MovieLens 100K Dataset Analysis with Cassandra and Spark

READY

MovieLens data sets were collected by the GroupLens Research Project at the University of Minnesota.

This data set consists of 100,000 ratings (1-5) from 943 users on 1682 movies and each user has rated at least 20 movies. It contains simple demographic info for the users (age, gender, occupation, zip)

The project aims to use Cassandra Query Language (CQL) and Spark2 Structured Query Language (SQL) to analyze the user info from the datasets to answer following questions:

- 1. Calculate the average rating for each movie.
- 2. Identify the top ten movies with the highest average ratings.
- 3. Find the users who have rated at least 50 movies and identify their favourite movie genres.
- 4. Find all the users who are less than 20 years old.
- 5. Find all the users whose occupation is "scientist" and whose age is between 30 and 40 years old.

This anlaysis uses Docker to run Zeppelin and Cassandra

1. Setting up the Docker

READY

a. Create a Docker network

docker network create zeppelin-net

b. Start Cassandra container

docker network create zeppelin-net

c. Start Zeppeling container

```
docker run -d --name zeppelin --network zeppelin-net -p 9090:8080 apache/zeppelin:0.12.0
```

d. Install Spark in Zeppelin container

```
docker exec -it zeppelin bash

cd /tmp

wget https://archive.apache.org/dist/spark/spark-3.1.2/spark-3.1.2-bin-hadoop3.2.tgz

tar -xvzf spark-3.1.2-bin-hadoop3.2.tgz
```

e. Reconfigure the interpreter

```
SPARK_HOME -> /tmp/spark-3.1.2-bin-hadoop3.2
spark.master -> local[*]
spark.cassandra.connection.host -> Cassandra
spark.jars.packages -> com.datastax.spark:spark-cassandra-connector_2.12:3.1.0
zeppeline.spark.enableSupportVersionCheck -> false
```

2. Uploading files from local into Zeppelin

```
docker cp "C:\Users\hajar\Downloads\ml-100k\u.user" zeppelin:/tmp/u.user

docker cp "C:\Users\hajar\Downloads\Assignemnt_3\ml-100k\u.data" zeppelin:/tmp/u.data

docker cp "C:\Users\hajar\Downloads\Assignment_3\ml-100k\u.item" zeppelin:/tmp/u.item
```

3. Create keyspaces and tables in Cassandra

READY

In Cassandra Shell

```
docker exec -it cassandra cqlsh
```

a) Create the keyspace

```
CREATE KEYSPACE movielens
WITH replication = {'class': 'SimpleStrategy', 'replication_factor': 1};
USE movielens;
```

b) Create the tables

```
CREATE TABLE users (
  user id int PRIMARY KEY,
  age int,
  gender text,
  occupation text,
  zip_code text
);
CREATE TABLE ratings (
    user id int,
   movie_id int,
   rating int,
   timestamp bigint,
   PRIMARY KEY (user_id, movie_id)
);
CREATE TABLE titles (
  movie_id int PRIMARY KEY,
  title text,
  release date text,
  imdb_url text,
  unknown int,
  action int,
  adventure int,
  animation int,
  children int,
  comedy int,
  crime int,
  documentary int,
  drama int,
  fantasy int,
  film_noir int,
  horror int,
  musical int,
  mystery int,
  romance int,
  sci_fi int,
  thriller int,
  war int,
```

```
western int
);
```

```
%pyspark

#test Spark
spark.range(1).show()

+---+
| id|
+---+
| 0|
+---+
```

4. Loading data from Zeppelin and write into Cassandra.

READY

READY

Here we create RDD (Resilient Distributed Dataset) objects

```
%pyspark
#libraries

from pyspark.sql import Row
from pyspark.sql import functions as F
```

