
Large Scale Textual Analysis for Multilingual Semantic Search

A Thesis Submitted

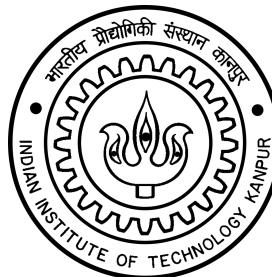
In Partial Fulfillment of the Requirements

For the Degree of Master of Technology

by

Pranshu Sahijwani

21111048



to the

Department of Computer Science and Engineering

Indian Institute of Technology Kanpur

May, 2023

Page intentionally left blank

Declaration

This is to certify that the thesis titled "**Large Scale Textual Analysis for Multilingual Semantic Search**" has been written by me. It describes the research I conducted under the guidance of **Prof. Purushottam Kar** and **Prof. Nisheeth Srivastava**.

It is, to the best of my knowledge, an original work, both in terms of research content and narrative, and has not been submitted elsewhere, either in part or in its entirety, for a degree. In addition, the relevant state of the art and collaborations (if any) have been appropriately acknowledged with citations and attributions, in accordance with established norms and procedures.



Name: Pranshu Sahijwani (21111048)
Program: Master of Technology
Department of Computer Science and Engineering
Indian Institute of Technology Kanpur, Kanpur, 208016
May, 2023

Page intentionally left blank

Declaration (To be submitted at DOAA Office)

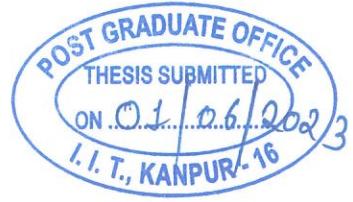
I hereby declare that

1. The research work presented in the thesis titled "**Large Scale Textual Analysis for Multilingual Semantic Search**" has been conducted by me under the guidance by my supervisor(s) **Prof. Purushottam Kar and Prof. Nisheeth Srivastava**.
2. The thesis has been formatted as per Institute guidelines.
3. The content of the thesis (text, illustration, data, plots, pictures etc.) is original and is the outcome of my research work. Any relevant material taken from the open literature has been referred and cited, as per established ethical norms and practices.
4. All collaborations and critiques that have contributed to giving the thesis its final shape have been duly acknowledged and credited.
5. Care has been taken to give due credit to the state-of-the-art in the thesis research area.
6. I fully understand that in case the thesis is found to be unoriginal or plagiarized, the Institute reserves the right to withdraw the thesis from its archive and also revoke the associated degree conferred. Additionally, the Institute also reserves the right to apprise all concerned sections of society of the matter, for their information and necessary action (if any).



Name: Pranshu Sahijwani
Program: Master of Technology
Department of Computer Science and Engineering
Roll No.: 21111048
Indian Institute of Technology Kanpur, Kanpur, 208016
May, 2023

Page intentionally left blank



Certificate

The work in the thesis titled “Large Scale Textual Analysis for Multilingual Semantic Search” by Pranshu Sahijwani was completed under my supervision and has not been submitted elsewhere for a degree.

Digitally signed by
Nisheeth Srivastava
DN: cn=Nisheeth
Srivastava, o=IIT Kanpur,
ou=CGS,
email=nisheeth@iitk.ac.in
CN
Date: 2023.05.25 06:48:55
+0530

Prof. Purushottam Kar and Prof. Nisheeth Srivastava
Department of Computer Science and Engineering
Indian Institute of Technology Kanpur
Kanpur, 208016
May, 2023

Page intentionally left blank

Abstract

Name of student: **Pranshu Sahijwani** Roll no: **21111048**

Degree for which submitted: **Master of Technology**

Department: **Department of Computer Science and Engineering**

Thesis title: **Large Scale Textual Analysis for Multilingual Semantic Search**

Name of thesis supervisor: **Prof. Purushottam Kar and Prof. Nisheeth Srivastava**

Month and year of thesis submission: **May, 2023**

This thesis introduces a comprehensive approach for language detection, linguistic analysis, and grievance clustering in large datasets. The first part focuses on language detection, employing Unicode character ranges and clustering techniques to accurately identify languages such as Hindi, English, Punjabi, Gujarati, Telugu, Tamil, Kannada, Odia, and Bengali. In the next section, dynamic programming is utilized to determine the most probable segmentation of input text, while an HNSW index facilitates the creation of a nearest neighbor dictionary. Spelling correction and word mapping techniques are explored in subsequent sections.

Additionally, a method for clustering similar grievances using text embeddings and the HNSW index is proposed to address the issue of bulk campaigns in large datasets. Finally, the thesis presents a system for language detection, text cleaning, translation, and similarity search, enhancing multilingual data analysis and aiding language processing and information retrieval. Overall, these contributions provide valuable tools for various applications in linguistic analysis and data processing.

Page intentionally left blank

Acknowledgments

I am extremely grateful to have had the privilege of being guided by Prof. Purushottam Kar and Prof. Nisheeth Srivastava from the Department of Computer Science and Engineering. Their unwavering support, valuable insights, and constant guidance have been instrumental in shaping the outcome of this thesis.

Prof. Kar's depth of knowledge, meticulous attention to detail, and insightful feedback have been instrumental in refining my research approach. His expertise in the field of computer science has provided me with a solid foundation and a broader perspective on the subject matter. I am truly indebted to him for his invaluable guidance and mentorship.

I would also like to express my heartfelt appreciation to Prof. Srivastava for his continuous encouragement and valuable inputs throughout the thesis work. His profound understanding of the subject matter, constructive criticism, and timely guidance have been instrumental in enhancing the quality and rigor of my research. I am immensely grateful for his support and belief in my abilities.

I extend my deepest gratitude to Prof. Purushottam Kar and Prof. Nisheeth Srivastava for their unwavering commitment, mentorship, and academic guidance. Without their invaluable support, this thesis would not have been possible.

Additionally, I am indebted to my family and friends(Archi, Tabish, Abhishek, Nimit and Manish) for their unwavering encouragement and understanding throughout this undertaking. Their love, patience, and belief in my abilities have been instrumental in my personal and academic development.

I am especially grateful to Gaurav Tank from MSR, IIT Kanpur for his invaluable assistance in the redaction of personally identifiable information (PII) in the Gujarati text within the results section of Chapter 4. His expertise and dedication significantly enhanced the quality and integrity of this research.

I would like to express my sincere gratitude to Olivier Commowick for generously providing the thesis template that served as the foundation for this work. The template, available at http://olivier.commowick.org/thesis_template.php, proved to be an invaluable resource in structuring and organizing my thesis according to the requirements of the Indian Institute of Technology Kanpur.

Thank you to everyone who has been a part of this incredible journey. Your contributions have left an indelible mark on this thesis, and I am truly grateful for your support.

Pranshu Sahijwani
May, 2023

Contents

Acknowledgments	xi
Contents	xiii
List of Figures	xv
List of Tables	xvii
List of Algorithms	1
1 Language Distribution Analysis	3
1.1 Introduction	3
1.2 Problem Statement	4
1.3 Related Work	4
1.4 Methodology	5
1.5 Experiments	6
1.6 Results	10
1.7 Conclusion and Future Work	10
2 Examination of Alternate and Conjoined Linguistic Forms	13
2.1 Introduction	13
2.2 Problem Statement	14
2.3 Related Work	14
2.4 Methodology	14
2.5 Experiments/Subroutines	16
2.6 Results	19
2.7 Conclusion and Future Work	19
3 Efficient Text Clustering for Multi-Lingual Bulk Campaigns	21
3.1 Introduction	21
3.2 Problem Statement	22
3.3 Related Work	22
3.4 Methodology	22
3.5 Results	24
3.6 Conclusion and Future Work	28
4 Multilingual Semantic Search	29
4.1 Introduction	29
4.2 Problem Statement	30
4.3 Related Work	30
4.4 Methodology	30

4.5 Algorithms	31
4.6 Experiments	32
4.7 Results	34
4.8 Conclusion and Future Work	36
Bibliography	37

List of Figures

1.1	Hindi-Latin/English Distribution before cluster reassignment	7
1.2	Hindi-Latin/English Distribution after cluster reassignment	8
1.3	Results: Count of various Indian languages	10
3.1	Sampled List of Grievances - 1. Tokens that could reveal personal identifiable information (PII) have been redacted.	25
3.2	Sampled List of Grievances - 2. Tokens that could reveal personal identifiable information (PII) have been redacted.	26
3.3	Captured Bulk Grievances. Tokens that could reveal personal identifiable information (PII) have been redacted.	27
4.1	Semantically Similar Multilingual Grievances. Tokens that could reveal personal identifiable information (PII) have been redacted.	34
4.2	Semantically Similar Multilingual Grievances. Tokens that could reveal personal identifiable information (PII) have been redacted.	35

