# BDPA winter 2017: Assignment 2

## Rajaa EL HAMDANI

## March 17, 2017

# 1 Pre-processing the input

## 1.1 Word counts job

### 1.1.1 Mapper

The Mapper tokenizes each line and ouput pairs where the keys are tokens and values are 1.

```java
public class WordCountMapper extends Mapper<LongWritable, Text, Text,
    IntWritable> {
    private Text word = new Text();
    private final static IntWritable ONE = new IntWritable(1);

    @Override
    public void map(LongWritable key, Text value, Context context)
            throws IOException, InterruptedException {

        String line = value.toString().toLowerCase();
        StringTokenizer tokenizer = new StringTokenizer(line);

        while (tokenizer.hasMoreTokens()) {
            word.set(tokenizer.nextToken());
            context.write(word, ONE);
        }
    }
}
```

### 1.1.2 Reducer

The reducer count the frequency of each word by summing over the received values.

```java
public class WordCountReducer extends Reducer<Text, IntWritable, Text,
    IntWritable> {

    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
        throws IOException, InterruptedException {

        int sum = 0;
```

```
 8          for (IntWritable val : values) {
 9              sum += val.get();
10          }
11          context.write(key, new IntWritable(sum));
12      }
13  }
```

### 1.1.3 Driver

```
 1  public class WordCount {
 2
 3    public static void main(String[] args) throws Exception {
 4
 5      if (args.length != 2) {
 6        System.out.printf("Usage: WordCount <input dir> <output dir>\n");
 7        System.exit(-1);
 8      }
 9
10      Job job = new Job();
11      job.setJarByClass(WordCount.class);
12      job.setJobName("Word count");
13
14      FileInputFormat.addInputPath(job, new Path(args[0]));
15      FileOutputFormat.setOutputPath(job, new Path(args[1]));
16
17      job.setOutputKeyClass(Text.class);
18      job.setOutputValueClass(IntWritable.class);
19
20      job.setMapperClass(WordCountMapper.class);
21      job.setCombinerClass(WordCountReducer.class);
22      job.setReducerClass(WordCountReducer.class);
23
24      System.exit(job.waitForCompletion(true) ? 0:1);
25    }
26  }
```

## 1.2 Preprocessing job

To carry the task of preprocessing, I wrote two MapReduce jobs. The first job counts the input words. The second job performs the preprocessing.

### 1.2.1 Mapper

The mapper performs the following tasks:
   - Loading the stop-words file (generated in the 1st assignment) in a HashSet.

```
 1  HashSet<String> stopwords = new HashSet<String>();
 2          BufferedReader Reader = new BufferedReader(new FileReader(new File("/
     home/cloudera/stopwords.csv")));
 3          String line;
 4          while ((line = Reader.readLine()) != null) {
 5              stopwords.add(line.split("\\s+")[0].toLowerCase());
 6          }
```

- Removing stopwords. - Removing words that contain special characters by matching each word with a regular expression.

```
String lineF = value.toString().toLowerCase();
StringTokenizer tokenizer = new StringTokenizer(lineF);

Pattern p = Pattern.compile("[^A-Za-z0-9]");
String token;
while (tokenizer.hasMoreTokens()) {
  token = tokenizer.nextToken().toLowerCase();
  if (!(stopwords.contains(token) || p.matcher(token).find() || token.
isEmpty())){
    word.set(token);
      context.write(key, word);
  }
    }
```

- Removing empty lines.

### 1.2.2  Reducer

The reducer performs the following tasks: - Remove repetition of words by loading the words of its line input in a HashSet.

```
ArrayList<String> tokens = new ArrayList<String>();

for (Text word : values) {
  tokens.add(word.toString());
}

HashSet<String> tokensU = new HashSet<String>(tokens);
```

   - Sort the words, of the line, by their global frequency. This is done by the following steps:

1. Loading the word count file (output of the WordCount job).

2. Storing the word counts (only for the words that appear in the current line) in a HashMap, where the keys are the word and values are their frequencies.

