RDD DEEP DIVE

Assignment 1

Task 1

1.Write a program to read a text file and print the number of rows of data in the document.

2. Write a program to read a text file and print the number of words in the document.

3. We have a document where the word separator is -, instead of space. Write a spark

code, to obtain the count of the total number of words present in the document.

```scala
package core

import org.apache.spark.sql.SparkSession

object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark =SparkSession
      .builder()
      .master(master="local")
      .appName(name="RDDExample")
      .config("spark.some.config.option","some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
```

using scpark context sc reading input file at location C:/Users/mypc/Desktop/input.txt & printing each line

```scala
    val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/input.txt");
    // data.count
```
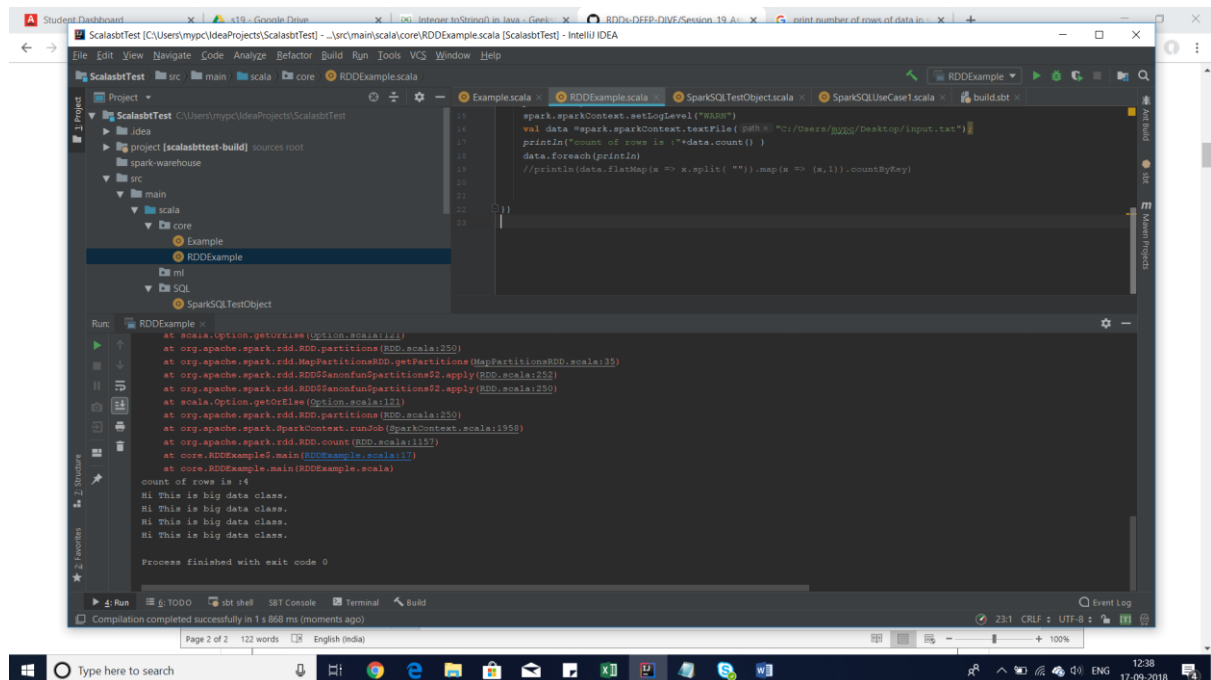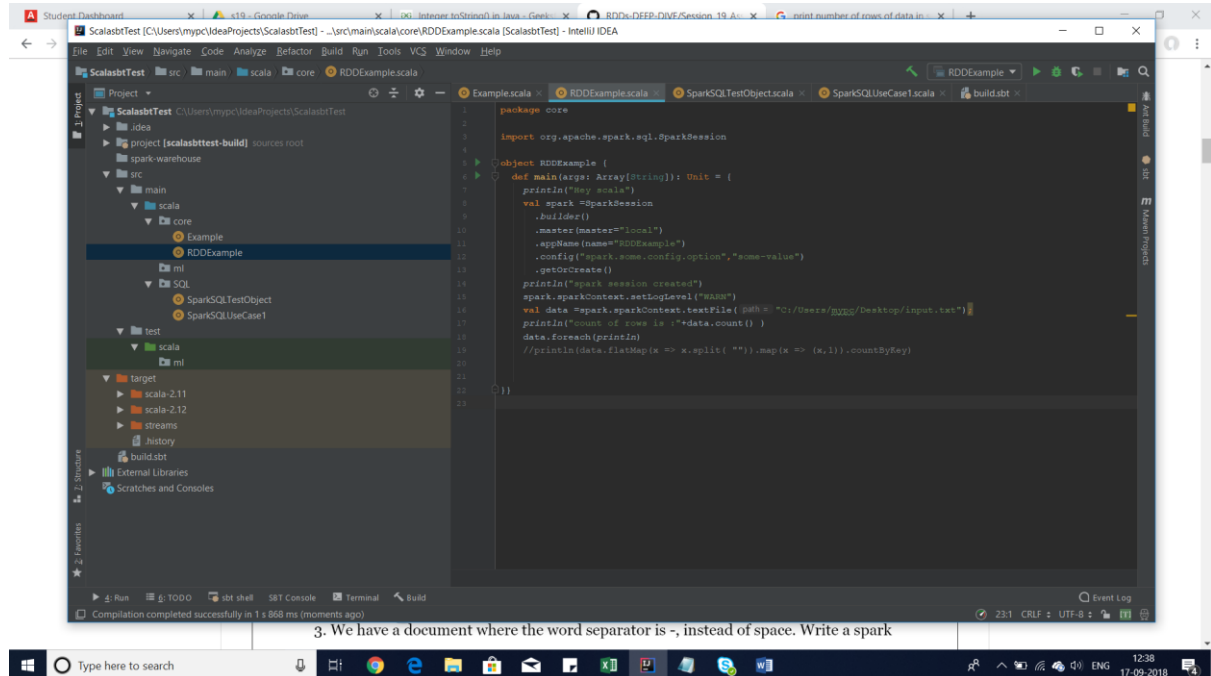
**(**Counting the total number of lines within the file**)**

```scala
    println("count of rows is :"+data.count() )
    data.foreach(println)
```

```
//println(data.flatMap(x => x.split( "")).map(x => (x,1)).countByKey)
```
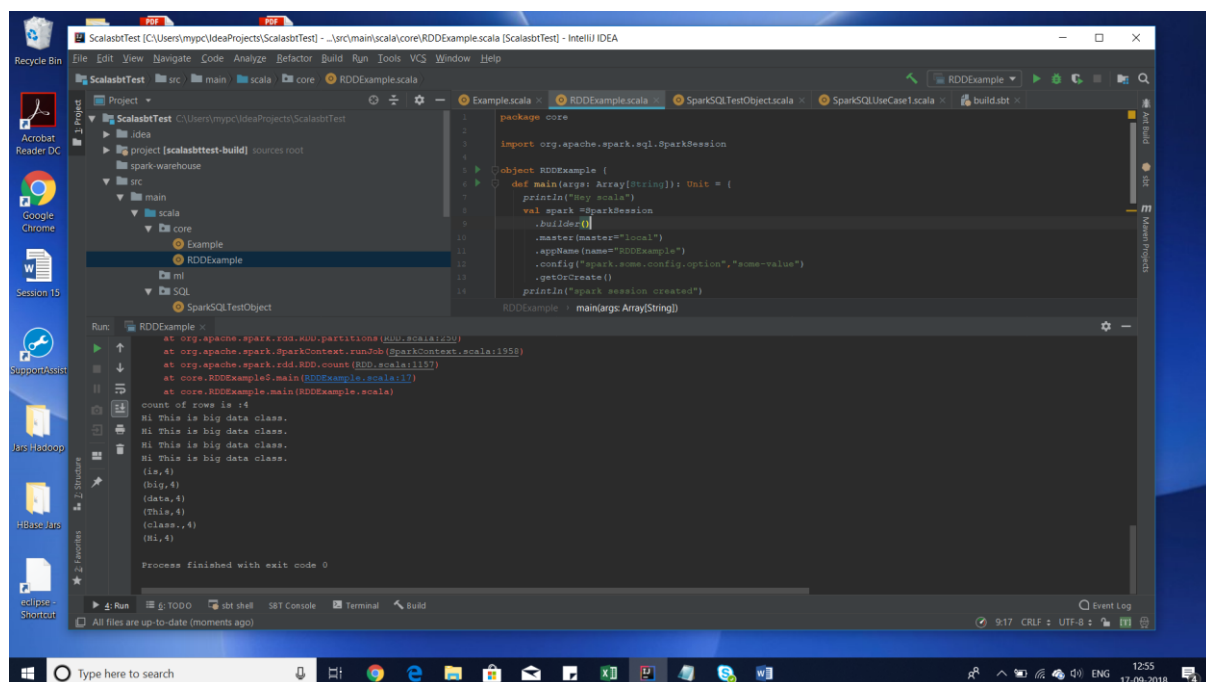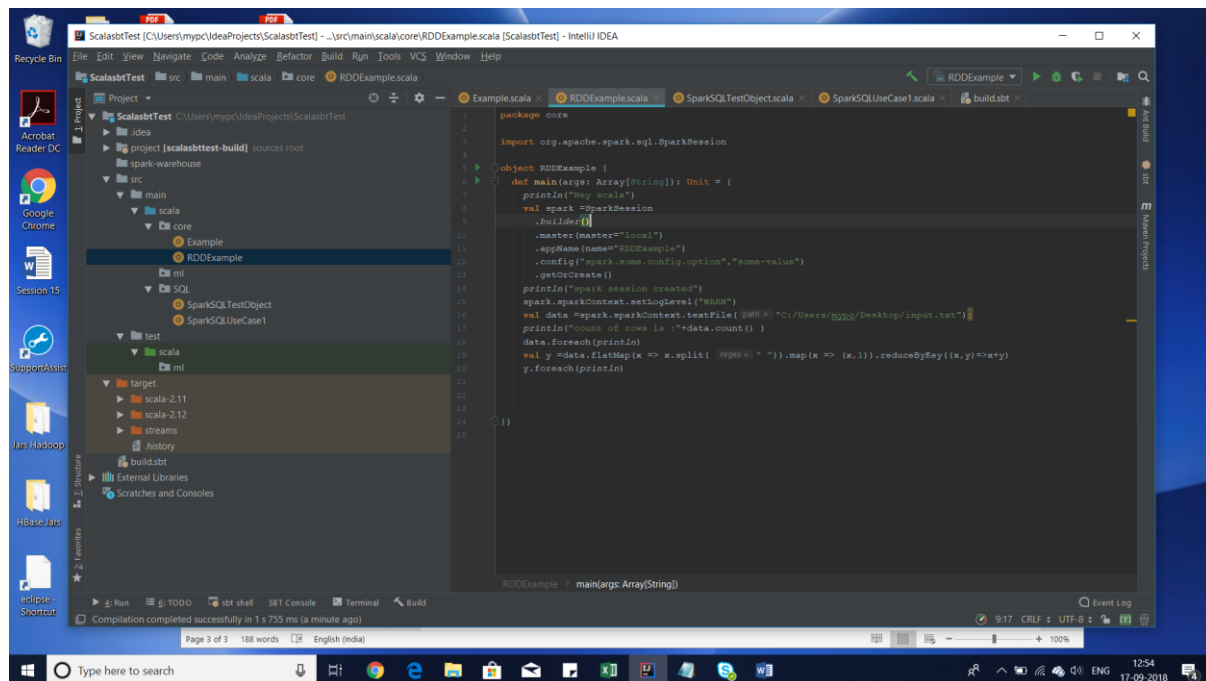

```
}}
```



