

CASE STUDY 5

Task1: You have to create a Spark Application which streams data from a file on local directory on your machine and does the word count on the fly. The word should be done by the spark application in such a way that as soon as you drop the file in your local directory, your spark application should immediately do the word count for you.

Sol: Step1: open a scala project and create a scala application

➔ Program for scala application

```
import org.apache.spark.{SparkConf, SparkContext}
import org.apache.spark.rdd.RDD
import org.apache.spark.streaming.{Seconds, StreamingContext, Time}
import org.apache.spark.sql.SparkSession
import org.apache.log4j.{Level, Logger}

object SqlNetworkWordCount {

  def main(args: Array[String]): Unit = {
    println("hey Spark SQL Streaming")
    val conf = new SparkConf().setMaster("local[2]").setAppName("SparkSteamingExample")
    val sc = new SparkContext(conf)
    val rootLogger = Logger.getRootLogger()
    rootLogger.setLevel(Level.ERROR)
    println("hey Spark Streaming ---> 1")
    //val sparkConf = new SparkConf().setAppName("NetworkWordCount")
    println("hey Spark Streaming ---> 2")
    val ssc = new StreamingContext(sc, Seconds(10))
    val lines = ssc.socketTextStream("localhost", 9999)
    println("hey Spark Streaming ---> 3")
    val words = lines.flatMap(_.split(" "))

    // Convert RDDs of the words DStream to DataFrame and run SQL query
    words.foreachRDD { (rdd: RDD[String], time: Time) =>
      val spark = SparkSessionSingleton.getInstance(rdd.sparkContext.getConf)
      import spark.implicits._

      // Convert RDD[String] to RDD[case class] to DataFrame
      val wordsDataFrame = rdd.map(w => Record(w)).toDF()

      // Creates a temporary view using the DataFrame
      wordsDataFrame.createOrReplaceTempView("words")
      // Do word count on table using SQL and print it
      val wordCountsDataFrame =
        spark.sql("select word, count(*) as total from words group by word")
    }
  }
}
```

```

println(s"===== $time =====")
wordCountsDataFrame.show()
}
ssc.start()
ssc.awaitTermination()
}
/** Case class for converting RDD to DataFrame */
case class Record(word: String)
/** Lazily instantiated singleton instance of SparkSession */
object SparkSessionSingleton {
  @transient private var instance: SparkSession = _
  def getInstance(sparkConf: SparkConf): SparkSession = {
    if (instance == null) {
      instance = SparkSession
        .builder
        .config(sparkConf)
        .getOrCreate()
    }
    instance
  }
}

```

```

eclipse-workspace - CaseStudyV/src/SparkFileStreamingWordCount.scala - Eclipse
File Edit Refactor Navigate Search Project Scala Run Window Help
Quick Access

Package Explorer
CaseStudyV
├─ Scala Library container [ 2.
├─ JRE System Library [JavaSt
├─ src
│   └─ (default package)
│       └─ SparkFileStreamingW
├─ Referenced Libraries
├─ SparkHiveInt
└─ SparkKafka

SparkFileStreamingWordCount.scala
1 import org.apache.spark.{SparkConf, SparkContext}
2 import org.apache.spark.rdd.RDD
3 import org.apache.spark.streaming.{Seconds, StreamingContext, Time}
4 import org.apache.spark.sql.SparkSession
5 import org.apache.log4j.{Level, Logger}
6
7 object SqlNetworkWordCount {
8
9   def main(args: Array[String]): Unit = {
10     println("hey Spark SQL Streaming")
11     val conf = new SparkConf().setMaster("local[2]").setAppName("SparkSteamingE
12     val sc = new SparkContext(conf)
13     val rootLogger = Logger.getRootLogger()
14     rootLogger.setLevel(Level.ERROR)
15     println("hey Spark Streaming ---> 1")
16     //val sparkConf = new SparkConf().setAppName("NetworkWordCount")
17     println("hey Spark Streaming ---> 2")
18     val ssc = new StreamingContext(sc, Seconds(10))
19     val lines = ssc.socketTextStream("localhost", 9999)
20     println("hey Spark Streaming ---> 3")
21     val words = lines.flatMap(_.split(" "))
22
23     // Convert RDDs of the words DStream to DataFrame and run SQL query
24     words.foreachRDD { (rdd: RDD[String], time: Time) =>
25       val spark = SparkSessionSingleton.getInstance(rdd.sparkContext.getConf)
26       import spark.implicits._
27
28       // Convert RDD[String] to RDD[case class] to DataFrame
29       val wordsDataFrame = rdd.map(w => Record(w)).toDF()
30

```

Quick Access

SparkFileStreamingWordCount.scala

```
30
31 // Creates a temporary view using the DataFrame
32 wordsDataFrame.createOrReplaceTempView("words")
33 // Do word count on table using SQL and print it
34 val wordCountsDataFrame =
35     spark.sql("select word, count(*) as total from words group by word")
36     println(s"===== $time =====")
37     wordCountsDataFrame.show()
38 }
39 ssc.start()
40 ssc.awaitTermination()
41 }
42
43
44
45 case class Record(word: String)
46 /** Lazily instantiated singleton instance of SparkSession */
47 object SparkSessionSingleton {
48     @transient private var instance: SparkSession =
49     def getInstance(sparkConf: SparkConf): SparkSession = {
50         if (instance == null) {
51             instance = SparkSession
52                 .builder
53                 .config(sparkConf)
54                 .getOrCreate()
55         }
56         instance
57     } }
58
59 }
```

Step2: Start nc server

➔ nc -lk 9999

```
ns_err_pid4625.log test.txt
install Videos
jarKafkaUtils
[acadgild@localhost ~]$ nc -lk 9999
hey bigdata
bigdata word word
hey hey tech tech
acadgild
acadgild acadgild
^C
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost ~]$
```

Step3: Run the program

eclipse-workspace - CaseStudyV/src/SparkFileStreamingWordCount.scala - Eclipse

File Edit Refactor Navigate Search Project Scala Run Window Help

Problems Tasks Console

```
<terminated> SqlNetworkWordCount$ [Scala Application] /usr/java/jdk1.8.0_151/bin/java (May 28, 2018, 8
18/05/28 08:56:40 INFO BlockManagerMasterEndpoint: Registering block manager 192.168.
18/05/28 08:56:40 INFO BlockManagerMaster: Registered BlockManager BlockManagerId(dri
18/05/28 08:56:40 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver,
hey Spark Streaming ----> 1
hey Spark Streaming ----> 2
hey Spark Streaming ----> 3
===== 1527478010000 ms =====
+-----+
| word|total|
+-----+
| word|    2|
|bigdata|  2|
|  hey|   3|
|  tech|   2|
+-----+

===== 1527478020000 ms =====
+-----+
|word|total|
```

Quick Access

Session */

on = {

eclipse-workspace - CaseStudyV/src/SparkFileStreamingWordCount.scala - Eclipse

File Edit Refactor Navigate Search Project Scala Run Window Help

Problems Tasks Console

```
<terminated> SqlNetworkWordCount$ [Scala Application] /usr/java/jdk1.8.0_151/bin/java (May 28, 2018, 8
|word|total|
+-----+
+-----+

===== 1527478070000 ms =====
+-----+
|word|total|
+-----+
+-----+

===== 1527478080000 ms =====
+-----+
| word|total|
+-----+
|acadgild|  3|
+-----+

===== 1527478090000 ms =====
+-----+
+-----+
```

