

# Gestão de Grandes Conjuntos de Dados

## 2º Trabalho Prático

-

Mestrado em Engenharia Informática  
Universidade do Minho

### **Grupo nº 8**

---

PG41080	João Ribeiro Imperadeiro
PG41081	José Alberto Martins Boticas
PG41091	Nelson José Dias Teixeira
PG41851	Rui Miguel da Costa Meira

14 de maio de 2020

# Conteúdo

<b>1</b>	<b>Introdução</b>	<b>3</b>
<b>2</b>	<b>Implementação</b>	<b>4</b>
2.1	Configuração . . . . .	5
2.2	1ª Tarefa . . . . .	5
2.2.1	<i>Log</i> . . . . .	5
2.2.1.1	Alternativa . . . . .	5
2.2.2	<i>Top3</i> . . . . .	5
2.2.2.1	Alternativa . . . . .	5
2.2.3	<i>Trending</i> . . . . .	5
2.2.3.1	Alternativa . . . . .	5
2.3	2ª Tarefa . . . . .	5
2.3.1	<i>Top10</i> . . . . .	5
2.3.1.1	Alternativa . . . . .	5
2.3.2	<i>Friends</i> . . . . .	5
2.3.2.1	Alternativa . . . . .	5
2.3.3	<i>Ratings</i> . . . . .	5
2.3.3.1	Alternativa . . . . .	5
2.4	3ª Tarefa . . . . .	5
<b>3</b>	<b>Conclusão</b>	<b>6</b>
<b>A</b>	<b>Observações</b>	<b>7</b>

# Lista de Figuras

# Capítulo 1

## Introdução

Neste trabalho prático é requerida a concretização e avaliação experimental de tarefas de armazenamento e processamento de dados através do uso da ferramenta computacional *Spark* (*batch* e *streaming*). Por forma a realizar estas tarefas, são utilizados os dados públicos do *IMDb*, que se encontram disponíveis em:

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

Para além destes dados, é também utilizado um gerador de *streams*, baseado nos mesmos, que simula uma sequência de votos individuais de utilizadores. Este utensílio foi desenvolvido pelo docente desta unidade curricular e encontra-se disponível na plataforma *Blackboard*.

Ao longo deste documento vão também ser expostos todos os passos tomados durante a implementação das tarefas pedidas neste projeto, incluindo as decisões tomadas pelos elementos deste grupo a nível de algoritmos e parâmetros de configuração. Para além disso são ainda apresentadas todas as instruções que permitem executar e utilizar corretamente os programas desenvolvidos. Por fim, na fase final deste manuscrito, são exibidos os objetivos atingidos após a realização das tarefas propostas.

De salientar também que durante os capítulos que se seguem são identificadas algumas alternativas para concretizar as tarefas indicadas neste trabalho prático.

## Capítulo 2

# Implementação

Para a realização com sucesso deste trabalho prático, é solicitada a elaboração de três tarefas. Apresentam-se de seguida as mesmas:

1. Desenvolver uma componente de processamento de *streams* que produza os seguintes resultados:
  - **Log**: armazenar todos os votos individuais recebidos, etiquetados com a hora de chegada aproximada ao minuto, em lotes de 10 minutos. Cada lote deve ser guardado num ficheiro cujo nome identifica o período de tempo;
  - **Top3**: exibir a cada minuto o top 3 dos títulos que obtiveram melhor classificação média nos últimos 10 minutos;
  - **Trending**: apresentar a cada 15 minutos os títulos em que o número de votos recolhido nesse período sejam superiores aos votos obtidos no período anterior, independentemente do valor dos votos.
2. Implementar uma componente de processamento em *batch* que permita realizar as seguintes tarefas:
  - **Top10**: calcular o top 10 dos atores que participaram em mais títulos diferentes;
  - **Friends**: computar o conjunto de colaboradores de cada ator (i.e., outros atores que participaram nos mesmos títulos);
  - **Ratings**: atualizar o ficheiro "*title.ratings.tsv*" tendo em conta o seu conteúdo anterior e os novos votos recebidos até ao momento.
3. Escolher a configuração e a implementação que, para o mesmo *hardware*, permite receber e tratar o maior débito de eventos. Esta tomada de decisão deve ser devidamente justificada com recurso a resultados experimentais.

Nas próximas secções são evidenciadas as implementações para cada uma destas tarefas bem como algumas sugestões alternativas que poderiam ser tomadas em consideração.

## 2.1 Configuração

### 2.2 1ª Tarefa

#### 2.2.1 *Log*

##### 2.2.1.1 Alternativa

#### 2.2.2 *Top3*

##### 2.2.2.1 Alternativa

#### 2.2.3 *Trending*

##### 2.2.3.1 Alternativa

### 2.3 2ª Tarefa

#### 2.3.1 *Top10*

##### 2.3.1.1 Alternativa

#### 2.3.2 *Friends*

##### 2.3.2.1 Alternativa

#### 2.3.3 *Ratings*

##### 2.3.3.1 Alternativa

### 2.4 3ª Tarefa

## Capítulo 3

## Conclusão

# Apêndice A

## Observações

- Documentação *Java* 8:  
`https://docs.oracle.com/javase/8/docs/api/`
- *Maven*:  
`https://maven.apache.org/`
- *Apache Spark*:  
`https://spark.apache.org/`
- *Docker*:  
`https://www.docker.com/`