

### Homework 3

**Technology :** Python: 2.7 Scala: 2.11 Spark: 2.2.1

#### **Directory:**

Surbhi\_Batra\_hw2

- Surbhi\_Batra
  - Surbhi\_Batra\_Description
- OutputFiles
  - Surbhi\_Batra\_ModelBasedCF\_Small.txt
  - Surbhi\_Batra\_ModelBasedCF\_Big.txt
  - Surbhi\_Batra\_SimilarMovies\_Jaccard.txt
  - Surbhi\_Batra\_UserBasedCF.txt
  - Surbhi\_Batra\_ItemBasedCF.txt
- Solution
  - Surbhi\_Batra\_task1\_Jaccard.py
  - Surbhi\_Batra\_task2\_ModelBasedCF.py
  - Surbhi\_Batra\_task2\_UserBasedCF.py
  - Surbhi\_Batra\_task2\_ItemBasedCF.py

#### **Algorithm:**

>>>Execute these Commands after : cd/to/Surbhi\_Batra\_hw3/Solution

#### **Python:**

#### **TASK : 1**

#### **LSH**

Command :

\$SPARK\_HOME/bin/spark-submit Surbhi\_Batra\_task1\_Jaccard.py MovieLens.Small.csv

Precision: 1.000000

Recall: 0.909214

Time : 184.732897997

## **TASK : 2**

	<b><u>Task 2A</u></b>		<b><u>Task2B</u></b>
	<b><u>SMALL</u></b>	<b><u>LARGE</u></b>	<b><u>SMALL</u></b>
<b><u>&gt;=0 and &lt;1</u></b>	13808	3223235	15052
<b><u>&gt;=1 and &lt;2</u></b>	4146	733558	4160
<b><u>&gt;=2 and &lt;3</u></b>	665	81967	890
<b><u>&gt;=3 and &lt;4</u></b>	105	7383	152
<b><u>&gt;=4</u></b>	9	188	2
<b><u>RMSE</u></b>	0.949470593765	0.819786213789	0.962587911664
<b><u>TIME (seconds)</u></b>	18.9599509239	1797.52702594	81.8034739494

## **MODEL - BASED**

rank = 10  
numIterations = 10  
lambdas = 0.1

### **MovieLens.Small.csv**

**Command :** \$SPARK\_HOME/bin/spark-submit Surbhi\_Batra\_task2\_ModelBasedCF.py  
MovieLens.Small.csv testing\_small.cs

```

small

No of rows in ground truth: 20256.000000
No of rows in Predictions: 18733.000000
(10, 10, 0.1)

>=0 and <1 : 13808more
>=1 and <2 : 4146 less
>=2 and <3: 665 less
>=3 and <4 : 105 less
>=4      : 9 less

RMSE : 0.949470593765 less

Endtime: 18.9599509239

```

MovieLens.Big.csv

**Command :**

\$SPARK\_HOME/bin/spark-submit Surbhi\_Batra\_task2\_ModelBasedCF.py MovieLens.Big.csv  
testing\_20m.csv

Output file not attached as of now.. Because too big

```

small

No of rows in ground truth: 4054451.000000
No of rows in Predictions: 4046331.000000
(10, 10, 0.1)

>=0 and <1 : 3223235more
>=1 and <2 : 733558more
>=2 and <3: 81967more
>=3 and <4 : 7383more
>=4      : 188more

RMSE : 0.819786213789 less

Endtime: 1797.52702594

```

## USER - BASED

Used Pearson co-relation in the algorithm.

\$SPARK\_HOME/bin/spark-submit Surbhi\_Batra\_task2\_UserBasedCF.py MovieLens.Small.csv  
testing\_small.csv

```
>=0 and <1 : 15052more  
>=1 and <2 : 4160 less  
>=2 and <3: 890 less  
>=3 and <4 : 152 less  
>=4      : 2 less  
  
RMSE : 0.962587911664 less  
Time : 81.8034739494
```

### **TASK : 3**

**Command :** \$SPARK\_HOME/bin/spark-submit Surbhi\_Batra\_task2\_ItemBasedCF.py  
MovieLens.Small.csv testing\_small.csv

### **JACCARD SIMILARITY**

#### **With LSH**

```
>=0 and <1 : 800 less  
>=1 and <2 : 1299 less  
>=2 and <3: 2681more  
>=3 and <4 : 6307more  
>=4      : 9169more  
  
RMSE : 3.54410355144more  
154.843222857
```

RMSE : 3.54410355144more

### **Pearson**

RMSE : 0.96

The result is very bad because Jaccard similarity doesn't consider the actual ratings  
But just 0 and 1s and hence we do not have an accurate similarity measure.  
While Pearson co - relation calculates the actual similarity using corated items.