# NLP
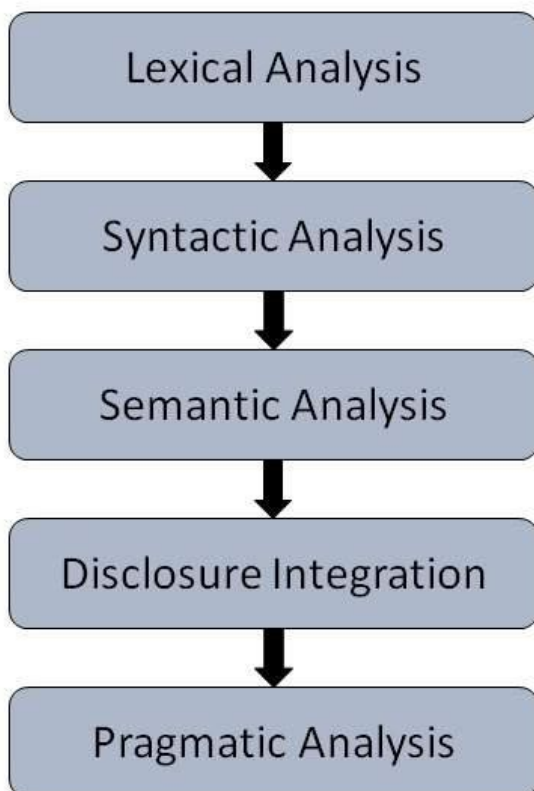
## Module 1

### 1. Explain about the phases of Natural Language processing ?

There are the following five phases of NLP:

Lexical Analysis − It involves identifying and analyzing the structure of words. Lexicon of a language means the collection of words and phrases in a language. Lexical analysis is dividing the whole chunk of txt into paragraphs, sentences, and words.

Syntactic Analysis (Parsing) − It involves analysis of words in the sentence for grammar and arranging words in a manner that shows the relationship among the words. The sentence such as "The school goes to boy" is rejected by English syntactic analyzer.



Semantic Analysis − It draws the exact meaning or the dictionary meaning from the text. The text is checked for meaningfulness. It is done by mapping syntactic structures and objects in the task domain. The semantic analyzer disregards sentence such as "hot ice-cream".

Discourse Integration − The meaning of any sentence depends upon the meaning of the sentence just before it. In addition, it also brings about the meaning of immediately succeeding sentence.

Pragmatic Analysis − During this, what was said is re-interpreted on what it actually meant. It involves deriving those aspects of language which require real world knowledge.

**2. What is Natural language processing? Explain generic NLP system?**

Natural Language Processing (NLP) refers to AI method of communicating with an intelligent systems using a natural language such as English.

Processing of Natural Language is required when you want an intelligent system like robot to perform as per your instructions, when you want to hear decision from a dialogue based clinical expert system, etc.

The field of NLP involves making computers to perform useful tasks with the natural languages humans use. The input and output of an NLP system can be −

- Speech
- Written Text

A generic NLP system consists of several components that work together to process and understand human language. Here's an overview of the key components of a generic NLP system:

- Input Text: The NLP system takes raw human language text as input. This text can be in the form of sentences, paragraphs, or entire documents.

- Tokenization: The input text is broken down into smaller units called tokens. Tokens can be words, phrases, or even characters. Tokenization is the initial step that divides the text into manageable units for further processing.

- Text Preprocessing: Text preprocessing involves cleaning and transforming the tokenized text. This may include converting text to lowercase, removing punctuation, special characters, and unwanted symbols, as well as handling issues like misspellings and contractions.

- Part-of-Speech Tagging: Part-of-speech tagging assigns grammatical labels (such as nouns, verbs, adjectives) to each token in the text. This helps in understanding the syntactic structure and grammatical relationships within the text.

- Parsing and Syntax Analysis: Parsing involves analyzing the grammatical structure of sentences to understand the relationships between words. This process helps in creating a structured representation of the text's syntax, such as a parse tree or a dependency graph.

- Named Entity Recognition (NER): NER identifies and classifies named entities in the text, such as people, places, organizations, and dates. This is crucial for extracting specific information and understanding the context.

- Sentiment Analysis: Sentiment analysis determines the emotional tone or sentiment expressed in the text. It can classify text as positive, negative, or neutral, providing insights into opinions and emotions.

- Feature Extraction: Feature extraction involves transforming the processed text into numerical or vector representations that machine learning models can work with. Techniques like word embeddings or other numerical encodings are used for this purpose.

- Machine Learning or Deep Learning Models: These models are trained on the extracted features to perform various NLP tasks. They can range from traditional machine learning algorithms to advanced neural networks, depending on the complexity of the task.

- Language Modeling and Context Understanding: Language models learn patterns and relationships between words to understand context. They can predict the next word in a sentence or generate coherent text based on a given input.

- Output Generation: The NLP system generates outputs based on the task. For instance, it might produce classifications, translations, summarizations, or even generate text in a conversational manner.

- Post-processing: The final output from the NLP system often undergoes additional processing to enhance its readability and coherence. This might involve reformatting, grammar correction, or smoothing out generated text.

