

Eduardo Darrazão - 1906399

Marcelo Guimarães da Costa - 1937570

Leandro Batista de Almeida - Professor

Big Data e Aplicações

06 de dezembro de 2023

Projeto 1: Aquisição e Processamento de um Dataset

Este projeto é a parte 1 de um projeto final de duas partes, sendo o objetivo final uma apresentação para a classe que envolveria todos os processos desde a seleção de um *dataset* até a utilização de modelos de *AutoML* sobre o *dataset* escolhido.

Nesta parte do projeto, os processos incluem:

- Selecionar *dataset* (fontes públicas ou privadas)
- Carregamento de dados (CSV, importação de bancos de dados, etc)
- Tratamento de informações
- Validação de erros
- Preenchimento de valores faltantes
- Limpeza de dados

1 – SELEÇÃO DO DATASET

O *dataset* escolhido foi do IMDB¹ (*Internet Movie Database*), que abrange séries e filmes de vários gêneros, bem como a equipe envolvida e a média das avaliações e números de votos. Os arquivos disponíveis são os seguintes:

¹ <https://www.imdb.com>

- *title.basics*: Informações principais sobre cada título, bem como gênero, e tipo (filme, série, vídeo, etc).
- *title.akas*: Traduções e localizações do nome do título para diferentes culturas.
- *title.principals*: Dados sobre cada membro da equipe de um título, bem como sua função naquele título (atores, escritores, diretores, etc).
- *title.crew*: Informações sobre a equipe de direção e escritores em cada título.
- *title.episode*: Dados sobre episódios no caso de séries, conectando a série com o episódio.
- *title.ratings*: Média de avaliações (1 a 10) e número de votos para certos títulos.
- *name.basics*: Mais informações sobre alguns registros de *title.principals*. Data de nascimento, nome, pelo que é conhecido, etc.

O *dataset* está disponível em <https://developer.imdb.com/non-commercial-datasets/>.

Os dados são atualizados diariamente, e foi realizada a aquisição na data 02/12/2023.

2 – CARREGAMENTO DOS DADOS

Para tratamento dos dados foi utilizada a biblioteca *pyspark*, mais especificamente utilizando a extensão Spark Dataframes. Segue um exemplo de carregamento:

```
title_basics = spark.read.load('../Data/title.basics.tsv', format='csv', sep='\t', inferSchema=True, header=True)
```

3 – TRATAMENTO, LIMPEZA E TRANSFORMAÇÃO DOS DADOS

A primeira etapa foi selecionar apenas os filmes, ou *movies*, que representam nosso corte de interesse. Dessa forma, começamos a trabalhar com um *dataset* de 664 mil linhas.

```
+-----+-----+
|  titleType|  count|
+-----+-----+
|  tvEpisode|7909714|
|      short| 964958|
|      movie| 664652|
|      video| 283511|
|   tvSeries| 253097|
|   tvMovie| 143596|
|tvMiniSeries|  51367|
|  tvSpecial|  44591|
| videoGame|  36627|
|   tvShort|  10084|
|   tvPilot|      1|
+-----+-----+
```

Em seguida foi mantido somente os dados que possuem gênero(s), utilizando a coluna *genres*. Ela contém informações que consideramos importantes para realizar a classificação de *rating* que planejamos. Assim, foram removidos mais de 70 mil registros, sobrando 590 mil.

Por fim, filtramos nosso *dataset* novamente pela junção interna (*inner join*) com a tabela *title.ratings*. Como *averageRating* será nosso campo *target* para os modelos construídos na parte 2, a existência do registro em *title.ratings* é essencial. Isso nos deixou com uma quantia de 290 mil linhas.

O *dataset* apresenta algumas colunas com dados relevantes faltantes, como no *title.basics*, as colunas *runtimeMinutes* e *startYear*. Estas contêm o valor \N em diversas

entradas, que neste *dataset* significa a ausência de valor, ou nulo. Estes foram substituídos por 0, para que não tenhamos mais perda de dados.

Todo o código está disponível em: <https://github.com/eduponto21/IMDB-BigData-Spark> .