

Data Structures and Algorithms

Coursework Assignment 2

Patryk Dobbek

stp18rsu

100023818

April 23, 2020

1 DictionaryMaker

This algorithm populates and returns a dictionary structure created from the the.

1.1 The 'formDictionary' algorithm

Algorithm 1 formDictionary(**doc**, n) **return** ($sortedDict, m$)

Require: list of strings **doc** of length n

Ensure: $sortedDict :=$, Dictionary of string keys and integer values of size m .

```
1:  $sortedDict$   $\triangleright$  empty dictionary with string keys and integer values.
2: for  $i \leftarrow 1$  to  $n$  do
3:    $temp \leftarrow doc_i$ 
4:   if  $sortedDict$  contains  $temp$  then
5:      $sortedDict$  put ( $key \leftarrow temp, value \leftarrow value + 1$ )
6:   else
7:      $sortedDict$  put ( $key \leftarrow temp, value \leftarrow 1$ )
8:   end if
9: end for
10: return ( $sortedDict, m$ )
```

1.2 Algorithm analysis for 'formDictionary' algorithm

Fundamental operation: (Line 3) **If** $sortedDict$ **contains** $temp$

Contain operations for the dictionary used in this algorithm (Java's red-black map-based tree, TreeMap) are $O(\log(m))$.

1.2.1 Worst case

The worst case for this algorithm is when every word in **doc** is unique, thereby making $m = n$ and **contains** operation on line 3 equal to $O(\log(n))$.

$$t(n) = \sum_{i=1}^n \log(n) = n \log(n)$$

$$t(n) = n \log(n)$$

Which gives us linearithmic running time and $\mathbf{O(n \log(n))}$ as the order of the algorithm.

1.2.2 Best case

The best case for this algorithm is for **doc** to contain repetitions of the same word, thereby making **m = 1** and **contains** operation on line 3 equal to $\mathbf{O(1)}$.

$$t(n) = \sum_{i=1}^n 1 = n$$

From here we can categorise the run-time function as linear and give $\mathbf{O(n)}$ as the order of the algorithm.

2 Trie

2.1 The 'add' algorithm

This function adds a string into the trie. It ensures space efficiency by initialising nodes and child arrays only when needed.

Algorithm 2 add(**key**, *n*, **root**) **return** *boolean*

Require: lower case alphabetical character array **key** of length *n*.

Require: reference to TrieNode **root** (root of a Trie structure)

Ensure: *boolean*, true if operation successful, false if key already in Trie.

```

1: curr ← root                                ▷ saving a reference to the root node in this Trie.
2: for i ← 1 to n do
3:   charIndex ← keyi − 98                      ▷ should give range between 1-26
4:   if curr.children equals null then
5:     curr.children := ▷ initialise space for references to other TrieNodes
6:   end if
7:   if curr.childrencharIndex equals null then
8:     curr.childrencharIndex := ▷ initialise TrieNode
9:   end if
10:  curr ← curr.childrencharIndex
11: end for
12: if curr.isWord equals false then
13:  curr.isWord ← true
14: else return false
15: end if
16: return true

```

2.2 Algorithm analysis for 'add' algorithm

Fundamental operation: (Line 8) $curr.children_{charIndex} :=$

Fundamental operation is constant time.

2.2.1 Worst case

None of the characters present in **key** are present in the Trie in the same order. This case would require n fundamental operations.

$$t(n) = \sum_{i=1}^n 1 = n$$

Which gives us $O(n)$ as the order of the algorithm and linear running time.

2.3 The 'contains' algorithm

This algorithm searches for the specified word and returns a boolean value depending on whether it is found within the structure.

Algorithm 3 contains(**key**, n , **root**) **return** *boolean*

Require: lower case ASCII character array **key** of length n .

Require: reference to a TrieNode **root** (root of a Trie structure)

Ensure: *boolean*, true if the word passed in is in the Trie as a whole word, not just prefix. Otherwise, false.

```
1:  $curr \leftarrow root$  ▷ saving a reference to the root node in this Trie.
2: for  $i \leftarrow 1$  to  $n$  do
3:    $charIndex \leftarrow key_i - 98$  ▷ should give range between 1-26
4:   if  $curr.children_{charIndex}$  equals null then
5:     return false
6:   end if
7:   if  $curr.children_{charIndex}$  equals null then
8:     return false
9:   end if
10:   $curr \leftarrow curr.children_{charIndex}$ 
11: end for
12: if  $curr.isWord$  equals false then
13:   return false
14: end if
15: return true
```

2.4 Algorithm analysis for 'contains' algorithm

Fundamental operation: (Line 7) **If** $curr.children$ **equals** *null*

Fundamental operation is constant time.

