1. Create a Spark project, to perform the data statistics about 10 abstracts.

   Use at least two transformations and actions.

   Input: Choose 10 of abstracts from your dataset.

   Data Statistics Tasks:

   a. Word Count Statistics
   b. Parts of Speech Statistics
   c. WordNet Statistics
   d. Medical Word Statistics

Project ▼

SparkWordCount.scala ×    build.sbt ×

```
Spark WordCount [SparkTransformationA
  .idea
  output
  project [spark-wordcount-build] source
  src
    main
      scala
        SparkWordCount
  target
  build.sbt
  input
External Libraries
Scratches and Consoles
```

```scala
1
2
3    import org.apache.spark.{SparkContext, SparkConf}
4
5    /**
6      * Created by Mayanka on 09-Sep-15.
7      */
8    object SparkWordCount {
9
10     def main(args: Array[String]) :Unit = {
11
12       //System.setProperty("hadoop.home.dir","D:\\winutils");
13
14       val sparkConf = new SparkConf().setAppName("SparkWordCount").setMaster("local[*]")
15
16       val sc=new SparkContext(sparkConf)
17
18       val input=sc.textFile( path = "input")
19
20       val wc=input.flatMap(line=>{line.split( regex = " ")}).map(word=>(word,1)).cache()
21
22       val output=wc.reduceByKey(_+_)
23
24       output.saveAsTextFile( path = "output")
25
26       val o=output.collect()
27
28       var s:String="Words:Count \n"
```

SparkWordCount

Run:    SparkWordCount ×

```
18/10/22 13:07:36 INFO DAGScheduler: Job 1 finished: collect at SparkWordCount.scala:26, took 0.037566 s
18/10/22 13:07:36 INFO SparkContext: Invoking stop() from shutdown hook
18/10/22 13:07:36 INFO SparkUI: Stopped Spark web UI at http://lakshmis-air:4040
18/10/22 13:07:36 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/10/22 13:07:36 INFO MemoryStore: MemoryStore cleared
18/10/22 13:07:36 INFO BlockManager: BlockManager stopped
18/10/22 13:07:36 INFO BlockManagerMaster: BlockManagerMaster stopped
18/10/22 13:07:36 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/10/22 13:07:36 INFO SparkContext: Successfully stopped SparkContext
18/10/22 13:07:36 INFO ShutdownHookManager: Shutdown hook called
18/10/22 13:07:36 INFO ShutdownHookManager: Deleting directory /private/var/folders/_j/z5xr05fj1sz7w5y9bsqydh4m0000gn/T/spark-0536ca81-1b89-4cff-8f5c-c3defba34d5d

Process finished with exit code 0
```