

# Reconhecimento de Padrões

## Análise de Jogos da Steam

**Grupo:**

- Guilherme Martiniano de Oliveira – 11215765
- Mateus Miquelino da Silva – 11208412
- Gustavo Fernandes Carneiro de Castro – 11369684

**Professor:**

Ricardo Zorzetto Nicolielo Vencio

## Introdução

O *dataset* selecionado para a análise possuía diversas informações sobre jogos distribuídos pela plataforma *Steam*, a maior plataforma de distribuição online de jogos eletrônicos. Esse *dataset*, de nome *steam.csv*, foi retirado do site *Kaggle*, uma comunidade online de pesquisas em análise de dados. O *dataset* contém dados retirados da própria plataforma, e do site *SteamSpy*, que disponibiliza informações diversas sobre os jogos.

## Objetivos

O objetivo da análise foi descobrir quais dos diferentes dados presentes no *steam.csv* eram relevantes para descobrir sua popularidade (número estimado de jogadores totais). Essa análise foi feita utilizando 3 métodos diferentes de análise supervisionada (KNN, com 3 e com 7 vizinhos mais próximos, e LDA). Além disso, os métodos foram comparados para descobrir o melhor método de análise possível para esse dataset.

## Dataset

O dataset era inicialmente composto por 18 colunas totais e um total de 27075 linhas, cada uma com um jogo diferente.

As colunas eram:

<i>appid</i>	Identificador diferenciado para cada jogo.
<i>name</i>	Título do jogo.
<i>release_date</i>	Data de lançamento do jogo no formato AAAA-MM-DD.
<i>english</i>	Suporte de idioma: 1 se o jogo é em inglês.
<i>developer</i>	Nome(s) do(s) desenvolvedor(es). Ponto e vírgula para delimitar múltiplos desenvolvedores.
<i>publisher</i>	Nome(s) do(s) distribuidora(s). Ponto e vírgula para delimitar múltiplas distribuidoras.
<i>platforms</i>	Lista delimitada por pontos e vírgulas das plataformas suportadas.
<i>required_age</i>	Idade mínima requerida para jogar, de acordo com o padrão da PEGI UK. 0 é livre para todos os públicos.
<i>categories</i>	Lista delimitada por pontos e vírgulas das categorias dos jogos.
<i>genres</i>	Lista delimitada por pontos e vírgulas dos gêneros dos jogos.
<i>steamspy_tags</i>	Lista delimitada por pontos e vírgulas das <i>tags</i> mais populares dos jogos (similar aos gêneros, mas dados pela comunidade).
<i>achievements</i>	Número de conquistas dentro do jogo.
<i>positive_ratings</i>	Número de análises positivas, dados do <i>SteamSpy</i> .
<i>negative_ratings</i>	Número de análises negativas, dados do <i>SteamSpy</i> .
<i>average_playtime</i>	Média de tempo jogado (minutos), dados do <i>SteamSpy</i> .
<i>median_playtime</i>	Mediana de tempo jogado (minutos), dados do <i>SteamSpy</i> .
<i>owners</i>	Número estimado de donos. Contém o limite superior e inferior.
<i>price</i>	Preço do jogo na época, em libras esterlinas.

Para tratar o dataset, nós analisamos cada uma das colunas para ver se havia discrepância. Após verificar uma a uma utilizando *summary(steam)*, nós verificamos q as

colunas *average\_playtime* e *median\_playtime* tinham: os valores de tempo médio de jogo iguais aos da mediana dos tempos de jogo; e/ou os valores de tempo médio de jogo maiores que os valores do jogo mais jogado da *Steam*, o *Dota 2*. Após um simples loop, esses dados foram retirados, deixando o dataset com 27059 linhas (jogos).

Além disso, *steam.csv* possuía diversas colunas com dados não numéricos e em muita quantidade, como gênero do jogo e plataformas de lançamento. Esses dados eram do tipo “string”, e separados por ponto e vírgula, como por exemplo: “Action; Adventure”. Para separá-los, foram criadas matrizes  $nrow(steam) \times 1$ , com informações binárias (0 e 1, “tem” e “não tem”, respectivamente) sobre a existência desses dados nas colunas do jogo, pela utilização da função *grepl*(“*dado*”, *steam[jogo]*, “*coluna*”)) para separar as “strings”.

