# CMP-5014Y Coursework 2 - Word Auto Completion with Tries

Student number: 100242165. Blackboard ID: dvd18scu

Wednesday 13th May, 2020 13:44

# Contents

# 1  Part 1: Form a Dictionary and Word Frequency Count

Part 1 deals with document file reading and dictionary writing, and creating dictionaries for the auto complete. Dictionaries are CSV files where each line has a word, as well as a number, the frequency that word appears in a document.

## 1.1  formDictionary Pseudocode

formDictionary is a function that reads a document file, a CSV file with words, and creates a dictionary and writes it to file. A Tree Map is used as it is efficient in both accessing from and adding items to it. The word itself is used as the key since it is unique, and the frequency the data. The function uses readWordsFromCSV, which returns an array of words.

---
**Algorithm 1** formDictionary($readPath, writePath$)

---
**Require:** A valid file to read located at $readPath$, and a file to write to located at $writePath$.
**Ensure:** The file located at $writePath$ contains a list of words from the file located at $readPath$, along with the frequency of each word's appearance.
1: $\mathbf{M} \leftarrow Map()$         ▷ $\boldsymbol{M}$ *is a map, where a string is the key, and the frequency is the data*
2: $\mathbf{L} \leftarrow readWordsFromCSV(readPath)$         ▷ $\boldsymbol{L}$ *is a list of words read from readPath*
3: **for** $i \leftarrow 1$ **to** $\mathbf{L}$.size() **do**
4:     $\mathbf{M}$.put($\mathbf{L}_i$,frequency($\mathbf{L}_i$)) ▷ $frequency()$ *linearly counts the frequency.* ▷ *Adds the words in* $\boldsymbol{L}$, *along with the word count for each word to* $\boldsymbol{M}$.
5: $saveToFile(\mathbf{M}, writePath)$         ▷ *saveToFile writes* $\boldsymbol{M}$ *to the file at writePath.*

---

## 1.2  formDictionary Analysis

1. The fundamental operation for formDictionary is:
   $\mathbf{M}$.put($\mathbf{L}_i$,frequency($\mathbf{L}_i$))

2. Treemap is $\mathcal{O}(log(n))$ when adding item to it. All cases, worst, best and average are the same. Where n is the size of $\mathbf{L}$

3.
$$f(g(n)) = \sum_{i=1}^{n} log(g(n)) \tag{1}$$

4. g(n) is the runtime complexity of $frequency()$, a linear scanning function, with a complexity of $\mathcal{O}(n)$. This makes the final runtime complexity

5.
$$f(n) = \sum_{i=1}^{n} log(n^2) \tag{2}$$

$$t(n) = nlog(n^2) \tag{3}$$

6. This simplifies further to 2nlog(n), log linear. Ignoring constants, the runtime complexity function is $\mathcal{O}(nlog(n))$

## 1.3  saveToFile Pseudocode

saveToFile is a simple function that writes a map with a string for its key and an integer for it's value to file. It saves the data as a CSV, with each line being a word and it's frequency.

---
**Algorithm 2** saveToFile($\mathbf{M}, writePath$)

---
**Require:** A map $\mathbf{M}$ with a string as the key, and an integer as the data, and a file located at $writePath$
**Ensure:** A file at $writePath$ that contains all keys (words) and data (frequency) to be written to file.
1: $openFile \leftarrow open(writePath)$         ▷ *Opens the file*
2: **for all i in M do**
3:     $tempString \leftarrow getKey(\mathbf{i}) + "," + getValue(\mathbf{i})$         ▷ *Joining two strings, with a comma in-between.*
4:     $writeline(openFile, tempString)$
5: $close(writePath)$         ▷ *Closes the file*

---

## 1.4    saveToFile Analysis

1. Fundamental operation:

   $tempString \leftarrow getKey(\mathbf{e}) + ", " + getValue(\mathbf{e})$

2. Getting items from a treemap is $\mathcal{O}(log(n))$. All cases, worst, best and average are the same. The number of times the loop loops, n is how many items/entries are in the treemap **M**.

3.

$$f(n) = \sum_{i=1}^{n} log(n) + \sum_{i=1}^{n} log(n) \tag{4}$$

$$t(n) = 2nlog(n) \tag{5}$$

4. Ignoring constants, the order of the runtime complexity is log linear, and the runtime complexity function is $\mathcal{O}(nlog(n))$.

# 2 Part 2: Implement a Trie Data Structure

This section is about implementing a Trie data structure. It's implementation is where each "node" is an array of size 26, where the index determines what letter it will be. For example, a is index 0, b is 1, and so on. An Ascii offset is mentioned a number of times in this section. This is for converting Ascii's representation of the letters into the node's representation for them. For example, the Ascii char "a" has the integer value of 97 instead of 0.

## 2.1 add Pseudocode

This function adds a string to the Trie structure. It returns true only if the word that was added was newly added. It works by traversing the trie, adding nodes if necessary, and returning the correct boolean value based on flags.

---

**Algorithm 3** add(**key**) return **true or false**

---

**Require:** $key$, a string.
**Ensure:** a boolean value, **true** or **false**
1:   $alreadyIn \leftarrow$ **true**
2:   $parent \leftarrow root$              $\triangleright$ *root is the the root node of the trie*
3:   **for** $level \leftarrow 1$ **to** $key.size()$ **do**       $\triangleright$ *goes through all letters of the key.*
4:      $index \leftarrow charToInt(level)$
5:      **if** $parent.getOffspring()_{index} =$ **null then**     $\triangleright$ *if the node for the current letter doesn't exist.*
6:         $parent.getOffspring()_{index} = TrieNode()$      $\triangleright$ *make a new one*
7:         $alreadyIn \leftarrow$ **false**
8:      $parent \leftarrow parent.getOffspring()_{index}$      $\triangleright$ *advanced parent node to one further down*
9:   **if** $alreadyIn$ && $parent.getIsEnd()$ **then**      $\triangleright$ *if everything was already done, return false*
10:      **return false**
11: $parent.setIsEnd($**true**$)$          $\triangleright$ *sets last node for key as end.*
12: **return true**

---

## 2.2 contains Pseudocode

This function checks if a word is present in the trie. It returns true only if the word is both in the trie and is marked as a key. It also works by traversing the trie.

---

**Algorithm 4** contains(**key**) return **true or false**

---

**Require:** $key$, a string.
**Ensure:** a boolean value, **true** or **false**
1:   $alreadyIn \leftarrow$ **true**
2:   $parent \leftarrow root$             $\triangleright$ *root is the the root node of the trie*
3:   **for** $level \leftarrow 1$ **to** $key.size()$ **do**      $\triangleright$ *goes through all letters of the key.*
4:      $index \leftarrow charToInt(level)$
5:      **if** $parent.getOffspring()_{index} =$ **null then**     $\triangleright$ *if the node for the current letter doesn't exist.*
6:         **return false**
7:      $parent \leftarrow parent.getOffspring()_{index}$      $\triangleright$ *advanced parent node to one further down*
8:   **return** $(parent \neq$ **null** && $parent.getIsEnd())$      $\triangleright$ *returns true only if parent exists AND is a key*

---

## 2.3 outputBreadthFirstSearch Pseudocode

This function explores and prints the contents of the trie in breadth first order. It uses a queue. It goes through the queue, originally starting with the root, adding it's offspring to the queue. It repeats this until there are no more nodes.

---
**Algorithm 5** outputBreadthFirstSearch() **return** ($result$)

---
**Require:** Access to a valid Trie root node.
**Ensure:** A string $result$, which contains all characters in a trie in breadth first order.
  1: $\mathbf{Q} \leftarrow Queue()$                                                ▷ $\mathbf{Q}$ *is a queue of trie nodes*
  2: $\mathbf{Q}.add(\mathbf{root})$                                      ▷ *Add the root node to the front of the queue*
  3: **while** $\mathbf{Q} \neq$ **empty do**
  4:     $\mathbf{tempNode} \leftarrow \mathbf{Q}.remove()$                         ▷ *remove() removes the top trie node*
  5:     $\mathbf{O} \leftarrow \mathbf{tempNode}.getOffspring()$         ▷ *getOffspring() returns an array of the node's children*
  6:     **for** $i \leftarrow 1$ **to** $\mathbf{O}.size()$ **do**
  7:         **if** $\mathbf{O}_i \neq$ **null then**
  8:             $\mathbf{Q}.add(\mathbf{O}_i)$
  9:             $result \leftarrow result + intToChar(i + 97)$     ▷ *append character to result. 97 is the offset for Ascii letters*
10: **return** ($result$)

---

## 2.4 outputDepthFirstSearch Pseudocode

This is the helper function to outputDepthFirstSearch. It provides the actual logic with a string to start appending to, as it is recursive.

---
**Algorithm 6** outputDepthFirstSearch() **return** ($result$)

---
**Require:** Access to a valid Trie root node.
**Ensure:** A string $result$, which contains all characters in a trie in depth first order.
  1: **return** $outputDepthFirstSearch(result, \mathbf{root})$

---

This is the actual logic of outputDepthFirstSearch. It linearly goes through every child of a node, and calls itself on the child as if it were the parent. It does this until there are no more nodes.

---
**Algorithm 7** outputDepthFirstSearch($result$, **root**) **return** ($result$)

---
**Require:** $result$, a string to build on, and **root**, a valid root node to a trie
**Ensure:** A string $result$, which contains all characters in a trie in depth first order.
  1: $\mathbf{O} \leftarrow \mathbf{root}.getOffspring()$             ▷ *getOffspring() returns an array of the node's children*
  2: **for** $i \leftarrow 1$ **to** $\mathbf{O}.size()$ **do**                           ▷ *For every child*
  3:     **if** $\mathbf{O}_i \neq$ **null then**
  4:         $result \leftarrow result + intToChar(i + 97)$     ▷ *append character to result. 97 is the offset for Ascii letters*
  5:         $outputDepthFirstSearch(result, \mathbf{O}_i)$
  6: **return** $result$

---

## 2.5  getSubTrie Pseudocode

This function gets all the words from a trie starting with a given prefix, and returns a new trie from it. It simply explores the trie to make sure the prefix is in the trie, then creates a trie with the node for the last letter in the prefix as the root.

---

**Algorithm 8** getSubTrie(**root**, $prefix$) return (**subTrie**)

---

**Require: root**, a valid node of a trie, and $prefix$, a prefix for a word/words.
**Ensure: subTrie** is a trie that starts from $prefix$.
 1: **tempNode** ← **root**
 2: **for** $i \leftarrow 1$ **to** $prefix.size()$ **do**
 3:     $index \leftarrow charToInt(prefix_i) - 97$          ▷ *charToInt(character) converts an Ascii character to its int version*
 4:     **if tempNode = null then**
 5:         **return null**
 6:     **tempNode** ← **tempNode**.$getOffSpring()_{index}$
 7: **return** $Trie($**tempNode**$)$

---

## 2.6  getAllWords Pseudocode

This is a helper function for getAllWords, providing the main recursive function with an ArrayList of strings as well as an empty StringBuilder.

---

**Algorithm 9** getAllWords() return **W**

---

**Require:** Access to a valid trie's **root**.
**Ensure: W**, an array of strings of varying length.
 1: $getAllWords($**W**$, builder,$ **root**$)$          ▷ *builder* is an empty string, and *root* is the root of the Trie to be used.
 2: $return$ **W**

---

This is the main function, where it is to build an array of all the words present in the Trie it is called recursively. It works quite similarly to outputDepthFirstSearch, but whilst building an array of the trie's words.

---

**Algorithm 10** getAllWords(**W**, $builder$, **root**)

---

**Require: W**, an array of strings of varying length, **builder**, a string to append to, and **root**, the root of the given Trie structure.
**Ensure: W** has been filled with all the words from the Trie.
 1: **if root**.$getIsEnd()$ **then**
 2:     **W**.$add(builder)$
 3: **O** ← **root**.$getOffspring()$
 4: **for** $i \leftarrow 1$ **to** **O**.$size()$ **do**
 5:     **if O**$_i \neq$ **null then**
 6:         $builder \leftarrow builder + intToChar(97 + i)$ ▷ *97 is the Ascii offset. intToChar converts the int to it's char equivalent*
 7:         $getAllWords($**W**$, builder,$ **O**$_i)$
 8:         $builder \leftarrow builder.shorten(1)$          ▷ *Removes last character from builder - "goes back up a level"*

---

# 3 Part 3: Word Auto Completion Application

This part is the autocomplete program. It uses part 1 and a modified part 2. The Trie structure now also implements the frequency count of each word. This affects add and contains slightly as well. It also features a nested static class, fullInfo. Each fullInfo object simply stores a word and its frequency, and getAllInfo acts like getAllWords but returning fullInfo objects instead. This is used to implement sorting to generate the top 3 results for a word. It also features another saveToFile, which simply saves the generated string output of the function to file. Finally, it implements a function addToTrie, used to add a dictionary to a trie, using add().

## 3.1 getAllInfo Pseudocode

This is a helper function, providing the main recursive function with an ArrayList of strings as well as an empty StringBuilder. The fullInfo object used consists of an int and a string, being frequency and the word itself.

---
**Algorithm 11** getAllInfo() **return I**

---
**Require:** Access to a valid **root** to a Trie
**Ensure: I**, an array of fullInfo objects of varying length.
1: $getAllInfo(\mathbf{I}, builder, \mathbf{root})$          ▷ *builder* is an empty string, and *root* is the root of the Trie to be used.
2: $\mathbf{I}.sortByFrequency()$     ▷ $sortByFrequency()$ is a function that sorts by fullInfo's frequency, followed by key.
3: **return I**

---

This is the main function, where it is to build an array of all the words present, as well as each word's frequency. The fullInfo object used consists of an int and a string, being frequency and the word itself. The function works similarly to getAllWords, which in turn is similar to outputDepthFirstSearch.

---
**Algorithm 12** getAllInfo(**I**, **builder**, **root**)

---
**Require: I**, an array of fullInfo objects of varying length, **builder**, a string to append to, and **root**, the root of the given Trie structure.
**Ensure: I** has been filled with all the words and information from the Trie.
  **if root**.$getIsEnd()$ **then**
    $tempInfo \leftarrow fullInfo(builder, \mathbf{root}.getFrequency())$
    $\mathbf{I}.add(tempInfo)$
  **else**
    $\mathbf{O} \leftarrow \mathbf{root}.getOffspring()$
    **for** $i \leftarrow 1$ **to** $\mathbf{O}.size()$ **do**
      **if** $\mathbf{O}_i \neq$ **null then**
        $builder \leftarrow builder + intToChar(97 + i)$ ▷ *97 is the Ascii offset. intToChar converts ints to it's char equivalent*
        $getAllInfo(\mathbf{I}, builder, \mathbf{O}_i)$
        $builder \leftarrow builder.shorten(1)$        ▷ *Removes last character from builder - "goes back up a level"*

---

## 3.2 addToTrie Pseudocode

This is quite a simple function. It reads in a dictionary formatted file and calls the add() function on each word and it's frequency. (add() in part 3 takes in both values)

---
**Algorithm 13** addToTrie(*readPath*)

---
**Require:** Access to a correctly formatted dictionary file, located with the string *readPath*. Access to a trie to add to.
**Ensure:** The contents of the dictionary file to be added to the trie.
  $openFile \leftarrow open(readPath)$
  **while** $openFile \neq end$ **do**          ▷ *while not the end of file*
    $line \leftarrow openFile.readLine()$
    $splitLine \leftarrow line.split(",")$          ▷ *Splits the line by the comma character*
    $add(splitLine_1, splitLine_2)$          ▷ *Calls the add function*
  $openFile.close()$

---

## 3.3 saveToFile Pseudocode

This is another saveToFile function that simply writes a string to a file as a line in the file.

---
**Algorithm 14** saveToFile($line$,$writePath$)
---
**Require:** Access to a file to write to, located at $writePath$.
**Ensure:** $line$ has been written to the file at $writePath$.
 1: $openFile \leftarrow open(writePath)$
 2: $openFile.writeLine(line)$
 3: $openFile.close()$

---

## 3.4 autoCompletion Pseudocode

This function uses the prior functions to generate the top three results for a prefix, then generates a formatted line to have it be written to file by saveToFile(). It makes a subTrie for every prefix, then calculates and generates the ideal output for it.

---
**Algorithm 15** completion(**masterTrie**, **P**,**n**, **writePath**)
---
**Require:** masterTrie, the complete trie structure, and P, an array of prefixes (strings), of length n. A location for the file to write to, $writePath$.
**Ensure:** The printing of the top three most frequent words for each prefix in $prefixes$, along with their frequency and probability.

  $infoString \leftarrow$ ""           ▷ *This is the string that will be written to file*
  **for** $i \leftarrow 1$ **to** $n$ **do**           ▷ *For every prefix*
     **currentSubTrie** $\leftarrow$ **masterTrie**.$getSubTrie(\mathbf{P}_i)$
     **I** $\leftarrow$ **currentSubTrie**.$getAllInfo()$     ▷ **I** *is an array of info about each word beginning with the prefix*
     $totalFreq \leftarrow 0$
     **for** $j \leftarrow 1$ **to** **I**.$size()$ **do**
       $totalFreq \leftarrow totalFreq + \mathbf{I}_j.getFreq()$

     **if** **I**.$size() < 3$ **then**      ▷ *If there are less results than the top 3*
       $printLimit \leftarrow \mathbf{I}.size()$
     **else**
       $printLimit \leftarrow 3$
     **for** $j \leftarrow 1$ **to** $printLimit$ **do**
       $prob \leftarrow \mathbf{I}_j.getFreq()/totalFreq$     ▷ *Calculates the probability ratio*
       $print(\mathbf{P}_i, \mathbf{I}_j.getKey() + "," + prob)$
       $infoString \leftarrow infoString + \mathbf{P}_i + \mathbf{I}_j.getKey() + "," + prob + ","$
     $infoString \leftarrow infoString + "/n"$      ▷ *End the line*
  $saveToFile(infoString, writePath)$

---

# 4 Code Listing

## 4.1 DictionaryMaker

Listing 1: DictionaryMaker.java

```java
1  /*
2  By/Modified by Robin Rai (100242165)
3  V.1.0.0
4  Created on 05/03/2020
5  */
6  package dsacoursework2;
7
8  import java.io.*;
9  import java.util.*;
10
11 public class DictionaryMaker {
12
13         public DictionaryMaker() {
14         }
15
16         public static ArrayList<String> readWordsFromCSV(String file, String delim)
17                         throws FileNotFoundException {
18                 Scanner sc = new Scanner(new File(file));
19                 sc.useDelimiter(delim);
20                 ArrayList<String> words = new ArrayList<>();
21                 String str;
22                 while (sc.hasNext()) {
23                         str = sc.next();
24                         str = str.trim();
25                         str = str.toLowerCase();
26                         words.add(str);
27                 }
28                 return words;
29         }
30
31         public static void saveToFile(Map<String, Integer> map, String file)
32                         throws IOException {
33                 //goes through any map and writes it to file
34                 FileWriter fileWriter = new FileWriter(file);
35                 PrintWriter printWriter = new PrintWriter(fileWriter);
36                 for (Map.Entry<String, Integer> entry : map.entrySet()) {
37                         //System.out.println(entry.getKey() + "," + entry.getValue());
38                         printWriter.println(entry.getKey() + "," + entry.getValue());
39                         //writes key/word, followed by its value/frequency
40                 }
41                 printWriter.close();
42         }
43
44         public static void saveToFile(String line, String file) throws IOException {
45                 //simply writes the string to file
46                 FileWriter fileWriter = new FileWriter(file);
47                 PrintWriter printWriter = new PrintWriter(fileWriter);
48                 printWriter.println(line);
49                 printWriter.close();
50         }
51
52         //form a set of words that exist and count the frequency of each word
53         public void formDictionary(String fileDirectory, String saveDirectory) {
54                 TreeMap<String, Integer> map = new TreeMap<String, Integer>();
55                 //a treeMap is the most efficient for this. The string is the key,
```

```
56                     // and the int is the value, since
57                     //all string's are unique
58                     try {
59                             ArrayList<String> temp = readWordsFromCSV(fileDirectory, ",");
60                             //takes in a big ol array of words
61                             for (int i = 0; i < temp.size(); i++) {
62                                     map.put(temp.get(i), Collections.frequency(temp,
                                        ↪ temp.get(i)));
63                                     //for every word, put them in the treeMap as the key,
                                        ↪ with the
64                                     // frequency as the value
65                             }
66                     } catch (Exception e) {
67                             System.out.println("formDictionary: readWordsFromCSV");
68                     }
69                     //System.out.println(map);
70                     try {
71                             saveToFile(map, saveDirectory);
72                             //save it to file
73                     } catch (Exception e) {
74                             System.out.println("formDictionary: saveToFile");
75                     }
76             }
77
78
79             public static void main(String[] args) throws Exception {
80                     dsacoursework2.DictionaryMaker df = new
                            ↪ dsacoursework2.DictionaryMaker();
81                     ArrayList<String> in = readWordsFromCSV(
82                                     "src\\TextFiles\\testDocument.csv", ",");
83                     System.out.println("Array: " + in);
84
85                     df.formDictionary("src\\TextFiles\\testDocument.csv",
86                                     "src\\TextFiles\\formDictionaryTest.csv");
87                     saveToFile("Text to write",
88                                     "src\\TextFiles\\formDictionaryTest2.csv");
89                     System.out.println("Saved to file (both methods) successfully.");
90             }
91
92 }
```

