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

Student number: 100238596. Blackboard ID: pxh18ksu

Thursday 14<sup>th</sup> May, 2020 03:52

## Contents

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

The method formDictionary takes in a List of strings and returns a TreeMap¡String, Integer¿ as an output. The string is the "word" from the list, and the value is the amount of times the word appears in the list

---

**Algorithm 1** formDictionary algorithm

---

**Require:** List of String words
1: **for** String *word* in *words* **do**
2:     **if** treeMap contains word **then**
3:         Add word, key + 1 to treeMap
4:     **else**
5:         Add word, 1 to treeMap
    **return** *treeMap*

---

## 1.1 Fundamental Operation

The fundamental operation for the algorithm is the "if treeMap contains" comparison as it is run every time the for loop runs

## 1.2 Run time complexity function

$$\sum_{i=1}^{n} n \tag{1}$$

$$\tag{2}$$

## 1.3 Worst case scenario

The worst case is that the words in the list passed are all words not already present within the treeMap

# 2 Part 2: Implement a Trie Data Structure

## 2.1 Add method for adding a key to a trie

---
**Algorithm 2** add algorithm

---
**Require:** String key
1: TrieNode $currentNode \leftarrow root$
2: **for** every char $c$ in $key$ **do**
3:      **if** $currentNode$ != null **then**
4:          $next \leftarrow currentNode$.getOffspring(c)
5:          **if** next == null **then**
6:              $next \leftarrow TrieNode.newNode(c)$
7:              $currentNode.addNode(next)$
8:          $currentNode \leftarrow next$
9: **if** $currentNode$.isComplete() **then return** false
10: **else**
11:      $currentNode$.setComplete(true)
12:      **return** true

---

## 2.2 Contains method to check if key passed in is a full word or a prefix

---
**Algorithm 3** contains algorithm

---
**Require:** String $key$
1: **for** char $c$ in $key$ do **do**
2:      **if** $currentNode$.getOffspring(c) == null **then return** false
3:      **else**
4:          $currentNode \leftarrow currentNode$.getOffspring(c)
     **return** $currentNode$.isComplete()

---

## 2.3 Output by Breadth First Search Method

---
**Algorithm 4** outputBreadthFirstSearch algorithm

---
**Require:** No Input
1: Queue $Nodes \leftarrow EmptyQueue$
2: ArrayList $characterArrayList \leftarrow EmptyArrayList$
3: $nodes$.add(root)
4: **while** $nodes$.isEmpty() = false **do**
5:      TrieNode $next \leftarrow nodes$.poll()
6:      **if** $next$.getOffspring != null **then**
7:          **for** each TrieNode node : $next$.getOffspring() **do**
8:              **if** $node$ != null **then**
9:                  $nodes$.add($node$)
10:          $characterArrayList$.add($next$.getCharacter())
11: String $output$
12: **for** char $c$ in $characterArrayList$ **do**
13:      $output \leftarrow output + c$
14: **return** output

---

## 2.4 Output by Depth First Search Method

**Algorithm 5** outputDepthFirstSearch algorithm

---

**Require:** TrieNode *next*, StringBuilder *stringBuilder*
1: *next.visited* ← *true*
2: **for** *i* ← 1 to *next*.offspring.length-1 **do**
3:     **if** *next*.offspring[i] != null and !*next*.offspring.visited **then**
4:         outputDepthFirstSearch(*next*.offspring[i], *stringBuilder*)
5: *stringBuilder*.append(*next*.getCharacter())
6: return result

---

## 2.5   getSubTrie method to return a trie rooted at the prefix passed

**Algorithm 6** getSubTrie algorithm

---

**Require:** String *prefix*
1: TrieNode *next* ← *root*
2: Trie *newTrie* ← *newTrie*()
3: **for** *i* ← 1 to *prefix*.charAt(i) **do**
4:     *index* ← 1
5:     **if** *next* == null **then return** null
6:     **if** *next.getOffspring(index)*! = *null* **then**
7:         *newTrie*.root ← *next*.getOffspring(index)
8:     *next* ← *next*.offspring[index]
    **return** *newTrie*

