

MapReduce w Sparku III

3 grudnia 2018

Opis pliku z zadaniami

Wszystkie zadania na zajęciach będą przekazywane w postaci plików `.pdf`, sformatowanych podobnie do tego dokumentu. Zadania będą różnego rodzaju. Za każdym razem będą one odpowiednio oznaczone:

- Zadania do wykonania na zajęciach oznaczone są symbolem \triangle – nie są one punktowane, ale należy je wykonać w czasie zajęć.
- Punktowane zadania do wykonania na zajęciach oznaczone są symbolem \diamond – należy je wykonać na zajęciach i zaprezentować prowadzącemu, w wypadku nie wykonania zadania w czasie zajęć lub nieobecności, zadania staje się zadaniem do wykonania w domu (\star).
- Zadania do wykonania w domu oznaczone są symbolem \star – są one punktowane, należy je dostarczyć w sposób podany przez prowadzącego i w wyznaczonym terminie (zwykle przed kolejnymi zajęciami).
- Zadania programistyczne można wykonywać w dowolnym języku programowania, używając jedynie biblioteki standardowej dostępnej dla tego języka.

1 Zapytania do danych MSDC

★

Treść

Wykorzystując pliki `.csv`, przygotowane na zajęciach dotyczących przetwarzania danych MSDC w powłocie `bash`, zapisz w Sparku poniższe zapytania (są to te same zapytania, jak na wcześniejszych zajęciach dotyczących danych MSDC):

- Ranking popularności piosenek,
- Ranking użytkowników ze względu na największą liczbę odsłuchanych unikalnych utworów,
- Artysta z największą liczbą odsłuchań,
- Sumaryczna liczba odsłuchań w podziale na poszczególne miesiące,
- Wszyscy użytkownicy, którzy odsłuchali wszystkie trzy najbardziej popularne piosenki zespołu Queen.

Można wykorzystać następujące instrukcje: `spark.read.csv`, `groupBy`, `count`, `join`, `select`, `orderBy`, `show`, `agg`, `countDistinct`, `col`, `toDF`, `limit`, `filter`.

Dla ułatwienia wykonania zadania poniżej przedstawiony jest kod dla pierwszego zapytania:

```
1 //Read data
2 val songs = spark.read.
3     option("delimiter", ",").
4     csv("songs").
5     toDF("song_id", "track_long_id", "
6         song_long_id", "artist", "song")
7 val facts = spark.read.
8     option("delimiter", ",").
9     csv("facts").
10    toDF("id", "user_id", "song_id", "
11        date_id")
12 //The most popular songs
13 facts.groupBy("song_id").
14    count.
15    join(songs, facts("song_id")===songs("song_id")).
16    select("song", "count").
17    orderBy(desc("count")).
18    show(10)
```

spark-msdc-1.scala

Za powyższe zadanie można otrzymać punkty bonusowe, po 1 punkcie za każde poprawnie wykonane zapytanie.