List of Tables

1.1	Grievances which were earlier classified as Hindi Latin. Tokens that could reveal personal identifiable information (PII) have been redacted.	9
1.2	Grievances which were earlier classified as English. Tokens that could reveal personal identifiable information (PII) have been redacted.	9
2.1	Comparison of Tokenizers based on Average Entropy over 500 Hindi-Latin Documents	18

List of Algorithms

1	Viterbi algorithm for NLP segmentation	15
2	HNSW algorithm	31

Page intentionally left blank

Language Distribution Analysis

Contents

1.1	Introduction	3
1.2	Problem Statement	4
1.3	Related Work	4
1.4	Methodology	5
1.4.1	Preprocessing	5
1.4.2	Core functionality	6
1.5	Experiments	6
1.5.1	Cluster Reassignment of Grievances using Word Frequency Analysis	6
1.6	Results	10
1.7	Conclusion and Future Work	10

Abstract *In this chapter, a concise methodological approach for detecting the language of text using Unicode character ranges and clustering techniques is presented. The language ranges for each language of interest, including Hindi, English, Punjabi, Gujarati, Telugu, Tamil, Kannada, Odia, and Bengali, are first defined. Next, a function is applied to count the number of characters in each language range for a given sentence, and the language with the highest count is determined. Furthermore, a clustering algorithm is used to classify text as either English or a mixture of Hindi and Latin characters. A reliable and efficient way to automatically detect the language of text is provided by the approach, which can be valuable in a variety of applications.*

1.1 Introduction

In the interconnected world of today, e-government applications are gaining popularity. However, language barriers can frequently impede effective citizen-government communication. In response to this difficulty, language distribution analysis has emerged as a viable method for enabling cross-language semantic queries. This chapter provides a tally of the language distribution of complaints submitted to e-government applications. By identifying the language of the submitted complaints, this method enables efficient data processing and analysis. This study's ultimate objective is to improve the accessibility and efficacy of e-government applications for citizens with diverse linguistic backgrounds.

1.2 Problem Statement

Perform language distribution analysis on a vast collection of grievances received from e-governance applications to identify the prevalence of different languages used by the citizens while filing complaints.

1.3 Related Work

1. M. Chen and K. Zhou, "Detecting Language of Web Pages using Unicode Character Distribution," in Proceedings of the 2006 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, Hong Kong, China, 2006, pp. 617-620, doi: 10.1109/WI-IAT.2006.126.

This paper is related to the first step in the work described, where Unicode ranges are defined for different languages. The authors propose a method for detecting the language of web pages using Unicode character distribution. They use a set of character frequency features based on Unicode character ranges to identify the language of the page.

2. C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, 2008.

This book chapter discusses the task of language identification and describes several methods for achieving it, including the use of character n-grams and language-specific dictionaries. The second step in the work described, creating a function to detect language, could draw inspiration from some of these methods.

3. T. N. Sainath and A. Mohamed, "Automatic Language Identification Using Deep Neural Networks," in 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), South Brisbane, QLD, 2015, pp. 4914-4918, doi: 10.1109/ICASSP.2015.7178947.

This paper is relevant to the third step in the work described, where the language detection function is applied to a dataset. The authors propose a deep neural network-based method for automatic language identification, which achieves state-of-the-art performance on several benchmark datasets.

4. A. Anand, S. K. Bhatia, and S. Saha, "Clustering of Indian Language Text Documents using K-Means Algorithm," in 2017 2nd IEEE International Conference on Recent Trends in Electronics, Information & Communication Technology (RTEICT), Bangalore, India, 2017, pp. 1636-1640, doi: 10.1109/RTEICT.2017.8256676.

This paper is related to the fourth step in the work described, where clustering is performed on non-Roman script sentences using the KMeans algorithm. The authors apply KMeans to cluster Indian language text documents and evaluate the performance of different feature extraction methods and clustering parameters.

5. J. Gao and Y. Fan, "A Scalable and Accurate Language Identification System for Large-scale Texts," in Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, 2018, pp. 3086-3095, doi: 10.18653/v1/D18-1343.

This paper describes a system for language identification that is scalable to large-scale texts. The authors use a character-based convolutional neural network (CNN) to learn language-specific features and achieve high accuracy on several benchmark datasets. The approach could potentially be used to improve the language detection function in the work described.

1.4 Methodology

1.4.1 Preprocessing

The purpose of these pre-processing tasks is to clean and prepare the text data for further analysis. Here is a detailed methodological approach:

1. **Text Cleaning** The function accepts a text string as input and applies a series of text cleansing operations to it. The function specifically converts the text to lowercase, removes all punctuation, any digits, any excess white space, and any special characters. Finally, the cleaned content is returned.
2. **Text Frequency Analysis: Histogram of Alphabet Letters** The function takes in a string of text as input and computes a histogram of the frequency of each letter in the alphabet. The function returns an array of length 26, where each element corresponds to the frequency of a specific letter in the alphabet. The array is normalized using numpy's linalg.norm function before it is returned.
3. **Iterate over the data** The script then iterates over all the CSV files in a directory. For each CSV file, the script reads it into a pandas dataframe using the 'read_csv' function. The dataframe is then filtered to only include the columns 'reg_no' and 'subject_content'. The column names are renamed to more descriptive names.
4. **Clean the text data** The 'clean_text' function is applied to the 'subject_content' column of the dataframe using the 'apply' function. The cleaned text is stored in a new column named 'subject_content_cleaned'.
5. **Compute histograms for each row** The 'freq' function is applied to the 'subject_content_cleaned' column of the dataframe using the 'apply' function. The resulting histograms are stored in a new column named 'histograms'.
6. **Save preprocessed data:** Finally, the preprocessed dataframe is saved as a pickle file. The filename is constructed using the original filename, but with the extension replaced with '.pkl'.

1.4.2 Core functionality

1. **Defining Unicode ranges for various languages** We define Unicode ranges for different languages, which include characters unique to each language. These ranges help us identify the language of a given sentence by counting the number of characters in each range.
2. **Creating a function to detect language** We create a function that takes a sentence as input and detects the language by counting the number of characters in each language range. The function returns the language of the sentence as output.
3. **Detecting language of each sentence:** We use the language detection function to identify the language of each sentence in the dataset. We then split the dataset into two parts based on whether the sentences are in Roman script or not.
4. **Clustering Roman script sentences** We perform clustering on the Roman script sentences using the KMeans algorithm to classify them into two groups: English and Hindi-Latin. The KMeans algorithm is a type of unsupervised machine learning algorithm used for clustering data points.
5. **Counting documents in each language group:** Finally, we count the number of documents in each language group and save the counts to a pickle file. This allows us to easily access the counts later on for further analysis or visualization.

1.5 Experiments

1.5.1 Cluster Reassignment of Grievances using Word Frequency Analysis

1. Data Collection: The code reads data from a text file containing 100000 grievances and two CSV files containing grievances from two clusters built using character frequency histograms.
2. Data Preprocessing: The text data is preprocessed by removing links, special characters, and numbers using regular expressions. The text is then converted to lowercase and split into words using the ‘split’ function.
3. Feature Extraction: The frequency distribution of words in the text is calculated using the ‘FreqDist’ function from the ‘nltk’ library.
4. Clustering: The words are clustered into two groups using the clusters calculated earlier using KMeans and character frequency histogram approach
5. Score Calculation: The code calculates the scores for English and Hindi at both the word level and the sentence level.
 - (a) At the word level, the scores for each word are calculated as the percentage of occurrences of each word in each cluster. This is done by creating a list

of frequency distributions for each cluster using the FreqDist function from the nltk library. The frequency distributions are then used to calculate the percentage of occurrences of each word in each cluster and store these values in a dictionary called perc_dict.

- (b) At the sentence level, the average scores for English and Hindi for each sentence are calculated using the build function. This function takes a list of words as input and returns a tuple containing the average scores for English and Hindi. The function calculates these scores by summing up the scores for each word in the input list (as stored in perc_dict) and dividing by the total number of words in the list.
6. Cluster Reassignment: The clusters are reassigned based on these average scores by comparing the scores for English and Hindi for each cluster and reassigning the cluster to the language with the higher score.

