

Machine Learning Model for Real-life Predictions



Overview

- Machine learning models were trained on data gathered from Steamspy
- Determine the most popular features that well-selling games provide



Steam Database

Objective of Project

- The objective of this project is to train an appropriate machine learning model to be able to predict the relative success of a video game
- Based on characteristics such as:
 - Genre
 - Categories
 - Tags
 - Price

Data Pre-processing

- No missing values were found in the dataset
- One-hot encoding used to interpret data
- Variables were not scaled
- Removed outliers

Feature Selection

- Chose to prioritize the following columns:
 - Tags
 - Categories
 - Genre
 - Price
 - Positive ratings
 - Negative ratings.

NARROW BY TAG	
Action 1,310	Indie 801
Shooter 785	First-Person 737
Singleplayer 595	Multiplayer 587
Adventure 461	Sci-fi 330
Early Access 303	Gore 299
Violent 296	Co-op 271

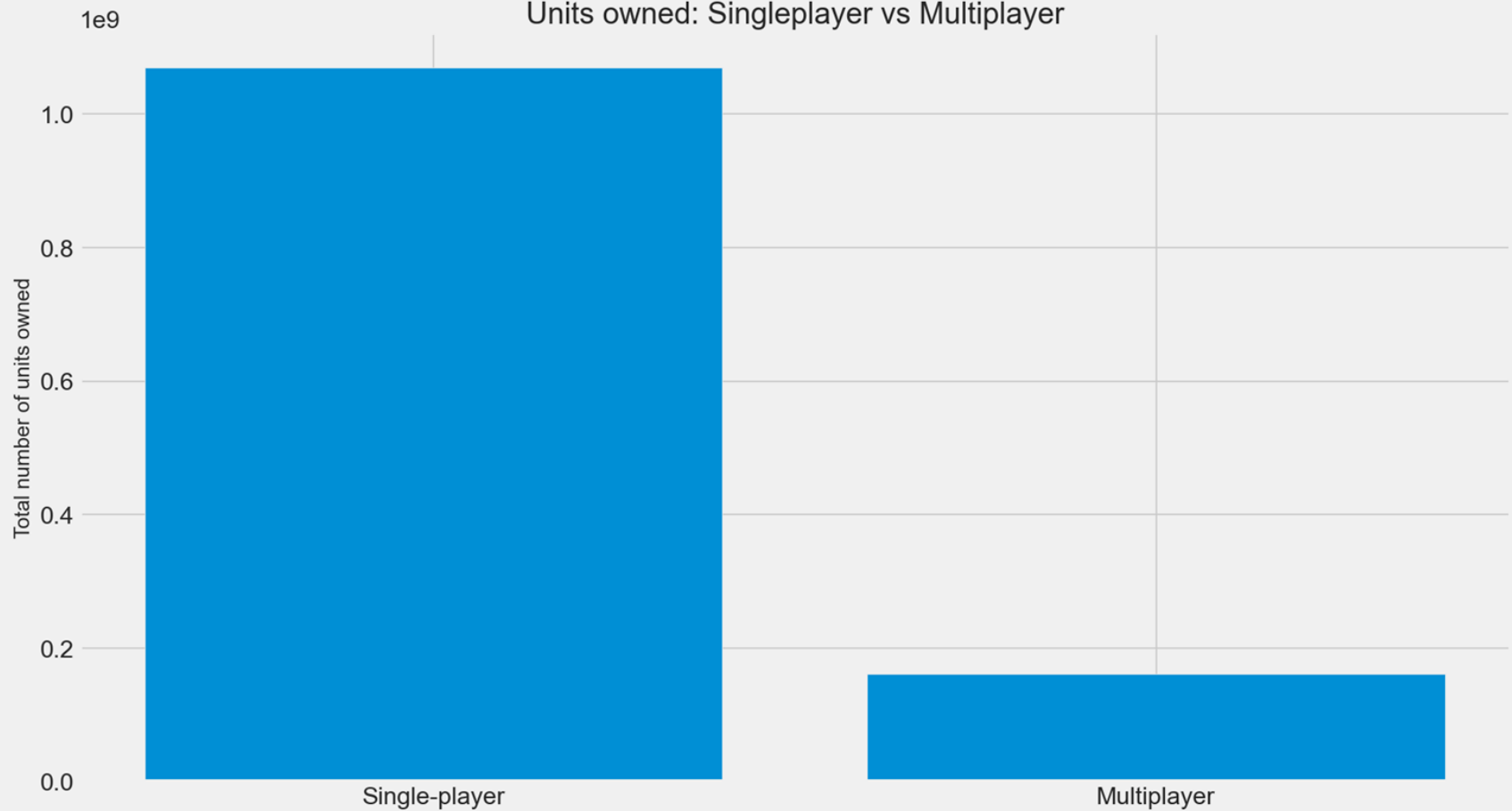
Feature Selection Continued

- Dropped achievements column
- Filtered out entries that were released before 2013
- Dropped entries that did not contain an ASCII value in the title

