

Assignment 3: Spell Checker Programming Report

s5343798 & s5460484

Algorithms and Data Structures 2024

1 Problem Description

As the title implies, we need to write a spell checker program. This program will need to scan a piece of text, find the words that are misspelt, print them out, and give a count of how many words in the text are wrong.

The program's input consists of two parts, separated by an exclamation mark "!". The first part of the input is a list of correctly spelt words that are separated by a new line. This list of words will be used to create a dictionary of correct words, which are then used for spell-checking the second part of the input that consists of the text that needs to be spell-checked. The text to be checked can be text of any sort, including capital letters and symbols.

2 Problem Analysis

Firstly, we must decide how to represent the data structure of the dictionary in the program. Here, a simple and easy implementation would be to create a list of all the correctly spelt words. However, this solution is crude at best, as it has to traverse the entire dictionary list to check every word in the text. This results in a time complexity of $O(n^2)$.

To reduce this time complexity to a more desirable $O(n)$, $O(\log(n))$, or even better $O(1)$, we need to narrow the search space (the number of words that need to be checked). By using a (Standard) Trie, we can employ a heuristic search due to the nature of the list. A Trie is based on a Linked List data structure where each element refers to the next element in the list. For the entire dictionary, the Trie allows us to define a root element (the first letter of a word) and multiple "children" of that root element (the second, third, etc. letters).

This structure allows us to narrow the check for every word by first checking if the first letter matches, then whether that letter has a "child" letter that matches or not and so forth until it reaches the end of a word. If no "children" letters match, the word is misspelt.

This brings our time complexity down to $O(n \log(n))$, which is much better than the original $O(n^2)$. Although creating the Trie data structure still only takes $O(n)$ time, the time it takes to check the text has greatly been reduced.

After implementing the Trie data structure, the rest of the program is reasonably simple to solve. First, we must add all the words in the first part of the input to the Trie. Thereafter, the program will go line by line and check if each word exists in the Trie; if they do not, that word is printed and the incorrect word counter is updated.

3 Program Design

For the implementation of the Trie data structure, we opted to use two classes; the first represents each node in the list (class name: *TrieNode*), and the second represents the list itself (class name: *StandardTrie*) found in the *standard_trie.py* file. This implementation removes the need to create the empty root node for the Trie manually.

The *insert* function works by checking if the letter of a word is not in the list already and adding it, then moving on till the end of the word is reached and an empty element is added.

The *term_exists* function works similarly to the *insert* function, but instead of adding a letter, it returns false. At the end of the word, if there aren't any empty elements, it returns false. This avoids words like "ap" being seen as current even though those letters are a part of the word "apple."

In the file *speller_trie.py*, an object `allowed_words` is created to which the words in the first part of the input are added. Thereafter, we use the *split* function to split the words on each line into a list following the splitter variable schema. The program then iterates over each word and checks if they are in the dictionary, if not, it prints the word and updates the `unknown_word_count` variable.

Finally, once the entire text has been checked, the program prints how many words were misspelt.

4 Evaluation of the Program

For testing, we used the given example inputs provided to us, as well as the test cases provided by Themis. For performance: In a small input text, the program was rapid, and in a large input text (for example, the English dictionary and a novel), the program was very efficient at processing the text quickly. We would imagine that in the best case, the program takes a few milliseconds, and in the worst case, a few seconds.

Our program passed all test cases on Themis.

5 Extension of the Program

Having only an insert function and `term_exists` function limits the usability of our custom standard trie class. Hence, we created an additional remove function: This takes in one string argument - `term` - and returns and removes the term from the trie.

Solving this problem required splitting the function into two processes: Iterating down the trie to reach the ending of the given term, and then working backwards to remove necessary characters. The reason we went with this approach was that in a standard trie, there may be words with similar roots, e.g. `app`, `apple`, `applet`. If we want to remove `"apple"` then, we would only want to remove up to `"app-"`, to allow for `"app"` and `"applet"` to still exist on the trie. Likewise, if we want to remove `"app"`, we should only remove the indicator (ending) that `"app"` exists on the trie, so that `"apple"` and `"applet"` still exist on the trie.

For the design: First, the program iterates down the trie, building a reference list for future backtracking. Next, the program removes the ending indicator (`None` key) for the given term and then begins iterating back up. In this loop, each key is removed until either the programme finds a child with an ending indicator, in which the function ends, or the loop ends, ending the function too.

For testing, there were no formal cases we used - we designed simple and short inputs for quick, iterative testing. So, it is difficult to evaluate the efficiency of this program. However, this program has a time complexity of $O(n)$, so it is somewhat efficient.

6 Process Description

We worked on the program together through an online call, bouncing back and forth in designing and testing. The work on the report was divided in two, with one handling the first half (up to 3) and the other the second half (up to 7). The distribution of labor was equal.

