

Bootcamp: Engenheiro(a) de Dados (Cloud)

Trabalho Prático

Módulo 2: Tecnologias de Big Data – Processamento de Dados Massivos

Objetivos de Ensino

Exercitar os seguintes conceitos trabalhados no Módulo:

1. Contexto de Big Data e das ferramentas de processamento massivo.
2. Fundamentos do Spark.
3. Funcionamento interno do Spark.
4. Manipulação de dados com Spark.

Enunciado

Durante a primeira metade do módulo, as aulas práticas utilizaram um conjunto de dados gratuito disponível no site do IMDB, com informações de filmes séries e outras produções audiovisuais. Neste trabalho prático, você deverá fazer o download das tabelas **title.basics** e **title.ratings** do site oficial (<https://datasets.imdbws.com/>) e realizar um processo de limpeza nos dados vistos em nas aulas, utilizando o Apache Spark. Será necessário alterar os tipos das colunas, tratar os valores nulos e realizar algumas análises com os dados.

Atividades

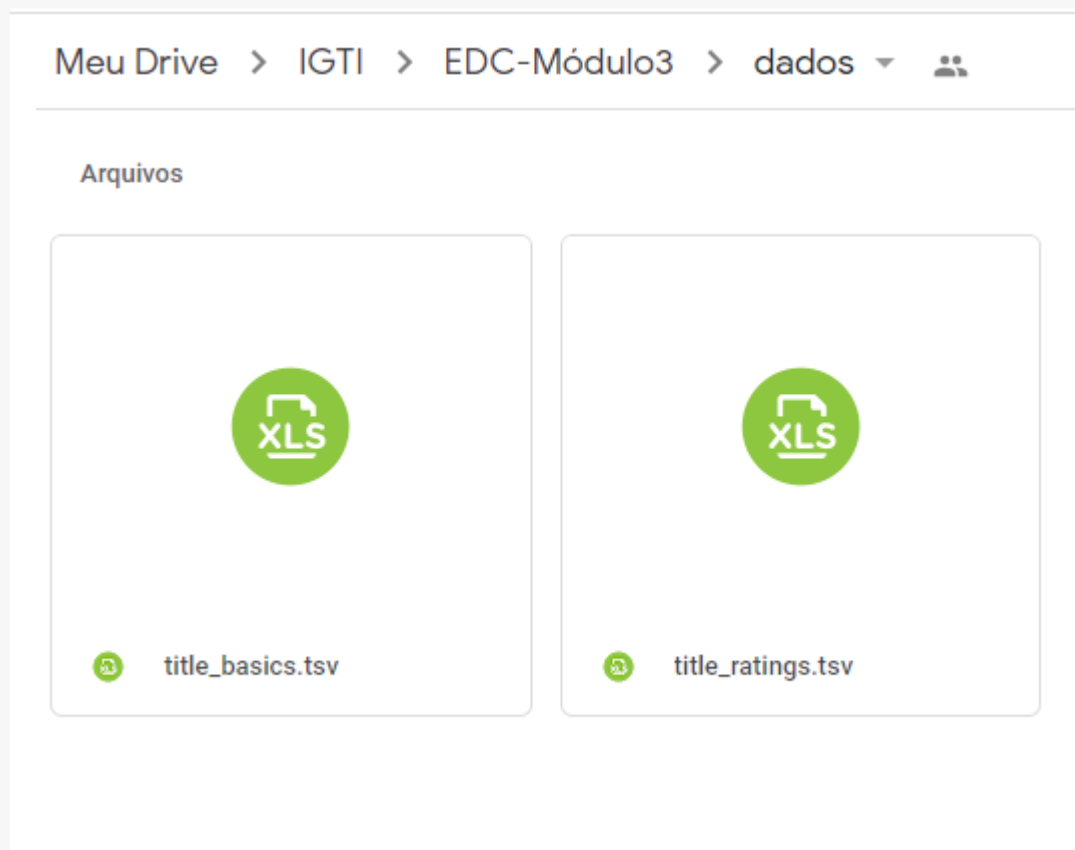
Os alunos deverão desempenhar as seguintes atividades:

Em seguida, será necessário fazer o download das tabelas no seguinte link:

<https://drive.google.com/drive/u/1/folders/1DfuJmIOsXU8hgCRUui-89FQhhHh3L5kS>

Obs.: Faça isso com antecedência, pois as tabelas a serem baixadas são grandes.

Figura 2 – Download dos Dados



Importante: É necessário que os dados estejam localizados na mesma pasta em que o shell do Spark esteja sendo executado para que eles possam ser lidos de forma mais fácil. Por isso, observe bem onde a shell está sendo executada ou mude o diretório da execução para a mesma pasta dos dados:

Figura 3 – Identificando o diretório de execução da Shell do Spark


```
df_ratings = spark.read.csv('title_ratings.tsv', header=True, sep='\t')
```

A partir desse momento, os dois DataFrames estarão disponíveis para manipulação no ambiente de execução. A prática consiste em manipular os DataFrames de forma a responder algumas das perguntas apresentadas. Além disso, temos algumas perguntas teóricas sobre o Spark e seu funcionamento.

Para um dicionário dos dados sendo trabalhados, acesse:

<https://www.imdb.com/interfaces/>

Respostas Finais

Os alunos deverão desenvolver a prática e, depois, responder às seguintes questões objetivas: