

**The Science of Gaming:  
Predicting the Success of Board Games**  
Andrew Riddle

## 1 Motivation

What makes a board game fun to play? Perhaps it depends on the number of players that can play the game or the complexity of the game, or some other factors. If you could predict how well certain games will do based on certain intrinsic characteristics of the game, measured by a numerical rating from many people, you could use that information to predict how well a new game will sell or how good an older game is that has no or very little review data available. Using sale price from historical sales as the target instead of game rating would allow for a pricing tool for new games going to market. With that, you could also identify segments of the board game market that lead to the highest priced games and focus on developing those types of games at a cheaper production cost to maximize profit.

## 2 Data Acquisition

The game data were taken from boardgamegeek.com, using their object API located at <https://www.boardgamegeek.com/xmlapi2/thing>. All games with a game ID between 1 and 110,317 were queried. About half of these (49,173) were board games, specified by the object type. In addition to the basic game features, the API can also return some statistics about the game based on user reviews as well as sale listings on their website (by users). The statistics include the average user rating on a scale from 1 to 10 as well as a bayesian average of the user rating (the "geek rating"), which the website uses to rank their games. If there are no user reviews, the game gets a geek rating of 0 and has no rank. Since the original goal was to use the geek rating as the target variable for a metric of the quality of a game, only games that have a defined rating (i.e. games that have user reviews) can be used. Of the 50,000 games that I had data for, only 10,000 had ratings. 8,000 of these games had available market data and it was this final subset of games that was used in the analysis. Table 1 shows a list of the column data returned by the API as well as a brief description of each column.

## 3 Preprocessing

The first step was to parse the market data into a usable form; the data coming out of the API consist of individual sale listings for each game by users, including the date, price, condition, and currency. I used historical exchange rates pulled from the API at <https://api.fixer.io> in conjunction with the year of the listing to convert each sale to USD. No adjustment was made for inflation. Once all prices were converted to USD, the median price was used as the feature for that game. The median price was used as it's less susceptible to outliers than the mean. There were a handful of games with very high prices ( $> \$500$ ), so I did 2 rounds of sigma clipping, setting values higher than 3 times the standard deviation from the mean equal to 3 times the standard deviation plus the mean each time. This reduced the

**Table 1.** Columns.

Column Name	Description	Data Type
age	Minimum suggested age for players	int
avg_rating	Average user rating (1-10)	float
avg_time	Average length of a game (min)	float
category	Category or genre the game belongs to	str
designer	Designer(s) of the game	str
game_id	ID used to identify the game	int
geek_rating	Bayesian average of the user ratings; used to rank games	float
market	Sales listings by users of the site	dict
max_players	Maximum number of players the game can support	int
max_time	Maximum time a game takes (min)	int
mechanic	Mechanic(s) by which the game progresses	str
min_players	Minimum number of players required to play	int
min_time	Minimum time required to play a game (min)	int
name	Name of the game	str
num_votes	Number of user ratings	int
numcomments	Number of comments left by users	int
numweights	Number of user reviews for the game weight (complexity)	int
owned	Number of users that own the game	int
rank	Game rank	int
trading	Number of users trading the game	int
wanting	Number of users actively seeking the game	int
wishing	Number of users wishing for the game	int
weight	Average weight (complexity) rating left by users (1-5)	float
year	Year the game was released	int

maximum price from \$10,000 to \$200, with a distribution pictured in Fig. 1. This was done for two reasons: the first being that linear models don't handle outliers well, and the second being that the prices on some of these games might not be reliable or they might be for a special edition of the game that doesn't reflect the price of the standard version.

