# ASSIGNMENT 20.1

## Problem Statement:

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.

## Solution:

Here is the Spark code in Scala to create a stream of strings over socket connection:

**EvenLines.scala:**

```
package org.spark

// import required Spark packages

import org.apache.spark.SparkConf
import org.apache.spark.storage.StorageLevel
import org.apache.spark.streaming.{Seconds, StreamingContext}

/*
 * Counts words in UTF8 encoded, '\n' delimited text received from the network every second.
 * Usage: EvenLines <hostname> <port>
 * <hostname> and <port> describe the TCP server that Spark Streaming would connect to receive data.
 * To run this on your local machine, you need to first run a Netcat server with this command:
 *    `$ nc -lk 9999`
 * and then run the example
 *    `$ bin/run-example org.apache.spark.examples.streaming.NetworkWordCount localhost 9999`
 */
object EvenLines {
  def main(args: Array[String]) {
// check whether user has passed sufficient command line arguments for execution of the program
    if (args.length < 2) {
```

```scala
      System.err.println("Usage: EvenLines <hostname> <port>")
      System.exit(1)
    }

    // create a streaming context with a 10 second batch size
    val sparkConf = new SparkConf().setAppName("EvenLines")
    val ssc = new StreamingContext(sparkConf, Seconds(10))

    // Create a socket stream on target ip:port and count the
    // words in input stream of \n delimited text
    // Note that no duplication in storage level only for running locally.
    // Replication necessary in distributed scenario for fault tolerance.
    val lines = ssc.socketTextStream(args(0), args(1).toInt, StorageLevel. MEMORY_AND_
DISK_SER)

   // retain only those lines in the current batch in which the sum of all numbers is even
    val linesFiltered = lines.filter { x => getLineSum(x)%2==0 }

     // compute sum for the retained lines with even sum
    val linesSum = linesFiltered.map { x => getLineSum(x) }

    // print the resulting lines followed by their sum
    println("Lines with even sum")
    linesFiltered.print()
    println("")
    print("Sum of numbers in even lines : ")
    linesSum.reduce( (c1, c2) => c1 + c2).print()

    // start the streaming operation until the program gets terminated
    ssc.start()
    ssc.awaitTermination()
  }

  // a routine to compute sum of numbers in a line passed by user over a socket connection
  def getLineSum(ln : String): Double={
```
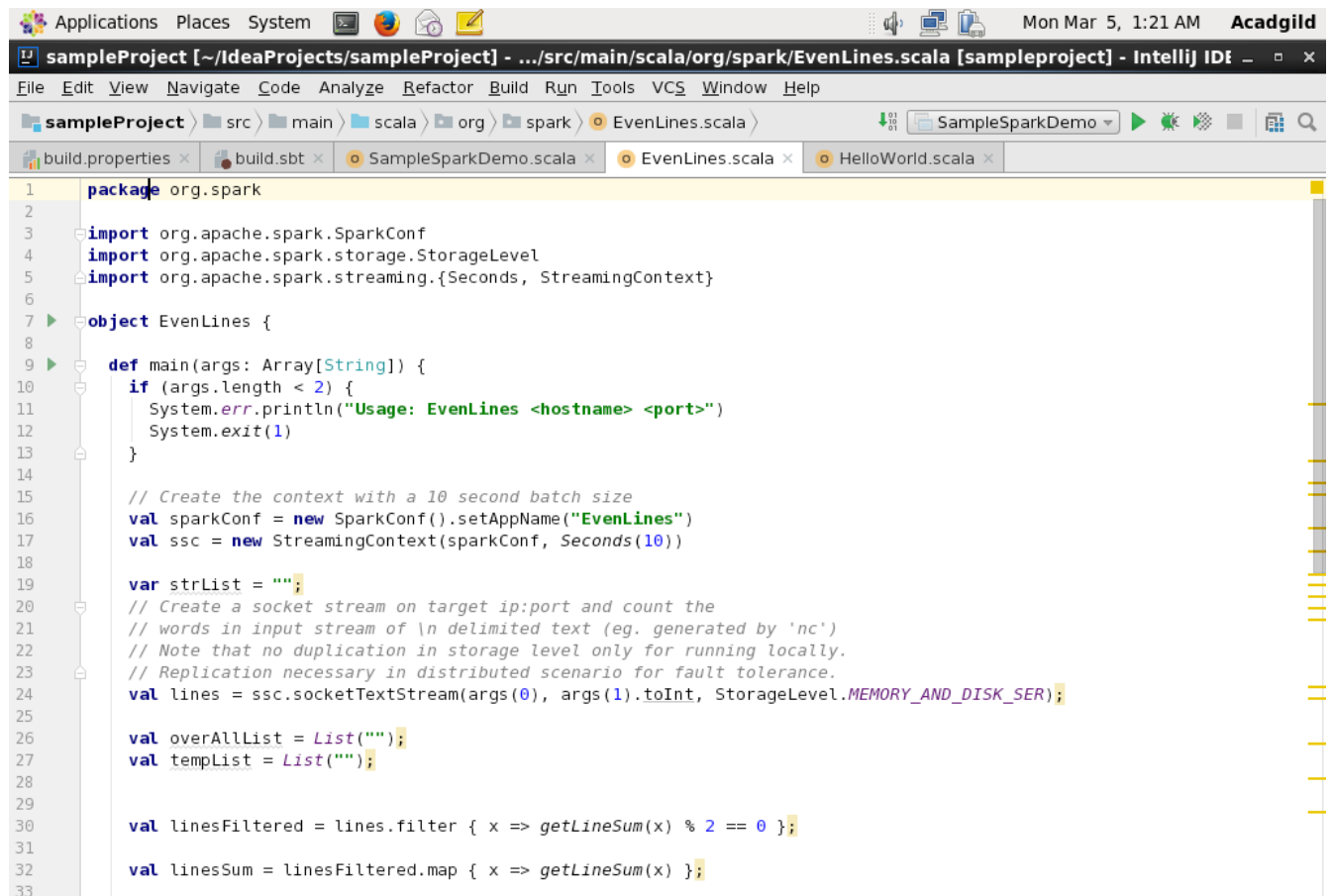
```scala
    val lineWords = ln.split(" ")

    var sum: Double = 0

    for(x <- lineWords)

    {

      try {

        val f = x.toDouble

        sum = sum + f

      } catch {

        case ex: Exception =>{    }

      }

    }

    sum                                // return sum of numbers present in a line

    }

}
```

