

Shayideep Sangam

Georgia state university – Fall 2018

ASSIGNMENT 8

In BDP-11-frequent items-1, we have discussed a-priori algorithm to identify frequent pairs. The Spark implementation we discussed has some flaw: we loop through all the elements contained by a RDD, which is not a good solution since RDD is a distributed collection of data elements, across different worker nodes in the cluster. In general, we should use transformations to specify what we would like to do for *each* element contained by the RDD, and let the cluster to manage the rest, we should never explicitly loop over its elements.

In this assignment, you will need to fix this issue. The details can be found in slide 42 – 45 from BDP-11-frequent items-1. As always, you can find our code discussed in the class from our course web site (the related text file is called `baskets.txt`). In the code, look for this comment, `===> the following has to be changed`, so you know where to replace/fix. Furthermore, *see slide 44 and slide 50 for hints*.

SOLUTION:

1. *all* your source code, with necessary comments so I know what is the part you have changed/fixed (90 points)

```
import java.util.*;
```

```
import org.apache.spark.mllib.fpm.AssociationRules;
```

```
import org.apache.spark.api.java.JavaPairRDD;
```

```
import org.apache.spark.api.java.JavaRDD;
```

```
import org.apache.spark.api.java.JavaSparkContext;
```

```
import org.apache.spark.api.java.function.FlatMapFunction;
```

```
import org.apache.spark.api.java.function.Function;
```

```
import org.apache.spark.api.java.function.Function2;
```

```
import org.apache.spark.api.java.function.PairFlatMapFunction;
```

```
import org.apache.spark.api.java.function.PairFunction;
```

```
import org.apache.spark.api.java.function.VoidFunction;
```

```

import org.apache.spark.broadcast.Broadcast;
import org.apache.spark.sql.Session;

import scala.Tuple2;

// this is the example in Chap 3, Example 3.6

public class SparkFrequency {

    private static final String FILE_URI = "C:/Users/Avinash/Downloads/baskets.txt";
    private static final double s = 0.5; // threshold = 50%

    public static void main(String[] args) {

        // initializing spark
        SparkSession spark = SparkSession.builder().config("spark.master", "local[*]").getOrCreate();
        JavaSparkContext sc = new JavaSparkContext(spark.sparkContext());
        sc.setLogLevel("WARN");

        // create RDD by using text files
        JavaRDD<String> baskets = sc.textFile(FILE_URI);
        System.out.println(baskets.take((int) baskets.count()).toString());

        // total number of baskets
        Broadcast<Long> basketCount = sc.broadcast(baskets.count());

        // organize basket content into integer array
        JavaRDD<Integer[]> basketContent = baskets.map(new Function<String, Integer[]>() {
            public Integer[] call(String line) throws Exception {
                String[] itemStr = line.substring(1, line.length() - 1).split(","); // get rid of { and }, and split
                Integer[] items = new Integer[itemStr.length];
            }
        });
    }
}

```

```

        for (int i = 0; i < itemStr.length; i++) {
            items[i] = new Integer(Integer.parseInt(itemStr[i].trim()));
        }
        return items;
    }
});

System.out.println("basketContent has [" + basketContent.count() + "] elements");
basketContent.foreach(new VoidFunction<Integer[]>() {
    public void call(Integer[] items) throws Exception {
        for (int i = 0; i < items.length; i++) {
            System.out.print(items[i] + " ");
        }
        System.out.println();
    }
});

// first pass: list all the items
JavaPairRDD<Integer, Integer> items = basketContent.flatMap(new FlatMapFunction<Integer[],
Integer>() {
    public Iterator<Integer> call(Integer[] items) throws Exception {
        return Arrays.asList(items).iterator();
    }
}).mapToPair(new PairFunction<Integer, Integer, Integer>() {
    public Tuple2<Integer, Integer> call(Integer item) throws Exception {
        return new Tuple2<Integer, Integer>(item, 1);
    }
});

System.out.println("items has [" + items.count() + "] elements");
System.out.println(items.take((int) items.count()).toString());

// first pass: count each item

