

**Modern Education Society's Wadia College of Engineering, Pune**  
**Department of Computer Engineering**

<b>Name of Student:</b>	<b>Class:</b>
<b>Semester/Year:</b>	<b>Roll No.:</b>
<b>Date of Performance:</b>	<b>Date of Submission:</b>
<b>Examined By:</b>	<b>Assignment No.:</b>

**Assignment No. 3**

**Title:** Write a program for Pre-processing of a Text Document: stop word removal

**Objectives:** To Implement stopword removal

**Problem Statement:** Write a program to Implement stop word removal

**Outcomes:** Student can understand how to Implement stop word removal

**Tools Required:**

**Hardware:**

**Software:** Open source operating system

**Theory:**

**Introduction:**

In natural language processing (NLP), stopwords are frequently filtered out to enhance text analysis and computational efficiency. Eliminating stopwords can improve the accuracy and relevance of NLP tasks by drawing attention to the more important words, or content words.

**Process of Stopword Removal**

The process generally follows these steps:

- **Tokenization:** First, the text is divided into individual words or tokens. For example, a sentence like *"The cat is sitting on the mat."* would be split into ['The', 'cat', 'is', 'sitting', 'on', 'the', 'mat'].
- **Stopword List:** A predefined list of stopwords is utilized. This list may vary depending on the language and specific application, but many programming libraries such as NLTK or SpaCy provide default lists for various languages.
- **Filtering:** Each word in the tokenized text is compared with the stopwords list, and if a match is found, the word is removed. After filtering, the sentence above would become ['cat', 'sitting', 'mat'], which better captures the key concepts.
- **Output:** The resulting text, with stopwords removed, is ready for further processing such as stemming, lemmatization, or vectorization for machine learning models.

## Tools for Stopword Removal:

Several libraries support stopwords removal:

- **NLTK (Natural Language Toolkit)**: Provides a comprehensive list of stopwords for multiple languages.
- **SpaCy**: A fast and modern NLP library that includes stopwords functionality.
- **Gensim**: Often used for topic modeling and also supports stopwords removal.
- **Scikit-learn**: Contains utilities for removing stop words during vectorization

**Conclusion:** Stopword removal is a fundamental preprocessing step in NLP and IR tasks. By eliminating words that do not contribute meaningful information, it enhances the efficiency and accuracy of downstream processes such as classification, clustering, and information retrieval. Customization of stopwords is essential for achieving optimal results, particularly in specialized domains.

## Questions:

Q1) What are stopwords?

Q2) Why is stopwords removal important in text preprocessing?

Q3) How does stopwords removal improve the performance of machine learning models?