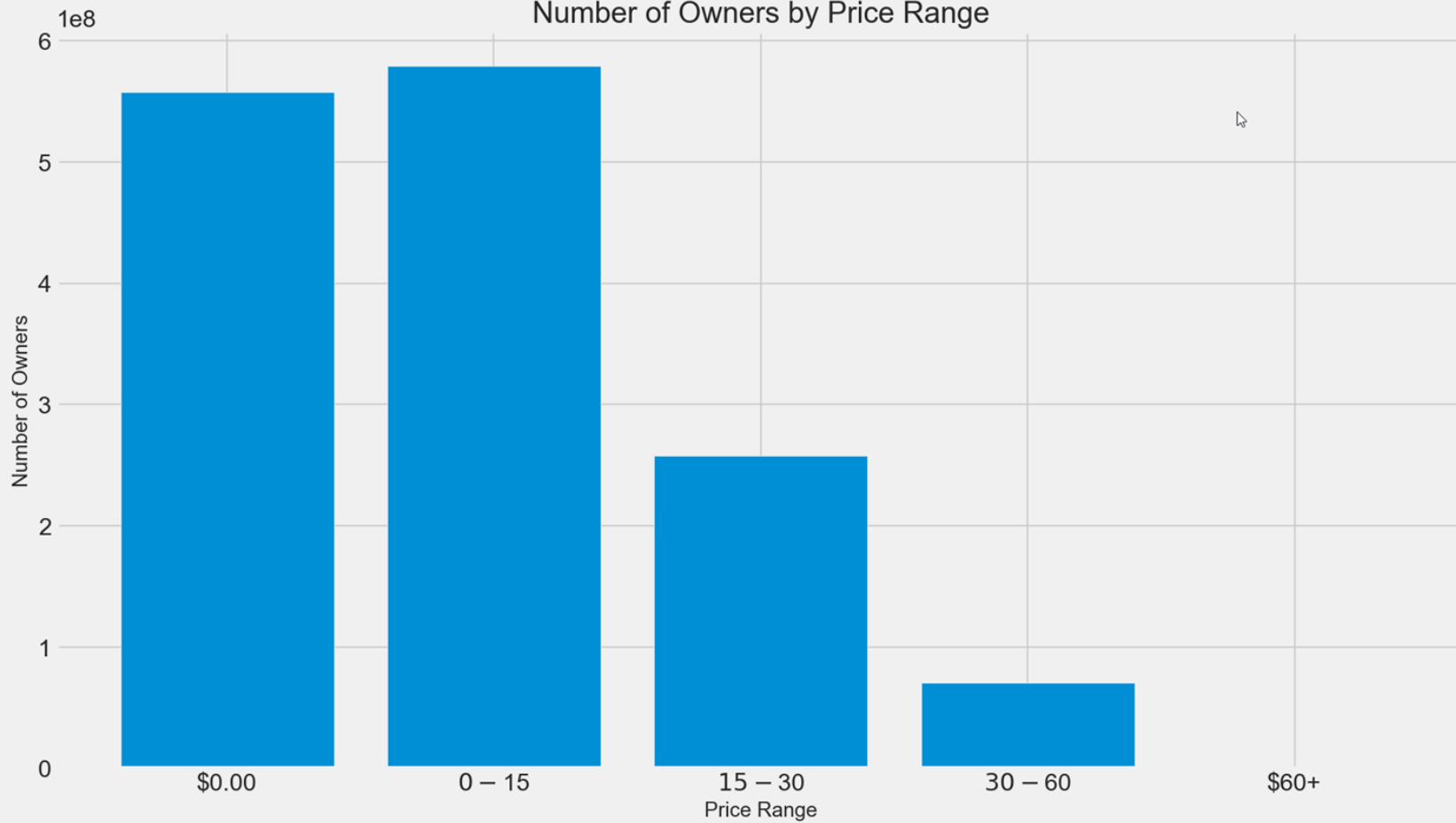
Data Visualization

- Process of creating graphs and charts to help you understand the patterns and trends in the data.