sampleProject [~/IdeaProjects/sampleProject] - .../src/main/scala/org/spark/EvenLines.scala [sampleproject] - IntelliJ IDE _ □ ×

File  Edit  View  Navigate  Code  Analyze  Refactor  Build  Run  Tools  VCS  Window  Help

sampleProject ⟩ src ⟩ main ⟩ scala ⟩ org ⟩ spark ⟩ EvenLines.scala ⟩                    SampleSparkDemo ▾  ▶ ※ ▓ ■  ⊞ Q

build.properties ×   build.sbt ×   SampleSparkDemo.scala ×   EvenLines.scala ×   HelloWorld.scala ×

```scala
1     package org.spark
2
3     import org.apache.spark.SparkConf
4     import org.apache.spark.storage.StorageLevel
5     import org.apache.spark.streaming.{Seconds, StreamingContext}
6
7  ▶  object EvenLines {
8
9  ▶    def main(args: Array[String]) {
10        if (args.length < 2) {
11          System.err.println("Usage: EvenLines <hostname> <port>")
12          System.exit(1)
13        }
14
15        // Create the context with a 10 second batch size
16        val sparkConf = new SparkConf().setAppName("EvenLines")
17        val ssc = new StreamingContext(sparkConf, Seconds(10))
18
19        var strList = "";
20        // Create a socket stream on target ip:port and count the
21        // words in input stream of \n delimited text (eg. generated by 'nc')
22        // Note that no duplication in storage level only for running locally.
23        // Replication necessary in distributed scenario for fault tolerance.
24        val lines = ssc.socketTextStream(args(0), args(1).toInt, StorageLevel.MEMORY_AND_DISK_SER);
25
26        val overAllList = List("");
27        val tempList = List("");
28
29
30        val linesFiltered = lines.filter { x => getLineSum(x) % 2 == 0 };
31
32        val linesSum = linesFiltered.map { x => getLineSum(x) };
33
```

```
34        println("Lines with even sum");
35        linesFiltered.print();
36        println("");
37        print("Sum of numbers in even lines : ");
38        linesSum.reduce((c1, c2) => c1 + c2).print();
39
40        ssc.start()
41        ssc.awaitTermination()
42    }
43    def getLineSum(ln : String): Double={
