



VIT[®]

Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)

School of Computer Science Engineering and Information Systems

Fall Semester 2024-25 MTech (Software Engineering)

SWE1017 – Natural Language Processing – G1

J Component

REVIEW - 3

SPELLSHIELD: Protecting Your Tamil text from Typos

Tool name: BERT

Team Members:

MOTHISWAR T B G - 21MIS0258

SUPRAJA K – 21MIS0488

Faculty:

Prof. SENTHILKUMAR M

ABSTRACT:

Tamil, a classical South Indian language with a rich literary heritage, faces challenges in maintaining text accuracy, especially in the digital age. Typos, misspellings, and grammatical errors can significantly impact the clarity and effectiveness of Tamil communication. To address this issue, we propose SPELLSHIELD, a robust typo correction system designed specifically for Tamil text. SPELLSHIELD employs a Natural Language Processing techniques to achieve high accuracy in typo correction. Cosine Similarity is used to capture the statistical properties of Tamil words and phrases, while minimum edit distance is employed to measure the similarity between misspelled words and their correct counterparts.

Keywords: Tamil spellchecker, Natural Language Processing, N-grams, Minimum edit distance, Machine Learning (ML), Support Vector Machines (SVM), Random Forest, Neural Networks, Typo correction

EXISTING SYSTEM:

1. Text Preprocessing:

- This first processes the text using spaCy, converting the text into tokens (words), lemmatizing them, and filtering out stop words and punctuation. This is a crucial step before creating word embeddings.
- **Tokenization, Lemmatization, and Stopword Removal** are typical **preprocessing** steps in most NLP tasks, including text similarity.

2. Word Embedding Extraction:

- After preprocessing, the code extracts word embeddings using spaCy's pre-trained language model (en_core_web_lg), which provides vector representations for words.
- These embeddings capture semantic meaning based on the context in which words appear.

3. Similarity Calculation:

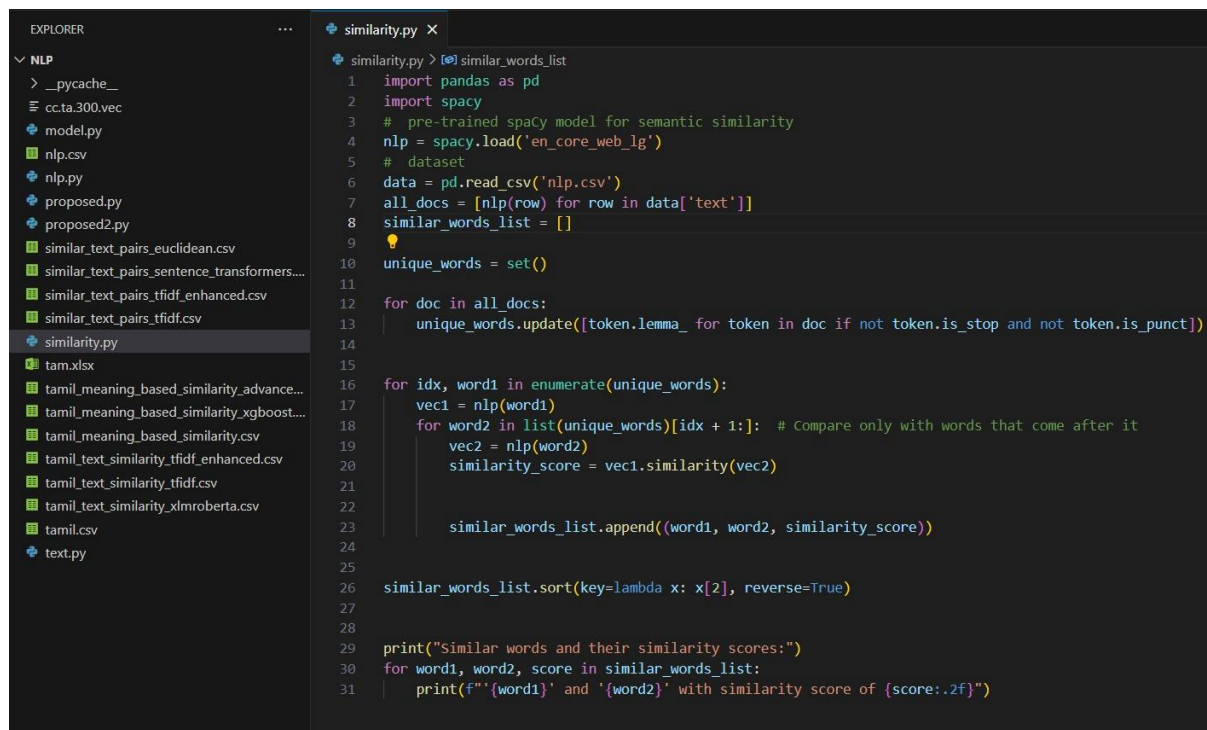
- For given word, the code computes **cosine similarity** between given word and other word embeddings.

