

Spell Correction

Project report submitted in partial fulfillment of the requirement for the degree
of Bachelor of Technology

in

Computer Science and Engineering/Information Technology

By

Anubhav Thapa(181286)
&
Aashish Chauhan(181244)

Under the supervision of
Dr. Rakesh Kanji

Department of Computer Science & Engineering and Information Technology

**Jaypee University of Information Technology Waknaghat, Solan-173234,
Himachal Pradesh**

Candidate's Declaration

I hereby declare that the work presented in this report entitled "**Spell Correction**" in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science and Engineering/Information Technology** submitted in the department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology Waknaghat is an authentic record of my own work carried out over a period from August 2021 to December 2021 under the supervision of **(Dr. Rakesh Kanji)** (Assistant Professor(SG), Department of Computer Science and Information Technology).

The matter embodied in the report has not been submitted for the award of any other degree or diploma.

Anubhav Thapa(181286)

&

Aashish Chauhan(181244)

This is to certify that the above statement made by the candidate is true to the best of my knowledge.

(Supervisor Signature)

Supervisor Name :Dr. Rakesh Kanji

Designation: Assistant Professor

Department name :Department of Computer Science and Information Technology

Dated:

ACKNOWLEDGEMENT

I would like to thank and express our gratitude to our Project supervisor Dr. Rakesh Kanji for the opportunity that he provided us with this wonderful project "Spell Correction". The outcome would not be possible without his guidance. This project taught me many new things and helped to strengthen concepts of Machine Learning. Next, I would like to express my special thanks to the Lab

Assistant for cordially contacting us and helping us in finishing this project within the specified time.

Lastly, I would like to thank my friends and parents for their help and support.

Anubhav Thapa(181286)
&
Aashish Chauhan(181244)

TABLE OF CONTENT

- 1) Abstract

- 2) Chapter 1- Introduction
 - 1.1 Introduction

 - 1.2 Problem Statement

 - 1.3 Objectives

 - 1.4 Methodology

- 3) Chapter 2- Literature Survey

- 4) Chapter 3- System Development
 - 3.1 Analysis
 - 3.2 Computational
 - 3.3 Experimental
 - 3.4 Mathematical
 - 3.4.1 N-gram Probability
 - 3.4.2 Perplexity
 - 3.4.3 Smoothing
 - 3.4.4 Laplace Smoothing

- 5) Chapter 4- Performance Analysis

- 6) Chapter 5- Conclusions

7) References

ABSTRACT: -

Text correction is an important task in document processing, permitting the automatic managing of sizeable streams of files in digital form. One issue in coping with some lessons of archives is the presence of specific sorts of textual errors, such as spelling and grammatical blunders in email, and persona cognizance blunders in files that come thru OCR. Text correction ought to work reliably on all input, and as a result need to tolerate some degree of these sorts of problems.

We describe right here an N-gram-based strategy to textual content correction that is accepting of textual errors. The gadget is small, quick and robust. This gadget labored very properly for language correction, accomplishing in one take a look at a 99.8% right classification charge on Usenet newsgroup articles written in unique languages. The gadget additionally labored fairly properly for classifying articles from a variety of one of a kind computer-oriented newsgroups in accordance to subject, accomplishing as excessive as an 80% right classification rate. There are additionally a number of apparent instructions for enhancing the system's classification overall performance in these instances the place it did now not do as well.

Natural language processing explains how machines can not only make sense of words but also make sense of words in their context. N-grams are one way to help machines understand a word in its context to get a better understanding of the meaning of a word. For example, we need to "book our tickets soon" versus "we need to read this book soon". The former 'book' is used as a verb and is therefore about the action of planning a trip somewhere. The latter 'book' is used as a noun and is therefore about a little book or object. How do I know this? How can we tell the difference between the verb book and the noun book? We take into account the context of the sentence and we do this innately as we humans have been attuned to language cues since we were born. Machines on the other hand have to learn these cues by looking at the surrounding context of the target word. Think of it like a context window of the before word and after word. This is what n-grams look at. They look at what came before the target word 'book' and what came after to then determine if the word is used as a noun or a verb or in another context. 'This book', 'a book', 'your book', 'my book', 'his book' 'her book', are all examples of by grams where the before word indicates 'book' is used as a noun. The 'n' in n-grams is just the number of words you want to look at.

Bi-grams are two pairs of words that occur together looking at the before word and afterward sliding over the words, for example "read this book soon" is split up into 'read this', 'this book', 'book soon' we could train a machine to analyze that when these phrases 'read this' and 'this book' manifest collectively in pairs the textual content is mainly discussing a literal book. You can additionally prolong the context window to make your n-grams a tri-gram, searching at three pairs of phrases at a time: 'read this book' 'this e book soon'. But endure in thought the longer your context window the more difficult it is to pick out up on phrases that often show up at some stage in the textual content when you are searching at pretty special units of words. I advocate taking the Goldilocks strategy to n-grams: no longer too long, now not too short, simply right. And through simply proper I suggest searching at two pairs of phrases as the earlier than phrase and after phrase is possibly all the context you want to seize the which means of the text. N-grams are additionally beneficial when attempting to seize phrases used in a poor context and vice versa for instance "the group of workers have been now not friendly, terribly really", 'not friendly' and 'friendly terrible' is sufficient context to be aware of that the phrase 'friendly' is used in a bad context. In isolation, the phrase pleasant is fine when we're searching at the earlier than and after word, 'not' and 'terrible' cancel out the high-quality meaning, reversing it to have a bad meaning. Another instance is shooting sarcasm such as "that's funny... not". When 'funny not' happens collectively it additionally cancels out 'funny' and reverses it to be the actual contrary in meaning. By searching at n-grams or pairs of phrases to seize the broader context of phrases to then teach machines to research these language queues and acquire a higher perception of the actual which means of the text. N-grams are a pretty easy but advantageous method to capture the context and that means of phrases in herbal language processing. And that sums up n-grams for you.

CHAPTER - 1

INTRODUCTION

1.1 INTRODUCTION

What is spell checking? Date back to 1980s, a spell checker is more like a "verifier"[1]. It has no corresponding suggestions to the spelling error detected. As many of the readers are using word processor nowadays, a spell checker will first mark a word as mistaken(Detection) and give a list of replacement of word(Suggestion). Therefore the definition of spell checking involve more than only *checking*, it is the process of *detecting* misspelled words in a document/sentence and *suggest* with a suitable word in the context. Therefore, to construct a spell checker, it needs to have the following features:

1. Spelling Detection: the ability to detect a word error
2. Spelling Suggestion(Correction): the ability to suggest a suitable word to users which matches their need in context

Spelling mistakes are collective, and furthestmost persons are discarded to software indicating if a fault was complete. From autocorrect on our receivers, to red emphasizing in text publishing supervisor, spell checking is an essential feature for many different products.

Python offers many modules to use for this purpose, making writing a simple spell checker an easy 20-minute ordeal.

The main aim is to develop a context delicate spell manager to solve the real world delivering errors.

1.2 PROBLEM STATEMENT

The major problem of the basic spell checker is about the spell detection stage. It is designed in the assumption that all word errors are the word that are NOT in the dictionary. These are classified as *non-word spelling error*. However, there are cases where spelling error is not simply a "spelling error", imagine the following case:

I would like a peace of cake as desert.

