# NETFLIX RECOMMENDATION SYSTEM

## CSYE7200 34322 BIG-DATA SYS ENGR USING SCALA SEC 01 – SPRING 2018

Team - 1

Sonali Chaudhari-001285029
Praneeth Reddy-001225762

# OUR PROPOSAL

- WITH THE HELP OF USER DATA AND RATINGS.

- GENERATE A RELEVANT SUGGESTION BASED ON PAST AVAILABLE

  DATASET.

- RECOMMENDING HIGHEST PREDICTED RATING TO A PARTICULAR USER.

- MEETING THE DEADLINE OF THE PROJECT.

# USE CASE/ ACTOR

**ACTOR**

- **USER**

- **THE SYSTEM**

**USE CASE**

- **USER WILL BELOW OPERATIONS:**

  - **PROVIDING RATING FOR MOVIES**

- **THE APPLICATION WILL PROVIDE LIST OF MOVIES WITH HIGHLY PREDICTED RATINGS**
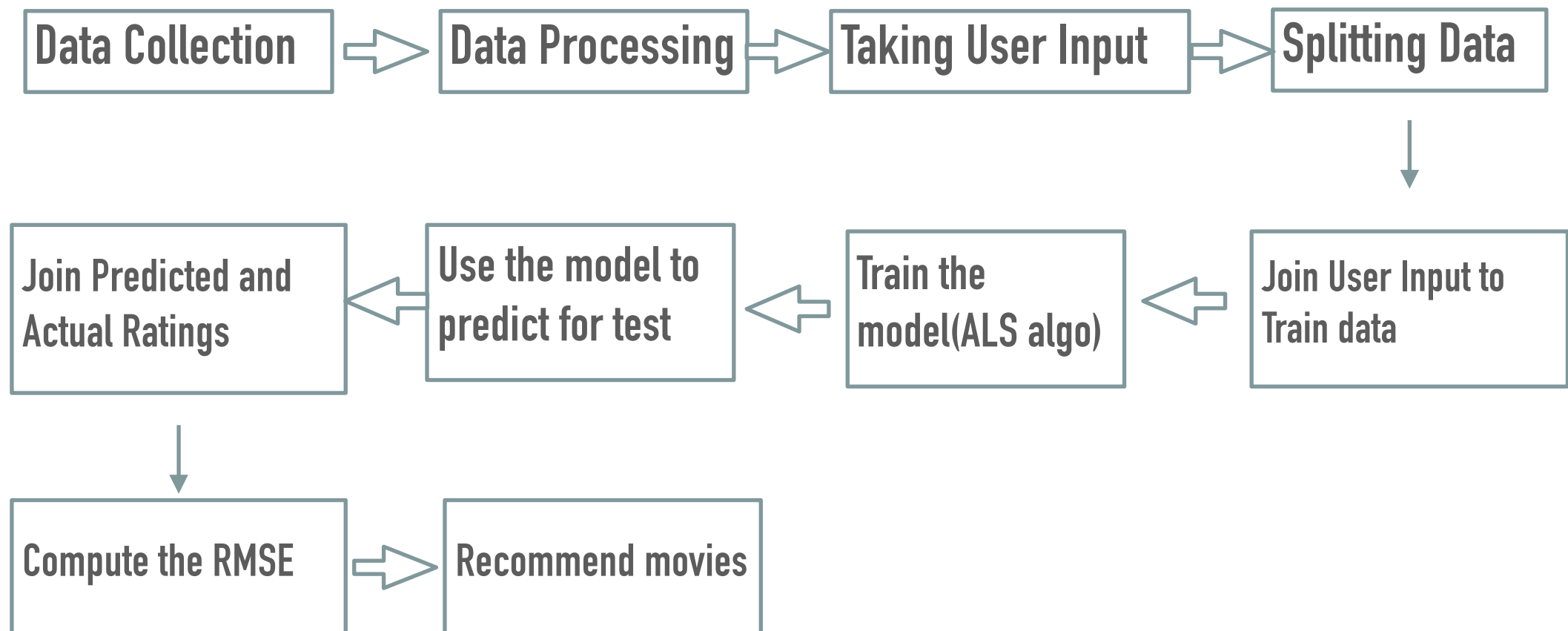
# DETAILS OF DATA