Regarding difficulty, there was not too much frustration in the assignment. Given that we had some pseudocode from the Reader, it was just a matter of implementation, although we did need to first figure out how to insert terms into the trie. After that though, the work was straightforward.

That being said, there was a small problem with interpreters: One member used Python 3.11, while the other used 3.9. While 3.11 had no issues with the given code, 3.9 had some difficulty, so this required the member to update their interpreter.

By the end of this assignment, we became more knowledgeable about data structures in how to create a custom one from scratch and how to iterate through a data structure. We also learnt how efficient tries can be when implemented properly. And finally, we learnt how to interpret pseudocode effectively.

7 Conclusions

This program efficiently solves the given problem. Even with a large input, the program will handle it in just a few seconds. Regarding whether the program is optimal, perhaps there is some room for improvement. Insertion has a time complexity of $O(n)$ since it iterates n number of times until the program can end, and the same is the case for searching for the term. However, modifying the program to aim for a time complexity of $O(1)$ or

$O(\log(n))$ may require redesigning the structure to allow a search similar to a binary search. In the end though, given the context of the problem - checking if a term exists in a trie - this may be irrelevant scope-wise.

A Program Text

A.1 standard_trie.py

This program text also includes the additional remove function.

```
1  """
2  File: standard_trie.py
3  Authors:
4      Marcus Persson (m.h.o.persson@student.rug.nl),
5      Marinus van den Ende (m.van.den.ende.1@student.rug.nl)
6
7  Description:
8      A custom class for standard tries.
9  """
10
11
12  class TrieNode:
13      """
14      Represents a standard trie node.
15      """
16
17      def __init__(self):
18          self.children: dict[str | None, TrieNode] = dict()
19
20
21  class StandardTrie:
22      """
23      Represents a standard trie.
24      """
25
26      def __init__(self):
27          self.root = TrieNode()
28
29      def insert(self, term: str) -> None:
30          """
31          Inserts a term into the trie.
32          :param term: Term to insert in the trie.
33          """
34          trie_depth = self.root
35
36          for char in term:
37              if not trie_depth.children.get(char):
38                  # If the character is not present at the current trie depth,
39                  # we create a node
40                  trie_depth.children[char] = TrieNode()
41                  # Moving downwards
42                  trie_depth = trie_depth.children[char]
43
44          # Making sure there is an ending None node
45          trie_depth.children[None] = TrieNode()
46
47      def remove(self, term: str) -> str:
48          """
49          Returns and removes a term from the trie.
50          :param term: Term to remove from the trie.
51          :return: The removed term, if it exists.
52          """
```

```

53     if not self.term_exists(term):
54         return "No such term exists in the trie."
55     trie_depth = self.root
56     parents = []
57
58     # Iterating down the trie.
59     for char in term:
60         parents.insert(0, trie_depth)
61         trie_depth = trie_depth.children[char]
62
63     trie_depth.children.pop(None)
64
65     reference = [*term]
66     # Iterating back up
67     for parent in parents:
68         if trie_depth.children.get(None):
69             # If the term iterates back up to a pre-existing term,
70             # end the function
71             return term
72         trie_depth = parent
73         trie_depth.children.pop(reference[-1])
74         reference.pop(-1)
75     return term
76
77 def term_exists(self, term: str) -> bool:
78     """
79     Checks if a term exists in the trie.
80     :param term: Term to check in the trie.
81     :return: True if the term exists in the trie, otherwise False.
82     """
83     trie_depth = self.root
84
85     for char in term:
86         if not trie_depth.children.get(char):
87             # If the character does not exist at that depth,
88             # return False
89             return False
90         # Moving downwards
91         trie_depth = trie_depth.children[char]
92
93     # Validate if the character is the ending
94     return True if trie_depth.children.get(None) else False

```

A.2 speller_trie.py

```

1  """
2  File: speller_trie.py
3  Authors:
4      Marcus Persson (m.h.o.persson@student.rug.nl),
5      Marinus van den Ende (m.van.den.ende.1@student.rug.nl)
6
7  Description:
8      A spell checker, created using a standard trie.
9  """
10 from standard_trie import *
11
12 # Step 1: Create the dictionary
13 allowed_words = StandardTrie()
14
15 word = input()

```

```

16 while word != "!":
17     allowed_words.insert(word)
18     word = input()
19
20 # Step 2: Validate if words exist in a given input.
21 # Handle non-alphabetical characters.
22 splitters = [".", ",", "\'", " ", "/", "-", "\"", ":", ";", "?",
23             "!", "(", ")", "1", "2", "3", "4", "5", "6", "7", "8",
24             "9", "0"]
25
26 unknown_word_count = 0
27
28 word = input()
29 while word != "":
30     for splitter in splitters:
31         word = " ".join(word.lower().split(splitter))
32     for w in word.split():
33         if not allowed_words.term_exists(w):
34             print(w)
35             unknown_word_count += 1
36     word = input()
37
38 # Step 3: Generate output.
39 print("There are", unknown_word_count, "unknown words.")

```