3. Sorting the HashMap by its values in ascending order.

```
BufferedReader reader = new BufferedReader(new FileReader(new File("/home/
    cloudera/workspace/SetSimJoins/wordcount.txt")));

    Map<String, Integer> wordcount = new HashMap<String, Integer>(); //
    HashMap table to store words as keys and their frequency as vales
    String line;
    while ((line = reader.readLine()) != null) {
      String[] word = line.split("\\s+");
      if (tokensU.contains(word[0])){
        wordcount.put(word[0], Integer.parseInt(word[1]));
      }
```

```
10
11          }
12          // Sort table by values
13          Map<String, Integer> sortedMap = sortByValue(wordcount);
14              /* sortByValue is a function that sort a hashMap by its values.
      Its code is in          the appendix
15              */
```

- Writing sorted words in a string buffer

```
1  // Write the ordered words in a StringBuffer
2          StringBuffer bf = new StringBuffer();
3          for (Entry<String, Integer> entry : sortedMap.entrySet()) {
4            if(bf.length()!=0){
5              bf.append(" ");
6            }
7            bf.append(entry.getKey());
8          }
```

- Outputing the line in the HDFS (only if it not empty), and incrementing the counter of output records.

```
1
2          if(bf.length()!=0){
3            context.getCounter(RECORDS_COUNTER.NB_RECORDS).increment(1);
4          context.write(key, new Text(bf.toString()));
5          }
```

The variable RECORDS_COUNTER is instantiated in the driver.

### 1.2.3   Driver

The driver performs the following tasks:
- Instantiation and configuration of the MapReduce Job.

```
1   public static void main(String[] args) throws IOException,
     ClassNotFoundException, InterruptedException {
2     /*
3        * Validate that two arguments were passed from the command line.
4        */
5     if (args.length != 2) {
6       System.out.printf("Usage: StopWords <input dir> <output dir>\n");
7       System.exit(-1);
8     }
9
10      /*
11       * Instantiate a Job object for your job's configuration.
12       */
13     Configuration conf = new Configuration();
14     Job job = new Job(conf);
15     job.setJarByClass(Preprocessing.class);
16     job.setJobName("Part1_Preprocessing");
17
18     FileInputFormat.addInputPath(job, new Path(args[0]));
19     FileOutputFormat.setOutputPath(job, new Path(args[1]));
20
```

```
21        job.setMapperClass(PreprocessingMapper.class);
22        job.setReducerClass(PreprocessingReducer.class);
23
24        job.setOutputKeyClass(LongWritable.class);
25        job.setOutputValueClass(Text.class);
26
27        job.setOutputFormatClass(TextOutputFormat.class);
28
29        job.getConfiguration().set(
30          "mapreduce.output.textoutputformat.separator", "; ");
31
32        FileSystem fs = FileSystem.get(new Configuration());
33        if (fs.exists(new Path(args[1]))) {
34          fs.delete(new Path(args[1]));
35        }
36
37        job.waitForCompletion(true);
```

- Instantiation of the counter of the output records.

```
1  public static enum RECORDS_COUNTER {
2        NB_RECORDS,
3    };
```

- Writing the value of the counter in HDFS.

```
1
2        long counter = job.getCounters().findCounter(RECORDS_COUNTER.NB_RECORDS)
3          .getValue();
4      Path counterFile = new Path("NB_RECORDS.txt");
5      BufferedWriter bf = new BufferedWriter(new OutputStreamWriter(
6        fs.create(counterFile, true)));
7      bf.write(String.valueOf(counter));
8      bf.close();
9
10     System.exit(0);
```

| | Job Overview |
|---|---|
| **Job Name:** | Part1_Preprocessing |
| **User Name:** | cloudera |
| **Queue:** | root.cloudera |
| **State:** | SUCCEEDED |
| **Uberized:** | false |
| **Submitted:** | Fri Mar 17 08:15:24 PDT 2017 |
| **Started:** | Fri Mar 17 08:15:41 PDT 2017 |
| **Finished:** | Fri Mar 17 08:16:08 PDT 2017 |
| **Elapsed:** | 27sec |
| **Diagnostics:** | |
| **Average Map Time** | 11sec |
| **Average Shuffle Time** | 7sec |
| **Average Merge Time** | 0sec |
| **Average Reduce Time** | 2sec |

| ApplicationMaster | | | |
|---|---|---|---|
| Attempt Number | Start Time | Node | Logs |
| 1 | Fri Mar 17 08:15:33 PDT 2017 | quickstart.cloudera:8042 | logs |

| Task Type | Total | Complete |
|---|---|---|
| **Map** | 1 | 1 |
| **Reduce** | 1 | 1 |

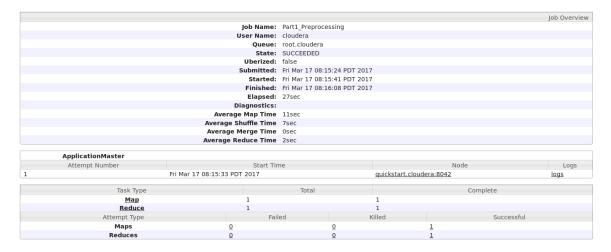| Attempt Type | Failed | Killed | Successful |
|---|---|---|---|
| **Maps** | 0 | 0 | 1 |
| **Reduces** | 0 | 0 | 1 |

Figure 1: Logs in Yarn of the pre-processing job

# 2 Set-similarity joins

## 2.1 First approach

In the first approach we perform all the pair-wise comparison between documents. To make sure that the same pair of documents is only compared once I implemented a custom WritablComparable class `DocPair` (the code is in the appendix).

### 2.1.1 Mapper

The mapper performs the following tasks:
- Read the pre-processed file.
- Associate the current document id with the ids of the rest of documents. The result is stored in a DocPair object and output as a key from the mapper.

```java
    @Override
    public void map(LongWritable key, Text value, Context context)
            throws IOException, InterruptedException {

        reader = new BufferedReader(new FileReader(new File("/home/cloudera/
    workspace/SetSimJoins/preprocessing_output_sample.txt")));
        String line;
        System.out.println(value.toString());
        String[] valueS = value.toString().split(";");
        while ((line = reader.readLine()) != null) {
            String key2 = line.split(";")[0];
            if (!key.toString().equals(key2)){
                pairKeys.set(new Text(valueS[0]), new Text(key2));
                context.write(pairKeys, new Text(valueS[1]));
            }
        }
    }
```

### 2.1.2 Reducer

The reducer performs the following tasks:—— - Computation of the Jaccard similarity between pairs of documents.

```java
public class SetsimjoinsReducer extends
    Reducer<DocPair, Text, Text, Text> {

    private BufferedReader reader;


    @Override
    public void reduce(DocPair key, Iterable<Text> values, Context context)
            throws IOException, InterruptedException {

        HashMap<String, String> allLines = new HashMap<String, String>();
```

```
12        reader = new BufferedReader(new FileReader(new File("/home/cloudera/
     workspace/SetSimJoins/preprocessing_output_sample.txt")));
13        String line;
14        while ((line = reader.readLine()) != null) {
15            String[] lineS = line.split(";");
16            allLines.put(lineS[0], lineS[1]);
17        }
18
19        HashSet<String> words1 = new HashSet<String>();
20
21        for (String word : values.iterator().next().toString().split(" ")) {
22          words1.add(word);
23        }
24
25        HashSet<String> words2 = new HashSet<String>();
26        String doc2 = allLines.get(key.getSecond()
27                  .toString());
28        for (String word : doc2.split(" ")) {
29          words2.add(word);
30        }
31        double sim = jaccardsim(words2, words2);
32
33
34 }
35 }
```

-Output the similar pairs in HDFS.

```
1          if (sim >= 0.8) {
2              context.write(new Text("(" + key.getFirst() + ", " + key.getSecond
     () + ")"),
3                        new Text(String.valueOf(sim)));
4          }
```

- Increment the number of comparisons. (The `COUNTER` variable is instantiated in the driver)

```
1  context.getCounter(COUNTER.NB_COMPARISIONS_I).increment(1);
```

### 2.1.3  Driver

The driver performs the following tasks:
- Instantiation and configuration of the MapReduce Job.

```
1 public class Setsimjoins {
2
3   public static void main(String[] args) throws IOException,
     ClassNotFoundException, InterruptedException {
4
5        if (args.length != 2) {
6          System.out.printf("Usage: StopWords <input dir> <output dir>\n");
7          System.exit(-1);
8        }
9
10        Configuration conf = new Configuration();
11        Job job = new Job(conf);
```

```
12        job.setJarByClass(Setsimjoins.class);
13        job.setJobName("Part2_Set Similarity Joins");
14
15        FileInputFormat.addInputPath(job, new Path(args[0]));
16        FileOutputFormat.setOutputPath(job, new Path(args[1]));
17
18        job.setMapperClass(SetsimjoinsMapper.class);
19        job.setReducerClass(SetsimjoinsReducer.class);
20
21      job.setMapOutputKeyClass(DocPair.class);
22      job.setMapOutputValueClass(Text.class);
23
24        job.setOutputKeyClass(Text.class);
25        job.setOutputValueClass(Text.class);
26
27        job.setOutputFormatClass(TextOutputFormat.class);
28
29
30        FileSystem fs = FileSystem.get(new Configuration());
31        if (fs.exists(new Path(args[1]))) {
32          fs.delete(new Path(args[1]));
33        }
34
35        job.waitForCompletion(true);
36
37    }
38 }
```

- Instantiation of the counter of the number of comparisons.

```
1    public static enum COUNTER {
2        NB_COMPARISIONS_I,
3    };
```

- Writing the value of the counter in HDFS.

```
1        long counter = job.getCounters()
2                .findCounter(COUNTER.NB_COMPARISIONS_I).getValue();
3        Path outFile = new Path("NB_COMPARISIONS_I.txt");
4        BufferedWriter br = new BufferedWriter(new OutputStreamWriter(
5                fs.create(outFile, true)));
6        br.write(String.valueOf(counter));
7        br.close();
8
9      System.exit(0);
```

| | Job Overview |
|---|---|
| **Job Name:** | Part2_Set Similarity Joins |
| **User Name:** | cloudera |
| **Queue:** | root.cloudera |
| **State:** | SUCCEEDED |
| **Uberized:** | false |
| **Submitted:** | Fri Mar 17 08:53:09 PDT 2017 |
| **Started:** | Fri Mar 17 08:53:40 PDT 2017 |
| **Finished:** | Fri Mar 17 08:54:12 PDT 2017 |
| **Elapsed:** | 31sec |
| **Diagnostics:** | |
| **Average Map Time** | 9sec |
| **Average Shuffle Time** | 5sec |
| **Average Merge Time** | 0sec |
| **Average Reduce Time** | 2sec |

**ApplicationMaster**

| Attempt Number | Start Time | Node | Logs |
|---|---|---|---|
| 1 | Fri Mar 17 08:53:14 PDT 2017 | quickstart.cloudera:8042 | logs |

| Task Type | Total | Complete | |
|---|---|---|---|
| **Map** | 1 | 1 | |
| **Reduce** | 1 | 1 | |

| Attempt Type | Failed | Killed | Successful |
|---|---|---|---|
| **Maps** | 0 | 0 | 1 |
| **Reduces** | 0 | 0 | 1 |

Figure 2: Logs in Yarn of the comparison job for the 1st approach

## 2.2 Second approach

In the second approach we compare only pairs of documents that have in common the first $|d| - \lceil t.|d| \rceil + 1$ words, where $|d|$ is the number of words in document d, and t is the Jaccard similarity threshold.

### 2.2.1 Mapper

The mapper implements an inverted index for the first $|d| - \lceil t.|d| \rceil + 1$ words of each document.

```java
public class Setsimjoins2Mapper extends
    Mapper<LongWritable, Text, Text, Text> {

    private Text word = new Text();

    @Override
    public void map(LongWritable key, Text value, Context context)
            throws IOException, InterruptedException {

        String doc = value.toString().split(";")[1];
        String docID = value.toString().split(";")[0];
        String[] words = doc.split(" ");
        long keptWordsNumber = Math.round(words.length - (words.length *
    0.8) + 1);
        String[] keptWords = Arrays.copyOfRange(words, 0,(int)
    keptWordsNumber);


        for (String keptWord : keptWords) {

            word.set(keptWord);
            //System.out.println(docID);
            context.write(word, new Text(docID));
        }
    }
```

9

```
24  }
```

### 2.2.2 Reducer

The reducer create pairs of documents the belong to the same line in the inverted index, ad then computes the Jaccard similarity.

```java
1  public class Setsimjoins2Reducer extends Reducer<Text, Text, Text, Text> {
2
3      private BufferedReader reader;
4
5      @Override
6    public void reduce (Text key, Iterable<Text> values, Context context)
7                throws IOException, InterruptedException {
8
9        HashMap<String, String> allLines = new HashMap<String, String >();
10           reader = new BufferedReader(new FileReader(new File("/home/cloudera/
    workspace/SetSimJoins/preprocessing_output_sample.txt")));
11           String line;
12           while ((line = reader.readLine()) != null) {
13               String[] lineS = line.split(";");
14               allLines.put(lineS[0], lineS[1]);
15           }
16
17           List<String> docSet = new ArrayList<String >();
18           for (Text id : values){
19             docSet.add(id.toString());
20           }
21
22        //System.out.println(docSet[0]);
23
24        if (docSet.size() > 1) {
25             ArrayList<String> pairs = new ArrayList<String >();
26             for (int i = 0; i < docSet.size(); ++i) {
27                 for (int j = i + 1; j < docSet.size(); ++j) {
28                     String pair = new String(docSet.get(i) + " "
29                             + docSet.get(j));
30                     pairs.add(pair);
31                 }
32             }
33             //System.out.println("*****************");
34             //System.out.println(pairs.size());
35             for (String pair : pairs) {
36               HashSet<String> words11 = new HashSet<String >();
37                 String words12 = allLines.get(pair.split(" ")[0].toString());
38                 for (String word : words12.split(" ")) {
39                     words11.add(word);
40                 }
41
42                 HashSet<String> words21 = new HashSet<String >();
43                 String words22 = allLines.get(pair.split(" ")[1].toString());
44                 for (String word : words22.split(" ")) {
45                     words21.add(word);
```

```
46                }
47
48                context.getCounter(COUNTER.NB_COMPARISIONS_II).increment(1);
49                double sim = jaccardsim(words11,
50                    words21);
51                System.out.println("*************");
52                if (sim >= 0.1) {
53                  System.out.println(pair.split(" ")[0]);
54                    context.write(new Text("(" + pair.split(" ")[0] + ", "
55                            + pair.split(" ")[1] + ")"),
56                            new Text(String.valueOf(sim)));
57                }
58            }
59        }
60    }
61 }
```

### 2.2.3 Driver

The driver of this second approach has the same structure as the first approach.

```
1 public class Setsimjoins2 {
2
3   public static enum COUNTER {
4       NB_COMPARISIONS_II,
5   };
6
7   public static void main(String[] args) throws IOException,
8     ClassNotFoundException, InterruptedException {
8      /*
9         * Validate that two arguments were passed from the command line.
10        */
11       if (args.length != 2) {
12         System.out.printf("Usage: StopWords <input dir> <output dir>\n");
13         System.exit(-1);
14       }
15
16       /*
17        * Instantiate a Job object for your job's configuration.
18        */
19       Configuration conf = new Configuration();
20       Job job = new Job(conf);
21       job.setJarByClass(Setsimjoins.class);
22       job.setJobName("Part2_Set Similarity Joins_2nd method");
23
24       FileInputFormat.addInputPath(job, new Path(args[0]));
25       FileOutputFormat.setOutputPath(job, new Path(args[1]));
26
27       job.setMapperClass(Setsimjoins2Mapper.class);
28       job.setReducerClass(Setsimjoins2Reducer.class);
29
30
31       job.setMapOutputKeyClass(Text.class);
32     job.setMapOutputValueClass(Text.class);
```

