

# Spell Correction Algorithm Report

## 1. Introduction

The spell correction algorithm plays a vital role in enhancing the quality and accuracy of text data by automatically correcting misspelled words. In this report, we present the methodology used for data cleaning, the approach taken for spell correction, experimental results, an accuracy analysis, and a comprehensive discussion of limitations and potential improvements of the algorithm.

## 2. Methodology for Data Cleaning

To ensure the spell correction algorithm operates on clean and relevant data, a thorough data cleaning process is conducted. The first step involves reading a corpus from the 'corpus.txt' file, which serves as the foundation for creating a dictionary of valid words. This corpus may be obtained from various sources such as books, articles, or other text repositories.

To create a clean and standardized dataset, regular expressions are used to remove any non-alphabetic characters, such as punctuation marks and numerical values, from the corpus. This process ensures that only words are retained in the dataset for further analysis.

Tokenization is a fundamental step in natural language processing (NLP) that involves breaking down a piece of text, such as a sentence or a document, into smaller units called tokens. Tokens are the individual components that make up the text, such as words or punctuation marks. The NLTK (Natural Language Toolkit) library provides various tools and functionalities for working with human language data, including tokenization. Tokenization is then performed using the NLTK library, which divides the cleaned corpus into individual words or tokens. These tokenized words become the basis for generating candidate words during the spell correction process.

Once the corpus is cleaned, the NLTK library is used to perform tokenization. The NLTK provides several tokenization methods, but in this algorithm, the `word_tokenize` function is used. The `word_tokenize` function breaks down the text into individual words, excluding punctuation and other non-word elements. For example, the sentence "Hello, how are you doing?" would be tokenized into the list of words ['Hello', 'how', 'are', 'you', 'doing'].

Finally, by using `FreqDist` from NLTK library frequency distribution object that stores the counts of all unique words in `tokenized_words` is calculated, regardless of the case (due to converting them to lowercase).

## 3. Correction Approach

The spell correction approach implemented in this algorithm relies on the Levenshtein distance metric. For each token in the input text, the algorithm generates a list of candidate words from the dictionary created from the cleaned corpus. The Levenshtein distance is calculated between the input token and each dictionary word. The Levenshtein distance represents the minimum number of edit operations, such as insertion, deletion, Transposition, or Replacement required to transform one word into another.

The algorithm selects the candidate with the smallest Levenshtein distance as the corrected word for the token and generates a list of candidate words which is sorted based on the number of repetition in the dictionary which cause to not be considered words by spelling errors as a best candidate words. If no suitable candidate is found, the original token is retained in the corrected text.

#### **4. Experimental Results**

To evaluate the performance of the spell correction algorithm, several experiments are conducted using different types of errors, such as insertion, deletion, and replacement.

**\*\*Insertion:\*\***

Input: "Affter reeeiving heer visitor"

Ground Truth: "After receiving her visitor"

Detected Errors: 4

Corrected Errors: 3

Accuracy: 75.0%

Output: "iffter receiving her visitor"

**\*\*Deletion:\*\***

Input: " you must returned fro by benefactor"

Ground Truth: " you must returned from by benefactor "

Detected Errors: 6

Corrected Errors: 1

Accuracy: 100%

Output: " you must returned from by benefactor "

**\*\*Replacement:\*\***

Input: "snly to this aim cam we always ftrive independentli"

Ground Truth: "Only to this aim can we always strive independently"

Detected Errors: 9

Corrected Errors: 7

Accuracy: 66.67%

Output: "only to his him can he always strive independently"

**\*\* Transposition:\*\***

Input: "have seen taht was represent"

ground truth: "have seen that was represent"

Detected Errors: 5

Corrected Errors: 2

Accuracy: 60.0

Output: "have been tht was represent"

## **5. Discussion of the Accuracy of the Correction**

The accuracy of the spell correction algorithm varies depending on the complexity and nature of errors present in the input text. In simple cases with minor errors, the algorithm demonstrates reasonable accuracy, achieving around 60% to 100% accuracy. However, in more challenging scenarios where multiple errors and complex patterns are present, the accuracy may decrease.