```

```

JavaPairRDD<Integer, Integer> itemCounts = items.reduceByKey(new Function2<Integer,
Integer, Integer>() {

    public Integer call(final Integer value1, final Integer value2) {

        return value1 + value2;

    }

}).sortByKey();

System.out.println("itemCounts has [" + itemCounts.count() + "] elements");

System.out.println(itemCounts.take((int) itemCounts.count()).toString());


// first pass: create frequent-items table

// ==> the following has to be changed since this is not the best way to do this (see slide 43 of
BDP-11)

Broadcast<Double> sValue = sc.broadcast(s);

int[] frequentItems = new int[(int) (itemCounts.count() + 1)];

int threshold = (int) Math.ceil(s * basketCount.value());


//    filter item count to to get frequent items

Map<Integer, Integer> collect = itemCounts.filter(new Function<Tuple2<Integer, Integer>,
Boolean>() {

    @Override

    public Boolean call(Tuple2<Integer, Integer> integerIntegerTuple2) throws Exception {

        return integerIntegerTuple2._2 >= threshold;

    }

}).collectAsMap();

//    Collect the Frequent items

for (int i = 1; i <= itemCounts.count(); i++) {

    frequentItems[i] = collect.containsKey(i) ? collect.get(i) : 0;

}

```

```

// broadcast the table to avoid network traffic

// the size should be fine

System.out.println("frequentItems[i]:");

for (int i = 1; i < frequentItems.length; i++) {

    System.out.print("count[" + i + "] = " + frequentItems[i]);

    if (i < frequentItems.length - 1) {

        System.out.print(", ");

    } else System.out.println();

}

Broadcast<int[]> frequentItemTable = sc.broadcast(frequentItems);

// ==> end of the changed part


// second pass: generate frequent-pairs, start from basketContent

class FrequentPairsChecker implements PairFlatMapFunction<Integer[], String, Integer> {

    public Iterator<Tuple2<String, Integer>> call(Integer[] items) throws Exception {

        List<Tuple2<String, Integer>> frequentPairs = new ArrayList<>();

        int[] intItems = Arrays.stream(items).mapToInt(Integer::intValue).toArray();

        String itemStr = SparkAPrioriUtils.generatePairs(intItems);

        if (itemStr == null) return frequentPairs.iterator();

        String[] itemPairs = itemStr.split(",");

        for (String itemPair : itemPairs) {

            String[] tmpPair = itemPair.split("-");

            if (tmpPair == null || tmpPair.length != 2) continue;

            if (Integer.parseInt(tmpPair[0]) > Integer.parseInt(tmpPair[1])) {

                String tmpStr = tmpPair[0];

                tmpPair[0] = tmpPair[1];

                tmpPair[1] = tmpStr;
            }
        }

        frequentPairs.addAll(itemPairs);

        return frequentPairs.iterator();
    }
}

```

```

        tmpPair[1] = tmpStr;
    } // make sure the item pair is ordered, such as 123-234, not the other way around

    // check frequent-items table, to make sure both are frequent
    if (frequentItemTable.value()[Integer.parseInt(tmpPair[0])] > 0 &&
        frequentItemTable.value()[Integer.parseInt(tmpPair[1])] > 0) {
        frequentPairs.add(new Tuple2<>(tmpPair[0] + "-" + tmpPair[1], 1));
    }

}

return frequentPairs.iterator();
}
}

```

```

JavaPairRDD<String, Integer> candidateDoubles = basketContent.flatMapToPair(new
FrequentPairsChecker());

System.out.println("candidateDoubles has [" + candidateDoubles.count() + "] elements");
System.out.println(candidateDoubles.take((int) candidateDoubles.count()).toString());

```

```

JavaPairRDD<String, Integer> frequentPairs =
    candidateDoubles.reduceByKey(new Function2<Integer, Integer, Integer>() {
        public Integer call(final Integer value1, final Integer value2) {
            return value1 + value2;
        }
    }).filter(new Function<Tuple2<String, Integer>, Boolean>() {
        public Boolean call(Tuple2<String, Integer> frequentPair) {
            if (frequentPair._2 >= (int) Math.ceil(sValue.value() * basketCount.value()))
                return true;
            else return false;
        }
    })

