

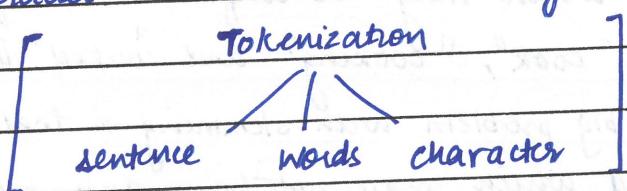


### Subject: Natural Language Processing

#### Module 2: Word level Analysis

Syllabus: Tokenization, Stemming, Segmentation, Lemmatization, Edit distance, collocations, Finite Automata, Finite state transducers, Porter Stemmer, Morphological analysis, Derivational and Inflectional morphology, Regular expressions and its types, N-Grams (Unigrams / Bigrams), language modelling, corpora, computing the probability of word sequence, Training and testing.

Tokenization: Tokenization is a common task of natural language processing. Tokenization is a way of separating a piece of text into smaller units called as Tokens. Here tokens can be either sentence, word or character. Tokenization is broadly classified into



Tokens are the building blocks of natural language. For example consider the sentence "Never give up". The most common way of forming tokens is based on space, assuming space as a delimiter. The tokenisation of above sentence results in 3 tokens - Never-Give-up.

Similarly tokens can be either characters or sub words for eg:-

Character tokens: S-m-a-r-t-e-r

Subword tokens: Smart-er.

→ Q) Why is tokenization required?

- :] Every sentence gets its meaning by the words present in it. So by analyzing the words present in the text we can easily interpret the meaning of the text.
- :] Once we have the list of words we can also use the standard tools and methods to get the insights into the text.



### Subject:Natural Language Processing

Stemming: Stemming usually refers to the crude heuristic process that chops off the end of words in hope of achieving the goal to reduce inflectional forms and sometimes derivationally related forms of word to a common base form.

With stemming words are reduced to their word stems. Stemming algorithms are typically rule based, one can view them as a heuristic process that sort of lops off the ends of words.

For example: we may have a suffix rule that based on a list of known suffixes cuts them off, In English language we have suffixes like "ed" and "-ing" which may be useful to cut off in order to map the words "cook", "Cooking" and "cooked" to the same stem of cook. The big problem with stemming is that it produces the root word which may not have any meaning.

Over stemming comes from when too much of the word is cut off. This can result in non-sensical stems where all the meaning of the word is lost/or muddled. Or it can result in words being resolved to the same stems, even though they probably should not be. Understemming is the opposite issue it comes from when we have several words that actually are the forms of one another.

Lemmatization ① It is the process of grouping together the different inflected forms of a word so that they can be analysed as a single item.

② Lemmatization usually refers to doing things normally with the use of vocabulary and morphological analysis of words aiming to remove inflectional endings and return the base or dictionary form of the word; also known as lemma.

Challenges in Stemming



### Subject:Natural Language Processing

Lemmatization is a more calculated process, it involves resolving words to their dictionary form. In fact lema of the word is its dictionary or canonical form. For lemmatization it requires an extra computational linguistic power such as part of speech. It is seen with a token 'saw' stemming might just return 's' however lemmatization would return see or saw depending on whether the use of token as a verb or a noun.

Stemming	Lemmatization
1) Stemming is faster because it chops words without knowing the context of word in a given sentence.	1] Lemmatization is slower as compared to stemming because it knows the context of word in a given sentence .
2) It is a rule based approach	2] It is a dictionary based approach
3] When we convert any word to its root form then stemming may create non existence meaning of word.	3] Lemmatization is more a calculated process, it involves resolving words to their dictionary form.
4) Accuracy is less	4] Accuracy is more
5) Stemming is preferred when the meaning of word is not important example: spam detection.	5] Lemmatization is preferred when the meaning of the word is important for analysis; example: Question answering system.
6] for example : "Studies" $\Rightarrow$ "Studi"	6] for example : 'Studies' $\Rightarrow$ "study".

Segmentation: ① segmentation refers to dividing text into larger, coherent units such as sentences, paragraphs or even topics.

Purpose: ① To structure text into meaningful sections that preserve the flow and context of original content.  
② To facilitate higher level text processing tasks.



