

NTNU

TDT4225 STORE, DISTRIBUERTE DATAMENGDER

ASSIGNMENT #4

Øving 4

Medvirkende:

Brede Yabo KRISTENSEN

March 13, 2019

1 RDD API Tasks

1.1 Task 1

The files we are reading are rows with columns separated by , .

I first split each column using , as a delimiter. This is done by using the `map()` function as it applies a function to all of the rows. The result is yet another RDD. When splitting the rows we get an array for each row, we then want the fourth column as it is the genre column. We could use `flatMap()`, but we want an RDD that has the same number of rows.

We then use the function `distinct()` to remove all duplicates from the RDD. Finally we count the number of rows using the `count()` function, which returns a number.

1.2 Task 2

We start the same way as last time, only now we get the fourth column from the `artist.csv` file (year). We also make sure that the rows are converted to integers as we will compare them later. We then reduce the RDD using the `reduce()` function. We want to return the oldest artist and do this by comparing all the rows.

Finally we print the oldest artist birth date.

1.3 Task 3

Like before, we split the lines but now as a tuple with the integer 1, used for counting later. We then use `reduceByKey()` which executes the function for each row that has the same key, in our case, we count each time and add it to the counter in the tuple. This leaves us with the number of artists for each country.

The task specifies that we sort the total number of artists for each country descending and artists alphabetically for countries with the same number of artists. This is done by first sorting alphabetically and then sorting by count.

When sorting the keys, we just flip the tuple and `sortByKey` (with `false` for descending sort) and then flip back. The result is outputted to `result3directory`

1.4 Task 4

We begin by mapping out $\text{artist}_i, \text{dc_column}$ and keeping_a_number1 in the tuple.

1.5 Coalesce vs repartition