By simply looking at the words on the sentence above, all of them are fine in terms of spelling. However, errors still occur as the word "peace" and "desert" are not suitable to the context. They are called *real-word spelling errors*. In a spell checker that uses dictionary check, this kind of error will go undetected and proceed. It is clear that dictionary check is not a optimal spelling detection method.

1.3 OBJECTIVE

In this project, the main aim is to develop a context-sensitive spell checker to solve the real-word spelling errors. For real-word spell checking, the spell checker will take a mixed part-of-speech trigram approach for spelling detection and uses confusion sets for spelling corrections. The spell checker will also attempt to combine the base spell checker with the context-sensitive spell checker hence it has both non-word and real-word spell error detections along with corrections. More about the system design will be mentioned in Chapter 5. The performance will be compared with the Google context-sensitive spell checker to measure the outcome of the program, which will be mentioned in Chapter 7.

The ultimate goal of the project is to create a spell checker which can detect and suggest on all type of real-word error made. In the evaluation, the performance will be analysed by the following question:

1. In what extent the real-word typing errors were detected?
2. In what extent a suitable suggestion was given to an user for each type of spelling error?
3. In what extent the performance of the context-sensitive spell checker improved from the base spell checker

1.4 METHODOLOGY:-

1. Correct() function and text blob.

The one easiest way to make spelling corrections is done by this method as it is known to find spelling mistakes by making a pre-existing library given by python itself.

2 N-Gram approach

-
Other approach is by analysing a set of n-grams derived from the context. A n-gram language model is widely used in natural language processing. N-gram means a set of n things, can be letter, words, symbols.... in which word n-gram is used in context-sensitive spell checker. For example, the sentence "I want to be a guy" will derived word trigrams(3-gram): {I, want, to}, {want, to, be}, {be, a, guy}. The use of this model often predicts the probability of item(word) i occurs in the item set j, the probability formula derived as follow[9]:

In an n -gram model, the probability $P(w_1, \dots, w_m)$ of observing the sentence w_1, \dots, w_m is approximated as

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_1, \dots, w_{i-1}) \approx \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1})$$

The formula varies with different n-gram (change of parameter n inside the formula) and it is the basis of many n-gram approach in context-sensitive spell checker.

In 2006, Google published a Web 1T 5-gram striped from their web crawling data. There are subsequence approaches based on that piece of data(e.g. [10], [11]). In general, the method involves matching the words in 5-gram descending to 1-gram, and suggests the suitable correction based on pre-defined confusion set. The conditional probability served as a measure if a word is suitable for the context, also the ranking of error suggestions. This method shows a good accuracy (over 90% average accuracy).

We will also be diving into a bilingual approach to n gram to better understand the fundamentals of the n gram concepts and how it manages to process the languages of different concepts in a given dataset with a dedicated corpus for the better.

We can use sampling to determine and illustrate what sort of information a language model embodies, and we can use it to sample from it. Sampling from a supply entails selecting random points based on their probability. Thus, sampling from a language model—which reflects a distribution of phrases—means generating some sentences and selecting each one based on the model's likelihood. As a result, we're more likely to generate statements with a high likelihood and less likely to generate sentences with a low probability, according to the perfect. This method of displaying a voiced model through specimen was proposed originally.

CHAPTER 2- LITERATURE SURVEY: -

The writing audit on the opinion examination shows the great exploration has been finished by the different specialist's dependent on feeling investigation on archive level.

In this paper was proposed a different multi-mark request on feeling investigation (Liu and Chen, 2015). They have used eleven staggered portrayal strategies appearing on two more limited size blog dataset moreover eight unmistakable appraisal networks for assessment. Besides that, they have also used three particular opinion dictionaries for staggered gatherings. As demonstrated by the analyst, the multi-name plan handle plays out the endeavor essentially in two phases i.e., issue change and estimation change (Zhang and Zhou, 2007). In the issue change stage, the issue is changed into various single-name issues. In the midst of the planning stage, the system gains from these changed single imprint data, and in the testing stage, the informed classifier makes an assumption at a singular name and after that makes a translation of it to a few names. In computation adaptation, the data is changed by the essential of the estimation.

As examined above, determination of right highlights and their scores is the way to work on the exhibition of AI based methodology. TF-IDF and count vectorizer are by and large utilized as highlights for the text order. A Few scientists use dictionary-based methodologies to include extraction and choose the scores in mix with a count vectorizer. Cross-area approval guarantees appropriateness of opinion examination to deal with this present reality of informational indexes were preparing designs is not accessible or costly to acquire. In such a manner, many endeavors have been made in the new past. In the cross-space learning issue, the preparation informational index and the objective informational index are from various sources. For instance, Medinas et al. utilized preparing information from Browser (Customer) dataset and testing information from Miscellaneous (Editor) dataset of CNETs programming download site.

n-gram ($n = 1$ to 5) measurements and different properties of the English language were inferred for applications in regular language comprehension and text handling. They were processed from a notable corpus made out of 1 million word tests. Comparative properties were additionally gotten from the most incessant 1000 expressions of three other corpuses. The positional conveyances of n-grams in the current review are talked about. Measurable examinations on word length and patterns of n-gram frequencies versus jargon are introduced. Notwithstanding an overview of n-gram insights found in the writing, an assortment of n-gram measurements acquired by different scientists is evaluated and thought about.

2.1 CONS:-

In this program n-gram we have evidence of earlier research that we used a program like n-gram in our systems sue to its nature that it has a huge corpus of words that is in bulk we used in the have to be very careful with the size of the load od words the system can handle.

1. Corpus of words.
2. Corpus of words leading to slower process time.
3. Words leading to big computational problem (the bigger the corpus more words the program has to handle in uni-, bi-, tri- gram etc.)
4. Accuracy is harmed because the system has to make the words easier to be used, the accuracy that the program will pick the best word is being damaged.

2.2 Earlier Research:-

Test that had been run by professors and student before were divided into categories because N-gram have a lot of functionality lime text characterization and test manipulation, spell correction etc. use paragraphs, newgroups, book and much more.

1. Become exercise sets for each language to be could be confidential. These are majorly the sets. They follow no particular manner of requirement of samples.
2. Calculated N-gram incidence shapes on the drill sets as mentioned above.

3. Computed a piece object's N-gram figure as labeled overhead.

4. Computed an general coldness amount between the sample's outline and the category outline for each language using the out of place amount, and then picked the sort with the smallest remoteness.

CHAPTER 3 - SYSTEM DEVELOPMENT :-

3.1 ANALYSIS :-

The essential benefit to this methodology is that of is undeniably appropriate for text awaiting from uproarious sources like email or OCR outlines. We initially created N-gram-based ways to deal with different report handling tasks to utilize extremely inferior quality pictures like those found in postal addresses. Albeit one may trust that filtered records that track down their direction into text assortments reasonable for recovery will be of to some degree better caliber, we expect that there will be a lot of changeability in the report information base. This changeability is to be expected to such factors as scanner contrasts, unique report printing quality, bad quality copies, and faxes, just as preprocessing and character acknowledgment contrasts. Our N-gram-based plan gives vigorous access notwithstanding such mistakes. This capacity might make it adequate to utilize an extremely quick yet bad quality person acknowledgment module for comparability examination.

It is conceivable that one could accomplish comparative results utilizing entire word insights. In this approach, one would utilize the recurrence insights for entire words. Notwithstanding, there are a few potential issues with this thought. One is that the framework turns out to be substantially more touchy to OCR issues—a solitary misrecognized character loses the insights for an entire word. A second conceivable trouble is that short sections (for example, Usenet articles) are basically excessively short to get agent subject word measurements. By definition, there are basically more N-grams in guaranteed sections than there are words, and there are thus more prominent freedoms to gather enough N-grams to be critical for coordinating. We trust to straightforwardly think about the presentation of N-gram based profiling with entire word-based profiling soon.

One more related thought is that by utilizing N-gram investigation, we get word stemming basically for free. The N-grams for related types of a word (e.g., 'advance', 'progressed', 'progressing', 'headway', and so on) consequently have a ton in normal when seen as sets of N-grams. To get identical outcomes with entire words, the framework would need to perform word stemming, which would necessitate that the framework have definite information about the specific language that the records were written in. The N-gram recurrence approach gives language autonomy for free.

Acquired preparing sets (class tests) for every language to be arranged. Commonly, these preparation sets were on the request of 20K to 120K bytes long. There was no specific configuration necessity, yet all the same each preparing set didn't contain tests of any language other than the one it should address.

- Figured N-gram recurrence profiles on the preparation sets as depicted previously.
- Figured each article's N-gram profile as depicted previously. The subsequent profile was on the request for 4K long.
- Figured a general distance measure between the example's profile and the classification profile for every language utilizing the awkward measure, and afterward picked the class with the littlest distance.

Such a framework has unobtrusive computational and capacity prerequisites, and is exceptionally successful. It requires no semantic or content investigation separated from the N-gram recurrence profile itself.

3.2 COMPUTATIONAL :-

We address the issue of foreseeing something from past words in an example of text. Specifically, we examine n-gram models dependent on classes of words. We likewise talk about a few factual calculations for relegating words to classes dependent on the recurrence of their co-event with different words. We observe that we can separate classes that have the kind of either linguistically based groupings or semantically based groupings, contingent upon the idea of the hidden insights.

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Pre-handling of text for language ID tasks essentially intends to eliminate messages which are language autonomous/breaking down elements and fuse the rationale which can improve the exactness of the ID task. We have considered the accompanying pre-handling ventures prior to making bi-gram language model.

Every one of the texts was changed over to bring down the case.

Every one of the digits were taken out from the message sentences.

Accentuation imprints and unique characters were taken out.

Every one of the sentences was connected with space in the middle.

Series of touching blank areas were supplanted by single space.

It is imperative to take note of that the text document should be perused in Unicode design which envelops the person set including every one of the dialects. The Python code for previously mentioned steps can be seen in the next segment.

3.3 EXPERIMENTAL:-

Test corpus contains almost 10,000 sentences for every language. To group a message sentence among the language models, the distance of the info sentence is determined with the bi-gram language model. The language with the negligible distance is picked as the language of the information sentence. When the pre-handling of the info sentence is done, the bi-grams are separated from the information sentence. Presently, the frequencies of every one of these bi-grams are determined from the language model and are summarized. The summarized recurrence incentive for every language is standardized by the amount of frequencies of the multitude of bi-grams in the separate language. This standardization is important to eliminate any predisposition because of the size of the preparation text corpus of every language. Likewise, we have duplicated the frequencies by an element of 10,000 to stay away from the situation when standardized recurrence ($f/\text{total}[i]$) becomes zero. The condition for the previously mentioned computation is given below (where is condition).

$$F(j) = \frac{\sum_{i=1}^k C(i,j) * 10000}{\sum_{i=1}^m C(i,j)}$$

, where $F(j)$ is the standardized recurrence amount of language, $C(i,j)$ is the recurrence count of the i^{th} bi-gram in j^{th} language. k is the quantity of bi-grams which happen in the test sentence, while m is the complete number of bi-grams in a similar language.

The full execution of shut set assessment of language recognizable proof undertaking on wortschatz test corpus is given beneath. tp and fp are valid up-sides and bogus up-sides separately. Genuine up-sides are the number of sentences which are identified effectively and bogus up-sides are the number of sentences which were wrongly distinguished as another dialect.

Size of preparing set – Larger preparing corpus will prompt estimation of exact measurements (recurrence counts) of bi-grams in the language. One can make the language models on 1 million sentences downloaded from wortschatz leipzig corpus.

There are loads of named substances (formal people, places or things) in the message sentence which debase the language model as these names are consistently language free. As I would like to think, the exactness of location errand will increase in the event that we can eliminate such words.

In the following approach of n-gram models, we have made models with $n = 2$. Exactness accomplished in the assessment interaction will positively increment as $n = 3$ or 4 (tri-grams and quad-grams) will be utilized.

The pairwise correlation of proteins depends on the substance normalities expected to extraordinarily portray each succession. These abnormalities are caught by n-gram based displaying methods and in the spin-off are differentiated by cross-entropy related measures. In this absolute first endeavor to intertwine hypothetical thoughts from computational etymology inside the field of bioinformatics, we tried different things with various executions having consistently as extreme objective the improvement of pragmatic, computational effective calculations. The exploratory investigation gives proof to the convenience of the new methodology and persuades the further improvement of etymology related devices as a way to unravel the natural arrangements.

Two central issues concern the treatment of huge n-gram language models: ordering, that is, compacting the n-grams and related satellite qualities without undermining their recovery speed, and assessment, that is, registering the likelihood circulation of the n-grams removed from an enormous literary source.

Playing out these two errands proficiently is essential for a considerable length of time in the fields of Information Retrieval, Natural Language Processing, and Machine Learning, for example, auto-fruitation in web search tools and machine interpretation.

Concerning the issue of ordering, we portray compacted, precise, and lossless information structures that all the while accomplishes high space decreases and no time debasement as for the cutting-edge arrangements and related programming bundles. Specifically, we present a compacted tire information

structure in which each expression of a n -gram following a setting of fixed length k , that is, its previous k words, is encoded as a whole number whose worth is relative to the quantity of words that follow such setting. Since the quantity of words following a given setting is ordinarily tiny in regular dialects, we bring down the space of portrayal to pressure levels that were never accomplished, permitting the ordering of billions of strings. In spite of the critical investment funds in space, our procedure presents a unimportant punishment at inquiry time.

In particular, the most space-proficient rivals in the writing, which are both quantized and lossy, don't take not exactly our trie information structure and are up to multiple times slower. On the other hand, our trie is just about as quick as the quickest contender yet additionally holds a benefit of up to 65% in outright space.

With respect to the issue of assessment, we present an original calculation for assessing changed Kneser-Ney language models that have arisen as the true decision for language displaying in both scholarly world and industry because of their generally low perplexity execution. Assessing such models from enormous literary sources represents the test of conceiving calculations that utilize the circle.

The best-in-class calculation utilizes three arranging steps in outside memory: we show a further development that requires just one arranging venture by taking advantage of the properties of the removed n -gram strings. With a broad exploratory investigation performed on billions of n -grams, we show a normal improvement of 4.5 occasions on the complete runtime of the past approach.

Clients interface with web-based media in various ways, giving an assortment of information, from evaluations and endorsements to amounts of text. Public

conversation for areas of interest specifically creates a huge volume and speed of client contributed text, much of the time inferable from a client identifier or pseudonym. It might very well be possible to decide the creation of different lots of text via online media utilizing n-gram examination on the piece level interpretation of the text. This paper investigates the office of spot level n-gram examination with other measurable arrangement approaches for deciding origin on two months of caught client postings from a web-based news and assessment site with direct conversation. The outcomes show that this methodology can accomplish a decent acknowledgment rate with a low bogus negative rate.

So, assuming that we are given a corpus of text and need to think about two diverse n-gram models, we partition the information into preparing and test sets, train the

boundaries of both models on the preparation set, and afterward look at how well the two prepared models fit the test set.

Be that as it may, what's the significance here to "fit the test set"? The appropriate response is basic: whichever model allots a higher likelihood to the test set—which means it all the more precisely predicts the test set—is a superior model. Given two probabilistic models, the better model is the one that throws a tantrum to the test information or that better predicts the subtleties of the test information, and thus will dole out a higher likelihood to the test information.

3.4 MATHEMATICAL: -

N-gram Probabilities: -

This framework proposes words which could be utilized next in a given sentence. Assume I give the framework the sentence "Thank you kindly for your" and anticipate that the system should foresee what the following word will be. Presently you and I both realize that the following word is "help" with an exceptionally high likelihood. However, how might the framework realize that?

Something significant to note here is that, concerning some other man-made brainpower or AI model, we want to prepare the model with an enormous corpus of information. When we do that, the framework, or the NLP model will have a very smart thought of the "likelihood" of the event of a word after a specific word. So, trusting that we have prepared our model with a colossal corpus of information, we'll accept that the model furnished us the right response.

I talked about the likelihood of a piece there, yet we should now expand on that. At the point when we're fabricating an NLP model for anticipating words in a sentence, the likelihood of the event of a word in a succession of words is what makes a difference. Also, how would we gauge that? Suppose we're working with a bigram model here, and we have the accompanying sentences as the preparation corpus:

$$\text{count}(w2\ w1) / \text{count}(w2)$$

- Example: -
1. I really like snow.
 2. We really get to stop.
 3. You don't know what it is really like.

```
count(really like) / count(really)
= 1 / 3
= 0.33
```

Also, assuming we needed to know the joint likelihood of a whole arrangement of words like its water is so straightforward, we could do it by inquiring "out of all conceivable successions of five words, the number of them is that its water is so straightforward?" We would need to get the count of its water so straightforward and partition by the amount of the counts of all conceivable five-word groupings. That appears to be fairly a ton to appraise! Therefore, we'll need to present more sharp methods of assessing the likelihood of a word w given a set of experiences h , or the likelihood of a whole word arrangement W . We should begin with a little formalization of documentation. To address the likelihood of a specific arbitrary variable X_i taking on the value "the", or $P(X_i = \text{"the"})$, we will utilize the improvement $P(\text{the})$. We'll address a succession of N words either as

$w_1 \dots w_n$ or $w_1 : n$ (so the articulation $w_1 : n-1$ implies the string w_1, w_2, \dots, w_{n-1}). For the joint likelihood of each word in a succession having a specific worth $P(X = w_1, Y = w_2, Z = w_3, W = w_n)$ we'll use $P(w_1, w_2, \dots, w_n)$.

Presently how might we figure probabilities of whole arrangements like $P(w_1, w_2, \dots, w_n)$?

One thing we can do is break down this likelihood utilizing the chain rule of likelihood:

$$\begin{aligned} P(X_1 \dots X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_{1:2}) \dots P(X_n|X_{1:n-1}) \\ &= \prod_{k=1}^n P(X_k|X_{1:k-1}) \end{aligned} \tag{3.3}$$

Applying the chain rule to words, we get: -

$$\begin{aligned} P(w_{1:n}) &= P(w_1)P(w_2|w_1)P(w_3|w_{1:2}) \dots P(w_n|w_{1:n-1}) \\ &= \prod_{k=1}^n P(w_k|w_{1:k-1}) \end{aligned} \tag{3.4}$$

The chain rule shows the connection between processing the joint likelihood of a grouping

What's more, registering the restrictive likelihood of a word given past words. Condition 3.4 proposes that we could assess the joint likelihood of a whole arrangement of

words by duplicating together various contingent probabilities.

Perplexity:-

By and by we don't utilize crude likelihood as our measurement for assessing language model perplexity els, however a variation called perplexity. The perplexity (at times called PP for short)

of a language model on a test set is the converse likelihood of the test set, standardized

by the quantity of words. For a test set $W = w_1w_2 \dots w_N$,:

$$\begin{aligned} \text{PP}(W) &= P(w_1w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{P(w_1w_2 \dots w_N)}} \end{aligned} \quad (3.14)$$

We can utilize the chain rule to grow the likelihood of W :

$$\text{PP}(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i|w_1 \dots w_{i-1})}} \quad (3.15)$$

In this way, in case we are processing the perplexity of W with a bigram language model,

we get:

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i|w_{i-1})}} \quad (3.16)$$

Note that in view of the opposite in Eq. 3.15, the higher the contingent likelihood of the word arrangement, the lower the perplexity. In this manner, limiting perplexity is comparable to amplifying the test set likelihood as per the language model.

What we for the most part use for word grouping in Eq. 3.15 or Eq. 3.16 is the whole grouping of words in some test set. Since this arrangement will cross many sentence limits, we really want to incorporate the start and end-sentence markers <s> and </s> in the likelihood calculation. We likewise need to incorporate the finish-of-sentence marker </s> (however not the start-of-sentence marker <s>) in the all out count of word tokens N.

There is one more method for considering perplexity: as the weighted normal stretching element of a language. The spreading component of a language is the quantity of conceivable next words that can follow any word. Think about the assignment of perceiving the digits in English (zero, one, two,..., nine), considering that (both in some preparation set and in a few test sets) every one of the 10 digits happens with equivalent likelihood $P = 1, 10$. The perplexity of this small scale language is truth be told 10. To see that, envision a test series of digits of length N, and accept that in the preparation set every one of the digits happened with equivalent likelihood.

By Eq. 3.15, the perplexity will be

$$\begin{aligned}
 \text{PP}(W) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\
 &= \left(\frac{1}{10}\right)^{-\frac{1}{N}} \\
 &= \frac{1}{10} \\
 &= 10
 \end{aligned}
 \tag{3.17}$$

In any case, assume that the number zero is truly successive and happens undeniably more frequently than different numbers. Suppose that 0 happens multiple times in the preparation set, and each of different digits happened 1 time each. Presently we see the accompanying test set: 0 0 0 0 0 3 0 0 0 0. We ought to expect the perplexity of this test set to be lower since more often than not the following number will be zero, which is truly unsurprising, for example has a high likelihood. Subsequently, albeit the stretching factor is as yet 10, the perplexity or weighted expanding factor is more modest. We leave this accurate estimation as exercise 12. We find in Section 3.8 that perplexity is additionally firmly identified with the information theoretic thought of entropy.

At last, how about we check out an illustration of how perplexity can be utilized to look at changed n-gram models. We prepared unigram, bigram, and trigram language structures on 38 million words (counting beginning-of-sentence tokens) from the Wall Street Journal, utilizing 19,979-word jargon. We then, at that point, figured the perplexity of each of these models on a test set of 1.5 million words with Eq. 3.16. The table beneath shows the perplexity of a 1.5-million-word WSJ test set by every one of these sentence structures.

	Unigram	Bigram	Trigram
Perplexity	962	170	109

As we see over, the more data the n-gram gives us about the word grouping, the lower the perplexity (since as Eq. 3.15 showed, perplexity is connected conversely to the probability of the test arrangement as indicated by the model).

Note that in registering perplexities, the n-gram model P should be developed with next to no information on the test set or any earlier information on the jargon of the test set. Any sort of information on the test set can make the perplexity be misleadingly low. The perplexity of two language models is just similar if they utilize indistinguishable vocabularies.

An (inherent) improvement in perplexity doesn't ensure an (outward) improvement in the presentation of a language handling task like discourse acknowledgment or then again machine interpretation. In any case, since perplexity regularly corresponds with such upgrades, it is normally utilized as a fast keep an eye on a calculation. Yet, a model's improvement in perplexity ought to consistently be affirmed by a start to finish assessment of a genuine undertaking prior to closing the assessment of the model.

Smoothing: -

How would we manage words that are in our jargon (they are not obscure words) but show up in a test set in an inconspicuous setting (for instance they show up after a word they never showed up after in preparing)? To keep a language model from relegating no likelihood to these inconspicuous occasions, we'll need to shave off a touch of likelihood mass from some more regular occasions and give it to the occasions we've won't ever see. This adjustment is called smoothing or limiting. In this part and the accompanying ones, we'll acquire an assortment of ways to do smoothing: Laplace (add-one) smoothing, add-k smoothing, moronic backoff, and Kinser-Ney smoothing.

Laplace Smoothing: -

The least difficult method for doing smoothing is to add one to all the n-gram counts, previously we standardize them into probabilities. Every one of the counts that used to be zero will now have a count of 1, the counts of 1 will be 2, etc. This calculation is called Laplace smoothing. Laplace smoothing doesn't perform all around ok to be utilized in present day n-gram models, yet it helpfully presents a considerable lot of the ideas that we see in other smoothing calculations, gives a valuable benchmark, and is likewise a functional smoothing calculation for different errands like text arrangement

How about we start with the use of Laplace smoothing to unigram probabilities. Review that the unsmoothed greatest probability gauge of the unigram likelihood of the word w_i is its count c_i standardized by the all-out number of word tokens N :

$$P(w_i) = \frac{c_i}{N}$$

Laplace smoothing only adds one to each count (thus its substitute name adds one smoothing). Since there are V words in the jargon and every one was increased, we additionally need to change the denominator to consider

the additional V perceptions. (What befalls our P esteems assuming we don't expand the denominator?)

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V} \quad (3.20)$$

Rather than changing both the numerator and denominator, it is advantageous to portray what a smoothing calculation means for the numerator, by characterizing a changed count c^* . This changed count is simpler to contrast straightforwardly and the MLE counts and can be transformed into a likelihood like a MLE count by normalizing by N . To characterize this count, since we are just changing the numerator as well as adding 1 we'll additionally need to increase by a standardization factor $N/N+V$:

$$c_i^* = (c_i + 1) \frac{N}{N + V} \quad (3.21)$$

We would now be able to turn c^*l into a likelihood P^*l by normalizing by N . A connected method for survey smoothing is as limiting (bringing down) some non-zero\ includes to get the likelihood mass that will be allotted to the zero counts. In this way, rather than alluding to the limited counts c , we may portray a smoothing calculation as far as a relative rebate d_c , the proportion of the limited counts to the first counts:

$$d_c = \frac{c^*}{c}$$

Since we have the instinct for the unigram case, how about we smooth our Berkeley Restaurant Project bigrams. Figure 3.6 shows the add-one smoothed counts for the bigrams.

Figure 3.7 shows the add-one smoothed probabilities for the bigrams. Review that ordinary bigram probabilities are figured by normalizing each column of sums by the unigram total:

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})} \quad (3.22)$$

For the method addone smoothed bigram count that we are observing and want to upsurge the unigram sum by the number of whole word that are written in the missed jaron V:

$$P_{\text{Laplace}}^*(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{\sum_w (C(w_{n-1}w) + 1)} = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + V} \quad (3.23)$$

	i	want	to	eat	chinese	food	lunch	spend
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

Method smoothed bigram counts for number of the observation (V that is 1484) in the major corpus of more than 10K words in sentences.

	i	want	to	eat	chinese	food	lunch	spend
i	0.0015	0.21	0.00025	0.0025	0.00025	0.00025	0.00025	0.00075
want	0.0013	0.00042	0.26	0.00084	0.0029	0.0029	0.0025	0.00084
to	0.00078	0.00026	0.0013	0.18	0.00078	0.00026	0.0018	0.055
eat	0.00046	0.00046	0.0014	0.00046	0.0078	0.0014	0.02	0.00046
chinese	0.0012	0.00062	0.00062	0.00062	0.00062	0.052	0.0012	0.00062
food	0.0063	0.00039	0.0063	0.00039	0.00079	0.002	0.00039	0.00039
lunch	0.0017	0.00056	0.00056	0.00056	0.00056	0.0011	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