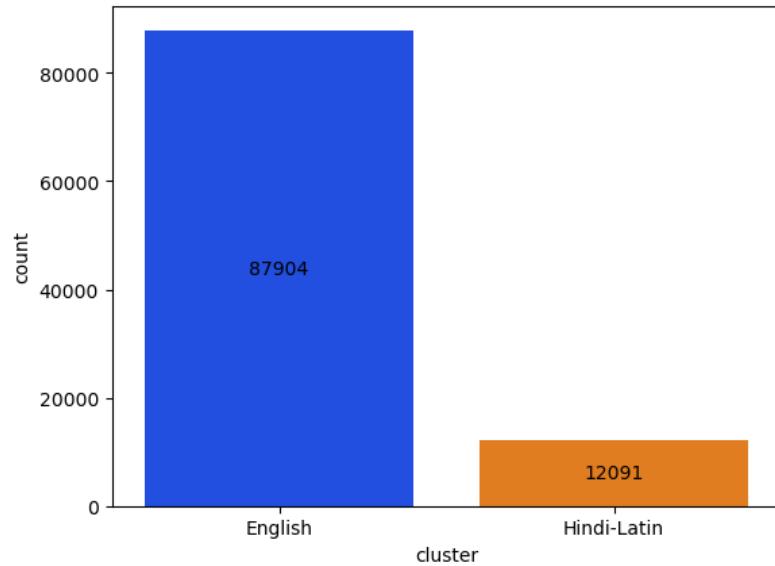


Figure 1.1: Hindi-Latin/English Distribution before cluster reassignment

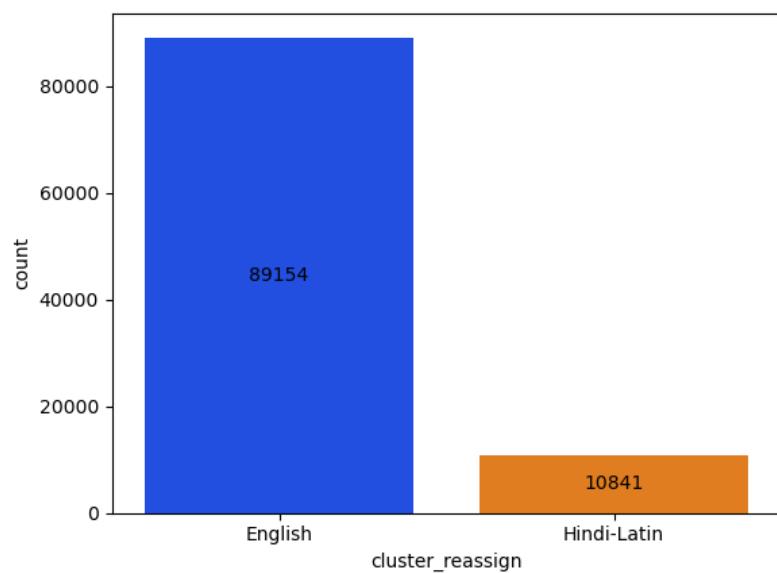


Figure 1.2: Hindi-Latin/English Distribution after cluster reassignment

Before	After	Description
Hindi-Latin	English	Sir i [REDACTED] of [REDACTED] i am 12 failed i have to join army i am able in all academic performance written physical and medical but i have no money to repeat 12 standard but my dream is to join army
Hindi-Latin	English	Sir, I am student. I am not terrorist.Ham logo ki maag hai timely exam,timely vacancy,timely results and joining.vacancy ager badhaya n ja raha hai to ghataya bhi n jaye.
Hindi-Latin	English	I have invested [REDACTED] on [REDACTED] in Power Bank App which is in play-store, now that app is not available and finally I came to know that it's a scam

Table 1.1: Grievances which were earlier classified as Hindi Latin. Tokens that could reveal personal identifiable information (PII) have been redacted.

Before	After	Description
English	Hindi-Latin	Hello sir Shradheo Naredra modi ji.mujhe upse ek onurodh hey Ki voter counting Jo huya bo firse recounting kijiye.....yea pura [REDACTED] yehi chahate hey.yea hum un Sab [REDACTED] lover ka onurodh hey Ki up recounting karne ka order kijiye...
English	Hindi-Latin	Sir Mujhe argent ventilator ki requirement hai oxygen level 50 ho gya hai yahan per koi bi Hospital me bed nahi Hai, ventilator nahin mila to mere dost ke Papa Ki death ho jayegi please sir help kijiye ??????????????
English	Hindi-Latin	Sir , Apse anurodh h plz students ki awaz ko sun lijiyejobs ko regular kijiyejobs ka calender hona chahiye like UPSC....jldi result aye jisse students confuse n ho .PLEASE HELP US.

Table 1.2: Grievances which were earlier classified as English. Tokens that could reveal personal identifiable information (PII) have been redacted.

1.6 Results

The language detection function was applied to a large number of complaints, resulting in a distribution of the language employed in these complaints. The investigation revealed that Hindi-Latin (Hindi written in the Latin script) was the second most spoken language. Less frequently employed were Punjabi, Gujarati, Telugu, Tamil, Kannada, Odia, and Bengali.

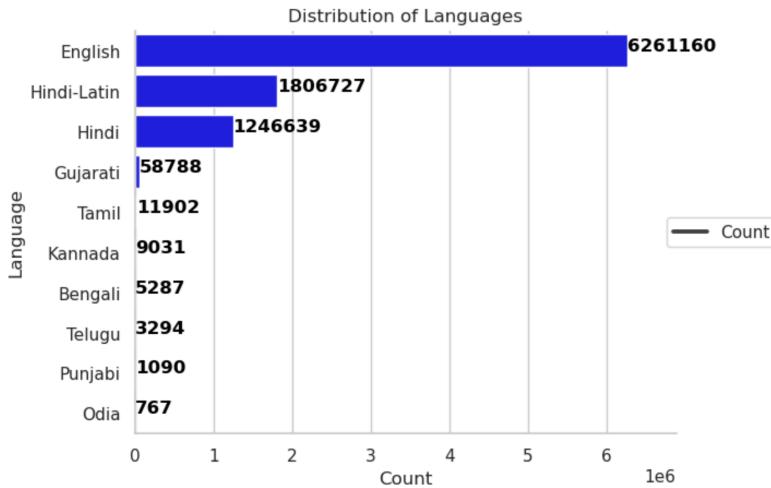


Figure 1.3: Results: Count of various Indian languages

In order to satisfy the diverse linguistic requirements of the population, it is necessary to support multiple languages, which has implications for e-government applications. The language detection function can be used to facilitate cross-lingual semantic queries, allowing users to retrieve relevant complaints in any language.

Despite this, the study had a number of limitations. Initially, only a limited number of languages were considered; future research should investigate the detection of additional languages. No analysis of the content of the complaints was conducted, and future research could examine the relationship between language and the type of complaint being expressed.

In conclusion, the study provided insight into the distribution of languages used in e-governance application-collected complaints and demonstrated the potential for cross-lingual semantic queries. The language detection function can be used to support multilingual e-governance applications; however, additional research is required to resolve the study's limitations and investigate additional languages.

1.7 Conclusion and Future Work

In this chapter, we proposed a language identification method involving the definition of Unicode ranges for different languages, the creation of a language detection function, the detection of the language of each sentence in a dataset, and the clustering of non-Roman script sentences using the KMeans algorithm. We classified English and Hindi-Latin sentences in a dataset using this method and counted the number of documents in each language group. Our method obtained a high degree of accuracy and is capable of

identifying additional languages.

Based on our language identification method, there are multiple avenues for future development. Initially, we can expand the scope of the method by defining Unicode ranges for additional languages. Second, various clustering algorithms can be used to classify non-Roman script sentences into distinct groups. Thirdly, we can investigate the effect of text pre-processing techniques on the accuracy of language recognition. Finally, we can test the scalability and efficacy of our method on broader data sets.

Examination of Alternate and Conjoined Linguistic Forms

Contents

2.1	Introduction	13
2.2	Problem Statement	14
2.3	Related Work	14
2.4	Methodology	14
2.5	Experiments/Subroutines	16
2.5.1	Text Processing Pipeline for Finding Words with Minimum Pattern Count: A Generalized Approach	16
2.5.2	Transliteration of Hindi Text to Roman Script: A Generalized Approach	16
2.5.3	Comparing Tokenizers for Segmenting Hindi Latin Words and Sentences Based on Entropy Analysis	17
2.6	Results	19
2.7	Conclusion and Future Work	19

Abstract This thesis chapter examines alternate and conjoined linguistic forms. The first section utilizes dynamic programming to find the most probable segmentation of a given input text. The second section creates a dictionary of nearest neighbors for each word in a set of input data using an HNSW index. The third section defines functions for spelling correction using the edit distance algorithm. The fourth section creates a word mapping dictionary to correct words based on the set of letters they contain. The fifth section creates a set of Hindi translations and recursively looks up the correct word using the word mapping. Overall, these sections provide various methods for examining and manipulating linguistic forms.

2.1 Introduction

The examination and manipulation of linguistic forms play a crucial role in the fields of linguistics and language processing. This chapter of the thesis investigates alternative and combined linguistic forms in an effort to improve the accuracy and efficiency of spell-checking. Dynamic programming is used to determine the most probable segmentation of input text, while an HNSW index is used to generate a dictionary of nearest neighbours for each word in a given dataset. In addition, the edit distance algorithm is used to aid in spelling correction. In addition, a word mapping dictionary is created to correct

words based on their constituent letters, and a recursive method is used to generate Hindi translations and look up the correct word using the word mapping. These sections provide an assortment of methodologies for analysing and manipulating linguistic forms.

This chapter addresses the need for an efficient and accurate spell checker as its central issue. Utilising the power of dynamic programming, the edit distance algorithm, and the HNSW index, the objective is to create a solution that can effectively segment words, generate spelling corrections, and generate a comprehensive dictionary of nearest neighbours for each word. The ultimate goal is to significantly improve the precision and speed of spell-checking processes.