**ACTUAL DATA:**

- **RATINGS DATA : 10048050**

- **MOVIES : 17770**

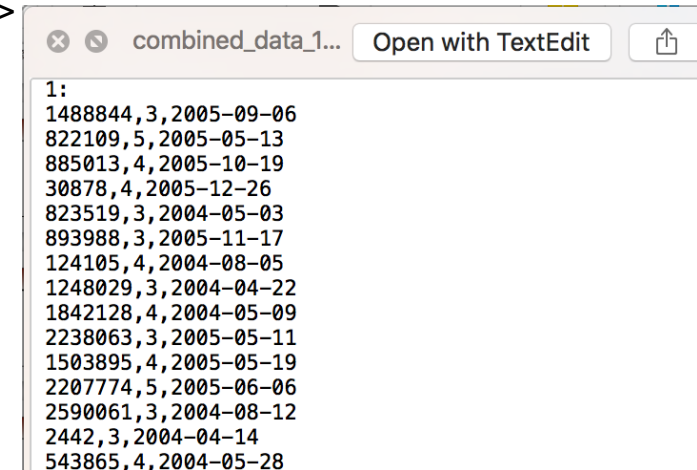- **USER: 480189**

**DATA USED  (HEAP SIZE ISSUE):**

- **RATING DATA : 5010199**

- **MOVIES : 1000**

- **USERS: 404555**

# WORKFLOW

Data Collection ➡ Data Processing ➡ Taking User Input ➡ Splitting Data

Join Predicted and Actual Ratings ⬅ Use the model to predict for test ⬅ Train the model(ALS algo) ⬅ Join User Input to Train data

Compute the RMSE ➡ Recommend movies

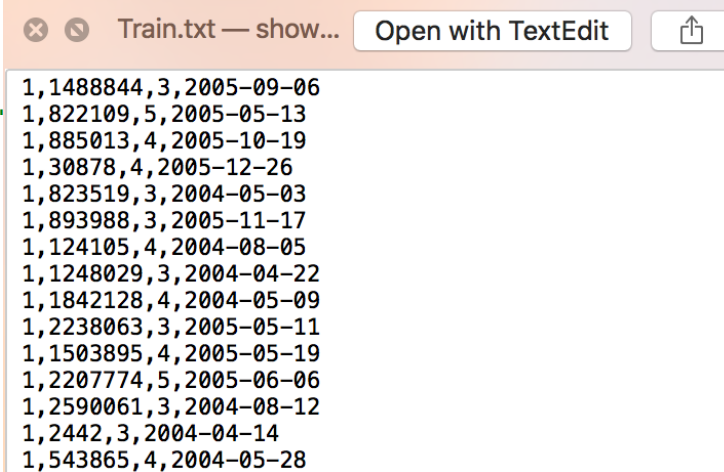# DATA CLEANING

- The train data is contained in 4 different text files in the format shown ->



```
1:
1488844,3,2005-09-06
822109,5,2005-05-13
885013,4,2005-10-19
30878,4,2005-12-26
823519,3,2004-05-03
893988,3,2005-11-17
124105,4,2004-08-05
1248029,3,2004-04-22
1842128,4,2004-05-09
2238063,3,2005-05-11
1503895,4,2005-05-19
2207774,5,2005-06-06
2590061,3,2004-08-12
2442,3,2004-04-14
543865,4,2004-05-28
```

- The four txt file data are converted into the required format and combine into Train.txt file

```scala
// Generating a combined file in the required format
val NtflixRecosFile = "/Users/sonalichaudhari/Desktop/netflix-prize-data/"
val file = new File("/Users/sonalichaudhari/Desktop/netflix-prize-data/Train.txt")
val bw = new BufferedWriter(new FileWriter(file))
var train_files = Array("combined_data_1.txt","combined_data_2.txt","combined_data_3.txt","combined_data_4.txt")
    for ( i <- 0 to (train_files.length - 1)) {
      var app = ""
      for (line <- Source.fromFile(NtflixRecosFile+train_files(i)).getLines) {
        if (line.contains(":")) {
          app = line.toString().stripSuffix(":")
        }
        else {
          var entry = app+","+line
          bw.write(entry+"\n")
        }
      }
    }
println(("Combined file generated!!"))
bw.close()
```

```
1,1488844,3,2005-09-06
1,822109,5,2005-05-13
1,885013,4,2005-10-19
1,30878,4,2005-12-26
1,823519,3,2004-05-03
1,893988,3,2005-11-17
1,124105,4,2004-08-05
1,1248029,3,2004-04-22
1,1842128,4,2004-05-09
1,2238063,3,2005-05-11
1,1503895,4,2005-05-19
1,2207774,5,2005-06-06
1,2590061,3,2004-08-12
1,2442,3,2004-04-14
1,543865,4,2004-05-28
```

