

### *Sentence Segmentation*

Process of identifying the boundaries between sentences in a piece of text, and it is a fundamental task in NLP.

Sentence segmentation can be carried out using a variety of techniques, including rule-based methods, statistical methods, and machine learning algorithms.

Rule-based methods - Pre-defined rules based on punctuation and other markers to segment sentences.

Machine learning algorithms - Annotated datasets to train models that can automatically identify sentence boundaries. These models learn to recognize patterns and features in text that signal the end of one sentence and the beginning of another .

### *Relevance and Application*

Sentence segmentation is the fundamental and the first step in the Natural Language Processing (NLP) pipeline.

Its primary relevance lies in breaking down large, unstructured text into meaningful, individual sentences, which are the basic units of thought and grammatical structure for subsequent analysis.

1. Contextual Understanding: Processing text sentence by sentence helps algorithms understand the local context and meaning more effectively than processing a continuous stream of words.

2. Data Preprocessing: It is a critical preprocessing step for many other NLP tasks, providing the necessary structure for downstream analysis.

### *Tokenization*

Tokenization is the process of breaking a stream of text into smaller, manageable units called tokens, which can be words, characters, or sentences. This is a fundamental first step in natural language processing (NLP) that structures raw text into a format that machines can process and analyze, enabling downstream tasks like text classification, sentiment analysis, and machine translation.

### *Relevance and Applications*

1. Enables machine learning: It converts raw text into a numerical representation, which is the only format neural networks can process.

2. Named Entity Recognition (NER): Identifying entities like names, dates, and locations within the text.

patterns in tokens that indicate spam.

3. Reduces complexity and improves efficiency: It breaks down large, unstructured text into discrete elements, making statistical and computational analysis more manageable and efficient.

4. Handles ambiguity and rare words: Break down rare or unknown words into smaller, more frequent units, reducing vocabulary size and improving model performance.

### *Stop Word Removal*

Stop word removal is a text processing step in NLP that removes common words like "the," "a," and "is" to reduce noise and improve efficiency for tasks. The process involves identifying and filtering out these insignificant words from a text.

How it works:

1. Tokenize the text
2. Convert to lowercase
3. Filter out stop words
4. Join the words

### *Relevance and Application*

Improves efficiency: Processing is faster, and the size of the corpus is decreased.

Increases model performance: Focuses algorithms on more meaningful, content-bearing words, which can improve accuracy in tasks like topic modelling.

Reduces data size: The total number of words to be processed is reduced

### *Lemmatization*

Lemmatization is a Natural Language Processing (NLP) technique that reduces a word to its base or dictionary form. It considers the word's context and part of speech to ensure the result is a valid dictionary word. This process improves text analysis for tasks like search and sentiment analysis.

Example: The words "running," "runs," and "ran" would all be lemmatized to "run". Similarly, "saw" could become "see" or "saw", depending on the context.

### *Relevance and Applications*

Accuracy: It provides more accurate results because the output is always a real, meaningful word.

Improved understanding: By grouping different inflected forms of a word, it helps NLP systems understand the core meaning of words, which is crucial for applications like search engines.

Effective text processing: It simplifies text analysis by reducing the number of unique words that need to be processed, which improves the performance of tasks like text classification and sentiment analysis.

### *Part-of-speech (POS)*

Part-of-speech (POS) tagging in NLP is the process of assigning a grammatical category, such as noun, verb, or adjective, to each word in a text. It is a foundational step that helps machines understand the structure and meaning of language, which is crucial for applications like machine translation, sentiment analysis, and named entity recognition. POS tagging also helps resolve ambiguity, as the same word can have different parts of speech depending on its context .

### *Relevance and Applications*

Understanding sentence structure: POS tagging provides the structure of a sentence.

Enabling downstream tasks: It serves as a first step for many other NLP applications by providing structured information that can be used as a feature in algorithms.

Improving information extraction: It is used in information extraction and named entity recognition (NER) to identify relationships between words and distinguish proper nouns from common nouns.

Analyzing patterns: POS tagging allows for searches based on grammatical patterns, such as finding all plural nouns that are not preceded by an article.

### *NER*

Named Entity Recognition (NER) is a subtask of Natural Language Processing (NLP) that locates and classifies named entities in unstructured text into predefined categories like person, organization, location, date, and money. It's a key process for information

extraction, question answering, and text summarization, which powers technologies like search engines and chatbots.

### *Relevance and Application*

Information extraction: Pulling specific pieces of information from a large volume of text.

Search engines: Making search results more relevant by understanding the entities in a query.

Chatbots: Enabling chatbots to understand and respond to user questions in a more human-like way.

Text summarization: Helping to identify key topics and entities for a concise summary.

Sentiment analysis: Providing more nuanced analysis by identifying entities associated with sentiment.

Knowledge management: Organizing and making information in documents like reports or customer feedback more accessible.