**Spark Project in Scala documentation**

I have chosen to use **IntelliJ IDEA** environment and reason is that I have been using this tools for 5 years. It is very easy to use and supports many languages such as Java, Scala, Kotlin, Node.js, Spring framework project and so on. So, let's get started.
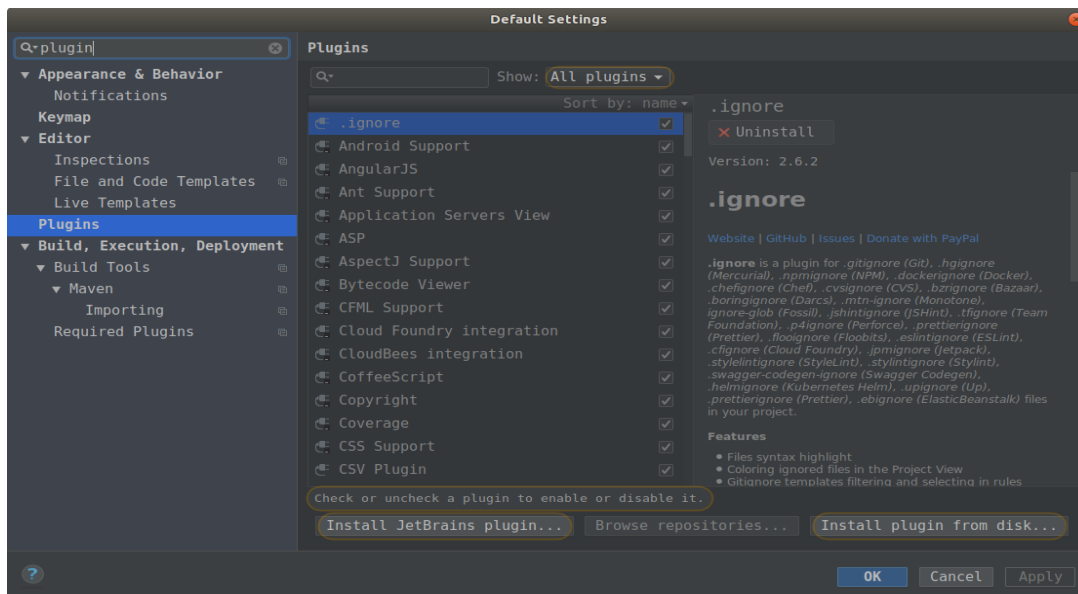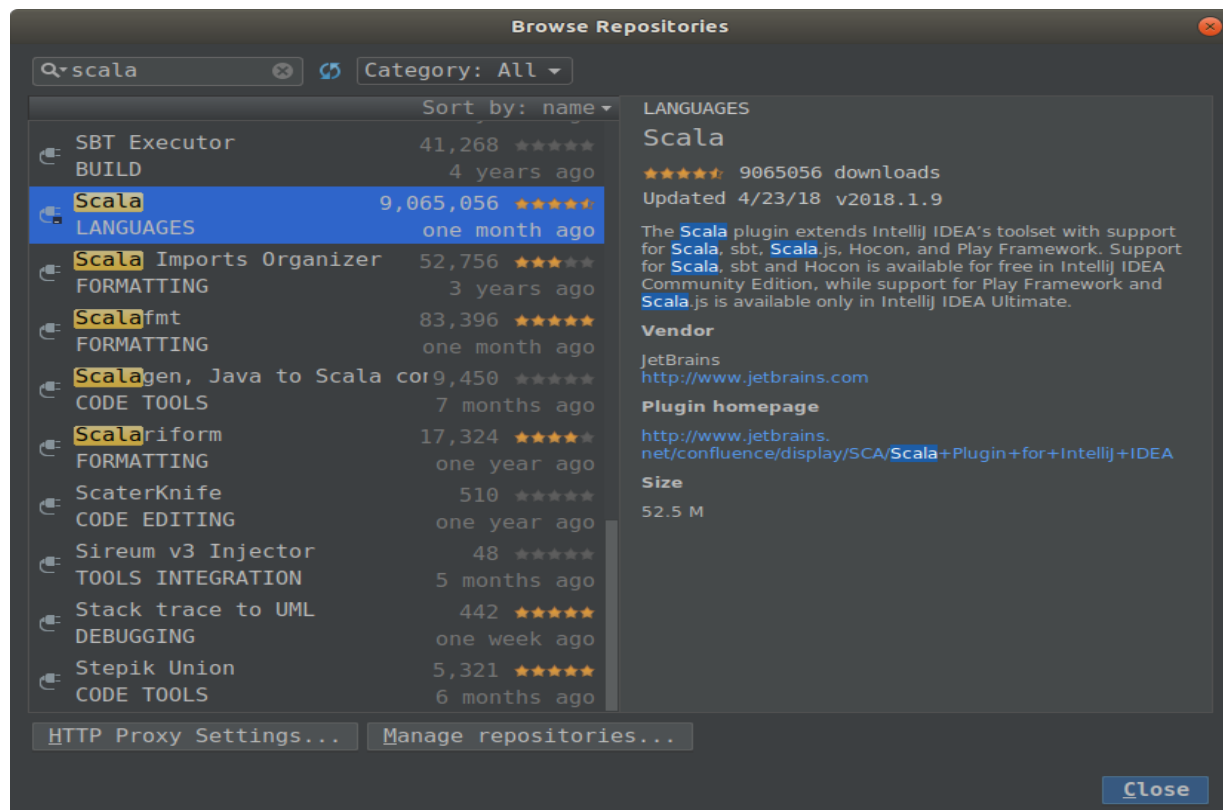First, let's open IntelliJ IDEA environment.



After opening the the tool, we need to configure IntelliJ IDEA in order to supports Scala project. So, click configure menu from the view.
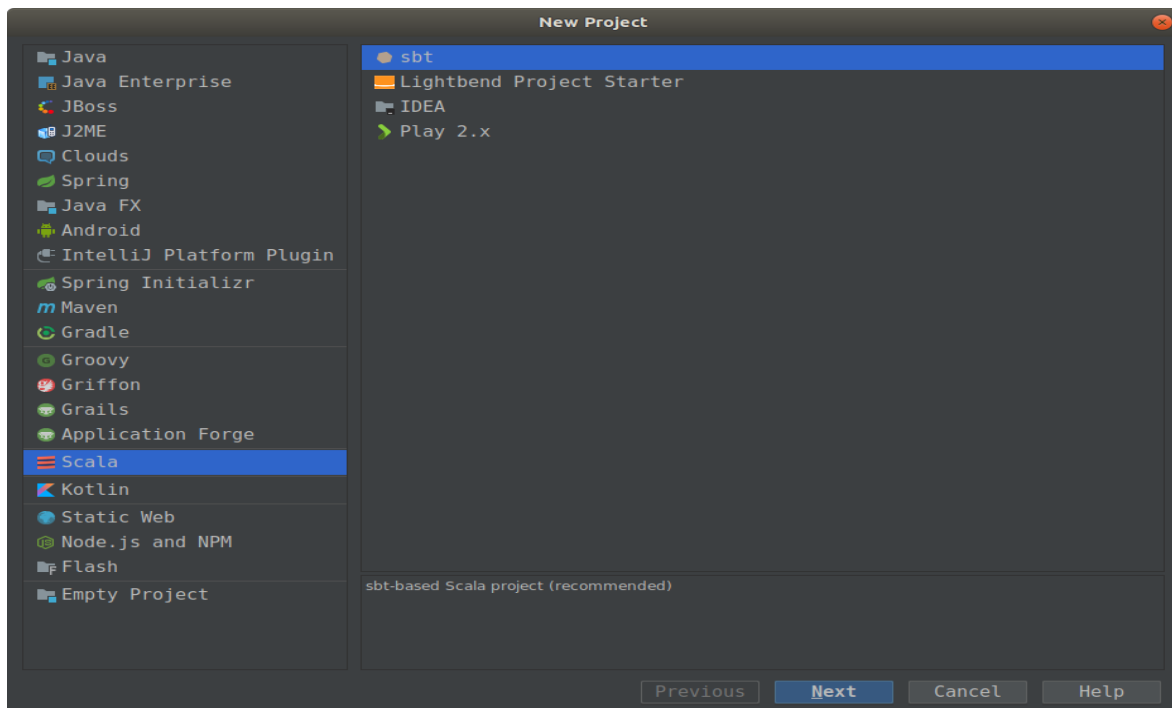
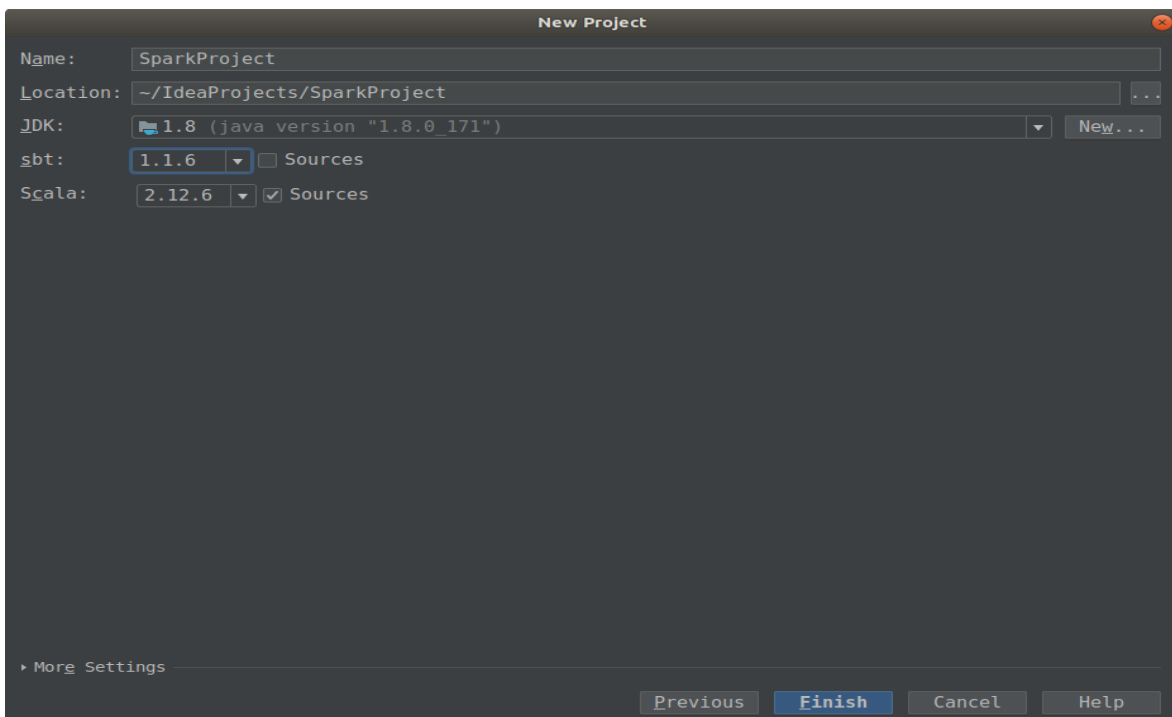Next step is to select plugins menu and type **Scala** keyword



In my machine it is already installed, if you haven't installed yet there is **install** to instal it. You just need to click it.
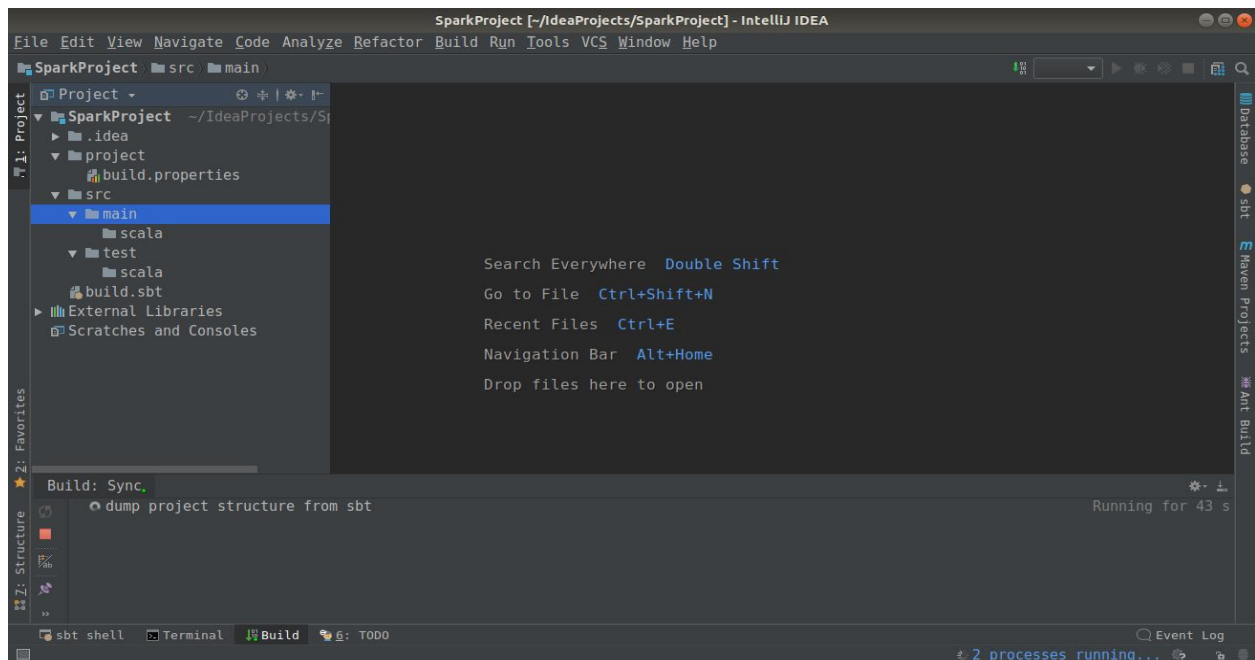
Now, everything is ready to start creating a project. Select Scala on the left menu and sbt on right menu as it is shown here.
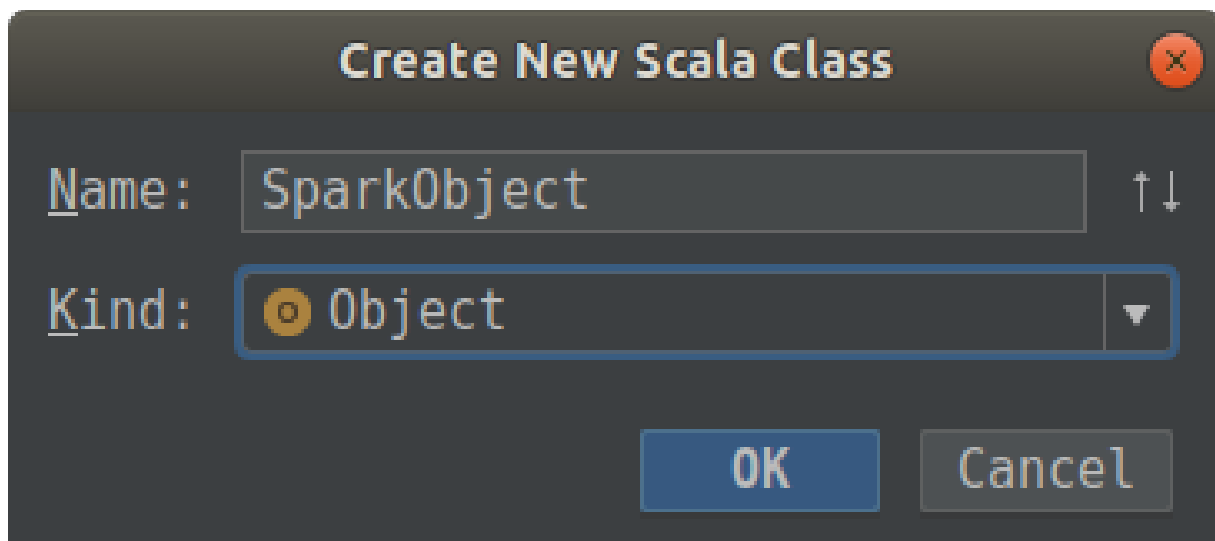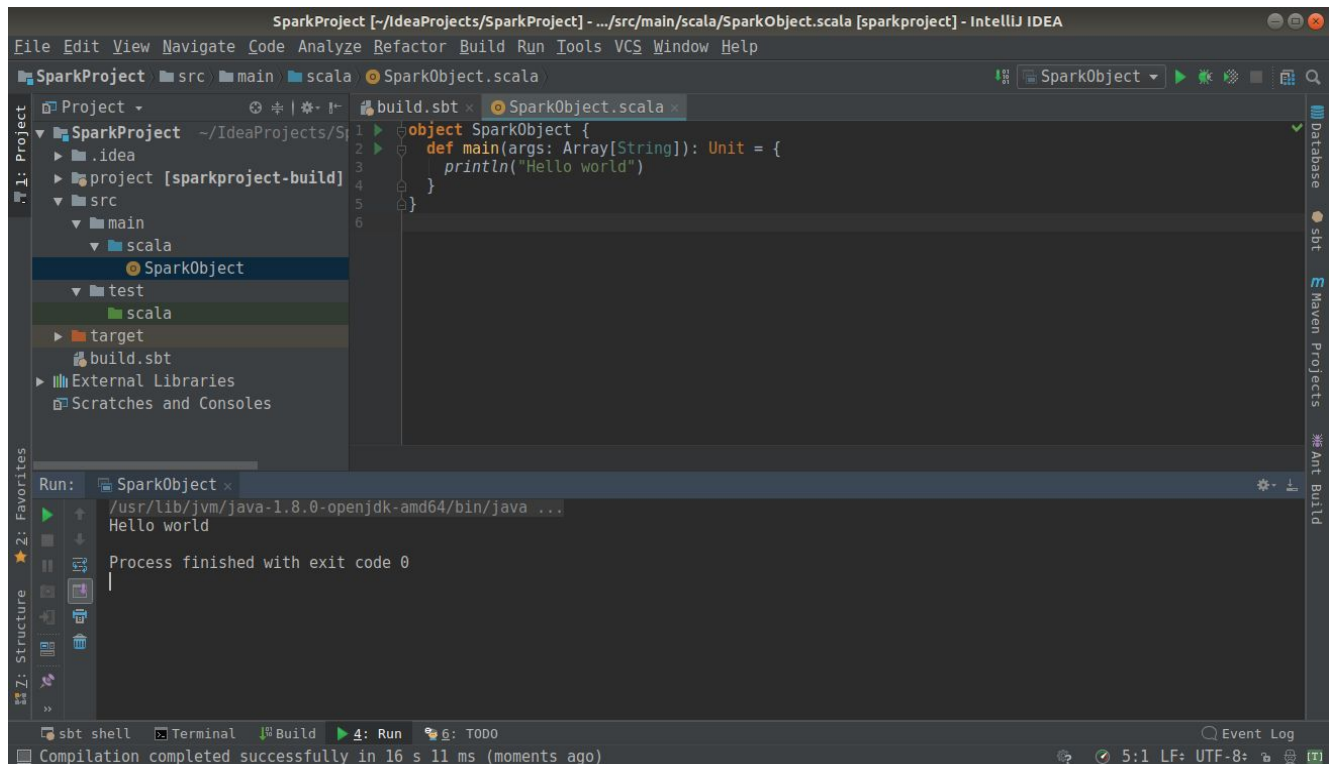


Select as it is shown here.

After then, this view shows up. Now, let's create test scala object and run it.



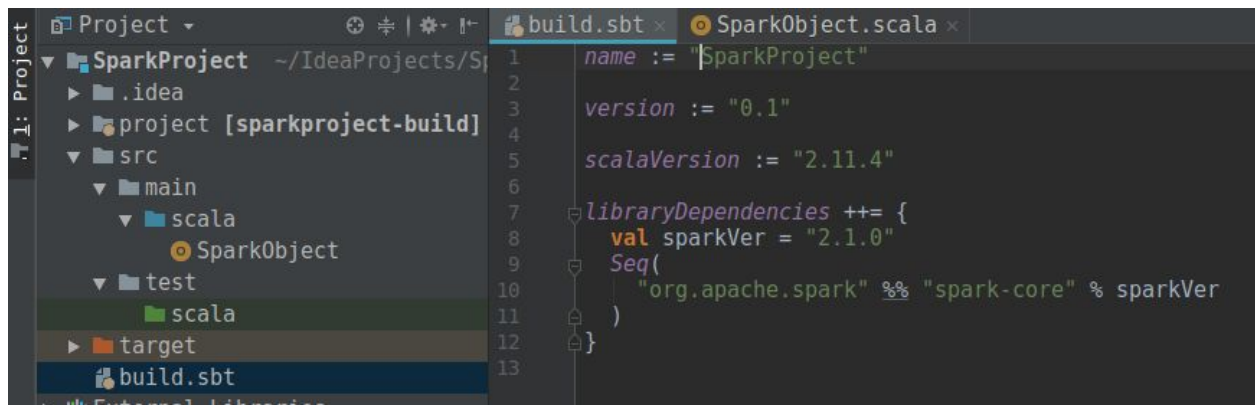Here we need to select to **Object** and click **OK**.

Here is sample example of Scala project. Work of the **SparkObject** is just printing **Hello World** as it is good example to begin learning a new language -:)). We can see the result from bottom view and it means scala project has successfully worked. The next step is to add **Spark** library to the project and use it.
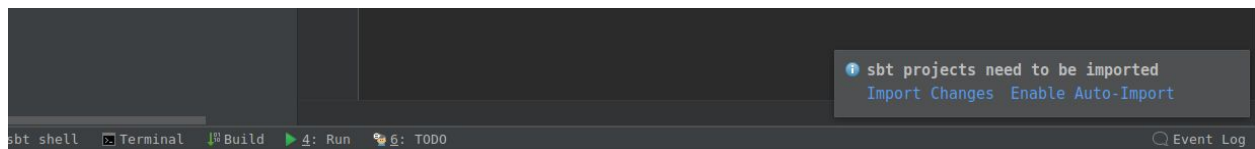


Here is sample code of the code of Spark. Spark's keywords are red as you can see. It means the project doesn't support **Spark** now. So let's add Spark library into the project.

We need to open **build.sbt** file to add Spark library. **Build.sbt** stores all the dependencies that are used in the project such as maven or gradle do.



Once we add the code of the library, it asks us if we want to import the library. We will choose yes.



After that we can import Spark codes.

It is the view when we import all the keywords and now we can test the code by running it.



Here is cvs sample file to check the project. We read the file, create a RDD, perform the output result. Let's run it.

SparkProject [~/IdeaProjects/SparkProject] - .../cane.csv [sparkproject] - IntelliJ IDEA

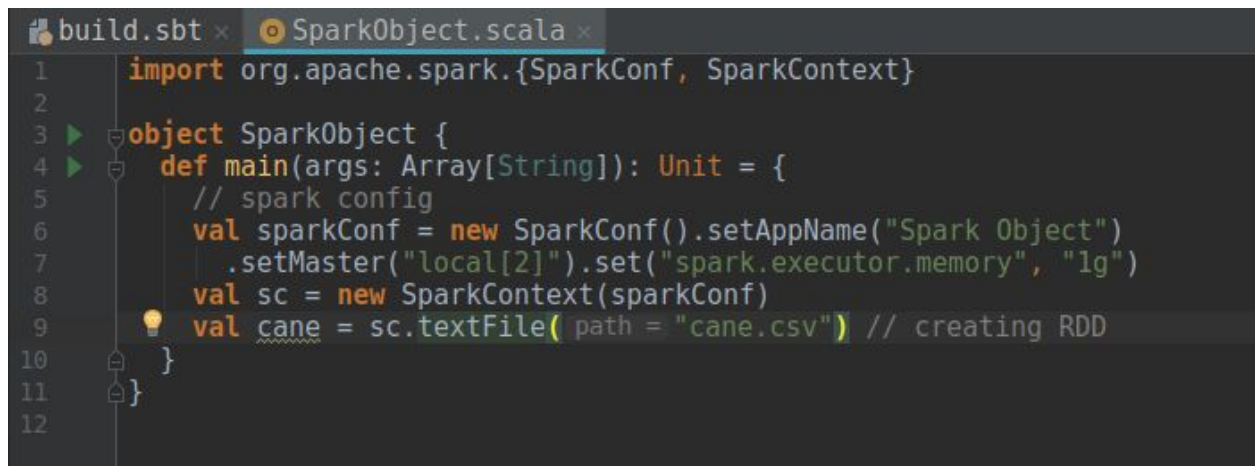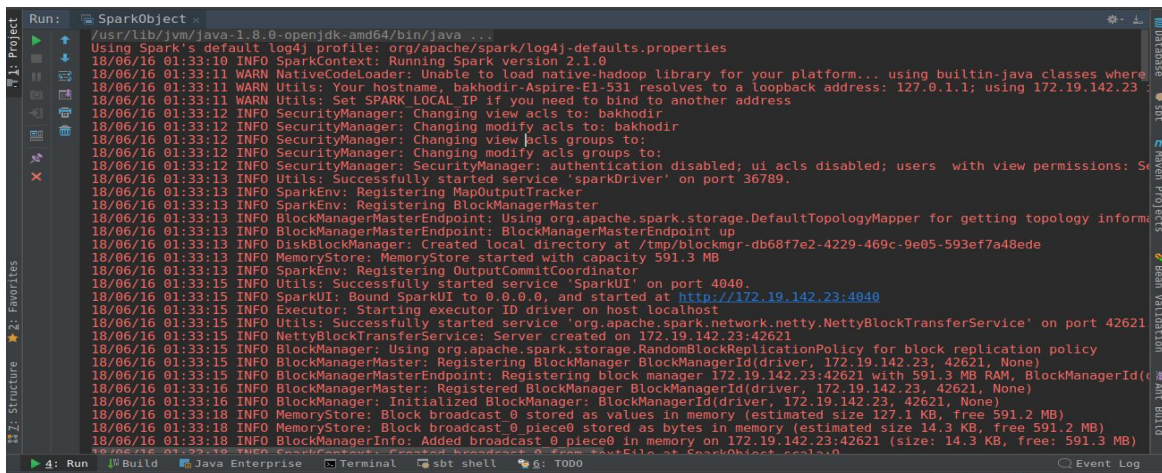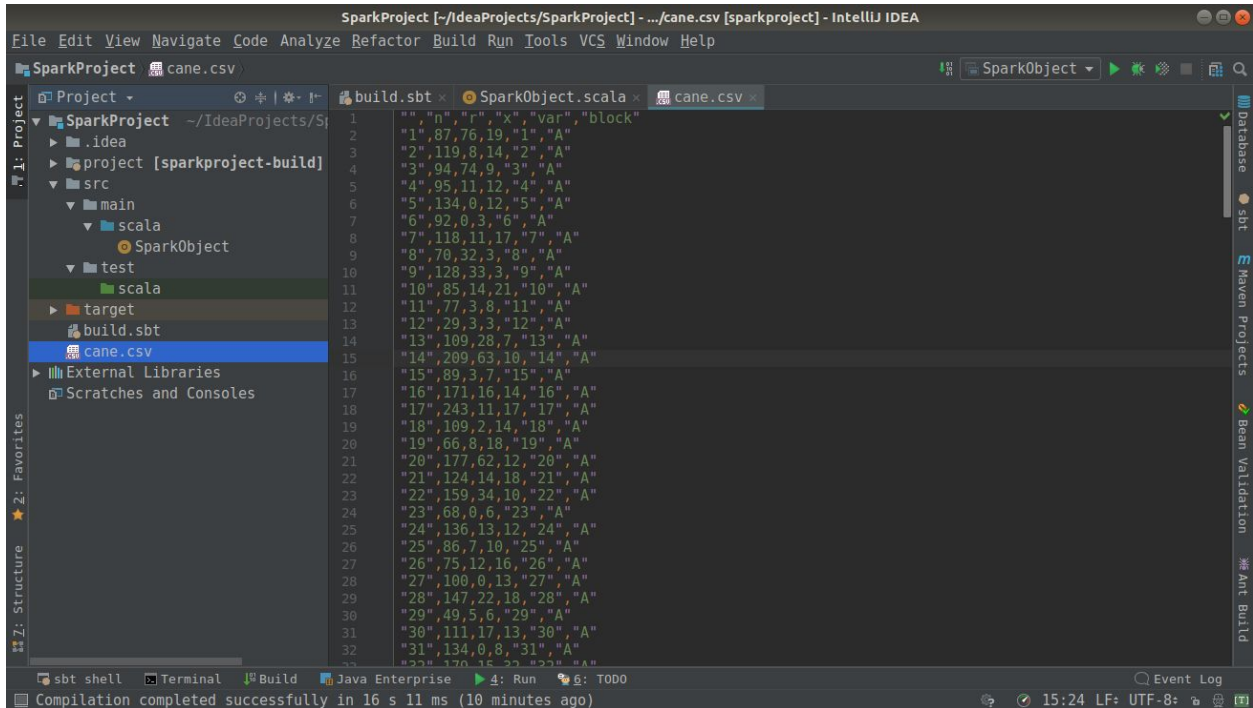File  Edit  View  Navigate  Code  Analyze  Refactor  Build  Run  Tools  VCS  Window  Help