- Helps to identify outliers and anomalies in the data.

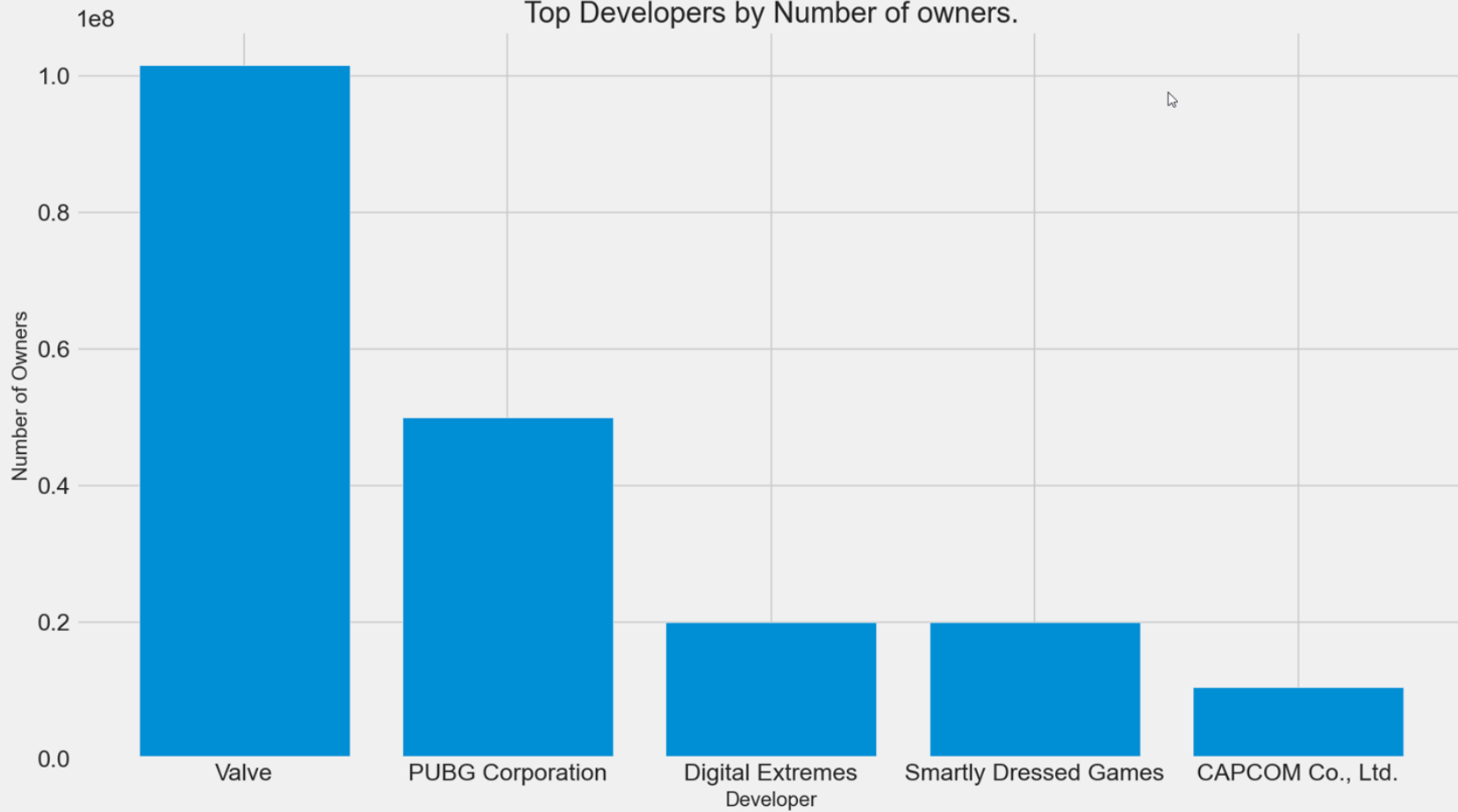
Units owned: Singleplayer vs Multiplayer



