

Part 1:

9	Flerken		0.85442	5	1d
Your Best Entry 					

Part 2:

1. fixed numIterations = 20

a. rank = 5

```
In [9]: 1 rank = 5
2 numIterations = 20
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [10]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.6187260909958737

b. rank = 10

```
In [11]: 1 rank = 10
2 numIterations = 20
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [12]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.472247470360593

c. rank = 20

```
In [13]: 1 rank = 20
2 numIterations = 20
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [14]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.29089591114428714

d. rank = 30

```
In [15]: 1 rank = 30
2 numIterations = 20
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [16]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.18314214236371854

As it shows, as rank increased, the MSE decreased.

2. fixed rank = 20
 - a. numIterations = 2

```
In [3]: 1 rank = 20
2 numIterations = 2
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [4]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.5235247099214314

- b. numIterations = 5

```
In [7]: 1 rank = 20
2 numIterations = 5
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [8]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.337948246332415

- c. numIterations = 10

```
In [9]: 1 rank = 20
2 numIterations = 10
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [10]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.30763536590666024

- d. numIterations = 20

```
In [13]: 1 rank = 20
2 numIterations = 20
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [14]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Mean Squared Error = 0.2908959114428714

- e. numIterations = 30

because numIteration is too large, the function cannot run properly.

As we see, as numIteration increased, the MSE decreased.

3. With fixed rank = 20, numIterations = 30

```
In [1]: 1 from pyspark import SparkContext
2 from pyspark.mllib.recommendation import ALS, MatrixFactorizationModel, Rating
3 sc = SparkContext.getOrCreate()
4 sc.setLogLevel("WARN")
5 train_data_pre = sc.textFile('re_u.data')
6 train_data = sc.parallelize(train_data_pre.take(2000))
7 train_ratings = train_data.map(lambda l: l.split(',')).map(lambda l: Rating(int(l[0]), int(l[1]), float(l[2])))
```

```
In [3]: 1 rank = 20
2 numIterations = 20
3 model = ALS.train(train_ratings, rank, numIterations)
```

```
In [4]: 1 testdata = train_ratings.map(lambda p: (p[0], p[1]))
2 predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
3 ratesAndPreds = train_ratings.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
4 MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
5 print("Mean Squared Error = " + str(MSE))
```

Size of data

2000

Mean Squared Error

0.291120758082147923

5000	0.291259682781938134
10000	0.29208892894787620
20000	0.291487303979932457
50000	0.291021837001804823
100000	0.291294923919096340

As data shows, the size of data won't affect MSE.