44        val lineWords = ln.split(" ");
45        var num: Double = 0;
46        for(x <- lineWords)
47        {
48            try {
49                val f = x.toDouble;
50                num = num + f;
51            } catch {
52                case ex: Exception =>{      }
53            }
54        }
55        return num;
56    }
57 }
```

**Output:**

Here is the command to submit the jar file for the execution of the above code:

$ spark-submit --master local[2] -class org.scala.EvenLines Assignment20_1.jar localhost 9999

```
[acadgild@localhost scala-2.10]$
[acadgild@localhost scala-2.10]$ spark-submit --master local[2] --class org.spark.EvenLines sampleproject_2.10-0.1.jar localh
ost 9999
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
18/03/05 01:06:05 INFO SparkContext: Running Spark version 1.6.0
18/03/05 01:06:07 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes
 where applicable
18/03/05 01:06:08 WARN SparkConf:
SPARK_WORKER_INSTANCES was detected (set to '2').
This is deprecated in Spark 1.0+.

Please instead use:
 - ./spark-submit with --num-executors to specify the number of executors
 - Or set SPARK_EXECUTOR_INSTANCES
 - spark.executor.instances to configure the number of instances in the spark config.

18/03/05 01:06:08 WARN Utils: Your hostname, localhost.localdomain resolves to a loopback address: 127.0.0.1; using 192.168.4
3.63 instead (on interface eth0)
18/03/05 01:06:08 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
18/03/05 01:06:08 INFO SecurityManager: Changing view acls to: acadgild
18/03/05 01:06:08 INFO SecurityManager: Changing modify acls to: acadgild
18/03/05 01:06:08 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissio
ns: Set(acadgild); users with modify permissions: Set(acadgild)
```

On another terminal we need to run netcat command to start entering some numbers which will be captured by the streaming application
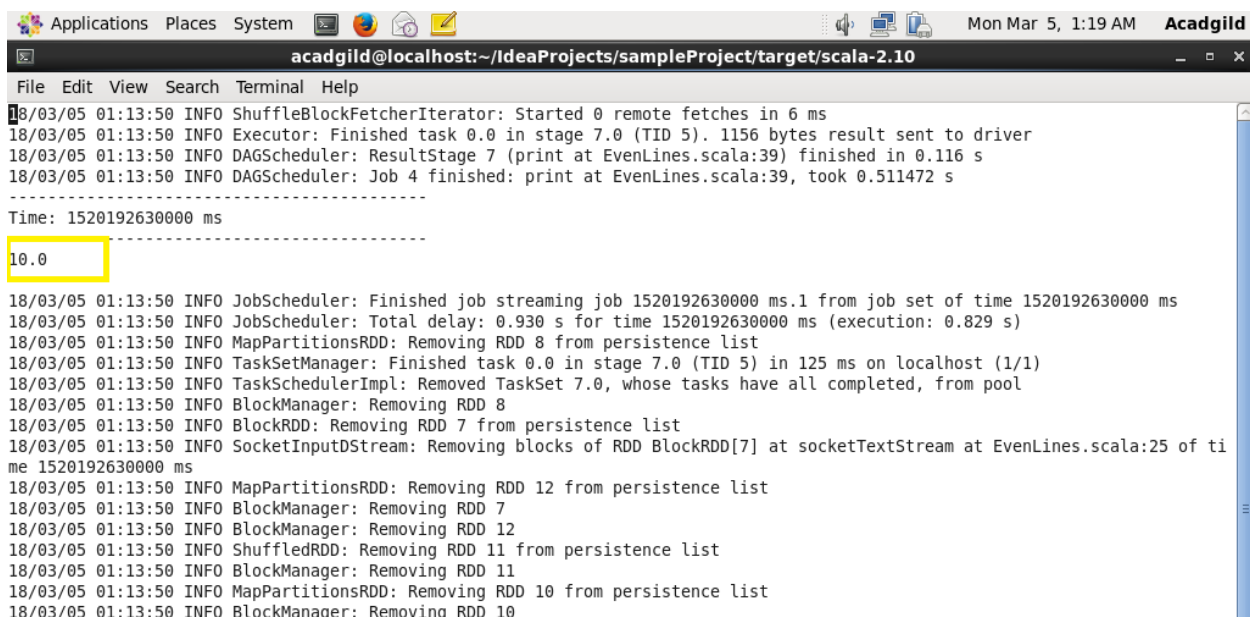
```
[root@localhost ~]# nc -lk 9999
test 2 3 5
test 7
test 10 20
35
^C
[root@localhost ~]#
```