## 4.2   Trie

Listing 2: Trie.java

```
 1 /*
 2 By Robin Rai (100242165)
 3 V.1.0.0
 4 Created on 05/03/2020
 5 */
 6 package dsacoursework2;
 7
 8 import java.util.ArrayList;
 9 import java.util.LinkedList;
10 import java.util.Queue;
11
12 public class Trie {
13         private TrieNode root;
14
15         public Trie() {
16                 root = new TrieNode();
```

```
17              }
18
19          public Trie(TrieNode input) {
20                  root = input;
21              }
22
23          boolean add(String key) {
24                  if (key.equals("")) {
25                          return false;
26                  }
27                  int index;
28                  boolean alreadyIn = true;
29                  TrieNode temp = root;
30                  for (int level = 0; level < key.length(); level++) {
31                          //goes through all letters of the key
32                          index = key.charAt(level) - 'a';
33                          if (temp.getOffspring()[index] == null) {
34                                  //if the node for the current letter doesn't exist
35                                  temp.getOffspring()[index] = new TrieNode();
36                                  //make a new one
37                                  alreadyIn = false;
38                                  //set the flag
39                          }
40                          temp = temp.getOffspring()[index];
41                          //advanced temp node to one further down
42                  }
43                  if (alreadyIn && temp.getIsEnd()) {
44                          //if everything was already done, return false
45                          return false;
46                  }
47                  temp.setIsEnd(true);
48                  //sets last node for key as end.
49                  return true;
50          }
51
52          boolean contains(String key) {
53                  int level;
54                  int length = key.length();
55                  int index;
56                  TrieNode temp = root;
57                  for (level = 0; level < length; level++) {
58                          //goes through all letters of the key
59                          index = key.charAt(level) - 'a';
60                          //index is current char's number
61                          if (temp.getOffspring()[index] == null) {
62                                  //if the node for the current letter doesn't exist
63                                  return false;
64                                  //return false, it isn't there
65                          }
66                          temp = temp.getOffspring()[index];
67                          //advance down a level
68                  }
69                  return (temp != null && temp.getIsEnd());
70                  //returns true only if it exists AND is a key - not just a prefix.
71          }
72
73
74          String outputBreadthFirstSearch() {
75                  Queue<TrieNode> queue = new LinkedList<>();
76                  StringBuilder builder = new StringBuilder();
```

```java
                    queue.add(root);
                    //starts with root
                    while (!queue.isEmpty()) {
                            TrieNode temp = queue.remove();
                            //takes first thing in queue
                            for (int i = 0; i < 26; i++) {
                                    if (!(temp.getOffspring()[i] == null)) {
                                            //for all of it's children
                                            queue.add(temp.getOffspring()[i]);
                                            //add them to the queue so it goes through
                                            //their children - width first
                                            builder.append((char) (i + 97));
                                            //adds it to the output
                                    }
                            }
                    }
                    return builder.toString();
        }

        public String outputDepthFirstSearch() {
                //helper function, provides string to append to
                StringBuilder builder = new StringBuilder();
                return outputDepthFirstSearch(builder, this.root);
        }

        String outputDepthFirstSearch(StringBuilder builder, TrieNode trieNode) {
                TrieNode[] children = trieNode.getOffspring();
                //makes an array of all children in the tree
                for (int i = 0; i < 26; i++) {
                        //for all of a node's children
                        if (children[i] != null) {
                                //if there's a child
                                builder.append((char) (i + 97));
                                //add it's character
                                outputDepthFirstSearch(builder, children[i]);
                                //call the function again on it's children
                        }
                }
                return builder.toString();
        }

        Trie getSubTrie(String key) {

                TrieNode temp = root;
                for (int level = 0; level < key.length(); level++) {
                        //for the key's length
                        int index = key.charAt(level) - 'a';
                        //goes to right character/node
                        if (temp.getOffspring()[index] == null) {
                                //checks if whole key is present
                                return null;
                        }
                        temp = temp.getOffspring()[index];
                }
                return new Trie(temp);
                //returns only if it exists AND is a key - not just a prefix.

        }

        public ArrayList<String> getAllWords() {
```

```
137                    //helper function, is used just for the ArrayList and
                          ↪ String/stringBuilder
138                    ArrayList<String> listOfWords = new ArrayList<>();
139                    getAllWords(listOfWords, new StringBuilder(), root);
140                    return listOfWords;
141            }
142
143        private void getAllWords(ArrayList listOfWords, StringBuilder builder,
144                                                    TrieNode root) {
145                    if (root.getIsEnd()) {
146                            //checks if it's reached the end
147                            //System.out.println(builder.toString());
148                            listOfWords.add(builder.toString());
149                            //print it and add it to the array
150                    }
151                    TrieNode[] offspring = root.getOffspring();
152                    //get's current node's children
153                    for (int i = 0; i < offspring.length; i++) {
154                            //for all the children
155                            if (offspring[i] != null) {     //if it's a valid node
156                                    getAllWords(listOfWords, builder.append((char)
157                                                    (97 + i)), offspring[i]);
158                                    //call itself, with the child as the root
159                                    builder.setLength(builder.length() - 1);
160                                    //goes up a level - resets stringbuilder to right
                                          ↪ place
161                            }
162                    }
163            }
164
165        public static void main(String args[]) {
166                    Trie potato = new Trie();
167                    System.out.println("Adding...");
168                    System.out.println(potato.add("bat"));
169                    System.out.println(potato.add("cat"));
170                    System.out.println(potato.add("chat"));
171                    System.out.println(potato.add("cheese"));
172                    System.out.println(potato.add("cheers"));
173
174                    System.out.println("Checking...");
175                    System.out.println(potato.contains("yeet"));
176                    System.out.println(potato.contains("ca"));
177                    System.out.println(potato.contains("cat"));
178
179                    System.out.println("Breadth...");
180                    System.out.println(potato.outputBreadthFirstSearch());
181                    System.out.println("Depth...");
182                    System.out.println(potato.outputDepthFirstSearch());
183
184                    System.out.println("getSubTrie \"ch\"");
185                    Trie potato2 = potato.getSubTrie("ch");
186
187                    System.out.println("Breadth...");
188                    System.out.println(potato2.outputBreadthFirstSearch());
189                    System.out.println("Depth...");
190                    System.out.println(potato2.outputDepthFirstSearch());
191
192                    System.out.println("Checking...");
193                    System.out.println(potato2.contains("eese"));
194                    System.out.println(potato2.contains("eers"));
```

```
195                    System.out.println(potato2.contains("at"));
196                    System.out.println(potato2.contains("yeet"));
197                    System.out.println(potato.getAllWords());
198            }
199  }
```

