



# Removing stop words with NLTK in Python

Last Updated : 03 Jan, 2024

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. The article aims to explore stopwords.

## Table of Content

- [What are Stop words?](#)
- [Need to remove the Stopwords](#)
- [Types of Stopwords](#)
- [Checking English Stopwords List](#)
- [Removing stop words with NLTK](#)
- [Removing stop words with SpaCy](#)
- [Removing stop words with Genism](#)
- [Removing stop words with SkLearn](#)

Certain words, like “the,” “and,” and “is,” are thought to be ineffective for communicating important information. The objective of eliminating words that add little or nothing to the comprehension of the text is to expedite text processing, even though the list of stopwords may differ.

## What are Stop words?

A [stop word](#) is a commonly used word (such as “the”, “a”, “an”, or “in”) that a search engine has been programmed to ignore, both when indexing entries for searching and when retrieving them as the result of a search query.

We would not want these words to take up space in our database or take up valuable processing time. For this, we can remove them easily, by storing a list

them in the nltk\_data directory.

`home/PratimaPython/nltk_data/corpora/stopwords` are the directory address. **or**  
(Do not forget to change your home directory name)

| Sample text with Stop Words                         | Without Stop Words                              |
|---|---|
| GeeksforGeeks – A Computer Science Portal for Geeks | GeeksforGeeks , Computer Science, Portal ,Geeks |
| Can listening be exhausting?                        | Listening, Exhausting                           |
| I like reading, so I read                           | Like, Reading, read                             |

## Need to remove the Stopwords

The necessity of removing stopwords in NLP is contingent upon the specific task at hand. For [text classification](#) tasks, where the objective is to categorize text into distinct groups, excluding stopwords is common practice. This is done to channel more attention towards words that truly convey the essence of the text. As illustrated earlier, certain words like “there,” “book,” and “table” contribute significantly to the text’s meaning, unlike less informative words such as “is” and “on.”

Conversely, for tasks like machine translation and text summarization, the removal of stopwords is not recommended. In these scenarios, every word plays a pivotal role in preserving the original meaning of the content.

## Types of Stopwords

Stopwords are frequently occurring words in a language that are frequently

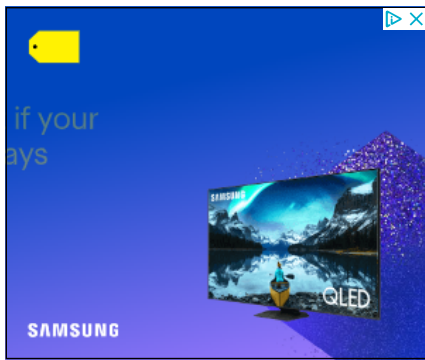
can change based on the language being studied and the context. The following is a broad list of stopword categories:

- **Common Stopwords:** These are the most frequently occurring words in a language and are often removed during text preprocessing. Examples include “the,” “is,” “in,” “for,” “where,” “when,” “to,” “at,” etc.
- **Custom Stopwords:** Depending on the specific task or domain, additional words may be considered as stopwords. These could be domain-specific terms that don’t contribute much to the overall meaning. For example, in a medical context, words like “patient” or “treatment” might be considered as custom stopwords.
- **Numerical Stopwords:** Numbers and numeric characters may be treated as stopwords in certain cases, especially when the analysis is focused on the meaning of the text rather than specific numerical values.
- **Single-Character Stopwords:** Single characters, such as “a,” “I,” “s,” or “x,” may be considered stopwords, particularly in cases where they don’t convey much meaning on their own.
- **Contextual Stopwords:** Words that are stopwords in one context but meaningful in another may be considered as contextual stopwords. For instance, the word “will” might be a stopword in the context of general language processing but could be important in predicting future events.

## Checking English Stopwords List

An English stopwords list typically includes common words that carry little semantic meaning and are often excluded during text analysis. Examples of these words are “the,” “and,” “is,” “in,” “for,” and “it.” These stopwords are frequently removed to focus on more meaningful terms when processing text data in natural language processing tasks such as text classification or sentiment analysis.

