# 福建师范大学
## FUJIAN NORMAL UNIVERSITY

## 计算机与网络空间安全学院学生实验报告

实验课程名称： 大数据导论 教师： 林鑫泓

| 实验名称 | 大数据软件引用与技术实践 | | | 实验成绩 | |
|---|---|---|---|---|---|
| 学生姓名 | 叶建行 | 学 号 | 116052020005 | 年级专业班级 | 2020 级软件工程一班 |
| 小组成员 | 无 | | | 实验日期 | 2022 年 12 月 17 |

## 一、实验要求

1. 结合 Netflix Dataset 或 ml-25m，按照实验步骤完成实验。

2. Netflix 数据集包含了 1999.12.31 至 2005.12.31 期间由网站用户提供的超过一亿条电影评价。Netflix Dataset.7z 压缩文件包含电影信息、training set（训练集）、probe set（探测集）和 qualifying set(评估集)。压缩文件的详细信息如下图 1 所示：
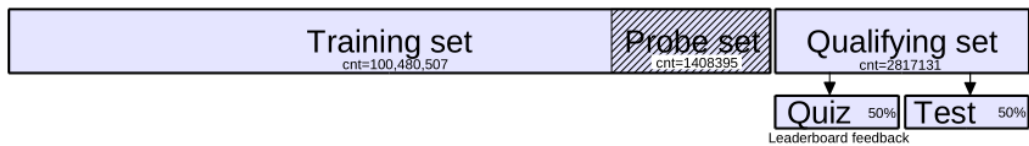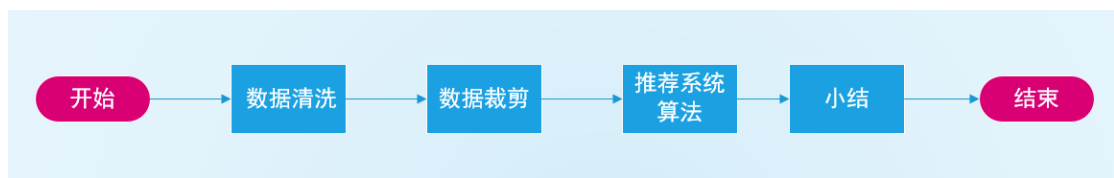


图 1 Netflix Dataset 说明

3. 实验成绩按照各步骤完成情况给分。

4. 大数据软件需包括教材当中有的，可以是介绍过的也可以是没介绍过的，尽可能地使用本学期学习并实验过的大数据软件和方法。

## 二、实验步骤

1. 独立设计一套切实可行的数据挖掘任务 task，并简要介绍该任务的各个步骤 step，使用的关键技术与软件应用。（10 分）

**任务流程图：**



数据清洗：通过 hive 实现，清除掉看的人少的，差评多的电影，全给好评或差评的用户。

数据裁剪：通过 hivesql 实现，随机抽样。

推荐系统算法：通过使用 mahout 推荐算法实现。

小结：对得到的结果进行分析和总结。

2. 详细介绍数据准备过程，例如：如何对数据进行预处理，如何做数据持久化存储。（10 分）

**（1） 导入数据（hive 导表过程过程相似，只展示 ratings 表的导表过程）**

（由于之前没有保存结果截图，因此只是语句的截图，希望老师不要在意）

● **建表：**

```
hive> create table ratings
    > (userId int,movieId int,rating double,`timestamp` BIGINT)
    > ROW FORMAT DELIMITED
    > fields terminated by ','
    > STORED AS TEXTFILE
    > tblproperties("skip.header.line.count"="1");_
```

● **导入本地数据：**

```
hive> load data local inpath '/root/big_data/ml-25m/ratings.csv' into table ratings;_
```

**（2） ratings 表筛选了评分人数<=1000 的电影**

**执行语句：**

```
insert overwrite table ratings select * from ratings where movieId in (select movieId from ratings group by movieId having count(*)>1000);
```

执行结果：



**（3）ratings 表筛选了电影评分的平均值<=3 的电影**

执行语句：

insert overwrite table ratings select * from ratings where movieId in (select movieId from ratings group by movieId having sum(rating)/count(*)>3);

执行结果：



**（4）ratings 表筛选了给差评多的用户的评价（评价平均分<=3）**

执行语句：

insert overwrite table ratings select * from ratings where userId in (select userId from ratings group by userId having sum(rating)/count(*)>3);

执行结果：



**（5）ratings 表筛选了全给好评的用户的评价（评价平均分==5）**

执行语句：

insert overwrite table ratings select * from ratings where userId in (select userId from ratings group by userId having sum(rating)/count(*)>3 and sum(rating)/count(*)<5);

