

# University of Virginia

## DS 5559: Big Data Analytics

### Music Recommendation

**Last updated: Feb 29, 2020**

#### Instructions

In this assignment, you will work with a recommendation algorithm based on user listening data from Autoscrobber.

The code is outlined below. Make the requested modifications, run the code, and copy all answers to the **ANSWER SECTION** at the bottom of the notebook. Note the *None* variable is a placeholder for code.

NOTE: For a given userID, some/many recommendation might come back as *None*.

These should be filtered out using a list comprehension as follows:

```
In [1]: 1 #print([x for x in recommendationsForUser if x is not None])
```

**TOTAL POINTS: 10**

---

#### About the Alternating Least Squares Parameters

rank

The number of latent factors in the model, or equivalently, the number of columns  $k$  in the user-feature and product-feature matrices. In nontrivial cases, this is also their rank.

iterations

The number of iterations that the factorization runs. More iterations take more time but may produce a better factorization.

lambda

A standard overfitting parameter. Higher values resist overfitting, but values that are too high hurt the factorization's accuracy.

```
In [2]: 1 # import modules
        2 import os
        3
        4 from pyspark import SparkContext
        5 from pyspark import SparkConf
        6 from pyspark.mllib import recommendation
        7 from pyspark.mllib.recommendation import *
        8 import pandas as pd
        9
       10 from pyspark.mllib.recommendation import *
       11 import random
       12 from operator import *
```

```
In [3]: 1 from pyspark import SparkContext, SparkConf
        2 spark = SparkContext.getOrCreate()
        3 spark.stop()
        4 spark = SparkContext('local', 'Recommender')
```

```
In [4]: 1 # set configurations
        2 conf = SparkConf().setMaster("local").setAppName("autoscrobbler")
```

```
In [5]: 1 # set context
        2 sc = SparkContext.getOrCreate(conf=conf)
```

```
In [6]: 1 # pathing and params
        2 user_artist_data_file = 'user_artist_data.txt'
        3 artist_data_file = 'artist_data.txt'
        4 artist_alias_data_file = 'artist_alias.txt'
        5
        6 numPartitions = 2
        7 topk = 10
```

```
In [7]: 1 # read user_artist_data_file into RDD (417MB file, 24MM records of users' plays of artists, along with count)
        2 # specifically, each row holds: userID, artistID, count
        3 rawDataRDD = sc.textFile(user_artist_data_file, numPartitions)
        4 rawDataRDD.cache()
```

Out[7]: user\_artist\_data.txt MapPartitionsRDD[1] at textFile at NativeMethodAccessorImpl.java:0

```
In [8]: 1 # userid, artistid, playcount
        2 rawDataRDD.take(5)
```

Out[8]: ['1000002 1 55',  
'1000002 1000006 33',  
'1000002 1000007 8',  
'1000002 1000009 144',  
'1000002 1000010 314']

```
In [9]: 1 # read artist_data_file using *textFile*
        2 # Import test files from location into RDD variables
        3 # YOUR CODE GOES HERE
        4 #import os
        5 #os.getcwd()
        6 artistData = sc.textFile('artist_data.txt').map(lambda s:(int(s.split("\t")[0]),s.split("\t")[1]))
```

```
In [10]: 1 # inspect some records
         2 #artistid, artist_name
         3 artistData.take(5)
```