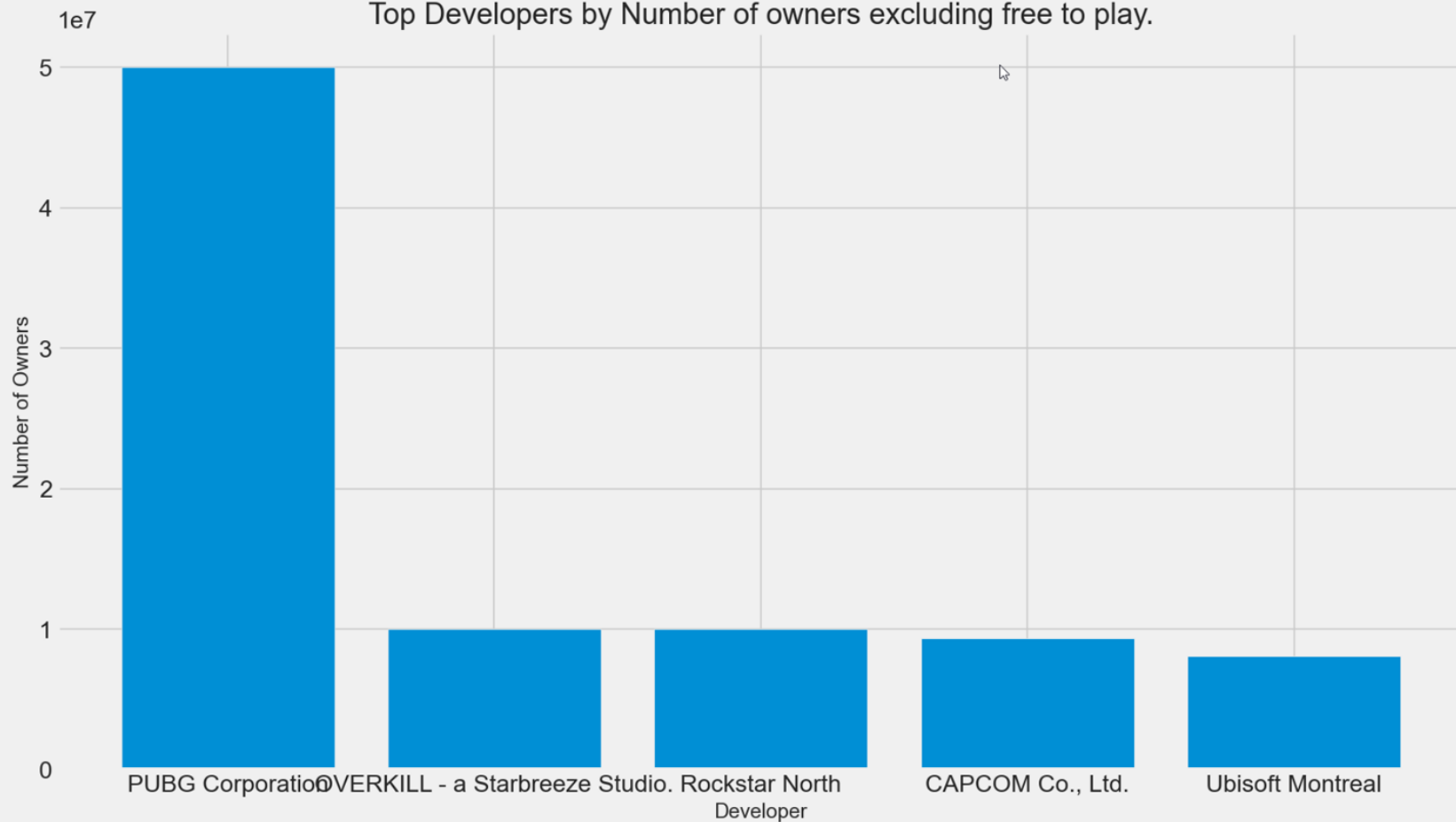
Number of Owners by Price Range

