Thejas Bharadwaj – SEC01 (NUID 002727189)

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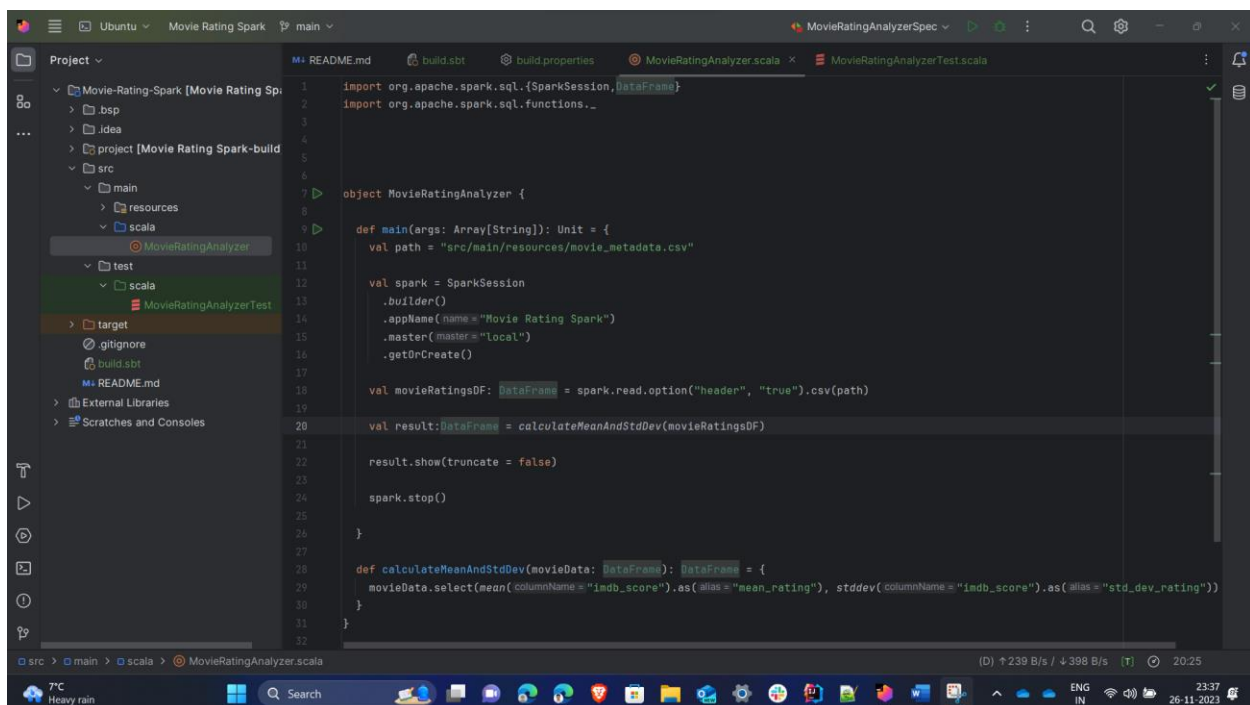
**-List of Tasks Implemented**