There were four other features with outliers: year, max\_players, avg\_time, and age. The year column had a massive range as the list of games includes very old game like Go and Backgammon (both "released" over 2000 years ago). To accomodate for this I set a minimum year value of 1950, imputing any values less than that to be 1950. I also created a "Classic" feature that equals 1 for games released before 1900 and 0 for games released after. A similar process was done for the other three columns. For max\_players a maximum cutoff of 20 was set as there are very few games beyond this point and the difference between 20 players and 20+ players isn't that significant. For the avg\_time a maximum cutoff of 480 minutes was used. I used 21 as the upper bound on age as a minimum suggested age beyond that doesn't make sense in terms of maturity/things you are legally allowed to do (e.g purchase alcohol). There are also some games that have no user reviews on the weight (complexity) of the game and thus have a value of 0 in that column. I imputed this missing data using the median weight of all games for which there was a non-zero weight.

Keeping in mind that the goal of the project is to be able to predict the game rating, which is a proxy of the popularity, of games that are brand new or have no user reviews, any

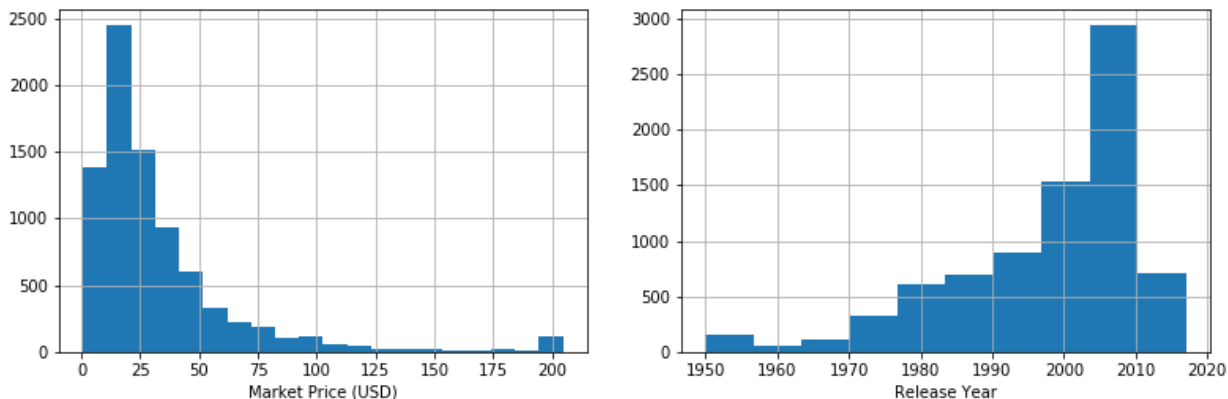


Figure 1: *Left*: Histogram of the available market data. *Right*: Histogram of release year after imputing very old games to 1950.

feature tied to user data must be thrown out. For this reason, I dropped `avg_rating`, `num_votes`, `numcomments`, `numweights`, `owned`, `trading`, `wanting`, and `wishing`. Also, because the market data consists of secondary market sales listings it is more representative of actual game value or popularity. Instead of using the market price as a feature, I did a separate analysis with it as the target variable instead of rating.

I created a column called `name_len` that is the number of characters in the game name under the theory that shorter names might be catchier and thus more popular. Once these step were taken, I used a heatmap to check for colinear features and to get an idea of which features might have the most predictive power (see Fig. 2). The features related to time spent playing the game are all colinear (Pearson correlation coefficient  $> 0.8$ ). I chose to keep just the `avg_time` of those three features. I also used a Seaborn pairplot of the remaining numerical columns to look for trends between features and the two target variables as well as any other interesting feature relationships. Some of these plots for the more important features are shown in Fig. 4. The game ID was left in temporarily so that the results can be matched up after the analysis, but will be removed before being fed into any models.

That leaves the four text fields to be dealt with: `name`, `mechanic`, `category`, and `designer`. Other than the `name_len`, no other features were engineered from the `name` column as the names are too specific to each game. For the other three I used a `CountVecotorizer` to convert the text data into numerical columns. Each game can have multiple entries in a column (it can have multiple mechanic types, categories, and designers), but each entry is only ever listed once (per game). The vectorizer turns each column into a distinct bag of words and then creates a column for each word that is a 1 if the game has that mechanic, category, or designer and 0 otherwise. Before the count vectorizer is used, missing data in each text field is changed to be “Designer Unknown” for the designer column and so forth. For the designer column, I used a minimum document frequency of 7. In other words, a designer had to be listed on 7 games before he or she was added to the bag of words for that column. This was done so that only designers with some pedigree would be included. A few different numbers were tried for this minimum document frequency, but 7 was found to produce the best results.

Each of the numeric columns (not counting the vectorized text columns) was passed

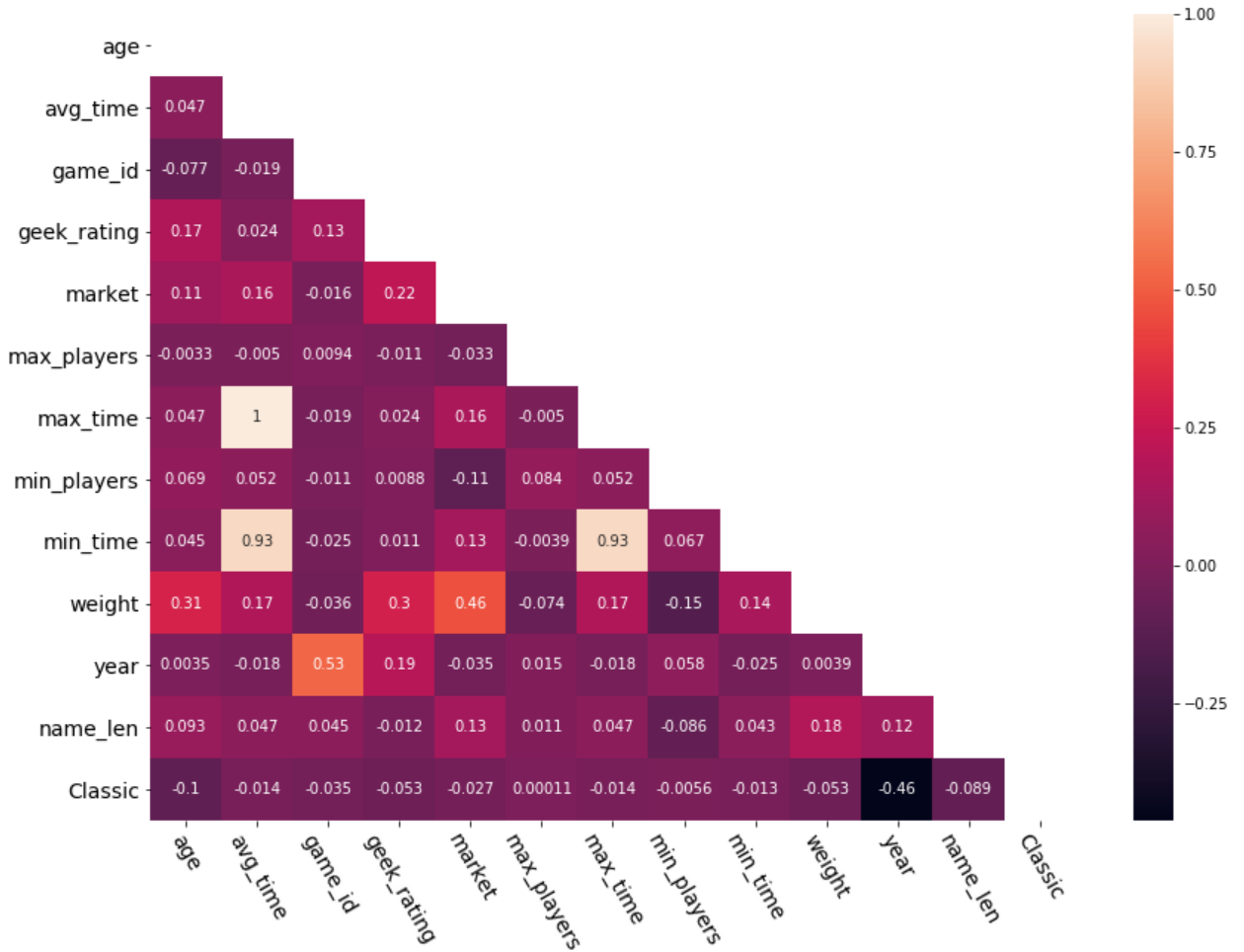


Figure 2: Heatmap of numerical columns showing the Pearson correlation coefficients.