## 4.3  TrieNode

Listing 3: TrieNode.java

```
1   /*
2   By Robin Rai (100242165)
3   V.1.0.0
4   Created on 05/03/2020
5   */
6   package dsacoursework2;
7
8   public class TrieNode {
9           private TrieNode[] offspring;     //a is 0, b is 1, etc
10          private boolean isEnd;     //if the node is the end of a word/key
11
12          public TrieNode() {
13                  isEnd = false;
14                  this.offspring = new TrieNode[26];
15                  //I don't think we're getting any extra letters soon
16          }
17
18          public TrieNode[] getOffspring() {
19                  return this.offspring;
20          }
21
22          public boolean getIsEnd() {
23                  return this.isEnd;
24          }
25
26          public void setIsEnd(boolean input) {
27                  this.isEnd = input;
28          }
29
30  }
```

## 4.4  AutoCompletionTrie

Listing 4: AutoCompletionTrie.java

```
1   /*
2   By Robin Rai (100242165)
3   V.1.0.0
4   Created on 05/03/2020
5   */
6   package dsacoursework2;
7
8   import java.io.File;
9   import java.io.FileNotFoundException;
10  import java.util.*;
11
12  public class AutoCompletionTrie {
13          private AutoCompletionTrieNode root;
14
15          public AutoCompletionTrie() {
16                  root = new AutoCompletionTrieNode();
17          }
```

```java
18
19          public AutoCompletionTrie ( AutoCompletionTrieNode input ) {
20                  root = input ;
21          }
22
23          boolean add ( String key , int frequency ) {
24                  if ( key . equals ( "" )) {
25                          return false ;
26                  }
27                  int index ;
28                  boolean alreadyIn = true ;
29                  AutoCompletionTrieNode temp = root ;
30                  for ( int level = 0; level < key . length (); level ++) {
31                          // goes through all letters of the key
32                          index = key . charAt ( level ) - 'a';
33                          if ( temp . getOffspring ()[ index ] == null ) {
34                                  // if the node for the current letter doesn't exist
35                                  temp . getOffspring ()[ index ] = new
                                      ↪ AutoCompletionTrieNode ();
36                                  // make a new one
37                                  alreadyIn = false ;
38                                  // set the flag
39                          }
40                          temp = temp . getOffspring ()[ index ];
41                          // advanced temp node to one further down
42                  }
43                  if ( alreadyIn && temp . getIsEnd ()) {
44                          // if everything was already done , return false
45                          return false ;
46                  }
47                  temp . setIsEnd ( true );
48                  // sets last node for key as end .
49                  temp . setFrequency ( frequency );
50                  return true ;
51          }
52
53          boolean contains ( String key ) {
54                  int level ;
55                  int length = key . length ();
56                  int index ;
57                  AutoCompletionTrieNode temp = root ;
58                  for ( level = 0; level < length ; level ++) {
59                          // goes through all letters of the key
60                          index = key . charAt ( level ) - 'a';
61                          // index is current char's number
62                          if ( temp . getOffspring ()[ index ] == null ) {
63                                  // if the node for the current letter doesn't exist
64                                  return false ;
65                          }
66                          temp = temp . getOffspring ()[ index ];
67                          // advance down a level
68                  }
69                  return ( temp != null && temp . getIsEnd ());
70                  // returns true only if it exists AND is a key - not just a prefix .
71          }
72
73          String outputBreadthFirstSearch () {
74                  if ( root == null ) {
75                          return null ;
76                  }
```

```java
                    Queue<AutoCompletionTrieNode> queue = new LinkedList<>();
                    StringBuilder builder = new StringBuilder();
                    queue.add(root);
                    //starts with root
                    while (!queue.isEmpty()) {
                            AutoCompletionTrieNode temp = queue.remove();
                            //takes first thing in queue
                            for (int i = 0; i < 26; i++) {
                                    if (!(temp.getOffspring()[i] == null)) {
                                            //for all of it's children
                                            queue.add(temp.getOffspring()[i]);
                                            //add them to the queue so it goes through
                                                ↪ their
                                            // children - width first
                                            builder.append((char) (i + 97));
                                            //adds it to the output
                                    }
                            }
                    }
                    return builder.toString();
            }

        String outputDepthFirstSearch() {
                //helper function, provides string to append to
                StringBuilder builder = new StringBuilder();
                return outputDepthFirstSearch(builder, root);
        }

        private String outputDepthFirstSearch(StringBuilder builder,
                                                                                AutoComple
                                                                                   ↪ trie
                                                                                   ↪ {
                AutoCompletionTrieNode[] children = trieNode.getOffspring();
                //makes an array of all children in the tree
                for (int i = 0; i < 26; i++) {
                        //for all of a node's children
                        if (children[i] != null) {
                                //if there's a child
                                builder.append((char) (i + 97));
                                //add its character
                                outputDepthFirstSearch(builder, children[i]);
                                //call the function again on it's children
                        }
                }
                return builder.toString();
        }

        AutoCompletionTrie getSubTrie(String key) {

