# Analysis of Predictors of Price and Popularity of Steam Games in 2016

## Introduction

This is a dataset on all steam games scraped from public Steam APIs and steamspy.com. The dataset comes from data.world at this url: https://data.world/craigkelly/steam-game-data. The dataset has 78 columns and includes things like price, descriptions and release date. According to the data on Github, the repository was created in Septembere 12, 2016 so we are assuming that the data also comes from the end of the year in 2016. The dataset has the following columns:

QueryID - (Integer) The original ID in idlist.csv

ResponseID - (Integer) The ID returned in the Steam response (should equal QueryID)

QueryName - (Text) The original name in idlist.csv

ResponseName - (Text) The name returned in the Steam response (should equal QueryName)

ReleaseDate - (Text) Appears to the be the initial release date for the game

RequiredAge - (Integer) list named required_age in JSON

DemoCount - (TextualCount) list named demos in JSON

DeveloperCount - (TextualCount) list named developers in JSON

DLCCount - (TextualCount) list named dlc in JSON

Metacritic - (Integer) numeric score from metacritic object in JSON

MovieCount - (TextualCount) list named movies in JSON (used object id for unique count)

PackageCount - (TextualCount) list named packages in JSON

RecommendationCount - (Integer) from recommendations.total in JSON

PublisherCount - (TextualCount) list named publishers in JSON

ScreenshotCount - (TextualCount) list named screenshots in JSON

AchievementCount - (Integer) achievements.total in JSON

AchievementHighlightedCount - (TextualCount) for achievements.highlighted in JSON

ControllerSupport - (Boolean) True if controller_support was full

IsFree - (Boolean) is_free in JSON

FreeVerAvail - (Boolean) True if is_free_license is True in package_groups list

PurchaseAvail - (Boolean) True if price_in_cents_with_discount greater than 0 in package_groups list

SubscriptionAvail - (Boolean) True if is_recurring_subscription is True in package_groups

PlatformWindows - (Boolean) True if platforms.windows is True

PlatformLinux - (Boolean) True if platforms.linux is True

PlatformMac - (Boolean) True if platforms.mac is True

PCReqsHaveMin - (Boolean) True if pc_requirements.minimum is non-empty string

PCReqsHaveRec - (Boolean) True if pc_requirements.recommended is non-empty string

LinuxReqsHaveMin - (Boolean) True if linux_requirements.minimum is non-empty string

LinuxReqsHaveRec - (Boolean) True if linux_requirements.recommended is non-empty string

MacReqsHaveMin - (Boolean) True if mac_requirements.minimum is non-empty string

MacReqsHaveRec - (Boolean) True if mac_requirements.recommended is non-empty string

CategorySinglePlayer - (Boolean) True if for any i, categories[i].description is "single-player"

CategoryMultiplayer - (Boolean) True if for any i, categories[i].description is one of: "cross-platform multi-player", "local multi-player", "multi-player", "online multi-player", "shared/split screen"

CategoryCoop - (Boolean) True if for any i, categories[i].description is one of: "co-op", "local co-op", "online co-op"

CategoryMMO - (Boolean) True if for any i, categories[i].description is "mmo"

CategoryInAppPurchase - (Boolean) True if for any i, categories[i].description is "in-app purchases"

CategoryIncludeSrcSDK - (Boolean) True if for any i, categories[i].description is "includes source sdk"

CategoryIncludeLevelEditor - (Boolean) True if for any i, categories[i].description is "includes level editor"

CategoryVRSupport - (Boolean) True if for any i, categories[i].description is "vr support"

GenreIsNonGame - (Boolean) True if for any i, genres[i].description is one of: "utilities", "design & illustration", "animation & modeling", "software training", "education", "audio production", "video production", "web publishing", "photo editing", "accounting"

GenreIsIndie - (Boolean) True if for any i, genres[i].description is "indie"

GenreIsAction - (Boolean) True if for any i, genres[i].description is "action"

GenreIsAdventure - (Boolean) True if for any i, genres[i].description is "adventure"

GenreIsCasual - (Boolean) True if for any i, genres[i].description is "casual"

GenreIsStrategy - (Boolean) True if for any i, genres[i].description is "strategy"

GenreIsRPG - (Boolean) True if for any i, genres[i].description is "rpg"

GenreIsSimulation - (Boolean) True if for any i, genres[i].description is "simulation"

GenreIsEarlyAccess - (Boolean) True if for any i, genres[i].description is "early access"

GenreIsFreeToPlay - (Boolean) True if for any i, genres[i].description is "free to play"

GenreIsSports - (Boolean) True if for any i, genres[i].description is "sports"

GenreIsRacing - (Boolean) True if for any i, genres[i].description is "racing"

GenreIsMassivelyMultiplayer - (Boolean) True if for any i, genres[i].description is "massively multiplayer"

PriceCurrency - (Text) price_overview.currency in JSON

PriceInitial - (Float) price_overview.initial in JSON, divided by 100.0 to converts cents to currency

PriceFinal - (Float) price_overview.final in JSON, divided by 100.0 to converts cents to currency

SteamSpyOwners - (steamspy.com) total owners, which includes free weekend trials and other possibly spurious numbers.

SteamSpyOwnersVariance - (steamspy.com) total owners, which includes free weekend trials and other possibly spurious numbers. Note that this is not technically variance: according to steamspy.com, "the real number... lies somewhere on... [value +/- variance]"

SteamSpyPlayersEstimate - (steamspy.com) best estimate of total number of people who have played the game since March 2009

SteamSpyPlayersVariance - (steamspy.com) errors bounds on SteamSpyPlayersEstimate. Note that this is not technically variance: according to steamspy.com, "the real number... lies somewhere on... [value +/- variance]"

SupportEmail - (Textual) support_info.email in JSON

SupportURL - (Textual) support_info.url in JSON

AboutText - (Textual) about_the_game in JSON

Background - (Textual) background in JSON

ShortDescrip - (Textual) short_description in JSON

DetailedDescrip - (Textual) detailed_description in JSON

DRMNotice - (Textual) drm_notice in JSON

ExtUserAcctNotice - (Textual) ext_user_account_notice in JSON

HeaderImage - (Textual) header_image in JSON

LegalNotice - (Textual) legal_notice in JSON

Reviews - (Textual) reviews in JSON

SupportedLanguages - (Textual) supported_languages in JSON

Website - (Textual) website in JSON

PCMinReqsText - (Textual) text of pc_requirements.minimum

PCRecReqsText - (Textual) text of pc_requirements.recommended

LinuxMinReqsText - (Textual) text of linux_requirements.minimum

LinuxRecReqsText - (Textual) text of linux_requirements.recommended

MacMinReqsText - (Textual) text of mac_requirements.minimum

MacRecReqsText - (Textual) text of mac_requirements.recommended

Our goal is to discover interesting trends that allow us to predict how a game is price and it's popularity after release. The price is divided into price initial and price final. We are going to use price initial as our metric to determine price and SteamSpyOwners as a metric to predict popularity.

## Data Cleanup

From the above, we can see that there is a lot of information that is hard to use to predict price. One example is the description. Unless we can extract some kind of metric from this, we won't be able to use it to predict price. As such, we choose to remove the following columns: QueryID, ResponseID, PCMinReqsText, PCRecReqsText, MacMinReqsText, MacRecReqsText, LinuxMinReqsText, LinuxRecReqsText, Reviews, LegalNotice, HeaderImage, DRMNotice, ExtUserAcctNotice, ShortDescrip, Background, AboutText, PriceCurrency.