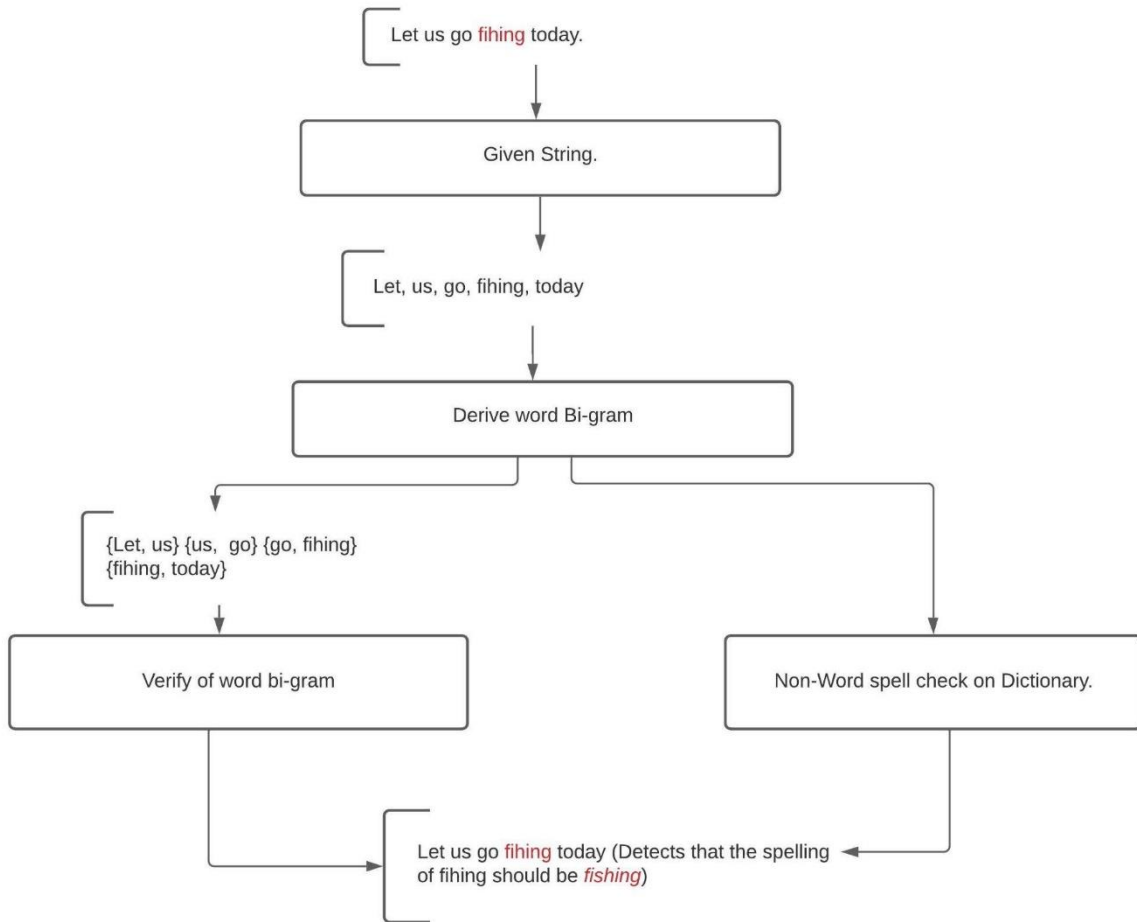
CHAPTER 4 - PERFORMANCE ANALYSIS :-

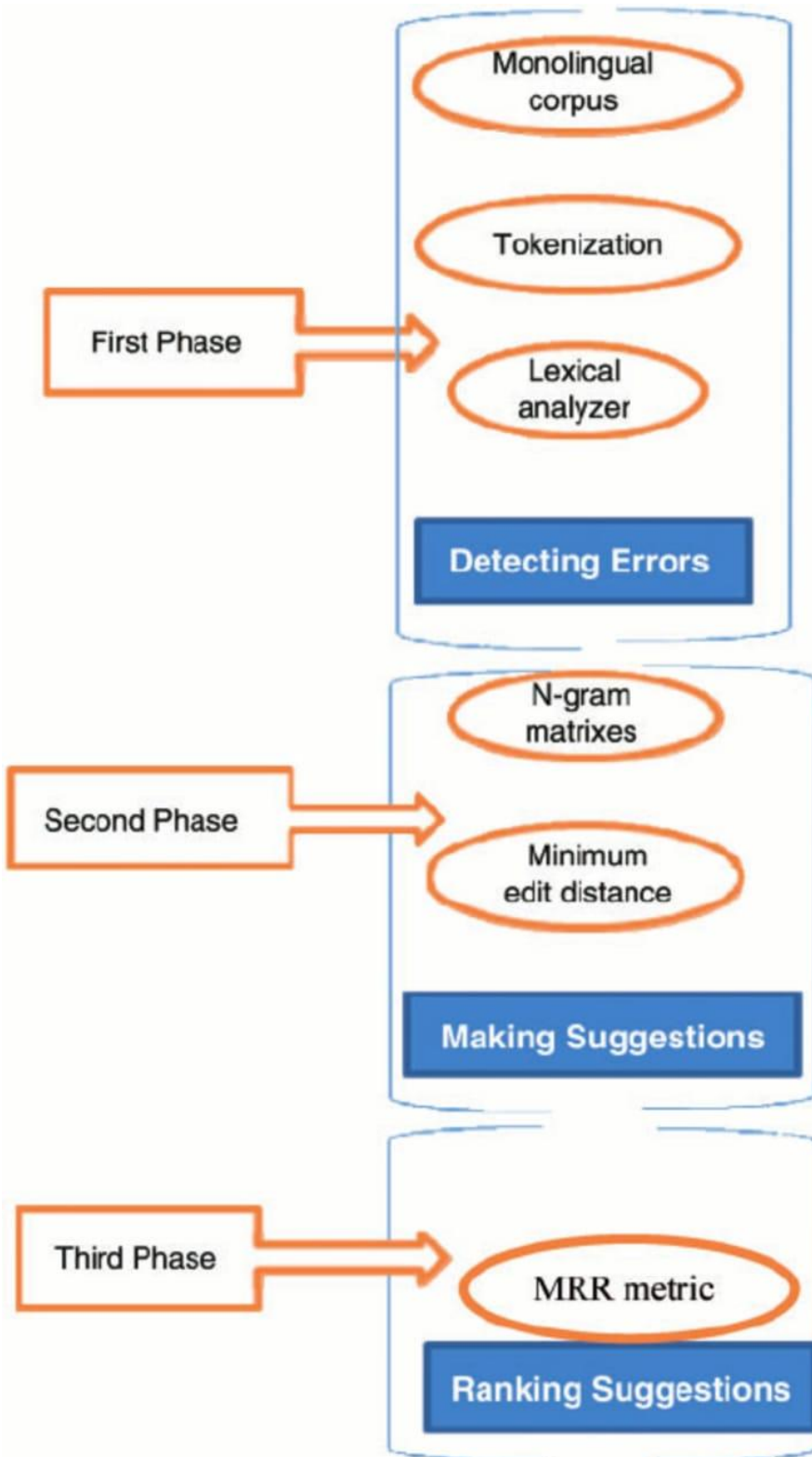
Neural Network language models (NNLMs) have as of late become a significant supplement to regular n-gram language models (LMs) in discourse-to-message frameworks. In any case, little is been called to be known that had some significant awareness of the conduct of NNLMs. The investigation introduced in this paper expects to comprehend which kinds of occasions are better displayed by NNLMs when contrasted with n-gram LMs, in what cases upgrades are generally generous and why this is the situation. Such an examination is critical to take further advantage from NNLMs utilized in blend with regular n-gram models. The investigation is completed for various sorts of neural organization (feed-forward and repetitive) LMs. The outcomes appearing for which kind of occasions NNLMs give better likelihood gauges are approved on two arrangements that are diverse in their size and the level of information homogeneity.

1. When we use the correct() function tells us that this technique succeeded to form an opinion that gets the spelling fault percentage from 60.6% to 15.9%.
2. The function is slow when compared to other alternatives.
3. Correct() func tends to be more accurate than with small amounts of data when it is given a bigger set of databases then it gives in.
4. The spelling mistakes that it corrects are accurate but when there are two words eg:- 'Ber' it can be replaced with 'beard', 'bear' or 'beer' making the meaning of the sentence inaccurate.
5. Domain Specific Features in the Corpus - Words that occur only under the given domain so the words belong to one given text like if we use API or links or slangs .
6. Use An Exhaustive Stop word List - The most common words that are guaranteed to occur like the, a, of, etc.
7. Noise Free Corpus - No extra words that end up littering the text using only the words that we need to use, not the links, punctuation marks

etc..

8. Eliminating features with extremely low frequency - the unused or non-repetitive words can be deleted before because it is not that often used therefore be a mess.
9. Normalized Corpus- Using the root form of the word only no nouns, pronouns, adjectives etc.