User Data

```
#to parse function
def parseInput(line):
    x = line.split('|') # the delimiter used is pipe
    return Row(user_id=int(x[0]), age=int(x[1]), gender=x[2], occupation=x[3], zip_code=x[4])
```

```
# Load raw u.user file
lines = spark.sparkContext.textFile("file:///tmp/u.user")
 # Apply transformation
 users = lines.map(parseInput)
 # Create DataFrame
 usersDataset = spark.createDataFrame(users)
 # Write to Cassandra
 usersDataset.write.format("org.apache.spark.sql.cassandra") \
      .mode("append") \
      .options(table="users", keyspace="movielens").save()
 # Read back from Cassandra
 readUsers = spark.read.format("org.apache.spark.sql.cassandra").options(table="users", keyspace="movielens").load()
|user id|age|gender| occupation|zip code|
+----+
    885 | 30 |
                 F
                                   95316
                          other
     791 31
                 Μİ
                        educator
                                   20064
     176 28
                      scientist
                                   07030
     393 | 19 |
                 Μİ
                         student
                                   83686
     582 | 17 |
                 M |
                        student
                                   93003
     371 | 36 |
                 Μ
                        engineer|
                                   99206
     561 23
                 Μĺ
                        engineer|
                                   60005
                                   68106
     89 | 43 |
                 F|administrator|
                 M|administrator|
     78 | 26 |
                                   61801
     265 | 26 |
                       executive|
                                   84601
     216 | 22 |
                 Μĺ
                        engineer|
                                   02215
     243 | 33 |
                 Μĺ
                        educator
                                   60201
                                   55337
     674 13
                 F|
                        student
     111 57
                        engineer
                                   90630
     721 401
                 Fladminictnatanl
```

Ratings Data READY

%pyspark

from pyspark.sql import Row
from pyspark.sql import functions as F

```
# ParseInput Function
 def parseInput2(line):
     x = line.split('\t') # the delimiter use is tab
     return Row(user id=int(x[0]), movie id=int(x[1]), rating=x[2], timestamp=x[3])
 # Load raw u.data file
lines = spark.sparkContext.textFile("file:///tmp/u.data")
 # mapping to row
 ratings = lines.map(parseInput2)
 # Create DataFrame
 ratingsDataset = spark.createDataFrame(ratings)
 # Write to Cassandra
ratingsDataset.write.format("org.apache.spark.sql.cassandra").mode("append").options(table="ratings", keyspace="movielens").save()
 # Read back from Cassandra
 readRatings = spark.read.format("org.apache.spark.sql.cassandra") \
     .options(table="ratings", keyspace="movielens").load()
 readRatings.show()
+----+
|user id|movie id|rating|timestamp|
+----+
                      4 | 891273683 |
     384
             258
                      4 | 891283502 |
     384
             271
                      5 | 891273509 |
     384
             272
                      4 | 891273649 |
     384
             286
             289
                      5 | 891283502 |
     384
                      4 | 891273809 |
     384
             300
                      5 | 891273509 |
     384
             302
                      5 | 891273683 |
     384
             313
     384
             316
                      5 | 891274055 |
             327
                      4 | 891273761 |
     384
    384
                      4 | 891274091 |
             328
     384
             329
                      3 | 891273761 |
     384
             333|
                      4 | 891273509 |
     384
              343
                      3 | 891273716 |
                      4 L001 272 E00 L
     2011
             2471
```

Movies data READY

```
%pyspark
from pyspark.sql import Row
from pyspark.sql import functions as F
# ParseInput Function
def parseInput3(line):
    fields = line.split('|') # the delimiter use is pipe
    return Row(
        movie id=int(fields[0]),
        title=fields[1],
        release date=fields[2],
        imdb url=fields[4],
        unknown=int(fields[5]),
        action=int(fields[6]),
        adventure=int(fields[7]),
        animation=int(fields[8]),
        children=int(fields[9]),
        comedy=int(fields[10]),
        crime=int(fields[11]),
        documentary=int(fields[12]),
        drama=int(fields[13]),
        fantasy=int(fields[14]),
        film noir=int(fields[15]),
        horror=int(fields[16]),
        musical=int(fields[17]),
        mystery=int(fields[18]),
        romance=int(fields[19]),
        sci fi=int(fields[20]),
        thriller=int(fields[21]),
        war=int(fields[22]),
        western=int(fields[23])
        );
# Load raw u.item file
lines = spark.sparkContext.textFile("file:///tmp/u.item")
movies = lines.map(parseInput3)
moviesDataset = spark.createDataFrame(movies)
#write to cassandra
moviesDataset.write.format("org.apache.spark.sql.cassandra").mode("append").options(table="titles", keyspace="movielens").save()
#read to check
readMovies = spark.read.format("org.apache.spark.sql.cassandra").options(table="titles", keyspace="movielens").load()
```

```
imdb url|musical|mystery|release
|movie id|action|adventure|animation|children|comedy|crime|documentary|drama|fantasy|film noir|horror|
date|romance|sci fi|thriller|
                                              title|unknown|war|western|
                                                                                                         0|http://us.imdb.co...|
                                                                                                                                                 0| 11-0ct-
     1507
               01
                                                                                                                                        0 l
1996
           0|
                   01
                            0|Three Lives and 0...|
                                                           0 0
               0 l
                                                                                       01
                                                                                                  01
                                                                                                         0|http://us.imdb.co...|
                                                                                                                                        11
                                                                                                                                                 0 | 21-Jun-
      596
                                              11
1996
           0|
                   01
                            0|Hunchback of Notr...|
                                                           0 0
     1145
               01
                                                            01
                                                                         0 I
                                                                               1|
                                                                                       01
                                                                                                  0|
                                                                                                         0|http://us.imdb.co...|
                                                                                                                                        0 l
                                                                                                                                                 0 | 01-Jan-
1994
           01
                   01
                            01
                                 Blue Chips (1994)
                                                               01
     1609
               0|
                                     01
                                              0|
                                                     1|
                                                            01
                                                                                       0 l
                                                                                                  0 l
                                                                                                         0|http://us.imdb.co...|
                                                                                                                                        0 l
                                                                                                                                                 0 28-Mar-
                                     B*A*P*S (1997)
1997
           01
                   0 l
                            0|
                                                           0 0
                                                                                                  0 l
                                                                                                         0|http://us.imdb.co...|
     1365
               1|
                                              01
                                                            01
                                                                         0 I
                                                                               1|
                                                                                       0 l
                                                                                                                                        0 l
                                                                                                                                                 0 | 01-Jan-
1993
           01
                   01
                            0|Johnny 100 Pesos ...|
                                                           0 0
                                                                        0|
                                                                                                         0|http://us.imdb.co...|
                                                                                                                                                0| 04-Oct-
     1364
               11
                                     0 l
                                                      0 l
                                                            0 l
                                                                         0 l
                                                                               0 l
                                                                                       01
                                                                                                  0 l
                                                                                                                                        0 l
10061
                            al pind of prov (1006) |
```

```
%pyspark READY
```