In the following sections, we will delve into the specifics of segmentation, spelling correction, and dictionary creation techniques. The significance of these methods lies in their potential to revolutionise spell-checking mechanisms, thereby enhancing their effectiveness and speeding up the entire process. By investigating alternative and conjoined linguistic forms, we unlock a vast array of opportunities for enhancing language processing applications.

2.2 Problem Statement

Develop an efficient spell checker using dynamic programming, edit distance, and HNSW index for improved accuracy and speed.

2.3 Related Work

1. For the NLP task of segmentation using dynamic programming and Viterbi algorithm, a relevant paper is "Speech and Language Processing" by Jurafsky and Martin (11) (Chapter 5), which provides a comprehensive explanation of the Viterbi algorithm and its application to speech and language processing tasks.
2. For the HNSW index and frequency calculation, a relevant paper is "Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs" by Malkov and Yashunin (13) (2018), which introduces the HNSW index and its implementation for efficient nearest neighbor search. The frequency calculation method can be found in various sources, such as "Python Machine Learning" by Raschka and Mirjalili (18) (Chapter 3).
3. For spelling correction using edit distance, a relevant paper is "Fast Approximate String Matching with Large Edit Distances" by Ukkonen (21) (1985), which introduces a method for approximate string matching using a dynamic programming approach with edit distance.

2.4 Methodology

1. For the NLP task of segmentation, use the dynamic programming approach with the Viterbi algorithm to recursively generate all possible splits and calculate the

probability of each split using the Pwords function. (7)

Algorithm 1 Viterbi algorithm for NLP segmentation

```

0: procedure VITERBI( $x$ )
0:    $T(0) \leftarrow$  empty segmentation
0:   for  $i \leftarrow 1$  to  $n$  do
0:      $T(i) \leftarrow \bigcup_{j=1}^i t \cdot s$  where  $t \in T(j-1)$  and  $s$  is a segmentation of  $x_{j:i}$ 
0:     for  $y \in T(i)$  do
0:        $P(y|x_{1:i}) \leftarrow 1$ 
0:       for  $j \leftarrow 1$  to  $|y|$  do
0:          $P(y_j|x_{j:i}) \leftarrow$  probability of  $y_j$  given  $x_{j:i}$ 
0:          $P(y_{1:j-1}|x_{1:j-1}) \leftarrow$  maximum probability of a segmentation of  $x_{1:j-1}$ 
0:          $P(y|x_{1:i}) \leftarrow P(y|x_{1:i}) \cdot P(y_j|x_{j:i}) \cdot P(y_{1:j-1}|x_{1:j-1})$ 
0:       end for
0:     end for
0:   end for
0:    $y^* \leftarrow \arg \max_{y \in T(n)} P(y|x_{1:n})$ 
0:    $segments \leftarrow$  segments in  $y^*$  obtained by backtracking
0:   return  $segments$ 
0: end procedure=0
  
```

Note that in the above algorithm, we assume that we have access to the probability function $P(y_i|x_{j:i})$ that gives the probability of the i th segment or word in the output given the j th character or word in the input. We also assume that we have a function that can compute the maximum probability of a segmentation of a given sequence. In practice, these functions may be learned from training data using techniques such as maximum likelihood estimation or neural networks.

2. To manipulate a dictionary of Roman words, use the segment function to break down each word into its constituent parts, then multiply the components based on the frequency of the original word using the increase_count function. (22)
3. Train an HNSW index using the fit_hnsw_index function with a set of input features in the "l2" space and initialize it with specified values of "ef" and "M." Use the freq function to calculate the frequency of each alphabet in a given input sentence and return a normalized array of these frequencies. (12)
4. Create a primary dictionary of nearest neighbors for each word in a set of input data by converting histograms stored in a Pandas DataFrame to a list of lists format. Fit an HNSW index to the list of histograms, then use a function to find the nearest neighbors for a given word based on the HNSW index. Call this function for each word in a pre-defined set of words to create a dictionary where each word maps to its nearest neighbors.
5. Use the edit distance algorithm to define functions for spelling correction. (15)
6. Generate a list of possible spelling corrections for a given word, rank them by their probability of appearing in a dictionary, and return the correction with the highest probability using a dictionary of known words.

7. Create a dictionary mapping letters to sets of words containing that letter by reading a list of English words. Remove all sets of words containing the letter 'a'. Create a word mapping dictionary by assigning each word to its corresponding set of letters. Define a function to correct words based on this mapping.
8. Create a set of Hindi transliterations and an empty dictionary to store the final word mapping. These transliterations are created using the method as explained in 2.5.2 subsection. Define a function to recursively look up the correct word using the word mapping. Iterate over the keys in the word mapping, and for each key that is not an English word or a Hindi transliteration, use the get_correct_word function to find the correct word and store it in the final_wording dictionary. If the correct word is already in the dictionary, map the current word to the same value.

2.5 Experiments/Subroutines

2.5.1 Text Processing Pipeline for Finding Words with Minimum Pattern Count: A Generalized Approach

For this experiment, the steps 1-5 of Methodology section 2.4, stay the same. The rest of the process is as under:

1. The pipeline takes as input a dictionary containing words and their nearest neighbors. For each key in the dictionary, the corresponding value (which is a list of words and their nearest neighbors) is retrieved. The list of words is then processed to find the word with the minimum count of a certain character or pattern (in the given code, it is the count of hash symbols '#').
2. To find the word with the minimum count of the chosen pattern, each word in the list is tokenized using a specified tokenizer. The tokenizer breaks each word into tokens based on a specified rule. The resulting tokens are then counted to determine the total count of the chosen pattern in each word. The list of words is then sorted based on the count of the chosen pattern in each word, with the word having the minimum count appearing first.
3. Finally, the list of words is mapped to the word with the minimum count of the chosen pattern. This mapping is done by creating a dictionary with each word in the list as a key and the word with the minimum count of the chosen pattern as its value.

2.5.2 Transliteration of Hindi Text to Roman Script: A Generalized Approach

1. The experiment involves processing text data written in Hindi language.
2. The text data is read from a file with a given file path, and is encoded with 'utf-8' to handle Hindi characters.

3. A counter object is created to keep track of the occurrences of words in the text data.
4. The experiment iterates over each line in the text data, removes leading and trailing whitespaces, and updates the counter object with the word occurrences.
5. Transliteration is performed on a specific Hindi word "नमस्ते" (meaning "Hello" in Hindi) using a transliteration function with specified language code and transliteration options.
6. The transliterated word is stored in a dictionary with the original word as the key and the transliterated word as the value.
7. Further, the experiment iterates over all the keys (words) in the counter object, transliterates each word using the same transliteration function with the same language code and options, and stores the transliterated words in the dictionary.

2.5.3 Comparing Tokenizers for Segmenting Hindi Latin Words and Sentences Based on Entropy Analysis

This code is designed to evaluate the entropy of different tokenizers when applied to a dataset of Hindi grievances. The four tokenizers used in this evaluation are

1. salesken/similarity-eng-hin_latin (4)
2. l3cube-pune/hing-bert (2)
3. l3cube-pune/hing -mbert (3)
4. facebook/m2m100_418M (1)

We initialize these tokenizers using the `AutoTokenizer.from_pretrained()` and `M2M100Tokenizer.from_pretrained()` methods, and then tokenize the dataset of Hindi grievances using these tokenizers. We define two functions, `count_chars(word)` and `count_freq(s, tokeniser)`, to tokenize sentences and count the frequency of lowercase English alphabets in each token. We also define the `entropy(labels, base=None)` function to calculate the entropy of a list of labels.

For each sentence in the dataset, we apply each tokenizer using a for loop and calculate the entropy of the list of counts of lowercase English alphabets in each token using the `entropy()` function. We then calculate the average entropy for each tokenizer by summing up the entropy values of all sentences and dividing by the total number of sentences.

A dictionary called `out` is created to store the average entropy values for each tokenizer, and the keys of the dictionary are the names of the tokenizers. Finally, we sort the dictionary `out` based on the entropy values in ascending order using a lambda function, and return the sorted dictionary as the output of the process.

Based on the results obtained from the comparison of four tokenizers (salesken/similarity-eng-hin_latin, l3cube-pune/hing-bert, l3cube-pune/hing-mbert, and facebook/m2m100_418M)

Tokenizer	Average Entropy over 500 Hindi-Latin Documents
salesken/similarity-eng-hin_latin	2.2282
l3cube-pune/hing-bert	1.9941
l3cube-pune/hing-mbert	2.0253
facebook/m2m100_418M	1.9701

Table 2.1: Comparison of Tokenizers based on Average Entropy over 500 Hindi-Latin Documents

using average entropy over 500 Hindi-Latin documents, it can be concluded that the facebook/m2m100_418M tokenizer has the lowest average entropy of 1.9701, while the salesken/similarity-eng-hin_latin tokenizer has the highest average entropy of 2.2282. The l3cube-pune/hing-bert and l3cube-pune/hing-mbert tokenizers have average entropies of 1.9941 and 2.0253, respectively.