## 3. What are different ambiguities which needs to handle by natural language processing ?

Natural language is inherently ambiguous, and resolving these ambiguities is a significant challenge in Natural Language Processing (NLP). Here are some common types of ambiguities that NLP systems need to handle:

- Lexical Ambiguity: This type of ambiguity arises when a word has multiple meanings. For example, the word "bank" can refer to a financial institution or the side of a river. Context is crucial in determining the correct meaning.

- Syntactic Ambiguity: Syntactic ambiguity occurs when the structure of a sentence allows for multiple interpretations. For instance, in the sentence "I saw the man with the telescope," it's unclear whether "with the telescope" modifies "I" or "the man."

- Semantic Ambiguity: Semantic ambiguity arises when a sentence can be interpreted in different ways due to the ambiguity of a word or phrase. For instance, the sentence "Time flies like an arrow; fruit flies like a banana" has different meanings for "flies."

- Anaphoric Ambiguity: Anaphoric ambiguity occurs when a pronoun refers back to a previous noun in a way that's unclear. For example, in the sentence "She saw her brother," it's unclear who "her" refers to.

- Referential Ambiguity: Referential ambiguity arises when it's unclear which entity a noun phrase refers to. In the sentence "The cat chased the dog," it's not clear which cat and which dog are being referred to.

- Pragmatic Ambiguity: Pragmatic ambiguity is related to the context and involves how language users interpret implied meaning. Sarcasm, irony, and indirect speech can all introduce pragmatic ambiguity.

## 4. Why is handling ambiguities important in Natural Language Processing applications?

Handling ambiguities is crucial in Natural Language Processing (NLP) applications for several reasons:

- Enhanced Understanding: Ambiguities in language can lead to different interpretations of the same text. Resolving these ambiguities accurately improves the NLP system's understanding of the intended meaning, allowing for more accurate and meaningful interactions.

- Precision in Results: NLP tasks often involve extracting information, making classifications, or generating text. Ambiguities can lead to incorrect results or misinterpretations, which can affect the quality and precision of the outcomes.

- Contextual Relevance: Ambiguities can create confusion and result in text that lacks contextually relevant information. By resolving ambiguities, NLP systems can provide more contextually appropriate responses and insights.

- User Experience: In applications like chatbots or virtual assistants, resolving ambiguities leads to more coherent and relevant responses. This improves the user experience by reducing misunderstandings and frustration.

- Translation and Localization: Ambiguities in one language can lead to errors during translation. Resolving these ambiguities ensures that the translated content accurately conveys the intended meaning.

- Natural Communication: Human communication relies on understanding and resolving ambiguities through context, tone, and shared knowledge. NLP systems aim to replicate this ability to achieve more natural and human-like communication.

- Information Retrieval: Ambiguities can impact the retrieval of relevant information from large text corpora or databases. Accurate ambiguity resolution aids in retrieving more relevant results.

- Decision Making: In applications like sentiment analysis or opinion mining, ambiguous language can lead to incorrect assessments of emotions or opinions. Resolving ambiguities is essential for making informed decisions based on accurate sentiment analysis.

## 5. What are the components of NLP?
There are the following two components of NLP -

1. Natural Language Understanding (NLU)
Natural Language Understanding (NLU) helps the machine to understand and analyse human language by extracting the metadata from content such as concepts, entities, keywords, emotion, relations, and semantic roles.
NLU mainly used in Business applications to understand the customer's problem in both spoken and written language.
NLU involves the following tasks -
It is used to map the given input into useful representation.
It is used to analyze different aspects of the language.

2. Natural Language Generation (NLG)

It is the process of producing meaningful phrases and sentences in the form of natural language from some internal representation.
It involves −

- Text planning − It includes retrieving the relevant content from knowledge base.
- Sentence planning − It includes choosing required words, forming meaningful phrases, setting tone of the sentence.
- Text Realization − It is mapping sentence plan into sentence structure.

Note: The NLU is difficult than NLG.

Difference between NLU and NLG

| NLU | NLG |
| --- | --- |

| NLU is the process of reading and interpreting language. | NLG is the process of writing or generating language. |
|---|---|
| It produces non-linguistic outputs from natural language inputs. | It produces constructing natural language outputs from non-linguistic inputs. |

## 6. Explain the application of NLP?

- Language Translation: NLP enables automated language translation, allowing text to be seamlessly translated from one language to another. Services like Google Translate use NLP algorithms to understand the structure and context of sentences in different languages and produce accurate translations.

- Sentiment Analysis: Sentiment analysis uses NLP to determine the emotional tone of text, whether it's positive, negative, or neutral. This application is used to analyze social media posts, customer reviews, and feedback, providing businesses with insights into public sentiment about their products or services.

- Chatbots and Virtual Assistants: Chatbots and virtual assistants, like those found on websites or platforms like Siri and Alexa, use NLP to understand and respond to user inquiries in natural language. NLP allows these systems to provide helpful responses and perform tasks like scheduling appointments or answering questions.

- Text Summarization: NLP-driven text summarization automatically condenses lengthy documents, articles, or reports into concise summaries. This application is valuable for quickly understanding the key points of a document without having to read the entire content.

- Named Entity Recognition (NER): NER involves identifying and classifying named entities, such as names of people, places, organizations, dates, and more. NLP-powered NER is used in various contexts, including information extraction from legal documents, medical records, and news articles.