#to allow DF view as SQL
usersDataset.createOrReplaceTempView("users")
ratingsDataset.createOrReplaceTempView("ratings")
moviesDataset.createOrReplaceTempView("movies")

4. Analysis

a. Average rating for each movie

```
%pyspark

from pyspark sql functions import avg
```

```
from pyspark.sql.functions import avg

joined_df = readRatings.join(readMovies, on="movie_id")
avg_ratings = joined_df.groupBy("title").agg(avg("rating").alias("avg_rating"))
avg_ratings.limit(20).show()
```

```
title
                              avg rating
          Cosi (1996)
        Psycho (1960) | 4.100418410041841
 Three Wishes (1995)|3.2222222222223|
 If Lucy Fell (1996) | 2.7586206896551726 |
|When We Were King...| 4.045454545454546|
   Annie Hall (1977) | 3.911111111111111 |
     Fair Game (1995) 2.1818181818181817
|Heavenly Creature...|3.6714285714285713|
|Paris, France (1993)|2.3333333333333335|
|Snow White and th...|3.7093023255813953|
|Night of the Livi...|
                                3.421875
|I'll Do Anything ...|
                                     2.6
|Spanking the Monk...| 3.074074074074074|
         Mondo (1996)
     Throngoma /1004\| 2 0207006774102EE
```

b. Top ten movies with the highest average ratings

```
%pyspark
                                                                                                                                    READY
top10 movies = avg ratings.orderBy("avg rating", ascending=False).limit(10)
top10 movies.show()
              title avg rating
+----+
|Someone Else's Am...|
                          5.0
|Saint of Fort Was...|
                          5.0
|Aiqing wansui (1994)|
                          5.0
|Marlene Dietrich:...|
                          5.0
     Star Kid (1997)
                          5.0
|Great Day in Harl...|
                          5.0
|Entertaining Ange...|
                          5.0
|They Made Me a Cr...|
                          5.0
|Santa with Muscle...|
                          5.0
 Prefontaine (1997)
                          5.0
```

+----+

c. Users who have rated at least 50 movies and what are their favourite movie genres

