

## Program – 5

AIM: Write a program for bloom-filtering in Hadoop.

### Introduction & Theory

---

#### Bloom Filter

A Bloom filter is a space-efficient probabilistic data structure, conceived by Burton Howard Bloom in 1970, that is used to test whether an element is a member of a set. False positive matches are possible, but false negatives are not – in other words, a query returns either "possibly in set" or "definitely not in set". Elements can be added to the set, but not removed (though this can be addressed with a "counting" filter); the more elements that are added to the set, the larger the probability of false positives

Unlike a standard hash table using open addressing for collision resolution, a Bloom filter of a fixed size can represent a set with an arbitrarily large number of elements; adding an element never fails due to the data structure "filling up". However, the false positive rate increases steadily as elements are added until all bits in the filter are set to 1, at which point all queries yield a positive result. With open addressing hashing, false positives are never produced, but performance steadily deteriorates until it approaches linear search.

Union and intersection of Bloom filters with the same size and set of hash functions can be implemented with bitwise OR and AND operations respectively. The union operation on Bloom filters is lossless in the sense that the resulting Bloom filter is the same as the Bloom filter created from scratch using the union of the two sets. The intersect operation satisfies a weaker property: the false positive probability in the resulting Bloom filter is at most the false-positive probability in one of the constituent Bloom filters, but may be larger than the false positive probability in the Bloom filter created from scratch using the intersection of the two sets. (interesting properties: source Wikipedia)

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.
- **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

## Program – 5

### Code

#### Bloom Filter Class

```
1  package utils;
2
3  import java.io.Serializable;
4  import java.nio.charset.Charset;
5  import java.security.MessageDigest;
6  import java.security.NoSuchAlgorithmException;
7  import java.util.BitSet;
8  import java.util.Collection;
9
10
11 public class FilterBloom<E> implements Serializable {
12     private BitSet bitset;
13     private int bitSetSize;
14     private double bitsPerElement;
15     private int expectedNumberOfFilterElements; // expected (maximum)
16     number of elements to be added
17     private int numberOfAddedElements; // number of elements actually
18     added to the Bloom filter
19     private int k; // number of hash functions
20
21     static final Charset charset = Charset.forName("UTF-8"); //
22     encoding used for storing hash values as strings
23
24     static final String hashName = "MD5"; // MD5 gives good enough
25     accuracy in most circumstances. Change to SHA1 if it's needed
26     static final MessageDigest digestFunction;
27     static { // The digest method is reused between instances
28         MessageDigest tmp;
29         try {
30             tmp = java.security.MessageDigest.getInstance(hashName);
31         } catch (NoSuchAlgorithmException e) {
32             tmp = null;
33         }
34         digestFunction = tmp;
35     }
36
37     /**
38      * Constructs an empty Bloom filter. The total length of the
39      * Bloom filter will be
40      * c*n.
41      *
42      * @param c is the number of bits used per element.
43      * @param n is the expected number of elements the filter will
44      * contain.
45      * @param k is the number of hash functions used.
46      */
47     public FilterBloom(double c, int n, int k) {
48         this.expectedNumberOfFilterElements = n;
49         this.k = k;
50         this.bitsPerElement = c;
51         this.bitSetSize = (int) Math.ceil(c * n);
52         numberOfAddedElements = 0;
53         this.bitset = new BitSet(bitSetSize);
```