执行结果：



**（6） tags 表筛选了评分人数<=1000 的电影**

执行语句：

insert overwrite table tags select * from tags where movieId in (select movieId from ratings group by movieId having count(*)>1000);

执行结果：



（7）**tags** 表筛选了电影评分的平均值<=3 的电影

执行语句：

insert overwrite table tags select * from tags where movieId in (select movieId from ratings group by movieId having sum(rating)/count(*)>3);

执行结果：



（8）**tags** 表筛选了给差评多的用户的评价（评价平均分<=3）

执行语句：

insert overwrite table tags select * from tags where userId in (select userId from ratings group by userId having sum(rating)/count(*)>3);

执行结果：



**（9）tags 表筛选了全给好评的用户的评价（评价平均分==5）**

执行语句：

insert overwrite table tags select * from tags where userId in (select userId from ratings group by userId having sum(rating)/count(*)>3 and sum(rating)/count(*)<5);

执行结果：

**（10）数据裁剪（数据随机排序，然后截取数据）**

（共有 18000000 条数据，没截取算法跑了很多次，内存调到 15G 都崩了）

（截取一半的数据）

```
hive> create table new_ratings as select * from ratings order by rand() limit 9000000;
Query ID = root_20221216173919_cd457ac3-cad1-4d4c-a322-48f3350e1fc2
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-16 17:39:30,393 Stage-1 map = 0%,   reduce = 0%
2022-12-16 17:40:05,831 Stage-1 map = 17%,  reduce = 0%
2022-12-16 17:40:48,223 Stage-1 map = 35%,  reduce = 0%
2022-12-16 17:40:54,322 Stage-1 map = 45%,  reduce = 0%
2022-12-16 17:40:58,376 Stage-1 map = 100%,  reduce = 0%
2022-12-16 17:41:04,456 Stage-1 map = 50%,  reduce = 0%
2022-12-16 17:41:34,735 Stage-1 map = 67%,  reduce = 0%
2022-12-16 17:42:11,052 Stage-1 map = 83%,  reduce = 0%
2022-12-16 17:42:23,150 Stage-1 map = 86%,  reduce = 0%
2022-12-16 17:42:29,181 Stage-1 map = 91%,  reduce = 0%
2022-12-16 17:42:34,262 Stage-1 map = 96%,  reduce = 0%
2022-12-16 17:42:40,303 Stage-1 map = 100%,  reduce = 0%
2022-12-16 17:42:53,435 Stage-1 map = 100%,  reduce = 17%
2022-12-16 17:42:59,568 Stage-1 map = 100%,  reduce = 67%
2022-12-16 17:43:04,622 Stage-1 map = 100%,  reduce = 68%
2022-12-16 17:43:10,682 Stage-1 map = 100%,  reduce = 69%
```

(还是崩溃，降到了 5000000 条数据)

```
hive> create table new_ratings1 as select * from new_ratings limit 5000000;
Query ID = root_20221216204925_102172e0-43fe-4ba1-941e-a38a6dcfc7df
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-16 20:49:33,657 Stage-1 map = 0%,   reduce = 0%
2022-12-16 20:49:50,901 Stage-1 map = 97%,  reduce = 0%
2022-12-16 20:49:52,003 Stage-1 map = 100%,  reduce = 0%
2022-12-16 20:50:04,173 Stage-1 map = 100%,  reduce = 67%
2022-12-16 20:50:05,189 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local1752748001_0001
Moving data to directory hdfs://localhost:9000/user/root/warehouse/new_ratings1
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 260620288 HDFS Write: 130217603 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 40.938 seconds
hive> _
```

(还是崩溃，降到了 4000000 条数据)

```
hive> create table new_ratings2 as select * from new_ratings limit 4000000;
Query ID = root_20221216211526_1db4f067-33fd-482f-947e-c5a54101218c
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-16 21:15:34,824 Stage-1 map = 0%,  reduce = 0%
2022-12-16 21:15:53,080 Stage-1 map = 100%,  reduce = 0%
2022-12-16 21:16:05,163 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local1146012397_0001
Moving data to directory hdfs://localhost:9000/user/root/warehouse/new_ratings2
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 208560128 HDFS Write: 104174306 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 40.136 seconds
```

(还是崩溃，降到了 2000000 条数据)

```
hive> create table new_ratings3 as select * from new_ratings limit 2000000;
Query ID = root_20221216213955_f336c5d9-ae13-42b6-a58f-092836102302
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-16 21:40:05,030 Stage-1 map = 0%,  reduce = 0%
2022-12-16 21:40:15,151 Stage-1 map = 100%,  reduce = 0%
2022-12-16 21:40:23,290 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local1955617406_0001
Moving data to directory hdfs://localhost:9000/user/root/warehouse/new_ratings3
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 104390656 HDFS Write: 52087525 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 29.315 seconds
hive> _
```