Let's observe our streaming application run to see the output:

acadgild@localhost:~/IdeaProjects/sampleProject/target/scala-2.10 ▫️ ◻️ ✕

File  Edit  View  Search  Terminal  Help

```
18/03/05 01:14:20 INFO MemoryStore: Block broadcast_10_piece0 stored as bytes in memory (estimated size 1122.0 B, free 104.3
KB)
18/03/05 01:14:20 INFO BlockManagerInfo: Added broadcast_10_piece0 in memory on localhost:41378 (size: 1122.0 B, free: 517.4
MB)
18/03/05 01:14:20 INFO SparkContext: Created broadcast 10 from broadcast at DAGScheduler.scala:1006
18/03/05 01:14:20 INFO DAGScheduler: Submitting 1 missing tasks from ResultStage 13 (MapPartitionsRDD[32] at filter at EvenLi
nes.scala:31)
18/03/05 01:14:20 INFO TaskSchedulerImpl: Adding task set 13.0 with 1 tasks
18/03/05 01:14:20 INFO TaskSetManager: Starting task 0.0 in stage 13.0 (TID 10, localhost, partition 0,NODE_LOCAL, 2085 bytes
)
18/03/05 01:14:20 INFO Executor: Running task 0.0 in stage 13.0 (TID 10)
18/03/05 01:14:20 INFO BlockManager: Found block input-0-1520192654800 locally
18/03/05 01:14:20 INFO Executor: Finished task 0.0 in stage 13.0 (TID 10). 928 bytes result sent to driver
18/03/05 01:14:20 INFO TaskSetManager: Finished task 0.0 in stage 13.0 (TID 10) in 22 ms on localhost (1/1)
18/03/05 01:14:20 INFO TaskSchedulerImpl: Removed TaskSet 13.0, whose tasks have all completed, from pool
18/03/05 01:14:20 INFO DAGScheduler: ResultStage 13 (print at EvenLines.scala:36) finished in 0.010 s
18/03/05 01:14:20 INFO DAGScheduler: Job 8 finished: print at EvenLines.scala:36, took 0.067723 s
-------------------------------------------
Time: 1520192660000 ms
-------------------------------------------
test 10 20

18/03/05 01:14:20 INFO JobScheduler: Finished job streaming job 1520192660000 ms.0 from job set of time 1520192660000 ms
18/03/05 01:14:20 INFO JobScheduler: Starting job streaming job 1520192660000 ms.1 from job set of time 1520192660000 ms
18/03/05 01:14:20 INFO SparkContext: Starting job: print at EvenLines.scala:39
18/03/05 01:14:20 INFO DAGScheduler: Registering RDD 34 (reduce at EvenLines.scala:39)
18/03/05 01:14:20 INFO DAGScheduler: Got job 9 (print at EvenLines.scala:39) with 1 output partitions
18/03/05 01:14:20 INFO DAGScheduler: Final stage: ResultStage 15 (print at EvenLines.scala:39)
18/03/05 01:14:20 INFO DAGScheduler: Parents of final stage: List(ShuffleMapStage 14)
18/03/05 01:14:20 INFO DAGScheduler: Missing parents: List(ShuffleMapStage 14)
```

