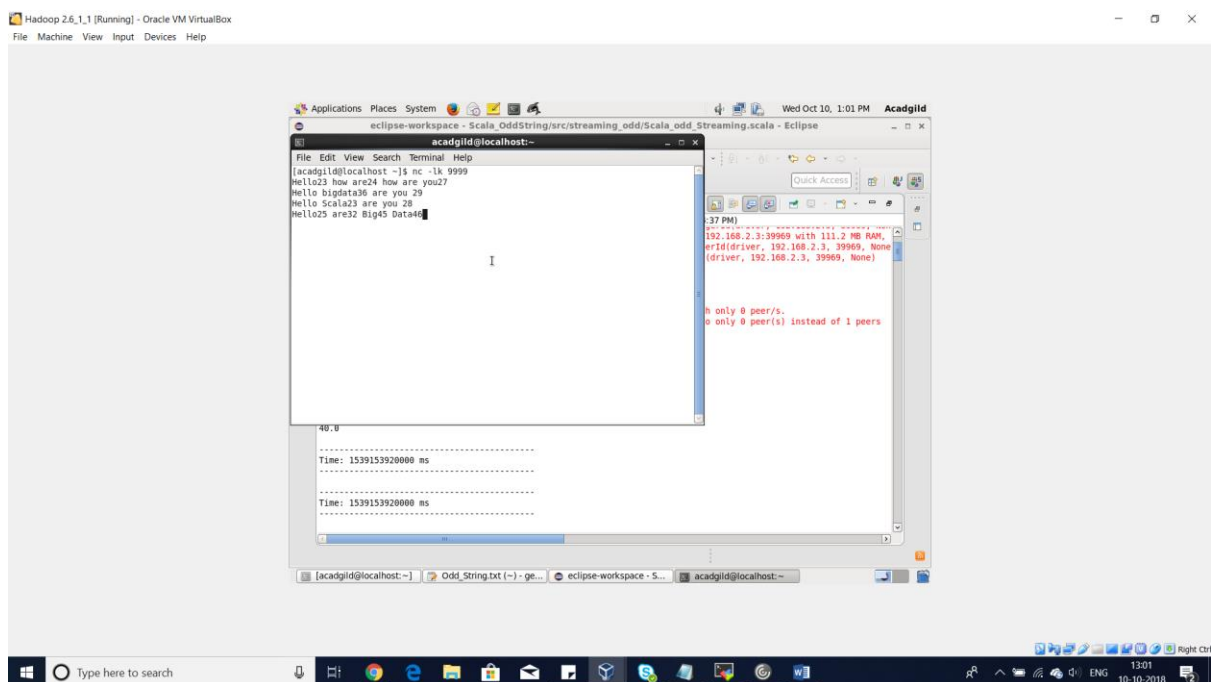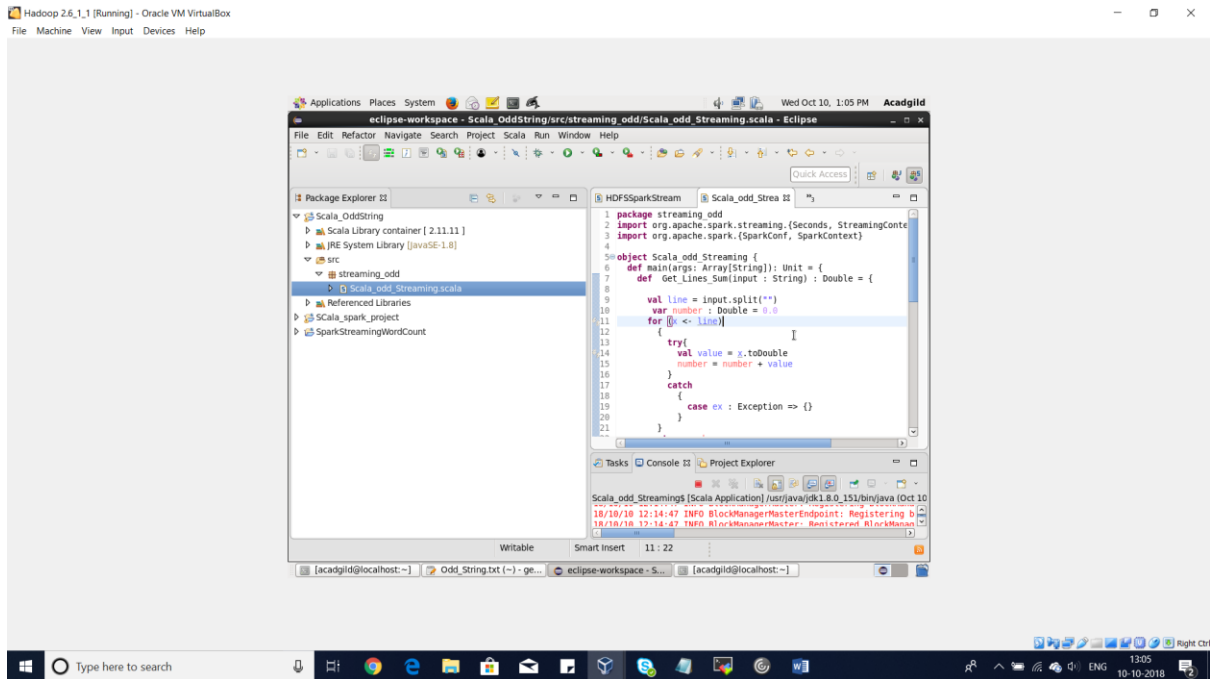**Session 24:**

**SPARK STREAMING**

**Assignment 1**

Task 1

Read a stream of Strings, fetch the words which can be converted to numbers. Filter out the rows,

where the sum of numbers in that line is odd.