11

```java
33
34        job.setOutputKeyClass(Text.class);
35        job.setOutputValueClass(Text.class);
36
37        job.setOutputFormatClass(TextOutputFormat.class);
38
39
40        FileSystem fs = FileSystem.get(new Configuration());
41        if (fs.exists(new Path(args[1]))) {
42          fs.delete(new Path(args[1]));
43        }
44
45        job.waitForCompletion(true);
46
47        long counter = job.getCounters()
48                .findCounter(COUNTER.NB_COMPARISIONS_II).getValue();
49        Path outFile = new Path("NB_COMPARISIONS_II.txt");
50        BufferedWriter br = new BufferedWriter(new OutputStreamWriter(
51                fs.create(outFile, true)));
52        br.write(String.valueOf(counter));
53        br.close();
54
55      System.exit(0);
56    }
57 }
```
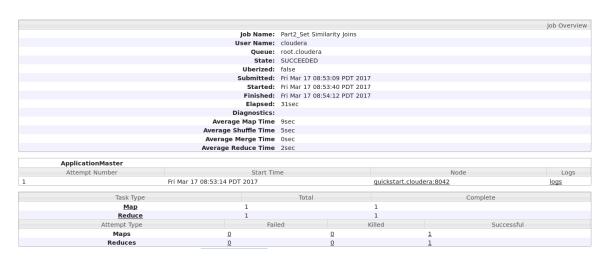
| | Job Overview |
|---|---|
| **Job Name:** | Part2_Set Similarity Joins |
| **User Name:** | cloudera |
| **Queue:** | root.cloudera |
| **State:** | SUCCEEDED |
| **Uberized:** | false |
| **Submitted:** | Fri Mar 17 08:53:09 PDT 2017 |
| **Started:** | Fri Mar 17 08:53:40 PDT 2017 |
| **Finished:** | Fri Mar 17 08:54:12 PDT 2017 |
| **Elapsed:** | 31sec |
| **Diagnostics:** | |
| **Average Map Time** | 9sec |
| **Average Shuffle Time** | 5sec |
| **Average Merge Time** | 0sec |
| **Average Reduce Time** | 2sec |

| **ApplicationMaster** | | | |
|---|---|---|---|
| Attempt Number | Start Time | Node | Logs |
| 1 | Fri Mar 17 08:53:14 PDT 2017 | quickstart.cloudera:8042 | logs |

| Task Type | Total | Complete |
|---|---|---|
| **Map** | 1 | 1 |
| **Reduce** | 1 | 1 |

| Attempt Type | Failed | Killed | Successful |
|---|---|---|---|
| **Maps** | 0 | 0 | 1 |
| **Reduces** | 0 | 0 | 1 |

Figure 3: Logs in Yarn of the comparison job for the 2nd approach

## 2.3   Comparison of the two approaches

| | | NB of comparisons | Execution time |
|---|---|---|---|
| Result on sample text | 1st method | 258 | 31 sec |
| | 2nd method | 212 | 36 sec |

As expected the second approach performs less comparisons than the 1st approach. However, the execution time of the 2nd method is larger than the 1st method. I think it is due to the

fact that the comparison is carried on a short sample file and therefore it won't be a large difference in elapsed time in addition to that the elapsed time of a same job can fluctuate by few seconds, therefore we obtain a larger elapsed time.

