

UNIT 1

Natural Language Processing

-Natural Language Processing (NLP) is a field of computer science that deals with developing algorithms and models that enable computers to understand, analyse, and generate human languages.

To better understand NLP, let's break down the definition:

- **Algorithms:** These are a set of rules or instructions that computers use to perform a specific task.
- **Models:** These are mathematical representations that computers use to learn patterns in data and make predictions.
- **Human languages:** These refer to any language that people speak, such as English, Spanish, Chinese, etc.

-NLP involves several stages or techniques that enable machines to process and understand natural language data. These techniques include:

- **Tokenization:** This involves breaking down a piece of text into smaller units called tokens, such as words or sentences.
- **Part-of-speech tagging:** This involves identifying the part of speech (noun, verb, adjective, etc.) of each word in a sentence.
- **Parsing:** This involves analyzing the grammatical structure of a sentence to identify the relationship between words.
- **Named entity recognition:** This involves identifying and extracting named entities, such as names of people, organizations, or locations, from a piece of text.
- **Sentiment analysis:** This involves analyzing the tone or emotion conveyed in a piece of text.
- **Machine learning:** This involves training computer models on large datasets to learn patterns and make predictions.

-NLP has many real-world applications, such as language translation, chatbots, speech recognition, and text-to-speech systems.

-For example, language translation tools like Google Translate use NLP techniques to translate text from one language to another.

-Chatbots use NLP techniques to understand and respond to natural language queries from users. Speech recognition systems like Siri and Alexa use NLP techniques to convert speech to text.

-Text-to-speech systems use NLP techniques to convert text into spoken language.

-field of computer science that deals with developing algorithms and models that enable computers to process and understand human languages.

- involves several techniques such as tokenization, part-of-speech tagging, parsing, named entity recognition, sentiment analysis, and machine learning, and has many real-world applications.

Why NLP is hard?

NLP is hard because natural language is inherently complex, ambiguous, and highly variable. Here are some of the reasons why NLP is a challenging field:

- **Ambiguity:** Natural language is inherently ambiguous, meaning that the same word or phrase can have multiple meanings depending on the context. For example, the word "bank" can refer to a financial institution or the side of a river.
- **Syntax:** Natural language has complex syntax and grammar, which can be difficult to parse and analyse. For example, the sentence "I saw the man with the telescope" can be interpreted in two different ways, depending on whether the man or the speaker had the telescope.
- **Colloquialisms:** Natural language includes many colloquialisms, idioms, and other expressions that can be difficult to understand for non-native speakers or for machines.
- **Variability:** Natural language is highly variable, meaning that people can express the same idea in different ways using different words, sentence structures, and idioms.
- **Domain-specific knowledge:** Natural language often requires background knowledge of the subject matter or context to be understood fully. For example, understanding medical terminology requires knowledge of medical concepts and vocabulary.
- **Data scarcity:** NLP requires large amounts of annotated data for training models, which can be scarce or difficult to obtain for certain languages or domains.

-To overcome these challenges, NLP researchers and practitioners use a combination of rule-based and machine learning-based approaches, and they continuously develop new techniques and models to improve the accuracy and effectiveness of NLP applications.

Programming languages Vs Natural Languages

-Programming languages and natural languages are fundamentally different in their structure, purpose, and use.

-Programming languages are designed for communicating with computers to give them specific instructions to execute tasks. They have strict grammatical rules and syntax, and each statement is precisely defined and must be followed exactly.

-Programming languages are used to create software, applications, and websites that perform specific functions.

-On the other hand, natural languages are used by humans to communicate with each other.

-They are flexible and adaptive, with grammar and syntax that can vary depending on the context and the speaker.

-Natural languages are used for a wide range of purposes, from everyday conversation to literature and art.

-Programming languages and natural languages also differ in their expressiveness.

-Programming languages are designed to express precise instructions for computers to execute, whereas natural languages are designed to express a wide range of ideas, emotions, and experiences.

-While both programming languages and natural languages involve a form of communication, they serve very different purposes and are designed to be used in different contexts.

-However, there are areas where these two types of languages overlap, such as in the field of natural language processing (NLP), where programming languages are used to develop algorithms and models that can understand and process natural language.

Are natural languages regular?

-Natural languages are not regular in the strict sense of the term.

-In computer science, a regular language is a formal language that can be described by a regular expression or a deterministic finite automaton (DFA).

-Regular languages are characterized by a set of simple rules that can be applied recursively to generate all possible valid strings in the language.

-However, natural languages are complex and often have irregularities and exceptions that cannot be captured by a simple set of rules.

-Natural languages have complex syntax and grammar that can vary depending on the context and the speaker, and they often involve idioms, slang, and other expressions that are not easily captured by a formal system.

For example, irregular verbs and noun plurals in English do not follow a predictable pattern, and there are many exceptions to the rules.

- In addition, natural languages can have multiple meanings for the same word or phrase, making it difficult to generate a finite set of valid strings.

-Despite these irregularities, researchers in the field of natural language processing (NLP) have developed techniques and models that can process and analyse natural language data.

- These techniques typically involve a combination of rule-based and machine learning-based approaches to handle the complexity and variability of natural language.

Finite automata for NLP

-Finite automata are a theoretical framework for studying formal languages, including natural languages.

-In NLP, finite automata can be used to model and analyse the structure and grammar of natural language.

-Finite automata are abstract machines that can recognize and generate strings in a formal language.

-They consist of a set of states, a set of input symbols, a transition function that maps a state and an input symbol to a new state, and a set of accepting states.

-The automaton starts in an initial state and reads input symbols from a string one at a time, following the transition function until it either reaches an accepting state or does not have any valid transitions left.

-In NLP, finite automata can be used to model and analyse various aspects of natural language, including:

- **Morphology:** Finite automata can be used to model the structure of words and how they are formed from smaller units such as roots, prefixes, and suffixes.
- **Syntax:** Finite automata can be used to model the structure of sentences and how they are formed from phrases and clauses.
- **Semantics:** Finite automata can be used to model the meaning of words and how they combine to form the meaning of sentences.
- **Text classification:** Finite automata can be used to classify texts based on their structure and content, for example, to identify spam emails or to categorize news articles by topic.

-However, while finite automata provide a useful theoretical framework for analysing the structure and grammar of natural language, they are limited in their ability to capture the full complexity and variability of natural language.

-As a result, researchers in NLP often use more advanced techniques, such as machine learning, to model and analyse natural language data.

Stages of NLP

The stages of natural language processing (NLP) generally involve several steps, which can vary depending on the specific application and techniques used.

However, the following are some of the common stages involved in most NLP tasks:

Text Pre-processing: This stage involves cleaning and preparing the text data for further analysis. This may include tasks such as removing punctuation and stop words, tokenizing the text into individual words, and stemming or lemmatizing the words to reduce them to their root form.

Part-of-Speech (POS) Tagging: This stage involves labeling each word in the text with its corresponding part-of-speech (e.g., noun, verb, adjective, etc.). POS tagging is important because it provides information about the grammatical structure of the text, which can be used for further analysis.

Parsing: This stage involves analyzing the syntactic structure of the text to identify the relationships between words and phrases. This may involve techniques such as constituency parsing or dependency parsing.

Named Entity Recognition (NER): This stage involves identifying and categorizing named entities in the text, such as people, organizations, and locations. NER is important for many NLP applications, such as information extraction and sentiment analysis.

Sentiment Analysis: This stage involves analyzing the text to determine the writer's overall sentiment or opinion. This may involve techniques such as identifying keywords, analyzing sentence structure, and using machine learning algorithms to classify the sentiment of the text.

Machine Translation: This stage involves translating text from one language to another. Machine translation involves techniques such as statistical machine translation, neural machine translation, or rule-based machine translation.

Text Generation: This stage involves generating new text based on existing text data. This may involve techniques such as language models, generative models, or deep learning algorithms.

These stages are often interconnected and may involve multiple techniques and algorithms. However, they provide a general framework for understanding the different tasks involved in NLP.

Challenges and Issues in NLP

-Natural language processing (NLP) is a complex and challenging field that involves analysing and understanding human language using computational methods.

-While NLP has made significant advances in recent years, there are still a number of challenges and issues that must be addressed to improve the accuracy and effectiveness of NLP systems.

-Some of the key challenges and issues in NLP are:

Ambiguity: Natural language is often ambiguous and can have multiple meanings. For example, the word "bank" can refer to a financial institution or the edge of a river. Resolving such ambiguities can be a significant challenge in NLP.

One approach to resolving ambiguity is to use context to disambiguate the meaning of words. For example, if the word "bank" appears in a sentence with other financial terms, such as "loan" or "investment," it is likely that it is referring to a financial institution.

Variations in language: Different people use language in different ways, which can create variations in language that can be challenging to process. Variations can include dialects, idioms, slang, and colloquialisms.

To address variations in language, NLP systems must be able to recognize and understand different forms of language use. This can involve training NLP models on diverse data sources, including social media, informal speech, and other forms of non-standard language use.

Context: Understanding language requires understanding the context in which it is used. Context can include both the surrounding words and the broader cultural and situational context. Interpreting context can be a significant challenge in NLP.

To address context, NLP systems must be able to analyze and understand the broader context in which language is used. This can involve incorporating knowledge of cultural and situational factors into NLP models, as well as using contextual information to disambiguate the meaning of words.

Data quality: NLP systems require large amounts of high-quality data to be effective. However, data can be noisy, incomplete, or biased, which can impact the accuracy of NLP systems.

To address data quality issues, NLP researchers must be diligent in collecting, cleaning, and pre-processing data to ensure that it is of high quality. Additionally, machine learning algorithms can be used to identify and remove noisy or biased data.

Machine learning challenges: Many NLP systems rely on machine learning algorithms, which can be challenging to train and optimize. Machine learning algorithms can be prone to overfitting, require large amounts of training data, and be sensitive to the quality of the data.

To address machine learning challenges, NLP researchers must be diligent in selecting appropriate training data and machine learning algorithms. Additionally, techniques such as regularization and

cross-validation can be used to improve the performance and generalizability of machine learning models.