through a sigma clipping function that changed values greater than 3 standard deviations from the mean to be at the edge of that boundary. Finally, they were passed through a MinMax scaler so that each column would be in the range from 0 to 1. This is done so that each feature starts with an equal weight in the models, although it's not necessary for tree-based models. As there are many features after using the count vectorizer, I thought using SVD might be helpful. However, when looking at the explained variance with each component (c.f. Fig. 3) there is no “elbow”, and thus I decided to not use SVD so as to preserve feature importances/weights in the models.

## 4 Modeling

I built the steps discussed in §3 into a pipeline for analysis. The full data set was split into 80-20 train/test partitions twice. The first split was to generate a holdout data set to use for final testing, while the second split was done for the traditional validation purposes. This

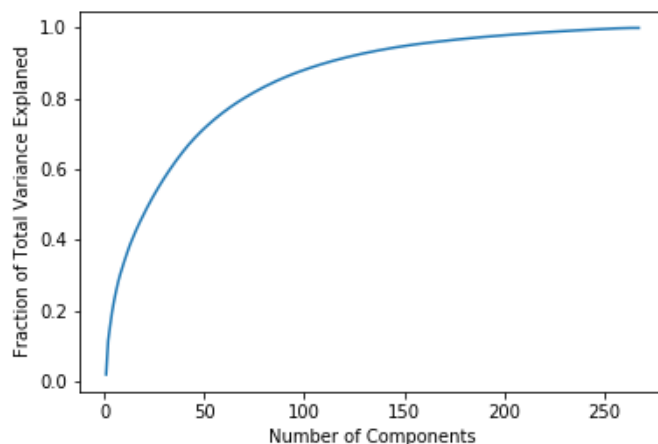


Figure 3: Cumulative explained variance ratio as you add more SVD components.

means that the models were trained on 64% of the total data set with 16% used as validation data and 20% used as holdout data.

A number of different regression models were tested on the training data and validated on the test set, including Random Forest, Elastic Net, Gradient Boosted Tree, K-Neighbors, SVM, and a Neural Network. For each model I did a grid search over relevant hyperparameters, optimizing on RMS error. The parameter grid I tried for each model is shown in Table 2.

## 5 Results

### 5.1 Geek Rating

The baseline score for the test set (the RMSE using the mean rating as the prediction for all games) is 0.47. The model that produced the best result on the test set was a Random Forest model with `min_samples_leaf=1`, `n_estimators=150`, and `max_features=1/3`. However, even the best model had a RMSE of 0.40 on the test set and the holdout set. Also, this model had a rather large difference between the scores on the training set and the test set (0.24). Changing the `min_samples_leaf` from 1 to 5 leads to a small decrease in test score (0.005), but a significant reduction in the score difference between the train/test sets (0.068). However, this model still performs about the same on the holdout set. Some of the top features from the RF model include `weight`, `year`, `name_len`, `avg_time`, and `age`.

### 5.2 Market Price

Similar results were found using the market price as the target variable. The baseline RMSE of the test set was \$33.76. The best score attained was still the RF model with a test score of \$27.77 and a holdout score of \$29.07. The top features were similar to that for the model of game rating, with few features having an importance greater than 0.01, but the top couple feature were stronger overall. Some of these top features are shown plotted against the two

target variables in Fig. 4.

