**Program – 6**

# AIM: Implement Flajolet-Martin Algorithm for counting distinct elements in stream data.

## Introduction & Theory

### Flajolet-Martin Algorithm

The Flajolet–Martin algorithm is an algorithm for approximating the number of distinct elements in a stream with a single pass and space-consumption which is logarithmic in the maximum number of possible distinct elements in the stream.

1. Create a bit vector (bit array) of sufficient length L, such that 2L>n, the number of elements in the stream. Usually a 64-bit vector is sufficient since 264 is quite large for most purposes.

2. The i-th bit in this vector/array represents whether we have seen a hash function value whose binary representation ends in 0i. So initialize each bit to 0.

3. Generate a good, random hash function that maps input (usually strings) to natural numbers.

4. Read input. For each word, hash it and determine the number of trailing zeros. If the number of trailing zeros is k, set the k-th bit in the bit vector to 1.

5. Once input is exhausted, get the index of the first 0 in the bit array (call this R). By the way, this is just the number of consecutive 1s (i.e. we have seen 0, 00, ..., $0^{R-1}$ as the output of the hash function) plus one.

6. Calculate the number of unique words as $2R/\phi$, where $\phi$ is 0.77351. A proof for this can be found in the original paper listed in the reference section.

7. The standard deviation of R is a constant: $\sigma(R)=1.12$. (In other words, R can be off by about 1 for 1-0.68=32% of the observations, off by 2 for about 1-0.95=5% of the observations, off by 3 for 1-0.997=0.3% of the observations using the Empirical rule of statistics). This implies that our count can be off by a factor of 2 for 32% of the observations, off by a factory of 4 for 5% of the observations, off by a factor of 8 for 0.3% of the observations and so on.

Below are the steps for MapReduce data flow:

- **Step 1:** One block is processed by one **mapper** at a time. In the mapper, a developer can specify his own business logic as per the requirements. In this manner, Map runs on all the nodes of the cluster and process the data blocks in parallel.

- **Step 2:** Output of Mapper also known as intermediate output is written to the local disk. An output of mapper is not stored on HDFS as this is temporary data and **writing on HDFS** will create unnecessary many copies.

- **Step 3:** Output of mapper is shuffled to **reducer** node (which is a normal slave node but reduce phase will run here hence called as reducer node). The shuffling/copying is a physical movement of data which is done over the network.

# Program – 6

- **Step 4:** Once all the mappers are finished and their output is shuffled on reducer nodes then this intermediate output is merged & sorted. Which is then provided as input to reduce phase.

- **Step 5:** Reduce is the second phase of processing where the user can specify his own custom business logic as per the requirements. An input to a reducer is provided from all the mappers. An output of reducer is the final output, which is written on HDFS

## Code

FM algo Class

```java
package algos;

import java.io.IOException;
import java.util.Collection;
import java.util.Collections;
import java.util.List;
import java.util.ArrayList;

public class FlajoletMartin {

    private static final double PHI = 0.77351D;
    private int hashFunctions;
    private HashFunction[] hashes;
    public FlajoletMartin(int numHashFunctions) {
        this.hashFunctions = numHashFunctions;
        hashes = new HashFunction[hashFunctions];
        generateHashes();

    }

    private void generateHashes() {
        for (int i = 0; i < hashFunctions; i++) {
            hashes[i] = new HashFunction(i);
        }
    }

    public int uniques(String str) {
        List<Integer> averageR = new ArrayList<Integer>();
        String[] stream = str.split(" ");
        for (int i = 0; i < hashFunctions; i++){
            int sumR = 0;
            for (String word : stream) {
                sumR =
Math.max(getFirstZero(hashes[i].getHash(word)), sumR);
            }
            averageR.add(sumR);
        }

        Collections.sort(averageR);
        double r = 0;
        int avgMid = averageR.size() / 2;
        if (averageR.size() % 2 == 0) {
            r = (averageR.get(avgMid) + averageR.get(avgMid+1))/2;
        } else {
```