- Question Answering: NLP enables systems to comprehend questions asked in natural language and provide relevant answers by extracting information from large datasets. This technology is used in chatbots, search engines, and virtual assistants to answer user queries.

- Speech Recognition and Text-to-Speech (TTS): NLP-driven speech recognition converts spoken language into text, enabling applications like voice commands, transcription services, and voice assistants. On the other hand, NLP-powered text-to-speech technology transforms written text into natural-sounding speech, used in audiobooks, navigation systems, and accessibility features.

These applications showcase how NLP enhances communication, information processing, and user interaction across a wide range of industries, from language translation and sentiment analysis to conversational agents and accessibility technologies.

## 7. What are challenges of NLP ?

- Ambiguity: Natural language is often ambiguous due to words having multiple meanings and sentence structures allowing for various interpretations. Resolving these ambiguities accurately is a significant challenge in NLP, as incorrect interpretations can lead to incorrect results.
- Lack of Context: Understanding context is essential for accurate language comprehension. However, capturing context across sentences or larger contexts remains a challenge, especially in complex conversations or documents.

- Data Limitations: NLP models require large amounts of high-quality data to learn effectively. However, obtaining labeled data for specific tasks, especially for niche domains or languages, can be challenging and time-consuming.

- Out-of-Distribution Data: NLP models can struggle when encountering data that is significantly different from what they were trained on. Handling out-of-distribution or unseen data is crucial for robust performance.

- Domain Adaptation: NLP models trained in one domain might not perform well in another due to differences in vocabulary, language style, and context. Adapting models to new domains while retaining their performance is a challenge.

- Bias and Fairness: NLP models can inadvertently learn and propagate biases present in training data. Ensuring fairness and addressing biases in language generation, sentiment analysis, and other tasks is a complex ethical challenge.

- Understanding Nuances: Human language contains nuances, cultural references, sarcasm, and idioms that are challenging for NLP systems to understand accurately. Capturing and interpreting these subtleties is an ongoing challenge.

**8. Write short notes of Indian language processing ?**
Indian Language Processing (ILP) refers to the field of Natural Language Processing (NLP) that focuses on developing language technologies for the diverse languages spoken in India. Given the linguistic diversity in the country, ILP presents its own set of challenges and opportunities. Here are some key points about Indian Language Processing:

1. Linguistic Diversity: India is home to hundreds of languages from different language families, including Indo-Aryan, Dravidian, Austroasiatic, and Tibeto-Burman. This linguistic diversity poses challenges for ILP, as each language has unique grammatical structures, phonetics, and vocabulary.

2. Resource Availability:Many Indian languages lack comprehensive language resources like large text corpora, labeled datasets, and linguistic tools. The scarcity of such resources hampers the development of robust ILP systems.

3. Low-Resource Languages: Several Indian languages are considered low-resource languages due to limited digital content and linguistic documentation. Building effective NLP models for these languages is challenging due to the scarcity of data.

4. Script Variations: India uses various scripts, including Devanagari, Tamil, Bengali, and more. Handling multiple scripts and their variations is a challenge when designing ILP systems that can process text in different scripts.

5. Code-Mixing and Code-Switching: Many Indians use a mix of English and their native language while communicating. Developing models that can understand and process code-mixed or code-switched text is important for accurate ILP.

6. Language Specifics: Different languages have their own linguistic characteristics, such as agglutination, inflection, and word order. These specifics impact the design of morphological analyzers, part-of-speech taggers, and other language processing tools.

7. Machine Translation: Machine translation is crucial for bridging language barriers in India. However, translating between Indian languages presents unique challenges due to structural differences, idiomatic expressions, and word order variations.

8. Speech Processing: ILP also includes speech recognition and synthesis for Indian languages. Speech processing systems need to account for variations in pronunciation, accents, and phonetics.

9. NLP Applications: ILP is applied to various tasks such as sentiment analysis, text classification, named entity recognition, and information retrieval in Indian languages. These applications have significant implications for regional content analysis and social media sentiment tracking.

10. Research and Development: Despite the challenges, ILP research and development are actively advancing. Academics, researchers, and organizations are working to create language resources, develop NLP models, and address language-specific issues.

In summary, Indian Language Processing focuses on adapting NLP techniques to the linguistic and cultural diversity of India. The field addresses challenges related to linguistic variations, resource scarcity, and multilingual communication to enable effective language technologies for Indian languages.

# Module 2

## 1.What is morphological analysis ?

Morphological analysis is a key component of Natural Language Processing (NLP) that focuses on breaking down words into their smallest meaningful units, known as morphemes. Morphemes are the smallest linguistic units that carry meaning, such as prefixes, suffixes, and root words. Morphological analysis involves analyzing the structure and composition of words to understand their grammatical and semantic properties.

Here's how morphological analysis works:

1. Morphemes: A word can be composed of one or more morphemes. For example, the word "unhappiness" consists of three morphemes: "un-" (a prefix indicating negation), "happy" (the root word), and "-ness" (a suffix indicating a state or quality).

2. Affixes: Morphemes can be classified into affixes and roots. Affixes are added to a root word to create new words or modify their meanings. Prefixes come before the root, while suffixes come after it.

3. Lemmatization and Stemming: Morphological analysis often involves lemmatization and stemming. Lemmatization aims to reduce words to their base or dictionary form, while stemming involves removing suffixes to obtain the word's root form.