We then try and extract some useful information out of the other columns. We first need to transform all of the true/false columns to factors. We then need to change the date time format so that it can be easily read by R.

```
##   [1]   306   600   669   682   688   744   905  1025  1181  1229  1254  1268
##  [13]  1269  1270  1597  1986  2011  2318  2845  2876  2999  3563  3651  3806
##  [25]  3830  3913  3972  4104  4125  4147  4298  4944  4986  5322  5643  6203
##  [37]  6323  6395  6844  6984  7003  7145  7344  7564  7753  7762  7884  7938
##  [49]  8065  8075  8240  8359  8685  8752  8807  8844  8930  8954  9121  9132
##  [61]  9152  9183  9319  9442  9481  9501  9545  9553  9564  9582  9671  9714
##  [73]  9863  9890 10019 10051 10215 10260 10293 10294 10302 10319 10325 10362
##  [85] 10468 10536 10598 10696 10754 10809 10853 10863 10961 10963 10970 11029
##  [97] 11085 11102 11169 11180 11501 11566 11585 11598 11652 11740 11747 11767
## [109] 11777 11778 11789 11791 11795 11826 11852 11858 11882 11885 11894 11909
## [121] 11924 11928 11941 11960 11961 11971 11999 12005 12013 12018 12019 12021
## [133] 12029 12034 12047 12052 12067 12068 12089 12097 12100 12103 12106 12124
## [145] 12160 12163 12174 12182 12195 12198 12199 12202 12204 12211 12233 12242
## [157] 12257 12278 12293 12303 12379 12382 12383 12393 12396 12424 12429 12482
## [169] 12487 12493 12506 12520 12521 12523 12532 12550 12556 12559 12562 12566
## [181] 12576 12580 12582 12601 12634 12649 12655 12665 12681 12712 12733 12767
## [193] 12789 12794 12809 12831 12844 12886 12900 12913 12924 12935 12941 13021
## [205] 13031 13102 13104 13116 13126 13180 13205 13207 13248 13263 13270 13280
## [217] 13289 13293 13297 13313 13325 13336 13341 13342 13343 13347
```

Now, we changed the Supported Languages column to contain the total number of supported languages, instead of the specific languages.

Next, we modified the Detailed Description field to contain the number of words.

We consolidated the Support Email and Support URL columns into a single column called Support, which holds a binary value based on whether a game has either a support email or a support URL.

We changed the Website column to indicate whether or not a website was provided

## Principle Components Analysis

Using principal component analysis, we tried to reduce the dimensions of the data given to something that is easier to interpret. The components extracted using the PCA function are as follows:

```
##               FreeVerAvailTrue           ControllerSupportTrue
##                    0.074508528                      0.070456149
##                    ReleaseDate               GenreIsNonGameTrue
##                    0.045623039                      0.043416614
##              PlatformWindowsTrue          GenreIsFreeToPlayTrue
##                    0.004339572                     -0.003198739
##              SubscriptionAvailTrue           GenreIsSportsTrue
##                   -0.008387476                     -0.013410307
##              CategoryVRSupportTrue            GenreIsCasualTrue
##                   -0.014164153                     -0.014987983
##                 CategoryMMOTrue     CategoryInAppPurchaseTrue
##                   -0.017014156                     -0.017425834
##               GenreIsRacingTrue      CategoryIncludeSrcSDKTrue
##                   -0.019477385                     -0.019621476
## GenreIsMassivelyMultiplayerTrue         GenreIsEarlyAccessTrue
##                   -0.022229404                     -0.029657569
##                    RequiredAge                        DLCCount
##                   -0.035270958                     -0.041101655
##              GenreIsSimulationTrue                    DemoCount
```

```
##                   -0.045608914                      -0.058909342
##                   GenreIsRPGTrue              ControllerSupportFalse
##                   -0.069529578                      -0.070456149
##             RecommendationCount    CategoryIncludeLevelEditorTrue
##                   -0.072258215                      -0.075454092
##             GenreIsStrategyTrue              GenreIsAdventureTrue
##                   -0.076956068                      -0.082023865
##          SteamSpyPlayersEstimate                SteamSpyOwners
##                   -0.083716051                      -0.087388876
##               CategoryCoopTrue                GenreIsActionTrue
##                   -0.094873905                      -0.099024698
##          CategoryMultiplayerTrue                      Metacritic
##                   -0.104781595                      -0.118210773
##               AchievementCount               SupportedLanguages
##                   -0.125752326                      -0.128166381
##               PCReqsHaveRecTrue                   PackageCount
##                   -0.138874371                      -0.149434076
##               GenreIsIndieTrue                PurchaseAvailTrue
##                   -0.155333942                      -0.164471179
##                 DetailedDescrip                 DeveloperCount
##                   -0.167448077                      -0.186315032
##                    WebsiteTRUE                      MovieCount
##                   -0.187502045                      -0.191829546
##                 ScreenshotCount              LinuxReqsHaveRecTrue
##                   -0.194203024                      -0.200554516
##     AchievementHighlightedCount                PlatformLinuxTrue
##                   -0.201618269                      -0.202732910
##               MacReqsHaveRecTrue          CategorySinglePlayerTrue
##                   -0.206223338                      -0.217801564
##                 PlatformMacTrue              LinuxReqsHaveMinTrue
##                   -0.220588387                      -0.238113499
##                    SupportTRUE                 PublisherCount
##                   -0.238153760                      -0.244145291
##               MacReqsHaveMinTrue                PCReqsHaveMinTrue
##                   -0.249825736                      -0.297218170
```