Subject: Natural Language Processing

granularity

Segmentation typically operates at sentence level, paragraph level or document level.

Examples of Sentence segmentation:

Input : "Mr Smith bought a car. Its new and shiny."

Output: ["Mr Smith bought a car", "Its new and shiny."]

Paragraph segmentation

Input : "This is the first paragraph. It talks about introduction. This is the second paragraph. It talks about details."

Output: ["This is the first paragraph. It talks about introduction", "This is a second paragraph. It talks about details"].

Challenges of segmentation :

1] correctly identifying sentence boundaries, especially when dealing with abbreviations, punctuations and quotes .

2] Determining paragraph boundaries in documents without explicit markers.

3] Segmenting text into topics in a way that each segment is coherent and contextually relevant.

Tools:

NLTK

Spacy

Differences

Stanford NLP.

Tokenisation

Segmentation

1] Tokenisation focuses on breaking text into meaningful units.

2] Tokenisation focuses on finer

granularity.

3] Tokenisation is foundation step for nlp tasks such as parsing & tagging.

1] Segmentation focuses on breaking text into larger coherent units.

2] Segmentation deals with coarser granularity.

3] Segmentation is essential for understanding structure and flow of large text bodies.



Subject: Natural Language Processing

### Edit Distance [MU 2023 Dec] - (Important) —

Edit distance also known as Levenshtein distance is the measure of similarity between the two strings. It quantifies how many character edits (insertions, deletions or substitutions) are required to transform one string to another.

key concepts:

1. Insertions: Adding a character to the string.
2. Deletions: Removing a character from the string
3. Substitutions: Replacing one character with another.

Example:

Consider the transformation of word kitten to sitting.

kitten → sitten (substitute 'k' with 's')

sitten → sittin ( substitute 'i' with 'e' )

sittin → sitting ('insertion of 'g').

The edit distance is '3' [two substitutions and one insertions].

Applications of Edit distance:-

1. Spell checking: Find the closest correct word to a misspelled word.
2. Approximate string matching: searching for a pattern in a text where the pattern may have errors.
3. Plagiarism Detection: comparing documents to find similarity.
4. DNA sequence analysis: comparing genetic sequences for similarity and differences.
5. Machine Translation Evaluation: comparing translated sentences to reference translations to measure translation quality.

Collocation: Collocations in NLP refer to the combinations of words that frequently occur together and exhibit a linguistic pattern or a relationship. These words pairs or groups often convey a particular meaning or context when used together which may not be evident when words are used separately.



Subject:Natural Language Processing

Examples of collocations:

1. Adjective + noun: "strong tea", "heavy rain"
2. Verb + noun: "make a decision", "catch a bus".
3. Noun + noun: "data analysis", "coffee machine"
4. Adverb + adjective: "deeply concerned", "highly effective"
5. Verb + adverb: "run quickly", "speak softly".

Importance of collocations in NLP:

1. Natural language understanding: Recognizing collocations helps in understanding the context
2. Language modelling: Collocations help to improve the performance of language models.
3. Machine Translations: Translating collocations is essential for preserving meaning across languages.
4. Text Generation: Generating natural sounding text requires an understanding of common word pairing.

Methods to identify collocations:-

1. Frequency based methods: ① Simply counting the frequency of word pairs occurring together in the corpus.  
Advantages: 1. Easy to implement  
Disadvantages: 2. High frequency alone does not account for significance common words may co-occur frequently by chance.
2. Statistical methods: Methods such as chi-square test, point wise mutual information, likelihood ratio test are used.



3. Machine learning approaches: Using supervised and unsupervised learning to identify collocations. Word embeddings and neural networks are extensively used.

Regular Expressions: R.E is a language for specifying text search strings or search patterns. The regular expression languages are used for searching texts in UNIX. Usually search patterns are used by searching algorithms for find or find/replace operations on strings or for input validation. A regular expression is a formula in a special language that is used for specifying simple class of strings. Formally, a regular expression is an algebraic notation for characterizing the set of strings. Basically a regular expression is a pattern describing a certain amount of text.

→ Sums on Regular Expressions

$0^* = \{ \text{E, } 0, \text{ } 00, \text{ } 000, \text{ } 0000 \} \rightarrow \text{zero or more occurrences}$

$0^+ = \{ 0, \text{ } 00, \text{ } 000, \text{ } 0000 \} \rightarrow \text{one or more occurrences of } 0.$

: Write a R.E for a language accepting all combinations of a.  $\Rightarrow a^* \rightarrow \{ \text{aa, } \text{aaa, } \text{aaaa, } \dots \}$

: Design a R.E for the language accepting all strings having any number of a's and b's  $\Rightarrow \{ a+b \}^*$

: Construct an R.E for the language accepting all strings which are ending with 00 over input 0,1  $(0|1)^* 00.$

: Input {a,b,c} in which any number of a's followed by any number of b's and c's  $- a^* b^* c^*.$

: Construct R.E for language which consist of atleast 2 b's over the string a,b. zero b  $\rightarrow (a)^* + \epsilon$ , one b  $\rightarrow (a)^* + b$ , two b  $\rightarrow (a)^* + bb.$

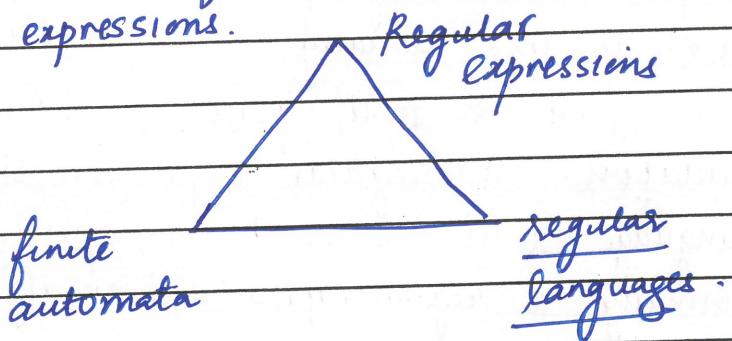
$$[(a^* + \epsilon) + (a^* + b) + (a^* + bb)].$$

$$a^* [\epsilon + b + bb].$$



### Finite Automata:-

The regular expression is more than just a convenient metalinguage for text searching. First a regular expression is one way of describing a finite state automaton (FSA). Both regular expressions and finite state automata can be used to describe regular expressions.



Finite automata (FA) also known as FSA are computational models used to represent and manipulate a set of states and transitions. In context of natural language processing, finite automata are employed to model and process sequences of characters, words or other linguistic units.

### Key concepts of finite automata:-

1. **States:** Finite automata consist of finite number of states, one of these states are initial state and one or more can be designated as final state.
2. **Alphabets:** The set of symbols that the automaton can read. In NLP, this could be characters, word or tokens.
3. **Transitions:** Rules that describe how to move from one state to the other.
4. **DFA [Deterministic finite automaton]:** In DFA for each state for each input symbol there is exactly one transition.
5. **NFA [Non deterministic finite automata]:** In NFA for each state for each input symbol there can be zero, one or more multiple transitions.



Subject: Natural Language Processing

Components of finite automaton:-

A finite automaton is defined by a 5 tuple  $(Q, \Sigma, \delta, q_0, F)$

Q: A finite set of states

$\Sigma$ : A finite set of alphabets

$\delta$ : A transition function  $\delta: Q \times \Sigma \rightarrow Q$  for DFA and  $\delta: Q \times \Sigma \rightarrow 2^Q$  for NFA.

$q_0$ : initial state

F: The set of final (accepting) states.

FST [Finite state transducers] extend the concept of finite automata by associating input symbols with output symbols, making them powerful tools for various tasks in natural language processing.

FSTs are particularly useful for modelling and implementing transformations in text processing such as morphological analysis, phonological rules and text to speech systems.

A FST is defined as a 6-tuple  $(Q, \Sigma, \Gamma, \delta, q_0, F)$

Q: finite set of states

$\Sigma$ : A finite set of ip alphabets

$\Gamma$ : A finite set of op symbols.

$\delta$ : A transition function  $Q \times \Sigma \rightarrow Q \times \Gamma$

$q_0$ : initial state

F: The set of final (accepting) states.