SparkProject  cane.csv                                                    SparkObject

Project

SparkProject  ~/IdeaProjects/Sp
  .idea
  project [sparkproject-build]
  src
    main
      scala
        SparkObject
    test
      scala
  target
  build.sbt
  cane.csv
External Libraries
Scratches and Consoles

build.sbt    SparkObject.scala    cane.csv

```
1    "","n","r","x","var","block"
2    "1",87,76,19,"1","A"
3    "2",119,8,14,"2","A"
4    "3",94,74,9,"3","A"
5    "4",95,11,12,"4","A"
6    "5",134,0,12,"5","A"
7    "6",92,0,3,"6","A"
8    "7",118,11,17,"7","A"
9    "8",70,32,3,"8","A"
10   "9",128,33,3,"9","A"
11   "10",85,14,21,"10","A"
12   "11",77,3,8,"11","A"
13   "12",29,3,3,"12","A"
14   "13",109,28,7,"13","A"
15   "14",209,63,10,"14","A"
16   "15",89,3,7,"15","A"
17   "16",171,16,14,"16","A"
18   "17",243,11,17,"17","A"
19   "18",109,2,14,"18","A"
20   "19",66,8,18,"19","A"
21   "20",177,62,12,"20","A"
22   "21",124,14,18,"21","A"
23   "22",159,34,10,"22","A"
24   "23",68,0,6,"23","A"
25   "24",136,13,12,"24","A"
26   "25",86,7,10,"25","A"
27   "26",75,12,16,"26","A"
28   "27",100,0,13,"27","A"
29   "28",147,22,18,"28","A"
30   "29",49,5,6,"29","A"
31   "30",111,17,13,"30","A"
32   "31",134,0,8,"31","A"
```

sbt shell    Terminal    Build    Java Enterprise    4: Run    6: TODO    Event Log