Privacy and ethical concerns: NLP systems can collect and process sensitive personal information, which can raise privacy and ethical concerns. For example, there may be concerns around data ownership, consent, and potential biases in the data.

To address privacy and ethical concerns, NLP researchers must be mindful of the potential risks and implications of their work. Additionally, techniques such as differential privacy and data anonymization can be used to protect sensitive personal information.

Multilingualism: NLP systems must be able to handle multiple languages and language variations. This can be challenging, as different languages have different grammatical structures and vocabularies. To address multilingualism, NLP researchers must develop models that are capable of handling multiple languages and language variations. This can involve using multilingual corpora, developing language-specific models, or using techniques such as transfer learning to adapt models to new languages.

Overall, NLP faces several challenges and issues that require ongoing research and development to address.

Tokenization

-Tokenization is the process of breaking down a text into smaller units, called tokens, which can be words, phrases, or even individual characters.

- Tokenization is an important pre-processing step in natural language processing (NLP) that enables machines to analyse and understand text data.

-Tokenization typically involves identifying and separating words and punctuation marks in a text. For example, consider the sentence: "The quick brown fox jumps over the lazy dog." Tokenization of this sentence would result in the following tokens: "The", "quick", "brown", "fox", "jumps", "over", "the", "lazy", "dog", "."

-In some cases, tokenization can be more complex.

-For example, in the sentence "I'm not sure if I've seen this movie before," the tokenization process would need to recognize the contraction "I'm" as a single token rather than two separate tokens ("I" and "m").

-There are several tokenization techniques that can be used in NLP, including:

Word-based tokenization: This involves breaking down a text into individual words. This is the most common type of tokenization and is often used for tasks such as sentiment analysis, topic modeling, and text classification.

Sentence-based tokenization: This involves breaking down a text into individual sentences. This is often used for tasks such as summarization and text generation.

Character-based tokenization: This involves breaking down a text into individual characters. This can be useful for tasks such as handwriting recognition and machine translation.

Subword-based tokenization: This involves breaking down a text into smaller units, such as character n-grams or syllables. This is often used for tasks such as text normalization and speech recognition.

Overall, tokenization is an important step in NLP that enables machines to analyze and understand text data. By breaking down text into smaller units, machines can more easily process and analyze text data, leading to improved accuracy and effectiveness of NLP models.

Stemming

-Stemming is a technique in natural language processing (NLP) that reduces words to their root or stem form.

-The goal of stemming is to group together different forms of the same word so that they can be treated as a single entity during analysis.

-For example, the words "run", "running", and "runner" are different forms of the same word, but they have different suffixes. Stemming would reduce these words to their common root form "run", which can help to simplify the analysis of text data.

-Stemming algorithms typically use a set of rules to identify the common root form of a word by removing suffixes and prefixes.

- These rules are based on linguistic and morphological principles and vary depending on the language and the stemming algorithm used.

-There are several popular stemming algorithms, including:

Porter stemming algorithm: This is one of the oldest and most widely used stemming algorithms. It uses a set of heuristic rules to remove common suffixes and reduce words to their stem form.

Snowball stemming algorithm: This is a more advanced stemming algorithm that was developed as an improvement to the Porter stemming algorithm. It uses a more complex set of rules and can stem words in multiple languages.

Lancaster stemming algorithm: This is a highly aggressive stemming algorithm that uses a set of rules to remove suffixes and prefixes from words. It can sometimes reduce words to very short and sometimes unrecognizable forms.

Overall, stemming is a useful technique in NLP that can help to simplify the analysis of text data by reducing different forms of words to their common root form. However, it is important to note that stemming is not perfect and can sometimes produce incorrect stems, especially for irregular or uncommon words.

Lemmatization

-Lemmatization is a technique in natural language processing (NLP) that involves reducing words to their base or dictionary form, called a lemma.

- The goal of lemmatization is to group together different forms of a word so that they can be treated as a single entity during analysis.

-For example, the word "better" can have different forms such as "good", "well", or "improve". By lemmatizing the word "better", all of its different forms can be reduced to its base form, "good", which can help to simplify the analysis of text data.

-Unlike stemming, which only removes the suffixes from a word to produce a root form, lemmatization takes into account the part of speech of the word and uses a detailed morphological analysis to derive the base form of a word.

-For instance, in the sentence "The dogs are barking", lemmatization would identify "dogs" as the plural form of "dog" and would reduce it to the base form "dog".

-Lemmatization algorithms typically use a lexical knowledge base, such as a dictionary or a thesaurus, to identify the correct lemma for each word.

-The algorithm may also take into account the context of the word in the sentence to determine its correct part of speech and lemma.

Some popular lemmatization algorithms include:

WordNet lemmatizer: This is a lemmatization algorithm that uses WordNet, a large lexical database of English words, to identify the base form of each word.

Stanford lemmatizer: This is a lemmatization algorithm that uses a combination of rule-based and statistical techniques to derive the base form of a word.

Spacy lemmatizer: This is a lemmatization algorithm that uses a statistical model to identify the correct lemma for each word based on its context in the sentence.

