

NTNU

TDT4225 STORE, DISTRIBUERTE DATAMENGDER

ASSIGNMENT #4

Øving 4

Medvirkende:

Brede Yabo KRISTENSEN

March 12, 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 have used `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