Após implementar os tratamentos de dado acima, nas colunas que foram selecionadas como relevantes, *steam.csv* ficou com um total de 27059 linhas (jogos) e 27 colunas (dados), nesta mesma ordem:

<i>english</i>	Dado binário. Possui, ou não, suporte de idioma (inglês).
<i>required_age</i>	Idade mínima requerida para jogar, de acordo com o padrão da PEGI UK.
<i>achievements</i>	Número de conquistas dentro do jogo.
<i>positive_ratings</i>	Número de análises positivas, dados do <i>SteamSpy</i> .
<i>negative_ratings</i>	Número de análises negativas, dados do <i>SteamSpy</i> .
<i>average_playtime</i>	Média de tempo jogado (minutos), dados do <i>SteamSpy</i> .
<i>median_playtime</i>	Mediana de tempo jogado (minutos), dados do <i>SteamSpy</i> .
<i>price</i>	Preço do jogo na época, em libras esterlinas.
<i>windows</i>	Dado binário. Possui, ou não, suporte para <i>Windows</i> .
<i>linux</i>	Dado binário. Possui, ou não, suporte para <i>Linux</i> .
<i>mac</i>	Dado binário. Possui, ou não, suporte para <i>Mac</i> .
<i>free_to_play</i>	Dado binário. É, ou não, de graça para instalar e jogar.
<i>action</i>	Dado binário. É, ou não, um jogo de ação.
<i>indie</i>	Dado binário. É, ou não, um jogo indie (criado por uma pequena empresa).
<i>adventure</i>	Dado binário. É, ou não, um jogo de aventura.
<i>strategy</i>	Dado binário. É, ou não, um jogo de estratégia.
<i>rpg</i>	Dado binário. É, ou não, um jogo de RPG.
<i>racing</i>	Dado binário. É, ou não, um jogo de corrida.
<i>simulation</i>	É, ou não, um jogo de simulação. Dado binário.
<i>casual</i>	Dado binário. É, ou não, um jogo casual.
<i>sports</i>	Dado binário. É, ou não, um jogo de esporte.
<i>violent</i>	Dado binário. É, ou não, um jogo violento.
<i>nudity</i>	Dado binário. Possui, ou não, nudez.
<i>single_player</i>	Dado binário. É, ou não, um jogo violento.
<i>multiplayer</i>	Dado binário. Pode, ou não, ser jogado sozinho.
<i>steam_cards</i>	Dado binário. Pode, ou não, ser jogado por múltiplos jogadores.
<i>created_by_valve</i>	Dado binário. É, ou não, um jogo criado pela própria <i>Steam</i> (a <i>Valve</i> é a criadora e distribuidora de jogos cuja <i>Steam</i> é dona).

Como os rótulos (labels) para a nossa análise, foi escolhida a coluna *owners*, que resume bem a popularidade de um jogo: quanto mais popular, mais pessoas tem o jogo.

## Métodos

Depois de ajeitar o dataset, deu-se espaço às análises. Primeiramente, foi escolhida uma *seed* para os testes, *set.seed(6969)*. Depois, foi estabelecido um grupo de *treino* com 80% dos jogos (escolhidos arbitrariamente pela função *sample()*). A comparação entre resultado das análises com os 20% restantes devolve a quantidade de acertos e erros e a precisão da análise, através da função da matriz de confusão, *confusionMatrix(resultado, labels[-treino])*, retirada da biblioteca *caret*.

Essa análise foi feita para o KNN, com os 3 e os 7 vizinhos mais próximos(KNN(3) e KNN(7)), e com o LDA. A quantidade de vizinhos foi escolhida propositalmente: 3 pois é a quantidade mais usual, e 7 pois foi a com a melhor precisão.

Primeiro, foi feita a análise dos KNNs e do LDA para o *treino* usual e anotada a precisão. Depois, 27 análises novas foram feitas para o KNN(3), o KNN(7) e para o LDA: a cada iteração, uma versão do *treino* sem uma coluna diferente era analisada. As precisões de cada iteração foram anotadas e comparadas.

Depois de comparadas, ao analisar que dado diminuía ou aumentava a precisão após ser retirado, foram descobertos os melhores dados (colunas) para cada método. Os *sub-treinos* de *treino* que possuíam apenas as melhores colunas foram chamados de *perfect\_Method\_* (*perfectKNN3*, *perfectKNN7* e *perfectLDA*).

A última análise foi feita após juntar os *perfect\_Method\_*, sendo chamados de *superPerfect\_Method\_*. Os dois *superPerfect\_Method\_* criados foram o *superPerfectKNN*, criado da junção do *perfectKNN3* e do *perfectKNN7*; e o *superPerfectAll*, criado da junção de todos os métodos de análise. A junção foi feita mantendo apenas os dados (colunas) presentes em todos os *perfect\_Method\_* que compunham o *superPerfect\_Method\_*.