Overall, lemmatization is a useful technique in NLP that can help to simplify the analysis of text data by reducing different forms of words to their base form. It is often used in tasks such as text classification, sentiment analysis, and information retrieval.

Part of Speech Tagging

-Part of speech (POS) tagging is a technique in natural language processing (NLP) that involves labelling the words in a text corpus with their corresponding part of speech, such as noun, verb, adjective, adverb, etc.

-POS tagging is an important pre-processing step in many NLP tasks, such as text classification, sentiment analysis, named entity recognition, and machine translation.

-POS tagging algorithms use statistical or rule-based methods to assign the correct part of speech to each word in a sentence.

-Statistical methods use machine learning algorithms to learn from large annotated datasets to predict the correct part of speech of each word in a new sentence.

-Rule-based methods, on the other hand, use a set of hand-crafted rules to assign the correct part of speech to each word based on its context in the sentence.

-For example, consider the sentence "The cat is sleeping on the couch". A POS tagging algorithm would label "The" and "the" as determiners, "cat" and "couch" as nouns, "is" as a verb, and "sleeping" as a present participle.

-There are several popular POS tagging algorithms, including:

Stanford POS Tagger: This is a widely-used POS tagging algorithm that uses a probabilistic model to assign the correct part of speech to each word in a sentence.

NLTK POS Tagger: This is a popular POS tagging algorithm that uses a rule-based approach to assign the correct part of speech to each word in a sentence.

Spacy POS Tagger: This is a POS tagging algorithm that uses statistical models to predict the correct part of speech of each word based on its context in the sentence.

Overall, POS tagging is an important technique in NLP that helps to extract meaningful information from text data by identifying the parts of speech of words in a sentence.

UNIT 2

What is Morphology?

- Morphology is the study of how words are formed and how they are related to each other.
- It is like breaking down words into their smallest parts and figuring out how those parts combine to create new words.
- For example, the word "unhappiness" can be broken down into two parts: "un-" and "happiness." ---
- The prefix "un-" means "not," and when added to "happiness," it changes the meaning to the opposite of happiness. So, "unhappiness" means "not happy."
- Morphology also looks at how words change their form to show different meanings.
- For instance, the word "walk" can become "walked" to show that the action happened in the past. - Similarly, the word "friend" can become "friends" to show that there is more than one friend.
- Overall, morphology helps us understand how words are created and how they can be changed to convey different meanings.

Types of Morphemes

Free Morphemes:

- Free morphemes are the smallest units of meaning in language that can stand alone as separate words and carry meaning on their own.
- are not attached to any other morphemes. In other words, free morphemes can exist on their own as words, without needing any other words to provide context or meaning.
- Examples of free morphemes include words like "cat," "dog," "happy," "run," "jump," and "talk." These words can be used by themselves to convey a meaning.

Bound Morphemes:

- Bound morphemes are morphemes that cannot stand alone as separate words and must be attached to other morphemes to create a word.
- They are usually added to the beginning or end of a word, and they change the word's meaning or function. Bound morphemes can be further divided into two categories:
 - Prefixes: Prefixes are bound morphemes that are added to the beginning of a word to change its meaning. For example, the prefix "un-" is added to the word "happy" to create "unhappy," which means the opposite of happy. Other examples of prefixes include "pre-" (as in "preview"), "dis-" (as in "disagree"), and "in-" (as in "inactive").
 - Suffixes: Suffixes are bound morphemes that are added to the end of a word to change its meaning or function. For example, the suffix "-ed" is added to the verb "walk" to create the past tense "walked." Other examples of suffixes include "-ing" (as in "running"), "-s" (as in "cats"), and "-able" (as in "comfortable").

-It's important to note that some morphemes can be both free and bound, depending on how they are used in a sentence.

-For example, the word "book" is a free morpheme when used as a standalone word, but it becomes a bound morpheme when combined with other morphemes to form a new word, such as "bookshelf" (where "book" is the bound morpheme).

Inflectional morphology

- branch of linguistics that deals with the study of how words are changed or inflected to indicate different grammatical features such as tense, aspect, mood, number, case, and gender.

-deals with how words change their form to convey grammatical information.

-Inflectional morphemes are a type of bound morphemes that are added to the end of a word to indicate these grammatical features.

-do not change the word's meaning or its lexical category, but rather they modify its grammatical function within a sentence.

- helps us understand the grammatical relationships between words and how they change depending on their role in a sentence.

- fundamental aspect of language that allows us to convey complex ideas through relatively simple structures.

- By adding inflectional morphemes to words, we can change their form to indicate a wide variety of grammatical features, which in turn can affect the meaning and function of a sentence.

- Without inflectional morphology, we would need to use additional words to convey these same ideas, which would make sentences longer and more complex.

- Inflectional morphology allows us to convey the same information in a more efficient and streamlined way.

- Inflectional morphology is also important for understanding how words are related to each other within a language.

Derivational morphology

-Derivational morphology is the branch of linguistics that deals with the study of how words are created or derived from other words through the addition of prefixes or suffixes.

- Unlike inflectional morphology, which adds morphemes to words to indicate grammatical features, derivational morphology changes the meaning or part of speech of a word.

