

Previsão do desempenho de estudantes ao jogar

Trabalho realizado no âmbito da Unidade Curricular Processamento de Big Data do 2º ano
em 2022/2023 da Licenciatura em Ciência de Dados

André Plancha, 105289
Andre_Plancha@iscte-iul.pt
Allan Kardec Rodrigues, 103380
aksrs@iscte-iul.pt

9 Abril 2023
Versão 1.0.0

Introdução

Jo Wilder and the Capitol Case é um jogo educacional sobre a história de Wisconsin, direcionado a crianças com idades entre os 8 e os 12 anos. O jogo de aventura e investigação segue a história de Jo Wilder, que descobre histórias sobre artefactos do passado do estado, investigando objetos, e encontrando pessoas. O jogo pode ser jogado no site oficial do PBD Wisconsin Education¹. Esta competição está a ser promovida pela plataforma kaggle².

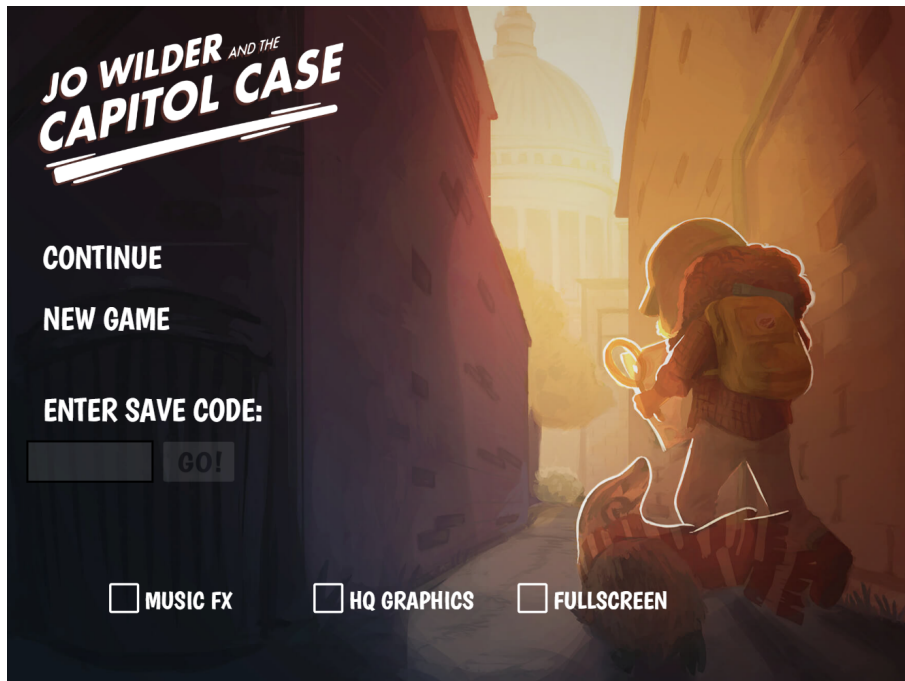


Figura 1: Imagem do menu principal

Este projeto tem como objetivo prever como os jogadores vão responder às 18 questões que o jogo apresenta, baseado na sua atividade durante o tal, implementando uma solução computacional para estudo e análise de dados de grande dimensão. Para isso, vamos usar Apache SparkTM e *PySpark* para processar os dados e a biblioteca *MLlib* para a criação de um modelo de aprendizagem supervisionada.

Análise exploratória dos dados

A competição disponibiliza um conjunto de ficheiros de dados, nos quais 2 são de interesse: `train.csv` e `train_labels.csv`, elas tendo as variáveis de interesse e as respetivas respostas, respetivamente. `train.csv` contém as seguintes colunas³, entre outras:

- `session_id`: Identificador da sessão do evento.
- `index`: Índice do evento na sessão.
- `elapsed_time`: Tempo decorrido desde o início da sessão em ms.
- `event_name`: Tipo do evento.
- `name`: Específicos do tipo de evento.

¹pbdwisconsineducation.org/jowilder/play-the-game/

²<https://www.kaggle.com/competitions/predict-student-performance-from-game-play>

³Mais informações sobre cada coluna está no site da competição, e no EDA em anexo

- `hover_duration`: O tempo de permanência do cursor sobre um objeto, em ms, se aplicável.
- `text`: O diálogo/monólogo que o jogador viu no evento, se aplicável
- `fqid`: O identificador único do evento.
- `fullscreen`: Se o jogador está em modo de ecrã inteiro.
- `hq`: Se o jogo está em alta qualidade.
- `music`: Se a música do jogo está ligada.
- `level_group`: O grupo de níveis a que o evento pertence.