---

## 2.6   getAllWords method to return all the words in the trie

**Algorithm 7** getAllWords algorithm

---

**Require:** String *wordSoFar*, TrieNode *currentNode*, List of Strings *nodes*
1: **for** TrieNode *temp* in *currentNode*.getOffspring() **do**
2:     **if** *temp* != null **then**
3:         String *currentPrefix* ← *wordSoFar* + *temp*.getCharacter()
4:         getAllWords(*currentPrefix*, *temp*, *nodes*);
5: **if** *currentNode*.isComplete() **then**
6:     *nodes*.add(*wordSoFar*)
    **return** output

---

# 3 Part 3: Word Auto Completion Application

## 3.1 Auto Completion Method

---
**Algorithm 8** AutoCompletion algorithm

---
1: List of Strings $dictWordsList \leftarrow readWordsFromCSV(file)$
2: List of Strings $LotrQueries \leftarrow readWordsFromCSV(file)$
3: TreeMap(String, Integer) $dictWordsMap \leftarrow formDictionary(dictWordsList)$
4: Trie $dictTrie \leftarrow dictWordsList$
5: NavigableMap(String, Integer) $navMap \leftarrow newTreemap$
6: **for** String $s$ in $LotrQueries$ **do**
7:     Trie $subTrie \leftarrow dictTrie$.getSubTrie(s)
8:     List of Strings $list \leftarrow temp.getAllWords()$
9:     **for** String $listString$ in $list$ **do**
10:        String $auto \leftarrow s.trim() + listString.trim()$
11:        **for** every $entry$ of Map(String, Integer) in $dictWordsMap$.entrySet() **do**
12:            **if** $auto$.equals($entry$.getKey() **then**
13:                $navMap$.put($entry$.getKey(), $entry$.getValue())
14:     LinkedHashMap(String, Float) $tempList \leftarrow sortByFrequency(navMap)$                 ▷ *Described below*
15:     LinkedHashMap(String, Float) $finalList$.putAll($tempList$)
16:     $navMap$.clear()
17: $DictionaryFinder$.saveToFile($finalList$, ”$lotrMatches.csv$”)

---

Forms a dictionary and word count of words from the file passed
Constructs a trie using this dictionary
Loads prefixes from LotrQueries into a list
For each prefix:
Recover all possible words
Choose 3 most frequent and display to standard output
Write results to file

sortByFrequency is a function created to mainly sort the list of possible words by how frequently they appear (their value in the Map). It has a running tally of frequency each time it is run, so it is possible to divide each key's value by that total to work out the probability. This uses the Java List built in sort and comparator.

Once the map has been sorted, it fetches the front 3 (most populous), adds them to a LinkedHashMap¡String, Float¿, and returns it to the main method.

## 3.2 Auto Completion Results

Listing 1: lotrMatches.csv

```
 1  about =0.56666666
 2  above =0.3
 3  able =0.1
 4  going =0.2777778
 5  go =0.24074075
 6  good =0.16666667
 7  the =0.626703
 8  they =0.15395096
 9  them =0.06811989
10  merry =0.94736844
11  merrily =0.02631579
12  merely =0.02631579
13  frodo =0.4909091
14  from =0.43636364
15  front =0.07272727
16  great =0.1969697
17  ground =0.18181819
18  grass =0.15151516
```

```
19  goldberry =0.6
20  golden =0.4
21  sam =1.0
```