3. We have a document where the word separator is -, instead of space. Write a spark



2. Write a program to read a text file and print the number of words in the document.

package core

import org.apache.spark.sql.SparkSession

**//creating RDD from a text file. Using split method on each line to split the line based on spaces. flatmap is used to map each line & then flatten those lines into isolated words. Using map again to associate an integer 1 with each word within file. reduceByKey treats all the words as key & adds up the ones for a particular word that is repeated**

```scala
object RDDExample {

  def main(args: Array[String]): Unit = {

    println("Hey scala")

    val spark =SparkSession

      .builder()

      .master(master="local")

      .appName(name="RDDExample")

      .config("spark.some.config.option","some-value")

      .getOrCreate()

    println("spark session created")

    spark.sparkContext.setLogLevel("WARN")

    val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/input.txt");

    println("count of rows is :"+data.count() )

    data.foreach(println)

    val y =data.flatMap(x => x.split( " ")).map(x => (x,1)).reduceByKey((x,y)=>x+y)

    y.foreach(println)



}}
```
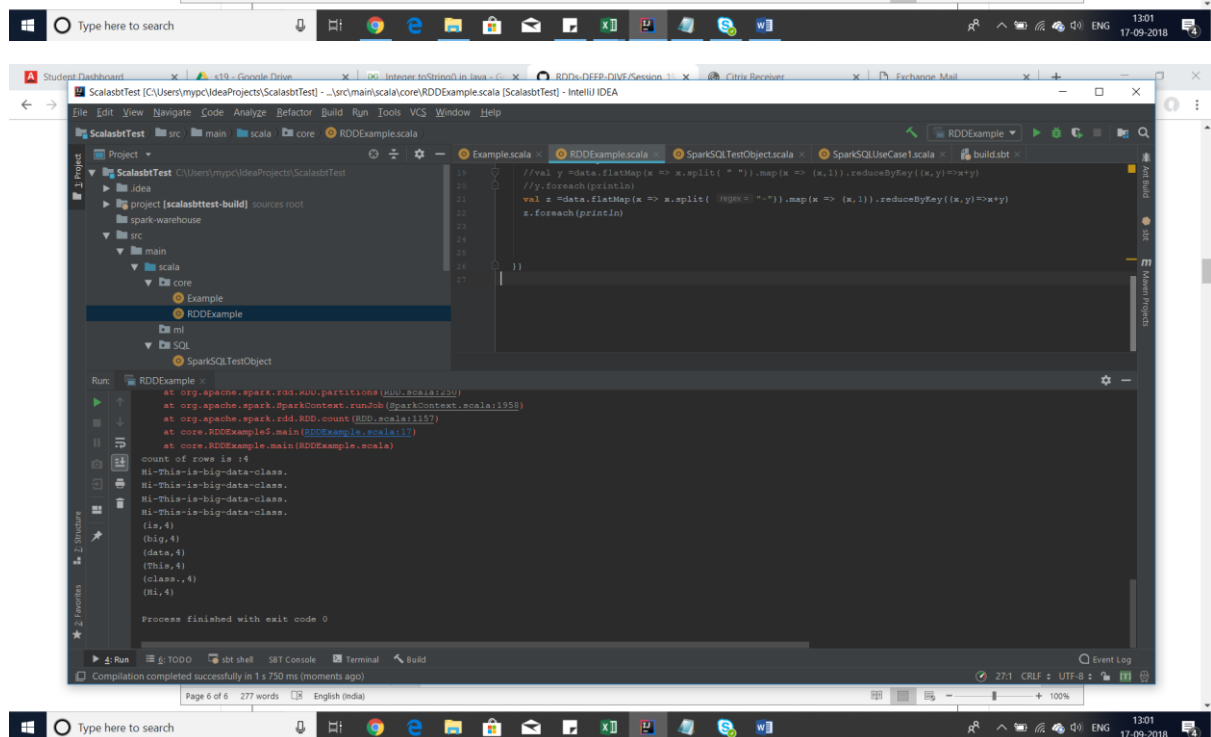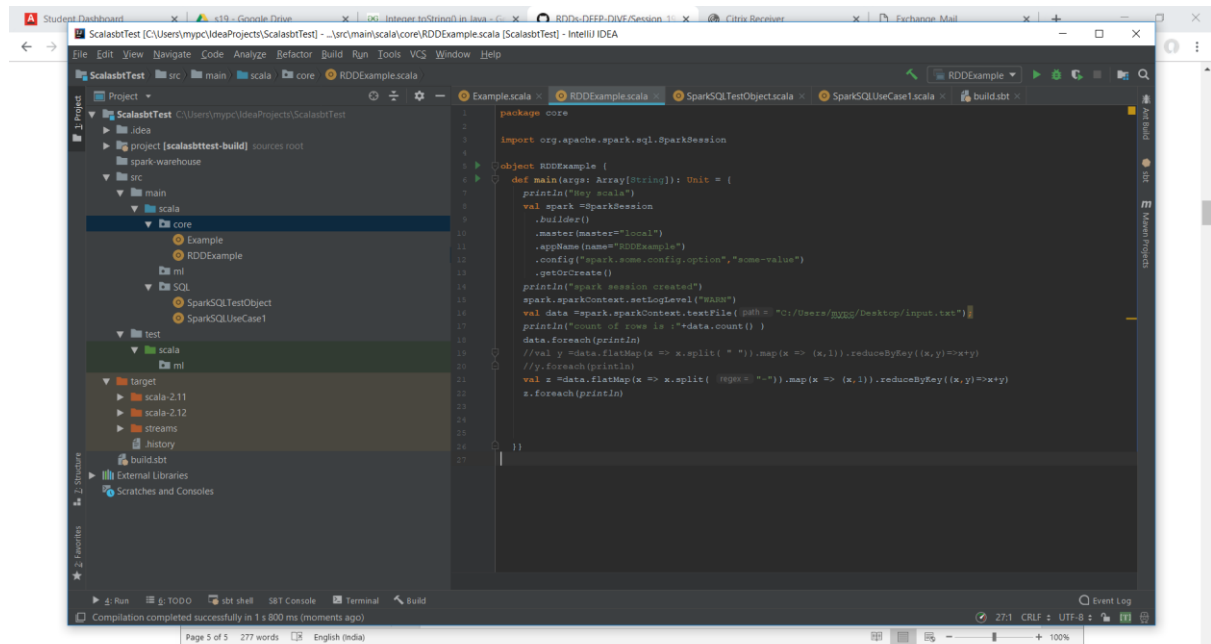