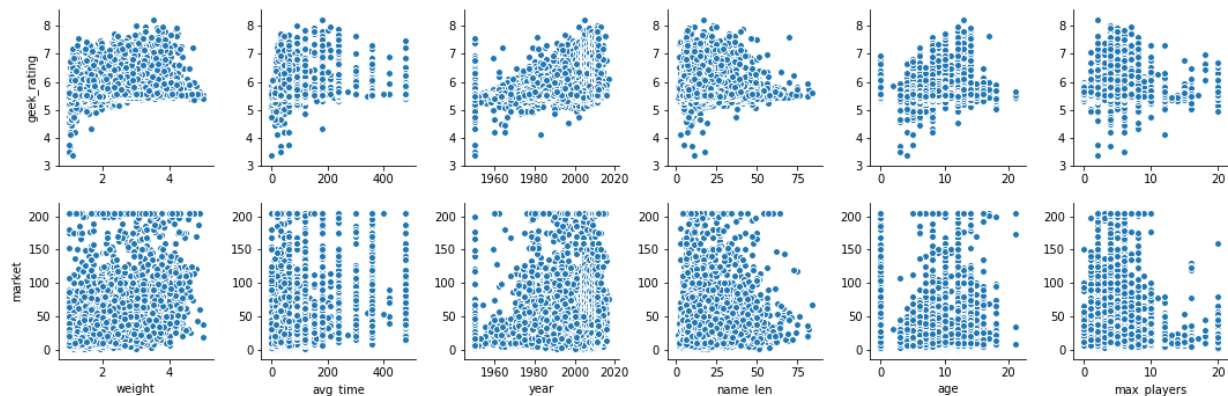


Figure 4: Scatter plots of the target variables against some of the top feature from the RF models.

### 5.3 Conclusions

There are several issues with this data set that lead to the models not being that much better than the baseline scores. Many of the features show a trend with the target variable over a small subset of the feature and then a more flat distribution over the rest of the range of that feature. Take `avg_time`, for example (left plot of Fig 5). Games with a lower `avg_time` tend to have a lower rating and almost all of the games with low ratings have short average times ( $<60$  min). However, for games with an `avg_time` longer than 60 minutes, the distribution is flatter and the feature has much less predictive power. The distribution of geek ratings and market prices are both heavily skewed (c.f. Fig. 1 and Fig. 5), which makes predictions tricky.

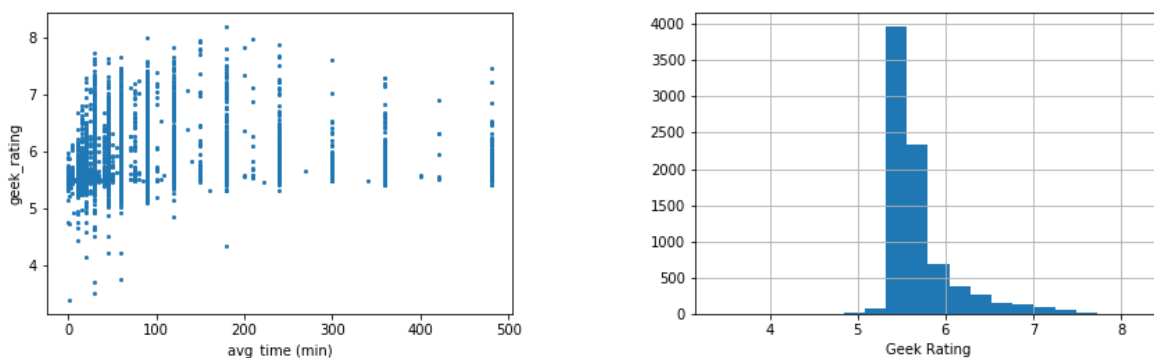


Figure 5: *Left*: Geek rating versus average game time (min). *Right*: Geek rating histogram.