4. Part-of-Speech Tagging: Morphological analysis is closely related to part-of-speech tagging, where words are assigned grammatical categories such as nouns, verbs, adjectives, etc. This is because the morphological structure of a word often indicates its grammatical role.

5. Morphological Analyzers: Morphological analysis is performed using morphological analyzers, which are tools or algorithms that apply linguistic rules and patterns to break down words into their constituent morphemes. These analyzers can be rule-based or statistical, and they help in identifying the morphemes, affixes, and their grammatical functions.

Morphological analysis is important in NLP for various reasons:

- It aids in understanding the grammatical structure of words, which is crucial for syntactic analysis and language understanding.
- Morphological information is valuable in tasks like part-of-speech tagging, where the grammatical role of a word determines its tag.
- Morphological analysis is essential for languages with rich morphology, where words can have complex inflections and multiple forms.

Overall, morphological analysis is a foundational step in NLP, helping systems to accurately process and understand the structural and grammatical properties of words within a given language.

## 2. What is significance of FST in morphological analysis ?

Finite State Transducers (FSTs) play a significant role in morphological analysis, especially for languages with complex morphology. FSTs are computational models that operate on sequences of symbols and are widely used in linguistics and Natural Language Processing (NLP) for tasks like morphological analysis and generation. Here's the significance of FSTs in morphological analysis:

1. Efficient Representation: FSTs provide a compact and efficient way to represent the complex relationships between words, morphemes, and affixes in a language. This is particularly important for languages with rich inflectional and derivational morphology.

2. Rule-Based Approach: FSTs enable a rule-based approach to morphological analysis. Linguists and language experts can define rules that capture the morphological patterns of a language. These rules can include information about affixes, stems, and grammatical categories.

3. Morpheme Segmentation: FSTs can segment words into their constituent morphemes, which is crucial for understanding the grammatical structure of words. This segmentation aids in tasks like part-of-speech tagging, lemmatization, and syntactic parsing.

4. Morphological Generation: FSTs can be used not only for analysis but also for generation. Given a root and a set of affixes, an FST can generate all the valid inflected forms of a word, helping in language generation tasks.

5. Flexibility: FSTs can accommodate different levels of linguistic analysis, from simple stemming to more complex derivational and inflectional processes. This flexibility allows FSTs to adapt to the specific morphological characteristics of a language.

6. Language Maintenance: FSTs are particularly useful for languages that lack extensive language resources. Linguists can create FST-based morphological analyzers to process and analyze the morphology of these languages, aiding in linguistic research and language maintenance efforts.

7. Resource-Light Approach: FST-based morphological analyzers can be built with relatively small amounts of annotated data, making them suitable for languages with limited linguistic resources.

8. Customization: FSTs can be easily customized and extended as linguists gain more insights into a language's morphology. This adaptability is crucial for accurately capturing the nuances of a language's morphological structure.

In summary, FSTs are a powerful tool in morphological analysis due to their efficiency, rule-based approach, and flexibility. They enable the development of morphological analyzers and generators that can accurately process and generate words, making them valuable assets in language technology for languages with complex morphological systems.


**3. Explain Tokenization Stemming and lemmatization ?**
Tokenization, stemming, and lemmatization are all fundamental techniques in Natural Language Processing (NLP) that deal with processing and simplifying words for analysis. Here's an explanation of each:

1. Tokenization:
   Tokenization is the process of breaking down a text into smaller units called tokens. Tokens are typically words, phrases, or symbols that are meaningful for analysis. Tokenization is the first step in text processing and is crucial for various NLP tasks. For example, the sentence "Chatbots are amazing!" can be tokenized into the tokens "Chatbots," "are," "amazing," and "!". Tokenization helps to prepare text for further analysis, such as part-of-speech tagging or sentiment analysis.

2. Stemming:
   Stemming is a technique to reduce words to their base or root form by removing suffixes and sometimes prefixes. The goal is to reduce inflected or derived words to their common root to group similar words together. For example, the stemming process might convert "running," "runner," and "runs" to the stem "run."

While stemming can help in basic text processing and information retrieval, it can sometimes result in non-words or words that are no longer valid. Stemming algorithms apply rules heuristically without considering the context, leading to potential errors.

3. Lemmatization:
  Lemmatization is a more advanced technique than stemming. It also aims to reduce words to their base form, but it takes into account the word's part of speech and context to ensure that the resulting word is valid. For example, the lemma of "running," "runner," and "runs" would all be "run." Lemmatization relies on linguistic databases and rules to determine the correct base form, making it more accurate than stemming. It's particularly useful for tasks that require a meaningful analysis of the words, such as language translation or sentiment analysis.

In summary, tokenization breaks down text into meaningful units (tokens), stemming reduces words to their base form using heuristics, and lemmatization reduces words to their base form while considering linguistic rules and context. These techniques are crucial for preparing text for various NLP tasks and improving the efficiency and accuracy of text analysis.

**4. Compare and contrast Lemmatization and stemming ?**
Lemmatization and stemming are both techniques used in Natural Language Processing (NLP) to reduce words to their base forms, but they have differences in terms of accuracy, complexity, and the linguistic aspects they consider. Here's a comparison and contrast between lemmatization and stemming:

Similarities:
1. Base Form: Both lemmatization and stemming aim to convert words to their base or root form. This helps in grouping words with the same meaning together and reducing the dimensionality of the vocabulary.

