

# Report of CS590 Challenge 5: A Basic Spell Checker

Author: Yunxiu Zhang, Xiaoping Yu, Tianye Zhao

## 1. Introduction

The objective of this challenge is to implement a basic spell checker, which recommend the likeliest known word from a dictionary built from a corpus file, for the input mistyped words.

The term basic means this spell checker only handles the most common mistakes which account for 60%-70% of all spelling mistakes people might make while typing. These most common spelling mistakes occur for the following reasons: Deletion, Replacement, Transposition, Insertion of only one character(2 consecutive characters for transposition).

## 2. Comparison of possible solutions

There are a few possible solutions to implement this basic spell checker.

### 2.1 Solution A: Brute Force

It is not difficult to come up a naïve brute force solution. For each incoming word, we compare it with each word in the dictionary and compute their edit distance. After computing, we get a list of candidates for each incoming word, and pick one word from the candidate as a recommendation as per the requirement. The time complexity of this solution is  $O(n*m)$  (  $n$  is the number of words in dictionary and  $m$  is the length of the incoming word ).

### 2.2 Solution B: Compute all possible recommendations and look up

Because the possible mistakes that this spell checker would handle are countable, and the time complexity of computing their possible original words is  $O(m)$  ( $m$  is the length of the incoming word). And for each possible original word, we look up it in the dictionary, which is a HashMap, to see if it exists. If yes, we put it into the candidate list. The time complexity of look-up is  $O(1)$ , so the overall time complexity of this solution is  $O(m)$ . Since the length of a word is very short,  $O(m)$  is almost  $O(1)$ , which is much better than above solution.

### 3. Solution description

As per above comparison, solution B is much better than solution A, so in this section let's talk about the implementation of solution B.

#### 3.1 Dictionary building

The dictionary can be built as a HashMap, in which the key is a word and the value is the frequency it occurs in corpus text file.

To read all the words from the text file, we need to remove some noises, including space, symbols, and so on. Only the strings with alphabetic characters and apostrophes are valid words. Even for a word that is connected with a hyphen, we split it into two words. We use a regular expression to remove the noises. After that, we can simply read the valid words one by one, count their frequency and save the data into the dictionary.

#### 3.2 Possible recommendations computation

As per the requirement, the most common spelling mistakes occur for 4 reasons: Deletion, Replacement, Transposition, and Insertion. If the word exists in dictionary, then it is not a mistype, we return it as it is. Otherwise, we compute its possible original words according to the possible ways it was mistyped.

##### 3.2.1 Deletion

One character in the string gets deleted incorrectly. If the mistype is caused by deletion, then we can get the possible recommendation by insert one character into the word. If the incoming word is of length  $m$ , then there are  $m+1$  places to insert a character, and for each place we have 26 characters to choose from. This is because only the letters a-z are involved, as per the requirement.

##### 3.2.2 Replacement

If the mistype is because one character in the string is incorrectly replaced by another one, we can get the possible original words by replacing each one character of the word with another 25 characters (a-z excludes itself).

##### 3.2.3 Transposition

If the mistype is caused by swapping one pair of consecutive characters, we can get the possible original words by swapping each pair of consecutive characters again in the incoming word. If the incoming word is of length  $m$ , then we can get  $m-1$  pairs of consecutive characters with first  $m-1$  characters and their subsequent character.

##### 3.2.4 Insertion

If the mistype is caused by inserting one extra character somewhere in the string, then each character in the word is possibly the extra one. we can get the possible original

words by deleting one of the characters in the string. If the incoming word is of length  $m$ , then there are  $m$  possible original words.

### 3.3 Candidates selection

If we got multiple candidates for one incoming word after above step of possible recommendations computation, we need to select one as a result.

First of all, we check their frequency in the dictionary (the key is the word, and the value is the frequency), and return the one occurred most frequently.

Then if we found multiple candidates with same frequency, we sort it alphabetically and return the first one in the list.

In our code, we combine above two steps into one, where we sort the candidates by their frequency and alphabetical order.

## 4. Conclusion

We have implemented the basic spell checker by computing all possible recommendations according to the 4 major reasons of mistakes in  $O(m)$  time ( $m$  is the length of a incoming word): Deletion, Replacement, Transposition, and Insertion, and looking up each of the recommendation words in the dictionary HashMap in  $O(1)$  time. The overall time complexity of this solution is  $O(m)$ . The number  $m$  is very small, so the time complexity of this solution is almost  $O(1)$ .