3. We have a document where the word separator is -, instead of space. Write a spark

code, to obtain the count of the total number of words present in the document.

**//Explanation: creating RDD from a text file. Using split method on each line to split the line based on hyphen(-). flatmap is used to map each line & then flatten those lines into isolated words. Using map again to associate an integer 1 with each word within file. reduceByKey treats all the words as key & adds up the ones for a word that is repeated**

```scala
package core

import org.apache.spark.sql.SparkSession

object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark =SparkSession
      .builder()
      .master(master="local")
      .appName(name="RDDExample")
      .config("spark.some.config.option","some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/input.txt");
    println("count of rows is :"+data.count() )
    data.foreach(println)
    //val y =data.flatMap(x => x.split( " ")).map(x => (x,1)).reduceByKey((x,y)=>x+y)
    //y.foreach(println)
    val z =data.flatMap(x => x.split( "-")).map(x => (x,1)).reduceByKey((x,y)=>x+y)
    z.foreach(println)



}}
```
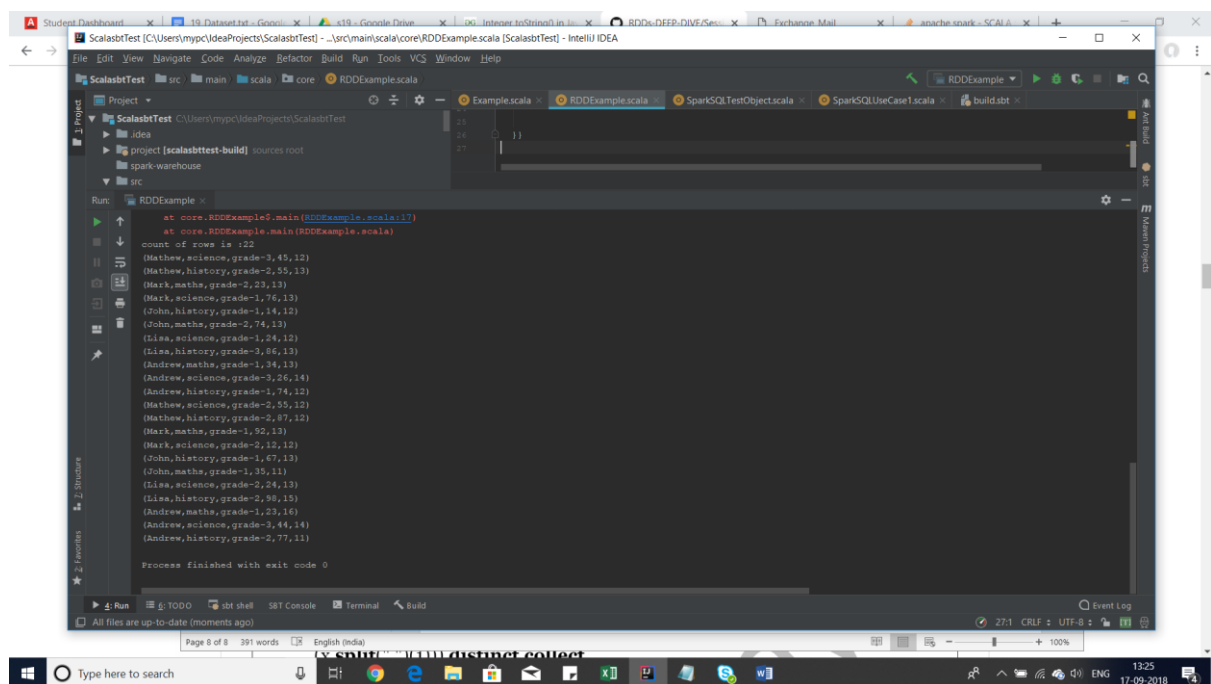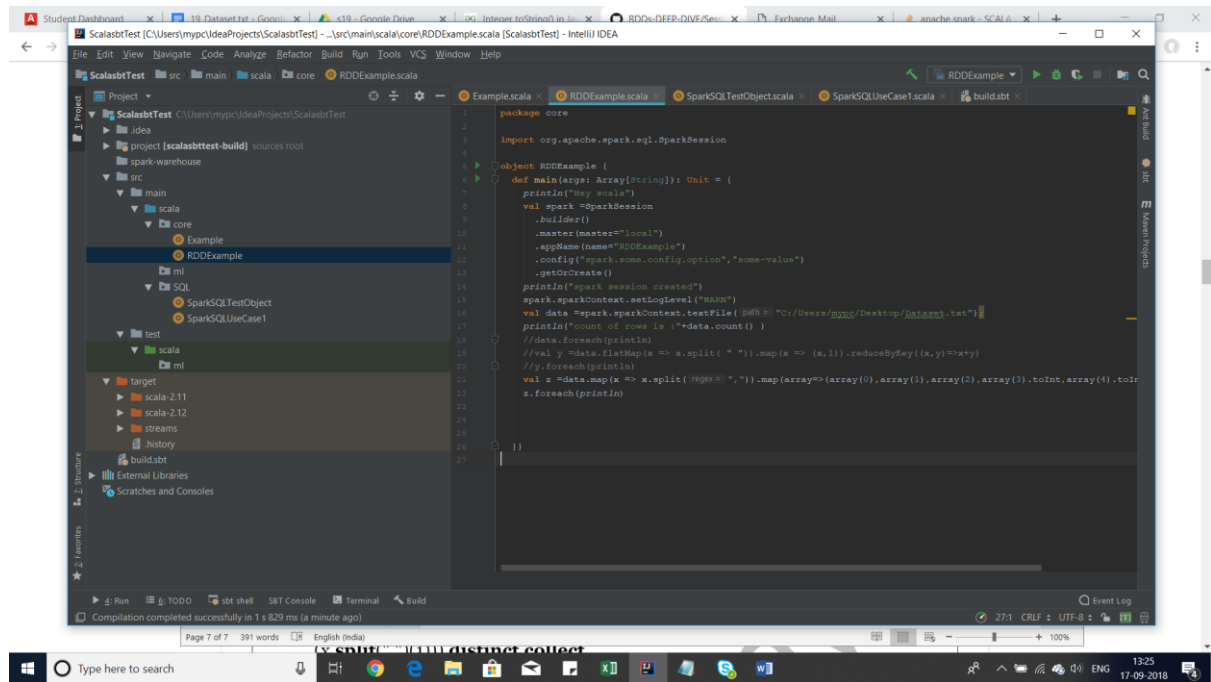
Task 2

Problem Statement 1:

1. Read the text file, and create a tupled rdd.

2. Find the count of total number of rows present.

3. What is the distinct number of subjects present in the entire school

4. What is the count of the number of students in the school, whose name is Mathew and

marks is 55

Solution:

**Explanation:**

**// creating tuple RDD from a text file. Using split method on each line to split the line based on ",". Using map again to create a tuple of each line from the file. toInt is used to convert the numeric columns from String to Integer**

```
package core

import org.apache.spark.sql.SparkSession

object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark =SparkSession
      .builder()
      .master(master="local")
      .appName(name="RDDExample")
      .config("spark.some.config.option","some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt");
    println("count of rows is :"+data.count() )
    //data.foreach(println)
    //val y =data.flatMap(x => x.split( " ")).map(x => (x,1)).reduceByKey((x,y)=>x+y)
    //y.foreach(println)
    val z =data.map(x =>
x.split(",")).map(array=>(array(0),array(1),array(2),array(3).toInt,array(4).toInt)).collect
    z.foreach(println)



}}
```
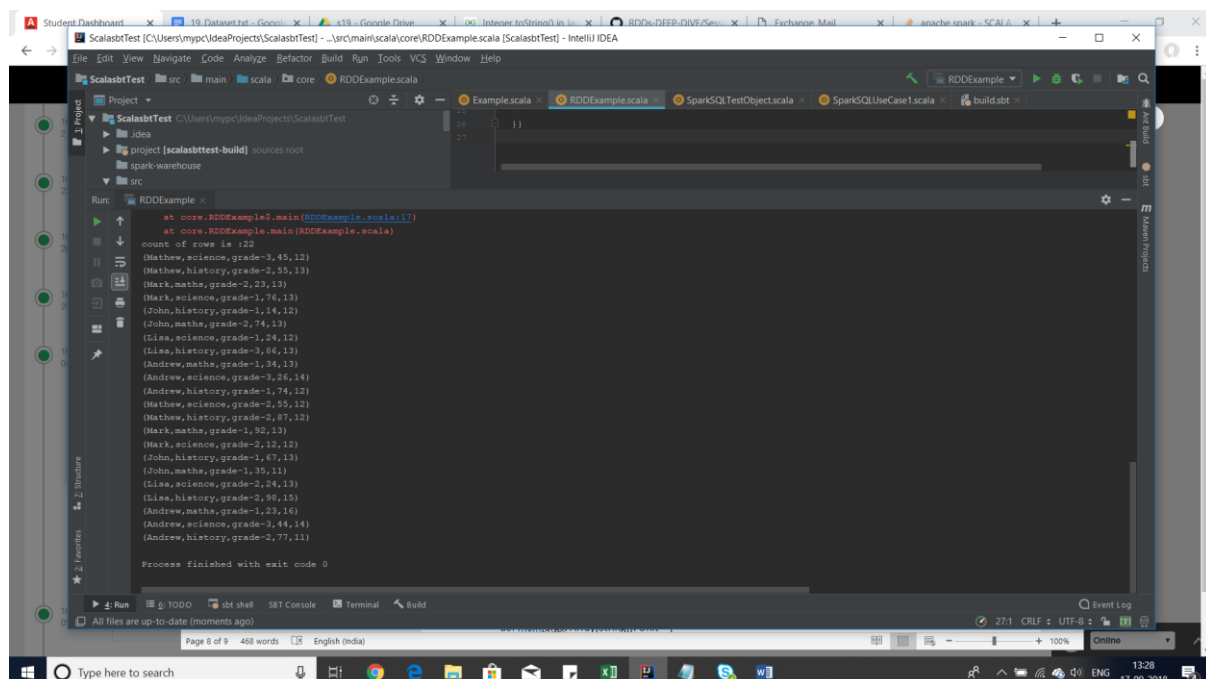
2. Find the count of total number of rows present.

Solution:We can find the count of rows using

"println("count of rows is :"+data.count() )" in the program

object RDDExample {

```scala
def main(args: Array[String]): Unit = {

  println("Hey scala")

  val spark =SparkSession

    .builder()

    .master(master="local")

    .appName(name="RDDExample")

    .config("spark.some.config.option","some-value")

    .getOrCreate()

  println("spark session created")

  spark.sparkContext.setLogLevel("WARN")

  val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt");

  println("count of rows is :"+data.count() )

  val z =data.map(x =>
x.split(",")).map(array=>(array(0),array(1),array(2),array(3).toInt,array(4).toInt)).collect

  z.foreach(println)




}}
```
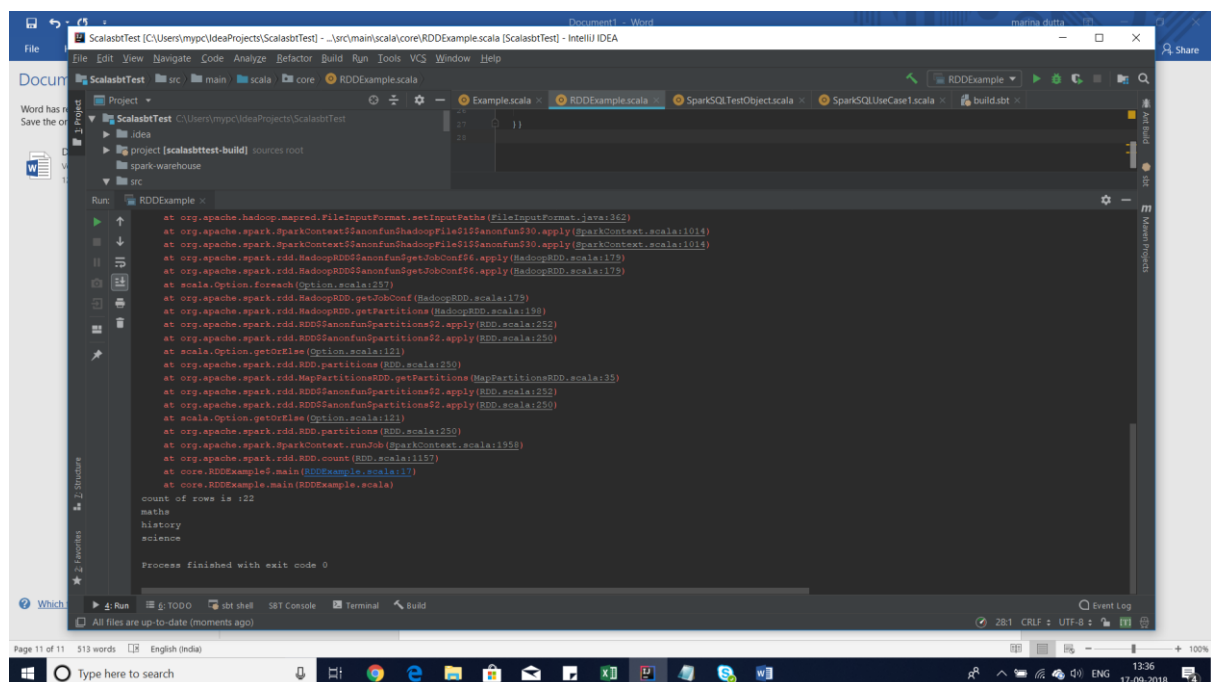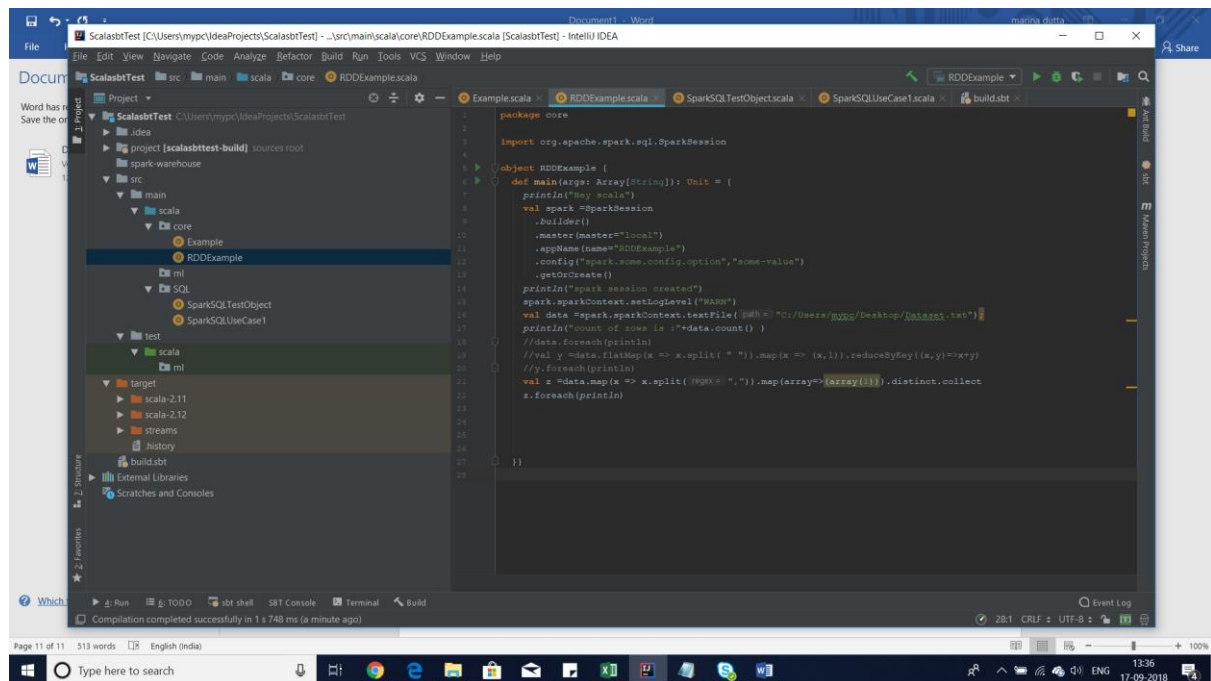
What is the distinct number of subjects present in the entire school

**Solution:**

**Explanation: Creating a n RDD by splitting an input dataset based on "," & selecting the subject field at index position 1. using distinct method over field at index position 1 to find the unique values within that column. Displaying them by calling the action collect.**

```scala
package core

import org.apache.spark.sql.SparkSession

object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark =SparkSession
      .builder()
      .master(master="local")
      .appName(name="RDDExample")
      .config("spark.some.config.option","some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt");
    println("count of rows is :"+data.count() )
    val z =data.map(x => x.split(",")).map(array=>(array(1))).distinct.collect
    z.foreach(println)



}}
```

4. What is the count of the number of students in the school, whose name is Mathew and marks is 55

**Solution:**

**Explanation:**

**// Using the filter function over the RDD to select all those first fields whose name is "Mathew" using the .equals() function to match the given**

**condition. filter returns all the matching columns when condition becomes True. Using logical && operator to**

**find all such fields where marks = 55. Using count method to to count all such lines that holds true for both the conditions**

```scala
package core

import org.apache.spark.sql.SparkSession

object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark =SparkSession
      .builder()
      .master(master="local")
      .appName(name="RDDExample")
      .config("spark.some.config.option","some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt");
    println("count of rows is :"+data.count() )

val z= data.filter(x=>(x.split(",")(0).equals("Mathew") && x.split(",")(3).toInt.equals(55))).count
    println("count of students whose name is  Mathew :" +z)


  }}
```