# TUNING THE MODEL

- ALS algorithm

- Parameters used in the model ALS(rank,iteration, regularization factor)

- Following is how the rise vary with different values of ALS parameters.

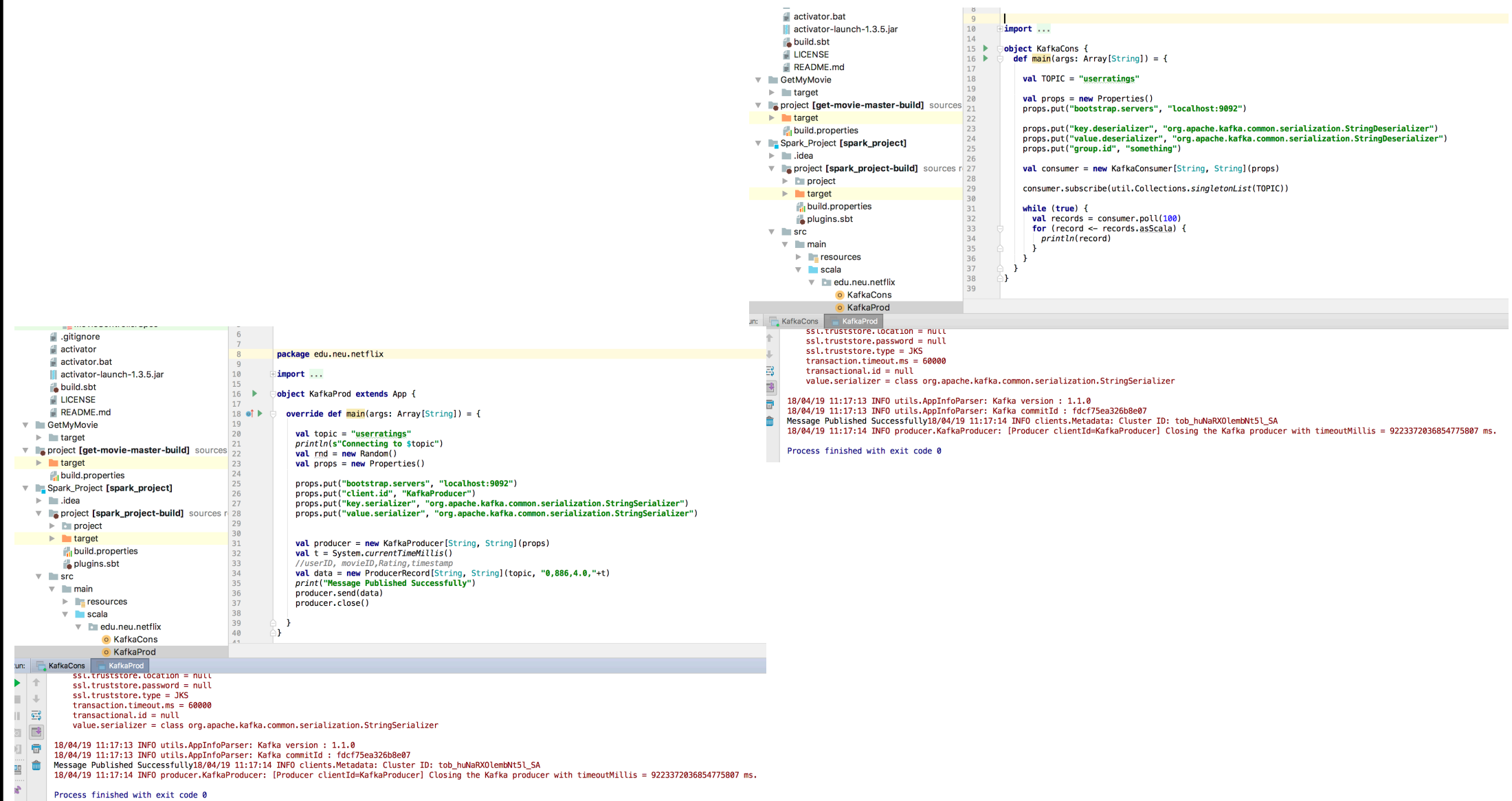| ALS(rank, iteration,regularization factor) | RMSE |
|---|---|
| 8,5,0.02 | 1.10 |
| 9, 4, 0.02 | 1.09 |
| 8,10,0.01 | 1.14 |
| 8, 20, 0.01 | 1.16 |
| 10, 5, 0.02 | 1.09 |
| 12, 5, 0.03 | 1.08 |
| 8, 5, 0.05 | 1.03 |
| 8,4,0.05 | 1.02 |
| 8, 4, 0.06 | 1.015 |
| 8, 4, 0.07 | 1.006 |
| 8,5,0.01 | 1.10 |
| **8,5,0.099** | **0.97** |