acadgild@localhost:~/IdeaProjects/sampleProject/target/scala-2.10 ▫️ ◻️ ✕

File  Edit  View  Search  Terminal  Help

```
18/03/05 01:14:20 INFO TaskSchedulerImpl: Removed TaskSet 15.0, whose tasks have all completed, from pool
18/03/05 01:14:20 INFO DAGScheduler: ResultStage 15 (print at EvenLines.scala:39) finished in 0.023 s
18/03/05 01:14:20 INFO DAGScheduler: Job 9 finished: print at EvenLines.scala:39, took 0.178395 s
-------------------------------------------
Time: 1520192660000 ms
-------------------------------------------
30.0

18/03/05 01:14:20 INFO JobScheduler: Finished job streaming job 1520192660000 ms.1 from job set of time 1520192660000 ms
18/03/05 01:14:20 INFO JobScheduler: Total delay: 0.387 s for time 1520192660000 ms (execution: 0.296 s)
18/03/05 01:14:20 INFO MapPartitionsRDD: Removing RDD 26 from persistence list
18/03/05 01:14:20 INFO BlockManager: Removing RDD 26
18/03/05 01:14:20 INFO BlockRDD: Removing RDD 25 from persistence list
18/03/05 01:14:20 INFO BlockManager: Removing RDD 25
18/03/05 01:14:20 INFO SocketInputDStream: Removing blocks of RDD BlockRDD[25] at socketTextStream at EvenLines.scala:25 of t
ime 1520192660000 ms
18/03/05 01:14:20 INFO MapPartitionsRDD: Removing RDD 30 from persistence list
18/03/05 01:14:20 INFO BlockManager: Removing RDD 30
18/03/05 01:14:20 INFO ShuffledRDD: Removing RDD 29 from persistence list
18/03/05 01:14:20 INFO BlockManager: Removing RDD 29
18/03/05 01:14:20 INFO MapPartitionsRDD: Removing RDD 28 from persistence list
18/03/05 01:14:20 INFO BlockManager: Removing RDD 28
18/03/05 01:14:20 INFO MapPartitionsRDD: Removing RDD 27 from persistence list
18/03/05 01:14:20 INFO BlockManager: Removing RDD 27
18/03/05 01:14:20 INFO ReceivedBlockTracker: Deleting batches ArrayBuffer(1520192640000 ms)
18/03/05 01:14:20 INFO InputInfoTracker: remove old batch metadata: 1520192640000 ms
18/03/05 01:14:29 INFO MemoryStore: Block input-0-1520192669600 stored as bytes in memory (estimated size 9.0 B, free 114.3 K
B)
18/03/05 01:14:29 INFO BlockManagerInfo: Added input-0-1520192669600 in memory on localhost:41378 (size: 9.0 B, free: 517.4 M
B)
18/03/05 01:14:29 INFO BlockGenerator: Pushed block input-0-1520192669600
18/03/05 01:14:30 INFO JobScheduler: Starting job streaming job 1520192670000 ms.0 from job set of time 1520192670000 ms
18/03/05 01:14:30 INFO JobScheduler: Added jobs for time 1520192670000 ms
18/03/05 01:14:30 INFO SparkContext: Starting job: print at EvenLines.scala:36
```

As we can see in the screenshots above, our program has picked up all the lines entered and filtered out those resulting in sum of odd number and has printed the rest.