## Program – 6

```
45              r = averageR.get(avgMid + 1);
46          }
47          return (int) (Math.pow(2, r));
48
49      }
50
51      private int getFirstZero(int val) {
52          if (val == 0)
53              return 0;
54          int pos = 0;
55          // counting the position of first set bit
56          for (int i = 0; i < 64; i++) {
57              if ((val & (1 << i))== 0)
58                  pos++;
59
60              else
61                  break;
62          }
63          return pos;
64      }
65      private static class HashFunction {
66          private int numLoop;
67
68          public HashFunction(int loop) {
69              this.numLoop = loop;
70          }
71
72          public int getHash(Object s) {
73              int hash;
74              int i = 0;
75              do {
76                  if (s instanceof String) {
77                      hash = s.hashCode() * (i+1);
78                      s = Integer.toString(hash);
79                  } else if (s instanceof Number) {
80                      hash = String.valueOf(s).hashCode() * (i+1);
81                      s = Integer.toString(hash);
82                  }
83                  else {
84                      String t = s.toString();
85                      hash = t.hashCode() *(i+1);
86                      s = Integer.toString(hash);
87                  }
88                  i += 1;
89              } while (i < numLoop);
90              return 13 + 7*Math.abs(hash);
91          }
92
93      }
94
95  }
96
```

## Program – 6

Mapper

```
1   package hadoop;
2
3   import java.io.IOException;
4
5   import algos.*;
6   import org.apache.hadoop.io.IntWritable;
7   import org.apache.hadoop.io.LongWritable;
8   import org.apache.hadoop.io.Text;
9   import org.apache.hadoop.mapreduce.Mapper;
10
11  public class FMMapper extends Mapper<LongWritable, Text, Text,
12  IntWritable>{
13      private static final int EMPTY = 0;
14      FlajoletMartin fm;
15
16      @Override
17      protected void setup(org.apache.hadoop.mapreduce.Mapper.Context
18  context)
19              throws IOException, InterruptedException {
20          super.setup(context);
21          fm = new FlajoletMartin(4);
22      }
23
24      @Override
25      public void map(LongWritable key, Text value, Context context)
26          throws IOException, InterruptedException {
27          String stream = value.toString();
28          String id = stream.substring(0, 4);
29          stream = stream.substring(5);
30          int distinct = fm.uniques(stream);
31          context.write(new Text(id), new IntWritable(distinct));
32
33      }
34  }
35
```