2. Text Processing: Both techniques are used to preprocess text before further analysis in NLP tasks. This can include tasks like text classification, sentiment analysis, information retrieval, and more.

Differences:
1. Accuracy:
  - Lemmatization: Lemmatization is generally more accurate than stemming because it considers the word's part of speech and context before reducing it to the base form. It uses linguistic databases and rules to determine the appropriate lemma.
  - Stemming: Stemming is a heuristic process that applies simple rules to chop off prefixes and suffixes. While it's faster, it can lead to inaccurate results, as the generated stems might not always be valid words.

2. Linguistic Aspects:
  - Lemmatization: Lemmatization takes into account the grammatical aspect of words. It considers the word's part of speech and applies linguistic rules to ensure that the resulting lemma is a valid word.
  - Stemming: Stemming is rule-based and doesn't consider part of speech or linguistic rules. This can result in generated stems that are not actual words and might not accurately represent the original word's meaning.

3. Complexity:
  - Lemmatization: Lemmatization is more complex due to its consideration of part of speech and linguistic rules. It requires linguistic databases or knowledge resources to accurately identify the correct lemmas.
  - Stemming: Stemming is simpler and faster since it follows a set of predefined rules to truncate affixes. However, this simplicity comes at the cost of potentially generating non-words.

4. Use Cases:

- Lemmatization: Lemmatization is suitable for tasks that require accurate analysis and understanding of words, such as machine translation, sentiment analysis, and language generation.
- Stemming: Stemming is often used in information retrieval tasks like search engines, where speed and simplicity are prioritized over linguistic accuracy.

In essence, lemmatization is more accurate but computationally heavier, while stemming is faster but less accurate. The choice between the two depends on the specific NLP task, the importance of linguistic accuracy, and the resources available for processing.


## 5. Explanation inflectional and derivational morphology ?
Inflectional and derivational morphology are two fundamental aspects of the way languages modify words to convey different meanings or grammatical information. They are key components of linguistic analysis and play a significant role in understanding the structure of words in various languages.

1. Inflectional Morphology:
   Inflectional morphology involves altering a word's form to indicate grammatical features like tense, number, gender, case, and person. Inflections provide additional information about the word's role in a sentence without changing its core meaning or part of speech. Inflectional morphemes are typically suffixes that are added to the base form of a word.

   For example, consider the English verb "run." By adding the "-ing" suffix, we get "running," indicating the present participle form of the verb. Similarly, "walk" becomes "walked" in the past tense.

   Inflectional morphology does not create entirely new words or change the part of speech; it modifies words to convey grammatical information within the same word category.

2. Derivational Morphology:
   Derivational morphology involves the creation of new words by adding prefixes or suffixes to a base word. Unlike inflectional morphology, derivational morphology can change the word's meaning, part of speech, or both. It often involves more substantial modifications to the word's structure.

   For instance, consider the English noun "nation." By adding the suffix "-al," we form the adjective "national." Similarly, adding the prefix "un-" to the adjective "happy" creates the word "unhappy," changing the meaning to its opposite.

   Derivational morphology is responsible for word formation processes like forming adjectives from nouns, verbs from adjectives, and so on. It plays a crucial role in expanding a language's vocabulary and expressing nuanced meanings.

In summary, inflectional morphology deals with modifying words to convey grammatical information within the same part of speech, while derivational morphology involves creating new words by adding affixes that can change the word's meaning or part of speech. Both aspects contribute to the richness and complexity of languages' morphological systems.


## 6. What is significance of porter stemmer in NLP applications ?
The Porter Stemmer is a widely used algorithm for stemming words in Natural Language Processing (NLP) applications. Developed by Martin Porter, this algorithm aims to reduce words to their base or root form by removing common suffixes. The Porter Stemmer holds significance in various NLP applications for several reasons:

1. Text Normalization: In NLP, text normalization involves transforming words to a standard or canonical form. The Porter Stemmer helps achieve this by reducing inflected words to their base form, allowing systems to treat different forms of the same word as equivalent. This is particularly useful for information retrieval and search engines.

2. Reducing Dimensionality: Stemming helps reduce the dimensionality of text data. Similar words with different inflections are converted to the same stem, enabling more effective analysis and comparison across documents or texts.

3. Information Retrieval: In applications like search engines, stemming enhances the retrieval of relevant documents. Users may use different forms of a word while searching, and stemming ensures that variations of the same root word are matched.

4. Text Classification: In tasks like sentiment analysis and text categorization, stemming simplifies the vocabulary by collapsing words with common roots. This reduces the number of features, making classification models more efficient and robust.

5. Document Clustering: Stemming can help improve the quality of document clustering and topic modeling by grouping similar words together. This enhances the accuracy of identifying themes within a collection of documents.

6. Speed and Memory Efficiency: The Porter Stemmer is relatively fast and requires minimal memory resources. This efficiency makes it suitable for large-scale text processing tasks where computational resources are a concern.

7. Basic Linguistic Analysis: The Porter Stemmer can be used as a basic linguistic analysis tool in scenarios where more advanced techniques like lemmatization are not required. It's particularly useful for tasks that don't demand high linguistic accuracy.