Compilation completed successfully in 16 s 11 ms (10 minutes ago)    15:24  LF:  UTF-8:



Run:    SparkObject

```
/usr/lib/jvm/java-1.8.0-openjdk-amd64/bin/java ...
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
18/06/16 01:33:10 INFO SparkContext: Running Spark version 2.1.0
18/06/16 01:33:11 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where
18/06/16 01:33:11 WARN Utils: Your hostname, bakhodir-Aspire-E1-531 resolves to a loopback address: 127.0.1.1; using 172.19.142.23
18/06/16 01:33:11 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
18/06/16 01:33:12 INFO SecurityManager: Changing view acls to: bakhodir
18/06/16 01:33:12 INFO SecurityManager: Changing modify acls to: bakhodir
18/06/16 01:33:12 INFO SecurityManager: Changing view acls groups to:
18/06/16 01:33:12 INFO SecurityManager: Changing modify acls groups to:
18/06/16 01:33:12 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users  with view permissions: Se
18/06/16 01:33:13 INFO Utils: Successfully started service 'sparkDriver' on port 36789.
18/06/16 01:33:13 INFO SparkEnv: Registering MapOutputTracker
18/06/16 01:33:13 INFO SparkEnv: Registering BlockManagerMaster
18/06/16 01:33:13 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology informa
18/06/16 01:33:13 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
18/06/16 01:33:13 INFO DiskBlockManager: Created local directory at /tmp/blockmgr-db68f7e2-4229-469c-9e05-593ef7a48ede
18/06/16 01:33:13 INFO MemoryStore: MemoryStore started with capacity 591.3 MB
18/06/16 01:33:13 INFO SparkEnv: Registering OutputCommitCoordinator
18/06/16 01:33:15 INFO Utils: Successfully started service 'SparkUI' on port 4040.
18/06/16 01:33:15 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://172.19.142.23:4040
18/06/16 01:33:15 INFO Executor: Starting executor ID driver on host localhost
18/06/16 01:33:15 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 42621
18/06/16 01:33:15 INFO NettyBlockTransferService: Server created on 172.19.142.23:42621
18/06/16 01:33:15 INFO BlockManager: Using org.apache.spark.storage.RandomBlockReplicationPolicy for block replication policy
18/06/16 01:33:15 INFO BlockManagerMaster: Registering BlockManager BlockManagerId(driver, 172.19.142.23, 42621, None)
18/06/16 01:33:15 INFO BlockManagerMasterEndpoint: Registering block manager 172.19.142.23:42621 with 591.3 MB RAM, BlockManagerId(d
18/06/16 01:33:16 INFO BlockManagerMaster: Registered BlockManager BlockManagerId(driver, 172.19.142.23, 42621, None)
18/06/16 01:33:16 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, 172.19.142.23, 42621, None)
18/06/16 01:33:18 INFO MemoryStore: Block broadcast_0 stored as values in memory (estimated size 127.1 KB, free 591.2 MB)
18/06/16 01:33:18 INFO MemoryStore: Block broadcast_0 piece0 stored as bytes in memory (estimated size 14.3 KB, free 591.2 MB)
18/06/16 01:33:18 INFO BlockManagerInfo: Added broadcast_0 piece0 in memory on 172.19.142.23:42621 (size: 14.3 KB, free: 591.3 MB)
18/06/16 01:33:18 INFO SparkContext: Created broadcast 0 from textFile at SparkObject.scala:9
```

4: Run    Build    Java Enterprise    Terminal    sbt shell    6: TODO    Event Log

Everything is fine. We can do all the project requirements that is given. After being done, let's run the project and see the result.

And here is the result file.

## Conclusion

From doing the project I have learnt a lot of knowledge that are about Scala language, how it works briefly, what is Apache Spark open-source cluster-computing framework and integrating it with Scala. At first, I had a problem with importing Spark library inside Scala project and it took me much time and finally I realized that the problem was library version compatibility of Scala and Spark framework. It was really good experience doing the project and learnt new things.

Here, I have also provided github url to go my project. You can copy and work with it easily.
**https://github.com/bakhodir10/Spark-Scala**