Out[10]: [(1134999, '06Crazy Life'),  
(6821360, 'Pang Nakarin'),  
(10113088, 'Terfel, Bartoli- Mozart: Don'),  
(10151459, 'The Flaming Sidebur'),  
(6826647, 'Bodenstandig 3000')]

```
In [11]: 1 # read artist_alias_data_file using *textFile*
         2 artist_alias= sc.textFile('artist_alias.txt')
```

```
In [12]: 1 # inspect some records
          2 # id, id
          3 artist_alias.take(5)
```

```
Out[12]: ['1092764\t1000311',
          '1095122\t1000557',
          '6708070\t1007267',
          '10088054\t1042317',
          '1195917\t1042317']
```

```
In [13]: 1 from pyspark.mllib.recommendation import *
          2 import random
          3 from operator import *
          4 def parser(s, delimiters=" ", to_int=None):
          5     s = s.split(delimiters)
          6     if to_int:
          7         return tuple([int(s[i]) if i in to_int else s[i] for i in range(len(s))])
          8     return tuple(s)
          9 artistData = sc.textFile("artist_data.txt").map(lambda x: parser(x, '\t', [0]))
         10 artistAlias = sc.textFile("artist_alias.txt").map(lambda x: parser(x, '\t', [0,1]))
         11 userArtistData = sc.textFile("user_artist_data.txt").map(lambda x: parser(x, ' ', [0,1,2]))
```

```
In [14]: 1 # 1) (1 PT) Print the first 10 records from rawDataRDD
          2 rawDataRDD
          3 rawDataRDD.top(topk)
```

```
Out[14]: ['9875 9973009 2',
          '9875 979 41',
          '9875 976 3',
          '9875 949 29',
          '9875 930 1',
          '9875 929 1',
          '9875 92 1',
          '9875 910 1',
          '9875 891 32',
          '9875 868 12']
```

```
In [15]: 1 def parseArtistIdNamePair(singlePair):
2         splitPair = singlePair.rsplit('\t')
3         # we should have two items in the list - id and name of the artist.
4         if len(splitPair) != 2:
5             #print singlePair
6             return []
7         else:
8             try:
9                 return [(int(splitPair[0]), splitPair[1])]
10            except:
11                return []
```

```
In [16]: 1 rawArtistRDD = sc.textFile(artist_data_file)
```

```
In [17]: 1 artistByID = dict(rawArtistRDD.flatMap(lambda x: parseArtistIdNamePair(x)).collect())
```

```
In [22]: 1 # 2) (1 PT) Print 10 values from artistByID, using topk variable
2         from collections import Counter
3         import collections
4         topk = 10
5         # Hint: the most_common() function may help
```

```
In [23]: 1 c = Counter(artist_vals)
2         c.most_common(topk)
```

```
Out[23]: [('06Crazy Life', 1),
('Pang Nakarin', 1),
('Terfel, Bartoli- Mozart: Don', 1),
('The Flaming Sidebur', 1),
('Bodenständig 3000', 1),
('Jota Quest e Ivete Sangalo', 1),
('Toto_XX (1977', 1),
('U.S Bombs -', 1),
('artist formally know as Mat', 1),
('Kassierer - Musik für beide Ohren', 1)]
```

```
In [24]: 1 def parseArtistAlias(alias):
2         splitPair = alias.rsplit('\t')
3         # we should have two ids in the list.
4         if len(splitPair) != 2:
5             #print singlePair
6             return []
7         else:
8             try:
9                 return [(int(splitPair[0]), int(splitPair[1]))]
10            except:
11                return []
```

```
In [25]: 1 rawAliasRDD = sc.textFile(artist_alias_data_file)
```

```
In [26]: 1 artistAlias = rawAliasRDD.flatMap(lambda x: parseArtistAlias(x)).collectAsMap()
```

```
In [27]: 1 # Create a dictionary of artist id's
2         # artist
```

```
In [28]: 1 # turn the artistAlias into a broadcast variable.
2         # This will distribute it to worker nodes efficiently, so we save bandwidth.
3         artistAliasBroadcast = sc.broadcast( artistAlias )
```

```
In [29]: 1 artistAliasBroadcast.value.get(2097174)
```

Out[29]: 1007797

```
In [30]: 1 # Print the number of records from the Largest RDD, rawDataRDD
2         print( rawDataRDD.count() )
```

24296858

```
In [31]: 1 # Sample 10% of rawDataRDD using seed 314, to reduce runtime. Call it sample.
2 seed = 314
3 weights = [.9, .1]
4 sample, _ = rawDataRDD.randomSplit(weights, seed)
5 sample.cache()
```

Out[31]: PythonRDD[26] at RDD at PythonRDD.scala:53

```
In [32]: 1 # take the first 5 records from the sample. each row represents userID, artistID, count.
2 sample.take(5)
```

Out[32]: ['1000002 1 55',  
'1000002 1000006 33',  
'1000002 1000007 8',  
'1000002 1000009 144',  
'1000002 1000010 314']

```
In [33]: 1 artistByIDBroadcast = sc.broadcast(artistByID)
```

```
In [34]: 1 # Based on sampled data, build the matrix for model training
2 def mapSingleObservation(x):
3     # Returns Rating object represented as (user, product, rating) tuple.
4     # [add line of code here to split each record into userID, artistID, count]
5     userID, artistID, count = map(lambda lineItem: int(lineItem), x.split())
6     # given possible aliasing, get finalArtistID
7     finalArtistID = artistAliasBroadcast.value.get(artistID)
8     if finalArtistID is None:
9         finalArtistID = artistID
10    return Rating(userID, finalArtistID, count)
```

```
In [35]: 1 trainData = sample.map(lambda x: mapSingleObservation(x))
2 trainData.cache()
```

