

CS5542 Big Data Apps and Analytics

LAB ASSIGNMENT #2

Name: Sidrah Junaid

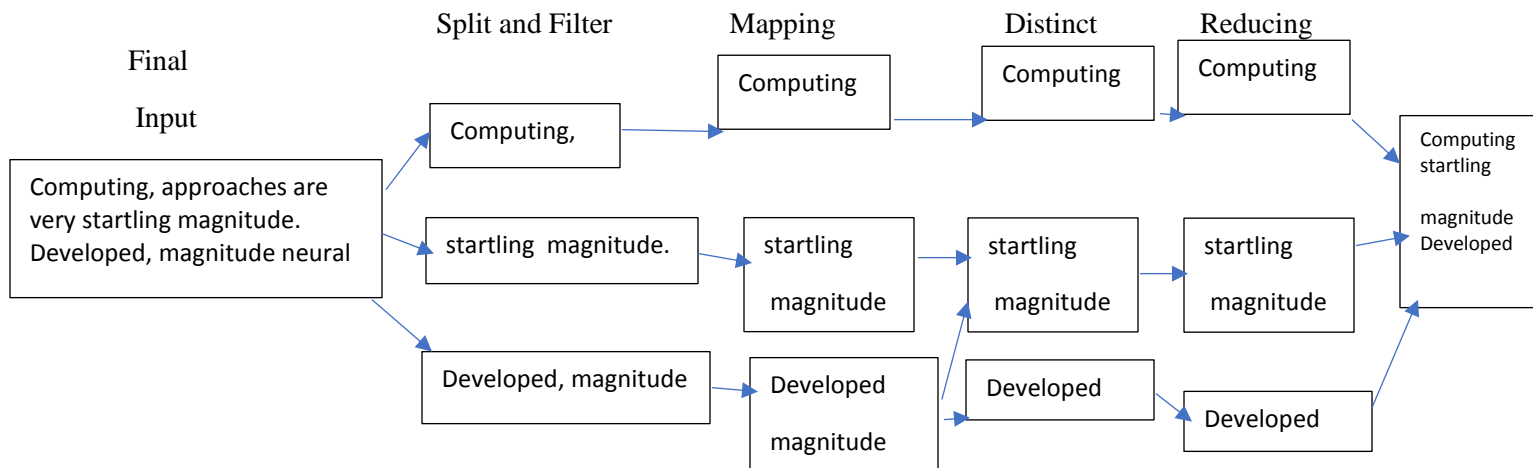
Class Id: 15

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).

USE CASE:



USE CASE IMPLEMENTATION:

The implemented use case shows the word having length of nine alphabets from the text present in input file.

1. The spark program is written in Scala using IntelliJ IDE and it is showing the words having 9 alphabets. The provided input has whitespaces and contain punctuation marks.
2. Program goes line by line and split the words based on empty characters and punctuation marks.
3. Then it will filter the words whose length is nine. In the provided text there is some redundant words from which only distinct word selected.
4. Then, the whole output would be collected and select in a separate output text file.

TRANSFORMATIONS:

1.Flat map:

Each line in the input file will be checked and split based on whitespaces and punctuation mark.

2.Distinct:

Here we select distinct words to avoid repetition.

3.Filter:

Here filter is applied so the words only having length 9 will be saved in main_filter.

Actions:

1.Collect:

Return all the elements of the input file as an array at the driver program.

2.saveAsTextFile:

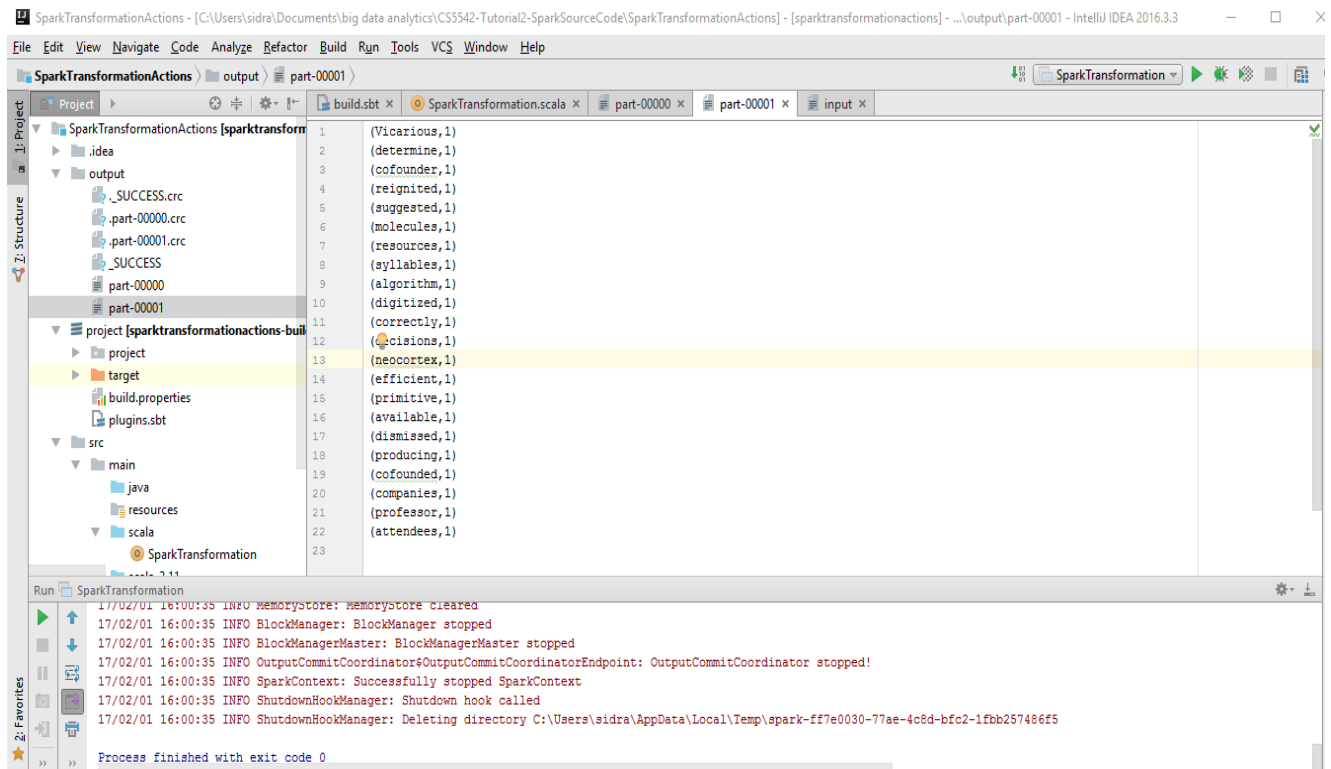
Write the elements of the input file as a text file (or set of text files) in a given directory in the local filesystem

Output:

File one:

```
SparkTransformationActions - [C:\Users\sidra\Documents\big data analytics\CS5542-Tutorial2-SparkSourceCode\SparkTransformationActions] - [sparktransformationactions] - ...\\output\part-00000 - IntelliJ IDEA 2016.3.3
File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window Help
SparkTransformationActions > output > part-00000
Project: SparkTransformationActions [sparktransform]
  .idea
  output
    _SUCCESS.crc
    .part-00000.crc
    .part-00001.crc
    _SUCCESS
    part-00000
    part-00001
  project [sparktransformationactions-build]
    project
    target
    build.properties
    plugins.sbt
  src
    main
      java
      resources
      scala
        SparkTransformation
Run: SparkTransformation
17/02/01 16:00:35 INFO MemoryStore: MemoryStore created
17/02/01 16:00:35 INFO BlockManager: BlockManager stopped
17/02/01 16:00:35 INFO BlockManagerMaster: BlockManagerMaster stopped
17/02/01 16:00:35 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
17/02/01 16:00:35 INFO SparkContext: Successfully stopped SparkContext
17/02/01 16:00:35 INFO ShutdownHookManager: Shutdown hook called
17/02/01 16:00:35 INFO ShutdownHookManager: Deleting directory C:\Users\sidra\AppData\Local\Temp\spark-ff7e0030-77ae-4c8d-bfc2-1fbb257486f5
Process finished with exit code 0
```

File two:



Input file:

