

WordCount.java

```
1  import java.io.File;
2  import java.io.FileNotFoundException;
3  import java.util.ArrayList;
4  import java.util.HashMap;
5  import java.util.Scanner;
6  public class WordCount {
7      //Data structure to map a word to its frequency i.e. a Map
8      HashMap<String,Integer> wordMap = new
HashMap<String,Integer>();// KEY VALUE ie word & its count
9      HashMap<String,Integer> pairMap = new
HashMap<String,Integer>();//KEY VALUE ie word & its count
10
11     public void readFile(String file) {
12         String word = "";
13         //var for pairMap
14         String prevWord = "";
15         String nextWord = "";
16         String tempKey = "";
17         Scanner sc;
18         try {
19             sc = new Scanner(new
File("C:\\Users\\bhupe\\Downloads\\"+file));
20             //Split the text into words using whitespace
character
21             sc.useDelimiter(" ");//("\\s +") for multiple spaces
ie one of more occurance of space
22             //Get a word at a time
23             while(sc.hasNext()) {
24                 word = sc.next(); //Get a word
25                 word = word.toLowerCase(); //Convert words to
lowercase
26                 word = word.replaceAll("[,;.]", ""); //Delete
punctuation
27                 //Store words as keys and frequencies as values
28                 if(!word.equals("")) {
29                     if(wordMap.containsKey(word))
30                         wordMap.put(word, wordMap.get(word) +
1);
```

```

31         else
32             wordMap.put(word, 1);
33     }
34 }
35
36 //for keyMap
37 sc = new Scanner(new
File("C:\\Users\\bhupe\\Downloads\\"+file));
38 //Split the text into words using whitespace
character
39 sc.useDelimiter(" ");//( "\\s +" ) for multiple spaces
ie one of more occurance of space
40 //Get a word at a time
41 while(sc.hasNext()) {
42     nextWord = sc.next(); //Get a word
43     nextWord = nextWord.toLowerCase(); //Convert
words to lowercase
44     nextWord = nextWord.replaceAll("[,;]", "");
//Delete punctuation
45     //Store words as keys and frequencies as values
46     if(!nextWord.equals("") && !prevWord.equals(""))
{
47         tempKey = prevWord + ":" + nextWord;
48         if(pairMap.containsKey(tempKey))
49             pairMap.put(tempKey,
pairMap.get(tempKey) + 1);
50         else
51             pairMap.put(tempKey, 1);
52     }
53     prevWord = nextWord;
54 }
55
56
57 sc.close();
58 } catch (FileNotFoundException e) {
59     System.out.println("File does not exist!" + e);
60 }
61 }
62 public void printPairFreq() {
63     System.out.print("check");
64     for(String pairWord: pairMap.keySet()) {
65         System.out.print("check");

```

```

66
67         int count = pairMap.get(pairWord);
68         System.out.print("{ "+pairWord+", "+count+"}"+ "\n");
69     }
70 }
71
72 public void printWordFreq() {
73
74     for(String word: wordMap.keySet()) {
75         int count = wordMap.get(word);
76         System.out.print("{ "+word+", "+count+"}"+ "\n");
77     }
78 }
79
80 double probWord(String word){//tf
81     double p;//probability
82     int total_count=0;
83     int word_count=0;
84     for(String wordInMap: wordMap.keySet()) {
85         total_count=total_count+1;
86     }
87     word_count = wordMap.get(word);
88     p = word_count*1.0/total_count*1.0; // type casting
89     return p;
90 }
91
92
93 double probOfWordPair(String first,String second){
94     String word = first + ":" +second;
95     double p;//probability
96     int total_count=0;
97     int wordPair_count=0;
98     for(String wordInMap: pairMap.keySet()) {
99         total_count=total_count+1;
100     }
101     wordPair_count = pairMap.get(word);
102     p = wordPair_count*1.0/total_count*1.0; // type casting
103     return p;
104 }
105

```

```

106
107     double idf(String[] arr,String term){
108         int length = arr.length;
109         double idft=0;
110         double p;
111         int dft=0;
112         for (int i = 0; i < length; i++) {
113             readFile(arr[i]+".txt");//file name
114             p=probWord(term);
115             if(p>0){
116                 dft++;
117             }
118         }
119
120         idft=Math.log(length/dft); //length is total number of
files & dft is the total number of files that have the
particular tern that is passed in the function

```

```

121
122         return idft;
123     }

```

```

124
125
126     double tf_idf(String[] arr,String term){
127         double overAllProb=0;
128         double p;
129         for(int i = 0; i < arr.length; i++) {
130             readFile(arr[i]+".txt");//file name
131             p=probWord(term);
132             if(p>0){
133                 overAllProb+=p;
134             }
135         }
136
137         double tf = overAllProb/arr.length;
138         double idf = idf(arr, term);//this logic is wrong
139
140         return tf*idf;
141     }

```

```
144 public void main(String[] args) {
145     WordCount unigram = new WordCount();
146     unigram.readFile("file.txt");
147     //unigram.printWordFreq();
148     //unigram.printPairFreq();
149     double pWord = unigram.probWord("that");
150     System.out.print(pWord+"\n");
151     double pWordPair = unigram.probOfWordPair("is", "the");
152     System.out.print(pWordPair+"\n");
153
154     String[] numbers = {"file1", "file2", "file3"};
155     double idf = idf(numbers, "the");
156     System.out.println(idf);
157
158     double id_idf = tf_idf(numbers, "the");
159     System.out.println(id_idf);
160
161 }
162 }
163
```