2.4.1 Worst case

In the worst case the series of characters passed in as a character array are present in the Trie. This case would require n fundamental operations.

$$t(n) = \sum_{i=1}^n 1 = n$$

Which gives us $O(n)$ as the order of the algorithm and linear running time.

2.5 The 'getSubTrie' algorithm

This algorithm returns a trie structure that is a subset of another trie. Returns null when prefix string cannot be found with the original trie.

Algorithm 4 getSubTrie(**prefix**, n , **root**) **return** ($subTrie, m$)

Require: lower case ASCII character array **key** of length n .

Require: reference to TrieNode **root** (root of a Trie structure).

Ensure: $subTrie$, Trie data structure of size m .

```
1:  $curr \leftarrow root$   $\triangleright$  saving a reference to the root node in this Trie.
2: for  $i \leftarrow 1$  to  $n$  do
3:    $charIndex \leftarrow key_i - 98$   $\triangleright$  should give range between 1-26
4:   if  $curr.children$  equals  $null$  then
5:     return  $null$ 
6:   end if
7:   if  $curr.children_{charIndex}$  equals  $null$  then
8:     return  $null$ 
9:   end if
10:   $curr \leftarrow curr.children_{charIndex}$ 
11: end for
12:  $subTrie :=$   $\triangleright$  initialising an empty Trie
13:  $subTrie.root \leftarrow curr$   $\triangleright$  setting the root node in subTrie
14: return  $subTrie$ 
```

2.6 The 'outputBreadthFirstSearch' algorithm

Algorithm 5 outputBreadthFirstSearch(**root**,**q**) **return** (*str*,*n*)

Require: reference to TrieNode **root** (root of a Trie structure). Every TrieNode can hold a fixed array of 26 references to other TrieNodes called **children**

Ensure: *str*, string of lower-case alphabetical characters of size *n* (*n* equals to the size of the Trie).

```
1: q                                ▷ empty queue structure that can hold TrieNodes
2: add root to q (enqueue)
3: str :=                             ▷ declaring a new empty string
4: TrieNode                          ▷ declaring a new TrieNode
5: while q is not empty do
6:   TrieNode ← dequeue q
7:   if TrieNode.children is not null then
8:     for i ← 1 to 26 do
9:       if TrieNode.childreni is not null then
10:        str add ASCII char equal to i + 98
11:        add child to q (enqueue)
12:      end if
13:    end for
14:  end if
15: end while
16: return (str,n)
```

2.7 The 'outputDepthFirstSearch' algorithm

The Trie function 'outputDepthFirstSearch' is not described here because it simply calls on a recursive function 'getDepthFirstString' of its root TrieNode.

Algorithm 6 getDepthFirstString() **return** (*str*,*n*)

Require: this TrieNode's fixed-sized array of 26 references to other TrieNodes, **children**.

Ensure: *str*, string of lower-case alphabetical characters of size *n* (*n* equals to the amount of offspring TrieNodes).

```
1: str :=                             ▷ declare a new string
2: if children is null then return (str,0)    ▷ return an empty string
3: end if
4: for i ← 1 to 26 do
5:   if childreni is not null then
6:     str add ASCII char equal to i + 98
7:     str add childreni.getDepthFirstString()    ▷ recursive call
8:   end if
9: end for
10: return (str,n)
```

2.8 The 'findAllWords' algorithm

The Trie function '*findAllWords*' is not described here because it simply calls on a recursive function '*findWords*' of its root TrieNode.

Algorithm 7 *findWords(strList, n, str, m)* **return void**

Require: **strList**, list of n complete word strings found so far.

Require: **str**, string of lower-case alphabetical characters of size m .

Require: this TrieNode's fixed-sized array of 26 references to other TrieNodes, **children**.

Require: **isWord**, a boolean indicating whether this node is marks a complete word.