## Program – 5

```
54     }
55
56     /**
57      * Constructs an empty Bloom filter. The optimal number of hash
58      functions (k) is estimated from the total size of the Bloom
59      * and the number of expected elements.
60      *
61      * @param bitSetSize defines how many bits should be used in
62      total for the filter.
63      * @param expectedNumberOElements defines the maximum number of
64      elements the filter is expected to contain.
65      */
66     public FilterBloom(int bitSetSize, int expectedNumberOElements) {
67         this(bitSetSize / (double)expectedNumberOElements,
68             expectedNumberOElements,
69             (int) Math.round((bitSetSize /
70 (double)expectedNumberOElements) * Math.log(2.0)));
71     }
72
73     /**
74      * Constructs an empty Bloom filter with a given false positive
75      probability. The number of bits per
76      * element and the number of hash functions is estimated
77      * to match the false positive probability.
78      *
79      * @param falsePositiveProbability is the desired false positive
80      probability.
81      * @param expectedNumberOfElements is the expected number of
82      elements in the Bloom filter.
83      */
84     public FilterBloom(double falsePositiveProbability, int
85 expectedNumberOfElements) {
86         this(Math.ceil(-(Math.log(falsePositiveProbability) /
87 Math.log(2))) / Math.log(2), // c = k / ln(2)
88             expectedNumberOfElements,
89             (int) Math.ceil(-(Math.log(falsePositiveProbability) /
90 Math.log(2)))); // k = ceil(-log_2(false prob.))
91     }
92
93     /**
94      * Construct a new Bloom filter based on existing Bloom filter
95      data.
96      *
97      * @param bitSetSize defines how many bits should be used for the
98      filter.
99      * @param expectedNumberOfFilterElements defines the maximum
100      number of elements the filter is expected to contain.
101      * @param actualNumberOfFilterElements specifies how many
102      elements have been inserted into the <code>filterData</code> BitSet.
103      * @param filterData a BitSet representing an existing Bloom
104      filter.
105      */
106     public FilterBloom(int bitSetSize, int
107 expectedNumberOfFilterElements, int actualNumberOfFilterElements,
108 BitSet filterData) {
109         this(bitSetSize, expectedNumberOfFilterElements);
110         this.bitset = filterData;
```

## Program – 5

```
111         this.numberOfAddedElements = actualNumberOfFilterElements;
112     }
113
114     /**
115      * Generates a digest based on the contents of a String.
116      *
117      * @param val specifies the input data.
118      * @param charset specifies the encoding of the input data.
119      * @return digest as long.
120      */
121     public static int createHash(String val, Charset charset) {
122         return createHash(val.getBytes(charset));
123     }
124
125     /**
126      * Generates a digest based on the contents of a String.
127      *
128      * @param val specifies the input data. The encoding is expected
129 to be UTF-8.
130      * @return digest as long.
131      */
132     public static int createHash(String val) {
133         return createHash(val, charset);
134     }
135
136     /**
137      * Generates a digest based on the contents of an array of bytes.
138      *
139      * @param data specifies input data.
140      * @return digest as long.
141      */
142     public static int createHash(byte[] data) {
143         return createHashes(data, 1)[0];
144     }
145
146     /**
147      * Generates digests based on the contents of an array of bytes
148 and splits the result into 4-byte int's and store them in an array.
149 The
150      * digest function is called until the required number of int's
151 are produced. For each call to digest a salt
152      * is prepended to the data. The salt is increased by 1 for each
153 call.
154      *
155      * @param data specifies input data.
156      * @param hashes number of hashes/int's to produce.
157      * @return array of int-sized hashes
158      */
159     public static int[] createHashes(byte[] data, int hashes) {
160         int[] result = new int[hashes];
161
162         int k = 0;
163         byte salt = 0;
164         while (k < hashes) {
165             byte[] digest;
166             synchronized (digestFunction) {
167                 digestFunction.update(salt);
```

