

# Intelligent Data Management - Exercise 1

Christoph Prinz

October 19, 2017

## Assignment 1

### Exercise a

Term  $i$  appears in  $n_i$  of the  $N$  documents

Appears in 40 documents:  $IDF_i = \log_2\left(\frac{N}{n_i}\right) = \log_2\left(\frac{10000000}{40}\right) = 18$

Appears in 10000 documents:  $IDF_i = \log_2\left(\frac{N}{n_i}\right) = \log_2\left(\frac{10000000}{10000}\right) = 10$

### Exercise b

Given the occurrence of a term  $i$  in document  $j$  is  $f_{ji}$  and  $\max_k f_{kj}$  the maximum number of occurrences of any term in this document, the term frequency  $TF$  is defined as:

$$TF_{ij} = \frac{f_{ij}}{\max_k f_{kj}}$$

Word  $w$  appears in 320 documents. In document  $d$  the maximum occurrence of any word is 15.

a)  $w$  appears once:

$$TF_{wd} = \frac{1}{15}$$

$$IDF_w = \log_2 \frac{10000000}{320} = 15$$

$$TF.IDF = \frac{1}{15} * 15 = 1$$

b)  $w$  appears five times:

$$TF_{wd} = \frac{5}{15} = \frac{1}{3}$$

$$IDF_w = *log_2 \frac{10000000}{320} = 15$$

$$TF.IDF = \frac{1}{3} * 15 = 5$$

### Exercise c

The basic rule is, that the possible hash-keys should not have any common factor with B (15 in this example). The population should therefore not be generated with a  $c$  of 3 or 5. To achieve an equal distribution of the generated hashes,  $c$  should be set to 1.

## Assignment 2

### Exercise a

The implementation consists out of a class `Document`, that counts words in a document and calculates the term-frequency of a given term in this document. The class `DocumentCollection` reads multiple input files and calculates IDF and TF.IDF for a given document and a given term. The source code is printed on the following pages.

### Exercise b

The choice of parameters is explained with comments in the source code.

```

1 import java.io.IOException;
2 import java.nio.charset.Charset;
3 import java.nio.file.Files;
4 import java.nio.file.Paths;
5 import java.util.Collections;
6 import java.util.HashMap;
7 import java.util.Map;
8
9 public class Document {
10     private Map<String, Integer> wordCount = new HashMap<String, Integer>();
11
12     public Document(String path) throws IOException {
13         countWordsFromFile(path, Charset.forName("Cp1252"));
14     }
15
16     public boolean doesTermAccur(String term){
17         return wordCount.containsKey(term);
18     }
19
20     /**
21      * calculate term frequency for a given term in document
22      * @param term: term w to calculate tf
23      * @return value for tf
24      */
25     public double get_tf(String term){
26         int f_ji = 0;
27         if (wordCount.containsKey(term)) {
28             f_ji = wordCount.get(term);
29         }
30         int max_k_fkj = Collections.max(wordCount.values());
31
32         return calculate_TF((double)f_ji, (double) max_k_fkj);
33     }
34
35     /**
36      * calculates tf
37      * @param f_ji: number of occurrences of a term
38      * @param max_k_fkj: maximum number of occurrences of any term
39      * @return value for tf
40      */
41     private static double calculate_TF(double f_ji, double max_k_fkj){
42         return f_ji/max_k_fkj;
43     }
44
45     /**
46      * Method to count words from a text file
47      * @param path: path of file to read
48      * @param encoding: encoding of file
49      * @throws IOException if file can't be loaded
50      */
51     private void countWordsFromFile(String path, Charset encoding)
52         throws IOException {
53         String text = new String(Files.readAllBytes(Paths.get(path)), encoding);
54         //Trim file in order to get only raw words
55         text = text.replaceAll("[\\s]", " ");
56         text = text.replaceAll("[^A-Za-zöüÄÖÜß ]", "");
57         text = text.toLowerCase();
58         String[] parts = text.split(" ");
59
60         //count words
61         for (String part : parts) {
62             if (wordCount.containsKey(part)) {
63                 wordCount.put(part, wordCount.get(part) + 1);
64             } else {
65                 wordCount.put(part, 1);
66             }
67         }
68     }
69 }

```

```

1 import java.util.HashMap;
2
3 public class DocumentCollection {
4     private HashMap<String, Document> documents = new HashMap<>();
5
6     public DocumentCollection(String[] paths) {
7         try {
8             for (String path : paths) {
9                 documents.put(path, new Document(path));
10            }
11        } catch (Exception e) {
12            System.out.println("Error Reading file!: " + e.getMessage());
13        }
14    }
15
16    /**
17     * calculates tf_idf
18     * @param term: term w to calculate tf_idf
19     * @param document: document d in collection to calculate tf_idf
20     * @return value for tf_idf
21     */
22    public double get_tf_idf(String term, String document) {
23        double tf = documents.get(document).get_tf(term);
24        try {
25            double idf = get_idf(term);
26            return calculate_tf_idf(tf, idf);
27        } catch (Exception e) {
28            System.out.println("Term does not occur: " + e.getMessage());
29            return 0;
30        }
31    }
32
33    public static double calculate_tf_idf(double tf, double idf) {
34        return tf * idf;
35    }
36
37    /**
38     * returns idf
39     *
40     * @param term: given term i to use to calculate idf
41     * @return value for idf
42     * @throws Exception
43     */
44    public double get_idf(String term) throws Exception {
45        int n_i = 0;
46        for (Document document : documents.values()) {
47            if (document.doesTermAccur(term)) {
48                n_i += 1;
49            }
50        }
51        return calculate_idf(documents.size(), n_i);
52    }
53
54    /**
55     * Method to calculate inverse document frequency
56     * @param N: Total number of documents
57     * @param n_i: Number of occurrences of a term
58     * @return value for idf
59     */
60    public static double calculate_idf(double N, double n_i) {
61        return log2(N / n_i);
62    }
63
64    private static double log2(double x) {
65        return (Math.log(x) / Math.log(2));
66    }
67 }

```