

Term Project Report Group63

Project

1. Recommend System
2. dataSet: comes from <https://grouplens.org/datasets/movielens/>
(<https://grouplens.org/datasets/movielens/>)
ml-latest-small.zip
3. output: (userId, movieId) predicate rate

implementation

1. 讀取資料，並且使用movies計算總共有幾個movie在計算最後的結果時才知道有幾個movies

```
1  movies = set()
2  def readData(inputFile):
3      global movies
4      with open(inputFile) as fp:
5          data = fp.readlines()
6
7      data = data[1:]
8      for i in range(len(data)):
9          data[i] = data[i].strip().split(',')[0:3]
10         data[i][0] = int(data[i][0])
11         data[i][1] = int(data[i][1])
12         movies.add(data[i][1])
13         data[i][2] = float(data[i][2])
14     return data
15
16 data = readData('ml-latest-small/ratings.csv')
```

2. Systme setting

```
1  import pyspark
2  from pyspark import SparkContext, SparkConf
3  import math
4  import os
5
6
7  conf = SparkConf().setMaster("local[*]").setAppName("Project").set("spark.
8  sc = SparkContext(conf=conf).getOrCreate()
```

3. 需要幾個最相似的item來計算推薦值

```
1 | topK = 10
```

4. input data是[userId, movieId, score]

透過第一個map去掉userId,只剩下movieId,再將相同的movieId聚集起來
再透過mapValue取的movies的平均分數

```
1 | inputMat = sc.parallelize(data)
2 | point_for_each_item = inputMat.map(lambda s:(s[1], s[2])).groupByKey() # type (movieId, Iterable(score))
3 | avg_point_for_each_item = point_for_each_item.mapValues(lambda s:sum(s)/len(s))
```

5. 將input 轉換為 (movieId, (userId, score))的形式

再將上述結果與前面計算的movieId平均值相減
就可以得到扣除平均值的結果

```
1 | point_minus_avg = inputMat.map(lambda s:(s[1], (s[0], s[2]))) # type (movieId, (userId, score))
2 | point_minus_avg = point_minus_avg.join(avg_point_for_each_item) # type (movieId, ((userId, score), avg_score))
3 | point_minus_avg = point_minus_avg.map(lambda s:((s[0], s[1][0][0]), s[1][0][1] - s[1][1]))
```

6. 計算cosim similarity的方法

首先將item1與item2有的value平方相加開根號
再根據暴力法找看過相同的movie並且相乘在相加

```
1 | # calculate cosine similarity
2 | def cosim(_input):
3 |     item1_L2_nor = 0
4 |     for ele in list(_input[0][1]):
5 |         item1_L2_nor += (ele[1]**2)
6 |     item1_L2_nor = item1_L2_nor ** 0.5
7 |
8 |     item2_L2_nor = 0
9 |     for ele in list(_input[1][1]):
10 |         item2_L2_nor += (ele[1]**2)
11 |     item2_L2_nor = item2_L2_nor ** 0.5
12 |
13 |     rxi_ryi = 0
14 |     for ele1 in list(_input[0][1]):
15 |         for ele2 in list(_input[1][1]):
16 |             if(ele1[0] == ele2[0]):
17 |                 rxi_ryi += (ele1[1] * ele2[1])
18 |             break
19 |     if item1_L2_nor * item2_L2_nor != 0:
20 |         return ((_input[0][0], _input[1][0]), rxi_ryi/(item1_L2_nor * item2_L2_nor))
21 |     else:
22 |         return ((_input[0][0], _input[1][0]), 0.0)
```

7. point_minus_avg.map將資料轉換為(movieId, (userId, score))

透過item_score.cartesian(item_score)得到所有的pair

利用filter刪除 (2, 1)這種case，防止(2, 1),(1, 2)重複計算

all_pair_for_item.map(cosim)得到所有的movie pair的similarity

```
1 item_score = point_minus_avg.map(lambda s:(s[0][0], (s[0][1],s[1]))).group
2 all_pair_for_item = item_score.cartesian(item_score).filter(lambda s: s[0]
3 all_pair_score_for_item = all_pair_for_item.map(cosim) # type((movieId1, n
```