## Program – 5

```
168         salt++;
169         digest = digestFunction.digest(data);
170     }
171
172     for (int i = 0; i < digest.length/4 && k < hashes; i++) {
173         int h = 0;
174         for (int j = (i*4); j < (i*4)+4; j++) {
175             h <= 8;
176             h |= ((int) digest[j]) & 0xFF;
177         }
178         result[k] = h;
179         k++;
180     }
181 }
182 return result;
183 }
184
185 /**
186  * Compares the contents of two instances to see if they are
187  * equal.
188  *
189  * @param obj is the object to compare to.
190  * @return True if the contents of the objects are equal.
191  */
192 @Override
193 public boolean equals(Object obj) {
194     if (obj == null) {
195         return false;
196     }
197     if (getClass() != obj.getClass()) {
198         return false;
199     }
200     final FilterBloom<E> other = (FilterBloom<E>) obj;
201     if (this.expectedNumberOfFilterElements !=
202 other.expectedNumberOfFilterElements) {
203         return false;
204     }
205     if (this.k != other.k) {
206         return false;
207     }
208     if (this.bitSetSize != other.bitSetSize) {
209         return false;
210     }
211     if (this.bitset != other.bitset && (this.bitset == null ||
212 !this.bitset.equals(other.bitset))) {
213         return false;
214     }
215     return true;
216 }
217
218 /**
219  * Calculates a hash code for this class.
220  * @return hash code representing the contents of an instance of
221  * this class.
222  */
223 @Override
224 public int hashCode() {
```

## Program – 5

```
225         int hash = 7;
226         hash = 61 * hash + (this.bitset != null ?
227 this.bitset.hashCode() : 0);
228         hash = 61 * hash + this.expectedNumberOfFilterElements;
229         hash = 61 * hash + this.bitSetSize;
230         hash = 61 * hash + this.k;
231         return hash;
232     }
233
234
235     /**
236      * Calculates the expected probability of false positives based
237 on
238      * the number of expected filter elements and the size of the
239 Bloom filter.
240      * <br /><br />
241      * The value returned by this method is the <i>expected</i> rate
242 of false
243      * positives, assuming the number of inserted elements equals the
244 number of
245      * expected elements. If the number of elements in the Bloom
246 filter is less
247      * than the expected value, the true probability of false
248 positives will be lower.
249      *
250      * @return expected probability of false positives.
251      */
252     public double expectedFalsePositiveProbability() {
253         return
254 getFalsePositiveProbability(expectedNumberOfFilterElements);
255     }
256
257     /**
258      * Calculate the probability of a false positive given the
259 specified
260      * number of inserted elements.
261      *
262      * @param numberOfElements number of inserted elements.
263      * @return probability of a false positive.
264      */
265     public double getFalsePositiveProbability(double
266 numberOfElements) {
267         // (1 - e^(-k * n / m)) ^ k
268         return Math.pow((1 - Math.exp(-k * (double) numberOfElements
269 / (double) bitSetSize)), k);
270     }
271 }
272
273     /**
274      * Get the current probability of a false positive. The
275 probability is calculated from
276      * the size of the Bloom filter and the current number of
277 elements added to it.
278      *
279      * @return probability of false positives.
280      */
281     public double getFalsePositiveProbability() {
```

## Program – 5

```
282         return getFalsePositiveProbability(numberOfAddedElements);
283     }
284
285
286     /**
287      * Returns the value chosen for K.<br />
288      * <br />
289      * K is the optimal number of hash functions based on the size
290      * of the Bloom filter and the expected number of inserted
291 elements.
292      *
293      * @return optimal k.
294      */
295     public int getK() {
296         return k;
297     }
298
299     /**
300      * Sets all bits to false in the Bloom filter.
301      */
302     public void clear() {
303         bitset.clear();
304         numberOfAddedElements = 0;
305     }
306
307     /**
308      * Adds an object to the Bloom filter. The output from the
309 object's
310      * toString() method is used as input to the hash functions.
311      *
312      * @param element is an element to register in the Bloom filter.
313      */
314     public void add(E element) {
315         add(element.toString().getBytes(charset));
316     }
317
318     /**
319      * Adds an array of bytes to the Bloom filter.
320      *
321      * @param bytes array of bytes to add to the Bloom filter.
322      */
323     public void add(byte[] bytes) {
324         int[] hashes = createHashes(bytes, k);
325         for (int hash : hashes)
326             bitset.set(Math.abs(hash % bitSetSize), true);
327         numberOfAddedElements++;
328     }
329
330     /**
331      * Adds all elements from a Collection to the Bloom filter.
332      * @param c Collection of elements.
333      */
334     public void addAll(Collection<? extends E> c) {
335         for (E element : c)
336             add(element);
337     }
338
```