                AutoCompletionTrieNode temp = root;
                for (int level = 0; level < key.length(); level++) {
                        //for the key's length
                        int index = key.charAt(level) - 'a';
                        //goes to right character/node
                        if (temp.getOffspring()[index] == null) {
                                //checks if whole key is present
                                return null;
                        }
                        temp = temp.getOffspring()[index];
                }
```

```java
134                    return new AutoCompletionTrie(temp);
135                    //returns only if it exists AND is a key - not just a prefix.
136            }
137
138        public ArrayList<String> getAllWords() {
139                    //helper function, is used just for the ArrayList
140                    // and String/stringBuilder
141                    ArrayList<String> listOfWords = new ArrayList<>();
142                    getAllWords(listOfWords, new StringBuilder(), root);
143                    return listOfWords;
144            }
145
146        private void getAllWords(ArrayList listOfWords, StringBuilder sb,
147                                               AutoCompletionTrieNode root)
                                                   ↪ {
148                    if (root.getIsEnd()) {
149                            //checks if it's reached the end
150                            System.out.println(sb.toString());
151                            listOfWords.add(sb.toString());
152                            System.out.println(root.getFrequency());
153                            //print it and add it to the array
154                    }
155                    AutoCompletionTrieNode[] children = root.getOffspring();
156                    //get's current node's children
157                    for (int i = 0; i < children.length; i++) {
158                            //for all the children
159                            if (children[i] != null) {    //if it's a valid node
160                                    getAllWords(listOfWords, sb.append((char)
161                                            (97 + i)), children[i]);
162                                    //call itself, with the child as the root
163                                    sb.setLength(sb.length() - 1);
164                                    //goes up a level - resets stringbuilder to right
                                        ↪ place
165                            }
166                    }
167            }
168
169        public static class fullInfo {
170                    //objects to store a word's full info - the word
171                    // as well as it's frequency. Done for the comparator.
172                    private String key;
173                    private int freq;
174
175                    fullInfo(String key, int freq) {
176                            this.key = key;
177                            this.freq = freq;
178                    }
179
180                    fullInfo(String key, String freq) {
181                            this.key = key;
182                            this.freq = Integer.parseInt(freq);
183                    }
184
185                    int getFreq() {
186                            return freq;
187                    }
188
189                    String getKey() {
190                            return key;
191                    }
```

```java
192              }
193
194         static class testComp implements Comparator<fullInfo> {
195                 public int compare(fullInfo m1, fullInfo m2) {
196                         //compares by frequency, then by string
197                         if (m1.getFreq() < m2.getFreq()) {
198                                 return 1;
199                         } else if (m1.getFreq() > m2.getFreq()) {
200                                 return -1;
201                         } else {
202                                 return m1.getKey().compareTo(m2.getKey());
203                         }
204                 }
205         }
206
207         public ArrayList<fullInfo> getAllInfo() {
208                 //helper function to provide recursive method with string and array
209                 ArrayList<fullInfo> listOfInfo = new ArrayList<>();
210                 getAllInfo(listOfInfo, new StringBuilder(), root);
211                 listOfInfo.sort(new testComp());
212                 return listOfInfo;
213         }
214
215         private void getAllInfo(ArrayList listOfInfo, StringBuilder sb,
216                                              AutoCompletionTrieNode root) {
217                 if (root.getIsEnd()) {
218                         listOfInfo.add(new fullInfo(sb.toString(),
219                             ↪ root.getFrequency()));
219                         //if the end of the word's been reached, add it to the list
220                 }
221                 AutoCompletionTrieNode[] children = root.getOffspring();
222                 //get's current node's children
223                 for (int i = 0; i < children.length; i++) {
224                         //for all the children
225                         if (children[i] != null) {
226                                 //if it's a valid node
227                                 getAllInfo(listOfInfo, sb.append((char)
228                                             (97 + i)), children[i]);
229                                 //call itself, with the child as the root
230                                 sb.setLength(sb.length() - 1);
231                                 //goes up a level - resets stringbuilder to right
232                                     ↪ place
232                         }
233                 }
234         }
235
236         public void addToTrie(String file) throws FileNotFoundException {
237                 //adds a dictionary file to a trie using trie's add function
238                 Scanner sc = new Scanner(new File(file));
239                 sc.useDelimiter("\n");
240                 String str;
241                 while (sc.hasNext()) {
242                         str = sc.next();
243                         str = str.trim();
244                         str = str.toLowerCase();
245                         //reads in and trims line
246                         String[] parts = str.split(",");
247                         //splits word and frequency, and adds them to the trie
248                         //System.out.println("adding: " + parts[0] + " " + parts[1]);
249                         //System.out.println(add(parts[0],
```

```
                                  ↪ Integer.parseInt(parts[1])));
250                          add(parts[0], Integer.parseInt(parts[1]));
251                  }
252          }
253 }
```

