# Machine Learning using Spark MLlib

Spark MLlib allows us to perform Machina Learning and Data Science operations for huge datasets on distributed clusters.

Before we start, we should realize that not all Machine Learning algorithms are easy to implement on distributed systems. This makes ML in Spark challenging. Sparks MLlib and the recent ML library has come far in creating algorithms, but it has a lot of ground to cover. Some of the algorithms covered in MLlib are:

1. Feature Extraction

2. Basic Statistics

3. Linear and Logistic Regression

4. Support Vector Machines

5. Naïve Bayes Classifier

6. Decision Trees

7. Decision Trees

MLlib comes with its specialized data types:

1. Vector – Vector is a way of representing large arrays of values, which may have missing values. For this we have Sparse and Dense vectors.

2. Labeled Point – Supervised ML is all about labelling data points based on some rules. Labeled Point helps us in making that association.

3. Rating – This data type is specialized for data which contains ratings for products given by consumers like in movies or online shopping.

**Getting the Data:**

We will be working on the movie rating dataset to predict movie recommendations. The data is read from the file and split using '\t' delimiter. We only take the columns that are needed by us (UserID, movieID and Rating) and store it in the final dataframe called 'ratings'.
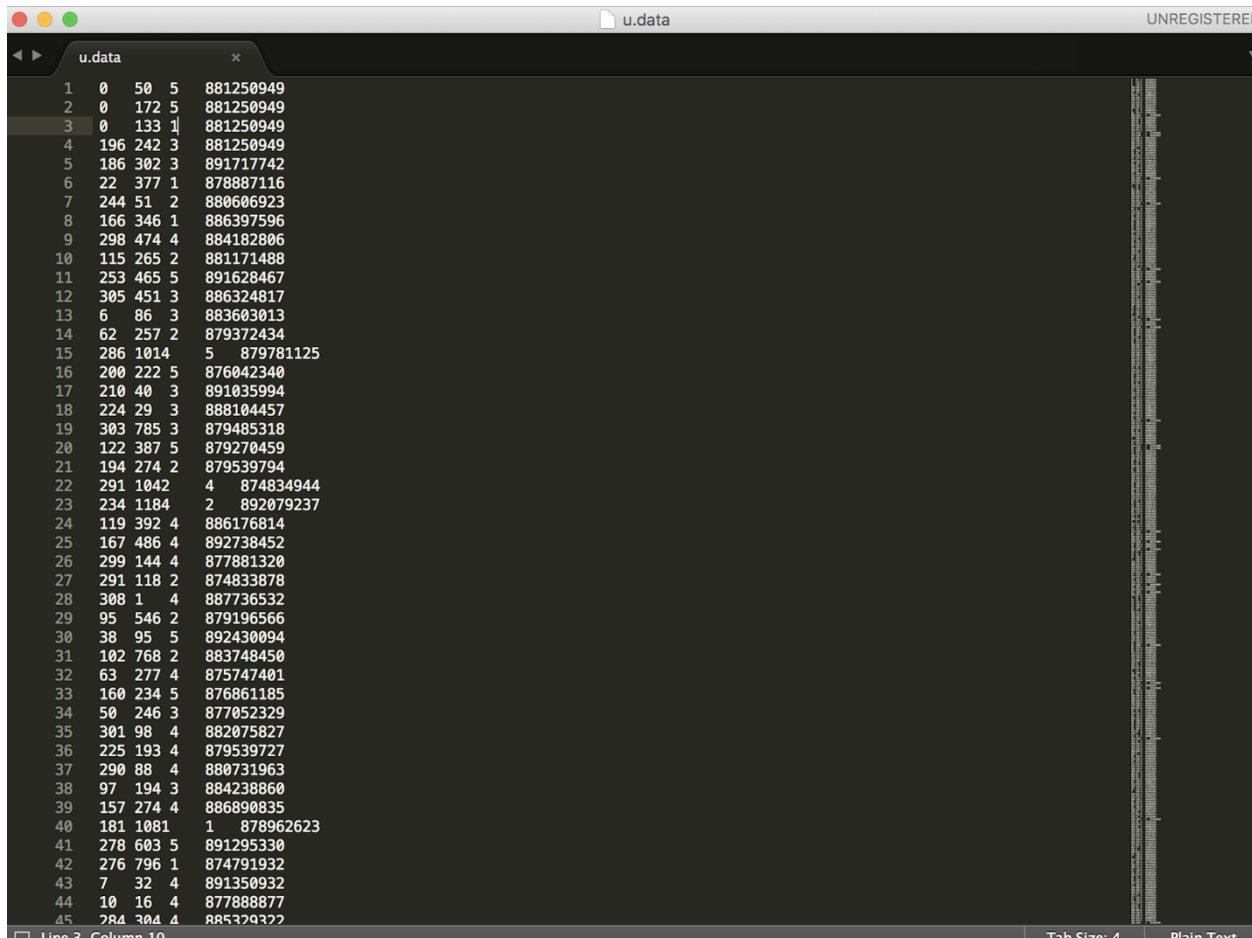
**Creating the ML Model:**

We will be using the Alternating Least Squares (ALS) model to learn from the ratings data. We will train this model using parameters like Rank and Number of Iterations on the training data (ratings dataframe). We store the trained model in the variable 'model'.