-Derivational morphemes are a type of bound morphemes that are added to the base form of a word to create a new word with a different meaning or part of speech.

- Unlike inflectional morphemes, derivational morphemes can change the lexical category of the word, creating a new word with a different part of speech.
- Derivational morphology is a powerful tool that allows us to create new words from existing ones, thereby expanding the vocabulary of a language.
- By adding prefixes and suffixes to words, we can create new words that convey shades of meaning that might not be possible with the original word alone.
- For example, the word "happy" can be modified with the prefix "un-" to create the word "unhappy," which conveys the opposite meaning. Similarly, the word "beauty" can be modified with the suffix "-ful" to create the word "beautiful," which describes something that has the quality of beauty.
- Derivational morphology also plays an important role in word formation in many languages.
- Derivational morphology is also important for understanding how words are related to each other within a language.

Morphological parsing with Finite State Transducers (FST)

- Morphological parsing is the process of analysing a word and breaking it down into its constituent morphemes.
- Finite state transducers (FST) are a type of computational model that can be used for morphological parsing.
- An FST is a type of automaton that consists of a set of states, a set of input symbols, a set of output symbols, a set of transition rules, and a set of final states.
- FSTs can be used to model regular languages, which are languages that can be generated by a regular grammar.
- In the context of morphological parsing, an FST can be used to model the set of possible morphological analyses for a given word.
- To create an FST for morphological parsing, we first need to define the set of possible morphemes in the language.
- We then define a set of rules that describe how these morphemes can be combined to form words. For example, in English, we might define rules for combining prefixes (e.g., "un-"), suffixes (e.g., "-able"), and roots (e.g., "write") to form words like "unwritable" or "writable."
- Once we have defined the set of possible morphemes and the rules for combining them, we can use an FST to model the set of possible analyses for a given word.
- The FST takes the word as input and generates a set of possible morphological analyses as output.
- Each analysis consists of a sequence of morphemes that can be combined to form the word.
- To use an FST for morphological parsing, we need to train it on a corpus of annotated data.
- This data consists of words and their corresponding morphological analyses.

-We use this data to learn the rules for combining morphemes and to determine the probabilities of different analyses for each word.

-Once the FST has been trained, we can use it to parse new words and generate morphological analyses.

- We input the word into the FST, and it generates a set of possible analyses.

-We can then use language-specific rules and heuristics to select the most likely analysis based on factors such as frequency, context, and semantic plausibility.

-Overall, morphological parsing with FSTs is a powerful tool for analysing the morphology of a language and generating new analyses for novel words.

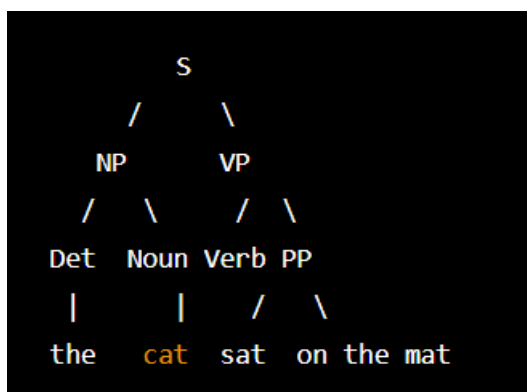
-It can be used for a wide range of applications, including language modelling, speech recognition, and machine translation.

SYNTACTIC ANALYSIS

Syntactic representation

Syntactic representation refers to the way in which the structure of a sentence or phrase is represented in a formal language. Syntactic representations of natural language can take many forms, ranging from simple tree structures to more complex formalisms based on logic or automata theory. Here are some examples:

Phrase Structure Trees: This is a common representation used in linguistics, where sentences are represented as trees where each node represents a constituent phrase (e.g., noun phrase, verb phrase, etc.), and the branches represent the relationships between these phrases. For example, the sentence "The cat sat on the mat" can be represented as the following tree:



-**Dependency Graphs:** This representation focuses on the relationships between words in a sentence, rather than their hierarchical structure. In a dependency graph, each word is represented as a node, and the edges represent the syntactic relationships between the words.

- **Logical Forms:** This representation uses formal logic to represent the meaning of a sentence. Each word or phrase is translated into a logical expression, and the relationships between these

expressions are represented using logical connectives. For example, the sentence "All cats are mammals" can be represented as the following logical form:

$\forall x [\text{cat}(x) \rightarrow \text{mammal}(x)]$

- **Abstract Syntax Trees:** This is a more formal representation used in computer science and programming languages. In an abstract syntax tree, each node represents a syntactic construct (e.g., function call, assignment, etc.), and the branches represent the relationships between these constructs. For example, the code fragment "x = y + 2" can be represented as the following abstract syntax tree:

```
=  
/  
 \  
x  +  
  / \  
  y  2
```

Parsing Algorithms

Recursive Descent Parsing:

Recursive descent parsing is a top-down parsing algorithm that uses a set of recursive functions to match the input sentence against a grammar. The grammar is typically defined in the form of production rules, which specify how different types of phrases can be combined to form sentences. The algorithm starts with the top-level symbol of the grammar (usually the sentence itself) and recursively applies the production rules to generate a parse tree. At each step, the algorithm chooses the next production rule based on the current input symbol and the context of the parse. If there is a match, the algorithm moves to the next symbol in the input and the next production rule. If there is no match, the algorithm backtracks to the previous symbol in the input and tries another production rule. Recursive descent parsing is simple to implement and easy to understand, but can be inefficient and may not handle all types of grammars.