CHAPTER - 5

Conclusions

Conclusions and Future Work

In this paper we introduce a concept of syntactic n-grams (syntactic n-grams). The difference between traditional n-grams and syntactic n-grams is related to the manner of what elements are considered neighbors. In case of syntactic n-grams, the neighbors are taken by following syntactic relations in syntactic trees, while traditional n-grams are formed as they appear in texts. The concept of syntactic n-grams allows bringing syntactic information into machine learning methods. Syntactic n-grams can be applied in all tasks when traditional n-grams are used. Any syntactic representation can be used for application of syntactic n-gram technique: dependency trees or constituency trees. In the case of dependency trees, we should follow the syntactic links and obtain syntactic n-grams.

In the case of constituency trees, some additional steps should be made, but these steps are very simple. We conducted experiments for authorship attribution tasks using SVM, NB, and J48 for several profile sizes.

Currently the system uses a number of different N-grams, some of which ultimately are more dependent on the language of the document than the words comprising its content. By omitting the statistics for those N-grams which are extremely common because they are essentially features of the language, it may be possible to get better discrimination from those statistics that remain. It is also possible that the system should include some additional statistics for rarer N-grams, thus gaining further coverage.

It seems clear that the quality of the document set affects the subject categorization performance. We would like to experiment with document sets that have a higher overall coherence and quality. For example, it would be interesting to test this technique on a set of technical abstracts for several different areas. By splitting the set for each area into training and testing portions, then computing the profile for each area from the training set, we could repeat this experiment in a more controlled way

In a related issue, the quality of the training set in general greatly affects matching performance. Although the FAQs were easy to obtain and work with, other training sets might have produced better results, even for these newsgroups. Of necessity, a FAQ lags the group it covers, since new “hot” topics of discussion have not yet made it into the FAQ. To test this, it would be interesting to compare the FAQ-based profiles with profiles derived from a separate set of articles from the appropriate newsgroups.

The raw match scores the system produces are largely useless by themselves except for imposing an overall relative ordering of matches for the various profiles. To correct this, we must devise a good normalization scheme, which would produce some sort of absolute measure of how good a particular match really is. This would allow the system to reject some documents on the grounds that their normalized scores were so low that the documents did not have good matches at all. Normalized scores would also let the system determine if a particular document lay between two classifications because of its interdisciplinary nature. A related idea would be to see how well the system could predict which articles get cross-posted to different groups precisely because of their interdisciplinary content.

This type of document similarity measure is ideally suited for document filtering and routing. All that a user needs to do is collect a representative set of documents that cover the relevant topics, then compute an overall profile. From that point on, it is simple and cheap to compute the profile of every incoming document, match it against the user’s overall profile, and accept those whose match scores are sufficiently good.

This system currently handles only languages that are directly representable in ASCII. The emerging ISO-6048/UNICODE standard opens up the possibility of applying the N-gram frequency idea to all of the languages of the world, including the ideographic ones.

Relatively massive corpus of works of three authors was once used. We used as a baseline characteristic ordinary n-grams of words, POS tags and characters. The consequences exhibit that sn-gram approach outperforms the baseline technique. The following instructions of future work can be mentioned: – Experiments with all characteristic units on large corpus (and greater authors, i.e., extra classes). – Analysis of the applicability of shallow parsing as a substitute of full parsing..

Investigation of the handiness of sn-grams of characters. – Analysis of the effect of parser mistakes on the presentation of sn-grams. – Analysis of conduct of sn-grams between dialects, e.g., in equal texts or similar texts. – Application of sn-grams in other NLP assignments. – Application of blended sn-grams.

Experiments that would reflect onconsideration on combos of the cited aspects in one characteristic vector. – Evaluation of the most useful range and measurement of sn-grams for a range of tasks. – Consideration of quite a number profile sizes with greater granularity. – Application of sn-grams in different languages.

This investigation suggested that spell-checkers have practically no effect by any stretch of the imagination on students' spelling botches on the scholarly level. It didn't help with fixing the botches, and the corrections are not masked; in like manner, allowing understudies to reiterate a comparative misstep. It is believed that future assessments contemplate the constraints of the survey and the thoughts for extra investigation. This is fundamental for future assessments to arrive at additional significant judgments due to spelling-checkers on students' abilities to deliver fixes. The individuals being understudies of an academically regarded school, it is typical that the understudies are most radically loath to commit blunder. Regardless, considering the disclosures, even awesome understudies for the most part disdain language capacity despite scoring An or B for English in their UPSR evaluations. As educators, we should understand that advancement can tragically do a restricted sum to help language understudies in additional fostering their language capacities. Albeit the spelling-checker was not created to help language understudies learn and deal with their

spelling, it tends to be utilized to fill that want with official course by way of the language teachers.

We see blunder identification, modification procedures, the phrase encouraged to the stop purchaser relies upon on two calculations one is Jaccard coefficient and 2nd is Levenshtein distance. These calculations sift thru the phrase reference phrases and provide the particular notion to the client, so the consumer enters textual content in the editorial supervisor ought to be a blunder free and it does not include any spelling botches.

REFERENCES:-

- Algoet, P. H. and T. M. Cover. 1988. A sandwich proof of the Shannon-McMillan-Breiman theorem. *The Annals of Probability*, 16(2):899–909.
- Bahl, L. R., F. Jelinek, and R. L. Mercer. 1983. A maximum likelihood approach to continuous speech recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 5(2):179–190.
- Baker, J. K. 1975a. The DRAGON system – An overview. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, ASSP-23(1):24–29.
- Baker, J. K. 1975b. Stochastic modeling for automatic speech understanding. In D. Raj Reddy, editor, *Speech Recognition*. Academic Press.
- Brants, T., A. C. Popat, P. Xu, F. J. Och, and J. Dean. 2007. Large language models in machine translation. *EMNLP/CoNLL*.
- Buck, C., K. Heafield, and B. Van Ooyen. 2014. N-gram counts and language models from the common crawl. *LREC*.
- Chen, S. F. and J. Goodman. 1998. An empirical study of smoothing techniques for language modeling. Technical Report TR-10-98, Computer Science Group, Harvard University.
- Jelinek, F. and R. L. Mercer. 1980. Interpolated estimation of Markov source parameters from sparse data. In Edzard S. Gelsema and Laveen N. Kanal, editors, *Proceedings, Workshop on Pattern Recognition in Practice*, pages 381–397. North Holland.
- Johnson, W. E. 1932. Probability: deductive and inductive problems (appendix to). *Mind*, 41(164):421–423.
- Jurafsky, D., C. Wooters, G. Tajchman, J. Segal, A. Stolcke, E. Fosler, and N. Morgan. 1994. The Berkeley restaurant project. *ICSLP*. Jurgens, D., Y. Tsvetkov, and D.
- Jurafsky. 2017. Incorporating dialectal variability for socially equitable language identification. *ACL*.
- King, S. 2020. From African American Vernacular English to African American Language: Rethinking the study of race and language in African Americans' speech. *Annual Review of Linguistics*, 6:285–300.