**（10）数据持久化存储（将所有的表导出存在本地,导出过程相似，只展示导出 movies 表的过程）**

```
hive> insert overwrite local directory '/root/downloads/ml-25m/movies'
    > row format delimited fields terminated by '\t'
    > COLLECTION ITEMS TERMINATED BY '\n'
    > select * from movies;
Query ID = root_20221215235355_70c6da9f-9bdf-451d-a459-1a12353d0be7
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Job running in-process (local Hadoop)
2022-12-15 23:53:57,598 Stage-1 map = 100%,  reduce = 0%
Ended Job = job_local1553092264_0007
Moving data to local directory /root/downloads/ml-25m/movies
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 456443138 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 2.479 seconds
hive> _
```

3. 详细介绍数据统计分析的过程，例如：统计、分析电影评价的数量、用户的数量、电影的数量等，统计、分析电影评价的平均值、方差等。（10 分）

（1）电影评价的数量：

```
hive> select count(*) from new_ratings3;
OK
2000000
Time taken: 9.141 seconds, Fetched: 1 row(s)
hive> _
```

（2）用户的数量：

```
hive> select count(t.userId) from (select userId from new_ratings3 group by userId) as t;
Query ID = root_20221217113506_2a02168d-ded6-48c9-95ce-096050628310
Total jobs = 2
Launching Job 1 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:35:15,960 Stage-1 map = 0%,  reduce = 0%
2022-12-17 11:35:27,360 Stage-1 map = 67%,  reduce = 0%
2022-12-17 11:35:33,664 Stage-1 map = 100%,  reduce = 0%
2022-12-17 11:35:42,325 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local154308396_0001
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:35:44,266 Stage-2 map = 100%,  reduce = 100%
Ended Job = job_local1457552082_0002
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 104174876 HDFS Write: 0 SUCCESS
Stage-Stage-2:  HDFS Read: 104174876 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
150526
Time taken: 37.325 seconds, Fetched: 1 row(s)
hive> _
```

（3）电影的数量：

```
hive> select count(t.movieId) from (select movieId from new_ratings3 group by movieId) as t;
Query ID = root_20221217113752_02808c8b-e5ea-482a-9b6e-eb9169fc8239
Total jobs = 2
Launching Job 1 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:37:54,712 Stage-1 map = 0%,  reduce = 0%
2022-12-17 11:37:56,725 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local203561829_0003
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:37:58,472 Stage-2 map = 100%,  reduce = 100%
Ended Job = job_local973204117_0004
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 208349752 HDFS Write: 0 SUCCESS
Stage-Stage-2:  HDFS Read: 208349752 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
2923
Time taken: 6.478 seconds, Fetched: 1 row(s)
hive> _
```

（4）电影评价的平均值

- 电影评价的每部电影的评分平均值（前 5 个）：

```
hive> select movieId,avg(rating) from new_ratings3 group by movieId limit 5;
Query ID = root_20221217114032_218f5537-a647-4953-b3ff-95ad2f6265e2
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:40:34,636 Stage-1 map = 0%,  reduce = 0%
2022-12-17 11:40:36,651 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local88080146_0005
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 312524628 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
1       3.941410129096326
2       3.336917562724014
3       3.195542472666106
5       3.11385737439222
6       3.909004196871423
Time taken: 4.569 seconds, Fetched: 5 row(s)
hive> _
```

- 电影评价的所有电影评分的评价值

```
hive> select avg(rating) from new_ratings3;
Query ID = root_20221217114202_6203fc11-5030-4b66-8b01-ccc2ccfa6137
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:42:05,276 Stage-1 map = 0%,  reduce = 0%
2022-12-17 11:42:07,289 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local1405769857_0006
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 416699504 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
3.7439385
Time taken: 5.133 seconds, Fetched: 1 row(s)
hive> _
```

（5）电影评价的方差

- 电影评价的每部电影的评分的方差（前 5 个）：

```
hive> select movieId,stddev(rating) from new_ratings3 group by movieId limit 5;
Query ID = root_20221217114341_e7b52e67-bc7a-497f-ad7f-a611da277bcd
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:43:43,829 Stage-1 map = 0%,  reduce = 0%
2022-12-17 11:43:45,855 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local1324484301_0007
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 520874380 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
1       0.8927749999558637
2       0.9184093876044268
3       0.9586867134545866
5       0.9371467087526948
6       0.8438903346811538
Time taken: 4.291 seconds, Fetched: 5 row(s)
hive> _
```

