

Extreme Computing

Assignment 2

1 Introduction

This is the second coursework assignment for the Extreme Computing course 2021/22. You need to use Apache Spark to solve problems you might encounter when working with collections. This section will give you administrative information and help with solving the assignment. This is followed by the actual tasks and finally a description of how to submit your solutions.

1.1 Administrative Information

Deadline The assignment is due at 4:00pm on November 26th.

Deadline Extension The School of Informatics has a policy on coursework deadlines, which applies across all taught courses. Further information can be found here:

<http://web.inf.ed.ac.uk/infweb/student-services/ito/admin/coursework-projects/late-coursework-extension-requests>

Questions All questions should go on Piazza

<https://piazza.com/class/ktlmf55uoar593>

in the “hw2” folder. Feel free to discuss general techniques amongst each other unless you would reveal an answer. If your question / discussion reveals an answer, ask privately.

Marking The assignment is worth 50 marks in total and counts for 20% of the final course mark. Marks are given for correctness, efficiency and proper use of tools.

Marks are indicated on the right margin of the page by two numbers, e.g. **1+3 marks**. The first number indicates achievable marks for correctness while the second number indicates achievable marks for efficiency.

Submission The submission process for your solution is described in Section 4. We start marking at the deadline and only mark once. If you are submitting on time, you can submit as many times as you want and the last one will be marked. If you are submitting late, you may only submit once in total (which implies that you should not submit before the deadline) or run the risk that we will mark an old version then penalise for the last submission time.

Marking Feedback You will receive your marks and feedback for your solution on LEARN. Once you have received your marks, you have three days to ask questions about them, e.g. feedback clarification. Please post any such questions in a private Piazza message.

Good Scholarly Practice Please remember the University requirement as regards all assessed work for credit. Details about this can be found at:

[http://web.inf.ed.ac.uk/infweb/admin/policies/
academic-misconduct](http://web.inf.ed.ac.uk/infweb/admin/policies/academic-misconduct)

Furthermore, you are required to take reasonable measures to protect your assessed work from unauthorised access. For example, if you put any such work on a public repository then you must set access permissions appropriately (generally permitting access only to yourself, or your group in the case of group practicals).

2 Tasks

For the dataset you will use in this assignment, there is a `imdb-small-data.zip` file that can be accessed at:

<https://amirsh.github.io/files/exc21/imdb-small-data.zip>

You can copy the extracted `tsv` files to `src/main/resources/imdb/` for local testing and debugging.

2.1 Internet Movie Database (IMDB)

This assignment will also be on processing the IMDB dataset. However, this time you need to use Apache Spark to efficiently process structured text. Note that you are only allowed to **use Spark RDDs** (e.g., `map`, `groupBy`, `join`). This means that you are NOT allowed to use other facilities provided by the Spark ecosystem including Spark SQL.

You can refer to the first assignment for the description of the IMDB dataset.

3 Tasks

Download `imdb-spark-src.zip` and extract it somewhere on your machine. You have to complete the missing implementations (specified by `???`) in `src/main/scala/imdb/ImdbAnalysis.scala`.

You are encouraged to look at the Apache Spark API documentation while solving this exercise:

<https://spark.apache.org/docs/3.0.3/api/scala/org/apache/spark/rdd/index.html>

Task 1

◀ Task

Calculate the average, minimum, and maximum runtime duration for all titles per movie genre.

Note that a title can have more than one genre, thus it should be considered for all of them. The results should be kept in minutes and titles with 0 runtime duration are valid and should be accounted for in your solution.

(5+6 marks)

```
return type: RDD[(Float, Int, Int, String)]
avg_runtime:Float
min_runtime:Int
max_runtime:Int
genre:String
```

Task 2

◀ Task

Return the titles of the movies which were released between 1990 and 2018 (inclusive), have an average rating of 7.5 or more, and have received 500000 votes or more.

For the titles use the `primaryTitle` field and account only for entries whose `titleType` is 'movie'.

(6+7 marks)

```
return type: RDD[String]
title:String
```

Task 3

◀ Task

Return the top rated movie of each genre for each decade between 1900 and 1999.

For the titles use the `primaryTitle` field and account only for entries whose `titleType` is 'movie'. For calculating the top rated movies use the `averageRating` field and for the release year use the `startYear` field.

The output should be sorted by decade and then by genre. For the movies with the same rating and of the same decade, print only the one with the title that comes first alphabetically. Each decade should be represented with a single digit, starting with 0 corresponding to 1900-1909.

(6+7 marks)

```
return type: RDD[(Int, String, String)]
decade:Int
genre:String
title:String
```

Task 4

◀ Task

In this task we are interested in all the crew names (`primaryName`) for whom there are at least two known-films released since the year 2010 up to and including the current year (2021). You need to return the crew name and the number of such films.

(6+7 marks)

```
return type: RDD[(String, Int)]
crew_name:String
film_count:Int
```

4 Submissions

To submit your work, please do the following:

1. To test the correctness of your program, you can use the `test` command of `sbt`. We will run your submissions against a bigger test suite than we provided to you. So if you want full points, be sure to add your own tests to double check that you have caught all corner cases.
2. We also check the performance of your implementation. Thus, make sure that you use the best possible implementation in terms of space and time complexity and communication time (e.g., not unnecessarily running actions).
3. To submit your code, please zip the entire `imdb-spark` directory and name it as `imdb-spark.zip`. Make sure that the data files are removed from the `resource` directory and remove the `target` directory.
4. Upload the zip file to Learn.