

Olasunkanmi Olayinka – SEC01 (NUID 001512266)

Big Data System Engineering with Scala

Fall 2022

Assignment No. 07

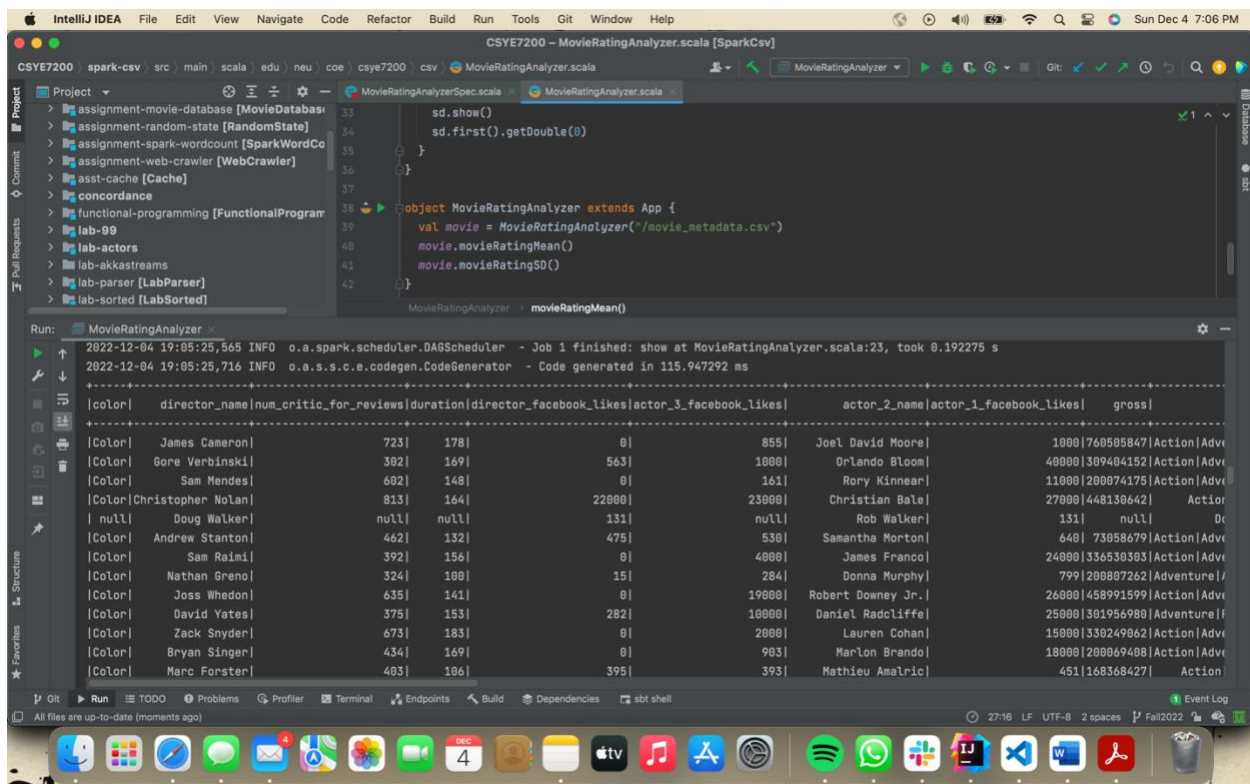


-List of Tasks Implemented

- Analyze a movie rating dataset stored in a CSV file
- Read the file using spark
- Calculate the mean and standard deviation for all movies
- Write a test case to ensure the programs works