## Program – 5

```
339     /**
340      * Returns true if the element could have been inserted into the
341      Bloom filter.
342      * Use getFalsePositiveProbability() to calculate the probability
343      of this
344      * being correct.
345      *
346      * @param element element to check.
347      * @return true if the element could have been inserted into the
348      Bloom filter.
349      */
350     public boolean contains(E element) {
351         return contains(element.toString().getBytes(charset));
352     }
353
354     /**
355      * Returns true if the array of bytes could have been inserted
356      into the Bloom filter.
357      * Use getFalsePositiveProbability() to calculate the probability
358      of this
359      * being correct.
360      *
361      * @param bytes array of bytes to check.
362      * @return true if the array could have been inserted into the
363      Bloom filter.
364      */
365     public boolean contains(byte[] bytes) {
366         int[] hashes = createHashes(bytes, k);
367         for (int hash : hashes) {
368             if (!bitset.get(Math.abs(hash % bitSetSize))) {
369                 return false;
370             }
371         }
372         return true;
373     }
374
375     /**
376      * Returns true if all the elements of a Collection could have
377      been inserted
378      * into the Bloom filter. Use getFalsePositiveProbability() to
379      calculate the
380      * probability of this being correct.
381      * @param c elements to check.
382      * @return true if all the elements in c could have been inserted
383      into the Bloom filter.
384      */
385     public boolean containsAll(Collection<? extends E> c) {
386         for (E element : c)
387             if (!contains(element))
388                 return false;
389         return true;
390     }
391
392     /**
393      * Read a single bit from the Bloom filter.
394      * @param bit the bit to read.
395      * @return true if the bit is set, false if it is not.
```



## Program – 5

```
396     */
397     public boolean getBit(int bit) {
398         return bitset.get(bit);
399     }
400
401     /**
402      * Set a single bit in the Bloom filter.
403      * @param bit is the bit to set.
404      * @param value If true, the bit is set. If false, the bit is
405 cleared.
406     */
407     public void setBit(int bit, boolean value) {
408         bitset.set(bit, value);
409     }
410
411     /**
412      * Return the bit set used to store the Bloom filter.
413      * @return bit set representing the Bloom filter.
414     */
415     public BitSet getBitSet() {
416         return bitset;
417     }
418
419     /**
420      * Returns the number of bits in the Bloom filter. Use count() to
421 retrieve
422      * the number of inserted elements.
423      *
424      * @return the size of the bitset used by the Bloom filter.
425     */
426     public int size() {
427         return this.bitSetSize;
428     }
429
430     /**
431      * Returns the number of elements added to the Bloom filter after
432 it
433      * was constructed or after clear() was called.
434      *
435      * @return number of elements added to the Bloom filter.
436     */
437     public int count() {
438         return this.numberOfAddedElements;
439     }
440
441     /**
442      * Returns the expected number of elements to be inserted into
443 the filter.
444      * This value is the same value as the one passed to the
445 constructor.
446      *
447      * @return expected number of elements.
448     */
449     public int getExpectedNumberOfElements() {
450         return expectedNumberOfFilterElements;
451     }
452
```

## Program – 5

```
453     /**
454      * Get expected number of bits per element when the Bloom filter
455      is full. This value is set by the constructor
456      * when the Bloom filter is created. See also
457      getBitsPerElement().
458      *
459      * @return expected number of bits per element.
460      */
461     public double getExpectedBitsPerElement() {
462         return this.bitsPerElement;
463     }
464
465     /**
466      * Get actual number of bits per element based on the number of
467      elements that have currently been inserted and the length
468      * of the Bloom filter. See also getExpectedBitsPerElement().
469      *
470      * @return number of bits per element.
471      */
472     public double getBitsPerElement() {
473         return this.bitSetSize / (double)numberOfAddedElements;
474     }
475 }
```