## Dados

Os dados retirados das análises acima seguem abaixo:

### KNN(3)

Número da coluna retirada	Nome da coluna retirada	Precisão	Comparação	Importante ou Problemático
Original	0	0,6833		
-1	<i>english</i>	0,6833	Igual	Irrelevante
-2	<i>required_age</i>	0,6807	Diminuiu	Importante
-3	<i>achievements</i>	0,6853	Aumentou	Problema
-4	<i>positive_ratings</i>	0,6787	Diminuiu	Importante
-5	<i>negative_ratings</i>	0,6779	Diminuiu	Importante
-6	<i>average_playtime</i>	0,6800	Diminuiu	Importante
-7	<i>median_playtime</i>	0,6726	Diminuiu	Importante
-8	<i>price</i>	0,6964	Aumentou	Problema
-9	<i>windows</i>	0,6857	Aumentou	Problema
-10	<i>linux</i>	0,6848	Aumentou	Problema
-11	<i>mac</i>	0,6842	Aumentou	Problema
-12	<i>free_to_play</i>	0,6772	Diminuiu	Importante
-13	<i>action</i>	0,6805	Diminuiu	Importante
-14	<i>indie</i>	0,6844	Aumentou	Problema
-15	<i>adventure</i>	0,6839	Aumentou	Problema
-16	<i>strategy</i>	0,6864	Aumentou	Problema
-17	<i>rpg</i>	0,6874	Aumentou	Problema
-18	<i>racing</i>	0,6861	Aumentou	Problema
-19	<i>simulation</i>	0,6807	Diminuiu	Importante
-20	<i>casual</i>	0,6892	Aumentou	Problema
-21	<i>sports</i>	0,6855	Aumentou	Problema
-22	<i>violent</i>	0,6866	Aumentou	Problema
-23	<i>nudity</i>	0,6813	Diminuiu	Importante
-24	<i>single_player</i>	0,6844	Aumentou	Problema
-25	<i>multiplayer</i>	0,6907	Aumentou	Problema
-26	<i>steam_cards</i>	0,6888	Aumentou	Problema
-27	<i>created_by_valve</i>	0,6848	Aumentou	Problema

Nome do sub-treino	Número das colunas que se mantiveram	Precisão
Perfect KNN(3)	2, 4, 5, 6, 7, 12, 13, 19, 23	0,7350
Perfect KNN(7)	1, 4, 5, 6, 7, 9, 10, 11, 12, 13, 16, 17, 20, 21, 22, 23, 27	0,7164
Perfect LDA	1, 3, 4, 5, 6, 7, 8, 10, 11, 18, 20	0,7199
Super Perfect KNN	4, 5, 6, 7, 12, 23	0,7385
Super Perfect All	4, 5, 6, 7, 23	0,7234

## KNN(7)

Número da coluna retirada	Nome da coluna retirada	Precisão	Comparação	Importante ou Problemático
Original	0	0,6949		
-1	<i>english</i>	0,6888	Diminuiu	Importante
-2	<i>required_age</i>	0,6962	Aumentou	Problema
-3	<i>achievements</i>	0,6977	Aumentou	Problema
-4	<i>positive_ratings</i>	0,6918	Diminuiu	Importante
-5	<i>negative_ratings</i>	0,6935	Diminuiu	Importante
-6	<i>average_playtime</i>	0,6949	Igual	Problema
-7	<i>median_playtime</i>	0,6912	Diminuiu	Importante
-8	<i>price</i>	0,6986	Aumentou	Problema
-9	<i>windows</i>	0,6929	Diminuiu	Importante
-10	<i>linux</i>	0,6942	Diminuiu	Importante
-11	<i>mac</i>	0,6938	Diminuiu	Importante
-12	<i>free_to_play</i>	0,6911	Diminuiu	Importante
-13	<i>action</i>	0,6992	Aumentou	Problema
-14	<i>indie</i>	0,6990	Aumentou	Problema
-15	<i>adventure</i>	0,6968	Aumentou	Problema
-16	<i>strategy</i>	0,6936	Diminuiu	Importante
-17	<i>rpg</i>	0,6944	Diminuiu	Importante
-18	<i>racing</i>	0,6960	Aumentou	Problema
-19	<i>simulation</i>	0,6959	Aumentou	Problema
-20	<i>casual</i>	0,6936	Diminuiu	Importante
-21	<i>sports</i>	0,6914	Diminuiu	Importante
-22	<i>violent</i>	0,6921	Diminuiu	Importante
-23	<i>nudity</i>	0,6949	Igual	Irrelevante
-24	<i>single_player</i>	0,6962	Aumentou	Problema
-25	<i>multiplayer</i>	0,6979	Aumentou	Problema
-26	<i>steam_cards</i>	0,6986	Aumentou	Problema
-27	<i>created_by_valve</i>	0,6914	Diminuiu	Importante

