Qinyun Lin – SEC01 (NUID 001582464)
# Big Data System Engineering with Scala
# Spring 2022
# Assignment No. 6

# Task

1. Read Movie dateset into spark
2. Calculate the mean rating and standard deviation for all movies

# Solution

The data I used is the one in the repository.
There are 1609 rows of records in the data, I calculate the "Ratings" based on the "imdb_score" column. The result I got is shown below:

```
+-----------------+------------------+
|             mean|           std_dev|
+-----------------+------------------+
|6.453200745804848|0.9988071293753289|
+-----------------+------------------+
```

### 1. Read Movie dateset into spark

```scala
val spark: SparkSession = SparkSession
  .builder()
  .appName( name = "AnalyzeMovieRating")
  .master( master = "local")
  .getOrCreate()

spark.sparkContext.setLogLevel("ERROR") // We want to ignore all of the INFO and WARN messages.

val resource = "src/main/resources/movie_metadata.csv"
val df:DataFrame = spark.read.option("header", "true").csv(resource)
```
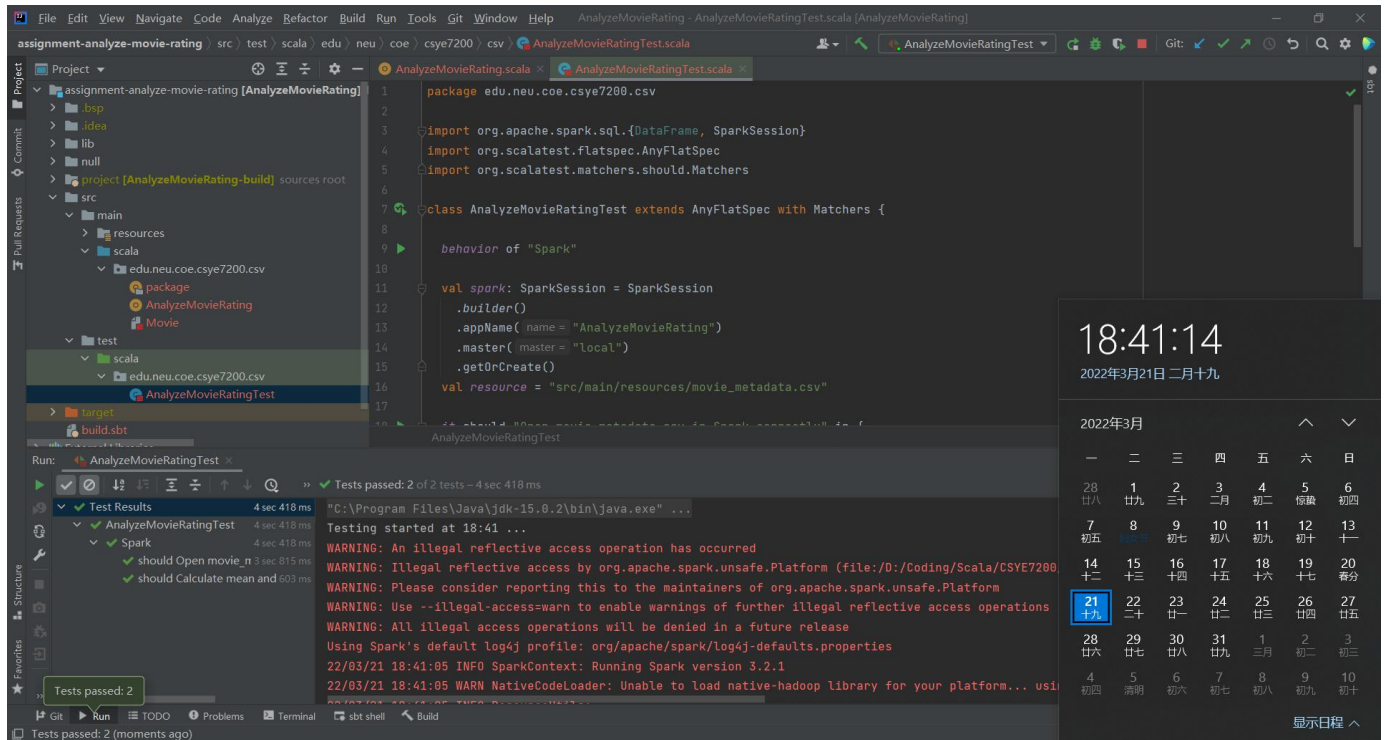
### 2. Calculate the mean rating and standard deviation for all movies

```scala
def calcMeanAndStd(df: DataFrame): DataFrame = {
  val colName = "imdb_score"
  df.select(
    mean(df(colName)).alias( alias = "mean"),
    stddev(df(colName)).alias( alias = "std_dev")
  )
}
```

# Unit Test Screenshot



# Project Source