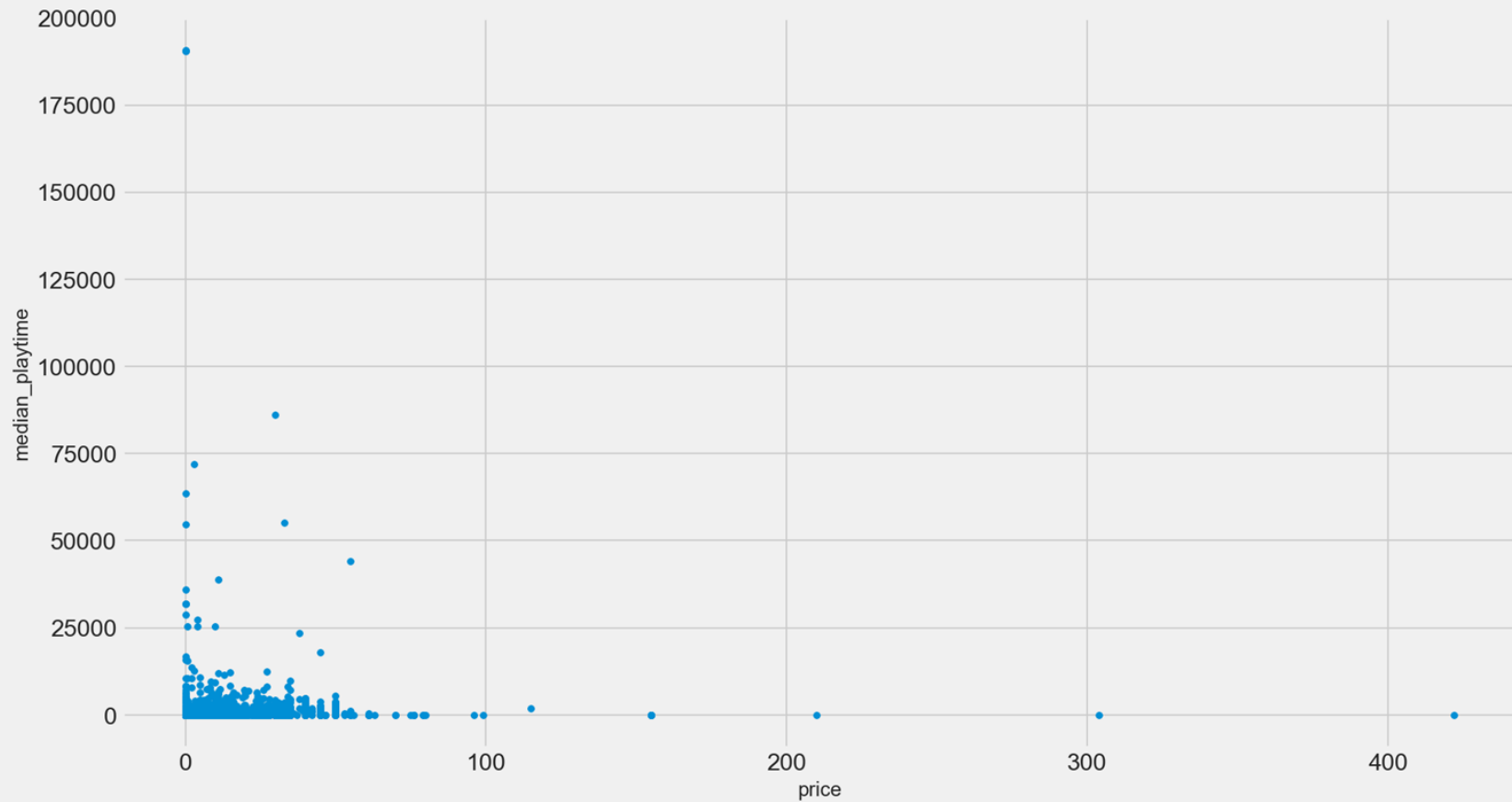


Top Developers by Number of owners.