We later on use the 'recommendProducts' function to create 10 recommendations for any given userID provided through the command line.
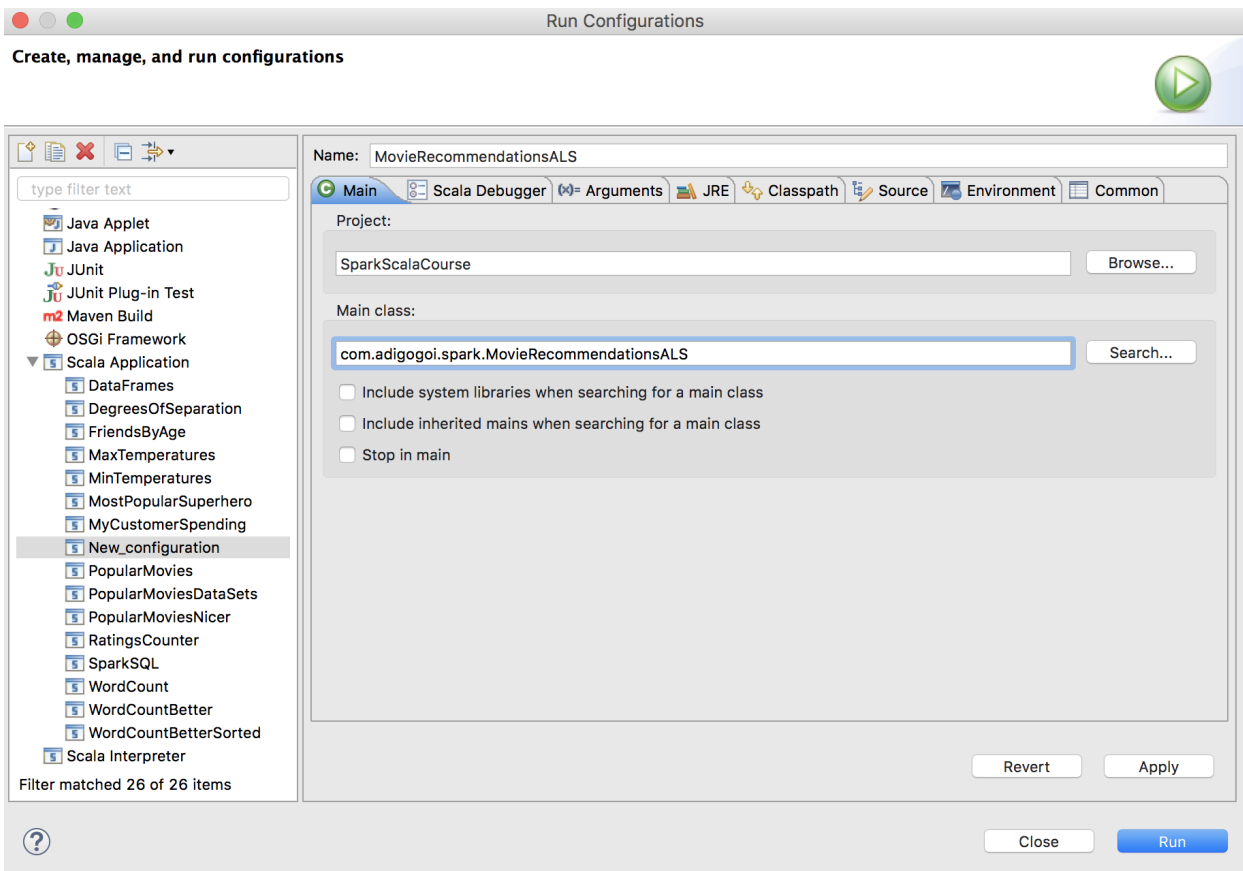
**Running the Program:**