Out[35]: PythonRDD[28] at RDD at PythonRDD.scala:53

```
In [36]: 1 # 3) (1 PT) Print the first 5 records from trainData
        2 trainData.take(5)
```

```
Out[36]: [Rating(user=1000002, product=1, rating=55.0),
          Rating(user=1000002, product=1000006, rating=33.0),
          Rating(user=1000002, product=1000007, rating=8.0),
          Rating(user=1000002, product=1000009, rating=144.0),
          Rating(user=1000002, product=1000010, rating=314.0)]
```

```
In [41]: 1 # Train the ALS model, using seed 314, rank 10, iterations 5, lambda_ 0.01. Call it model.
        2 from pyspark.mllib.recommendation import *
        3 model = ALS.trainImplicit(trainData, rank=10, iterations = 5, alpha = 0.01)
```

```
In [43]: 1 # Model Evaluation
        2
        3 # fetch artists for a test user
        4 testUserID = 1000002
        5
        6 # broadcast artistByID for speed
        7 artistByIDBroadcast = sc.broadcast( artistByID )
        8
        9 # from trainData, collect the artists for the test user. Call the object artistsForUser.
       10 # hint: you will need to apply .value.get(x.product) to the broadcast artistByID, where x is the Rating RDD.
       11 # if you don't do this, you may see artistIDs. you want artist names.
       12 artistsForUser = (trainData
       13                   .filter(lambda observation: observation.user == testUserID)
       14                   .map(lambda observation: artistByIDBroadcast.value.get(observation.product))
       15                   .collect())
```

```
In [48]: 1 res = [i for i in artistsForUser if i]
        2 print(res)
```

```
['Mallrats', 'Kerrang', 'Brian Hughes', 'Joshua Redman', 'The Mystick Krewe of Clearlight', 'Benny Goodman Orchestra',
'YMC', 'Brant Bjork and The Operators', 'Firebird', 'Elvis Costello', 'Café Del Mar', 'Eric Clapton', 'Enigma', 'Eurythm
ics', 'Armand Van Helden', 'Echo & the Bunnymen', 'George Duke']
```



```
In [49]: 1 # 4) (1 PT) Print the artist listens for testUserID = 1000002
        2 c = Counter(artist_vals)
        3 c.most_common(topk)
```

```
Out[49]: [('06Crazy Life', 1),
          ('Pang Nakarin', 1),
          ('Terfel, Bartoli- Mozart: Don', 1),
          ('The Flaming Sidebur', 1),
          ('Bodenstandig 3000', 1),
          ('Jota Quest e Ivete Sangalo', 1),
          ('Toto_XX (1977', 1),
          ('U.S Bombs -', 1),
          ('artist formaly know as Mat', 1),
          ('Kassierer - Musik für beide Ohren', 1)]
```

```
In [59]: 1 # 5) (2 PTS) Make 10 recommendations for testUserID = 1000002
        2 num_recomm = 500
        3 recommendationsForUser_rank10 = map(lambda observation: artistByID.get(observation.product), model.call("recommendPro
        4 print([x for x in recommendationsForUser_rank10 if x is not None])
```

```
['Eric Clapton', 'Elvis Costello', 'Eurythmics', 'Scorpions', 'Enigma', 'Gary Jules', '植松伸夫', 'Nena', 'Joss Stone']
```

```
In [56]: 1 # Train a second ALS model with rank 20, iterations 5, Lambda 0.01.
        2 model_2 = ALS.trainImplicit(trainData, rank= 20, iterations = 5, alpha = 0.01)
```

```
In [60]: 1 # 6) (2 PTS) Using the rank 20 model, make 10 recommendations for the same test user
        2 num_recomm = 500
        3 recommendationsForUser_rank20 = map(lambda observation: artistByID.get(observation.product), model_2.call("recommendP
        4 print([x for x in recommendationsForUser_rank20 if x is not None])
```

```
['Eric Clapton', 'Eurythmics', 'Scorpions', 'Elvis Costello', 'Enigma', 'Gary Jules', 'Nena', 'Joss Stone']
```