Earley Parsing:

Earley parsing is a bottom-up parsing algorithm that uses dynamic programming to build up a set of possible parse trees for a given sentence. The algorithm works by maintaining a set of items, which represent partially completed parse trees. Each item consists of a production rule, a position in the input, and a set of tokens that have already been matched. The algorithm starts by adding an initial item to the set, and then iteratively expands the set by applying production rules to the current items. The algorithm keeps track of all the possible paths through the set of items, and outputs the ones that generate a complete parse tree. Earley parsing is more powerful than recursive descent parsing and can handle any context-free grammar, but can be slower and may require more memory.

CYK Parsing:

CYK parsing is a bottom-up parsing algorithm that uses dynamic programming to determine if a given sentence can be generated by a given grammar. The algorithm works by constructing a table that represents all possible combinations of non-terminals that can generate a substring of the input sentence. The table is filled in using a recursive formula that combines smaller substrings to generate larger ones. Once the table is filled in, the algorithm checks whether the top-level symbol of the grammar can be generated by the entire sentence. CYK parsing is more efficient than Earley parsing for certain types of grammars, such as context-free grammars in Chomsky normal form.

Chart Parsing:

Chart parsing is a generalization of Earley parsing that uses a chart data structure to store partial parse trees. The chart is a table that keeps track of all the possible parse trees that have been generated so far, and allows for efficient sharing of sub-trees between different parts of the parse. The algorithm works by adding a set of initial items to the chart, and then iteratively expanding the chart by applying production rules to the current items. Chart parsing can handle a wide variety of grammars, including those with ambiguous or recursive rules.

Probabilistic context-free grammars

- type of formal grammar used in natural language processing (NLP) to model the syntax of a language.
- consists of a set of production rules, which describe how a sentence can be constructed by combining words and other linguistic structures.
- particularly useful for generating new sentences that are similar in structure and style to those in a training corpus
- can also be used for parsing, which is the process of analyzing a sentence to determine its syntactic structure.
- In parsing, a PCFG is used to generate a parse tree for a given sentence by applying the production rules in a way that maximizes the probability of the parse tree.
- do not capture semantic information, meaning that they cannot represent the meaning of a sentence.
- they sometimes produce sentences that are grammatically correct but do not make sense semantically

Statistical parsing

- type of natural language processing (NLP) technique that uses statistical models to analyze the structure of sentences in a language.
- involves building probabilistic models of how words in a sentence are related to each other, based on a large corpus of annotated text data.

- a parse tree is generated for a given sentence by selecting the most probable parse tree from a set of possible parse trees.

- Each parse tree represents a possible syntactic structure for the sentence, and the probability of a parse tree is determined by the probability of the production rules used to generate the tree.

-Types:

- Probabilistic context free grammars
- Dependency parsing

Probabilistic context-free grammar (PCFG) parsing:

PCFG parsing is a type of statistical parsing that uses a probabilistic context-free grammar to model the syntax of a language and generate parse trees. A context-free grammar (CFG) is a formal grammar that consists of a set of production rules that define how strings of symbols can be generated. In a PCFG, each production rule has a probability associated with it, indicating how likely it is to be used in generating a sentence.

PCFG parsing works by taking a sentence and generating all possible parse trees for that sentence using the PCFG model. Each parse tree represents a possible syntactic structure for the sentence, and the probability of a parse tree is determined by the probability of the production rules used to generate the tree. The most probable parse tree is selected as the final result.

Dependency parsing:

Dependency parsing is a type of statistical parsing that uses a dependency grammar to model the relationships between words in a sentence. In a dependency grammar, each word in a sentence is linked to one or more "head" words, which represent the words that the linked word is dependent on. The links between words are represented as directed edges in a graph, with the head words at the top and the dependent words at the bottom.

Dependency parsing works by taking a sentence and generating a dependency graph that represents the relationships between words in the sentence. The most probable graph is selected as the final result. Dependency parsing is often used for tasks such as information extraction and sentiment analysis, which require a deeper understanding of the relationships between words in a sentence.

SEMANTIC ANALYSIS

Lexical Semantic

Lexical semantics is a subfield of linguistics and natural language processing (NLP) that focuses on the study of the meaning of words and their relationships to each other. It is concerned with the way in which the meanings of individual words combine to create the meanings of phrases, sentences, and texts.

Lexical semantics involves the analysis of the meaning of words at different levels, including:

Word sense: This refers to the various meanings that a word can have, depending on the context in which it is used. For example, the word "bank" can refer to a financial institution, a river bank, or a place where one can sit.

Word similarity and relatedness: This involves the identification of words that are similar in meaning (synonyms) or related in meaning (antonyms, hyponyms, hypernyms, etc.). For example, "car" and "automobile" are synonyms, while "car" and "bicycle" are not.