; e `train_labels.csv` contém as seguintes colunas:

- `session_id`: Identificador da sessão do evento, em conjunto da pergunta que se pretende responder.
- `correct`: Se a pergunta está correta ou não.

Perante as colunas, como há uma incompatibilidade entre perguntas e respostas, nós decidimos criar uma tabela com cada sessão, características de tal (recursos) e as suas respostas, de forma usar classificação para prever as respostas.

Para criar estes recursos, foi necessário de uma análise dos tipos de eventos disponíveis. Uma análise mais aprofundada encontra-se em anexo, mas para resumir, os eventos estão divididos em 3 categorias: Exploração, exposição e revisão. Exploração são os eventos onde a personagem se move, interage com objetos opcionais, tenta entrar em zonas ainda não acedidas, etc... Exposição são os eventos onde a personagem interage com personagens, interage com objetivos principais, e é onde se encontra o que o utilizador aprende. Revisão são os eventos de quando o utilizador está perdido e portanto precisa de apoio de qual é o próximo passo.

A única exceção a estes eventos é o evento `checkpoint`, representando quando o utilizador começa a responder às perguntas (os eventos das respostas às perguntas não se encontra na tabela). Há 3 zonas onde o utilizador responde às perguntas, e na base de dados as 3 zonas são divididas pelos `checkpoints`, e cada zona está rotulada pelo seu `level_group`.

Análise e limpeza

Nós usámos como base de limpeza a nossa análise, e a análise de erros de um utilizador do kaggle⁴. Na Tabela 1 encontra-se os problemas encontrados e as soluções que nós usámos para os resolver.

Algumas das limpezas foram úteis na criação dos recursos, enquanto que outras apenas limpavam a base de dados. Nós não fizemos (nem registámos) algumas das limpezas, pois não seria benéfico para a criação dos recursos necessários.

Algo importante a notar também é que os roteiros dos jogos são diferentes para cada utilizador, entre 4 roteiros diferentes: *dry*, *nohumor*, *nosnark*, e *normal*. Um exemplo pode ser encontrado na Figura 2. A diferença entre roteiros é mínima na opinião dos autores, mas decidimos usar esta informação na mesma.

⁴<https://www.kaggle.com/code/abaojiang/eda-on-game-progress>

Tabela 1: Resolução dos problemas

Erros	Problemas	Solução	Notas
Tempo recorrido recuava no tempo	Cálculo de tempos entre eventos pode ficar inválido	Transformar em diferenças entre os tempos e tornar zero quando negativo	Máximo desta coluna continua afetado
Salto de <code>index</code>	O índice do evento não é sequencial	-	Suspeitamos que os eventos aconteceram na mesma, só não foram anotados
Sessões com menos ou mais de 3 checkpoints	As sessões são inválidas e dificultam análise	Foram removidas essas sessões	
Tempos de jogo demasiado elevados comparado com a maioria dos tempos	Dificulta o processo de aprendizagem	Tempos de jogo foram transformados e normalizados	

