# Spark ALS on EMR

In this document I'm recording critical steps on achieving movie recommendation on spark on AWS EC2 instances. Totally there are 3 primary steps as following:

**1.Set up AWS EMR**

**2. Go to Zeppelin on port 8890**

**3. Achieving movie recommendation with Spark Mllib ALS algorithm**

Here we go.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*


**1.Set up AWS EMR**

1.1 set AWS EMR  and connect it in terminal

## 2. Go to Zeppelin on port 8890



## 3. Achieving movie recommendation with Spark Mllib ALS algorithm

### 3.1 download data and save it on HDFS

## 3.2 read data from HDFS



## 3.3 Data stats

## 3.4 model and part of the code

**Zeppelin**    Notebook ▾    Job                  🔍 Search                    ● anonymous ▾

**ALS** ▷ ⣏ 📖 ✎ ⧉ ⬆   📄 ⊙ ⇄ Head ▾   🔍   🗑                    ⌨ ⚙ 🔒 default ▾

```
%pyspark
# ALS model                                                        ≡ SPARK JOBS  FINISHED  ▷ ⣏ 📖 ⚙
# the hyper parameter could be tweaked here
model = ALS.train(ratingsRDD, 10, 10, 0.01)
```

```
%pyspark
# recommend top 5 most related products to user_id 100              ≡ SPARK JOBS  FINISHED  ▷ ⣏ 📖 ⚙
model.recommendProducts(100,5)

[Rating(user=100, product=2192, rating=5.27332673384497), Rating(user=100, product=106, rating=5.021703887059571), Rating(user=100, product=632, rating=4.881118612831959), Rating(user=100, pr
oduct=3245, rating=4.802627726423592), Rating(user=100, product=958, rating=4.798462260132511)]
```

```
%pyspark
# recommend product_id 200 to top 5 users who might be interested   ≡ SPARK JOBS  FINISHED  ▷ ⣏ 📖 ⚙
model.recommendUsers(product=200,num=5)

[Rating(user=5760, product=200, rating=9.10564602910767), Rating(user=89, product=200, rating=8.048368649349037), Rating(user=1459, product=200, rating=7.564286080564312), Rating(user=3897, p
roduct=200, rating=7.478246775292424), Rating(user=2697, product=200, rating=7.3095822341429875)]
```