Word associations: This involves the identification of words that tend to occur together in a particular context. For example, "coffee" is associated with "mug", "caffeine", "aroma", and "morning".

Word ambiguity: This refers to situations where a word has multiple meanings, and it is unclear which meaning is intended in a particular context. For example, the word "bat" can refer to a flying mammal or a piece of sports equipment.

Lexical semantics is important in many NLP applications, such as text classification, information retrieval, and machine translation. It can be used to improve the accuracy of these applications by ensuring that the correct meaning of a word is identified in a particular context.

Relations among lexemes & their senses

-Lexemes are the basic units of meaning in a language, and they are typically associated with multiple senses or meanings.

-The relationships between lexemes and their senses are important in lexical semantics and NLP, as they help to capture the nuances of meaning in a language.

Homonymy

-type of lexical ambiguity where two or more words have the same form (spelling and/or pronunciation) but have different meanings.

- can cause confusion in language understanding, especially in speech recognition and natural language processing (NLP) applications.

-types:

- **Homophones:** These are words that have the same pronunciation but different meanings and spellings. Examples include "to", "too", and "two"; "there", "their", and "they're"; and "write" and "right".
- **Homographs:** These are words that have the same spelling but different meanings and pronunciations. Examples include "tear" (to rip) and "tear" (water from the eye); "wind" (to twist) and "wind" (moving air); and "lead" (to guide) and "lead" (a metal).

-In NLP, homonymy can present a challenge in language understanding, particularly in speech recognition and machine translation

Polysemy

- Polysemy is a type of lexical ambiguity where a single word has multiple related meanings.
- In other words, polysemous words have different senses that are related to each other in some way, rather than being completely distinct meanings.
- For example, the word "bank" can refer to a financial institution where people can deposit and withdraw money, or it can refer to the side of a river. Both meanings of the word are related to the concept of a place where something is held, stored or contained.
- Another example is the word "book," which can refer to a physical object made of paper and ink that contains printed or written information, or it can refer to a record of scheduled events or appointments. Both of these senses are related to the concept of information being organized and recorded.
- In natural language processing (NLP), identifying the correct sense of a polysemous word is important for accurate language understanding and communication.
- Techniques such as word sense disambiguation and context analysis can be used to determine the most appropriate sense of a word based on the surrounding text.
- Polysemy is a common phenomenon in natural languages, and it allows for a richer and more flexible communication.
- However, it can also lead to confusion and misunderstandings if the intended sense of a word is not clear from the context.

Synonymy

- Synonymy is a linguistic phenomenon where two or more words have the same or similar meanings.
- Synonyms are words that can be used interchangeably in a particular context, without affecting the meaning of the sentence.
- For example, the words "big" and "large" are synonyms because they have a similar meaning of referring to something that is not small in size.
- Other examples include "happy" and "glad," "pretty" and "beautiful," and "fast" and "quick."
- In natural language processing (NLP), recognizing synonyms is important for a range of applications, such as information retrieval, text classification, and machine translation.
- Techniques such as word embeddings and lexical databases can be used to identify and represent synonyms in a computational context.
- However, it is important to note that synonyms are not always completely interchangeable, and their use can depend on the specific context and connotations associated with each word.
- For example, the words "home" and "house" are synonyms in the sense that they both refer to a place where people live, but they can have different connotations and associations depending on the context in which they are used.

Hyponymy

-Hyponymy is a linguistic phenomenon where a word or phrase is a specific example or subset of a more general word.

-In other words, a hyponym is a word that is a type of another word, which is called its hypernym. Hyponyms are used to describe more specific or specialized meanings of a word.

-For example, "dog" is a hypernym that refers to a general category of four-legged animals that are kept as pets or used for hunting, while "poodle," "beagle," and "labrador" are hyponyms that refer to specific types of dogs.

-Similarly, "fruit" is a hypernym that refers to a general category of sweet and edible produce, while "apple," "orange," and "banana" are hyponyms that refer to specific types of fruit.

-In natural language processing (NLP), recognizing hyponyms and hypernyms is important for tasks such as information retrieval, text classification, and sentiment analysis.

-Techniques such as word embeddings and lexical databases can be used to identify and represent hyponyms and hypernyms in a computational context.

-Hyponymy is a useful linguistic tool for describing the relationships between words and the ways in which they are used in language.

-By understanding the hyponyms and hypernyms of a word, we can better understand its meaning and use it more effectively in communication.

WordNet

-WordNet is a large lexical database for the English language that was created by researchers at Princeton University.

-It is designed to provide a comprehensive resource for English words and their meanings, and is organized around the notion of "synsets," or sets of synonyms.

-Each synset in WordNet consists of a group of words that are considered synonymous, meaning that they share a similar meaning. For example, the synset for "dog" contains words such as "puppy," "hound," and "canine," all of which are related to the concept of a four-legged animal that is kept as a pet or used for hunting.

-In addition to providing synonyms for words, WordNet also includes information on the semantic relationships between words.

-These relationships include:

Hyponymy/hypernymy: This relationship connects synsets that are related by a hierarchical structure, where the hypernym is the more general concept, and the hyponym is a more specific instance of that concept. For example, "beagle" is a hyponym of "dog," which is a hypernym.