# 4   Code Listing

Listing 2: DictionaryFinder.java

```java
1
2  package com.company;
3
4  import java.io.*;
5  import java.util.*;
6
7  /**
8   *
9   * @author ajb
10  */
11 public class DictionaryFinder {
12
13     //String is the word, Int is the number of occurrences
14     TreeMap<String, Integer> dict;
15
16     public DictionaryFinder(){
17     }
18     /**
19      * Reads all the words in a comma separated text document into an Array
20      * @param file
21      */
22     public static ArrayList<String> readWordsFromCSV(String file) throws
          ↪ FileNotFoundException {
23         Scanner sc=new Scanner(new File(file));
24         sc.useDelimiter(" |,");
25         ArrayList<String> words=new ArrayList<>();
26         String str;
27         while(sc.hasNext()){
28             str=sc.next();
29             str=str.trim();
30             str=str.toLowerCase();
31             words.add(str);
32         }
33         return words;
34     }
35
36     public static ArrayList<String> readWordsFromCSVNewLine(String file) throws
          ↪ FileNotFoundException {
37         Scanner sc=new Scanner(new File(file));
38         sc.useDelimiter("\n");
39         ArrayList<String> words=new ArrayList<>();
40         String str;
41         while(sc.hasNext()){
42             str=sc.next();
43             str=str.trim();
44             str=str.toLowerCase();
45             words.add(str);
46         }
47         return words;
48     }
49     public static void saveCollectionToFile(Collection<?> c,String file) throws
          ↪ IOException {
50         FileWriter fileWriter = new FileWriter(file);
51         PrintWriter printWriter = new PrintWriter(fileWriter);
52         for(Object w: c){
53             printWriter.println(w.toString());
```

```java
54              }
55              printWriter.close();
56          }
57      public TreeMap<String, Integer> formDictionary(List<String> words) throws
            ↪ FileNotFoundException {
58          Collections.sort(words);
59          dict = new TreeMap<String, Integer>();
60          for(String word : words){
61              //If word already exists, increment counter
62              if(dict.containsKey(word)){
63                  dict.put(word, dict.get(word) + 1);
64              }
65              //Else, add word and set occurrences to 1
66              else{
67                  dict.put(word, 1);
68              }
69          }
70          return dict;
71      }
72
73  public static void saveToFile(LinkedHashMap map,String file) throws IOException {
74      FileWriter fileWriter = new FileWriter(file);
75      PrintWriter printWriter = new PrintWriter(fileWriter);
76      for(Object w: map.entrySet()){
77          printWriter.println(w);
78      }
79      printWriter.close();
80  }
81
82      public static void main(String[] args) throws Exception {
83          DictionaryFinder df=new DictionaryFinder();
84          //ArrayList<String> in=readWordsFromCSV("C:\\Teaching\\2017-2018\\Data
                ↪ Structures and Algorithms\\Coursework 2\\test.txt");
85          //DO STUFF TO df HERE in countFrequencies
86          //df.formDictionary(in);
87          //df.saveToFile();
88
89          Trie test = new Trie();
90          System.out.println("Testing adding words to the trie - all should be T apart
                ↪ from the last");
91          System.out.println(test.add("cheers"));
92          System.out.println(test.add("cheese"));
93          System.out.println(test.add("chat"));
94          System.out.println(test.add("cat"));
95          System.out.println(test.add("bat"));
96          System.out.println(test.add("bat"));
97          System.out.println("\n");
98          System.out.println("Testing the contains method - should produce F,F,F,T,T");
99          System.out.println(test.contains("chee"));
100         System.out.println(test.contains("afc"));
101         System.out.println(test.contains("ba"));
102         System.out.println(test.contains("cheese"));
103         System.out.println(test.contains("bat"));
104         System.out.println("\n");
105         System.out.println("Testing the breadth first search method - should produce
                ↪ bcaahttaetersse");
106         System.out.println(test.outputBreadthFirstSearch());
107         System.out.println("\n");
108         System.out.println("Testing the depth first search method - should produce
                ↪ batcathateersse");
```

```
109        System.out.println(test.outputDepthFirstSearch());
110        System.out.println("\n");
111        Trie subtrie = test.getSubTrie("ch");
112        System.out.println("Testing the subTrie method via breadth search - should
             ↪ produce haetersse");
113        System.out.println(subtrie.outputBreadthFirstSearch());
114        System.out.println("\n");
115        System.out.println("Testing getAllWords - should produce bat, cat, chat,
             ↪ cheers, cheese");
116        System.out.println(test.getAllWords());
117
118        System.out.println("\n");
119        System.out.println("Now running AutoCompletion");
120        AutoCompletion.main();
121
122    }
123
124 }
```

```
1   package com.company;
2
3   public class TrieNode {
4
5       private char character;
6       TrieNode[] offspring = new TrieNode[26];
7       private boolean isComplete;
8       Boolean visited;
9
10      public TrieNode(){
11          this.offspring = new TrieNode[26];
12          visited = false;
13      }
14
15      public TrieNode(char character){
16          this.character = character;
17          this.offspring = new TrieNode[26];
18          for(int i  = 0; i < offspring.length; i++){
19              offspring[i] = null;
20          }
21          visited = false;
22      }
23
24
25      public boolean isComplete() {
26          return this.isComplete;
27      }
28
29      public char getCharacter() {
30          return character;
31      }
32
33      public TrieNode getNode(char c) {
34          return this.offspring[getCharIndex(c)];
35      }
36
37      public void setCharacter(char character) {
38          this.character = character;
39      }
40
41      public void setOffspring(TrieNode[] offspring) {
42          this.offspring = offspring;
43      }
44
45      public static int getCharIndex(char c) {
46          return c - 'a';
47      }
48
49      public static TrieNode newNode(char c){
50          TrieNode newNode = new TrieNode();
51          newNode.isComplete = false;
52          newNode.character = c;
53          for (int i = 0; i < newNode.offspring.length; i++){
54              newNode.offspring[i] = null;
55          }
56          return newNode;
57      }
58      public void addNode(TrieNode next){
59          int index = (int)next.character -97;
```

```java
            offspring[index] = next;
        }

    public void setComplete(boolean complete) {
        isComplete = complete;
    }

    public TrieNode[] getOffspring(){
        return this.offspring;
    }

    public TrieNode getOffspring(char c){
        for(int i = 0; i < offspring.length; i++){
            if (offspring[i] != null && offspring[i].character == c){
                return offspring[i];
            }
        }
        return null;
    }

}
```