To test how ALS model of MLlib works, we will have to make a few adjustments to the dataset and then run the code in a specific way.
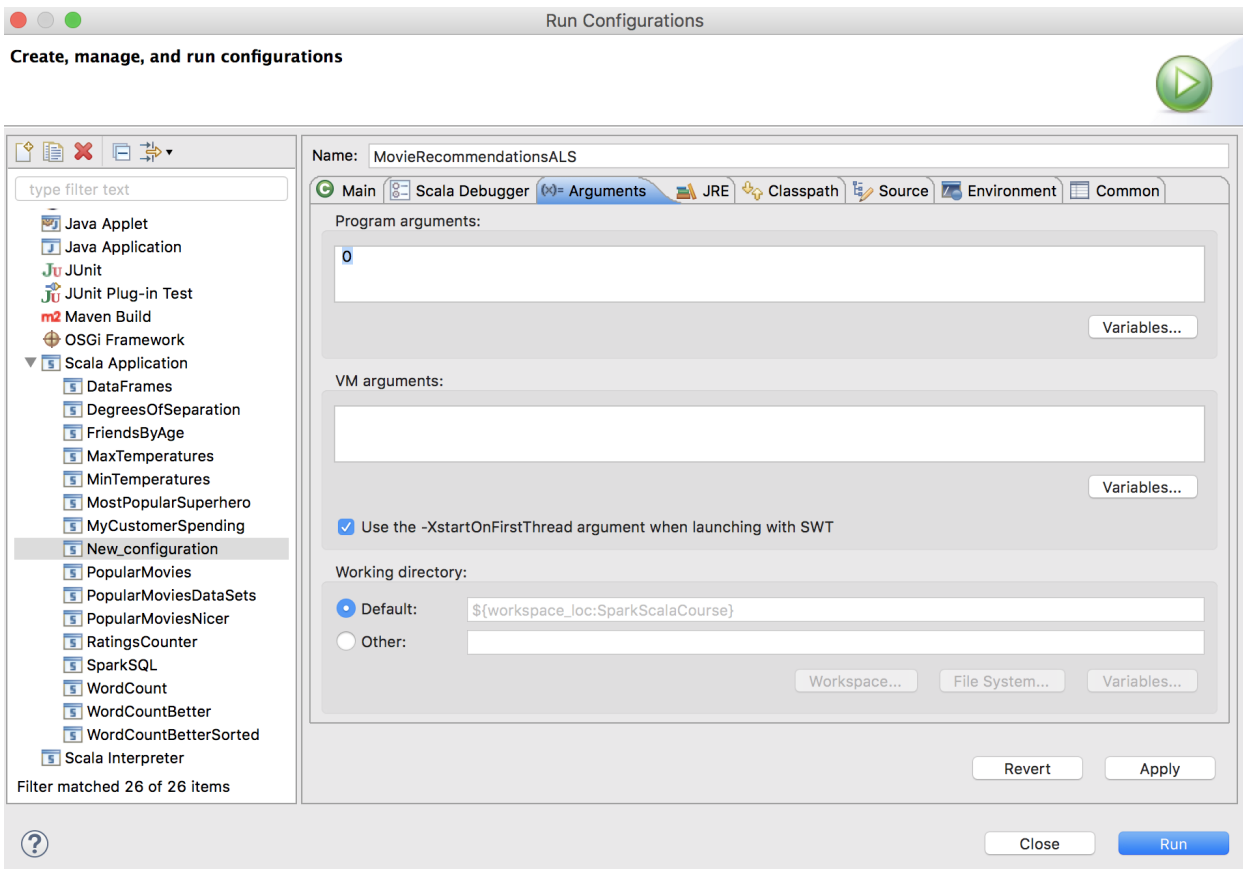
1. We will create a fictional user. This user loves Sci-fi movies and has hence rated Star Wars (Movie ID 50) a 5. This user hates romantic movies and has given Gone With The Wind (Movie ID 133) a 1 rating. This will be stored in our data file as User ID 0.



2. After that adjustment, we will run the code as if through a command line.  For this we have to go to the 'Run Configuration' option in the Scala IDE for Eclipse that we are using and then give the project name and the Main class of my project.

3. Next we go to the 'Arguments' window and enter '0', which is the User ID we want to predict our recommendations for.

4. The program will the train the model for the number of stages specified by us.



5. According to my code, the next thing to be displayed will be the movies rated by the User ID 0 and the rating given for each movie.



6. The ALS model will then make predictions and provide 10 recommendations for User ID 0.

## Conclusions:

After watching the model work and seeing its inaccurate predictions, I have come to the following conclusions:

1. The ALS algorithm is very sensitive to the parameters provided to it. It takes a little time to find the ideal parameters for a dataset. But due to the black-box nature of the algorithms, it becomes difficult to determine the correct parameters.

2. Complicated is not always better. Sometimes a simple approach to the problem may provide better results.

3. Never blindly trust results when analyzing Big Data. Small problems in algorithms magnify when dealing with such large datasets. Quality of data is the real issue here.