# CSEE5590 Big Data Programming

**In Class Programming –8  Report**
**(Jongkook Son)**

## Project Overview:

Apache Spark is a unified analytics engine for big data processing, with built-in modules for streaming, SQL, machine learning and graph processing.

## Requirements/Task(s):

1. Spark Programming:

Write a spark program with an interesting use case using text data as the input and program should have at least Two Spark Transformations and Two Spark Actions.

Present your use case in map reduce paradigm as shown below (for word count).

2. Secondary Sorting in Map Reduce

Secondary sorting is used to sort the values in the reducer phase.

Take any input of your interest and perform secondary sorting on it.

## What I learned in ICP:

I could have learned how to set up Spark on my local machine and installing scala as a plugin for IntelliJ. It was a little tricky to install spark in my local but I did it by following steps. By fulfilling the first task, I could have learned some transformations and actions in spark and how can I implement it with scala. It seems better than I did on Hadoop. By Fulfilling the second task, I could have understood what secondary sort in map-reduce. The secondary sorting technique will enable us to sort the values (in ascending or descending order) passed to each reducer. The goal of the Secondary Sort pattern is to give some order to the values received by a reducer.

# INSTALLATION OF Spark AND Scala

**(I installed it on my local maschine and intellij. )**

**&lt;Task1&gt; Write a spark program with an interesting use case using text data as the input and program should have at least Two Spark Transformations and Two Spark Actions.**

=> Here I used transformations map,flatmap,reducebykey,filter and in actions I used count, foreach, take.



&lt;Overall Code&gt;



&lt;Transformation1&gt; FlatMap

```
21/03/16 19:43:03 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 5 ms
(are,3)
(you,2)
(am,3)
(hi,1)
(i,3)
(how,1)
(very,1)
(hello,2)
(good,4)
(so,1)
```

<Transformation2> Map and Reduce

```
21/03/16 19:43:03 INFO DAGScheduler: ResultStage 5 (foreach at WordCount.scala:29) finished in 0.017 s
21/03/16 19:43:03 INFO DAGScheduler: Job 2 finished: foreach at WordCount.scala:29, took 0.044575 s
(am,3)
(are,3)
(good,4)
(hello,2)
(hi,1)
(how,1)
(i,3)
(so,1)
(very,1)
(you,2)
```

<Transformation3> Sort

```
21/03/16 19:43:03 INFO ShuffleBlockFetcherIterator: Getting 1 non-empty blocks out of 1 blocks
21/03/16 19:43:03 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 0 ms
(am,3)
(are,3)
(good,4)
(hello,2)
(i,3)
21/03/16 19:43:03 INFO Executor: Finished task 0.0 in stage 17.0 (TID 10). 1090 bytes result sent to driver
21/03/16 19:43:03 INFO TaskSetManager: Finished task 0.0 in stage 17.0 (TID 10) in 5 ms on localhost (executor driver) (1/1)
21/03/16 19:43:03 INFO TaskSchedulerImpl: Removed TaskSet 17.0, whose tasks have all completed, from pool
21/03/16 19:43:03 INFO DAGScheduler: ResultStage 17 (count at WordCount.scala:41) finished in 0.008 s
21/03/16 19:43:03 INFO DAGScheduler: Job 6 finished: count at WordCount.scala:41, took 0.009907 s
21/03/16 19:43:03 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-IHI3GNU.mshome.net:4040
21/03/16 19:43:03 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
(Count Unique words:,6)
21/03/16 19:43:03 INFO MemoryStore: MemoryStore cleared
21/03/16 19:43:03 INFO BlockManager: BlockManager stopped
21/03/16 19:43:03 INFO BlockManagerMaster: BlockManagerMaster stopped
21/03/16 19:43:03 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
```
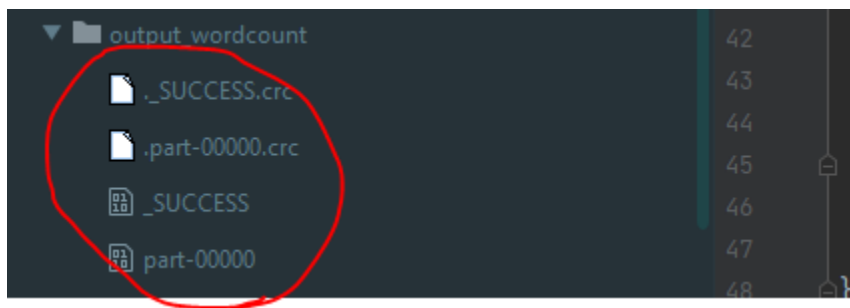
<Transformation4> Fliter

```scala
//Action 1
wordsFilter.take( num = 5).foreach(outputLIst=>println(outputLIst))
//Action 2
println("Count Unique words:",wordsFilter.count())
```