FST as a recognizer, FST as a generator, FST as a translator, FST as a set relator, these are roles of finite state transducers.



### Subject: Natural Language Processing

Morphological parsing with finite state transducers:

- The objective of morphological parsing is to produce output lexicons for a single input lexicon
  - The term morphological parsing is related to parsing of morphemes.
  - We define morphological parsing as the problem of recognizing that a word breaks down into smaller meaningful units called as morphemes producing some sort of linguistic structure for it.
  - For example we can break the word foxes into two fox and es
- In other sense morphology is the study of:
- The formation of words
  - The origin of words
  - grammatical forms of words.

Types of morphemes:—

- Stems
- Word order ]

Stems: It is core meaningful unit of word. we can also say it is root of the word. for example in the word foxes the stem is fox.

Affixes: As the name suggest they add some additional meaning and grammatical functions to the words.

Affixes can be divided into four types:—

- Prefixes — Prefixes precede the stem, for example in the word unbuckle (un → prefix)

Suffixes — suffixes follow the stem for example in the word cats , -s is suffix.



Subject: Natural Language Processing

Morphological parsing with FST: →

The FST is a multi function device and can be viewed in following ways:-

- Translator: It reads one string on one tape and outputs another string
- Recognizer: It takes a pair of strings as two tapes, and accepts/rejects based on matching.
- Generator: It outputs a pair of strings on two tapes along with yes/no result.
- Relator: It compares the relations between two strings

Input	Parsed output	Lexical	surface
cat	cat+N+SG	c   a   t	n + pl
cats	cat+N+PL	c   a   t   s	
geese	goose+N+PL		
reading	read+N+present part		

FST-1 translates between lexical and intermediate level.

Lexical      f | o | x | + N | + Pl

Intermediate      f | o | x | ^ | s | #

Orthographic rules for FST. — Language English often requires spelling changes at morpheme boundaries by introducing spelling rules (or orthographic rules)

Name	Description of rule	Example
consonant doubling	I letter consonant doubled	beg/begging
E deletion	silent 'e' dropped before -ing	make/making
E inserbing	e added after -s, -z, -x, -ch	watch/watches.
y replacement	y changes to ie before	try/tres
K insertion	verb ending with vowel + c add K	panic/panicked.



### Subject:Natural Language Processing

Porter's Stemming Algorithm [Imp, MU'23 (May, Dec), MU'22 (May)]

- The porter's stemming Algorithm is used to remove the suffixes from an English word and obtain its stem which is very useful in field of NLP
- This process reduces the number of terms kept by an NLP system and will be very advantageous.
- A consonant is a letter other than vowel and other than 'y' preceded by a consonant.
- If a letter is not a consonant it is a vowel
- A consonant is denoted by letter 'c' and a vowel is denoted by 'v'.
- Any word or a part of word therefore has one of the four forms given below :-
  - cvcv ... c → management, collection
  - cvcv ... v → conclude, revise
  - vcvc ... c → entertainment, illumination
- All of these forms can be represented using a single form as  $[c]vcvc \dots [v]$

Here square brackets denote the arbitrary presence of consonants or vowels.

$(vc)^m$  → denotes vc repeated m times so we can have expression written as

$[c] (vc)^m [v]$

- The value 'm' found in the above expression is called as the measure of any word or word part when represented in the form  $[c](vc)^m[v]$ . Some examples for different values of m are:



Subject: Natural Language Processing

<b>TREE</b> $\downarrow \downarrow \downarrow \downarrow$ $CC\ VV$ $[c][cv] [v]$ $m=0$	<b>T R O U B L E</b> $\downarrow \uparrow \downarrow \downarrow \downarrow \downarrow \downarrow$ $c\ c\ v\ v\ c\ c\ v$ $[c][vc][v]$ $m=1$
--	--

The rules for replacing a suffix will be given in the following form:—

$$(condition) S1 \xrightarrow{ } S2$$

This means that if a word ends with  $S_1$  and stem before  $S_1$  satisfies the given condition  $S_1$  is replaced by  $S_2$ .  
 This condition is usually given in terms of ' $m$ ' in regards to  $S_1$  ( $m > 1$ )ement  $\rightarrow S_1$  is ement and  $S_2$  is null.  
 This would map REPLACEMENT to REPLACE, since REPLACE is a word part for which  $m=2$ .

Conditions:

1.  $m \rightarrow$  measure of the stem,  $m=k$  or  $m>k$  where  $k$  is an integer.
2.  $*X \leftarrow$  the stem ends with a given letter  $X$ .
3.  $*V \leftarrow$  The stem contains a vowel.
4.  $*d \leftarrow$  The stem ends in a double consonant.
5.  $*o \leftarrow$  The stem ends in a consonant-vowel-consonant.
- Rules are divided into sets and in each successive step one set of rules are applied.



### Subject: Natural Language Processing

#### Porter steps:

- Each step corresponds to the set of rules. The rules in the step are examined in a sequence and only one rule from a step can apply.
  - Step 1a (word);
  - Step 1b (stem);  
 if (step second or third rule of step 1b was used)
    - Step 1b1 (stem);
    - Step 1c (stem);
    - Step 2 (stem);
    - Step 3 (stem);
    - Step 4 (stem);
    - Step 5a (stem);
    - Step 5b (stem);

Y

#### Porter rules:

##### Step 1a (rules)

conditions	suffix	Replacement	Examples
NULL	sses	ss	care sses → caress
NULL	ies	i	pones → poni
NULL	ss	ss	carress → carres
NULL	s	NULL	cats → cat

##### Step 1b (rules)

conditions	suffix	Replacement	Examples
(m>o)	eed	ee	fed → fed
(+v+)	ed	Null	plastered → plaster
(+v+)	ing	Null	motoring → motor



Subject: Natural Language Processing

Step 2 (rules).

Conditions	Suffix	Replacement	Examples
(m>0)	ational	ate	relationat → relati
(m>0)	enci	ence	valenci → valence
(m>0)	tional	tion	conditional → condition
(m>0)	izer	ize	digitiz → digitized
(m>0)	alism	al	fudalism → fudal
(m>0)	senti	ve	sensitivity → sensitive

+++  
(Refer ppt for more  
rules of step 2)

Step 3, Step 4, (Refer ppt for more rules of step 2, step 3, step 4).

Conditions	Suffix	Replacement	Examples
(m>0)	icate	ic	triplicate → tripli
(m>0)	ative	Null	formativ → form
(m>0)	alize	al	formalize → formal

Step 4 (Refer ppt for rules 4).

Conditions	Suffix	Replacement	Examples
(m>1)	al	null	revival → reviv
(m>1)	anc	null	allowanc → allow

Step 5 (Refer ppt for more rules 5).

Conditions	Suffix	Replacement	Examples
(m>1)	e	null	probate ↕ probat
(m>1)	e	null	cease → ceas



### Subject: Natural Language Processing

Example:

computers  $\xrightarrow{\text{step 1a}}$  computer  $\xrightarrow{\text{step 4}}$  comput  
 singing  $\xrightarrow{\text{step 1b}}$  sing.

Porter mishaps:-

gas (noun)  $\rightarrow$  ga

gases (plural)  $\rightarrow$  gase

gasses (verb, present tense)  $\rightarrow$  gass

gaseous (adjective)  $\rightarrow$  gaseou

\* This is not good — all these words should ideally reduce to the same stem.

Tradeoff: More rules (accurate but slow) vs Less rules (efficient but sometimes wrong).

\* Google does give different results for gas and gases, so maybe they use these porter rules...

Morphology: (Imp).

- Morphology is the study of <sup>way</sup> words are build up from smaller meaning bearing units, morphemes
- A morpheme is defined as minimal meaning bearing unit in a language.

For example the word 'fox' consist of a single morpheme while the word 'cats' consist of two morphemes morpheme 'cat' and morpheme 's'.

- Morphemes are divided into two types

Bound  
morphemes

Free  
morphemes



### Subject:Natural Language Processing

**Bound morphemes:** These are ones that cannot appear as a word by itself. eg) -s (dogs), (quickly)-ly, -ed (walked).

**Free morphemes:** These can appear as words by themselves... and can often combine with other morphemes.

eg) house (house-s), walk (walk-ed), of, or ...