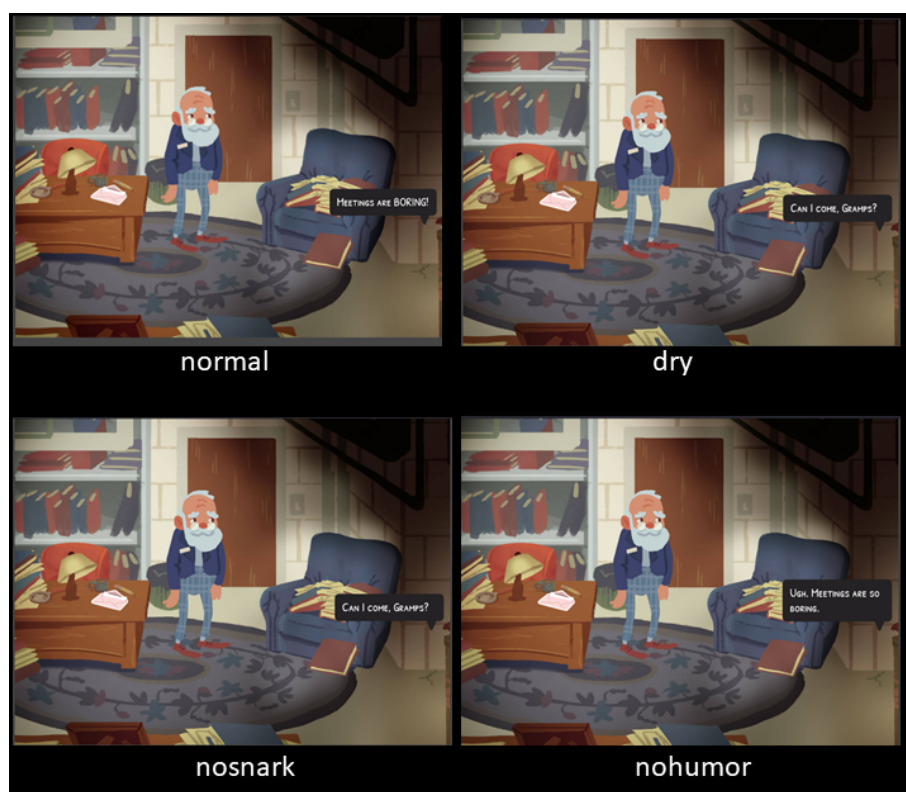


Figura 2: Diferenças entre roteiros

Processo

O processo da implementação da solução computacional pode ser descrito como caótico, sendo que como os dados não vieram com uma estrutura clara, foi necessário uma representação dos dados adaptada previamente ao análise de relações entre os dados; ou seja a criação de recursos teóricos foi executada *à priori*, e a análise de colineariadade e de relação foi feita *ad hoc*. Por outras palavras, a análise de dados foi feita ao longo do processo de criação de recursos, e a criação destes foram feitos apenas com análises comportamentais potências invés de análises estatísticas. Ainda assim, a Figura 3 pode ser visto o processo que foi seguido.

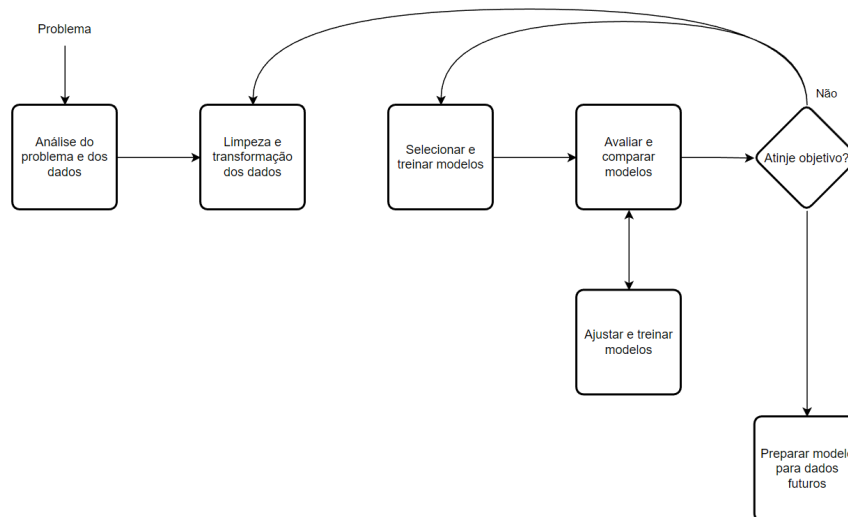


Figura 3: Processo

Recursos

As features que nós fizemos para cada sessão são entre as seguintes:

- O index máximo, que indica o número de eventos de cada sessão;
- Se a sessão esteve em tela cheia;
- Se a sessão estava em alta qualidade;
- Se a sessão tinha música ligada;
- Quantos objetivos opcionais a sessão interagiu com;
- Quantas salas inacessíveis a sessão tentou entrar em;
- Quantas vezes a sessão reviu o objetivo;
- Tempo médio de leitura do diálogo;
- Tempo médio por movimento do jogador;
- Tipo de roteiro.

Os promenores de como foram criadas encontram-se em anexo.

Exploração

Devido à natureza dos nossos dados a exploração não foi executada de uma forma tradicional, sendo que as análises menos teóricas foram feitas após a criação de cada recurso. De qualquer forma, estas encontram-se em anexo. Ainda assim, anotaremos aqui a análise de probabilidades das respostas.