Quick Access

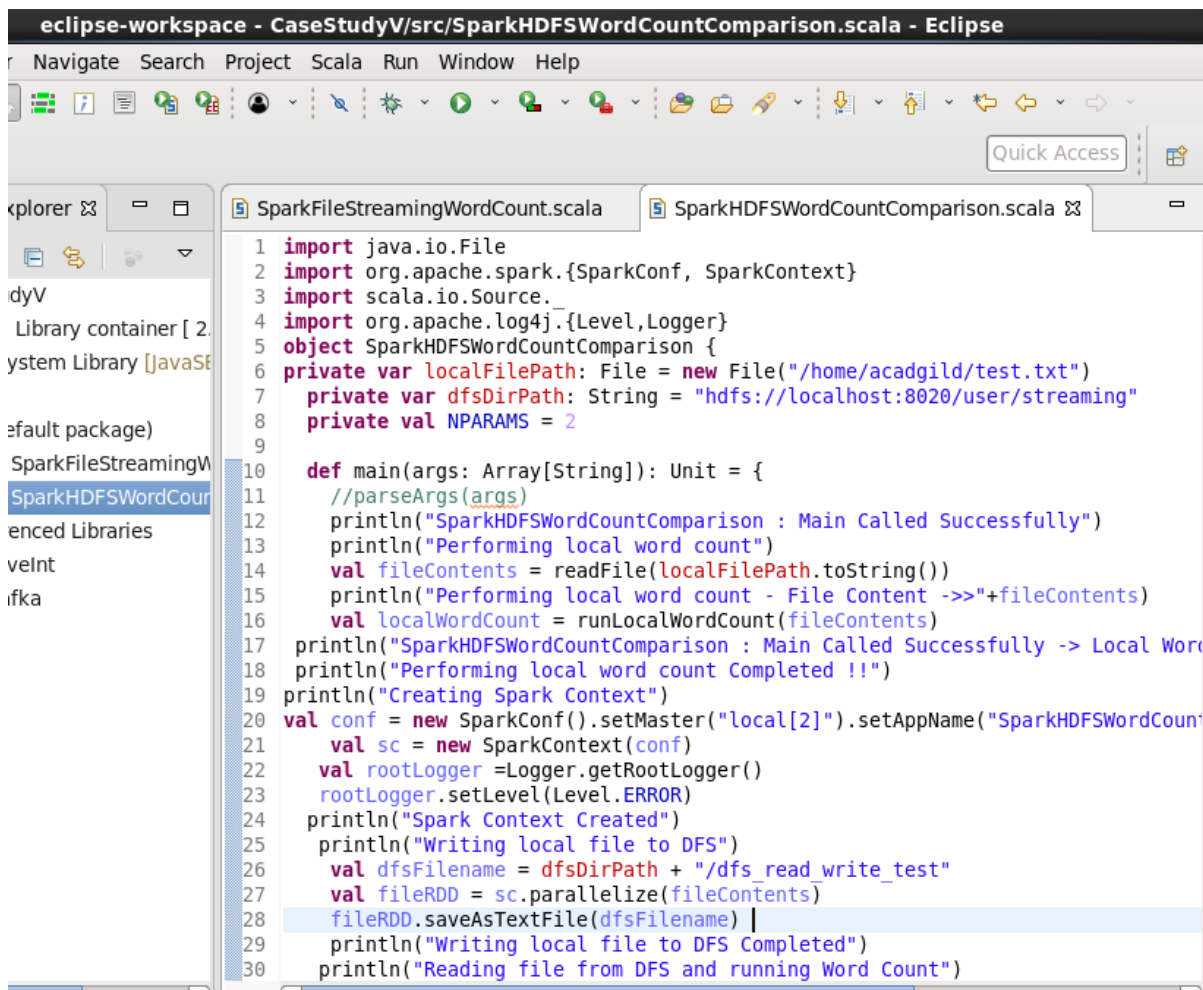
Session */

on = {

Task2: In this part, you will have to create a Spark Application which should do the following

1. Pick up a file from the local directory and do the word count
2. Then in the same Spark Application, write the code to put the same file on HDFS.
3. Then in same Spark Application, do the word count of the file copied on HDFS in step 2
4. Lastly, compare the word count of step 1 and 2. Both should match, other throw an error