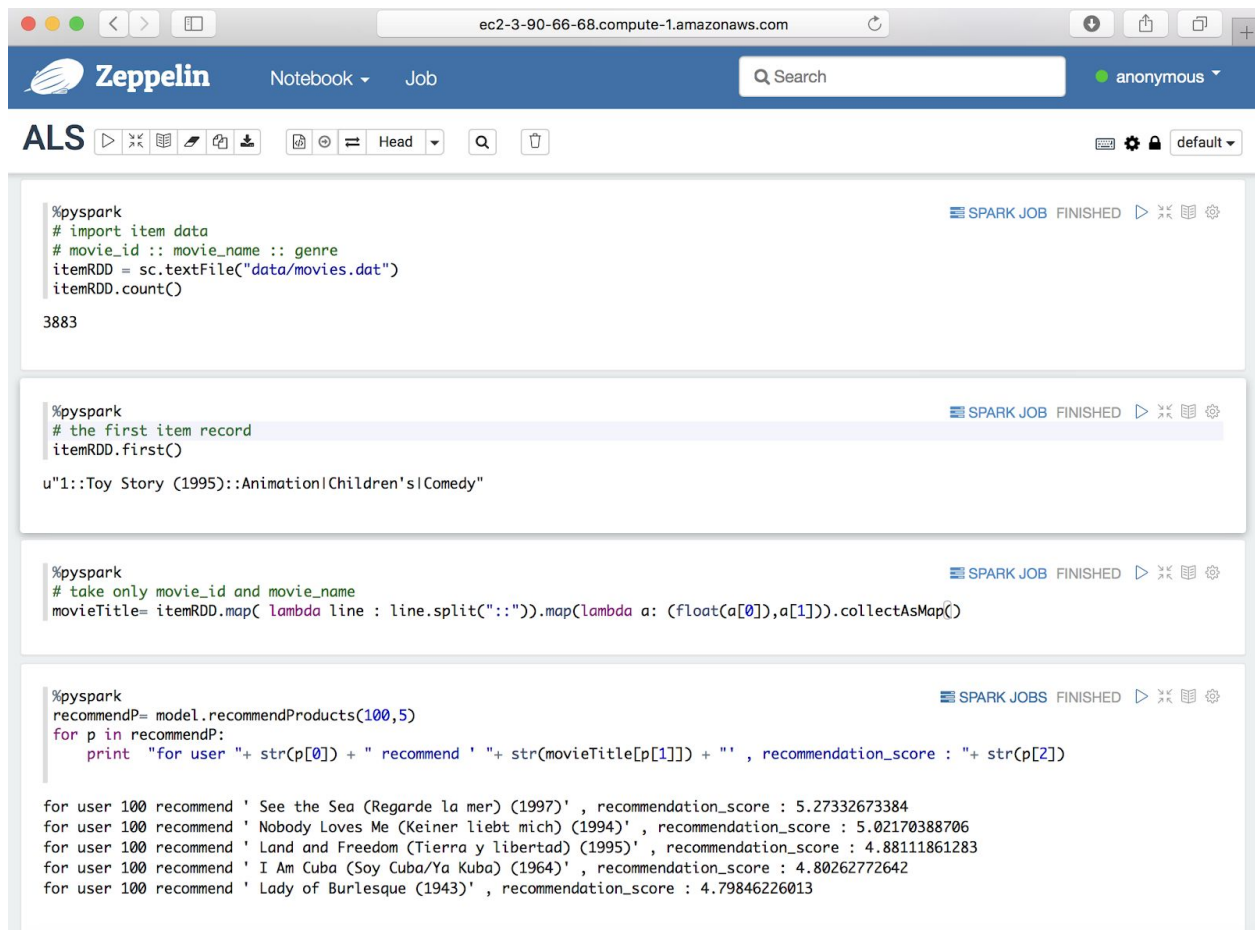**Zeppelin**    Notebook ▾    Job                  🔍 Search                    ● anonymous ▾

**ALS** ▷ ⣏ 📖 ✎ ⧉ ⬆   📄 ⊙ ⇄ Head ▾   🔍   🗑                    ⌨ ⚙ 🔒 default ▾

```
%pyspark
# import item data                                                 ≡ SPARK JOB  FINISHED  ▷ ⣏ 📖 ⚙
# movie_id :: movie_name :: genre
itemRDD = sc.textFile("data/movies.dat")
itemRDD.count()

3883
```

```
%pyspark
# the first item record                                            ≡ SPARK JOB  FINISHED  ▷ ⣏ 📖 ⚙
itemRDD.first()

u"1::Toy Story (1995)::Animation|Children's|Comedy"
```

```
%pyspark
# take only movie_id and movie_name                                ≡ SPARK JOB  FINISHED  ▷ ⣏ 📖 ⚙
movieTitle= itemRDD.map( lambda line : line.split("::")).map(lambda a: (float(a[0]),a[1])).collectAsMap()
```

```
%pyspark
recommendP= model.recommendProducts(100,5)                         ≡ SPARK JOBS  FINISHED  ▷ ⣏ 📖 ⚙
for p in recommendP:
    print  "for user "+ str(p[0]) + " recommend ' "+ str(movieTitle[p[1]]) + "' , recommendation_score : "+ str(p[2])

for user 100 recommend ' See the Sea (Regarde la mer) (1997)' , recommendation_score : 5.27332673384
for user 100 recommend ' Nobody Loves Me (Keiner liebt mich) (1994)' , recommendation_score : 5.02170388706
for user 100 recommend ' Land and Freedom (Tierra y libertad) (1995)' , recommendation_score : 4.88111861283
for user 100 recommend ' I Am Cuba (Soy Cuba/Ya Kuba) (1964)' , recommendation_score : 4.80262772642
for user 100 recommend ' Lady of Burlesque (1943)' , recommendation_score : 4.79846226013
```