The corrections are primarily based on the Levenshtein distance, which considers only the minimum number of edits required to transform one word into another. While this approach works well for simple errors, it may lead to incorrect corrections when the input text has multiple valid suggestions with similar edit distances.

The Spell Correction process successfully achieved an overall accuracy of 75% when dealing with four types of errors: Insertion, Deletion, Replacement, and Transposition. Notably, the highest accuracy of 100% was attained for the Deletion errors, indicating that the system effectively

rectified all instances where a single character was incorrectly omitted from the original text. On the other hand, the Transposition errors exhibited a slightly lower accuracy of 60%, indicating that there is still room for improvement in addressing cases where adjacent characters are swapped erroneously.

The accuracy achieved for Replacement errors stood at a commendable 66.66%, reflecting the system's ability to accurately identify and correct instances where one character was incorrectly replaced by another. These results demonstrate the efficacy of the spell correction approach, especially considering the complexity and variations present in natural language data.

Further analysis of the errors revealed that the Insertion type performed satisfactorily, with an accuracy of 75%. This indicates that the system effectively corrected cases where an additional character was mistakenly inserted into the original text. The success of this approach bodes well for the system's overall performance and shows its potential to enhance the quality and accuracy of textual data.

In conclusion, the spell correction system showcased promising results in addressing various types of errors. While the Deletion errors achieved perfection, the overall accuracy remained high for the other types as well. However, the Transposition errors demonstrated room for improvement. Further refinements and enhancements to the spell correction algorithm could potentially lead to even better results and greater accuracy across all error types. The system's versatility and adaptability make it a valuable tool for enhancing the reliability and precision of text processing tasks.

## **6. Limitations and Possible Improvements**

While the spell correction algorithm has demonstrated promising performance in addressing various types of errors, it is essential to acknowledge its limitations and identify potential areas for further improvement:

- ✓ **Limited Edit Distance:** The current algorithm operates within a maximum edit distance of 1, which effectively corrects simple errors. However, more complex errors, such as multiple consecutive insertions or deletions, might require a higher edit distance threshold. Expanding the algorithm to handle larger edit distances would enable it to rectify more extensive changes in the input words and improve correction accuracy for intricate errors.
- ✓ **Dictionary Coverage:** The algorithm's efficacy heavily relies on the size and quality of the dictionary generated from the corpus. While the current implementation offers a reasonable coverage of valid words, a more substantial and diverse corpus could lead to a more comprehensive dictionary. Integrating additional linguistic resources or incorporating

specialized domain-specific vocabularies would further enrich the dictionary, enhancing the system's ability to handle a broader range of words and errors.

- ✓ **Contextual Consideration:** Presently, the spell correction algorithm primarily focuses on individual words and their respective edit distances. While effective for many cases, it lacks contextual awareness of the input text. Introducing context-based approaches, such as utilizing N-grams or advanced language models like neural networks, would enable the system to consider neighboring words and their relationships. This contextual consideration could significantly boost correction accuracy, particularly in scenarios where the meaning of a word depends on its surrounding words.
- ✓ **Evaluation on Diverse Datasets:** While the algorithm has been evaluated on the provided dataset, it would be beneficial to test its performance on more diverse datasets. Evaluating the system on texts from different domains, languages, and writing styles would provide a comprehensive understanding of its generalizability and robustness across various real-world scenarios.

## **7. Conclusion**

In conclusion, the spell correction algorithm has proven to be a valuable tool in addressing various types of errors in natural language text. With its utilization of the Levenshtein distance metric and a well-constructed dictionary, it effectively rectifies spelling mistakes and typos, enhancing the overall accuracy and readability of the text.

However, there is room for improvement to further elevate its performance. By considering a higher edit distance threshold, the algorithm can handle more complex errors that involve multiple insertions, deletions, or replacements. Additionally, expanding the dictionary's coverage by incorporating a more diverse and extensive corpus of words will enable the algorithm to provide more comprehensive and contextually relevant candidate corrections. Addressing these areas of improvement will undoubtedly reinforce the algorithm's utility, making it a powerful asset in refining text data across various applications and domains.