- Two broad classes of morphemes are
  - stem
  - Affixes

stem → main morpheme of the word.

Affixes → add additional meanings of various kinds

Affixes are further classified into

→ prefixes

→ suffixes

→ circumfixes

→ infixes.

prefixes → precede the stem.

suffixes → follow the stem

circumfixes → precede as well as follow the stem

infixes → inserted inside the stem.

**Types of morphology (V.V.Imp) –**

- Morphology can be classified into two types

Inflectional

Derivational

morphology

morphology .

**Inflectional morphology:** It is a combination of a word stem with a grammatical morpheme which results in a word of same class as original stem. In English inflection is simple only nouns, verbs, adjectives can be inflected. Eg cat → cats, mouse → mice, walk → walking etc.



**Derivational morphology:** It is a combination of a word stem with a grammatical morpheme which results in a word of different class. In English it is very hard to predict the meaning of stem from its derived structure.  
Eg appoint → appointee clue → clueless, kill → killer etc

### Inflectional Morphology

1. Inflection may or may not be affected by affixation, reduplication, modification of bases and other morphological process

2. Inflection prototypically serves to modify the lexemes to fit a different grammatical context.

3. Inflection typically adds grammatical information

4. Inflection prefix is to change the category but does not add a new meaning eg ~~un~~appreciate

### Derivational Morphology

1. Derivation may be affected by formal means like affixation, reduplication or modification of bases and other morphological process.

2. Derivation serves to create new lexemes.

3. Derivation changes the category (takes a verb like employee and make it a noun employment).

4. Derivation adds a substantial new meaning. ej) unhappy.



Subject: Natural Language Processing

N-gram model. (v.v. imp). — (IDM) [MU'23 (May, Dec), MU'22 (Dec)].

- In natural language processing n gram is a contiguous sequence of n items generated from a given sample of text
- A statistical language model is a probability distribution  $p(s)$  over all possible word sequences. The dominant approach in a statistical language modelling is the N-gram model.
- The goal of a statistical language model is to estimate the probability of a sentence.
- This is achieved by decomposing sentence probability into a product of conditional probabilities using chain rules as follows -  $P(s) = P(w_1, w_2, w_3 \dots w_n)$

$$\begin{aligned} &= P(w_1) \cdot P(w_2|w_1) \cdot P(w_3|w_2, w_1) \cdot P(w_4|w_3, w_2, w_1) \\ &= \prod_{i=1}^n P(w_i|h_i) \end{aligned}$$

where  $h_i \rightarrow$  history of word  $w_i$

- In order to calculate sentence probability we need to calculate the probability of word, given sequence of words preceding it.
- In n gram model simplifies the task by approximately calculating the probability of a word given all the previous words by conditional probability given previous  $(n-1)$  words only.

$$P(w_i|h_i) = P(w_i|w_{i-n+1}, \dots, w_{i-1})$$

n gram model calculates  $P(w_i|h_i)$  by modelling language as Markov model of order  $n-1$  i.e. by looking at previous  $n-1$  words only.

A model that limits the history to previous one word is termed as bigram model and previous two words is termed as trigram model.



Subject: Natural Language Processing

Using bigram and trigram estimate the probability of a sentence :-

$$P(S) = \prod_{l=1}^n P(w_l | w_{l-1})$$

$$P(S) = \prod_{l=1}^n P(w_l | w_{l-2} w_{l-1})$$

- . A special <s> tag is used to mark the beginning of the sentence in a bigram estimation.
- . The probability of first word in a sentence is a condition on <s>. Similarly in trigram estimation we introduce two pseudo words <s1> <s2>

$$\left[ P(w_l | w_{l-n+1} \dots w_{l-1}) = \frac{C(w_{l-n+1} \dots w_{l-1}, w_l)}{\sum_w C(w_{l-n+1} \dots w_{l-1}, w)} \right]$$