# RMSE FOR DIFFERENT RUNS

# PREDICTION ACCURACY

**ROOT MEAN SQUARE ERROR(RMSE)**

**RMSE IS THE PARAMETER USED TO MEASURE THE DIFFERENCE BETWEEN THE PREDICTED VALUES AND THE ACTUAL VALUES.**

$$RMSE = \sqrt{\sum \frac{(y_{pred} - y_{ref})^2}{N}}$$

# APPLICATION AND USER INTERFACE

# KAFKA PRODUCER AND CONSUMER

## Movie Database

<------Please rate the movies (1(Low) to 5(High)) and Get Suggestions-->

| Movie | Rating |
|---|---|
| Ray | |
| | |
| Speed | |
| Reservoir Dogs | |
| Mean Girls | |
| Something's Gotta Give | |
| X-Men | |
| American Beaty | |
| Rush Hour | |
| Pay it forward | |

Show 10 entries    Search:

| name | expectedRating |
|---|---|
| Loading... | |

Showing 0 to 0 of 0 entries

Previous    Next

NETFLIX

# ACCEPTANCE CRITERIA

- **APPLICATION SHOULD BE ABLE TO HANDLE AT LEAST 2500 REQUESTS SIMULTANEOUSLY AND THE MODEL SHOULD BE SCALABLE TO ADD NEW DATA SOURCES AS AND WHEN REQUIRED.**

- **ACHIEVE >90% ACCURACY**



```
// Model training
val model = ALS.train(training,8,5,0.099)

// Implementing trained model on the test
val prediction = model.predict(test.map(x

// Joining predicted values and actual va

DataProcessing  ›  main(args: Array[String])

taProcessing
18/04/18 18:14:25 INFO executor.Executor:
 (TID 165). 945 bytes result sent to driv
18/04/18 18:14:25 INFO scheduler.TaskSetMa
 stage 32.0 (TID 165) in 990 ms on localh
18/04/18 18:14:25 INFO scheduler.TaskSche
 whose tasks have all completed, from poo
18/04/18 18:14:25 INFO scheduler.DAGSched
 DataProcessing.scala:111) finished in 0.
18/04/18 18:14:25 INFO spark.SparkContext
 DataProcessing.scala:111, took 3.6811863
RMSE: 0.9750691117807623
```

# CHALLENGES FACED

- SPARK, SCALA, KAFKA AND PLAY FRAMEWORK COMPATIBILITY ISSUE

- RESTRUCTURING DATA TO USE INTO CORRECT  FORMAT

- OVERFITTING OF THE MODEL

- TUNING THE MODEL

- HEAP SIZE ISSUE

# USING PLAY FRAMEWORK

- IMPLEMENTED MVC

- INTEGRATING SPARK

- IMPLEMENTED MOCKITO FOR APPLICATION ,MOVIE
  CONTROLLER SPEC TO TEST VARIOUS FEATURES

# USING PLAY FRAMEWORK

- **TASK COMPLETED**

  - DATA PROCESSING

  - TAKING USER INPUT

  - PREDICTION GENERATING

- **TASK TO BE COMPLETED**

  - USER INTERFACE

# GITHUB REPOSITORY LINK

https://github.com/reddyse/Big-Data-Engineering-Using-Scala

THANK YOU...