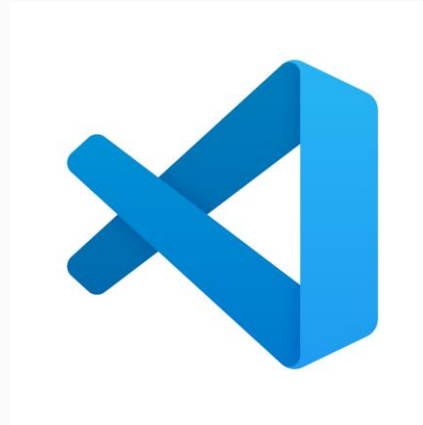
Top Developers by Number of owners excluding free to play.





Platform and Machine Configuration Used

- A Visual Studio Code environment was used with the Anaconda distribution. Every python script was run on a local machine using the environment.



Model Planning

- Different machine learning regression models were used when training on the dataset.
 - Linear Discriminant Analysis
 - Logistic Regression
 - K-Nearest Neighbour
 - Decision Trees
 - Lasso Regression
 - Gaussian Naive Bayes

Model Training

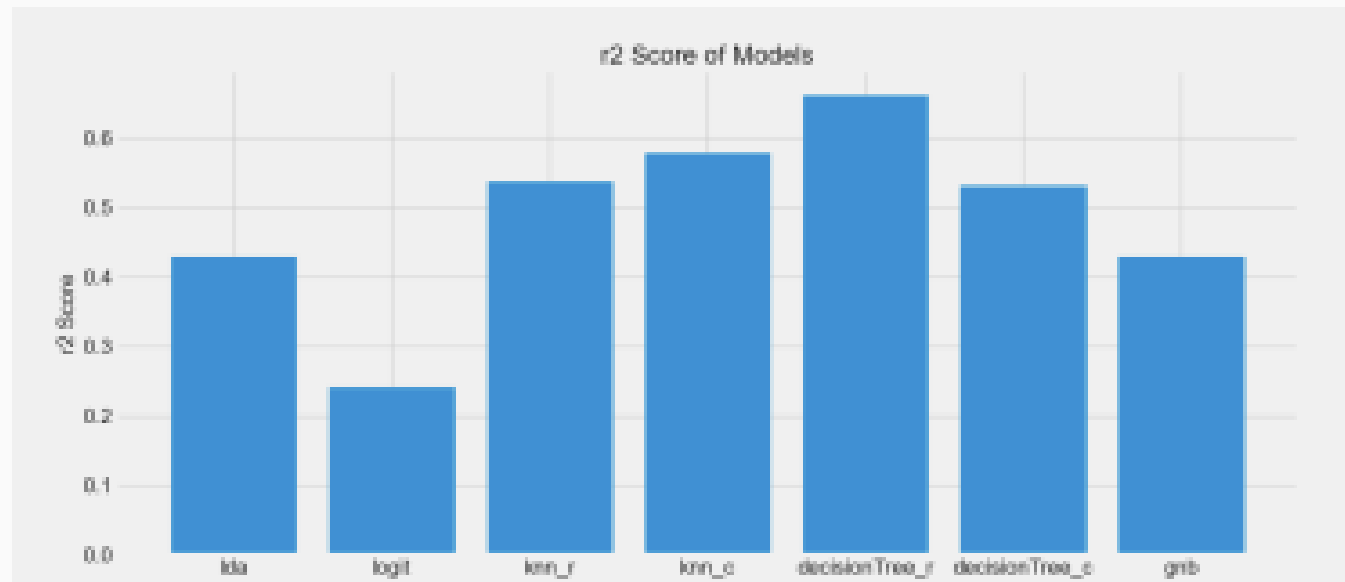
- Used sklearn
- Models were fitted and trained to predict for the “owners” metric which is an estimated number of the users that own the game on the platform.