## Mapper

```
1  package bfcode;
2  import java.lang.System.*;
3  import java.io.IOException;
4  import org.apache.hadoop.io.LongWritable;
5  import org.apache.hadoop.io.NullWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Mapper;
8  import utils.*;
9  public class FilterMapper extends Mapper<LongWritable, Text, Text,
10 NullWritable> {
11     FilterBloom<String> filter;
12     @Override
13     protected void setup(org.apache.hadoop.mapreduce.Mapper.Context
14 context)
15         throws IOException, InterruptedException {
16         super.setup(context);
17         double falsePositiveProbability = 0.1;
18         int expectedNumberOfElements = 100;
19         filter = new FilterBloom<String>(falsePositiveProbability,
20 expectedNumberOfElements);
21         filter.add("bad service");
22         filter.add("iron man");
23         filter.add("marvel");
24         filter.add("end game");
25     }
26     protected void map(LongWritable key, Text value, Context context)
27         throws java.io.IOException, InterruptedException {
28
29         String[] tokens = value.toString().split(",");
30         for(String token :tokens){
```

## Program – 5

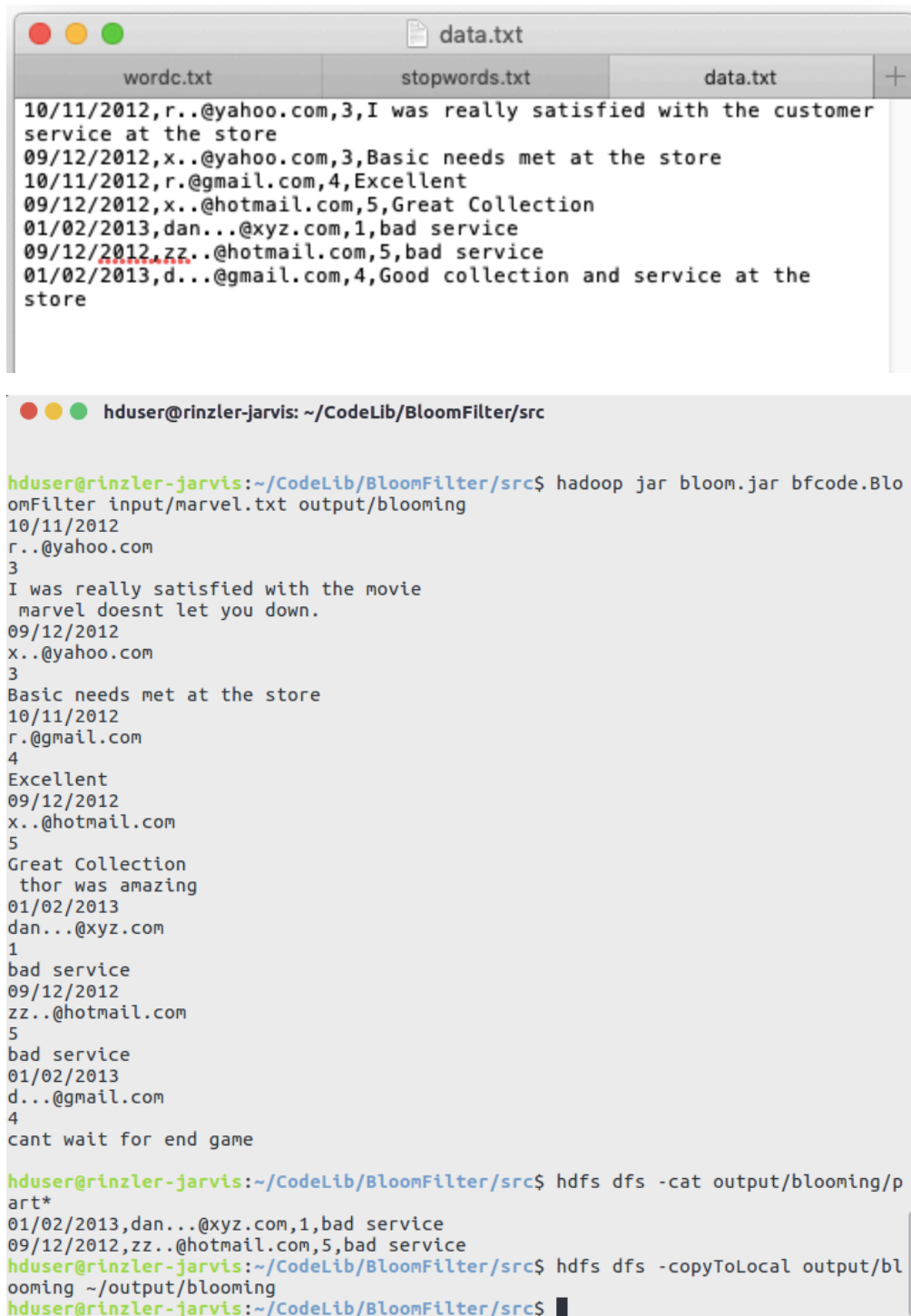
```
31         System.out.println(token);
32         if(filter.contains(token)){
33             context.write(value, NullWritable.get());
34         }
35     }
36 }
37 }
```