Problem Statement 2:

1. What is the count of students per grade in the school?

    import org.apache.spark.sql.SparkSession

// **Explanation:** Creating an RDD & selecting the field grade at the 2nd index i.e. 3rd position. Assigning 1 to each grade & thereby performing reduceByKey to count the value of each grade.

```scala
object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark =SparkSession
      .builder()
      .master(master="local")
      .appName(name="RDDExample")
      .config("spark.some.config.option","some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data =spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt");
    println("count of rows is :"+data.count() )

 val z= data.map(x=>x.split(",")(2)).map(x=>(x,1)).reduceByKey((x,y)=>x+y)

    println("count of students per grade is :" + " val y = z.foreach(println)")
    val y = z.foreach(println)


  }}
```
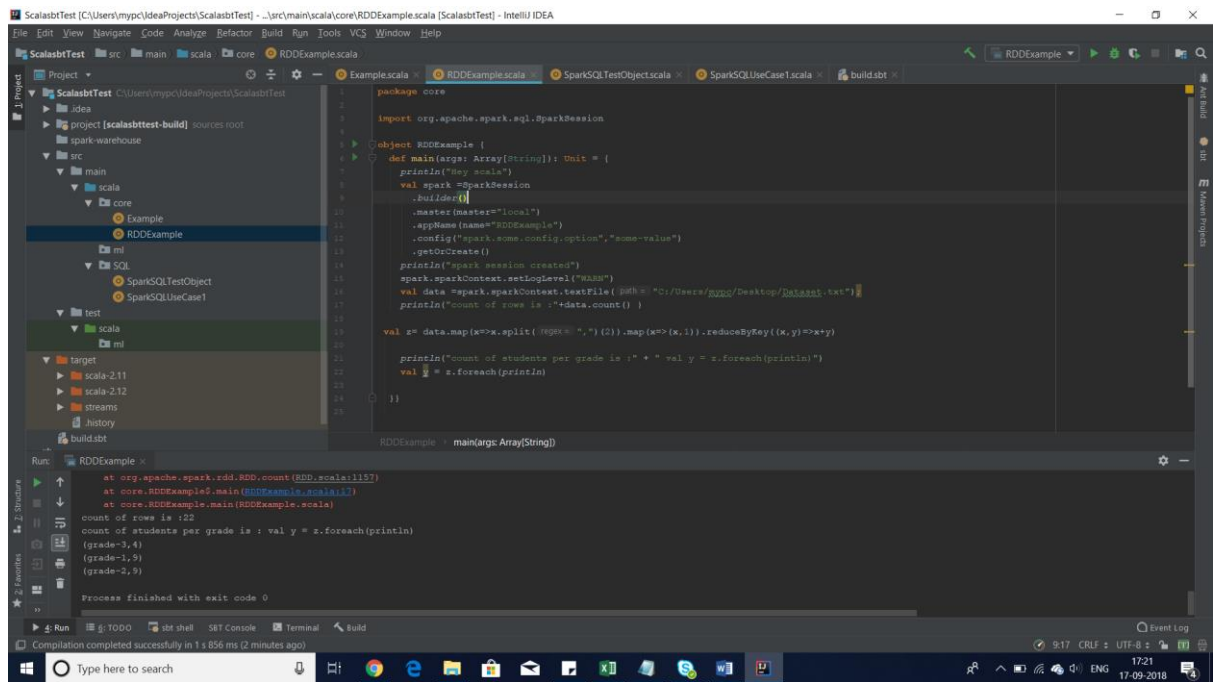
2. Find the average of each student (Note - Mathew is grade-1, is different from Mathew in some other grade!)

**Explanation:**

Creating an RDD from textFile. Selecting name & grade as key, marks as value. Converting marks into double datatype just to avoid any truncation after decimal. groupByKey returns a tuple of ((name, grade), marks). Applying map function & selecting name, grade as first value of outer tuple & marks as second value. Adding the marks for respective grade & for each distinct grade fetching the count of marks as x._2.size to divide the sum & find the average on it. Here grouping key is student & grade.

```scala
import org.apache.spark.sql.SparkSession


object RDDExample {

 def main(args: Array[String]): Unit = {

  println("Hey scala")

  val spark = SparkSession

    .builder()
```

```scala
      .master(master = "local")

      .appName(name = "RDDExample")

      .config("spark.some.config.option", "some-value")

      .getOrCreate()

    println("spark session created")

    spark.sparkContext.setLogLevel("WARN")

    val data = spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt")

    println("count of rows is :" + data.count())

    val z = data.map(x => ((x.split(",")(0), x.split(",")(2)), x.split(",")(3).toDouble)).groupByKey().map(x =>
(x._1, x._2.sum / x._2.size)).foreach(println)



  }}
```

2. What is the average score of students in each subject across all grades?

> **Explanation:** Creating an RDD from textFile. Selecting name & subject as key, marks as value. Converting marks into double datatype just to avoid any truncation after decimal. groupByKey returns a tuple of ((name, subject), marks). Applying map function & selecting name, subject as first value of outer tuple & marks as second value. Adding the marks for respective subject & for each distinct subject fetching the count of marks as x._2.size to divide the sum & find the average on it. Here grouping key is student & subject.

```scala
import org.apache.spark.sql.SparkSession


object RDDExample {
 def main(args: Array[String]): Unit = {

  println("Hey scala")

  val spark = SparkSession

    .builder()

    .master(master = "local")

    .appName(name = "RDDExample")

    .config("spark.some.config.option", "some-value")

    .getOrCreate()
```

```scala
println("spark session created")

spark.sparkContext.setLogLevel("WARN")

val data = spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt")

println("count of rows is :" + data.count())

val z=
data.map(x=>((x.split(",")(0),x.split(",")(1)),x.split(",")(3).toDouble)).groupByKey().map(x=>(x._1,x._2.sum / x._2.size)).foreach(println)

}}
```

4. What is the average score of students in each subject per grade?

**Explanation:** Creating an RDD from textFile. Selecting name, subject & grade as key, marks as value. Converting marks into double datatype just to avoid any truncation after decimal. groupByKey returns a tuple of ((name, subject, grade), marks). Applying map function