8. Search Engine Optimization (SEO): In web content optimization, stemming can be used to identify relevant keywords and improve the visibility of content in search engine results.

However, it's important to note that the Porter Stemmer has limitations. It's a rule-based algorithm that doesn't consider linguistic context or part-of-speech information. This can result in over-stemming (removing too many characters) or under-stemming (not removing enough characters), leading to incorrect or non-words. In situations where high accuracy is required, lemmatization, which considers linguistic context, might be a better choice.

In summary, the Porter Stemmer is significant in NLP applications for its ability to simplify words, reduce dimensionality, and improve information retrieval and text analysis. It's especially useful in cases where speed and efficiency are important and where linguistic accuracy can be sacrificed to some extent.


**7. What is regular expression ? What is is significance in word level analysis ?**
A regular expression (regex or regexp) is a sequence of characters that defines a search pattern. It is a powerful tool for string manipulation, text searching, and pattern matching. Regular expressions provide a concise and flexible way to define and identify patterns within strings of text.

In word-level analysis, regular expressions are significant for several reasons:

1. Pattern Matching: Regular expressions allow you to define specific patterns in text, such as finding words that start with a certain prefix, end with a particular suffix, or match a specific sequence of characters. This is useful for tasks like extracting specific types of words or phrases from a larger text.

2. Text Extraction: Regular expressions are commonly used to extract information from unstructured text. For example, you can use regex to extract email addresses, phone numbers, dates, or other structured data from a document.

3. Tokenization: Regular expressions are employed in tokenization, which is the process of breaking text into smaller units (tokens). By defining patterns for word boundaries, punctuation, and whitespace, regular expressions help split text into meaningful tokens.

4. Text Cleaning: Regular expressions are valuable for cleaning and preprocessing text. You can use them to remove special characters, punctuation, or unwanted formatting from a text corpus, improving the quality of the data before analysis.

5. Stemming and Lemmatization: Regular expressions can aid in basic stemming or lemmatization tasks by identifying and manipulating word forms. While they might not provide the same linguistic accuracy as dedicated stemmers or lemmatizers, they can be helpful in certain scenarios.

6. Search and Retrieval: In information retrieval, regular expressions are used to search for specific patterns or keywords within a document or dataset. This is essential for tasks like full-text search and document retrieval.

7. Data Validation: Regular expressions are used for validating user input in applications such as web forms. They help ensure that data entered by users follows a specific pattern or format.

8. Language Analysis: Regular expressions can assist in identifying linguistic patterns, such as finding all occurrences of adjectives followed by nouns, or detecting certain grammatical constructions.

9. Named Entity Recognition (NER): Regular expressions can aid in basic NER tasks by identifying patterns associated with named entities like dates, locations, and organizations.

While regular expressions are powerful, they can also be complex and require careful crafting to match the desired patterns accurately. Additionally, for more advanced linguistic analysis tasks or handling complex linguistic phenomena, other techniques like part-of-speech tagging, lemmatization, and syntactic parsing may be needed.

**8. Difference between free and bound morpheme ?**
Free morphemes and bound morphemes are two fundamental concepts in linguistics that describe different types of morphemes in a language. Morphemes are the smallest units of meaning in a language. Here's the difference between free and bound morphemes:

1. Free Morphemes:
   - Definition: Free morphemes are morphemes that can stand alone as independent words and carry meaning by themselves.
   - Examples: In English, words like "cat," "run," "book," and "happy" are free morphemes. Each of these words can exist on its own and convey meaningful content.
   - Independence: Free morphemes do not need to be attached to other morphemes to convey their intended meaning.

2. Bound Morphemes:
   - Definition: Bound morphemes are morphemes that cannot stand alone as independent words and need to be attached to other morphemes to convey meaning.
   - Examples: In English, prefixes like "un-" (as in "undo") and suffixes like "-ed" (as in "walked") are bound morphemes. They cannot function as words by themselves; they modify the meaning of other words.
   - Attachment: Bound morphemes must be affixed (attached) to a free morpheme to create a complete word.

Examples to Illustrate the Difference:
- "Cat" (free morpheme): A standalone word with a specific meaning.
- "Cats" (free morpheme + free morpheme): Plural form created by adding a free morpheme.
- "Undo" (bound morpheme + free morpheme): The prefix "un-" (bound morpheme) is added to the verb "do" (free morpheme) to change its meaning.
- "Books" (free morpheme + bound morpheme): Plural form created by adding the bound morpheme "-s" to the noun "book."
- "Happiness" (free morpheme + bound morpheme): The noun "happy" (free morpheme) is modified by adding the bound morpheme "-ness" to convey a state or quality.

In summary, the distinction between free morphemes and bound morphemes lies in their ability to function independently as words. Free morphemes can stand alone and carry meaning, while bound morphemes need to be attached to other morphemes to convey meaningful content. The combination of both types of morphemes allows for the creation of a wide variety of words and expressions in a language.