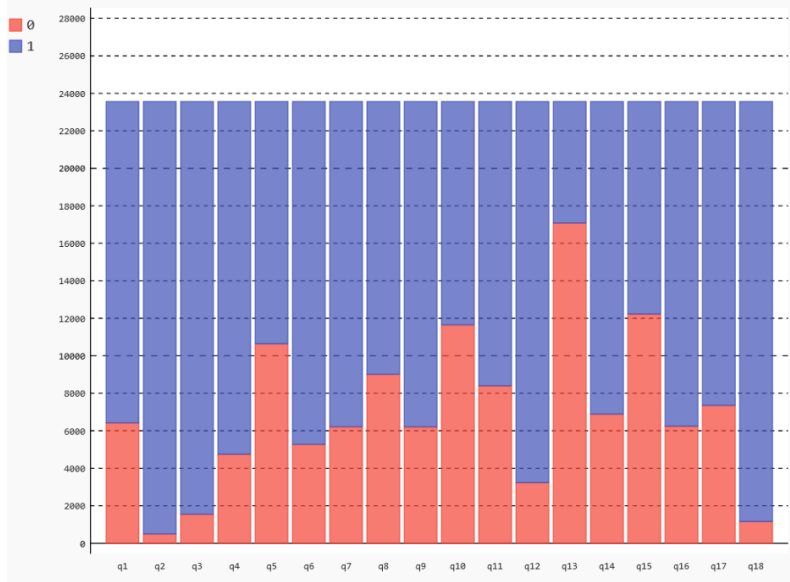


Figura 4: Probabilidades das respostas

Como se pode observar na Figura 4, algumas respostas são muito prováveis de acertar, e poucas são mais prováveis de errar. Notamos que há um acréscimo de perguntas erradas nas questões 5 e 13, sendo que essas são as primeiras do grupo. Ainda assim, não parece haver um padrão claro.

Modelação

Como processo de modelação, nós usamos aprendizagem supervisionada de classificação multi-rótulo, usando como rótulos cada pergunta certa, ou seja, cada sessão podia ter entre 0 e 18 rótulos. Nós consideramos várias opções de classificação, mas decidimos que a melhor forma para o nosso projeto foi com uma transformação para um problema de classificação binária, usando o método de *binary relevance*.

De forma simplificada, para cada rótulo, um modelo binário é treinado, independente dos outros rótulos; se 1 o rótulo é incluído, se 0 não. Esta estratégia fazia sentido, sendo que os nossos rótulos já são binários, o que facilita a modelação, e como temos uma grande quantidade de rótulos, outras estratégias (como *Label Powerset*) não faziam tanto sentido. Ainda assim, pode haver outros métodos mais adequados, como por exemplo *Classifier Chains*, sendo que não deve haver uma independência de rótulos.

Para a análise dos nossos modelos, usamos métricas de classificação multi-rótulo, como o *subset accuracy*, *hamming loss*, e *micro-F1*. Usamos também a métrica de avaliação *F1*, estrapulada e adaptada.

Subset accuracy é uma métrica que mede a percentagem de resultados que todos os rótulos estão certos. Previmos que esta métrica não seira muito útil, sendo que 18 rótulos corretos é muito difícil do nosso modelo conseguir, especialmente com o método de *binary relevance*.

$$\uparrow \text{subset-accuracy} = \frac{1}{|Y|} \sum_{i=1}^{|Y|} [\hat{Y}_i = Y_i]$$

Hamming loss é uma métrica que mede a percentagem de rótulos que estão errados. Sendo que este mede de forma mais individual, esta métrica é mais útil e direta.

$$\downarrow \text{hamming-loss} = \frac{1}{|Y| \cdot |L|} \sum_{i=1}^{|Y|} \sum_{l=1}^{|L|} [\hat{Y}_{i,l} \neq Y_{i,l}]$$

Micro-F1 é uma métrica que mede a percentagem de rótulos que estão certos, mas penaliza mais os rótulos que estão errados. Esta métrica é útil mais útil para comparar modelos, em vez de ser útil para avaliar o modelo individualmente.

$$\uparrow \text{micro-}F_1 = \frac{2TP}{2TP + FP + FN}$$

F1 é uma métrica usada para avaliar modelos de classificação binária, e é uma média ponderada entre a precisão e a revocação. A nossa métrica adaptada faz a média dos vários F1 scores de cada modelo que compoz o algoritmo. Esta métrica é também útil para comparar modelos, pois dá igual importância a cada modelo, enquanto que a métrica anterior tem em conta a quantidade de rótulos que cada X tem, no nosso caso isso equivaler às respostas certas, o que pode induzir em erro a métrica.

$$\uparrow F_1 = \frac{2}{|Y|} \sum_{i=1}^{|Y|} \frac{|Y_i \cap \hat{Y}_i|}{|Y_i| + |\hat{Y}_i|}$$

Nós começámos por treinar 4 modelos, de forma a compará-los e decidir que método usar, sendo eles regressão logística, máquina de vetores de suporte linear, árvore de decisão e *random forest*, todos eles com o método de *binary relevance*. Os promenores do executado e dos parametros estão em anexo. Os resultados podem ser observados na Tabela 2.