# 3 Appendix

## 3.1 sortByValue code

```
 1    private static Map<String, Integer> sortByValue(Map<String, Integer>
      wordcount) {
 2
 3      List<Map.Entry<String, Integer>> list =
 4              new LinkedList<Map.Entry<String, Integer>>(wordcount.entrySet
      ());
 5
 6        Collections.sort(list, new Comparator<Map.Entry<String, Integer>>() {
 7            public int compare(Map.Entry<String, Integer> o1,
 8                               Map.Entry<String, Integer> o2) {
 9                return (o1.getValue()).compareTo(o2.getValue());
10            }
11        });
12
13        Map<String, Integer> sortedMap = new LinkedHashMap<String, Integer>();
14        for (Map.Entry<String, Integer> entry : list) {
15            sortedMap.put(entry.getKey(), entry.getValue());
16        }
17        return sortedMap;
18    }
```

## 3.2 Implementation of WritableComparable Class

```
 1
 2  public class DocPair implements WritableComparable<DocPair> {
 3
 4    private Text key1;
 5      private Text key2;
 6
 7      public DocPair(Text key1, Text key2) {
 8          set(key1, key2);
 9      }
10      public DocPair() {
11          set(new Text(), new Text());
12      }
13
14      public DocPair(String key1, String key2) {
15          set(new Text(key1), new Text(key2));
16      }
17
18      public Text getFirst() {
19          return key1;
20      }
21
22      public Text getSecond() {
23          return key2;
24      }
25
26      public void set(Text key1, Text key2) {
```

```java
27          this.key1 = key1;
28          this.key2 = key2;
29      }
30
31      @Override
32      public void readFields(DataInput in) throws IOException {
33        key1.readFields(in);
34        key2.readFields(in);
35      }
36
37      @Override
38      public void write(DataOutput out) throws IOException {
39        key1.write(out);
40        key2.write(out);
41      }
42
43      @Override
44      public String toString() {
45          return key1 + " " + key2;
46      }
47
48      @Override
49      public int compareTo(DocPair other) {
50          int cmpFirstFirst = key1.compareTo(other.key1);
51          int cmpSecondSecond = key2.compareTo(other.key2);
52          int cmpFirstSecond = key2.compareTo(other.key2);
53          int cmpSecondFirst = key2.compareTo(other.key1);
54
55          if (cmpFirstFirst == 0 && cmpSecondSecond == 0 || cmpFirstSecond == 0
56                  && cmpSecondFirst == 0) {
57              return 0;
58          }
59
60          Text thisSmaller;
61          Text otherSmaller;
62
63          Text thisBigger;
64          Text otherBigger;
65
66          if (this.key1.compareTo(this.key2) < 0) {
67              thisSmaller = this.key1;
68              thisBigger = this.key2;
69          } else {
70              thisSmaller = this.key2;
71              thisBigger = this.key1;
72          }
73
74          if (other.key1.compareTo(other.key2) < 0) {
75              otherSmaller = other.key1;
76              otherBigger = other.key2;
77          } else {
78              otherSmaller = other.key2;
79              otherBigger = other.key1;
80          }
```

```
81
82          int cmpThisSmallerOtherSmaller = thisSmaller.compareTo(otherSmaller);
83          int cmpThisBiggerOtherBigger = thisBigger.compareTo(otherBigger);
84
85          if (cmpThisSmallerOtherSmaller == 0) {
86              return cmpThisBiggerOtherBigger;
87          } else {
88              return cmpThisSmallerOtherSmaller;
89          }
90      }
91      @Override
92      public int hashCode() {
93          return key1.hashCode() * 163 + key2.hashCode();
94      }
95
96      @Override
97      public boolean equals(Object o) {
98          if (o instanceof DocPair) {
99            DocPair tp = (DocPair) o;
100              return key1.equals(tp.key1) && key2.equals(tp.key2);
101          }
102          return false;
103      }
104 }
```

## 3.3 Jaccard similarity code

```
1      public double jaccardsim(HashSet<String> v1, HashSet<String> v2) {
2
3        HashSet<String> intersect1 = v1;
4        intersect1.retainAll(v2);
5        int intertsect = intersect1.size();
6
7          if (v1.size() < v2.size()) {
8            HashSet<String> unionSet = v1;
9            unionSet.addAll(v2);
10            int union = unionSet.size();
11              return (double) intertsect / union;
12          } else {
13            HashSet<String> unionSet = v2;
14              unionSet.addAll(v1);
15              int union = unionSet.size();
16              return (double) intertsect / union;
17          }
18 }
```