### Main

```
1  package bfcode;
2  import java.io.IOException;
3  import org.apache.hadoop.fs.Path;
4  import org.apache.hadoop.io.IntWritable;
5  import org.apache.hadoop.io.NullWritable;
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 BloomFilter {
15
16     public static void main(String[] args)
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
25         Job job = new Job();
26         job.setJarByClass(bfcode.FilterMapper.class);
27         job.setJobName("Customer Complaint Filter");
28
29         FileInputFormat.addInputPath(job, new Path(args[0]));
30         FileOutputFormat.setOutputPath(job, new Path(args[1]));
31
32         job.setMapperClass(bfcode.FilterMapper.class);
33         job.setOutputKeyClass(Text.class);
34         job.setOutputValueClass(NullWritable.class);
35         job.waitForCompletion(true);
36     }
37 }
```

## Program – 5

### Output



The image shows a text editor window with three tabs: wordc.txt, stopwords.txt, and data.txt. The data.txt tab is active, displaying a list of customer reviews. Below the text editor is a terminal window showing the execution of Hadoop and HDFS commands to process the data using a Bloom filter.

```
data.txt
wordc.txt  stopwords.txt  data.txt  +
10/11/2012,r..@yahoo.com,3,I was really satisfied with the customer
service at the store
09/12/2012,x..@yahoo.com,3,Basic needs met at the store
10/11/2012,r.@gmail.com,4,Excellent
09/12/2012,x..@hotmail.com,5,Great Collection
01/02/2013,dan...@xyz.com,1,bad service
09/12/2012,zz..@hotmail.com,5,bad service
01/02/2013,d...@gmail.com,4,Good collection and service at the
store

hduser@rinzler-jarvis: ~/CodeLib/BloomFilter/src

hduser@rinzler-jarvis:~/CodeLib/BloomFilter/src$ hadoop jar bloom.jar bfcode.Blo
omFilter input/marvel.txt output/blooming
10/11/2012
r..@yahoo.com
3
I was really satisfied with the movie
marvel doesnt let you down.
09/12/2012
x..@yahoo.com
3
Basic needs met at the store
10/11/2012
r.@gmail.com
4
Excellent
09/12/2012
x..@hotmail.com
5
Great Collection
thor was amazing
01/02/2013
dan...@xyz.com
1
bad service
09/12/2012
zz..@hotmail.com
5
bad service
01/02/2013
d...@gmail.com
4
cant wait for end game

hduser@rinzler-jarvis:~/CodeLib/BloomFilter/src$ hdfs dfs -cat output/blooming/p
art*
01/02/2013,dan...@xyz.com,1,bad service
09/12/2012,zz..@hotmail.com,5,bad service
hduser@rinzler-jarvis:~/CodeLib/BloomFilter/src$ hdfs dfs -copyToLocal output/bl
ooming ~/output/blooming
hduser@rinzler-jarvis:~/CodeLib/BloomFilter/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 bloom filters
4. We have successfully implemented bloom filters