This first component explains the most variance out of all the components. The most significant coefficient indicates whether a Free version is avaiable. This makes sense because games that have a frere or trial version will get players hooked to purchase the full, more costly version. Release Date is another significant Coefficient, and suggests that newer games have higher prices.

```
##     ControllerSupportFalse         PCReqsHaveMinTrue              PublisherCount
##              0.26799885                0.23861955                0.19613841
##          ScreenshotCount           DetailedDescrip CategorySinglePlayerTrue
##              0.14893131                0.14595983                0.12773316
##           DeveloperCount              PackageCount          PurchaseAvailTrue
##              0.12171039                0.11257952                0.11030653
##  CategoryMultiplayerTrue
##              0.09904693
##  [ reached getOption("max.print") -- omitted 44 entries ]
```

The second component seems indicative of PC games due to the highest component being lack of controller support. Thus, we can conclude this category contains older AAA high budget games due to the low emphasis on release date as well as a high emphasis on number of Publishers and Developers.

```
## CategorySinglePlayerTrue          PurchaseAvailTrue          PCReqsHaveMinTrue
##            0.25858586                 0.24076940                 0.14102728
##        GenreIsCasualTrue          GenreIsIndieTrue     ControllerSupportFalse
##            0.13685124                 0.12831029                 0.11018263
##         PublisherCount       GenreIsAdventureTrue                SupportTRUE
##            0.10219650                 0.09127584                 0.08888365
##         DeveloperCount
##            0.06325220
##  [ reached getOption("max.print") -- omitted 44 entries ]
```

The third component describes niche indie games. We can conclude this because of the positive emphasis on the coefficients of Single Plater Games, Casual Games, and Indie Games as well as a negative emphasis on Number of owners and Number of Players.

```
## GenreIsMassivelyMultiplayerTrue                          ReleaseDate
##                      0.2963920                            0.2948865
##                 CategoryMMOTrue                 GenreIsFreeToPlayTrue
##                      0.2843971                            0.2840576
##                     SupportTRUE                 GenreIsEarlyAccessTrue
##                      0.2262849                            0.1853767
##          ControllerSupportFalse             CategoryInAppPurchaseTrue
##                      0.1646891                            0.1537688
##                GenreIsIndieTrue                       GenreIsRPGTrue
##                      0.1485697                            0.1083563
##  [ reached getOption("max.print") -- omitted 44 entries ]
```

The fourth component seems to describe newer multiplayer games. This is due to the high positive coefficients of Massive Multiplayer Genre, Release Date, and MMO Category.

We couldn't accurately determine the subset of the fifth component, so we decided to stop any further PCA analysis.