Meronymy/holonymy: This relationship connects synsets that are related by a part-whole relationship, where the meronym is a part of the holonym. For example, "wheel" is a meronym of "car," which is a holonym.

Antonymy: This relationship connects synsets that are opposites in meaning. For example, "hot" is an antonym of "cold."

-WordNet is widely used in natural language processing (NLP) applications such as information retrieval, text classification, and machine translation.

-It is particularly useful for tasks that involve identifying semantic relationships between words and for word sense disambiguation, where the correct meaning of a word is determined in a given context.

-Overall, WordNet is an important resource for understanding the relationships between words and for facilitating accurate and effective communication in natural language processing.

Word Sense Disambiguation (WSD)

-Word Sense Disambiguation (WSD) is the process of determining the correct sense or meaning of a word in a given context.

-Many words in natural language have multiple senses, and it can be difficult for computer programs to accurately determine which sense is intended in a particular sentence or passage.

-WSD is an important problem in natural language processing (NLP), as it is crucial for many applications such as machine translation, information retrieval, and text classification.

-For example, in machine translation, it is important to accurately translate a word based on its correct meaning in context, rather than relying on a single translation for all possible senses.

-There are various approaches to WSD, including knowledge-based, supervised, and unsupervised methods.

-Knowledge-based approaches rely on resources such as WordNet, which contains information on the meanings of words and their relationships.

-These approaches typically use rules and heuristics to match the context of a given word with its possible meanings.

-Supervised approaches, on the other hand, require annotated data where the correct sense of a word is provided for each instance in a given context.

-Machine learning algorithms are trained on this data to learn patterns and features that can help identify the correct sense of a word in new contexts.

-Unsupervised methods do not rely on annotated data or pre-existing resources, and instead use statistical techniques to identify patterns in large datasets of text.

-These methods often involve clustering or clustering-like algorithms to group similar words and identify patterns in their usage.

-Despite significant progress in WSD research, it remains a challenging problem in NLP due to the complexity of language and the ambiguity of words in context.

-However, accurate WSD is critical for improving the performance of many NLP applications, and continued research in this area is essential for advancing the field.

Dictionary based approach

-The dictionary-based approach to natural language processing (NLP) involves the use of pre-existing dictionaries or lexicons to perform various NLP tasks.

-These dictionaries typically contain collections of words or phrases along with their associated meanings or definitions.

-One of the most common uses of a dictionary-based approach is in part-of-speech tagging, which involves assigning each word in a piece of text with its appropriate part of speech, such as noun, verb, adjective, or adverb.

-To perform this task, a pre-existing lexicon is used that maps each word to its corresponding part of speech. This lexicon can be created manually, or it can be generated automatically using machine learning algorithms.

-Another use of the dictionary-based approach is in sentiment analysis, which involves determining the overall sentiment of a piece of text.

-In this task, a sentiment lexicon is used that assigns each word in the text a sentiment score, such as positive or negative.

-The overall sentiment of the text is then determined by aggregating the sentiment scores of each word.

-Dictionary-based approaches can also be used in text classification, which involves categorizing a piece of text into one or more predefined categories, such as news articles or customer reviews.

- In this task, a pre-existing lexicon of keywords is used to identify the most relevant category or categories for the text.

-While dictionary-based approaches can be useful in many NLP tasks, they are often limited by the coverage and accuracy of the lexicon.

-If the lexicon does not contain a particular word or phrase, or if the meaning of a word is ambiguous, then the approach may not perform well.

-Additionally, dictionary-based approaches may not be effective for tasks that require a deeper understanding of language, such as machine translation or natural language generation.

-In recent years, dictionary-based approaches have been complemented by more advanced machine learning techniques, such as deep learning and neural networks.

-These techniques can learn from large amounts of data and can perform well even in situations where a pre-existing lexicon may be incomplete or inaccurate.

Latent Semantic Analysis

-Latent Semantic Analysis (LSA) is a way of analyzing text to understand the relationships between words and documents.

-It works by identifying patterns of word usage across a large set of documents and then using those patterns to infer the meanings of words.

-For example, LSA might analyze a large collection of news articles and discover that the words "election" and "vote" often appear together in the same articles.

-Based on this pattern, LSA would infer that these words are related and have similar meanings.

-LSA uses a mathematical technique called singular value decomposition (SVD) to identify these patterns and relationships.

-Essentially, it creates a matrix that represents the frequency of words across documents, and then breaks that matrix down into smaller matrices that capture the most important patterns.

-LSA can be useful for a variety of natural language processing tasks, such as clustering similar documents, identifying the most important words in a document, or even generating new text based on existing documents.

However, LSA has some limitations. It requires a large corpus of text to work well, and it may not perform as well on short or noisy texts.

-Additionally, more advanced deep learning techniques may be more effective for certain tasks.

LSA is based on the idea that words that are used in similar contexts are likely to have similar meanings.

-It represents words and documents as vectors in a high-dimensional space, where the dimensions correspond to different words.

- The similarity between two words or documents can then be measured by calculating the cosine of the angle between their vectors.