Tabela 2: Resultados dos primeiros modelos

Modelo	Subset accuracy \uparrow	Hamming loss \downarrow	Micro-F1 \uparrow	F1 \uparrow
Regressão Logística	0.0222	0.2477	0.8412	0.8239
Máquina de Vetores de Suporte Linear	0.0229	0.2517	0.8424	0.8285
Árvore de Decisão	0.0237	0.2440	0.8429	0.8268
Random Forest	0.0232	0.2442	0.8444	0.8295

Como se pode observar, estes modelos não mostram muita diferença nos resultados, mesmo que a árvore de decisão tenha um resultado melhor nas métricas individuais, e o Random Forest tenha obtido um melhor resultado nas métricas de comparação. Decidimos tentar melhorá-los, tirando proveito destes resultados. Logo, treiná-mos mais 4 modelos *Random Forest*, com parâmetros diferentes, especificamente na estratégia de divisão em cada nó. Mais promenores sobre este processo estão na documentação da classe *RandomForest* ⁵. Os resultados pdoem ser observados na Tabela 3.

Tabela 3: Resultados dos segundos modelos

Modelo	Subset accuracy \uparrow	Hamming loss \downarrow	Micro-F1 \uparrow	F1 \uparrow
RF - all	0.0247	0.2420	0.84481	0.8288
RF - sqrt/log2	0.0232	0.2442	0.84442	0.8295
RF - onethird	0.0242	0.2435	0.84480	0.8298

⁵<https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.mllib.tree.RandomForest.html>

Como se pode observar, as diferenças são minúsculas o suficiente que dificilmente se considera um melhor que o outro. Ainda assim, o modelo *RF* com estratégia *all* tem um resultado individual melhor, e *RF* com estratégia *onethird* tem um resultado de comparação melhor. Como os resultados não parecem melhorar tanto assim, e como temos já bons resultados com os modelos apresentados, decidimos não continuar a modelar.

Pela a análise das matrizes de confusão (em anexo), podemos observar muitas das perguntas sem positivos ou negativos. Isto é devido à própria amostra que tínhamos acesso não ter muitas respostas positivas ou negativas para as perguntas em questão. Isto devia ter sido verificado antes de começar, e devia ter-se adaptado a forma de divisão em conjunto de treino ou teste, ou devia ter-se usado validação cruzada.

Para melhorar o resultado, além do referido podia-se também ter-se melhorado a limpeza dos dados, sendo que há muitas sessões com valores inválidos, que podiam ter sido tratados e/ou as sessões removidas do conjunto de treino.

```
[[ 0, 392],
 [ 0, 18311]],

[[ 10, 1190],
 [ 0, 17503]],
```

Figura 5: Exemplos de matrizes de confusão resultadas por rótulo

Tempo de execução

De forma a comparar as ferramentas que tinha-mos disponíveis, decidimos medir o tempo de execução dos 4 primeiro modelos apresentados, uma primeira vez com o ambiente local, outra com 5 clusters da Amazon EMR (AWS). O ambiente local usado é composto por um AMD Ryzen 7 5800X com 8 cores (16 Threads). Os resultados podem ser observados na Tabela 4.

Tabela 4: Tempo de execução

Modelo	Local (ms)	AWS (ms)	Diferença (%)
Regressão Logística	16.743	51.014	304.7
SVM	41.470	112.787	272.0
Árvore de Decisão	9.962	35.154	352.9
Random Forest	15.247	46.656	306.0

Como podemos observar, o uso do AWS no nosso caso piorou o tempo de execução. No entanto, não descartamos a ferramenta para ambientes pessoais menos potentes, e a sua utilidade também pode ser observada para outros casos específicos. Também não tiramos proveito de todas as funcionalidades do serviço, portanto estes números podem não estar otimizados.