Ensure: **strList** populated with words found.

```
1: str := ▷ declare a new string
2: if children is null then
3:   for  $i \leftarrow 1$  to 26 do
4:     if children $i$  is not null then
5:       temp  $\leftarrow$  str add ASCII char equal to  $i + 98$ 
6:       if children $i$  isWord then
7:         strList add temp
8:       end if
9:       children $i$ .findWords(strList, temp)
10:    end if
11:  end for
12: end if
13: return
14:
```

3 AutoCompleteTrie

3.1 The 'populateFrequencyMap' algorithm

This is a recursive function retrieves all words and word frequencies within the AutoCompleteTrie using a depth first search. Those words are stored in a map of integer keys (frequencies) and string values (words).

This function is called by AutoCompleteTrie's '*getFrequencyMap()*' method, which is not described here because it simply declares a Sorted Dictionary and calls its root node's '*populateFrequencyMap*' function.

Algorithm 8 populateFrequencyMap(**sortedDict**, n ,**str**, m) **return** void

Require: this TrieNode's fixed-sized array of 26 references to other TrieNodes, **children**.

Require: **str**, string of lower-case alphabetical characters of size m .

Require: dictionary with integer keys and values of list of strings (representing words found so far), **sortedDict**, of size n .

Ensure: *sortedDict* populated with integer keys (corresponding to words frequency) and list of strings values.

```
1: if children is not null then
2:   for  $i \leftarrow 1$  to 26 do
3:     if children $_i$  is not null then
4:       if children $_i$  isWord then
5:         childFrequency  $\leftarrow$  children $_i$ .frequency
6:         list  $\leftarrow$  sortedDict $_{childFrequency}$ 
7:         if list is null then
8:           list := ▷ initialise list
9:         end if
10:        list add (str add ASCII char equal to  $i + 98$ )
11:        sortedDict $_{childFrequency}$   $\leftarrow$  list
12:      end if
13:      str add ASCII char equal to  $i + 98$ 
14:      children $_i$ .populateFrequencyMap(sortedDict, str) ▷ recursive
      call
15:    end if
16:  end for
17: end if
18: return
```

3.2 The 'writeAutoComplete' algorithm

Algorithm 9 writeAutoComplete(**prefix**,*n*,**writer**) **return** *void*

Require: reference to this AutoCompleteTrie, **trie**.

Require: **prefix**, string of lower-case alphabetical characters of size *n*, representing the query.

Require: **writer**, file writer object.

```

1: subTrie  $\leftarrow$  trie.getSubTrie(prefix)
2: allWords  $\leftarrow$  subTrie.getFrequencyMap()  $\triangleright$  map of frequencies and words
3: totalWordFreq  $\leftarrow$  0
4: for i  $\leftarrow$  1 to allWords.pairsCount do  $\triangleright$  count all frequencies
5:   totalWordFreq  $\leftarrow$  totalWordFreq + allWordsi
6: end for
7:
8:    $\triangleright$  If the prefix is a complete word, insert it as one of the results
9: subTrieRoot  $\leftarrow$  subTrie.root
10: entryValue  $\triangleright$  declare list of strings
11: if subTrieRoot.isWord then
12:   entryValue  $\leftarrow$  allWordssubTrieRoot.frequency
13:   if entryValue is null then
14:     entryValue :=  $\triangleright$  initialise
15:   end if
16:   entryValue push empty string
17:   allWordsubTrieRoot.frequency  $\leftarrow$  entryValue
18: end if
19: writer write prefix and ", "
20:
21: entry  $\triangleright$  declare an integer/list of string pair
22: probability  $\triangleright$  declare a float
23: writeCount  $\leftarrow$  0
24: while allWords has entries do
25:   entry  $\leftarrow$  allWords.pop  $\triangleright$  integer/list pair
26:   entryValue  $\leftarrow$  entry.value  $\triangleright$  list of strings
27:   if entryValue has more than 1 element then
28:     sort entryValue in lexicographical order
29:   end if
30:
31:   for i  $\leftarrow$  to entryValue.size do
32:     probability  $\leftarrow$  (entry.key/totalWordFreq)
33:     print prefix and entryValuei and probability in format "word
34:     (probability 0.5)"
35:
36:     writer write prefix and entryValuei and probability in format
37:     "word,0.5,"
38:
39:     wordCount  $\leftarrow$  wordCount + 1
40:     if wordCount is 3 then
41:       print new line
42:       writer write new line
43:       return
44:     end if
45:   end for 8
46: end while
47: print new line
48: writer write new line
49: return

```
