

## Jiayue Shi's Description

### 1. Version

Spark version: 2.2.1

Scala version: 2.11

Hadoop version: 2.7

### 2. How to run

**First of all, the input files should under the data folder, and output file will also be under data folder.**

Task1\_1

```
./bin/spark-submit --class JaccardLSH Jiayue_Shi_hw3.jar data/video_small_num.csv  
data/Jiayue_Shi_SimilarProducts_Jaccard.txt
```

Task2\_1

```
./bin/spark-submit --class UserBasedCF Jiayue_Shi_hw3.jar data/video_small_num.csv  
data/video_small_testing_num.csv data/Jiayue_Shi_UserBasedCF.txt
```

Task2\_2

```
./bin/spark-submit --class ModelBasedCF Jiayue_Shi_hw3.jar data/video_small_num.csv  
data/video_small_testing_num.csv data/Jiayue_Shi_ModelBasedCF.txt
```

Task2\_3

```
./bin/spark-submit --class ItemBasedCF --driver-memory 2g Jiayue_Shi_hw3.jar  
data/video_small_num.csv data/video_small_testing_num.csv  
data/Jiayue_Shi_SimilarProducts_Jaccard.txt data/Jiayue_Shi_ItemBasedCF.txt
```

### 3. Task1 Result

Task1	Jaccard	Cosine
Precision	1	0.9410651799018869
Recall	0.9235976599949131	0.7374794987863281
Time	44s	25min

### 4. Task2 Result

Task2	ModelBased	UserBased	ItemBased - USE LSH
-------	------------	-----------	---------------------

Accuracy	$\geq 0$ and $< 1$ : 4601 $\geq 1$ and $< 2$ : 2194 $\geq 2$ and $< 3$ : 568 $\geq 3$ and $< 4$ : 329 $\geq 4$ : 8 RMSE: 1.2954668	$\geq 0$ and $< 1$ : 4452 $\geq 1$ and $< 2$ : 2308 $\geq 2$ and $< 3$ : 650 $\geq 3$ and $< 4$ : 286 $\geq 4$ : 4 RMSE: 1.320513	$\geq 0$ and $< 1$ : 4117 $\geq 1$ and $< 2$ : 2310 $\geq 2$ and $< 3$ : 810 $\geq 3$ and $< 4$ : 400 $\geq 4$ : 63 RMSE: 1.4189949
Time	10s	25s	17s

**Compare:** If we use LSH, we will use less time to get our result ( $17s < 25s$ ), because after LSH's filter, we will use less movies to calculate their pearson similarity. But we will have larger RMSE ( $1.41 > 1.32$ ), because LSH may filter some valuable item pairs, which may make influence on result.

## 5. Improvement

In UserBased task, there are lots of ratings which algorithm cannot predict. So, I calculated RMSE of use certain user's average rating or certain movie's average rating, Then I find the previous RMSE is smaller than the later. But when I try the average of both to fill missing value, I find the average of both of them will get lower RMSE. Therefore, I use  $(0.6 * \text{certain user's average rating} + 0.4 * \text{certain movie's average rating})$  to fill missing value, and this one get best result.