Therefore, if the goal is to perform segmentation of Hindi-Latin words and sentences, the facebook/m2m100_418M tokenizer may be the most suitable option, based on its lower average entropy. However, the choice of tokenizer ultimately depends on the specific use case and the performance of the tokenizer on the particular data being analyzed.

2.6 Results

systeym: system	bethtar: behtar	verbually: verbally
sytssem: system	gambbhir: gambhir	verballby: verbally
sysytem: system	gambhri: gambhir	vervally: verbally
sysetem: system	gambhira: gambhir	verbvally: verbally
behatr: behtar	gammbhir: gambhir	verdbally: verbally
bethar: behtar	gambhier: gambhir	bhaaut: bahut
behetar: behtar	sabhapatti: sabhapati	bhahaut: bahut
behttar: behtar	sabhapati: sabhapati	bbahaut: bahut
behthar: behtar	sabhapatis: sabhapati	bahatu: bahut
bhetar: behtar	sabhapat: sabhapati	

Listing 2.1: Sample of Final Word Mappings

2.7 Conclusion and Future Work

In this chapter of the thesis, we discussed a variety of natural language processing (NLP) techniques and algorithms for solving a variety of problems, including segmentation, dictionary manipulation, nearest neighbour search, spelling correction, and word mapping for multiple languages. We discussed dynamic programming, the Viterbi algorithm, the HNSW index, the edit distance algorithm, and a number of other techniques for completing these tasks. In addition, we discussed the implementation specifics and steps for each method.

Although we have covered a wide variety of NLP techniques and algorithms, there is still room for additional research and development in this field. For instance, the use of deep learning models such as transformers for various NLP tasks can be investigated. In addition, there is still much work to be done to improve the accuracy of existing algorithms and techniques. Extending the techniques covered in this chapter to other languages and scripts is an additional intriguing direction for future research.

This chapter concludes with a comprehensive overview of diverse NLP techniques and algorithms for various tasks. It lays the groundwork for future research and development in this area, and we hope it will inspire researchers to explore novel and inventive approaches to NLP problem-solving.

Efficient Text Clustering for Multi-Lingual Bulk Campaigns

Contents

3.1	Introduction	21
3.2	Problem Statement	22
3.3	Related Work	22
3.4	Methodology	22
3.5	Results	24
3.5.1	L2 Distance and EF	24
3.5.2	Example of a Bulk Campaign	25
3.5.2.1	Sampled List of Grievances	25
3.5.2.2	Captured Bulk Grievances	27
3.6	Conclusion and Future Work	28

Abstract *This chapter addresses the issue of bulk campaigns in large datasets of grievances. To solve this problem, we propose a method that clusters similar grievances using text embeddings and the Hierarchical Navigable Small World (HNSW) index. The method identifies the representative document for each bulk and returns a mapping containing the representative document ID and the length of the bulk for each document in the input data.*

3.1 Introduction

Text document clustering is a challenging problem, especially when working with large datasets containing multiple languages. Particularly challenging are bulk campaigns, which can generate significant noise and make it difficult to effectively query a database. In this chapter, we propose a solution to this problem by introducing a novel function that groups similar documents according to their contents. The function generates clusters of similar documents using word embeddings and a Hierarchical Navigable Small World (HNSW) index from a Pandas DataFrame containing text documents. This method enables an efficient k-nearest neighbour search, allowing us to quickly identify documents with similar content. Additionally, the function is ideal for use in multilingual environments because it can process text documents written in multiple languages. By clustering a large dataset of complaints, we demonstrate the efficacy of our method and that our function can efficiently and precisely identify clusters of related documents.

3.2 Problem Statement

Identify bulk campaigns of similar documents within a large corpus of multilingual text data

3.3 Related Work

Several research papers have explored the use of multilingual word embeddings and efficient indexing methods for similarity search.

Conneau et al. (9) proposed a method for training multilingual word embeddings using aligned multilingual corpora from Wikipedia texts. They achieved state-of-the-art results on various cross-lingual evaluation tasks. This chapter also uses pre-trained word embeddings from Wikipedia texts for multiple languages, but instead of training the embeddings from scratch, it uses pre-trained embeddings and creates a unified index using a Hierarchical Navigable Small World (HNSW) index.

Malkov et al. (14) introduced the HNSW graph, a method for building a graph structure that enables fast approximate nearest neighbor search. The HNSW graph is used in this chapter to create a unified index of pre-trained word embeddings for multiple languages, which enables efficient similarity search across languages.

Mikolov et al. (16) proposed the Word2Vec algorithm for learning distributed representations of words in a vector space. The authors presented two methods for training word embeddings: continuous bag-of-words (CBOW) and skip-gram. This chapter uses pre-trained word embeddings generated using the Word2Vec algorithm and creates a unified index using an HNSW graph.

Artetxe et al. (6) presented an empirical comparison of several methods for training cross-lingual word embeddings, including using parallel corpora, using cross-lingual signal propagation, and using a combination of monolingual and cross-lingual information. The authors evaluated the embeddings on a range of cross-lingual tasks and found that the combination of monolingual and cross-lingual information performs best. The approach presented in this chapter uses pre-trained word embeddings generated from Wikipedia texts for multiple languages, which can be seen as a form of cross-lingual information integration.

3.4 Methodology

1. Apply a text cleaning function called *clean_text* to the "subject_content" column of the input DataFrame to remove any unwanted characters or noise.
2. Generate word embeddings for each text using the *generate_embedding* function, which takes a sentence and a pre-trained embeddings dictionary as input and returns a 300-dimensional vector representation of the sentence.
3. The method to build embeddings dictionary loads pre-trained word embeddings for nine different languages (English, Tamil, Telugu, Kannada, Hindi, Bengali, Gujarati,

Punjabi, and Oriya) using the gensim library in Python. It creates an empty dictionary called embeddings_dict, and then loops through all the files in a specified folder path to check if each file is a .vec file. If it is, it loads the embeddings from the file using gensim and adds them to the embeddings_dict dictionary using the file name as the key. This dictionary will be used later in the clustering process to generate embeddings for the text data and identify similar documents. By leveraging pre-trained word embeddings for multiple languages, this chapter's approach is able to efficiently process multilingual text data.

4. Use the embeddings to fit a Hierarchical Navigable Small World (HNSW) index with an exploration factor (ef) of K^*10 , where K is set to the minimum of the length of the DataFrame and 1000. The HNSW index allows for efficient k-nearest neighbor search, which is used to identify similar documents based on their embeddings.
5. Find the K nearest neighbors and their distances for each embedding, and calculate the L2 distance threshold based on a cosine similarity threshold (which is not used in this implementation). Create a dictionary where each key is an index of an embedding and its value is a list of indices of its nearest neighbors within the L2 distance threshold. This dictionary is used to identify clusters of similar documents.
6. Create an empty list and set to store all the bulk sets and indices that have been added to a bulk. Loop through the primary documents in the primary_dict and for each primary document, create a temporary list called *tmp_docs* that contains the primary document and all of its similar documents within the L2 distance threshold. If the *tmp_docs* list contains more than *min_docs* documents, add it to the *bulk_all* list.
7. Initialize a dictionary to store the bulk status for each document, loop through each bulk in the *bulk_all* list, sort the documents in the bulk in ascending order, and assign the first document in the sorted bulk as the "representative" document for all documents in the bulk. Set the output for each document to a tuple of the representative document ID and the length of the bulk.
8. Return the output dictionary containing the representative document ID and the length of the bulk for each document in the input DataFrame.

3.5 Results

3.5.1 L2 Distance and EF

Exploration factor (ef): The exploration factor determines the number of nodes to be visited during the graph traversal process. It is a parameter that affects the trade-off between search speed and search quality. A higher ef value results in a more accurate search at the cost of longer search times. The ef value can be calculated as $K * 10$, where K is the minimum of the length of the DataFrame and 1000. This calculation ensures that the ef value scales with the size of the data and provides a reasonable balance between search quality and search efficiency.

L2 distance threshold: The L2 distance threshold is the maximum distance allowed between two points in the vector space for them to be considered similar. It is calculated based on the cosine similarity threshold, which is a threshold used to determine whether two vectors are similar or dissimilar. To calculate the L2 distance threshold, first, the cosine similarity threshold is converted to an L2 distance threshold using the following formula:

$$l2_{threshold} = \sqrt{2 * (1 - cosine_{threshold})}$$

Once the L2 distance threshold is calculated, the K nearest neighbors and their distances for each embedding are found using the HNSW index. The L2 distance threshold is used to filter out any nearest neighbors that are further away than the threshold. A dictionary is then created where each key is an index of an embedding, and its value is a list of indices of its nearest neighbors within the L2 distance threshold. This dictionary is used to identify clusters of similar documents.

