Qinyun Lin – SEC01 (NUID 001582464)
# Big Data System Engineering with Scala
# Spring 2022
# Assignment No. 1
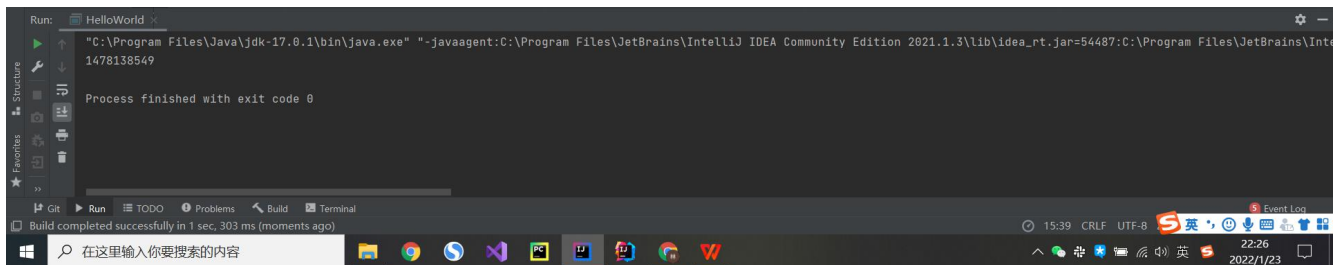
# Task

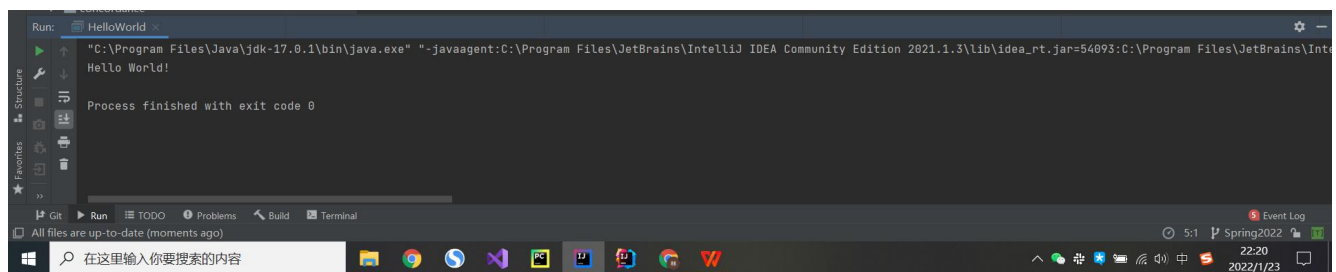**Ensure a development environment can compile, test and run Scala programs**

1. **Check the result of "scala.util.Random.nextInt"**
2. **Run "HelloWord"**
3. **Run "Ingest.scala" and view the list of movies**
4. **Modify "Ingest.scala" and find movies of "New Zealand"**
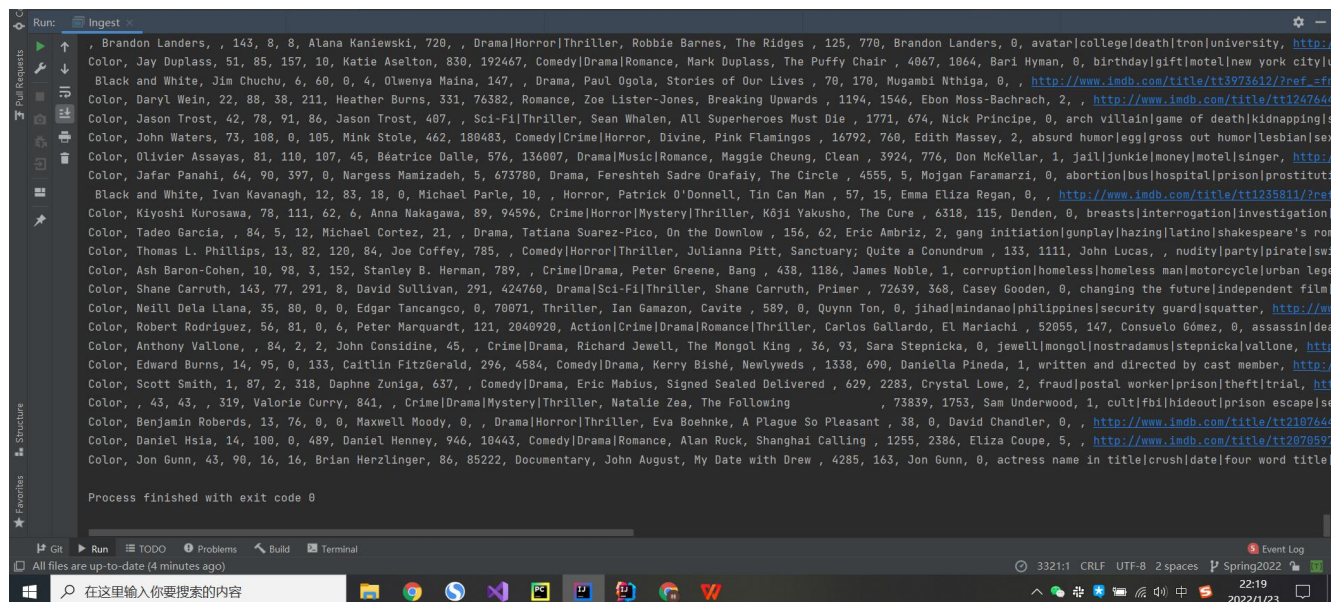5. **Get "Spark ready" and familiar with interacting with spark program**

# Solution

## 1. Check the result of "scala.util.Random.nextInt"
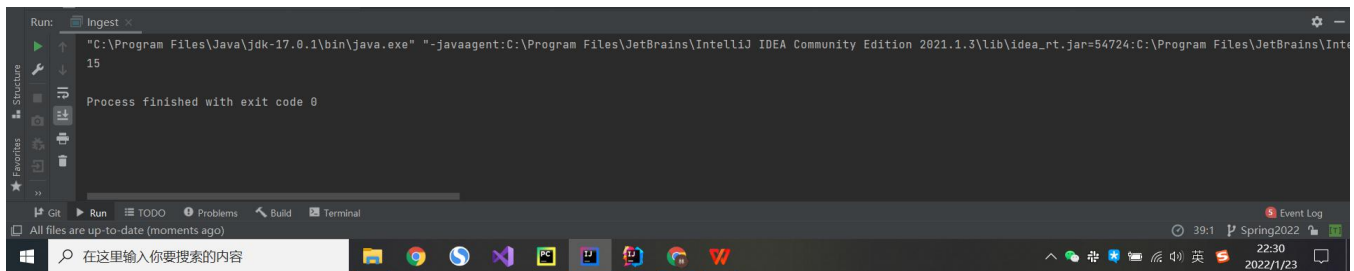


## 2. Run "HelloWord"



## 3. Run "Ingest.scala" and view the list of movies

## 4. Modify "Ingest.scala" and find number of movies of "New Zealand"



There are 15 movies are from "New Zealand" in the list.

## 5. Run simple spark program(word count) in Shell

**Project Source**

https://github.com/MrNiro/CSYE7200/tree/Spring2022/assignment-helloworld/src/main/scala/edu/neu/coe/csye7200/assthw