READY

```
%sql
WITH active users AS ( -- filter users who have rated at least 50 movies
  SELECT user id
  FROM ratings
  GROUP BY user id
  HAVING COUNT(*) >= 50
user genres AS ( -- joins ratings with movies to get genre for each movie
  SELECT r.user id, m.*
  FROM ratings r
  JOIN movies m ON r.movie id = m.movie id
  WHERE r.user id IN (SELECT user id FROM active users)
genre counts AS (
  SELECT user id, genre, COUNT(*) AS count
  FROM (
    SELECT user id, stack(19, -- to unpivot genre columns to row
      'unknown', unknown, 'Action', action, 'Adventure', adventure, 'Animation', animation,
      'Children', children, 'Comedy', comedy, 'Crime', crime, 'Documentary', documentary,
      'Drama', drama, 'Fantasy', fantasy, 'Film-Noir', film noir, 'Horror', horror,
      'Musical', musical, 'Mystery', mystery, 'Romance', romance, 'Sci-Fi', sci fi,
      'Thriller', thriller, 'War', war, 'Western', western
    ) AS (genre, is genre)
    FROM user genres
  WHERE is genre = 1 -- the genre is belong to the original genre which 1
  GROUP BY user id, genre
ranked genres AS ( -- to rank genres based on count
  SELECT *,
         ROW NUMBER() OVER (PARTITION BY user id ORDER BY count DESC) AS rank
  FROM genre counts
SELECT user id, genre, count
FROM ranked genres
WHERE rank = 1 -- top genre for that user
ORDER BY count DESC
```

 \sim

 \blacksquare

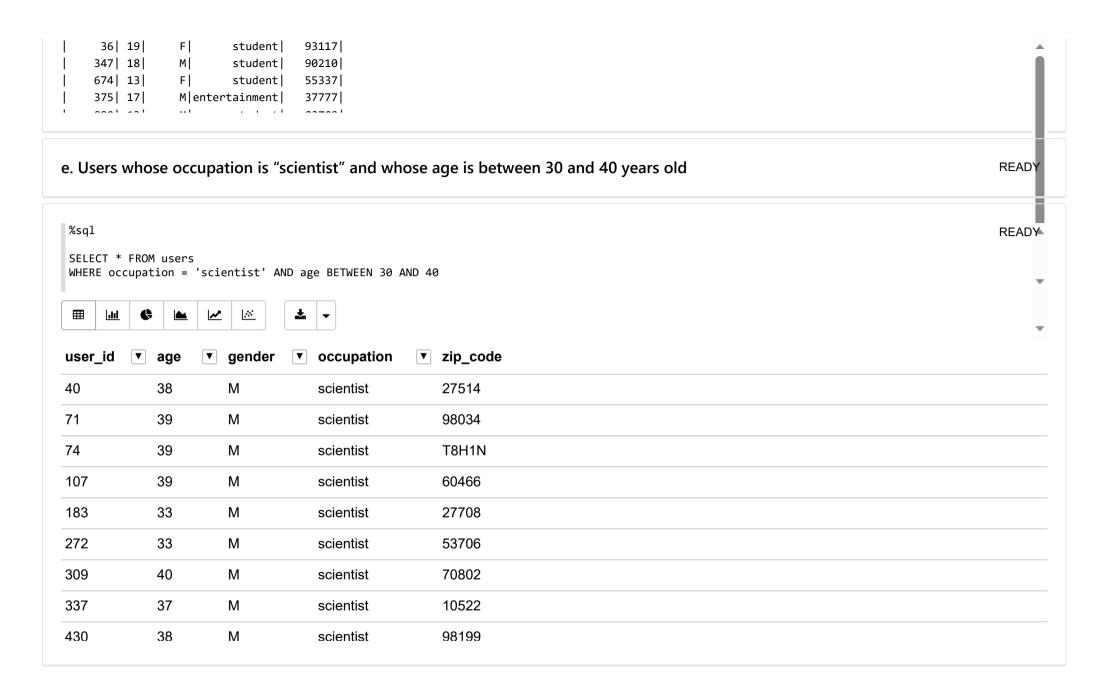
user_id	genre	▼ count
655	Drama	410
405	Drama	309
537	Drama	251
450	Drama	237
13	Drama	218
234	Drama	213
416	Drama	212
279	Comedy	211
201	Drama	196

d. Users who are less than 20 years old

%pyspark

READY

```
under20age = readUsers.filter("age < 20")</pre>
under20age.show()
|user_id|age|gender| occupation|zip_code|
    262 | 19 |
                                    78264
                 F|
                         student
    142 | 13 |
                 M|
                           other
                                    48118
                 F|
                                    47906
    223 | 19 |
                         student
                 M|
                                    94619
    289 11
                            none
    618 | 15 |
                 F|
                         student
                                    44212
    471 10
                 M|
                         student
                                    77459
    281 | 15 |
                 F|
                         student
                                    06059
    621 | 17 |
                 M
                         student
                                    60402
    887 | 14 |
                 F|
                         student|
                                    27249
    270 | 18 |
                 F|
                         student
                                    63119
```



%pyspark

READY

sc.version