Listing 4: Trie.java

```java
package com.company;

import java.util.ArrayList;
import java.util.LinkedList;
import java.util.List;
import java.util.Queue;

public class Trie extends TrieNode {

    private TrieNode root;

    public Trie() {
        root = new TrieNode();
    }


    /* let current node = root node
    for each letter in the key
        find the child node of current node associated with that letter
        if there is no child node associated with that letter, create a new node and
            ↪ add it to current node as a child associated with the letter
        set current node = child node
    add value to current node */
    public boolean add(String key) {
        TrieNode currentNode = root;
        for (int i = 0; i < key.length(); i++) {
            char currentChar = key.charAt(i);
            if (currentNode != null) {
                TrieNode next = currentNode.getOffspring(currentChar);
                if (next == null) {
                    next = TrieNode.newNode(key.charAt(i));
                    currentNode.addNode(next);
                }
                currentNode = next;
            }
        }

        // returns false if key is already in the trie
        if (currentNode.isComplete()) {
            return false;
        }

        currentNode.setComplete(true);
        return true;
    }


    public boolean contains(String key) {
        TrieNode currentNode = root;
        for (int i = 0; i < key.length(); i++) {
            char currentChar = key.charAt(i);
            TrieNode next = currentNode.getOffspring(currentChar);
            if (next == null) {
                return false;
            } else {
                currentNode = next;
            }
        }
        return currentNode.isComplete();
```

```java
59      }
60
61      public String outputBreadthFirstSearch() {
62          //Breadth first done using a queue
63          Queue<TrieNode> nodes = new LinkedList();
64          ArrayList<Character> characterArrayList = new ArrayList();
65          nodes.add(root); //add root node to the queue
66          while (!nodes.isEmpty()) { //while queue isn't empty
67              TrieNode next = nodes.poll(); //set next node to item front of queue
68
69              if (next.getOffspring() != null) {
70                  for (TrieNode node : next.getOffspring()) {
71                      //if node is not null
72                      if (node != null) {
73                          //add a node to nodes linkedlist
74                          nodes.add(node);
75
76                      }
77                  }
78                  characterArrayList.add(next.getCharacter()); //add character value of
                          ↪ node to an arraylist
79              }
80
81          }
82
83          //Adding the arraylist to a string
84          String outputBreadth = "";
85          for (char character : characterArrayList) {
86              outputBreadth = outputBreadth + character;
87          }
88          return outputBreadth;
89      }
90
91
92      public String outputDepthFirstSearch(){
93          StringBuilder builder = new StringBuilder();
94          outputDepthFirstSearch(root, builder);
95          return builder.toString();
96      }
97
98      private static void outputDepthFirstSearch(TrieNode next, StringBuilder str){
99          next.visited = true;
100         for(int i = 0; i < next.offspring.length-1; i++){
101             if(next.offspring[i] != null && !next.offspring[i].visited){
102                 outputDepthFirstSearch(next.offspring[i],str);
103             }
104         }
105         str.append(next.getCharacter());
106     }
107
108     public Trie getSubTrie(String prefix){
109         TrieNode next = root;
110         Trie newTrie = new Trie();
111         for(int i = 0; i < prefix.length(); i++){
112             int index = (int)prefix.charAt(i)-97;
113             if(next == null){
114                 return null;
115             }
116             if(next.getOffspring(prefix.charAt(i)) != null){
117                 newTrie.root = next.getOffspring(prefix.charAt(i));
```

```java
118              }
119              next = next.offspring[index];
120          }
121          return newTrie;
122      }
123
124      public List<String> getAllWords() {
125          List<String> output = new LinkedList<>();
126          getAllWords("", root, output);
127          return output;
128      }
129
130      private void getAllWords(String wordSoFar, TrieNode currentNode,
131                              List<String> nodes) {
132          for (TrieNode temp : currentNode.getOffspring()) {
133              if (temp != null) {
134                  String currentPrefix = wordSoFar + temp.getCharacter();
135                  getAllWords(currentPrefix, temp, nodes);
136              }
137          }
138
139          if (currentNode.isComplete()) {
140              nodes.add(wordSoFar);
141          }
142      }
143
144  }
```