-Code

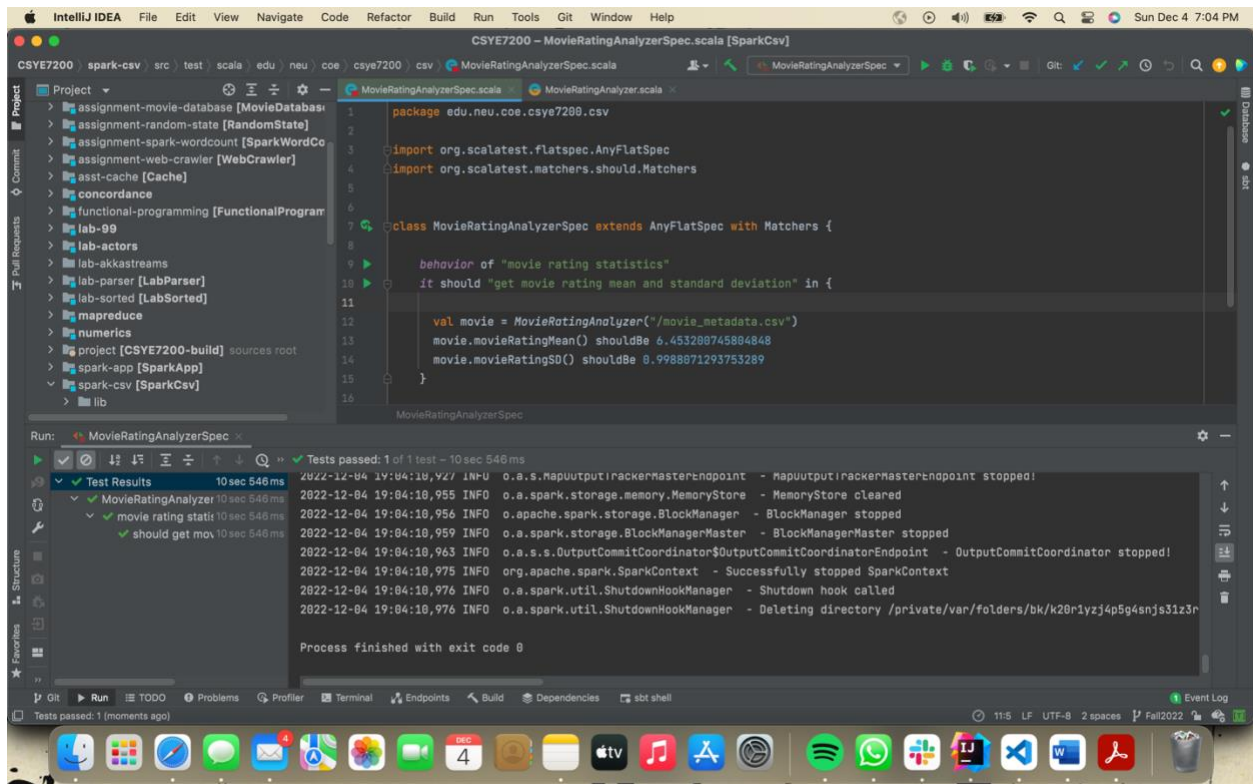
```
9  case class MovieRatingAnalyzer(resource: String) {
10
11    val spark: SparkSession = SparkSession
12      .builder()
13      .appName( name = "MovieRating")
14      .master( master = "local[*]")
15      .getOrCreate()
16
17    spark.sparkContext.setLogLevel("ERROR")
18
19    def apply(resource: String): MovieRatingAnalyzer = new MovieRatingAnalyzer(resource)
20
21    val df: DataFrame = spark.read.format( source = "csv").option("header", "true")
22      .load(getClass.getResource(resource).getPath)
23    df.show()
24
25    def movieRatingMean(): Double = {
26      val mean = df.select(functions.avg( columnName = "imdb_score"))
27      mean.show()
28      mean.first().getDouble(0)
29    }
30
31    def movieRatingSD(): Double = {
32      val sd = df.select(functions.stddev( columnName = "imdb_score"))
33      sd.show()
34      sd.first().getDouble(0)
35    }
36  }
37
```



```
2022-12-04 19:05:26,529 INFO o.a.s.s.c.e.codegen.CodeGenerator - Code generated in 20.165542 ms
+-----+
| avg(imdb_score)|
+-----+
| 6.453200745804848|
+-----+
```

```
2022-12-04 19:05:27,150 INFO o.a.spark.scheduler.DAGScheduler - Job 4 finished: show at MovieRatingAnalyzer.scala:33, took 0.095737 s
+-----+
| stddev_samp(imdb_score)|
+-----+
| 0.9988071293753289|
+-----+
```

-Unit tests



- Result

Mean = 6.453200745804848

Standard Deviation = 0.9988071293753289