed out in - is some kind of secret cloning research government program ugh		
Run faster - RUN		
And that's exactly what you will be doing for hours, days even - running around in a forest, like Gunga Himself, and no one knows where you get that world class running power from, because you don't run out of stamina at all. You run like a champion		
So, I run! Now what?		
If you've managed to figure out where all the PIGGIES are (ab being the most obvious one, with a not so obvious way to get OUT of the tab [inventory] once you click on something there but once you paste that little 200 IQ challenge, then you n		
Well about that....		
You realized soon enough that you need stones, and stones are literally everywhere around you, but you can't just pick these up like a regular person. You have to pick up special BUCKY APPROPRIATE STONES and they differ in color slightly and is a bit r		
And when you're done doing that, and being pretty proud of yourself, you're going to meet the locals, that is hunter NECTS Bears and Beasr. Don't worry - the Bears will either hug you in a surprising way and make grunting sounds while you enjoy i		
NECTs? This game has NECTs?		
Yes, and theyre armed, but don't worry here either- the NECTs are contemplating life in the game just like you are, they will go from spot to spot, oh, it's a deer, ah dear...it's another deer, oh there's a rabbit, there's another rabbit, oh there you ar		
...And it goes on.		
The animals are constantly on the run, just like you. And so are the GHOST animals! YES Ghost animals, because before you know it, you'll hear something ROAR past you and you didn't see a thing. You check your speakers and your surround		
Speaking of sounds...		
Sounds are plentiful and random. If you see a bear or a boar - you will hear it roar instantly - or not, because even though the coder has put some effort in sound distance, you'd be amazed of on how far or close the animals really are to you before		
Building stuff...		
Depends on how much wood you collect, and you will run for hours just to get the satisfaction of owning a few sticks, and collecting sticks is just a matter of finding a half pint sized bush where you need to aim to the left and at an EXACT pick-up-sp		
You will be equally surprised how much wood you need to collect, and there's plenty of trees, right? Well...yes, but not trees you can just chop down, as a lazy city dweller, with world class machron running powers, your arms are as weak as roos		
When you finally cut some branches up, and maybe even throw them onto some place, you will notice that every year basically kill all the same...and you chose barbaricly moved just to find our well...you are just an axe other off-line stone...are all th		

After:

[illegible]

OUTPUT:
Game analysis:



casual



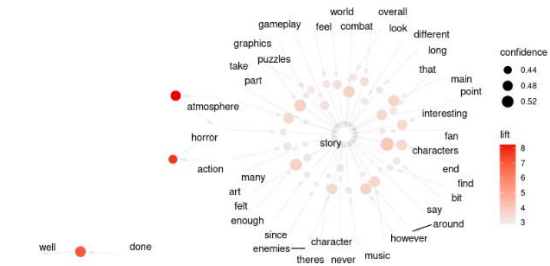
indie



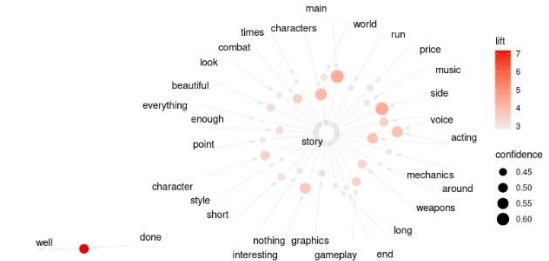
massively_multiplayer



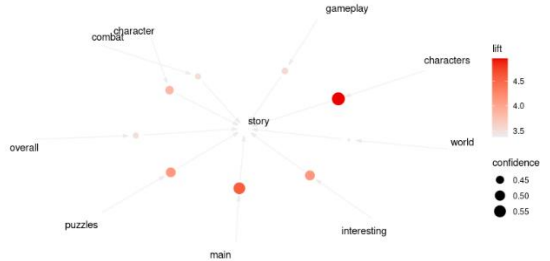
racing



rpg



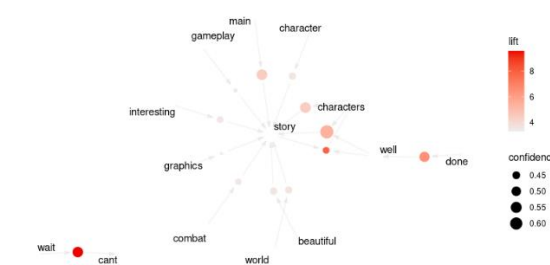
simulation



sports

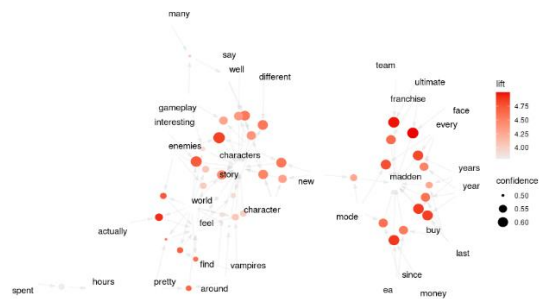


strategy

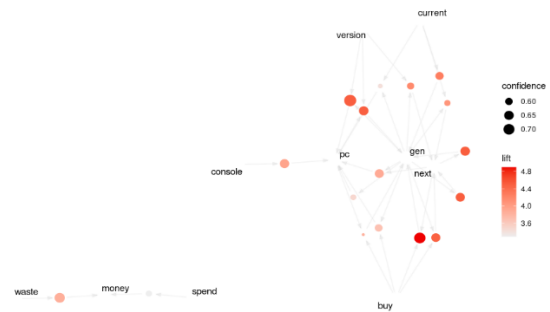


Game reviews analysis:

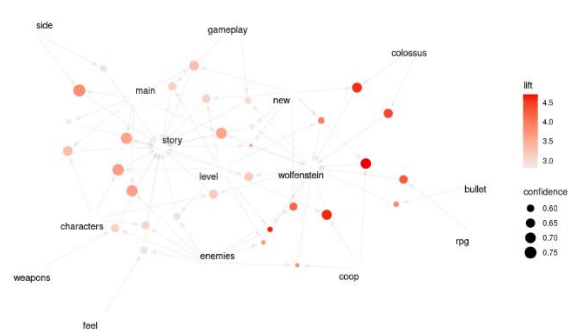
action



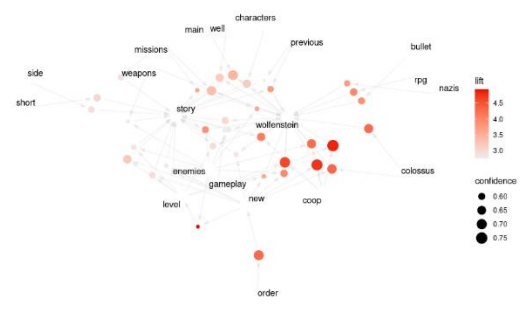
casual



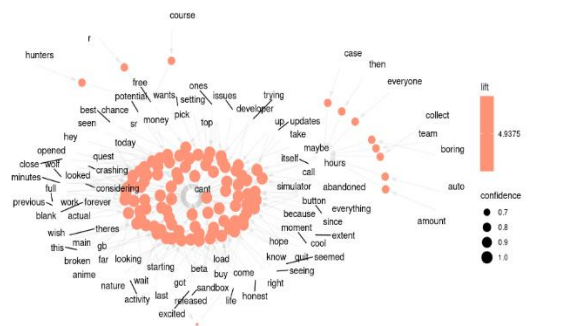
adventure



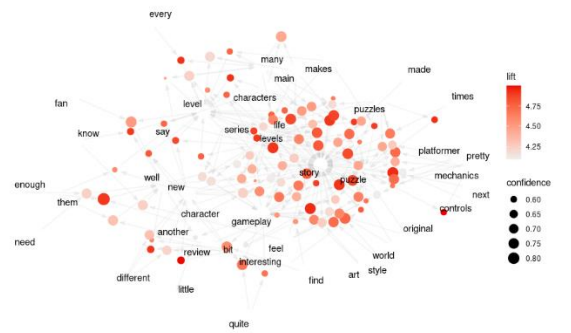
indie



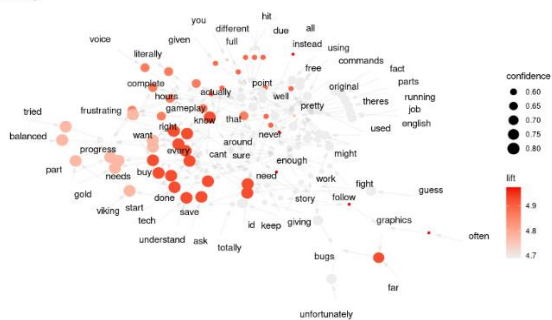
massively_multiplayer



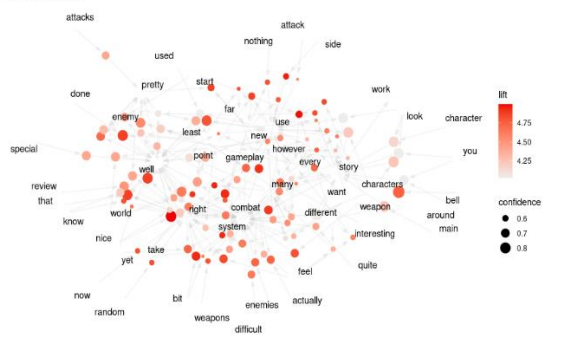
rpg

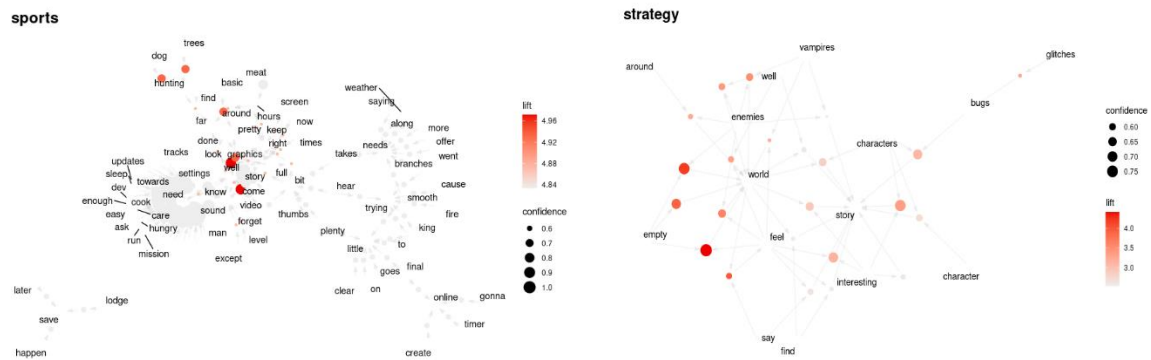


racing



simulation





WHAT WE LEARNED: Negative reviews are much more scattered about. While RPGs center around story (which in this case might indicate a bad story) to a degree, the other plots are less centralized around a common theme which could indicate that some elements of a game could be good, but one bad element of a game can ruin the whole thing.

4. **K-MEANS CLUSTERING** This analysis uses the K-Means clustering algorithm to group Steam apps into distinct clusters based on the features total_reviews, total_positive_reviews, total_negative_reviews, and price. The goal is to identify patterns and similarities among apps.

DATA PREPROCESSING:

Converts publisher into numerical columns using One-Hot Encoding. Scales numerical features (Total_Reviews and Price) to have a mean of 0 and a standard deviation of 1.

WHY WE USED IT:

1. Clustering: With an optimal $k=3$, clusters group products into categories like high review/high price (premium), moderate review/mid-range price (popular mid-tier), and low review/low price (budget).
2. Visualizations: A scatterplot visualizes clusters by total reviews and price, with each cluster color-coded.
3. Key Insights: Products are segmented into distinct price and popularity tiers, aiding in market understanding.

MODEL EVALUATION (Accuracy Scores):

Inertia (Within-Cluster Sum of Squares):

– Final inertia: 248.36, indicating a reasonable grouping of data points.

• Cluster Centers:

– Average prices for clusters:

Cluster 0: \$13.04

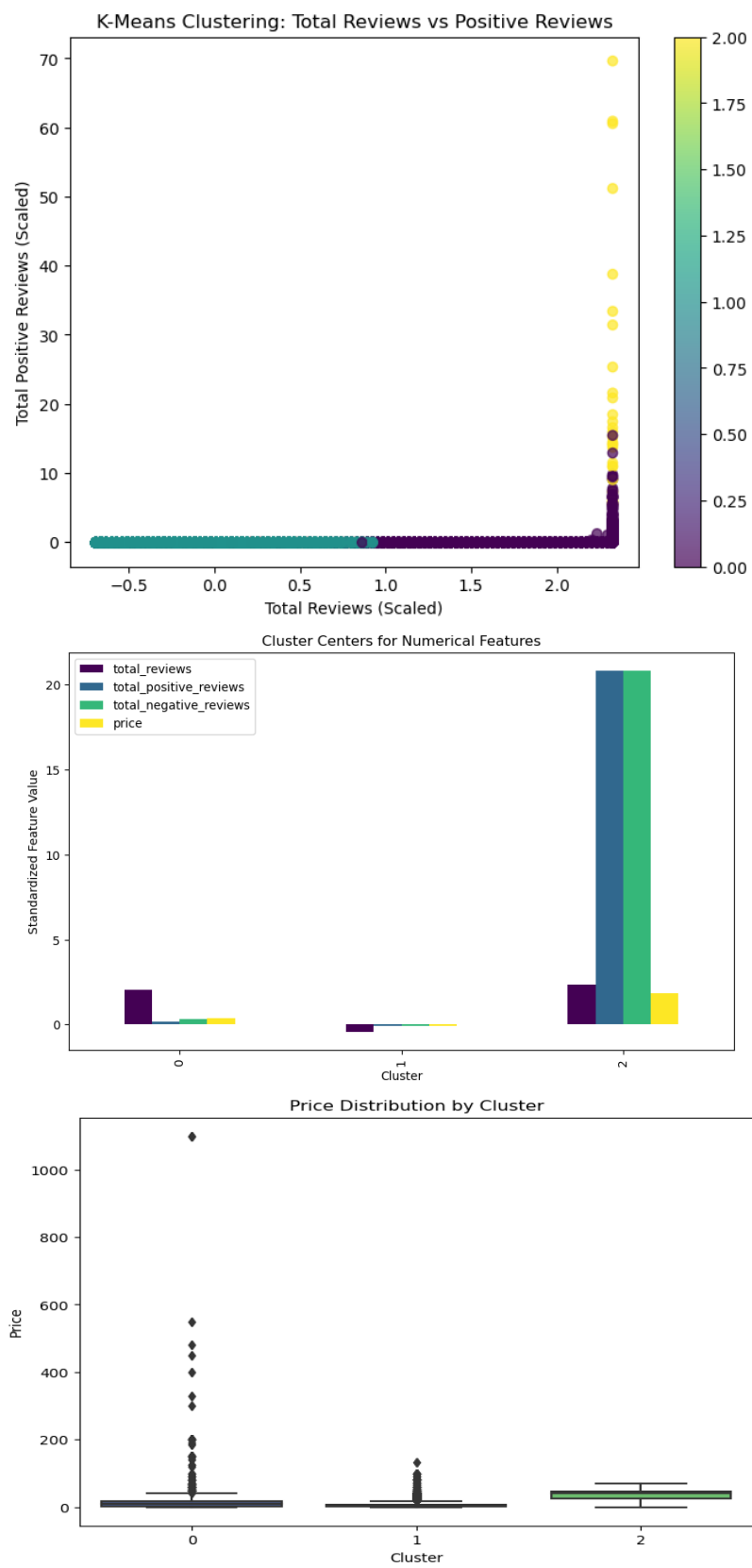
Cluster 1: \$5.94

Cluster 2: \$35.34

• Silhouette Score:

– An optional metric to evaluate clustering quality. Example: 0.52, indicating moderately well-separated clusters.

OUTPUT:



CONCLUSION

This project effectively studied patterns in the video game sector by gathering, refining, and presenting extensive data. The knowledge acquired can help developers, marketers, and stakeholders in making educated choices, matching product offerings with consumer preferences, and recognizing possible market opportunities. Future research might include a more in-depth examination of player involvement statistics and a wider investigation of upcoming gaming technologies.

In this project, we evaluated multiple models, including K-means, association rule mining, gradient boosting, CLIP model, to address various aspects of the game data analysis. Each model was assessed using metrics like accuracy, precision, recall, and F1-score to determine its effectiveness. The results demonstrated different outputs. While the models provided valuable insights, challenges such as data imbalance were addressed through techniques like hyperparameter tuning and data preprocessing. These findings highlight the importance of model selection and fine-tuning in deriving actionable insights from data.