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

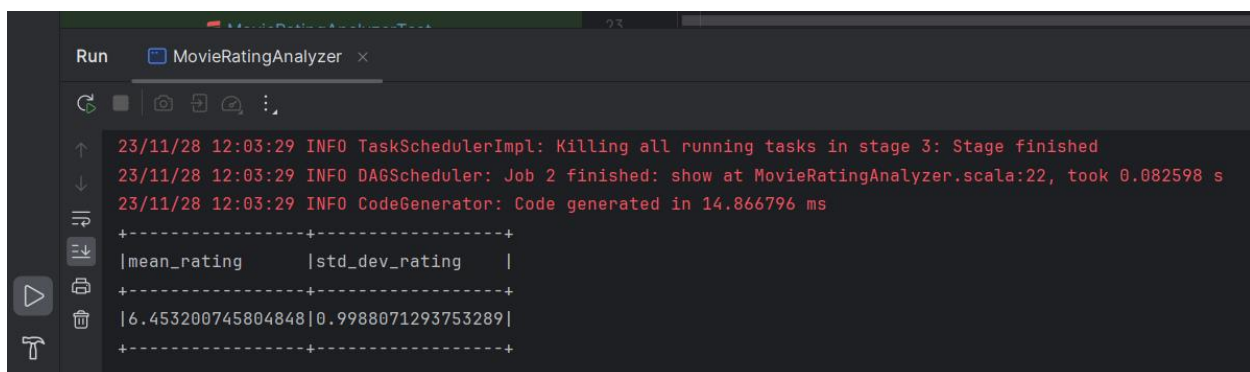4) The test file contains 2 test cases – One creates a sample df and tests if the function from main works correctly
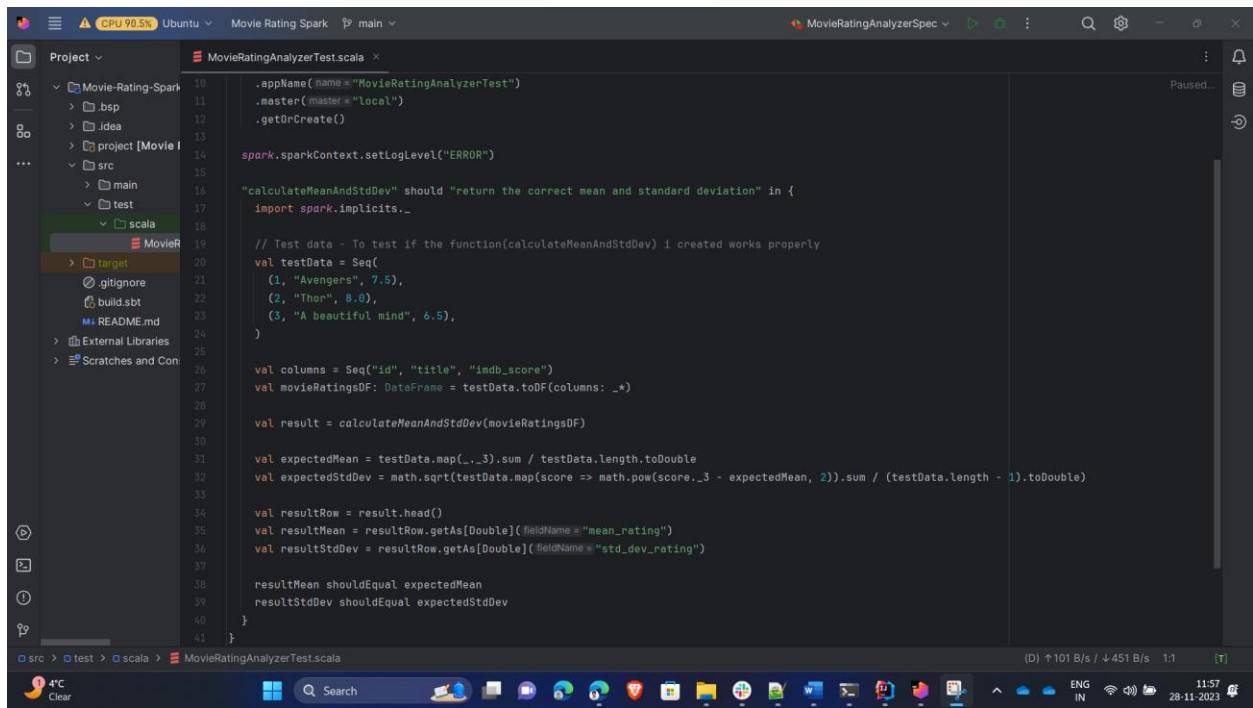
**-Code**

**1) Main file**



**Output –**

## Test file



## -Unit tests

MovieRatingAnalyzerTest.scala

```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)
```

Run    MovieRatingAnalyzerSpec ×

- ✓ Test Results
  - ✓ MovieRatingAnalyzerSpec
    - ✓ calculateMeanAndStdDev
      - ✓ should return the correct mean and standard deviation

✓ Tests passed: 1 of 1 test – 3 sec 749 ms

```
23/11/28 11:59:05 INFO Executor: Starting executor ID driver on host 172.30.203.81
23/11/28 11:59:06 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferSe
23/11/28 11:59:06 INFO NettyBlockTransferService: Server created on 172.30.203.81:44395
23/11/28 11:59:06 INFO BlockManager: Using org.apache.spark.storage.RandomBlockReplicationPolicy for block repl
23/11/28 11:59:06 INFO BlockManagerMaster: Registering BlockManager BlockManagerId(driver, 172.30.203.81, 44395
23/11/28 11:59:06 INFO BlockManagerMasterEndpoint: Registering block manager 172.30.203.81:44395 with 2.2 GiB R
23/11/28 11:59:06 INFO BlockManagerMaster: Registered BlockManager BlockManagerId(driver, 172.30.203.81, 44395,
23/11/28 11:59:06 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, 172.30.203.81, 44395, Non
```

src > test > scala > MovieRatingAnalyzerTest.scala