**ANSWER SECTION (COPY ALL ANSWERS HERE)**

```
In [220]: 1 # ANSWER 1 (1 PT)
          2 # Print the first 10 records from rawDataRDD
          3 rawDataRDD.top(topk)
```

```
Out[220]: ['9875 9973009 2',
           '9875 979 41',
           '9875 976 3',
           '9875 949 29',
           '9875 930 1',
           '9875 929 1',
           '9875 92 1',
           '9875 910 1',
           '9875 891 32',
           '9875 868 12']
```

```
In [86]: 1 # ANSWER 2 (1 PT)
          2 # Print topk values from artistByID
          3 c = Counter(artist_vals)
          4 c.most_common(topk)
```

```
Out[86]: [('06Crazy Life', 1),
          ('Pang Nakarin', 1),
          ('Terfel, Bartoli- Mozart: Don', 1),
          ('The Flaming Sidebur', 1),
          ('Bodenstandig 3000', 1),
          ('Jota Quest e Ivete Sangalo', 1),
          ('Toto_XX (1977', 1),
          ('U.S Bombs -', 1),
          ('artist formaly know as Mat', 1),
          ('Kassierer - Musik für beide Ohren', 1)]
```

```
In [62]: 1 # ANSWER 3 (1 PT)
          2 # Print the first 5 records from trainData
          3 trainData.take(5)
```

```
Out[62]: [Rating(user=1000002, product=1, rating=55.0),
          Rating(user=1000002, product=1000006, rating=33.0),
          Rating(user=1000002, product=1000007, rating=8.0),
          Rating(user=1000002, product=1000009, rating=144.0),
          Rating(user=1000002, product=1000010, rating=314.0)]
```

```
In [219]: 1 # ANSWER 4 (1 PT)
          2 # Print the artist listens for testUserID = 1000002
          3 res = [i for i in artistsForUser if i]
          4 print(res)
```

['Mallrats', 'Kerrang', 'Brian Hughes', 'Joshua Redman', 'The Mystick Krewe of Clearlight', 'Benny Goodman Orchestra', 'YMC', 'Brant Bjork and The Operators', 'Firebird', 'Elvis Costello', 'Café Del Mar', 'Eric Clapton', 'Enigma', 'Eurythmics', 'Armand Van Helden', 'Echo & the Bunnymen', 'George Duke']

```
In [63]: 1 # ANSWER 5 (2 PTS)
          2 # Make 10 recommendations for testUserID = 1000002
          3 num_recomm = 500
          4 recommendationsForUser_rank10 = map(lambda observation: artistByID.get(observation.product), model.call("recommendPro
          5 print([x for x in recommendationsForUser_rank10 if x is not None])
```

['Eric Clapton', 'Elvis Costello', 'Eurythmics', 'Scorpions', 'Enigma', 'Gary Jules', '植松伸夫', 'Nena', 'Joss Stone']

```
In [64]: 1 # ANSWER 6 (2 PTS)
          2 # Using the rank 20 model, make 10 recommendations for testUserID = 1000002
          3 num_recomm = 500
          4 recommendationsForUser_rank20 = map(lambda observation: artistByID.get(observation.product), model_2.call("recommendP
          5 print([x for x in recommendationsForUser_rank20 if x is not None])
```

['Eric Clapton', 'Eurythmics', 'Scorpions', 'Elvis Costello', 'Enigma', 'Gary Jules', 'Nena', 'Joss Stone']

```
In [ ]: 1 # ANSWER 7 (2 PTS)
          2 # How does the rank 10 model seem to perform versus the rank 20 model?
          3 # The contents of artistsForUser may help answer the question.
```

```
In [70]: 1 list1 = ['Mallrats', 'Kerrang', 'Brian Hughes', 'Joshua Redman', 'The Mystick Krewe of Clearlight', 'Benny Goodman Or  
2 list2_r10 = ['Eric Clapton', 'Elvis Costello', 'Eurythmics', 'Scorpions', 'Enigma', 'Gary Jules', '植松伸夫', 'Nena',  
3 list3_r20 = ['Eric Clapton', 'Eurythmics', 'Scorpions', 'Elvis Costello', 'Enigma', 'Gary Jules', 'Nena', 'Joss Stone  
4 comment_elements_r10 = []  
5 comment_elements_r20 = []  
6 set1=set(list1)  
7 common_elements_r10= set1.intersection(list2_r10)  
8 common_elements_r20 = set1.intersection(list3_r20)
```

```
In [71]: 1 print(common_elements_r10)  
  
{'Eric Clapton', 'Enigma', 'Eurythmics', 'Elvis Costello'}
```

```
In [73]: 1 print(common_elements_r20)  
  
{'Eric Clapton', 'Enigma', 'Eurythmics', 'Elvis Costello'}
```

- 1 The rank controls the number of internal parameters that must be fit from the data, too many and you get overfitting your training set.
- 2 Since we might not know the underlying factor. The more you use, the better the results up to a point, but the more memory and computation time you will need. I compare the rank 10 to artistForusers to see the common element. I also compare the rank 20 to artistForusers to see the common element. I found the common users are {'Eric Clapton', 'Enigma', 'Eurythmics', 'Elvis Costello'}. Again, since we might need to guess to see the underlying factors, the chosen rank higher should be better since the ranks refers to the presumed latent or hidden factors. However, we also need to avoid the overfitting issue.