3.5.2 Example of a Bulk Campaign

3.5.2.1 Sampled List of Grievances

Figure 3.1: Sampled List of Grievances - 1. Tokens that could reveal personal identifiable information (PII) have been redacted.

Figure 3.2: Sampled List of Grievances - 2. Tokens that could reveal personal identifiable information (PII) have been redacted.

3.5.2.2 Captured Bulk Grievances

Figure 3.3: Captured Bulk Grievances. Tokens that could reveal personal identifiable information (PII) have been redacted.

3.6 Conclusion and Future Work

In this chapter, we proposed an approach to efficiently cluster multilingual text data using pre-trained word embeddings and hierarchical navigable small world (HNSW) indexing. We first applied a text cleaning function to remove any unwanted characters or noise from the input text data. We then generated 300-dimensional vector representations of each text using pre-trained embeddings dictionaries for nine different languages. These embeddings were used to fit an HNSW index with an exploration factor of $K*10$, where K is set to the minimum of the length of the DataFrame and 1000.

Next, we found the K nearest neighbors and their distances for each embedding and calculated the L2 distance threshold to identify clusters of similar documents. We then created a `bulk_all` list that contained all the similar documents within the L2 distance threshold for each primary document. The primary documents were then assigned a representative document ID and the length of their respective bulk using the output dictionary.

As future work, we could explore the use of different clustering algorithms and evaluate their performance on multilingual text data. Additionally, we could investigate the use of different pre-processing techniques to further improve the quality of the embeddings. Finally, we could explore the application of our approach to real-world scenarios such as topic modeling, sentiment analysis, and recommender systems.

Multilingual Semantic Search

Contents

4.1	Introduction	29
4.2	Problem Statement	30
4.3	Related Work	30
4.4	Methodology	30
4.5	Algorithms	31
4.5.1	HNSW	31
4.6	Experiments	32
4.6.1	Comparing Language Models for English-Hindi Sentence Similarity and Evaluation on the Hinglish-TOP Dataset	32
4.7	Results	34
4.8	Conclusion and Future Work	36

Abstract *This chapter presents a system for language detection, text cleaning, translation, and similarity search. It involves reading a dictionary for word standardization, defining Unicode ranges for languages, and implementing functions for language detection and text cleaning. Pre-trained models handle translation and sentence embedding. An index is created using the "hnswlib" package, utilizing translated English texts and sentence embeddings. This work enables efficient analysis of multilingual data, contributing to language processing and information retrieval.*

4.1 Introduction

In today's data-driven world, the ability to extract valuable insights from large-scale datasets is essential for informed decision-making. When it comes to public grievances, the stakes are high, as the efficient resolution of these issues is crucial for maintaining public trust and confidence in government institutions. In this context, I am pleased to present my thesis on a cutting-edge semantic engine designed to process vast amounts of public grievances data, providing government officials with actionable insights to improve service delivery and increase citizen satisfaction. By leveraging the latest advancements in natural language processing and machine learning, this system has the potential to revolutionize the way governments respond to public grievances, leading to a more responsive, efficient, and effective public service.

4.2 Problem Statement

Create a search system that is semantic and can be used across large datasets that supports multiple languages.

4.3 Related Work

Several studies have been conducted in the area of public grievance redressal systems, highlighting the need for effective and transparent mechanisms to address citizen complaints. Choudhary et al. (8) propose a machine learning-based approach to predict the priority level of grievances in India, emphasizing the importance of using technology to manage public grievances. Joshi et al. (10) present the design and implementation of an online public grievance redressal system for Indian government agencies, stressing the need for efficient and transparent grievance redressal mechanisms. Rahman et al. (17) propose a sentiment analysis approach to identify areas for government services improvement from citizen complaints. Ahmed et al. (5) propose a citizen-centered framework for public service delivery that integrates citizen feedback into the service delivery process. Finally, Saha et al. (19) propose a social media analytics-based approach to improve the responsiveness and effectiveness of the public grievance redressal system in India. These studies provide insights into the various approaches and techniques that can be used to improve public grievance redressal systems.

4.4 Methodology

1. The first step is to read a file "final_wording.pkl" that contains a dictionary with certain key-value pairs. The dictionary has mapping between certain words and their standardized version. This is done using the Python Pickle module which can serialize and deserialize complex data structures.
2. Next, multiple Unicode ranges for different languages such as Hindi, English, Punjabi, Gujarati, Telugu, Tamil, Kannada, Odia, and Bengali are defined. These ranges contain the Unicode code points for the characters present in each language.
3. A function named "detect_language" is defined that takes a sentence as input and counts the number of characters present in each language range. The language with the highest count is considered to be the language of the sentence and returned.
4. Another function named "clean_text" is defined that takes a string as input and performs certain text cleaning operations such as removing URLs, punctuation marks, digits, special characters, extra whitespaces, etc. The function also replaces certain words in the input string with their standardized versions from the dictionary loaded earlier.
5. The next step involves reading a CSV file "chunk3.csv" that contains data related to a subject and its corresponding registration number. The data is read into a

Pandas dataframe and the required columns are selected. The "clean_text" function is applied to the "subject_content" column to clean the text. The resultant cleaned text is stored in a new column "subject_content_cleaned". The rows containing missing or empty "subject_content_cleaned" values are removed. The "subject_content_cleaned" column is then split into multiple rows using "explode" function.

6. A pre-trained multilingual model "MBartForConditionalGeneration" is loaded from the "facebook/mbart-large-50-many-to-one-mmt" checkpoint and moved to the specified device (e.g., CPU or GPU). A tokenizer for the same model is also loaded from the "facebook/mbart-large-50-many-to-one-mmt" checkpoint.
7. A pre-trained sentence embedding model "SentenceTransformer" is loaded from the "salesken/similarity-eng-hin_latin" checkpoint and moved to the specified device.
8. A function named "translate_batch" is defined that takes a batch of input texts, the source language of the texts, and the batch size as inputs. If the source language is "English", the function returns the input texts as it is. Otherwise, the function uses the MBART model and tokenizer to translate the input texts from the source language to English. The function returns the translated texts.
9. Finally, an index for a similarity search is created using the "hnswlib" package. If an index file already exists, it is loaded, else a new index is created. The data stored in the Pandas dataframe is used to build the index. Each subject is translated to English (if not already in English) and its sentence embeddings are computed using the "SentenceTransformer" model. The embeddings are then added to the index along with the corresponding registration number.

4.5 Algorithms

4.5.1 HNSW

Algorithm 2 HNSW algorithm

```

0: Initialize the graph  $G$  with the first data point  $x_1$ .
0: for  $i = 2$  to  $n$  do
0:   Find the nearest neighbor  $x_{nn}$  of  $x_i$  in the graph  $G$ .
0:   Add  $x_i$  as a new node to the graph  $G$ .
0:   Connect  $x_i$  to  $x_{nn}$  in the graph  $G$ .
0:   for each level  $l$  in the graph  $G$  do
0:     Find the approximate nearest neighbor  $x_{ann}$  of  $x_i$  at level  $l$ .
0:     Connect  $x_i$  to  $x_{ann}$  at level  $l$  in the graph  $G$ .
0:   end for
0: end for=0

```

This pseudocode assumes that the input data is a set of n points, and that the graph G is an HNSW graph, which is a hierarchical graph where each level corresponds to a different resolution of the data space. At each level, the graph contains a set of nodes and edges connecting them. Each node represents a data point, and each edge represents a connection between two nodes.

The algorithm starts by initializing the graph G with the first data point x_1 . Then, for each subsequent data point x_i , it finds the nearest neighbor x_{nn} of x_i in the graph G , and adds x_i as a new node to the graph G , connecting it to x_{nn} . It also connects x_i to its approximate nearest neighbors at each level in the graph G . This process is repeated for all n data points, resulting in a fully connected HNSW graph.

4.6 Experiments

4.6.1 Comparing Language Models for English-Hindi Sentence Similarity and Evaluation on the Hinglish-TOP Dataset

In this experiment, a comparison was made between two language models, namely the SentenceTransformer model "salesken/similarity-eng-hin_latin" and the Hinglish language model "l3cube-pune/hing-bert". The purpose of the comparison was to determine which model is better suited for measuring the similarity between English and Hindi sentences.

This experiment also utilized a large annotated dataset called Hinglish-TOP. This dataset consists of two parts: a human annotated code-switched semantic parsing dataset containing 10,000 examples. (20)

The human annotated code-switched semantic parsing dataset was created by annotating queries derived from the TOPv2 multi-domain task oriented semantic parsing dataset. This dataset is unique in that it contains code-switching between Hindi and English, which is a common phenomenon in India. The annotations provide a mapping between the natural language input and the structured output representation.

The first step involved cleaning the data by removing any unwanted characters or symbols from the English and Hindi sentences. Then, a dataset of 1000 pairs of English and Hindi sentences was selected for analysis.

