1. Generate the output (changes or transformations in the data) manually when the following tasks are applied on the input text. Show your output in details.
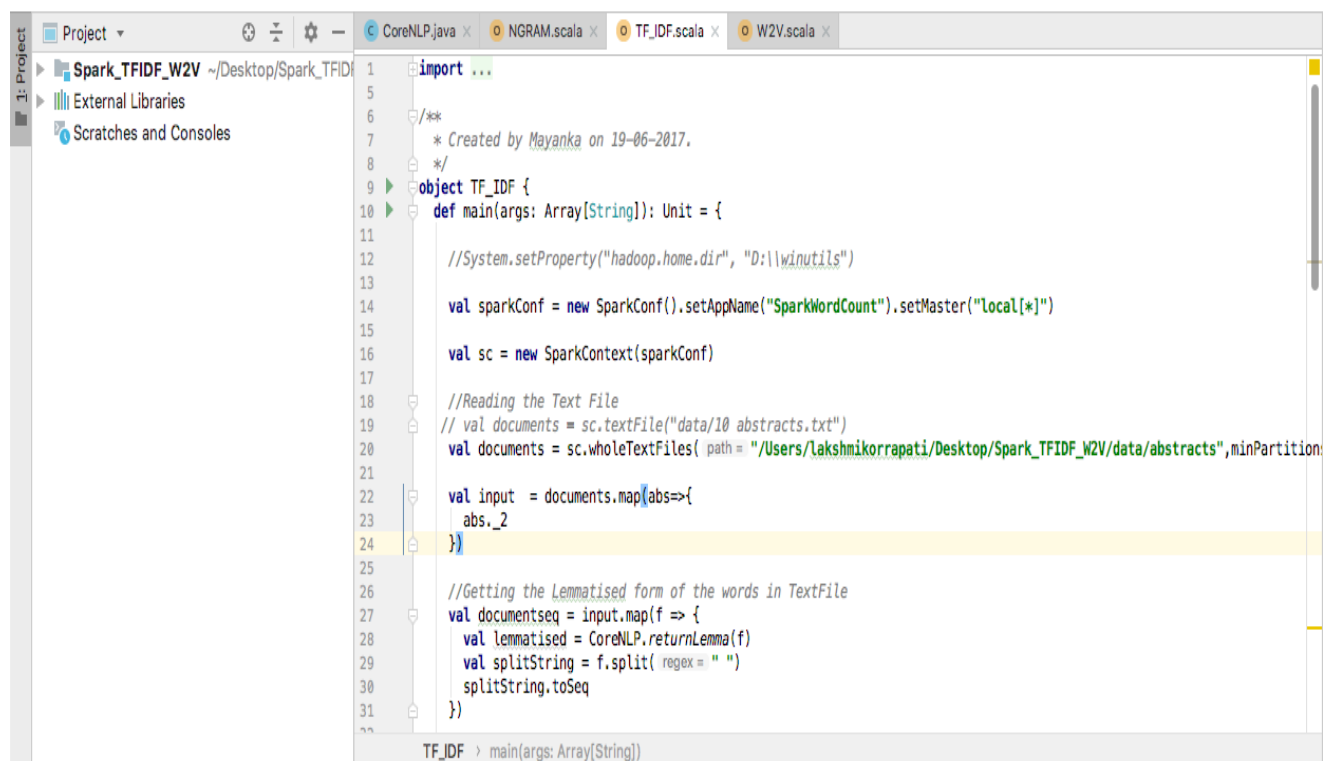
   Input: 5 Abstracts saved in separate files

Tasks:

   a. Find out the top TF-IDF words for the above input.
   b. Find out the top TF-IDF words for the lemmatized input
   c. Find out the top TF-IDF words for the n-gram based input.
2. Write a simple spark program to read a dataset and find the W2V Synonyms for the Top TF-IDF Words
   a. Try without NLP
   b. Try with Lemmatization
   c. Try with NGrams

   Compare the results from (a) , (b) and (c)

   TF-IDF:

```scala
Project ▼          C CoreNLP.java ×   O NGRAM.scala ×   O TF_IDF.scala ×   O W2V.scala ×
Spark_TFIDF_W2V ~/Desktop/Spark_TFIDF  1   import ...
External Libraries                      5
Scratches and Consoles                  6   /**
                                        7    * Created by Mayanka on 19-06-2017.
                                        8    */
                                        9   object TF_IDF {
                                       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       //Reading the Text File
                                       19       // val documents = sc.textFile("data/10 abstracts.txt")
                                       20       val documents = sc.wholeTextFiles( path = "/Users/lakshmikorrapati/Desktop/Spark_TFIDF_W2V/data/abstracts",minPartition:
                                       21
                                       22       val input  = documents.map(abs=>{
                                       23         abs._2
                                       24       })
                                       25
                                       26       //Getting the Lemmatised form of the words in TextFile
                                       27       val documentseq = input.map(f => {
                                       28         val lemmatised = CoreNLP.returnLemma(f)
                                       29         val splitString = f.split( regex = " ")
                                       30         splitString.toSeq
                                       31       })

TF_IDF > main(args: Array[String])
```

Output:



W2V for TF-IDF words:



```scala
import ...

/**
  * Created by Mayanka on 19-06-2017.
  */
object W2V {
  def main(args: Array[String]): Unit = {

    // System.setProperty("hadoop.home.dir", "D:\\winutils")

    val sparkConf = new SparkConf().setAppName("SparkWordCount").setMaster("local[*]")
      .set("spark.driver.memory", "6g").set("spark.executor.memory", "6g")

    val sc = new SparkContext(sparkConf)

    val input = sc.textFile( path = "data/sample").map(line => line.split( regex = " ").toSeq)

    val modelFolder = new File( pathname = "myModelPath")

    if (modelFolder.exists()) {
      val sameModel = Word2VecModel.load(sc,  path = "myModelPath")
      val synonyms = sameModel.findSynonyms( word = "zero",  num = 40)

      for ((synonym, cosineSimilarity) <- synonyms) {
        println(s"$synonym $cosineSimilarity")
      }
    }
    else {
      val word2vec = new Word2Vec() setVectorSize(1000)
```

W2V › main(args: Array[String])

Output:

```
nine 0.9945034980773926
class 0.9939331412315369
six 0.9934068322181702
eight 0.9930217266082764
late 0.9892733097076416
four 0.989136278629303
international 0.9888012409210205
modern 0.9885989427566528
european 0.9872747659683228
two 0.9852036833763123
votes 0.9833505749702454
civil 0.9820547103881836
theory 0.9796767234802246
front 0.9795065522193909
```