& selecting name, subject & grade as first value of outer tuple & marks as second value. Adding the marks for respective subject, grade & for each distinct subject, grade

```scala
package core

import org.apache.spark.sql.SparkSession

object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark = SparkSession
      .builder()
      .master(master = "local")
      .appName(name = "RDDExample")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data = spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt")
    println("count of rows is :" + data.count())
    val z=
data.map(x=>((x.split(",")(0),x.split(",")(1),x.split(",")(2)),x.split(",")(3).toDouble)).groupByKey().map(
x=>(x._1,x._2.sum / x._2.size)).foreach(println)
}}
```
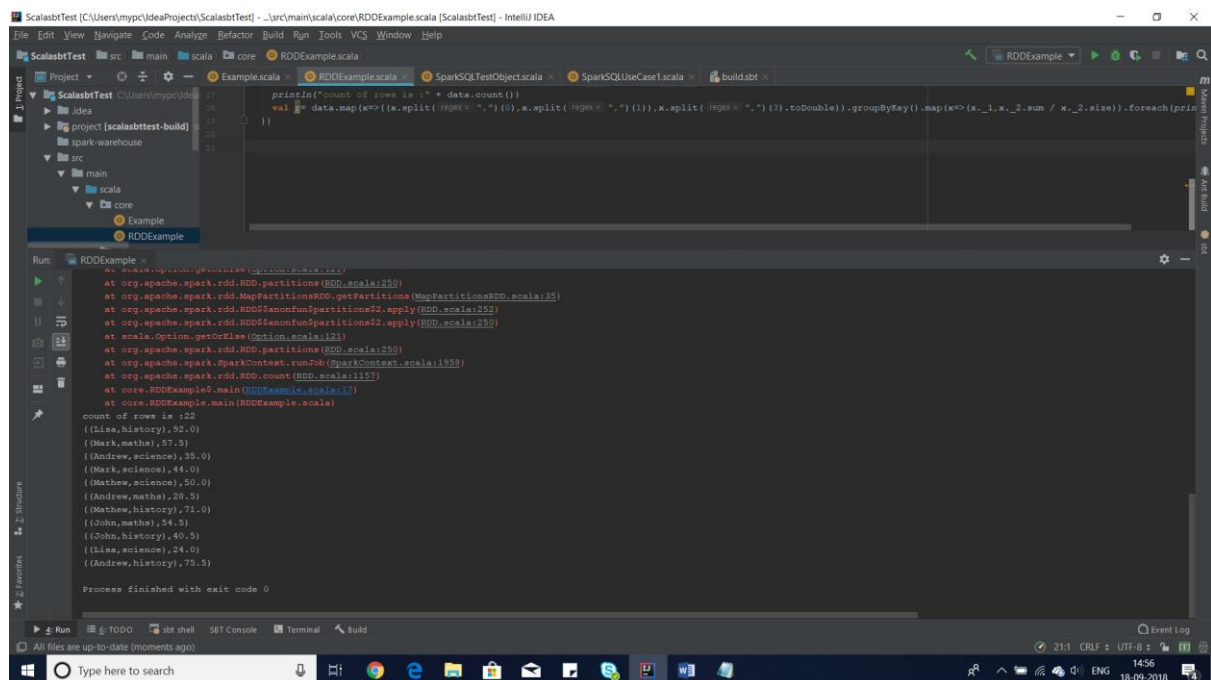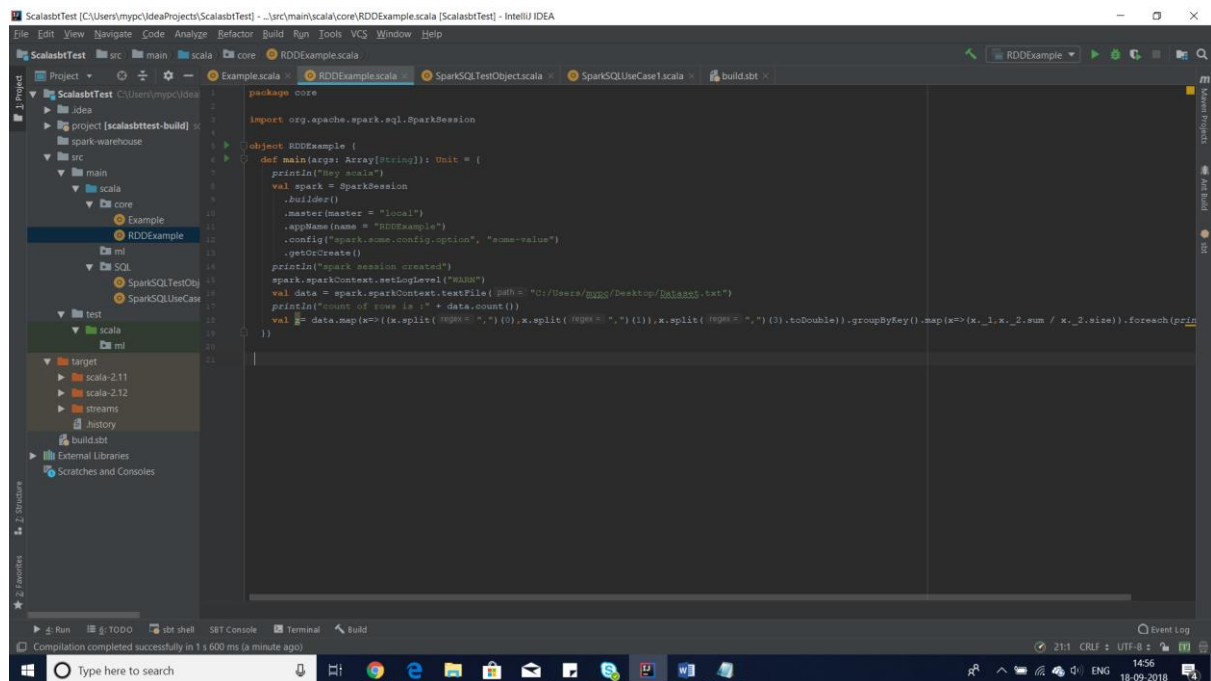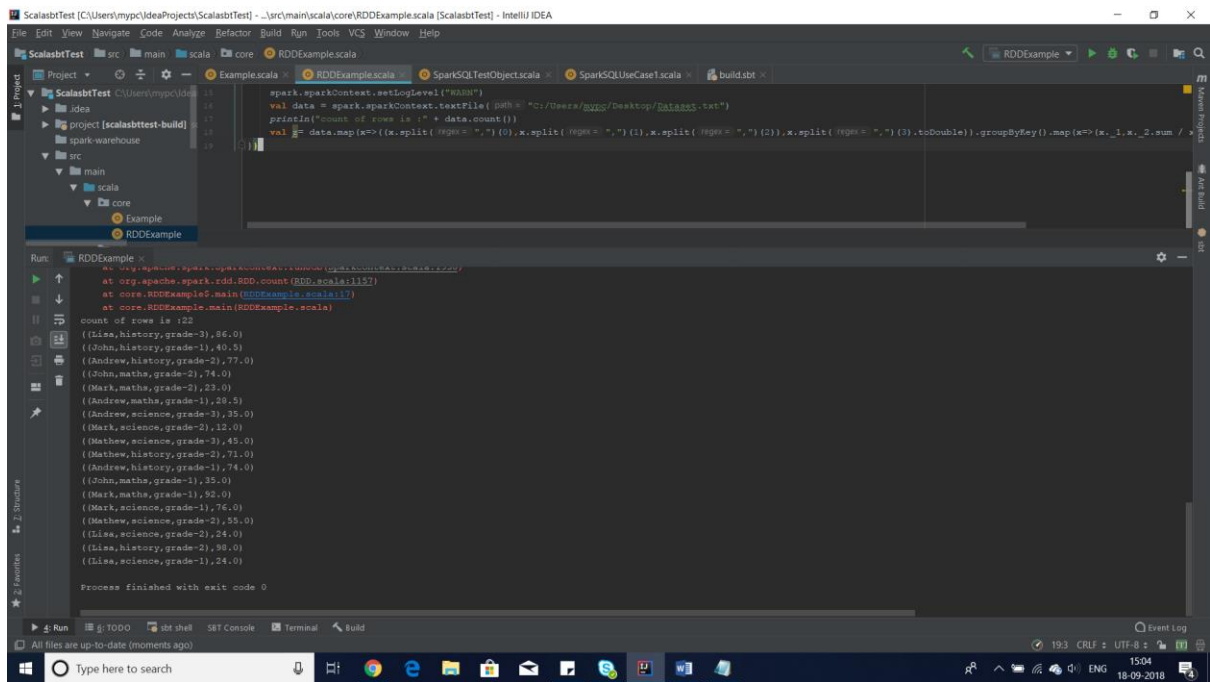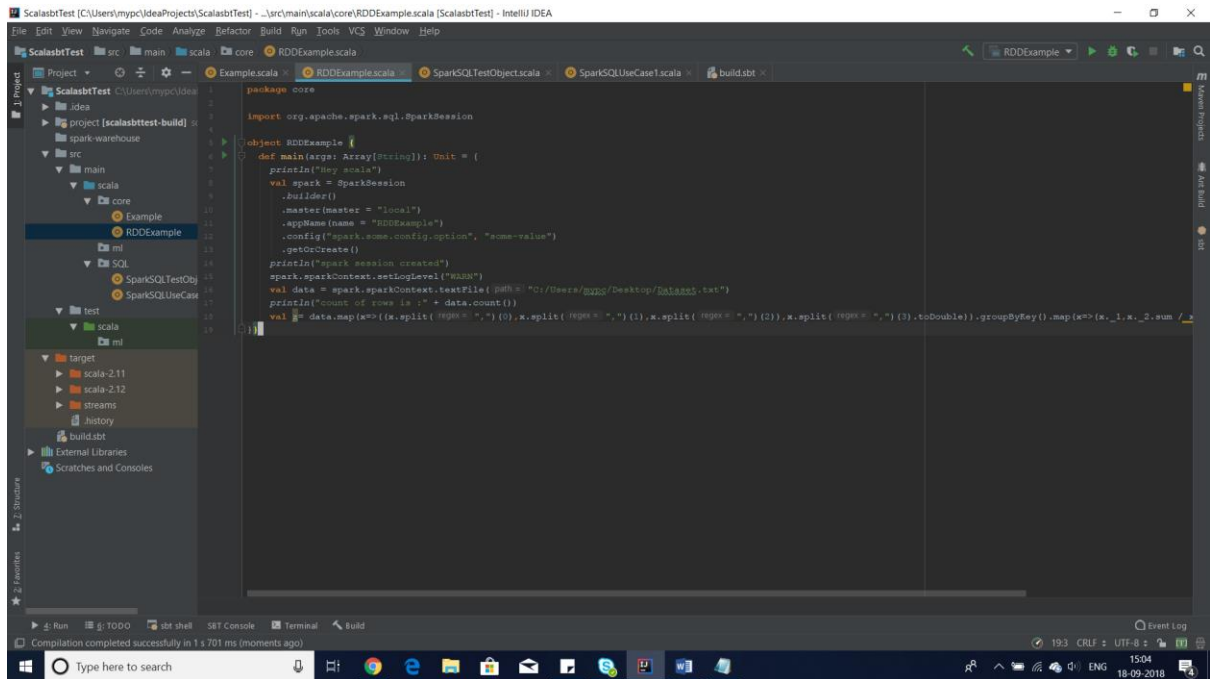
Run output:

```
        at org.apache.spark.rdd.RDD.count(RDD.scala:1157)
        at core.RDDExample$.main(RDDExample.scala:17)
        at core.RDDExample.main(RDDExample.scala)
count of rows is :22
((Lisa,history,grade-3),86.0)
((John,history,grade-1),40.5)
((Andrew,history,grade-2),77.0)
((John,maths,grade-2),74.0)
((Mark,maths,grade-2),23.0)
((Andrew,maths,grade-1),28.5)
((Andrew,science,grade-3),35.0)
((Mark,science,grade-2),12.0)
((Mathew,science,grade-3),45.0)
((Mathew,history,grade-2),71.0)
((Andrew,history,grade-1),74.0)
((John,maths,grade-1),35.0)
((Mark,maths,grade-1),92.0)
((Mark,science,grade-2),76.0)
((Mathew,science,grade-1),55.0)
((Lisa,science,grade-2),24.0)
((Lisa,history,grade-2),98.0)
((Lisa,science,grade-1),24.0)


Process finished with exit code 0
```