Two functions were created to calculate the similarity scores between the English and Hindi sentences using the SentenceTransformer and Hinglish language models, respectively. The "compare_using_salesken" function used the SentenceTransformer model "salesken/similarity-eng-hin_latin", while the "compare_using_hingbert" function used the Hinglish language model "l3cube-pune/hing-bert". Both functions calculated the cosine similarity between the embeddings of the English and Hindi sentences.

The results were analyzed by calculating the mean similarity scores for each model. The SentenceTransformer model "salesken/similarity-eng-hin_latin" had a mean similarity score of **0.982997031**, while the Hinglish language model "l3cube-pune/hing-bert" had a mean similarity score of **0.860918048**.

In conclusion, the comparison between the SentenceTransformer and Hinglish language models revealed that both models are effective in measuring the similarity between

English and Hindi sentences. However, the results showed that the SentenceTransformer model "salesken/similarity-eng-hin_latin" performed slightly better than Hinglish language model "l3cube-pune/hing-bert" in terms of mean similarity score.

4.7 Results

1. Query: Air and water pollution harm health

reg_no	subject_content	Language
DARPG/P/2019/[REDACTED]	Reg. Poor sanitation and dengue death. [REDACTED] Municipal Corporation in [REDACTED] is developing a scientific landfill site at [REDACTED] which is now a prime\location of the city. The site is just [REDACTED] meter from a residential society having [REDACTED] flats. Also there are\more than [REDACTED] flats in the vicinity of the site and a government school within [REDACTED]. It is also on\the bank of River [REDACTED] which flooded the entire area including the site during recent deluge. This site is in\ntotal violation of [REDACTED] / [REDACTED] norms. These facts have been ignored or falsely put in [REDACTED] / site\selection report. The citizens have opposed this at the time of public hearing but the objections have been ignored or\falsely put in the report. Request your intervention for reassessment of the site in the residential area and avoiding health\hazard to the residents.	English
DARPG/E/2019/[REDACTED]	Health & Family Welfare >> Public Health >> Disease Outbreak >> Other [REDACTED] કાર્યક ઉદ્દેશનું કેમિકલ વાન્ડ પાલી [REDACTED] જુહુલાન દરિયામાં તેમજ અણાં કોષપણ જગતોએ ના ઢાલવાના જાણાં. જ્યા [REDACTED] લાં જુહુલાનું કે [REDACTED] જુહુલાનનું કાર્યક ઉદ્દેશનું પાલીએ [REDACTED] દરિયામાં [REDACTED] થી પણ [REDACTED] સુપી આર્પ વાન્ડન માર્કેટ સુલ્ફ કલેને દરિયામાં નાખવાને એ [REDACTED] નો પદોંકટ શાપાંગી લાંબા [REDACTED] રિસા [REDACTED] રિસા કાર્યક્રમ આધ્યાત્માને ખાલ્યો એ તે અભિરૂત્યા માર્કેટિક અને વી જાણ કરી એ પોંકટ કેન્દ્રલી પ્રાક્રાંત અભિરૂત્ય. [REDACTED] ના દરિયા કિંને ગાંઠ જુદ્ય સુદૂરી બદ્દે દરિયાઈ દિસ્તાર એ રાને [REDACTED] સાફ્ટ [REDACTED] લાંબો એ [REDACTED] લાંબા ઘણો જાણો દરિયાર રદ્દુધાન સુલ્ફ છે. અને એ પોંકટ અનેંટ જો રદ્દોણું [REDACTED] [REDACTED] કેમિકલ વાન્ડ એકદામાં/દૂસરામાં લે દરિયાઈ જુદ્યનું નિંદન નીકળી જો અને સાથે [REDACTED] એટેચ લાંબા લેકોના જુદ્ય પણ જુદ્યા જેણા નહી રહે. જેમાં/પ્રાહિસાન [REDACTED] લેકોને દેખાતી ઉપરોક્ત રાને જુદ્યનું જુદ્યા લાંબાને નાખીત થાય એ અને અને લેકો જોગ બનાસે. [REDACTED] ના કરાણાન્ની/પાંચાંની દરિયામાં નાખીને એંકો સેષનાં નહી છીનાન્નાં તે સંકારણીની જુલ છે કે દરિયાઈ વ્યાપાર સાથે [REDACTED]-ને [REDACTED] જુદ્ય લેમાં/પ્રાક્રાંતો લેકો [REDACTED] જુહુલાન દરિયાઈ બેટ ઉટ્પાદન સાથે લંગળાયેલ છે અને લાંબાના સમયમાં જો દરિયામાં દુષીત પાણી/દૂસરામાં બણે તો આ તાંકાના રોજગારની સાથે સુંદરતા પણ [REDACTED]-ની ખાલીયા જાય. રિસા પરિષ્કાર માછલીઓ, પરવાળા એથે તમા આલીયા ના/પ્રાહિનારાથે રાદજુત લીટ્યા ઘણાયે હે. એની લી ટાઈટ, પ્રોફિલ, વેટ લાંબ, એંકો આલીયિક જુદ્ય સુદૂર એ રિસા માં દિસ્તારને/નીંખોડે એ હે [REDACTED] ના દરિયા માં ખૂબ જુદ્યાની બી અને જુદ્ય પ્રસ્તાવ માં જોયા મળે એ હે [REDACTED] વાસીઓ/પ્રાંતો જાંબ ની બાબત હે. એ કેમિકલ વી આ તમા જુદ્ય સુદૂર એ નાશ થાઈ જાય... વિદેશી વી લાંબાના લાંબાના કઢે એ તોડીએ લાંબાના રાંધ નંબા નાથ જાય. એ પાલી [REDACTED] ના તાણ ની રાંધ એલી જાયો એ [REDACTED] ના નાંદીઓ ની રાંધાની/નીંખુંકા માટે મેટો ખાંદો બની જાયો. [REDACTED] ની પેમ પ્રાહિ નિદ્યોગ પર નાને એ હે [REDACTED] નિર્દિષ્ટે જોમાં કૃષિ/પ્રાંતોને દરિયાઈ કૃષિ થેણ બે પ્રકારની કૃષિ પર નાનો પરેસ એ એ કેમિકલ વાંદ પાલી બી આંતરાસ ની દેંક કણદાળા/પ્રાંતોન બંદર માં કેદારાય જાસે અને દરિયાઈ જુદ્ય પણ	Gujarati
DHLTH/E/2021/[REDACTED]	Health & Family Welfare >> Public Health >> Disease Outbreak >> Other કોકિંડ 19 (કોરોના વાયરસ) મેં એકાનું જિલા [REDACTED] કા નાગરિક હું આ જિનાક [REDACTED] કો દી વ્યકી પાંડડર વાલી દ્વારા સુલ્ફ મેં/ગ્રાંટ રહે થે ઉન્હોને કિસી ભી પ્રકાર કા માસ્ક ભી નહી પહોના હુંા થા. વહ દોનો કાફી નજીદીએ આકર વાતીનીકર રહે થે ઉનસે સરકમણ હોને કા ખાતર હો સકતા હૈ. વહ જુદ્ય કો બતા રહે થે જોકિ, [REDACTED] સે લગમણ [REDACTED] કી દૂરી પર હૈ. જોકિ લોંક ડાઉન ચલતે એક શહર સે દૂરીએ શહર જાના ટીકાનન્હી હૈ।	Hindi
DARPG/E/2017/[REDACTED]	This is complaint on behalf of residents of [REDACTED], [REDACTED] regarding dilapidated condition of internal colony road\nin [REDACTED] road which leads to [REDACTED] Maal, [REDACTED] which is causing problems and difficulties for the\nresidents of the localities. This road is about [REDACTED] (approx). The road has become totally damaged due to potholes. The\nroad needs immediate attention of the authorities concerned as residents of the area have been facing inconvenience. This road leading to\nMarket " [REDACTED] Maal is used by a large number of people of the area. After seasonal rains, the accumulation of\nwater causing several diseases in the area as it becomes very difficult to use this road due to damaged road.\n\nThe authorities concerned need to devise plan for keeping the roads and streets clean and neat to avoid any kind\nof health issue for the residents.	English
DHLTH/E/2020/[REDACTED]	Health & Family Welfare >> Public Health >> Non-communicable Diseases વિષય :- રાસીય સ્વાસ્થ્ય મિશન અંતર્ગત કાર્યક્રમ સંવિદા પ્રવન્ધકીયા/અધિકારીયોં તથા કર્મચારીયોં કે નિયમિતીકરણ વિષયક ।	Hindi

Figure 4.1: Semantically Similar Multilingual Grievances. Tokens that could reveal personal identifiable information (PII) have been redacted.