- 电影评价的所有电影评分的方差

```
hive> select stddev(rating) from new_ratings3;
Query ID = root_20221217114423_50b1da47-19ba-42d8-a19f-cf2785a39afb
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2022-12-17 11:44:26,400 Stage-1 map = 0%,  reduce = 0%
2022-12-17 11:44:27,404 Stage-1 map = 100%,  reduce = 100%
Ended Job = job_local1088386712_0008
MapReduce Jobs Launched:
Stage-Stage-1:  HDFS Read: 625049256 HDFS Write: 0 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
0.9386925791854063
Time taken: 3.633 seconds, Fetched: 1 row(s)
hive> _
```
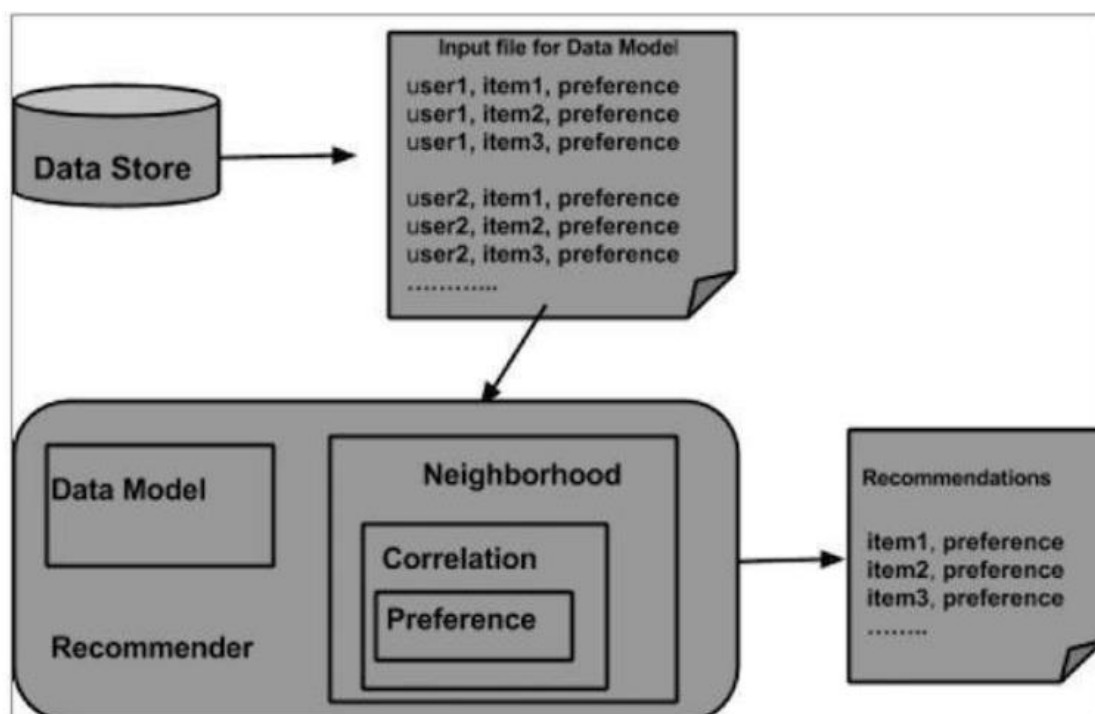
4. 详细介绍数据挖掘的算法、过程和结果。例如：使用了某个数据挖掘的算法，或编辑了某段代码，接着对多少用户和多少电影进行了协同过滤，得到了什么结果或什么成果，有什么结论。（10分）

（1）mahout 的推荐算法：

mahout 的推荐算法是基于 userbase 的,通过输入传递具有用户对项目首选项的文本文档，输出特定用户对其他项目的估计偏好。

从数据存储中，准备数据模型，并将其作为输入传递到推荐引擎。推荐器引擎为特定用户生成推荐。

推荐引擎的架构：



（2）将处理过后导出到本地的 ratings 表上传到 HDFS

```
[root@master ml-25m]# hdfs dfs -mkdir mahout/cf/input1
[root@master ml-25m]# hdfs dfs -put new_ratings3/000000_0 mahout/cf/input1
[root@master ml-25m]# _
```

（3）使用 mahout 的推荐算法处理上述步骤中上传的文件

**实验代码：**

```
hadoop jar mahout-examples-0.13.0-job.jar

org.apache.mahout.cf.taste.hadoop.item.RecommenderJob -i mahout/cf/input1/* -o

mahout/cf/output -s SIMILARITY_LOGLIKELIHOOD --tempDir /tmp/mahout/cf
```

**实验执行过程：**



**实验结果：**



（4）结论：

根据上面'3'的数据统计分析可知，推荐算法对 150526 名用户和 2923 部电影进行了协同过滤，得到了为 150526 名用户推荐的不同类型和不同个数的电影列表。