Nome do sub-treino	Número das colunas que se mantiveram	Precisão
Perfect KNN(3)	2, 4, 5, 6, 7, 12, 13, 19, 23	0,7428
Perfect KNN(7)	1, 4, 5, 6, 7, 9, 10, 11, 12, 13, 16, 17, 20, 21, 22, 23, 27	0,7413
Perfect LDA	1, 3, 4, 5, 6, 7, 8, 10, 11, 18, 20	0,7313
Super Perfect KNN	4, 5, 6, 7, 12, 23	0,7493
Super Perfect All	4, 5, 6, 7, 23	0,7408

## LDA

Número da coluna retirada	Nome da coluna retirada	Precisão	Comparação	Importante ou Problemático
Original	0	0,6883		
-1	<i>english</i>	0,6881	Diminuiu	Importante
-2	<i>required_age</i>	0,6896	Aumentou	Problema
-3	<i>achievements</i>	0,6883	Igual	Irrelevante
-4	<i>positive_ratings</i>	0,6874	Diminuiu	Importante
-5	<i>negative_ratings</i>	0,6855	Diminuiu	Importante
-6	<i>average_playtime</i>	0,6877	Diminuiu	Importante
-7	<i>median_playtime</i>	0,6866	Diminuiu	Importante
-8	<i>price</i>	0,6883	Igual	Irrelevante
-9	<i>windows</i>	0,6885	Aumentou	Problema
-10	<i>linux</i>	0,6875	Diminuiu	Importante
-11	<i>mac</i>	0,6864	Diminuiu	Importante
-12	<i>free_to_play</i>	0,6896	Aumentou	Problema
-13	<i>action</i>	0,6888	Aumentou	Problema
-14	<i>indie</i>	0,6916	Aumentou	Problema
-15	<i>adventure</i>	0,6888	Aumentou	Problema
-16	<i>strategy</i>	0,6887	Aumentou	Problema
-17	<i>rpg</i>	0,6885	Aumentou	Problema
-18	<i>racing</i>	0,6883	Igual	Irrelevante
-19	<i>simulation</i>	0,6888	Aumentou	Problema
-20	<i>casual</i>	0,6881	Diminuiu	Importante
-21	<i>sports</i>	0,6877	Aumentou	Problema
-22	<i>violent</i>	0,6881	Aumentou	Problema
-23	<i>nudity</i>	0,6885	Aumentou	Problema
-24	<i>single_player</i>	0,6885	Aumentou	Problema
-25	<i>multiplayer</i>	0,6874	Aumentou	Problema
-26	<i>steam_cards</i>	0,6896	Aumentou	Problema
-27	<i>created_by_valve</i>	0,6894	Aumentou	Problema

Nome do sub-treino	Número das colunas que se mantiveram	Precisão
Perfect KNN(3)	2, 4, 5, 6, 7, 12, 13, 19, 23	0,6864
Perfect KNN(7)	1, 4, 5, 6, 7, 9, 10, 11, 12, 13, 16, 17, 20, 21, 22, 23, 27	0,6912
Perfect LDA	1, 3, 4, 5, 6, 7, 8, 10, 11, 18, 20	0,6911
Super Perfect KNN	4, 5, 6, 7, 12, 23	0,6909
Super Perfect All	4, 5, 6, 7, 23	0,6912

## Conclusão

Após analisar os dados, é possível concluir que o melhor método, encontrado pelas nossas análises, para definir quão popular é um jogo é o KNN(7). Além disso, os dados (colunas) que fazem um jogo se aproximar mais dos jogos com mais usuários são: *positive\_ratings*, *negative\_ratings*, *average\_playtime*, *median\_playtime*, *free\_to\_play* e *nudity*, retirado do Super Perfect KNN, o *sub-treino* com maior precisão no KNN(7).

Sem levar em conta o *free\_to\_play* e o *nudity*, todos os outros dados do SuperPerfectKNN são gerados quando o jogo já está em mercado. Isso nos leva a concluir que não há preferência de gênero dos usuários da *Steam* no geral, e que a qualidade do jogo é mais importante que seu estilo, gênero, local de lançamento e coletáveis.

Além disso, o fato de o jogo ser *free\_to\_play* é muito relevante, pois qualquer um pode adicionar um jogo grátis na sua biblioteca de jogos, significando que muito mais gente possui o jogo.

## Referências

Dataset:

- <https://www.kaggle.com/nikdavis/steam-store-games/>

Auxílio no tratamento de dados:

- <https://steamspy.com/>

- <https://store.steampowered.com/>

Auxílio com a linguagem:

- <https://www.rdocumentation.org/>

- <https://www.r-bloggers.com/installing-r-packages/>