20/5/25  Lab-8.

Aggregate

scala> val numbers = SC.parallelize (1 to 20, 4).
numbers: org.apache.spark.rdd.RDD[Int] =
Parallel collectionRDD[18] at parallelize at <console>: 24

scala> numbers.glom().collect()
res22. Array[Array[Int]] = Array (Array (1, 2, 3, 4, 5),
Array(6, 7, 8, 9, 10), Array(11, 12, 13, 14, 15), Array(16, 17, 18, 19, 20))

scala> numbers.aggregate (0) (_ + _, _ + _).
res23 : Int = 210

Fold operation.

Scala> numbers.fold (0) ((x, y) => x + y).

Scala> numbers.aggregate (1) (_ + _, _ + _).
res37 : Int = 215

scala> val numbers = SC.parallelize (1 to 4, 2).
numbers: org.apache.spark.rdd.RDD[Int] =
ParallelcollectionRDD[20] at parallelize at <console>: 24

Scala> numbers.glom.collect
res39 : Array[Array[Int]] = Array (Array (1, 2), Array(3, 4))

Scala> numbers.aggregate (0) (_ + _, _ * _).

Scala> val func = (x: Int) => x + x

```
Scala> 2func(3)
Res0 : Int = 6
```

**Lab 9.**

```
Scala> val sampleRDD = SC.parallelize (List(6,2,5,3,10,25,9,20,35))
sampleRDD: org.apache.spark.rdd.RDD[Int] = ParallelCollection
RDD[26] at parallelize at <console>:24
```

```
Scala> sampleRDD.collect
res55: Array[Int]: Array (6,2,5,3,10,25,9,20,35)
```

**Lab-10.**

D Write a scala program to print numbers from 1 to 100
using for loop.

```
-D sudo apt install scala.
-D nano PrintNumbers.sala
    object PrintNumbers {
        def main(args: Array[String]): Unit = {
            for (i <- 1 to 100){
                println(i) }
            }
        }
```

-D scalac PrintNumbers.Scala
-D Scala PrintNumbers.

-D using RDD and flatmap count how many times each word
appears in a file and write out a list of words whose
count is strictly greater than 4 using spark

→ mkdir -p wordcounterSout/src/main/scala
cd wordcountsout

→ nano build.sbt

name := "wordcountsout"
version := "0.1"
scalaversion := "2.12.17"
library Dependencies += Seq(
  "org.apache.spark" %% "spark-core" %, "3.3.1"
)
mainclass in compile := some("wordcountsout")

→ nano src/main/scala/wordcountSout.scala

→ echo "hello spark hello spark hello world spark
spark "> /home/bhoom/Desktop/wc.txt

→ sbt run.

→③ import org.apache hadoop.io.IntWritable
import org.apache.hadoop.io.Text
import org.apache.hadoop.map.reduce.Reducer
import java.util.*,
public class wordcount Reducer {
  private final IntWritable result new IntWritable()
  public void reduce() throws IO exception {
    int sum=0.
    for (Intwritable val:values){
      sum+: val.get();
    }

```
result.set(sum)
context.write(key, result);
    }
```

20/5/25