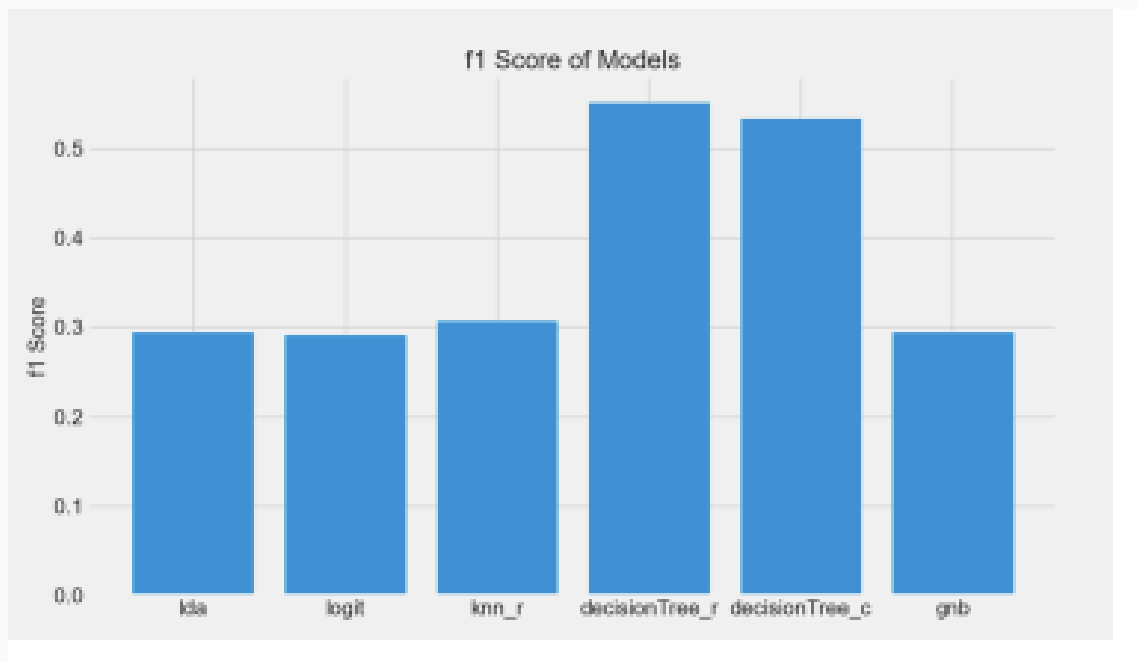
Model Evaluation

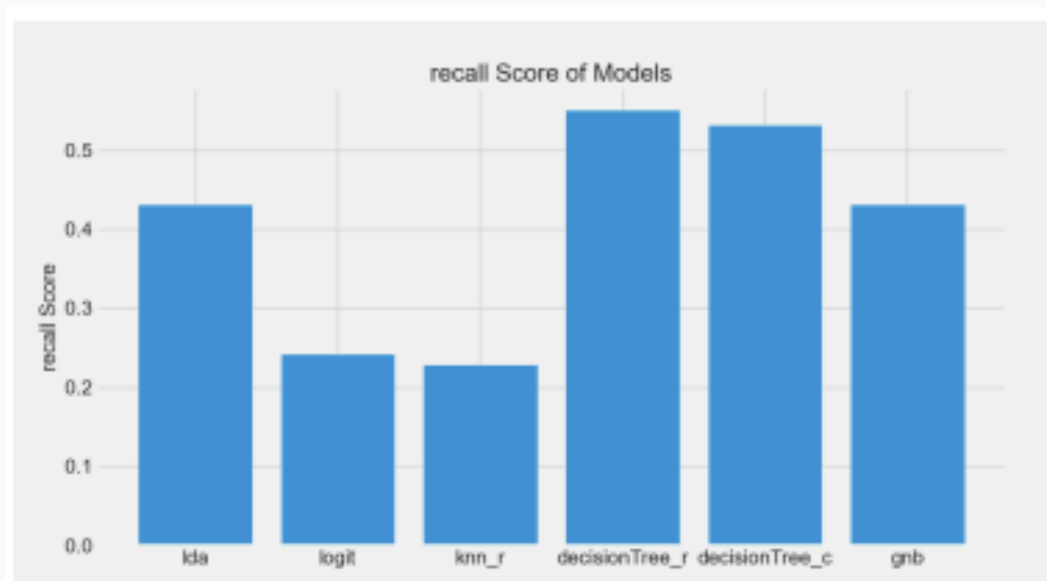
- Model evaluation was done using 4 statistics:
 - Accuracy / r^2 score
 - F1 score
 - Precision
 - Recall
- Model optimization was done during the data processing and feature selection phases.