**9.** Explain following Terms :
– Morpheme
– Morphotactics
– Orthographic rule

1. Morpheme:
   A morpheme is the smallest grammatical unit in a language that carries meaning. Morphemes are the building blocks of words and can be a whole word or a part of a word. They can't be further divided into smaller meaningful units without losing their meaning. Morphemes can be free (able to stand alone as words) or bound (attached to other morphemes to convey meaning). For example, in the word "unhappiness," "un-" and "happiness" are morphemes. "Un-" is a bound morpheme that indicates negation, while "happiness" is a free morpheme representing a state of being happy.

2. Morphotactics:
   Morphotactics refers to the rules and patterns that govern the arrangement and combination of morphemes within a word. It encompasses the way morphemes are ordered, attached, and combined to create words in a language. Morphotactics are language-specific and dictate the allowable sequences of morphemes to form grammatically correct words. For example, English has morphotactic rules that determine the order of prefixes and suffixes in words.

3. Orthographic Rule:
   Orthography refers to the conventional system of spelling and writing a language. Orthographic rules are the established conventions that dictate how words are spelled and written in a particular language. These rules encompass aspects like letter sequences, capitalization, punctuation, and more. Following orthographic rules ensures consistent and standardized written communication. For example, the orthographic rule in

English states that the first letter of a sentence should be capitalized, and question sentences should end with a question mark.

In summary, a morpheme is the smallest meaningful unit in a language, morphotactics are the rules governing how morphemes are combined, and orthographic rules define the correct way to spell and write words in a language. These concepts are essential in linguistics and language analysis to understand how languages structure words and communicate in written form.

**10. Explain n-gram model ?**
The n-gram model is a statistical language model used in Natural Language Processing (NLP) to predict the next word in a sequence of words based on the previous (n-1) words. It's a type of probabilistic model that captures the likelihood of seeing certain word sequences in a given language. The "n" in n-gram refers to the number of words in the sequence.

For example, let's consider an n-gram model with n=2, also known as a bigram model. Given the sentence "The cat is on the", the model would predict the next word by looking at the probability distribution of words that follow the sequence "The cat is on the." If "mat" occurs frequently after that sequence, the model might predict "mat" as the next word.

Here's how the n-gram model works:

1. Tokenization: The text is first tokenized into words or other meaningful units (e.g., characters or subwords) to create a sequence of tokens.

2. Building N-Grams: The text is then divided into overlapping sequences of n tokens. These sequences are the n-grams. For example, with n=3, the sentence "The cat is on the mat" would be split into trigrams: "The cat is," "cat is on," "is on the," "on the mat."

3. Frequency Counting: The frequency of each n-gram in the text is counted. This involves tallying how often specific sequences of words appear.

4. Calculating Probabilities: To predict the likelihood of a particular word following a sequence, the model calculates conditional probabilities based on the frequency of the n-gram and its subsequent words.

5. Prediction: When given a sequence of n-1 words, the model predicts the next word by selecting the most probable word according to the calculated probabilities.

N-gram models have several applications in NLP:

- Language Modeling: N-gram models are used to model the probabilities of word sequences in a language. They help in speech recognition, machine translation, and text generation.

- Text Prediction: N-gram models are used in autocomplete suggestions and predictive text input on devices like smartphones.

- Spell Checking: N-grams can be used to suggest correct spellings based on context. For example, if "they're" frequently follows "you're," the model can suggest "they're" when "you're" is mistyped.

- Information Retrieval: In search engines, n-gram models help identify relevant documents by predicting common word sequences.

- Document Classification: N-gram features are used in text classification tasks to capture patterns in documents.

However, n-gram models have limitations, especially for larger values of n. They struggle with capturing long-range dependencies and can become memory-intensive due to the large number of possible combinations. Nevertheless, they remain a simple yet effective tool in many NLP applications.

**11. How n-gram model work for spelling correction/ word suggestion ?**
The n-gram model can be leveraged for spelling correction and word suggestion by analyzing the frequency and probability of word sequences in a given text corpus. Here's how the n-gram model works for these tasks:

1. Building N-Gram Frequencies:
   - A text corpus is used to build n-gram frequencies. N-grams can be of varying sizes (unigrams, bigrams, trigrams, etc.), but for spelling correction and word suggestion, typically bigrams (n=2) or trigrams (n=3) are used.
   - The frequency of each n-gram (sequence of n words) is counted in the corpus.

2. Calculating Conditional Probabilities:
   - For each n-gram, the conditional probability of the next word (or the word following the sequence) is calculated. This involves dividing the frequency of the n-gram by the frequency of the preceding (n-1)-gram.
   - The calculated probabilities reflect how often specific sequences of words occur in the corpus.

3. Spelling Correction:
   - Given a misspelled word, the n-gram model identifies possible corrections by considering word sequences that are similar in terms of n-grams.
   - The model looks for n-grams that match parts of the misspelled word. It then suggests corrections based on words that frequently follow those matching n-grams in the corpus.

4. Word Suggestion:
   - When a user types a partial word, the n-gram model suggests the most probable completions based on the preceding n-gram.
   - The model suggests words that often follow the input sequence in the corpus, using their conditional probabilities to rank the suggestions.

5. Contextual Analysis:
   - The n-gram model considers the context in which a word appears. For spelling correction, it suggests words based on similar sequences, and for word suggestion, it considers the input's context to predict likely completions.

6. Ranking and Presentation:
   - The suggested corrections or word completions are ranked based on their conditional probabilities. The most probable suggestions are presented to the user.