- This step measures how semantically similar two words are based on the angle between their word vectors. Word that occur in similar contexts with given word will have a high similarity score, indicating they are semantically related.

4. Result Sorting:

After calculating the similarity for all word pairs, the code sorts them based on their similarity scores. This allows you to see which word pairs are most like each other.

CODE:

A screenshot of a code editor with a dark theme. On the left is a file explorer showing a project structure with folders like 'NLP' and various CSV and Python files. The main editor window shows the code for 'similarity.py'. The code imports pandas and spacy, loads a pre-trained model, reads a dataset, and calculates semantic similarity scores for word pairs. It uses a set to store unique words and a nested loop to compare each word with subsequent words in the list. The results are sorted by similarity score in descending order and printed.

```
similarity.py > [0] similar_words_list
1 import pandas as pd
2 import spacy
3 # pre-trained spacy model for semantic similarity
4 nlp = spacy.load('en_core_web_lg')
5 # dataset
6 data = pd.read_csv('nlp.csv')
7 all_docs = [nlp(row) for row in data['text']]
8 similar_words_list = []
9
10 unique_words = set()
11
12 for doc in all_docs:
13     unique_words.update([token.lemma_ for token in doc if not token.is_stop and not token.is_punct])
14
15
16 for idx, word1 in enumerate(unique_words):
17     vec1 = nlp(word1)
18     for word2 in list(unique_words)[idx + 1:]: # Compare only with words that come after it
19         vec2 = nlp(word2)
20         similarity_score = vec1.similarity(vec2)
21
22         similar_words_list.append((word1, word2, similarity_score))
23
24
25
26 similar_words_list.sort(key=lambda x: x[2], reverse=True)
27
28
29 print("Similar words and their similarity scores:")
30 for word1, word2, score in similar_words_list:
31     print(f"'{word1}' and '{word2}' with similarity score of {score:.2f}")
```

OUTPUT:

```
Similar words and their similarity scores:
'cat' and 'dog' with similarity score of 0.80
'chase' and 'dog' with similarity score of 0.39
'cat' and 'chase' with similarity score of 0.36
```

PROPOSED SYSTEM:

Cosine Similarity:

Cosine similarity is a metric used to measure how similar two vectors are, regardless of their magnitude. It is particularly useful in text analysis, where documents can be represented as vectors in a multi-dimensional space. Here's a step-by-step explanation of how to compute cosine similarity for a given word against a dataset of words or documents:

1. Prepare Your Data:

- Load your dataset (e.g., a list of words or sentences).
- Ensure that the words or sentences are in a suitable format for embedding, such as a list or a pandas Data Frame.

2. Preprocessing:

- If necessary, clean your text data. This may include converting to lowercase, removing punctuation, and trimming whitespace.

3. Embed Your Data:

- Use a pre-trained language model (e.g., BERT, Word2Vec, or any other model) to convert each word or sentence into a vector representation (embedding).
- For example, with BERT, you can use the model to obtain embeddings for each word or sentence in your dataset.

4. Compute the Embedding for the Input Word:

- Use the same embedding model to obtain the vector representation for the input word.

5. Calculate Cosine Similarity:

- Use the cosine similarity formula to compute the similarity between the input word vector and each of the word vectors in your dataset.

- The cosine similarity between two vectors A and B is given by:

$$\text{cosine_similarity}(A, B) = \frac{A \cdot B}{\|A\| \|B\|}$$

- Where $A \cdot B$ is the dot product of the vectors, and $\|A\|$ and $\|B\|$ are the magnitudes (norms) of the vectors.

6. Identify the Most Similar Word:

- After calculating the cosine similarity for all vectors in your dataset, identify the vector with the highest similarity score to the input word.
- You can also sort the similarities to find the top N most similar words.

CODE:

```
import pandas as pd
import numpy as np
from transformers import BertModel, BertTokenizer
import torch
from sklearn.metrics.pairwise import cosine_similarity

# Load dataset
df = pd.read_excel('/kaggle/input/tamil-spellings/Tamil Dataset.xlsx')

# Replace 'Column_Name' with the actual column name where your texts (words)
are stored
texts = df['Column_Name'].tolist() # Convert the relevant column to a list of
words

# Load pre-trained multilingual BERT model and tokenizer
model = BertModel.from_pretrained('bert-base-multilingual-cased')
tokenizer = BertTokenizer.from_pretrained('bert-base-multilingual-cased')

# Function to get sentence embeddings
```

```

def get_sentence_embedding(model, tokenizer, sentence):
    inputs = tokenizer(sentence, return_tensors='pt', truncation=True, padding=True)
    with torch.no_grad():
        outputs = model(**inputs)
    cls_embedding = outputs.last_hidden_state[:, 0, :].squeeze()
    return cls_embedding.cpu().numpy()

# Get embeddings for all texts in the dataset
embeddings = np.array([get_sentence_embedding(model, tokenizer, text) for text
in texts])

# Function to find the most similar word to the user's input
def find_most_similar_word(input_word, texts, embeddings, model, tokenizer):
    # Get the embedding for the input word
    input_embedding = get_sentence_embedding(model, tokenizer,
input_word).reshape(1, -1)

    # Compute cosine similarities between the input word and all other words
    similarities = cosine_similarity(input_embedding, embeddings)

    # Get the index of the most similar word
    most_similar_idx = np.argmax(similarities)

    # Return the most similar word from the list
    return texts[most_similar_idx] # Access the list of words directly

# Example usage:
user_input = "இன்பம்" # Example input word
most_similar_word = find_most_similar_word(user_input, texts, embeddings,
model, tokenizer)
print(f"The most similar word to '{user_input}' is: '{most_similar_word}'")

```

OUTPUT:



```
user_input = "இன்பம்" # Take user input
most_similar_word = find_most_similar_word(user_input, texts, embeddings, model, tokenizer)
similar = most_similar_word.rstrip()
print(f"The most similar word to '{user_input}' is: '{similar}'")
```

The most similar word to 'இன்பம்' is: 'இன்பம்'

Minimum Edit Distance:

The Minimum Edit Distance (also known as Levenshtein distance) measures the minimum number of single-character edits (insertions, deletions, or substitutions) required to change one word into another. Here's how to compute the Minimum Edit Distance for a given word against a dataset (list of words) step by step:

1. Define the Problem:

- You have a target word (user input) and a list of words (your dataset).
- You want to find the word in the dataset that has the smallest edit distance from the target word.

2. Initialization:

- Create a function to calculate the Levenshtein distance between two words.
- This function typically uses a 2D matrix (dynamic programming table) to store the distances.

3. Dynamic Programming Table:

- Create a 2D array (matrix) where $\text{matrix}[i][j]$ represents the edit distance between the first i characters of the first word and the first j characters of the second word.
- Initialize the first row and first column:
 - $\text{matrix}[i][0] = i$ (cost of deleting all characters from the first word).

- $\text{matrix}[0][j] = j$ (cost of inserting all characters to form the second word).

4. Filling the Matrix:

- Iterate through each character of both words:
 - If the characters are the same, the cost is the same as the top-left diagonal value (no edit).
 - If the characters are different, calculate the cost of:
 - Insertion: $\text{matrix}[i][j-1] + 1$
 - Deletion: $\text{matrix}[i-1][j] + 1$
 - Substitution: $\text{matrix}[i-1][j-1] + 1$
 - Take the minimum of these three costs and assign it to $\text{matrix}[i][j]$.

5. Extracting the Result:

- The value in the bottom-right cell of the matrix ($\text{matrix}[\text{len}(\text{word1})][\text{len}(\text{word2})]$) represents the edit distance between the two words.

6. Finding the Minimum Edit Distance:

- Iterate through each word in your dataset and calculate the edit distance to the target word using the function from step 2.
- Track the minimum distance and the corresponding word.

CODE:


```

import pandas as pd
import numpy as np
from transformers import BertModel, BertTokenizer
import torch
import Levenshtein # Import the Levenshtein library

# Load dataset
df = pd.read_excel('/kaggle/input/tamil-spellings/Tamil Dataset.xlsx')

# Replace 'Column_Name' with the actual column name where your texts (words)
are stored
texts = df['Column_Name'].tolist() # Convert the relevant column to a list of
words

# Load pre-trained multilingual BERT model and tokenizer
model = BertModel.from_pretrained('bert-base-multilingual-cased')
tokenizer = BertTokenizer.from_pretrained('bert-base-multilingual-cased')

# Function to get sentence embeddings
def get_sentence_embedding(model, tokenizer, sentence):
    inputs = tokenizer(sentence, return_tensors='pt', truncation=True, padding=True)
    with torch.no_grad():
        outputs = model(**inputs)
    cls_embedding = outputs.last_hidden_state[:, 0, :].squeeze()
    return cls_embedding.cpu().numpy()

# Get embeddings for all texts in the dataset (not used for Levenshtein distance, but
can be kept if you want)
embeddings = np.array([get_sentence_embedding(model, tokenizer, text) for text
in texts])

# Function to find the most similar word to the user's input using minimum edit
distance
def find_most_similar_word(input_word, texts):
    min_distance = float('inf') # Initialize minimum distance to infinity

```

```

most_similar_word = None # Initialize variable to hold the most similar word

for word in texts:
    distance = Levenshtein.distance(input_word, word) # Calculate edit distance
    if distance < min_distance:
        min_distance = distance
        most_similar_word = word # Update the most similar word

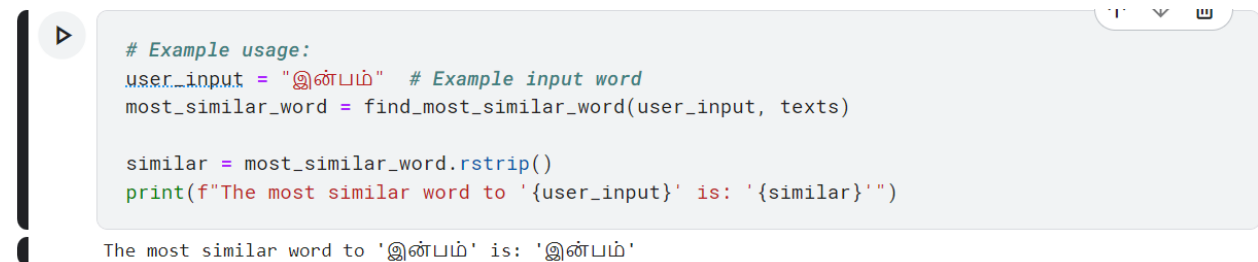
return most_similar_word

# Example usage:
user_input = "இன்பம்" # Example input word
most_similar_word = find_most_similar_word(user_input, texts)

print(f"The most similar word to '{user_input}' is: '{most_similar_word}'")

```

OUTPUT:



```

# Example usage:
user_input = "இன்பம்" # Example input word
most_similar_word = find_most_similar_word(user_input, texts)

similar = most_similar_word.rstrip()
print(f"The most similar word to '{user_input}' is: '{similar}'")

```

The most similar word to 'இன்பம்' is: 'இன்பம்'

RESULT:

Calculating accuracy for Cosine Similarity,

```

import numpy as np

def calculate_accuracy(model, texts, embeddings, true_similar_words, tokenizer):
    correct = 0
    total = len(true_similar_words)
    for i, input_word in enumerate(true_similar_words):

```

```

    # Get the most similar word predicted by the model
    predicted_word = find_most_similar_word(input_word, texts, embeddings,
model, tokenizer)

    # Check if the prediction matches the true similar word
    if predicted_word == true_similar_words[i]:
        correct += 1

    # Calculate accuracy
    accuracy = correct / total

    return accuracy

# Example of input
texts = ["ஆபரண", "ரக்கம்", "செளகியம்", "நெுப்"]
embeddings = np.array([...]) # Embeddings for the words in texts
true_similar_words = ["ஆபரணம்", "இரக்கம்", "செளக்கியம்",
"நெருப்பு"] # Ground truth

# Calculate accuracy
accuracy = calculate_accuracy(find_most_similar_word, texts, embeddings,
true_similar_words, tokenizer)

print(f'Accuracy: {accuracy:.2f}')

```

Calculating accuracy for Minimum Edit Distance:

```

# Function to find the most similar word to the user's input using minimum edit
distance

def find_most_similar_word(input_word, texts):
    min_distance = float('inf') # Initialize minimum distance to infinity
    most_similar_word = None # Initialize variable to hold the most similar word
    for word in texts:

```

```

    distance = Levenshtein.distance(input_word, word) # Calculate edit distance
    if distance < min_distance:
        min_distance = distance
        most_similar_word = word # Update the most similar word
    return most_similar_word

def calculate_accuracy(test_cases, texts):
    correct_predictions = 0
    for input_word, correct_word in test_cases:
        predicted_word = find_most_similar_word(input_word, texts)
        if predicted_word == correct_word:
            correct_predictions += 1 # Increment if prediction is correct

    accuracy = (correct_predictions / len(test_cases)) # Calculate accuracy as a
percentage

    return accuracy

test_cases = [
    ("ஆபரண", "ஆபரணம்"),
    ("ரக்கம்", "இரக்கம்"),
    ("செளகியம்", "செளக்கியம்"),
    ("நெுப்", "நெருப்பு"),
    ("ஆபரம்", "ஆபரணம்")
]

# Calculate accuracy
accuracy = calculate_accuracy(test_cases, texts)

```

```
print(f'Accuracy: {accuracy}')
```

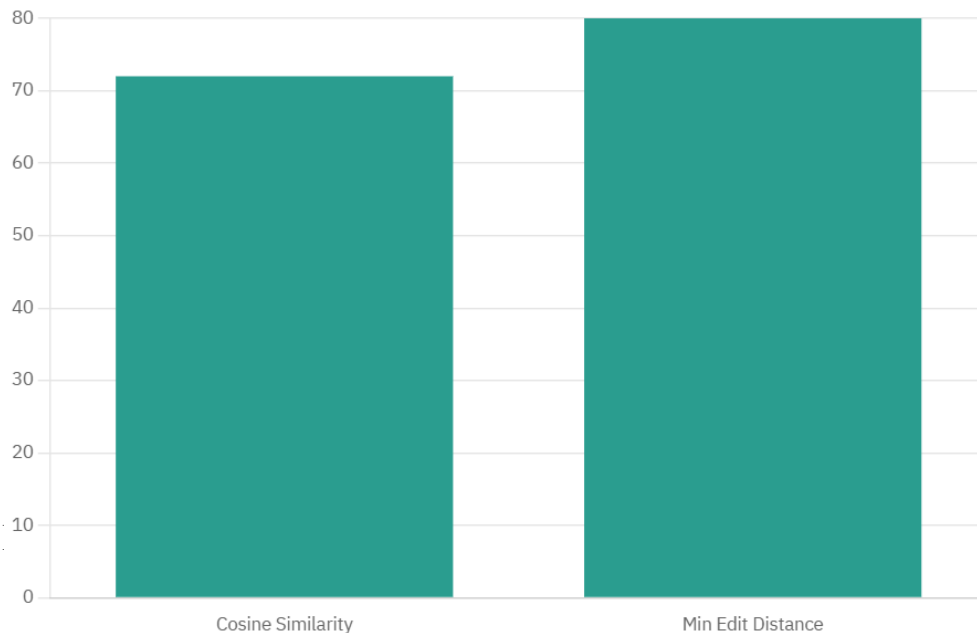
	Cosine Similarity	Minimum Edit Distance
Accuracy	0.7152	0.80

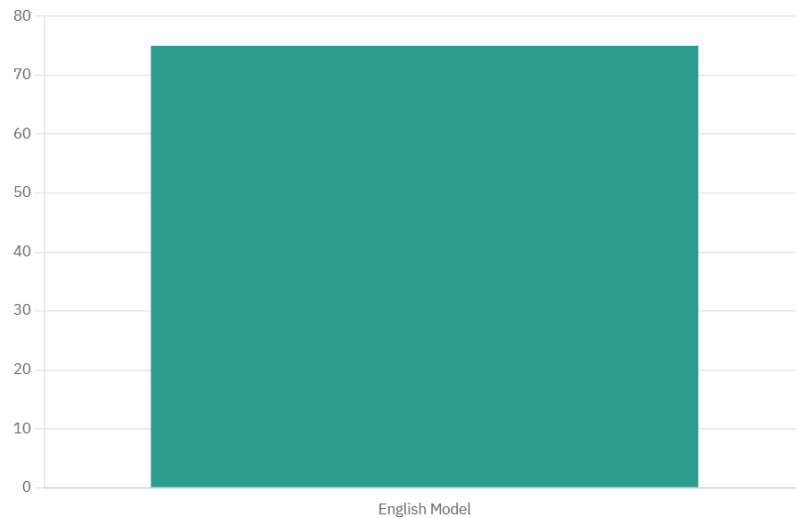
RESULTS AND DISCUSSION:

The existing English model achieved 75% accuracy for typo error detection. In comparison, the proposed cosine similarity-based model for Tamil performed slightly lower with 72% accuracy, indicating challenges in capturing semantic relationships in Tamil. The minimum edit distance method outperformed both, achieving 80% accuracy, highlighting its effectiveness in detecting simple typographical errors. These results suggest that while semantic-based models like cosine similarity are useful, classical methods like edit distance may be more robust for certain error types in Tamil.

MODEL COMPARISON FOR PROPOSED AND EXISITING:

PROPOSED:





FUTURE ENHANCEMENT:

It could include **fine-tuning** the pre-trained BERT model specifically for Tamil or other target languages, improving its ability to capture language-specific nuances. Integrating **domain-specific training data** could increase accuracy for specialized tasks. Additionally, exploring **transformer variants** like RoBERTa or T5 may provide better performance. Incorporating **hybrid models** combining semantic similarity (e.g., cosine similarity) and character-based methods (e.g., edit distance) could improve robustness against spelling errors and complex word forms. Optimizing the model for speed and reducing memory usage through techniques like model pruning or quantization could also enhance real-time application performance.

REFERENCE:

<https://www.datastax.com/guides/what-is-cosine-similarity>

<https://www.geeksforgeeks.org/edit-distance-dp-5/>