- . The sum of all n grams that share (n-1) words is equal to the count of common prefix,  $w_{l-n+1}$  hence the previous expression can be rewritten as follows:-

$$\left[ P(w_l | w_{l-n+1} \dots w_{l-1}) = \frac{C(w_{l-n+1}, \dots, w_{l-1}, w_l)}{\sum_w C(w_{l-n+1}, \dots, w_{l-1})} \right]$$



Subject: Natural Language Processing

Numericals on Ngrams (Imp)  $\star\star$  [ MU22'(May, Dec), MU23'(May, Dec) ]

- ① Find the probability of : The Arabian knights are fairy tales of east.  
Training set:

<S> The Arabian knights

<S> These are the fairy tales of the east

<S> The stories of the Arabian knights are translated in many languages.

$$P(\text{The}|\text{<S>}) = 2/3 = 0.67$$

$$P(\text{are}|\text{knights}) = 1/2 = 0.5$$

$$P(\text{Arabian}|\text{the}) = 2/5 = 0.4$$

$$P(\text{translated}|\text{are}) = 1/2 = 0.5$$

$$P(\text{knights}|\text{Arabian}) = 2/2 = 1.0$$

$$P(\text{in}|\text{translated}) = 1/1 = 1$$

$$P(\text{are}|\text{these}) = 1/1 = 1$$

$$P(\text{many}|\text{in}) = 1/1 = 1$$

$$P(\text{fairy}|\text{the}) = 1/5 = 0.2$$

$$P(\text{languages}|\text{many}) = 1/1 = 1$$

$$P(\text{tales}|\text{fairy}) = 1/1 = 1$$

$$P(\text{the}|\text{of}) = 2/2 = 1$$

$$P(\text{east}|\text{the}) = 1/5 = 0.2$$

$$P(\text{stones}|\text{the}) = 1/5 = 0.2$$

$$P(\text{of}|\text{stones}) = 1/1 = 1$$

Test sentence: The Arabian knights are the fairy tales of the east

$$P(\text{the}|\text{<S>}) \times P(\text{Arabian}|\text{the}) \times P(\text{knights}|\text{Arabian}) \times P(\text{are}|\text{knights}) \\ \times P(\text{the}|\text{are}) \times P(\text{fairy}|\text{the}) \times P(\text{tales}|\text{fairy}) \times P(\text{of}|\text{tales}) \times \\ P(\text{the}|\text{of}) \times P(\text{east}|\text{the}).$$

$$= 0.67 \times 0.5 \times 1 \times 0.5 \times 0.2 \times 1 \times 1 \times 1 \times 2 = \underline{\underline{0.0067}}$$



Subject:Natural Language Processing

Example on Naram model.

$\langle S \rangle I \text{ am Sam} \langle |S \rangle$

$$(do|\langle S \rangle) = 1/5 = 0.2$$

$\langle S \rangle Sam I \text{ am} \langle |S \rangle$

$$(I|do) = 1/2 = 0.5$$

$\langle S \rangle Sam I \text{ like} \langle |S \rangle$

$$(like|I) = 2/5 = 0.4$$

$\langle S \rangle Sam I \text{ do like} \langle |S \rangle$

$$(Sam|like) = 1/3 = 0.33$$

$\langle S \rangle do I \text{ like Sam} \langle |S \rangle$

$$(\langle S \rangle |Sam) = 2/5 = 0.4$$

$$(I|\langle S \rangle) = 1/5 = 0.2$$

$$(am|I) = 2/5 = 0.4$$

$$(Sam|am) = 1/2 = 0.5$$

$$(\langle S \rangle |Sam) = 2/5 = 0.4$$

$$(Sam|\langle S \rangle) = 3/5 = 0.6$$

$$(I|Sam) = 3/5 = 0.6$$

$$(am|I) = 2/5 = 0.4$$

$$(\langle S \rangle |am) = 1/2 = 0.5$$

$$(Sam|\langle S \rangle) = 3/5 = 0.6$$

$$(I|Sam) = 3/5 = 0.6$$

$$(like|I) = 2/5 = 0.4$$

$$(\langle S \rangle |like) = 2/3 = 0.66$$

$$(Sam|\langle S \rangle) = 3/5 = 0.6$$

$$(I|Sam) = 3/5 = 0.6$$

$$(do|I) = 1/5 = 0.2$$

$$(like|do) = 1/2 = 0.5$$

$$(\langle S \rangle |like) = 2/3 = 0.66$$

Prediction of the next word

1]  $\langle S \rangle Sam$  —

$(Sam|\langle S \rangle)$  has the probability of 0.6, the next highest probability of bigram having Sam on right hand side is  $(I|Sam) = 0.6$ ,  $\therefore$  The next predicted word after Sam will be I.