Sol: Step1: create a scala application in the previous scala project



```
eclipse-workspace - CaseStudyV/src/SparkHDFSWordCountComparison.scala - Eclipse
Navigate Search Project Scala Run Window Help
Quick Access
SparkFileStreamingWordCount.scala SparkHDFSWordCountComparison.scala
1 import java.io.File
2 import org.apache.spark.{SparkConf, SparkContext}
3 import scala.io.Source
4 import org.apache.log4j._.{Level, Logger}
5 object SparkHDFSWordCountComparison {
6 private var localFilePath: File = new File("/home/acadgild/test.txt")
7 private var dfsDirPath: String = "hdfs://localhost:8020/user/streaming"
8 private val NPARAMS = 2
9
10 def main(args: Array[String]): Unit = {
11 //parseArgs(args)
12 println("SparkHDFSWordCountComparison : Main Called Successfully")
13 println("Performing local word count")
14 val fileContents = readFile(localFilePath.toString())
15 println("Performing local word count - File Content ->" + fileContents)
16 val localWordCount = runLocalWordCount(fileContents)
17 println("SparkHDFSWordCountComparison : Main Called Successfully -> Local Word Count")
18 println("Performing local word count Completed !!")
19 println("Creating Spark Context")
20 val conf = new SparkConf().setMaster("local[2]").setAppName("SparkHDFSWordCountComparison")
21 val sc = new SparkContext(conf)
22 val rootLogger = Logger.getRootLogger()
23 rootLogger.setLevel(Level.ERROR)
24 println("Spark Context Created")
25 println("Writing local file to DFS")
26 val dfsFilename = dfsDirPath + "/dfs_read_write_test"
27 val fileRDD = sc.parallelize(fileContents)
28 fileRDD.saveAsTextFile(dfsFilename)
29 println("Writing local file to DFS Completed")
30 println("Reading file from DFS and running Word Count")
```


eclipse-workspace - CaseStudyV/src/SparkHDFSWordCountComparison.scala - Eclipse

File Edit Refactor Navigate Search Project Scala Run Window Help

Problems Tasks Console

```
<terminated> SparkHDFSWordCountComparison$ [Scala Application] /usr/java/jdk1.8.0_151/bin/java (May 18/05/28 09:16:41 INFO DiskBlockManager: Created local directory at /tmp/blockmgr-ffe
18/05/28 09:16:41 INFO MemoryStore: MemoryStore started with capacity 111.2 MB
18/05/28 09:16:41 INFO SparkEnv: Registering OutputCommitCoordinator
18/05/28 09:16:43 INFO Utils: Successfully started service 'SparkUI' on port 4040.
18/05/28 09:16:43 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://192.168.0.7:4040/
18/05/28 09:16:44 INFO Executor: Starting executor ID driver on host localhost
18/05/28 09:16:44 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 41337.
18/05/28 09:16:44 INFO NettyBlockTransferService: Server created on 192.168.0.7:41337
18/05/28 09:16:44 INFO BlockManager: Using org.apache.spark.storage.RandomBlockReplicationLevel2 of 2 blocks
18/05/28 09:16:44 INFO BlockManagerMaster: Registering BlockManager BlockManagerId(driver, 192.168.0.7, 41337)
18/05/28 09:16:44 INFO BlockManagerMasterEndpoint: Registering block manager 192.168.0.7:41337 with 2 blocks
18/05/28 09:16:44 INFO BlockManagerMaster: Registered BlockManager BlockManagerId(driver, 192.168.0.7, 41337)
18/05/28 09:16:44 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, 192.168.0.7, 41337)
Spark Context Created
Writing local file to DFS
Writing local file to DFS Completed
Reading file from DFS and running Word Count
Success! Local Word Count (44) and DFS Word Count (44) agree.
```

gild/test
8020/user/
ed Success
g()
nt ->"+fi
ts)
Successful
ppName("Sp