```

```

    });

    System.out.println("final frequent pairs =>");

    System.out.println(frequentPairs.take((int) frequentPairs.count()).toString());

    frequentItemTable.unpersist();

    frequentItemTable.destroy();

    basketCount.unpersist();

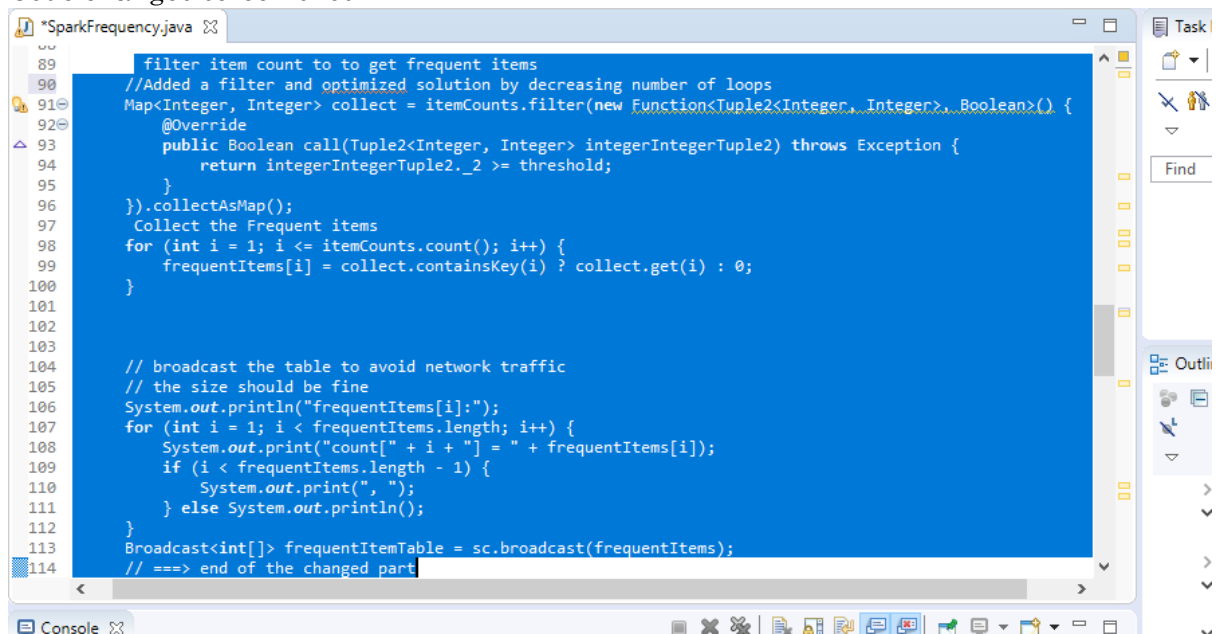
    basketCount.destroy();

    sc.close();
}

}

```

1. Code changes screen shot



```

*SparkFrequency.java
89      filter item count to to get frequent items
90      //Added a filter and optimized solution by decreasing number of loops
91      Map<Integer, Integer> collect = itemCounts.filter(new Function<Tuple2<Integer, Integer>, Boolean>() {
92          @Override
93          public Boolean call(Tuple2<Integer, Integer> integerIntegerTuple2) throws Exception {
94              return integerIntegerTuple2._2 >= threshold;
95          }
96      }).collectAsMap();
97      Collect the Frequent items
98      for (int i = 1; i <= itemCounts.count(); i++) {
99          frequentItems[i] = collect.containsKey(i) ? collect.get(i) : 0;
100      }
101
102
103
104      // broadcast the table to avoid network traffic
105      // the size should be fine
106      System.out.println("frequentItems[i]:");
107      for (int i = 1; i < frequentItems.length; i++) {
108          System.out.print("count[" + i + "] = " + frequentItems[i]);
109          if (i < frequentItems.length - 1) {
110              System.out.print(", ");
111          } else System.out.println();
112      }
113      Broadcast<int[]> frequentItemTable = sc.broadcast(frequentItems);
114      // ==> end of the changed part

```

2. an *acceptable* screen copy to show the results of running your code (10 points)