Provide the sum of all the remaining numbers in that batch.

Output: – Before running the application, start "netcat" in terminal using the command "nc -lk 9999"
– Run the spark application – Put some string with numbers on the shell as shown below

Given below screenshot – we are able to see that lines containing sum of odd numbers are filter, even number is displayed, and sum of the number is displayed in the next line.



Task 2: Read two streams 1. List of strings input by user 2. Real-time set of offensive words Find the word count of the offensive words inputted by the user as per the realtime set of offensive words Solution: Note: Source code file is provided along with this assignment report.

Solution:

**In the spark application, we have a set of words that we considered as offensive words "idiot", "fool", "bad", "nonsense", "shit", "damn", "stupid", "dash", "bloody", "rascal",thief,thug**

In screenshot below, we are able to see that spark application has counted the number of offensive words occur as per the input string provided by the user.

Hadoop 2.6_1_1 [Running] - Oracle VM VirtualBox

File   Machine   View   Input   Devices   Help

Applications   Places   System          Wed Oct 10, 1:50 PM   Acadgild

eclipse-workspace - Scala_OddString/src/streaming_odd/Scala_bad_words.scala - Eclipse

File   Edit   Refactor   Navigate   Search   Project   Scala   Run   Window   Help

Quick Access

Package Explorer ⊠

Scala_OddString
  ▷ Scala Library container [ 2.11.11 ]
  ▷ JRE System Library [JavaSE-1.8]
  ▽ src
    ▽ streaming_odd
      ▷ Scala_bad_words.scala
      ▷ Scala_odd_Streaming.scala
  ▷ Referenced Libraries
▷ SCala_spark_project
▷ SparkStreamingWordCount

Tasks   Console ⊠   Project Explorer

Scala_bad_words$ [Scala Application] /usr/java/jdk1.8.0_151/bin/java (Oct 10, 2018, 1:47:
Spark Session Object Created!
Set(rascal, thief, bad, thug, bloody, shit, stupid, nonsense, idiot)
hey Spark Streaming!
Spark Streaming Context Created!
18/10/10 13:47:24 WARN RandomBlockReplicationPolicy: Expecting 1 replic
18/10/10 13:47:24 WARN BlockManager: Block input-0-1539159444400 replic
-------------------------------------------
Time: 1539159460000 ms
-------------------------------------------

(thug,2)
(thief,2)
(bad,4)
(idiot,4)

-------------------------------------------
Time: 1539159480000 ms
-------------------------------------------

```
23        //Let us create a spark session object
24        val conf = new SparkConf().setMaster("local[2]").
25        val sc = new SparkContext(conf)
26
27
28        sc.setLogLevel("WARN")
29        println("Spark Session Object Created!")
```

eclipse-workspace - S...   acadgild   [import_org.txt (~) - g...   acadgild@localhost:~

Right Ctrl

Type here to search        13:50  10-10-2018   ENG

---

Hadoop 2.6_1_1 [Running] - Oracle VM VirtualBox

File   Machine   View   Input   Devices   Help

Applications   Places   System          Wed Oct 10, 1:51 PM   Acadgild

eclipse-workspace - Scala_OddString/src/streaming_odd/Scala_bad_words.scala - Eclipse

File   Edit   Refactor   Navigate   Search   Project   Scala   Run   Window   Help

Quick Access

Package Explorer ⊠

Scala_OddString
  ▷ Scala Library container [ 2.11.11 ]
  ▷ JRE System Library [JavaSE-1.8]
  ▽ src
    ▽ streaming_odd
      ▷ Scala_bad_words.scala
      ▷ Scala_odd_Streaming.scala
  ▷ Referenced Libraries
▷ SCala_spark_project
▷ SparkStreamingWordCount

Scala_odd_Strea   Scala_bad_words ⊠

```
1  package streaming_odd
2  import org.apache.spark.streaming.{Seconds, StreamingConte
3      import scala.collection.mutable.ArrayBuffer
4
5  import org.apache.spark.streaming.{Seconds, StreamingConte
6      import scala.collection.mutable.ArrayBuffer
7      import org.apache.spark.SparkConf
8      import org.apache.spark.SparkContext
9
10
11
12 object Scala_bad_words {
13
14
15        //ArrayBuffer to store list of offensive words in memc
16        val wordList: ArrayBuffer[String] = ArrayBuffer.empty[
17
18
19        def main(args: Array[String]) {
20          println("hey Scala, Streaming Offensive Words!")
21
```

Tasks ⊠   Project Explorer

0 items

✓   !   Description                          Resource

Writable        Smart Insert      35 : 1

eclipse-workspace - S...   acadgild   [import_org.txt (~) - g...   [acadgild@localhost:~]

Right Ctrl

Type here to search        13:51  10-10-2018   ENG