Main

```
1   package hadoop;
2   import java.io.IOException;
3
4   import org.apache.hadoop.fs.Path;
5   import org.apache.hadoop.io.IntWritable;
6   import org.apache.hadoop.io.Text;
7   import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
8   import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
9   import org.apache.hadoop.mapreduce.Counter;
10  import org.apache.hadoop.mapreduce.CounterGroup;
11  import org.apache.hadoop.mapreduce.Counters;
12  import org.apache.hadoop.mapreduce.Job;
13
14  public class FMAlgo {
15      public static void main(String[] args)
16
```
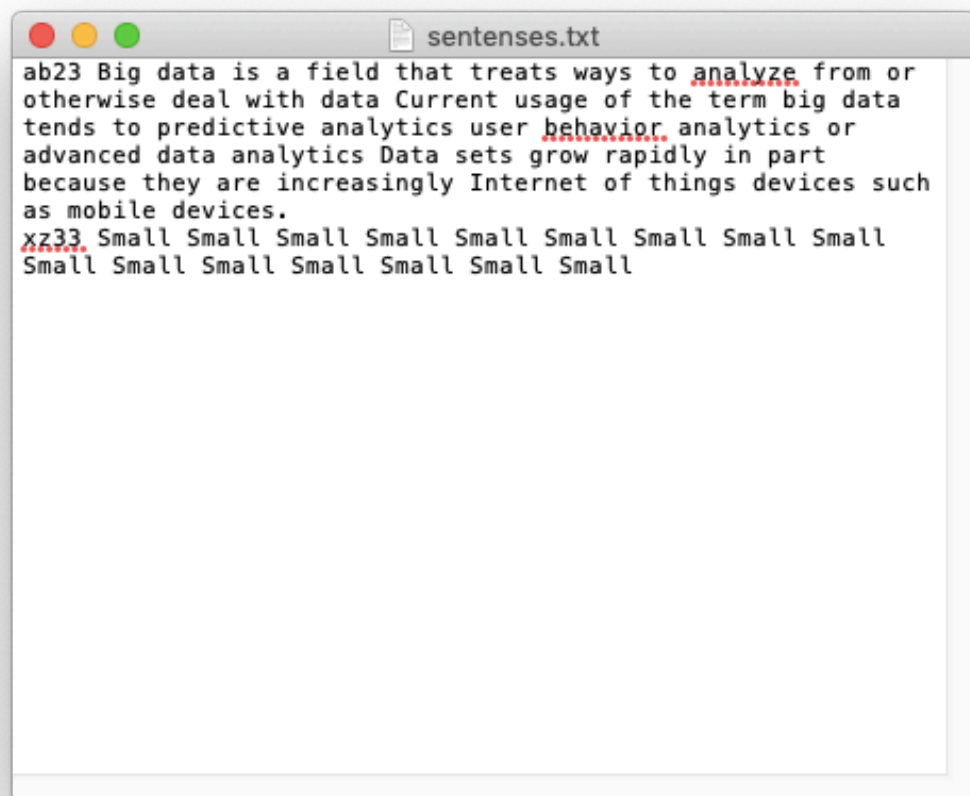
## Program – 6

```
17              throws IOException, ClassNotFoundException,
18   InterruptedException {
19          if (args.length != 2) {
20              System.err.println("Usage: FilterJob <input path> <output
21   path>");
22              System.exit(-1);
23          }
24          Job job = new Job();
25          job.setJarByClass(hadoop.FMMapper.class);
26          job.setJobName("Multi stream FM algo");
27
28          FileInputFormat.addInputPath(job, new Path(args[0]));
29          FileOutputFormat.setOutputPath(job, new Path(args[1]));
30
31          job.setMapperClass(hadoop.FMMapper.class);
32          job.setOutputKeyClass(Text.class);
33          job.setOutputValueClass(IntWritable.class);
34          job.waitForCompletion(true);
35      }
36   }
```

## Output



sentenses.txt

ab23 Big data is a field that treats ways to analyze from or
otherwise deal with data Current usage of the term big data
tends to predictive analytics user behavior analytics or
advanced data analytics Data sets grow rapidly in part
because they are increasingly Internet of things devices such
as mobile devices.
xz33 Small Small Small Small Small Small Small Small Small
Small Small Small Small Small Small Small

# Program – 6

```
● ● ●   hduser@rinzler-jarvis: ~/CodeLib/FMAlgo/src
hduser@rinzler-jarvis:~/CodeLib/FMAlgo/src$ hadoop jar FMA.jar hadoop.FMAlgo inp
ut/sentenses.txt output/fmalgo^C
hduser@rinzler-jarvis:~/CodeLib/FMAlgo/src$ hdfs dfs -copyToLocal output/fmalgo
~/output/fmalgo^C
hduser@rinzler-jarvis:~/CodeLib/FMAlgo/src$ cat ~/output/fmalgo/part-r-00000
ab23    16
xz33    2
hduser@rinzler-jarvis:~/CodeLib/FMAlgo/src$ █
```

## Findings and Learnings:

1. We learned how map-reduce works.
2. We learned how to code using hadoop in Java.
3. We have learned the working of the flajolet martin algorithm.

We have successfully implemented flajolet martin algorithm.