Also of interest is the plot of game ID versus release year shown in Figure 6. The `game_id` isn't exactly in chronological order of release, but there is a strong correlation. This plot shows that shortly after the year 2000, the board game industry started to explode as many more games were being released per year. This can also be seen in the histogram of release year shown in the right plot of Fig. 1.

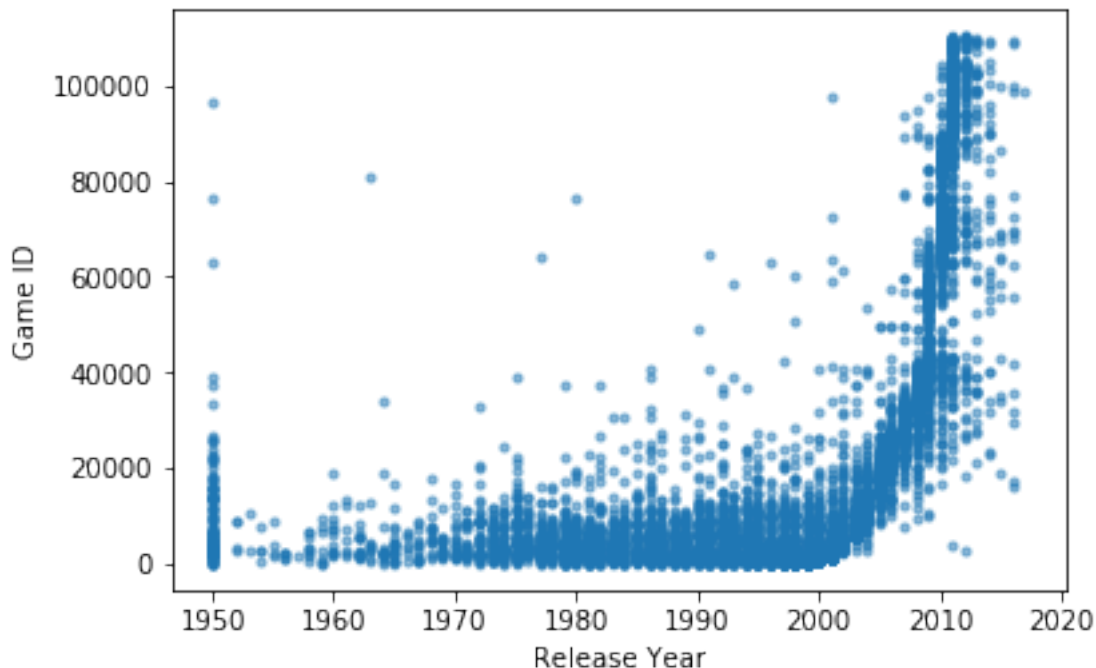


Figure 6: *Scatter plot of game\_id versus year*

## 6 Future Plans

### 6.1 Classification

One thing I'd like to try is to change the goal of the project from trying to predict a particular market price or rating for a game to instead predict whether the game will be good or bad (cheap or expensive). This will allow the model to leverage the small regions in parameter space that lead to lower average ratings and prices and not be as hurt by the other regions where the target variables have a more even distribution.

### 6.2 Logarithmic Targets

One of the problems mentioned in §5.3 was that the target variables themselves have very skewed distributions. Moving the target variable into log space will reduce the skew of the target distributions and might allow the models to be more predictive.

### 6.3 Consider Subsets of Games

Using SVD with two components, the data set is separable into two relatively distinct clusters (see Fig. 7). There appears to be a separation along the axis of the second components, with higher average prices in the upper cluster. This is born out by choosing a (somewhat arbitrary) cutoff of 0.2 for the second component and calculating the mean prices for each group. The mean prices of the lower and upper cluster are \$25.26 and \$45.18, respectively. Doing a K-Means clustering analysis with  $k = 2$  yields the clusters shown in Fig. 9. Of note, some of the top feature weights of the second SVD component ( $y$ -axis of Fig. 9) have to do