Listing 5: AutoCompletion.java

```java
package com.company;

import java.io.FileNotFoundException;
import java.io.IOException;
import java.util.*;

import static com.company.DictionaryFinder.readWordsFromCSV;
import static com.company.DictionaryFinder.readWordsFromCSVNewLine;

public class AutoCompletion {

    public static void main() throws FileNotFoundException,
            IOException {

        DictionaryFinder df = new DictionaryFinder();

        NavigableMap<String, Integer> storeStringMap
                = new TreeMap<String, Integer>();

        LinkedHashMap<String, Float> finalList = new LinkedHashMap<>();

        String auto;

        //list of words from dictionary file
        List<String> dictWordsList = readWordsFromCSV("lotr.csv");
        TreeMap<String, Integer> dictWordsMap = df.formDictionary(dictWordsList);

        // dictionary trie
        Trie dictTrie = new Trie();
        for (String s : dictWordsList){
            dictTrie.add(s);
        }


        // list of words from query file
        List<String> LotrQueries = readWordsFromCSVNewLine("lotrQueries.csv");


        //for each prefix in query file
        for(String s : LotrQueries){
            //get a subtrie rooted at the prefix
            Trie temp = dictTrie.getSubTrie(s);
            //create a list of words from the new trie - will be missing the prefix
            //    though as new root: e.g. eese
            List<String> list = temp.getAllWords();
            //re-add the prefix onto the list: e.g. ch + eese = cheese
            for (String listString : list){
                auto = s.trim() + listString.trim();

                for (Map.Entry<String, Integer> entry : dictWordsMap.entrySet()) {
                    //if words that are in auto equal to the words in entry map
                        if (auto.equals(entry.getKey())) {
                            //store those words in a storeAuto map
                            storeStringMap.put(entry.getKey(), entry.getValue());
                    }
                }

            }
```

```
59              LinkedHashMap <String, Float > tempList  = sortByFrequency(storeStringMap);
60              finalList.putAll(tempList);
61              storeStringMap.clear();
62          }
63          DictionaryFinder.saveToFile(finalList, "lotrMatches.csv");
64      }
65
66
67      public static LinkedHashMap<String, Float> sortByFrequency(NavigableMap<String,
            ↪ Integer > dictionary) throws IOException {
68          List <Map.Entry<String, Integer >> dictList = new
                ↪ LinkedList<>(dictionary.entrySet());
69          LinkedHashMap<String, Float> map = new LinkedHashMap<>();
70
71          dictList.sort(Comparator.comparingInt(Map.Entry::getValue));
72          Collections.reverse(dictList);
73
74          float totalFreq = 0;
75          for(Map.Entry<String, Integer> item : dictList){
76              totalFreq = totalFreq + item.getValue();
77          }
78
79          for(int i=0; i < 3; i++){
80              try{
81                  float probability = (float)dictList.get(i).getValue() / totalFreq;
82                  map.put(dictList.get(i).getKey(), probability);
83                  System.out.println(dictList.get(i).getKey() + " (probability " +
                        ↪ probability + ")");
84              }
85              catch (Exception ignored){
86              }
87
88          }
89
90          return map;
91      }
92
93
94 }
```