- Kneser, R. and H. Ney. 1995. Improved backing-off for Mgram language modeling. ICASSP, volume 1.
- Lin, Y., J.-B. Michel, E. Aiden Lieberman, J. Orwant, W. Brockman, and S. Petrov. 2012. Syntactic annotations for the Google books Ngram corpus. ACL.
- Markov, A. A. 1913. Essai d'une recherche statistique sur le texte du roman "Eugene Onegin" illustrant la liaison des epreuve en chain ('Example of a statistical investigation of the text of "Eugene Onegin" illustrating the dependence between samples in chain'). Izvestia Imperatorski Akademii Nauk (Bulletin de l'Académie Impériale des Sciences de St.-Petersbourg) , 7:153–162.
- Mikolov, T. 2012. Statistical language models based on neural networks. Ph.D. thesis, Ph. D. thesis, Brno University of Technology.
- Shannon, C. E. 1948. A mathematical theory of communication. Bell System Technical Journal, 27(3):379–423. Continued in the following volume.
- Shannon, C. E. 1951. Prediction and entropy of printed English. Bell System Technical Journal, 30:50–64.
- Stolcke, A. 1998. Entropy-based pruning of backoff language models. Proc. DARPA Broadcast News Transcription and Understanding Workshop.
- Stolcke, A. 2002. SRILM – an extensible language modeling toolkit. ICSLP .
- Talbot, D. and M. Osborne. 2007. Smoothed Bloom filter language models: Tera-scale LMs on the cheap. EMNLP/CoNLL .
- Witten, I. H. and T. C. Bell. 1991. The zero-frequency problem: Estimating the probabilities of novel events in adaptive text compression. IEEE Transactions on Information Theory, 37(4):1085–1094.

APPENDICES: -

The fundamental science of the n-gram was first proposed by Markov (1913), who utilized what are presently called Markov chains (bigrams and trigrams) to foresee regardless of whether a forthcoming letter in Pushkin's Eugene Onega would be a vowel or a consonant. Markov arranged 20,000 letters as V or C and processed the bigram and trigram likelihood that a given letter would be a vowel given the past one or two letters. Shannon (1948) applied n-grams to process approximations to English word arrangements. In view of Shannon's work, Markov models were generally utilized in designing, semantic, and mental work on displaying word groupings by the 1950s. In a progression of amazingly persuasive papers beginning with Chomsky (1956) and counting Chomsky (1957) and Miller and Chomsky (1963), Noam Chomsky contended that "limited state Markov processes", while a conceivably valuable designing heuristic, were unequipped for being a finished intellectual model of human syntactic information. These contentions drove numerous etymologists and computational language specialists to overlook work in factual displaying for quite a long time. The resurgence of n-gram models came from Jelinek and partners at the IBM Thomas J. Watson Research Center, who were impacted by Shannon, and Baker at CMU, who were affected by craft by Baum and associates. Freely these two labs effectively utilized n-grams in their discourse acknowledgment frameworks (Baker 1975b, Jelinek 1976, Baker 1975a, Buhl et al. 1983, Jelinek 1990). Add-one smoothing gets from Laplace's 1812 law of progression and was first applied as a designing answer for the zero-recurrence issue by Jeffreys (1948) in view of a previous add idea by Johnson (1932). Issues with the add-on calculation are summed up in Gale and Church (1994). A wide range of language demonstrating and smoothing strategies were proposed during the 80s and 90s, including Good-Turing limiting—first applied to the n-gram smoothing at IBM by Katz (Nada' ' 1984, Church and Gale 1991)—Witten-Bell limiting (Witten and Bell, 1991), and assortments of class-based n-gram models that pre-owned data about word classes.

Beginning in the last part of the 1990s, Chen and Goodman played out various cautiously controlled trials looking at changed limiting calculations, reserve models, class-based models, and other language model boundaries (Chen and Goodman 1999, Goodman 2006, entomb alia). They showed the

benefits of adapted Inserted Kneser-Ney, which was the same old gauge for n-gram linguistic presentation, exactly in minor of the certainty that they established that assets and class-based modes gave just minor extra development. these papers are cautioned for any peruser with extra curiosity in n-gram language displaying. SRILM (Stolcke, 2002) and KenLM (Heafield 2011, Heafield et al. 2013) are openly offered stratagem mass for making n-gram verbal ways.

Contemporary dialectal showing is all of the extra frequently finished with neural corporation philological mockups, which cope with the thrilling quandaries with n-grams: the number of barriers increments dramatically as the n-gram request increments, and n-grams need any tactic for tallying up from fixing to check set. Neuronic dialectal representations instead challenge phrases right into a nonstop area in which idioms with equal situations have comparable interpretations.

Here, the cleaned text input

People have travelled through and inhabited the Toronto area, located on a broad sloping plateau interspersed with rivers, deep ravines, and urban forest, for more than 10,000 years. After the broadly disputed Toronto Purchase, when the Mississauga surrendered the area to the British Crown, the British established the town of York in 1793 and later designated it as the capital of Upper Canada. During the War of 1812, the town was the site of the Battle of York and suffered heavy damage by American troops. York was renamed and incorporated in 1827 as the city of Toronto. It was designated as the capital of the province of Ontario in 1827 during Canadian Confederation. The city proper has since expanded past its original borders through both annexation and amalgamation to its current area of 630.8 km² (243.5 sq mi). The diverse population of Toronto reflects its current and historical role as an important destination for immigrants to Canada. More than 50 percent of residents belong to a visible minority population group, and over 100 distinct ethnic origins are represented among its inhabitants. While the majority of Torontonians speak English as their primary language, over 100 languages are spoken in the city. Toronto is a prominent center for music, theatre, motion picture production, and television production, and is home to the headquarters of Canada's major national broadcast networks and media outlets. Its varied cultural institutions, which include numerous museums and galleries, festivals and public events, entertainment districts, national historic sites, and sports activities, attract over 10 million tourists each year. Toronto is known for its many skyscrapers and high rise buildings, in particular the tallest free standing structure in the Western Hemisphere, the CN Tower.

Output:

People have travelled through and inhabited the Toronto area, located on a broad sloping plateau interspersed with rivers, deep ravines, and urban forest, for more than 10,000 years. After the broadly disputed Toronto Purchase, when the Mississauga surrendered the area to the British Crown, the British established the town of York in 1793 and later designed it as the capital of Upper Canada. During the War of 1812, the town was the site of the Battle of York and suffered heavy damage by American troops. York was renamed and incorporated in 1827 as the city of Toronto. It was designated as the capital of the province of Ontario in 1827 during Canadian Confederation. The city proper has since expanded past its original borders through both annexation and amalgamation to its current area of 630.2 km² (243.3 sq mi). The diverse population of Ontario reflects its current and historical role as an important destination for immigrants to Canada. More than 40 percent of residents belong to a visible minority population group, and over 100 distinct ethnic origins are represented among its inhabitants. While the majority of Torontonians speak English as their primary language, over 100 languages are spoken in the city. Toronto is a prominent center for music, theatre, motion picture production, and television production, and is home to the headquarters of Canada's major national broadcast network and media outlets. Its varied cultural institutions, which include numerous museums and galleries, festival and public events, entertainment districts, national historic sites, and sports activities, attract over 100 million tourists each year. Toronto is known for its many skyscrapers and high rise buildings, in particular the tables free standing structure in the Western Hemisphere, the CN Tower.