2]  $\langle S \rangle Sam I do$

$(do|I)$  have the probability 0.2, the next highest probability of bigram having do on right hand side is  $(I|do) = 0.5$  and  $(like|do) = 0.5$  so next predicted word is like or I.



41

Subject: Natural Language Processing

Numerical (3) on Ngram (Mu 2023, Dec) given the following corpus find the  
for probability of occurrence of test data  $\langle \text{S} \rangle$  Students are from Vellore  $\langle \text{S} \rangle$   
Corpus:

$\langle \text{S} \rangle$  I am from Vellore  $\langle \text{S} \rangle$

$\langle \text{S} \rangle$  I am a teacher  $\langle \text{S} \rangle$

$\langle \text{S} \rangle$  Students are good and are from various cities  $\langle \text{S} \rangle$

$\langle \text{S} \rangle$  Students from Vellore do engineering  $\langle \text{S} \rangle$ .

$$P(\text{student}|\langle \text{S} \rangle)^* P(\text{are}|\text{student})^* P(\text{from}|\text{are})^* P(\text{vellore}|\text{from})^* P(\langle \text{S} \rangle|\text{vellore})$$

$$P(\text{student}|\langle \text{S} \rangle) = \frac{\text{count}(\langle \text{S} \rangle \text{ student})}{\text{count}(\langle \text{S} \rangle)} = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{are}|\text{student}) = \frac{\text{count}(\text{student are})}{\text{count}(\text{student})} = \frac{1}{2}$$

$$P(\text{from}|\text{are}) = \frac{\text{count}(\text{are from})}{\text{count}(\text{are})} = \frac{1}{2}$$

$$P(\text{vellore}|\text{from}) = \frac{\text{count}(\text{from vellore})}{\text{count}(\text{from})} = \frac{2}{3}$$

$$P(\langle \text{S} \rangle|\text{vellore}) = \frac{\text{count}(\text{vellore}|\langle \text{S} \rangle)}{\text{count}(\text{vellore})} = \frac{1}{2}$$

$$P(\text{Test data}) = \frac{1}{2} * \frac{1}{2} * \frac{1}{2} * \frac{2}{3} * \frac{1}{2} = 0.0416.$$

Numerical (4) on Ngram: Corpus: The girl bought a chocolate  
The boy ate a chocolate

The girl bought a toy

(choose the option that fits) The girl played with the toy.

Input: The girl \_\_\_\_\_



Subject: Natural Language Processing

a) The girl bought

b) The girl played

out of these ② which has the highest probability.

$$P(\text{bought} \mid \text{girl}) = \frac{\text{count}(\text{The girl bought})}{\text{count}(\text{The girl})} = \frac{2}{3} = 0.67$$

$$P(\text{played} \mid \text{the girl}) = \frac{\text{count}(\text{The girl played})}{\text{count}(\text{the girl})} = \frac{1}{3} = 0.33$$

probability of  $P(\text{bought} \mid \text{the girl}) >$  probability of  $P(\text{played} \mid \text{the girl})$

$0.67 > 0.33$

$\therefore$  option a) The girl bought fills the blank.

List of previous years university questions —

- Q1. Explain in brief Inflectional and Derivational morphology with suitable examples?
- Q2. Explain Porter's Stemming Algorithm?
- Q3. Explain Edit distance algorithm with an example? Show the working of minimum number of operations required to transform "kitten" to sitting?
- Q4. Explain difference between Stemming and Lemmatization?
- Q5. Consider the following corpus:

<S> She asks you to wait patiently </S>

<S> He wants me to help him </S>

<S> They expect us to arrive early </S>

List all possible bigrams. Compute the conditional probabilities and predict the next word to -