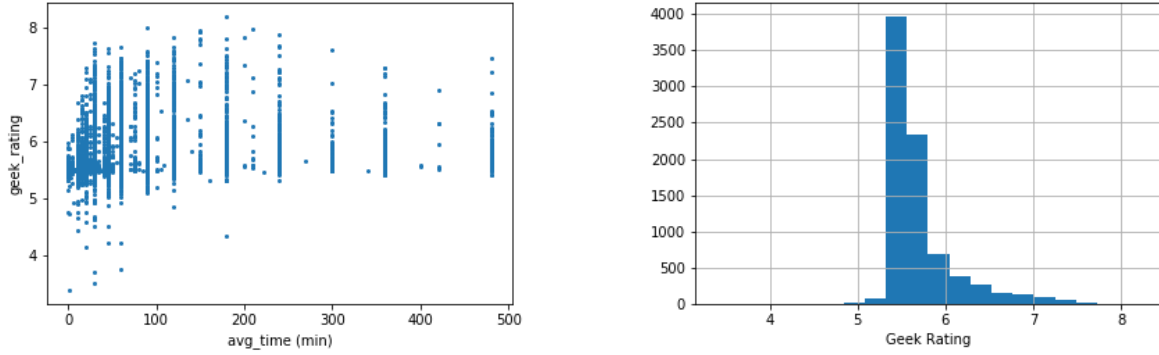


Figure 7: *Left*: Geek rating versus average game time (min). *Right*: Geek rating histogram.

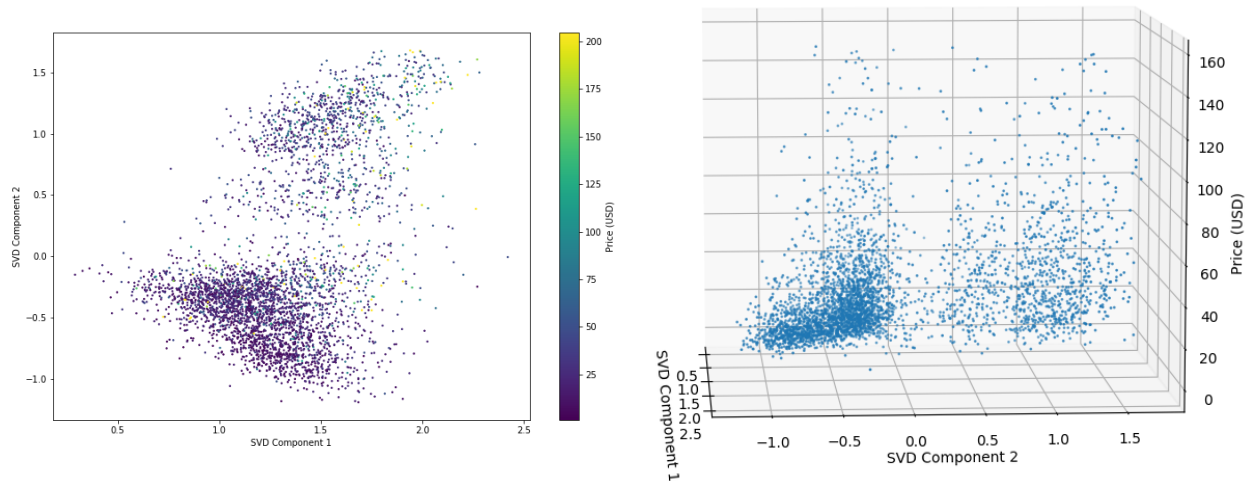


Figure 8: *Left*: Scatter plot of the top 2 SVD components with market price indicated by color. *Right*: Same as the left plot, but with market price plotted on the  $z$  axis in a 3-D plot.

with the game category of Wargame and game mechanics that are popular in such games, such as Hex-and-Counter (c.f. <https://boardgamegeek.com/boardgamemechanic/2026/hex-and-counter>). An interesting analysis might be to separate all the wargames and analyze just those. A similar thought might be to analyze the two clusters found via the K-Means clustering separately, with a majority of the wargames in the upper (and more expensive) cluster.