7. Error Handling:
   - To handle situations where a misspelled word doesn't have an exact n-gram match, the model can explore alternative strategies like phonetic similarity, edit distance, or incorporating larger context.

By analyzing the probabilities of word sequences, the n-gram model provides spelling correction and word suggestion capabilities that account for context and frequent patterns in the given language. While the n-gram

model is effective for simple errors and suggestions, more advanced techniques are used for handling complex linguistic phenomena and context-dependent corrections.

# Module 3

**1. What is Penn Tree bank tagset? How do we assign pos tags to words in sentence**
The Penn Treebank Tag Set is a widely used set of part-of-speech (POS) tags that provides a standardized way to annotate words in a sentence with their grammatical categories. It was developed as part of the Penn Treebank project, which aimed to annotate a large corpus of English text with syntactic and grammatical information. The tag set consists of various POS tags that represent different parts of speech and grammatical features.

Here are some common POS tags from the Penn Treebank Tag Set:

- Noun (NN): Common noun
- Verb (VB): Base form verb
- Adjective (JJ): Adjective
- Adverb (RB): Adverb
- Pronoun (PRP): Personal pronoun
- Preposition (IN): Preposition or subordinating conjunction
- Conjunction (CC): Coordinating conjunction
- Determiner (DT): Determiner
- Numeral (CD): Cardinal number
- Interjection (UH): Interjection or exclamation

Assigning POS tags to words in a sentence involves using linguistic rules, context, and sometimes machine learning algorithms. Here's a general approach to assigning POS tags:

1. Rule-Based Tagging:
   - Linguistic rules based on patterns and context are used to assign POS tags. For example, words ending in "-ing" are often verbs, while words ending in "-ly" are frequently adverbs.
   - Contextual information like nearby words and sentence structure can also influence tag assignment.

2. Dictionary-Based Tagging:
   - Some words have a limited set of possible POS tags. A dictionary can map words to their likely POS tags based on their definitions and typical usage.

3. Statistical Tagging:
   - Machine learning algorithms, like Hidden Markov Models (HMM) or Conditional Random Fields (CRF), can be trained on labeled corpora to learn the most likely POS tags for words based on their context.
   - These models consider the probabilities of transitions between POS tags and use them to make predictions for unseen sentences.

4. Combination of Approaches:
   - Many POS taggers combine rule-based and statistical approaches to achieve higher accuracy. For instance, a rule-based tagger might be followed by a statistical tagger to catch cases that don't fit well with rules.

5. Manual Corrections:
   - In some cases, manual corrections are made to fix errors made by automated taggers. Human annotators review and adjust the assigned POS tags to ensure accuracy.

Assigning POS tags accurately is crucial for various NLP tasks, such as syntactic analysis, sentiment analysis, and information retrieval. While POS tagging is highly effective, it can be challenging due to word ambiguity, context sensitivity, and irregular linguistic patterns.

**2. Explain sequence labeling using hidden markov model**
Sequence labeling is a natural language processing (NLP) task where the goal is to assign a label to each element in a sequence of input data. Hidden Markov Models (HMMs) are a class of probabilistic models commonly used for sequence labeling tasks, especially when there's a sequential dependence between the elements. One of the prominent applications of HMMs in NLP is part-of-speech tagging, where each word in a sentence is assigned a part-of-speech label.

Here's how sequence labeling using a Hidden Markov Model works:

1. State Space:
   - In the context of NLP, the state space represents the set of possible labels or tags that can be assigned to each element in the sequence. For example, in part-of-speech tagging, the state space includes tags like nouns, verbs, adjectives, etc.

2. Observation Space:
   - The observation space consists of the actual input elements (e.g., words) in the sequence.

3. Transition Probabilities:
   - HMMs model the transition probabilities between different states (tags) as the sequence progresses. These probabilities represent the likelihood of moving from one state to another.
   - In part-of-speech tagging, transition probabilities capture the likelihood of certain parts of speech following others based on linguistic patterns.

4. Emission Probabilities:
   - Emission probabilities represent the likelihood of observing a particular input element (word) given a certain state (tag).
   - For part-of-speech tagging, emission probabilities capture the likelihood of a word being associated with a particular part of speech.

5. Initial State Probabilities:
   - Initial state probabilities represent the likelihood of starting in a certain state (tag) at the beginning of the sequence.

6. Decoding:
   - The decoding process involves finding the most likely sequence of states (tags) that generated the observed sequence of input elements.
   - The Viterbi algorithm is commonly used for decoding in HMMs. It efficiently computes the most likely state sequence by considering both transition and emission probabilities.

7. Training:
   - HMMs can be trained using labeled data, where the correct state sequences are known. The Baum-Welch algorithm (a type of Expectation-Maximization algorithm) is used to estimate the model's parameters (transition and emission probabilities) from the training data.

8. Prediction and Labeling:
   - Once the HMM is trained, it can be used to predict labels for unseen sequences by finding the most likely state sequence using the Viterbi algorithm.

Sequence labeling using Hidden Markov Models is versatile and applicable to various NLP tasks beyond part-of-speech tagging, such as named entity recognition, chunking, and more. However, HMMs assume that the sequence of labels depends only on a fixed number of previous labels (Markov assumption), which might not capture long-range dependencies well.