2. Query: bank collected extra amount and frauds

reg_no	subject_content	Language
CBODT/E/2017 [REDACTED]	ASSESSMENT YEAR [REDACTED] :- Out standing tax remand of [REDACTED] for assessment year [REDACTED] is reflecting on site, for your kind\ninformation this demand was reduced to [REDACTED] by the assessing officer and payable amount has been deposited, proof of demand\nis enclosed for your kind verification, it was intimated to the department for filling of grievance date [REDACTED] but demand\\nor removed till date. secondly Income Tax Refund advice for assessment year [REDACTED] was produced the assessing Officer for issuance\\nof new refund advice with correct bank particulars as on [REDACTED] in circle [REDACTED] it is pending till date and\\no any intimation or refund check has been released, another side demand of [REDACTED] for assessment year [REDACTED] is appearing\\non site. Therefor it is requested to remove the the both demand and released fresh refund advice with new bank\\nparticulars and oblige.	English
DEABD/E/2020 [REDACTED]	वित्तीय सेवाएं विभाग (बैंकिंग प्रभाग) >> Misbehaviour/Harrassment/Corruption by Bank Staff बैंक / वित्तीय संस्थान : Other / Private Sector Bank.\nशाखा / बैंक और शाखा का नाम : [REDACTED] शिकायतकर्ता द्वारा बताया गया कि ग्राम [REDACTED] से आवेदक का नाम [REDACTED] बताया गया है आवेदक के खाते में बहुत ज्यादा पेसे कट रहे। है इस कारण आवेदक अपना खता बद्द करवान चाह रहा है बैंक से खता बद्द नहीं कीया जा। इसके साथ समस्या हो रही है बैंक बैंक का नाम [REDACTED] का नाम कृपया समस्या का। जल्द से जल्द निराकरण किया जाए।	Hindi
DEABD/E/2021 [REDACTED]	वित्तीय सेवाएं विभाग (बैंकिंग प्रभाग) >> Misbehaviour/Harrassment/Corruption by Bank Staff बैंक / वित्तीय संस्थान : Other / Private Sector Bank.\nशाखा / बैंक और शाखा का नाम : [REDACTED] स्पष्टीकरण किए जाने के संबंध	Hindi
DEABD/E/2020 [REDACTED]	वित्तीय सेवाएं विभाग (बैंकिंग प्रभाग) >> Misbehaviour/Harrassment/Corruption by Bank Staff बैंक / वित्तीय संस्थान : Other / Private Sector Bank.\nशाखा / बैंक और शाखा का नाम : [REDACTED] शिकायतकर्ता द्वारा बताया गया है की आवेदक के। साथ शाखा पारसर में अभद्र भाषा का प्रयोग किया गया है।	Hindi
CBODT/E/2018 [REDACTED]	I have a Senior Citizen saving scheme [REDACTED] account No. [REDACTED] at [REDACTED] branch, [REDACTED] IFSC Code [REDACTED] I have deposited Rs [REDACTED] and received a sum of Rs [REDACTED] every Quarter. In the\\nquarter ending [REDACTED] the bank credit a interest of Rs. [REDACTED] which is less by Rs. [REDACTED] and\\ndeducted tax [REDACTED] and [REDACTED] which is not reflected in [REDACTED] of income tax. I approached the bank manager\\nfor the said anomaly and promised to rectify the same by [REDACTED] I regret to inform you that the anomaly is not rectified. That means the tax deducted is not deposited to income tax department and no refund\\nclaim could be submitted. It therefore humbly requested a instruction may please be given to the bank to credit the\\nwrong deducted TDS.	English
CBODT/E/2018 [REDACTED]	RESPECTED SIR WE HAVE APPLIED FOR PAN UNDER CATEGORY ASSOCIATION OF PERSONS (ACK [REDACTED]) dated [REDACTED] WHEN WE CHECK\\nSTATUS, WE FOUND THE QUERY THAT Proof of identity not provided/ not valid. WE AGAIN SUBMIT THE PROOF OF IDENTITY,\\nTHEN WE TRIED TO CALL TIN FACILITY CENTER, THEY PROVIDED US MANY OF NUMBERS LIKE [REDACTED], [REDACTED] BUT NONE OF NUMBER IS PROVIDED ANY SOLUTION. THEY ARE NOT ATTENDING CALL, IF\\nATTENDING THEN TRANSFER TO OTHERS, WE ARE FOLLOWING LAST [REDACTED] TO GET STATUS OF PAN BUT WE HAVE NOT\\nFIND ANY SOLUTIONS. IT LOOK LIKE ONLY HARASSMENT FROM LAST [REDACTED]. ALMOST [REDACTED] PASSED, YOU ARE REQUESTED TO\\nISSUE PAN NUMBER OR INFORM TO US FOR QUERY AS EARLIEST AS. PLEASE TAKE NECESSARY ACTION WE REQUIRED URGENT PAN\\nTO OPEN BANK ACCOUNT AND OTHER FORMALITIES	English

Figure 4.2: Semantically Similar Multilingual Grievances. Tokens that could reveal personal identifiable information (PII) have been redacted.

4.8 Conclusion and Future Work

In conclusion, this chapter presents a pipeline for building a similarity search system for subject-content matching in a multilingual setting. The pipeline includes various steps such as text cleaning, language detection, machine translation, and sentence embedding. The system uses a pre-trained multilingual model for machine translation and a pre-trained sentence embedding model for computing sentence embeddings. The hnswlib package is used to build the index for similarity search.

Future work can be done in several directions. Firstly, the system can be extended to support more languages and their corresponding Unicode ranges. Secondly, more advanced techniques can be employed for text cleaning and preprocessing to improve the accuracy of the system. Thirdly, the system can be evaluated and compared with other existing systems for similarity search in a multilingual setting. Lastly, the system can be integrated into an application for practical use in the industry.

Bibliography

- [1] facebook/m2m100_418M. https://huggingface.co/facebook/m2m100_418M, 2021.
- [2] l3cube-pune/hing-bert. <https://huggingface.co/l3cube-pune/hing-bert>, 2021.
- [3] l3cube-pune/hing-mbert. <https://huggingface.co/l3cube-pune/hing-mbert>, 2021.
- [4] salesken/similarity-eng-hin_latin. https://huggingface.co/salesken/similarity-eng-hin_latin, 2021.
- [5] Faisal Ahmed and Abu Sayed Md Choudhury. A framework for citizen engagement in public service delivery. *International Journal of Public Administration in the Digital Age*, 5(1):41–57, 2018.
- [6] Mikel Artetxe, Gorka Labaka, and Eneko Agirre. Cross-lingual models of word embeddings: An empirical comparison. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1301–1311, 2018.
- [7] Debswapna Chakrabarti and Pabitra Mitra. A dynamic programming approach to segmentation of roman texts. *Journal of King Saud University-Computer and Information Sciences*, 29(2):171–178, 2017.
- [8] Shikha Choudhary, Rohit Singh, Arpan Dutta, and Gurmukh Singh. A machine learning approach to address public grievances in india. *2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT)*, pages 301–305, 2018.
- [9] Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. Multilingual word embeddings from wikipedia. *arXiv preprint arXiv:1801.06126*, 2018.
- [10] Deepali Joshi, Nitesh Yadav, Atul Dixit, and Gaurav Singh. Design and implementation of an online public grievance redressal system. *International Journal of Advanced Research in Computer Science and Software Engineering*, 5(6):697–700, 2015.
- [11] Daniel Jurafsky and James H. Martin. *Speech and Language Processing*. Prentice Hall, 2008.
- [12] Yinan Liu, Mingdong Wang, Guangxiang Wang, Yong Chen, and Jun Li. Efficient k-nearest neighbor search using hnsw. In *Proceedings of the 2018 ACM SIGMOD International Conference on Management of Data*, pages 649–664. ACM, 2018.
- [13] Yu. A. Malkov and D. A. Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *Information Systems*, 75:41–56, 2018.
- [14] Yu A Malkov and D A Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. In *Advances in Neural Information Processing Systems*, pages 3847–3857, 2018.
- [15] James Martin. Spelling correction and the edit distance algorithm. *Linguistic Data Consortium, University of Pennsylvania*, 4(2):1–10, 1999.
- [16] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

-
- [17] Md Saifur Rahman, Md Habibul Islam, and Tasnim Mahmud. Sentiment analysis of citizen complaints for government services improvement. *Journal of Information Technology and Economic Development*, 11(2):14–29, 2020.
 - [18] Sebastian Raschka and Vahid Mirjalili. *Python Machine Learning*. Packt Publishing, 2017.
 - [19] Partha Saha, Sourish Bhattacharya, Anirban Das, and Anirban Basu. Improving public grievance redressal system using social media analytics. *2019 2nd International Conference on Inventive Systems and Control (ICISC)*, pages 1321–1326, 2019.
 - [20] Vikas Singh, Prakhar Gupta, Harsh Shrivastava, Ravi Kiran Sarvadevabhatla, and Manish Shrivastava. Hinglish-top: A large scale code-switched semantic parsing dataset, 2022.
 - [21] Esko Ukkonen. Fast approximate string matching with large edit distances. *Algorithmica*, 6(3):251–260, 1985.
 - [22] George Kingsley Zipf. A study of frequency distributions of english letters and some words. *Journal of the American Statistical Association*, 30(201):28–33, 1935.