5. For all students in grade-2, how many have average score greater than 50?

**Explanation:** After the RDD creation. Filtering the records where grade equals grade-2 & then splitting the record where key is name value is marks. Thereby grouping by Name to calculate average of marks by calculating total marks for each name & dividing by count of times marks were awarded to each student. At last, applying filter & counting whoever has average marks greater than 50.

```scala
package core


import org.apache.spark.sql.SparkSession


object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark = SparkSession
      .builder()
      .master(master = "local")
      .appName(name = "RDDExample")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data = spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt")
    println("count of rows is :" + data.count())
    val z= data.filter(x=>(x.split(",")(2)).equals("grade-
2")).map(x=>(x.split(",")(0),x.split(",")(3).toDouble)).groupByKey().map(x=>(x._1,x._2.sum /
x._2.size)).filter(x=>(x._2>50)).foreach(println)


}}
```

Problem Statement 3:

Are there any students in the college that satisfy the below criteria:

1. Average score per student_name across all grades is same as average score per stud_name per grade

package core

**Solution: Here we are using the intersection property to find the above criteria. Finding the common data in both the lists by using the intersection property.**

```scala
import org.apache.spark.sql.SparkSession

object RDDExample {
  def main(args: Array[String]): Unit = {
    println("Hey scala")
    val spark = SparkSession
      .builder()
      .master(master = "local")
      .appName(name = "RDDExample")
      .config("spark.some.config.option", "some-value")
      .getOrCreate()
    println("spark session created")
    spark.sparkContext.setLogLevel("WARN")
    val data = spark.sparkContext.textFile("C:/Users/mypc/Desktop/Dataset.txt")
    println("count of rows is :" + data.count())
```
**//Finding the average score per student_name across all grades**

```scala
 val z = data.map(x=>(x.split(",")(0),x.split(",")(3).toDouble)).groupByKey().map(x=>(x._1,x._2.sum /
x._2.size))
```

```scala
    println("Average score per student_name across all grades is given below")
    val b=z.foreach(println)
```
**//average score per student per grade(**

```scala
    val y
=data.map(x=>((x.split(",")(0),x.split(",")(2)),x.split(",")(3).toDouble)).groupByKey().map(x=>(
x._1,x._2.sum / x._2.size))
    println("Average score each student  is given below")
    val y_1 =y.foreach(println)
```
**//average score per student per grade**

```scala
    val  per_grade
=data.map(x=>((x.split(",")(0),x.split(",")(2)),x.split(",")(3).toDouble)).groupByKey().map(x>(x._1._1
,x._2.sum / x._2.size))
```

```scala
    println("Average score per student_name per grade is given below")
    val pergrade_1 = per_grade.foreach(println)
    println("criteria for finding the average score per student_name across all grades is same
as average score per stud_name per grade is given below ")
    val z_1 =z.intersection(per_grade).collect()
}}
```

Note: Comparing the results using the intersection property we found that there was no common data in both the lists. There is no record that matches the criteria. After using the intersection property we did not get any output result hence no common data that matches the criteria.