8. all_pair_score_for_item.map(...)將所有資料轉為(similarity, ((movie1, movie2), (movie2, movie1)))

透過all_pair_score_for_item.flatMapValues 將資料轉為(similarity, (movie1, movie2))，如此一來也能得到(2, 1)的similarity

最後在map成 ((movie1, movie2), score)的形式

```
1 # type (score, ((movie1, movie2), (movie2, movie1)))
2 all_pair_score_for_item = all_pair_score_for_item.map(lambda s:(s[1], ((s|
3 # type ((movie1, movie2), score)
4 all_pair_score_for_item = all_pair_score_for_item.flatMapValues(lambda s:
```

9. 將movie依照similarity排序

all_pair_score_for_item.map將item轉為(movieId, sorted similarity)

sorted similarity是list每個element為(item, score)並且根據score的大小進行排序

ele當中為(movie, similarity)

```
1 def sort_similarity(_input):
2     ele = list(_input[1])
3     ele = sorted(ele, key=lambda s:s[1], reverse=True)
4     return (_input[0], ele)
5
6
7 similarity_score_for_each_item = all_pair_score_for_item.map(lambda s:((s|
```

10. 得到user沒看過的movieId

user對movie的pred

```

1  def get_unRating(_input):
2      global movies
3      movie_list = list(movies)
4      for ele in list(_input):
5          movie_list.pop(movie_list.index(ele[0]))
6
7      return movie_list
8
9  def get_rate(_input):
10     label = dict()
11     for ele in list(_input[1][0][1]): ## user have seen movie
12         label[ele[0]] = ele[1]
13     count = 0
14     pred = 0
15     weight_list = []
16     for j in range(len(_input[1][1])):
17         if _input[1][1][j][0] in label.keys() and _input[1][1][j][1]>0:
18             pred += (_input[1][1][j][1] * label[_input[1][1][j][0]])
19             weight_list.append(_input[1][1][j][1])
20             count += 1
21             if(count == topK):
22                 break
23     weight_sum = sum(weight_list)
24     if weight_sum == 0:
25         return ((_input[1][0][0], _input[0]), 0)
26     else:
27         return ((_input[1][0][0], _input[0]), pred/sum(weight_list))

```

11. `inputMat.map(...)`得到user看過的所有電影 type為`type(user ,list of (movie, score))`
 透過`get_unRating`得到user沒看有過的movies，透過`flatMapValues`將type轉為 (user, unseen movie)
`user_unrating.join(user_rating)`得到user沒看過move和看過的movie list與rating
`user_rating_and_unrating.map(...)`將unseen movie當成key
`user_rating_and_unrating.join(...)` 得到了unseen movie的similarity score
 最後透過`total_info.map(get_rate)`得到((userId, movieId), pred)

```

1  user_rating = inputMat.map(lambda s:(s[0],(s[1], s[2]))).groupByKey() # ty
2  user_unseen = user_rating.mapValues(get_unRating).flatMapValues(lambda s:s
3  user_rating_and_unrating = user_unseen.join(user_rating) # type(userId, (u
4  user_rating_and_unrating = user_rating_and_unrating.map(lambda s:(s[1][0],
5  total_info = user_rating_and_unrating.join(similarity_score_for_each_item)
6  user_unseen_rate = total_info.map(get_rate) # type(userId, (movieId, pred)
7  res = user_unseen_rate.sortByKey().collect()

```

12. 將檔案寫成file

```
1 | with open("Outputfile.txt", 'w') as fp:
2 |     for i in range(len(res)):
3 |         fp.write("{0},{1},{2}\n".format(res[i][0][0], res[i][0][1], res[i][0][2]))
```