## 4.5 AutoCompletionTrieNode

Listing 5: AutoCompletionTrieNode.java

```java
1  /*
2  By Robin Rai (100242165)
3  V.1.0.0
4  Created on 05/03/2020
5  */
6  package dsacoursework2;
7
8  public class AutoCompletionTrieNode {
9          private AutoCompletionTrieNode[] offspring;
10         //a is 0, b is 1, etc
11         private boolean isEnd;
12         //if the node is the end of a word/key
13         private int frequency;
14         //added frequency of word - is to be used at the end of a word like isEnd
15
16
17         public AutoCompletionTrieNode() {
18                 isEnd = false;
19                 this.offspring = new AutoCompletionTrieNode[26];
20                 frequency = 0;
21         }
22
23         public AutoCompletionTrieNode[] getOffspring() {
24                 return this.offspring;
25         }
26
27         public boolean getIsEnd() {
28                 return this.isEnd;
29         }
30
31         public void setIsEnd(boolean input) {
32                 this.isEnd = input;
33         }
34
35         public int getFrequency() {
36                 return frequency;
37         }
38
39         public void setFrequency(int input) {
40                 this.frequency = input;
41         }
42  }
```

## 4.6 AutoComplete

Listing 6: AutoComplete.java

```java
1  /*
2  By Robin Rai (100242165)
3  V.1.0.0
4  Created on 05/03/2020
5  */
```

```java
package dsacoursework2;

import java.io.IOException;
import java.util.ArrayList;

import static dsacoursework2.DictionaryMaker.*;
import static dsacoursework2.AutoCompletionTrie.*;

public class AutoComplete {

        public static void autoCompletion(AutoCompletionTrie masterTrie,
            ↪ ArrayList<String> prefixes, String saveLocation) throws IOException {
                //generates top three results for each query, prints them to console
                    ↪ and file
                String infoString = "";
                //string to be written to file
                for (int prefix = 0; prefix < prefixes.size(); prefix++) {
                        //for every prefix
                        System.out.println("\nResults for: " + prefixes.get(prefix));
                        AutoCompletionTrie currentSubTrie =
                            ↪ masterTrie.getSubTrie(prefixes.get(prefix));
                        //gets the subTrie for the prefix
                        ArrayList<fullInfo> info = currentSubTrie.getAllInfo();
                        //gets the information for the subTrie
                        int totalFreq = 0;
                        //counter
                        for (int i = 0; i < info.size(); i++) {
                                totalFreq += info.get(i).getFreq();
                                //sums up frequencies to calculate ratios
                        }
                        infoString += prefixes.get(prefix) + ",";
                        //starts string to write to file
                        for (int i = 0; i < 3 && i < info.size(); i++) {
                                //gets top three or all available sub words,
                                    ↪ whichever's smallest
                                infoString += prefixes.get(prefix) +
                                    ↪ info.get(i).getKey() + "," +
                                    ↪ info.get(i).getFreq() + "," + (double)
                                    ↪ info.get(i).getFreq() / totalFreq + ",";
                                //creates a line with all words under that prefix
                                    ↪ along with their info to write to file
                                System.out.println(prefixes.get(prefix) +
                                    ↪ info.get(i).getKey() + " (probability " +
                                    ↪ (double) info.get(i).getFreq() / totalFreq +
                                    ↪ ")");
                                //prints said info to console as well
                        }
                        infoString = infoString.substring(0, infoString.length() - 1);
                        //removes comma from last result, unnecessary.
                        infoString += "\n";
                        //new line for new prefix
                }
                saveToFile(infoString, saveLocation);
                //calls saveToFile from DictionaryMaker to save to file
        }

        public static void main(String[] args) throws IOException {
                dsacoursework2.DictionaryMaker df = new
                    ↪ dsacoursework2.DictionaryMaker();
                df.formDictionary("src\\TextFiles\\lotr.csv",
```

```
                                 ↪ "src\\TextFiles\\lotrDic.csv");
54              AutoCompletionTrie completionTest = new AutoCompletionTrie();
55              completionTest.addToTrie("src\\TextFiles\\lotrDic.csv");
56              autoCompletion(completionTest,
                         ↪ readWordsFromCSV("src\\TextFiles\\lotrQueries.csv", "\n"),
57                                  "src\\TextFiles\\lotrMatches.csv");
58
59          }
60
61  }
```

## 4.7 AutoComplete output

Listing 7: lotrMatches.csv

```
1  ab,about,17,0.5666666666666667,above,9,0.3,able,3,0.1
2  go,going,15,0.2777777777777778,go,13,0.24074074074074073,good,9,0.16666666666666666
3  the,the,460,0.6267029972752044,they,113,0.1539509536784741,them,50,0.0681198910081744
4  mer,merry,36,0.9473684210526315,merely,1,0.02631578947368421,merrily,1,0.02631578947368421
5  fro,frodo,27,0.4909090909090909,from,24,0.43636363636363634,front,4,0.07272727272727272
6  gr,great,13,0.19696969696969696,ground,12,0.18181818181818182,grass,10,0.15151515151515152
7  gol,goldberry,3,0.6,golden,2,0.4
8  sam,sam,18,1.0
```