Thejas Bharadwaj – SEC01 (NUID 002727189)

# Big Data System Engineering with Scala
# Spring 2023
# Assignment No. 7
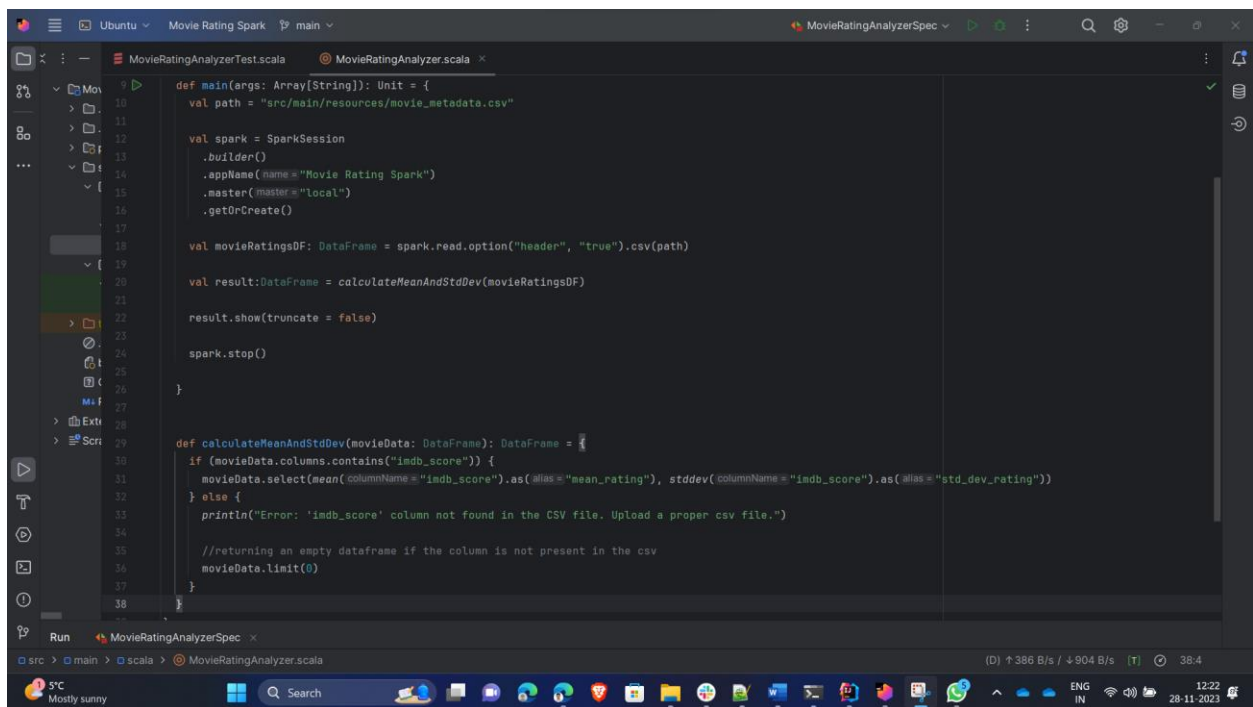
## -List of Tasks Implemented

1) Created a Github repo - [thejas98/Movie-Rating-Spark (github.com)](https://github.com)

2) Created two code files – main and test

3) The main file has the code to ingest the csv and a function which calculates the mean and the standard deviation of the 'imdb_score' column.

4) The test file contains 2 test cases – One creates a sample df and tests if the function calculatemean from main works correctly and the other makes sure if you upload a csv without a 'imdb_score' column, it handles the error gracefully.

## -Code

### 1) Main file



**Output –**

## 2) Test cases

### Test case #1

```scala
"calculateMeanAndStdDev" should "return an empty DataFrame if 'imdb_score' column is not present" in {
  import spark.implicits._

  // Test data without 'imdb_score' column
  val testData = Seq(
    (1, "Avengers"),
    (2, "Thor"),
    (3, "A beautiful mind")
  )

  val columns = Seq("id", "title")
  val movieDataWithoutImdbScore: DataFrame = testData.toDF(columns: _*)

  val result = calculateMeanAndStdDev(movieDataWithoutImdbScore)


  //result should be 0 since imdb_score column in not present
  result.count() shouldBe 0
}
```

### Test case #2

```scala
"calculateMeanAndStdDev" should "return the correct mean and standard deviation" in {
  import spark.implicits._

  // Test data - To test if the function(calculateMeanAndStdDev) i created works properly
  val testData = Seq(
    (1, "Avengers", 7.5),
    (2, "Thor", 8.0),
    (3, "A beautiful mind", 6.5),
  )

  val columns = Seq("id", "title", "imdb_score")
  val movieRatingsDF: DataFrame = testData.toDF(columns: _*)

  val result = calculateMeanAndStdDev(movieRatingsDF)

  val expectedMean = testData.map(_._3).sum / testData.length.toDouble
  val expectedStdDev = math.sqrt(testData.map(score => math.pow(score._3 - expectedMean, 2)).sum / (testData.length - 1).toDouble)

  val resultRow = result.head()
  val resultMean = resultRow.getAs[Double](fieldName = "mean_rating")
  val resultStdDev = resultRow.getAs[Double](fieldName = "std_dev_rating")

  resultMean shouldEqual expectedMean +- 0.001
  resultStdDev shouldEqual expectedStdDev +- 0.001
}
```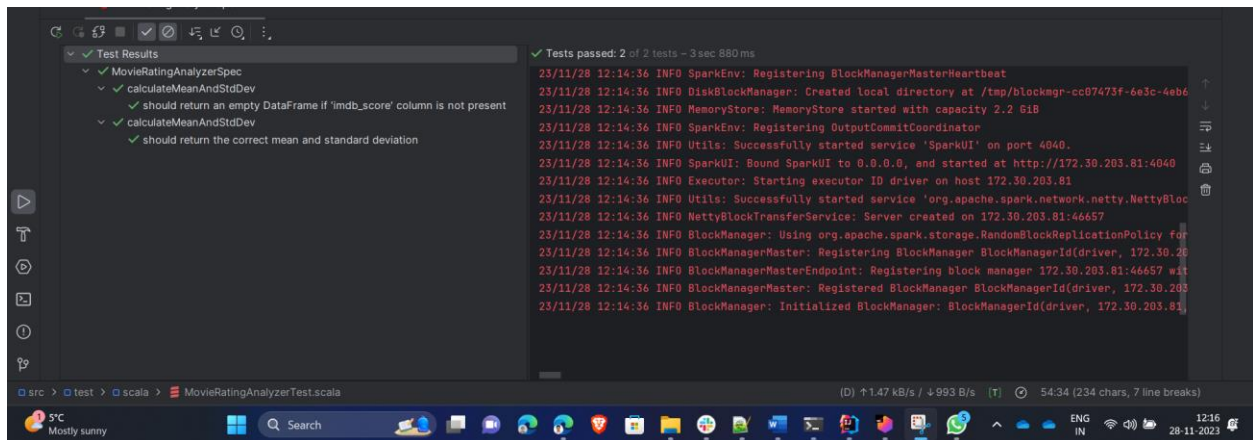