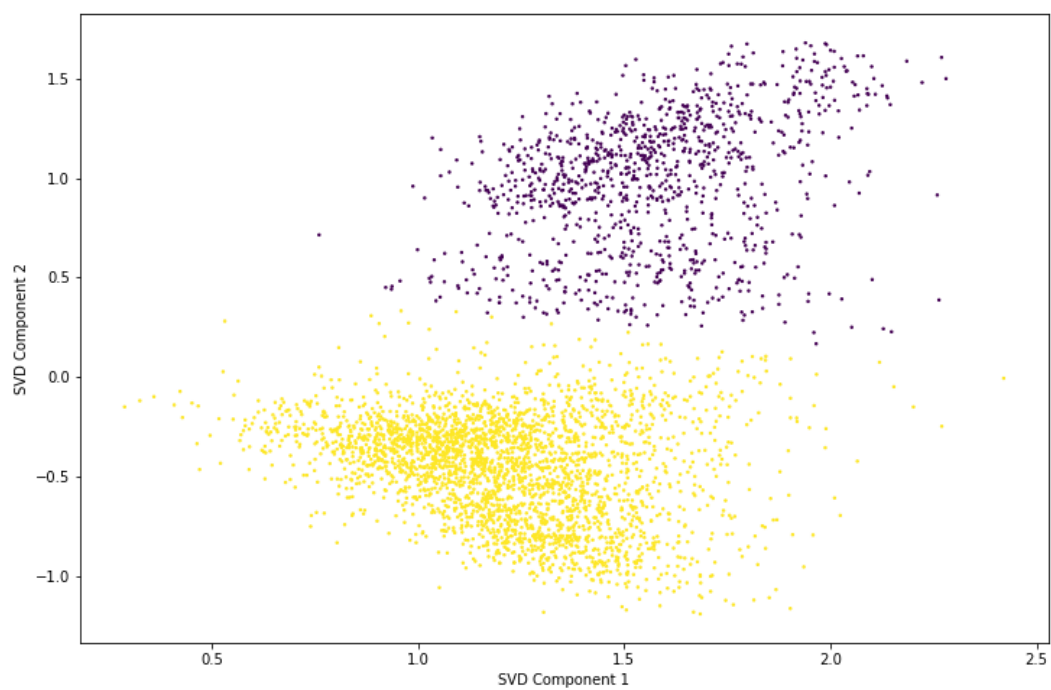


Figure 9: Clustering from K-Means with  $k = 2$ .

**Table 2.** Model Parameter Grids

Random Forest (RF)	
Parameter	Grid
min_samples_leaf	1, 3, 5, 7
n_estimators	25, 50, 75, 100, 125, 150
max_features	None, 1/3
Elastic Net (EN)	
Parameter	Grid
$\alpha$	0.001, 0.01, 0.1, 1, 10
l1_ratio	0.01, 0.05, 0.1, 0.25, 0.5, 0.8, 0.9
fit_intercept	True, False
Gradient Boosted Tree (GBT)	
Parameter	Grid
learning rate	0.05, 0.1, 0.3
max_depth	3, 5, 7, 9
max_features	None, 1/3
n_estimators	100, 150, 200
subsample	0.5, 0.8, 1
K-Nearest Neighbors (KNN)	
Parameter	Grid
n_neighbors	3-21
Support Vector Machine (SVM)	
Parameter	Grid
kernel	rbf, linear, poly
degree	2, 3, 4
$\gamma$	auto, 0.1, 1, 10
C	0.1, 1, 10
$\epsilon$	0.05, 0.1, 0.2, 0.4
Neural Network	
Parameter	Grid
Dropout	0, 0.25, 0.5
Nodes in 1st hidden layer	25, 50, 75, 100, 150
Nodes in 2nd hidden layer	0, 25, 50