<Actions>

=>I take 5 out for the filtered list and count the number of words filter

<Output of the Task>

<Task2>Performing the partitioning of the map values which will do group by with comparing the common values as a KEY and the tuple values will be taken as VAlUES

```scala
package com.jk.spark.SecondTask

import org.apache.spark.{HashPartitioner, SparkConf, SparkContext}


object SortTemperatureByDate {

  // function to parse the input file
  def parseLine(line: String): ((String, Int), Int) = {
    val fields = line.split( regex = ",")
    val year = fields(0).toInt
    val month = fields(1).toInt
    val day = fields(2).toInt
    val temperature = fields(3).toInt

    val k = (year + "-" + month)
    val k2 = (k, temperature)
    (k2, temperature)
  }
}
```

```scala
def main(args: Array[String]): Unit = {
    // Create a Scala Spark Configuration.
    val conf = new SparkConf().setMaster("local[*]").setAppName("Sort Temperature By Date")

    // Create a Scala Spark Context.
    val sc = new SparkContext(conf)

    // Turn off all the warnings but ERROR
    sc.setLogLevel("ERROR")

    // Load our input data.
    val lines = sc.textFile( path = "temperature.txt")

    // Split up lines into key and value pairs
    val rdd = lines.map(parseLine)

    // To partition the rdd into 2 clusters
    val partitionedRDD = rdd.partitionBy(new HashPartitioner( partitions = 2))
    println("partitionedRDD")
    partitionedRDD.foreach {println}

    //  To map the rdd into key, value pairs and grouping values by date (key) and sorting the values after putting them into List
    val listRDD = partitionedRDD.map(l => (l._1._1, l._2)).groupByKey().mapValues(k => k.toList.sortBy(r => r))

    // printing on the console
    println("listRDD")
    listRDD.collect().foreach(println)

    // Save the output back out to a text file, causing evaluation.
    listRDD.saveAsTextFile( path = "output_secondarySorting")
    sc.stop()

}
```

<Overall Code>

```
21/03/16 19:52:09 INFO BlockManager: Initialized BlockManager: BlockManagerId(driver, DESKTOP-IHI3GNU.mshome.net, 57576, None)
partitionedRDD
((2000-11,20),20)
((2000-12,10),10)
((2000-12,-20),-20)
((2000-11,30),30)
((2000-11,-40),-40)
```
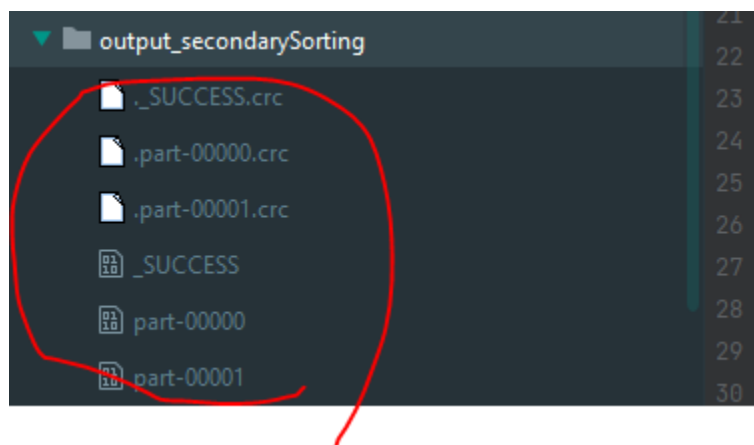
<Partitioned Rdd>

```
listRDD
(2000-12,List(-20, 10))
(2000-11,List(-40, 20, 30))

Process finished with exit code 0
```

<List RDD >

output_secondarySorting
    ._SUCCESS.crc
    .part-00000.crc
    .part-00001.crc
    _SUCCESS
    part-00000
    part-00001

21
22
23
24
25
26
27
28
29
30

<Output for Task2>