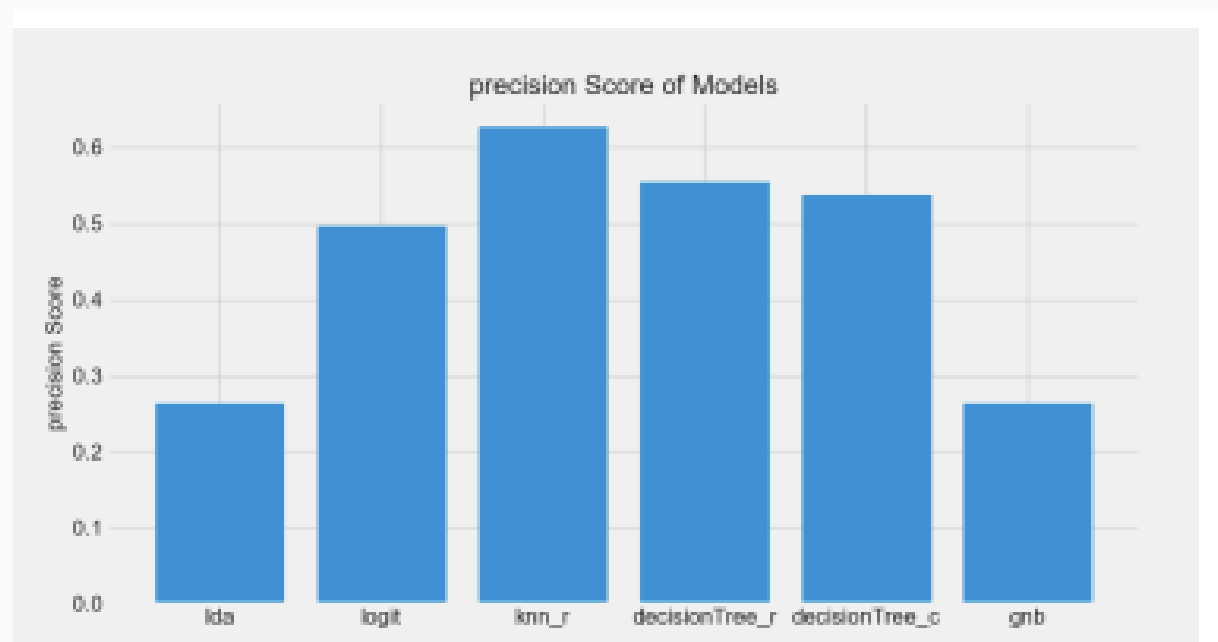
Final Model Building

- Comparison graphs were plotted after all model evaluation was done









Result Table

Model	R2 / Accuracy	F1	Recall	Precision
-----	-----	-----	-----	-----
Logistic Regression	0.2421	0.2918	0.2421	0.4976
LDA	0.4307	0.2949	0.4307	0.2642
KNN Regressor	0.5386	0.3071	0.2282	0.6281
KNN Classifier	0.5786	0.5603	0.5786	0.5520
DT Regressor	0.6623	0.5525	0.5501	0.5562
DT Classifier	0.5320	0.5345	0.5320	0.5383
Gaussian NB	0.4307	0.2949	0.4307	0.2642
-----	-----	-----	-----	-----
Mean	0.4879	0.4051	0.4275	0.4715
Std	0.1254	0.1249	0.1325	0.1359

Interpretation of Results:

- Decision tree model yields the best results of performance metrics.
- The final model did not achieve the specified accuracy
 - Number of owners was an estimated value between a broad range.
 - An accurate value of sales may provide better performance metrics.
- Strong decision tree results may suggest:
 - Non-linear relationships between variables
 - Mixture of both categorical and numerical data.

Sensitivity Analysis

- Various attempts were made at changing certain aspects of the process in order to obtain the best performance metrics:
 - Drop entries based on their release date
 - Experiment with the value of SelectKBest

Future Research

- More emphasis should be put on obtaining accurate sales figures.
- An alternative solution could be to replace the number of owners column with a sales column.