**To check the list of stopwords you can type the following commands in the python shell.**



[Click Here for More Info](#)

Ad By **Sponsor**

## Python3

```
import nltk
from nltk.corpus import stopwords

nltk.download('stopwords')
print(stopwords.words('english'))
```

### Output:

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours',
 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your',
 'yours', 'yourself',
 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her',
 'hers', 'herself',
 'it', "it's", 'its', 'itself', 'they', 'them', 'their', 'theirs',
 'themselves',
 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these',
 'those', 'am', 'is', 'are',
 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
 'do', 'does', 'did', 'doing',
 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until',
 'while', 'of', 'at', 'by', 'for',
 'with', 'about', 'against', 'between', 'into', 'through', 'during',
 'before', 'after', 'above', 'below', 'to',
 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under',
```

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#)

```
'each', 'few', 'more', 'most',
'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same',
'so', 'than', 'too', 'very',
's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've",
'now', 'd', 'll', 'm', 'o', 're', 've', 'y',
'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
'doesn', "doesn't",
'hadn', "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't",
'ma', 'mightn',
"mightn't", 'mustn', "mustn't", 'needn', "needn't", 'shan', "shan't",
'shouldn', "shouldn't",
'wasn', "wasn't", 'weren', "weren't", 'won', "won't", 'wouldn',
"wouldn't"]
```

**Note:** You can even modify the list by adding words of your choice in the English .txt. file in the stopwords directory.

## Removing stop words with NLTK

The following program removes stop words from a piece of text:

### Python3

```
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

example_sent = """This is a sample sentence,
                  showing off the stop words filtration."""

stop_words = set(stopwords.words('english'))

word_tokens = word_tokenize(example_sent)
# converts the words in word_tokens to lower case and then checks whether
#they are present in stop_words or not
filtered_sentence = [w for w in word_tokens if not w.lower() in stop_words]
#with no lower case conversion
filtered_sentence = []

for w in word_tokens:
    if w not in stop_words:
        filtered_sentence.append(w)
```

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#)

## Output:

```
['This', 'is', 'a', 'sample', 'sentence', ',', 'showing',  
'off', 'the', 'stop', 'words', 'filtration', '.']  
['This', 'sample', 'sentence', ',', 'showing', 'stop',  
'words', 'filtration', '.']
```

The provided Python code demonstrates stopword removal using the [Natural Language Toolkit](#) (NLTK) library. In the first step, the sample sentence, which reads “This is a sample sentence, showing off the stop words filtration,” is tokenized into words using the `word_tokenize` function. The code then filters out stopwords by converting each word to lowercase and checking its presence in the set of English stopwords obtained from NLTK. The resulting `filtered_sentence` is printed, showcasing both lowercased and original versions, providing a cleaned version of the sentence with common English stopwords removed.

## Removing stop words with SpaCy

---

### Python3

```
import spacy  
  
# Load spaCy English model  
nlp = spacy.load("en_core_web_sm")  
  
# Sample text  
text = "There is a pen on the table"  
  
# Process the text using spaCy  
doc = nlp(text)  
  
# Remove stopwords  
filtered_words = [token.text for token in doc if not token.is_stop]  
  
# Join the filtered words to form a clean text  
clean_text = ' '.join(filtered_words)  
  
print("Original Text:", text)
```

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#)

## Output:

Original Text: There is a pen on the table

Text after Stopword Removal: pen table

The provided Python code utilizes the [spaCy](#) library for natural language processing to remove stopwords from a sample text. Initially, the spaCy English model is loaded, and the sample text, "There is a pen on the table," is processed using spaCy. Stopwords are then filtered out from the processed tokens, and the resulting non-stopword tokens are joined to create a clean version of the text.

## Removing stop words with Genism

---

### Python3

```
from gensim.parsing.preprocessing import remove_stopwords

# Another sample text
new_text = "The majestic mountains provide a breathtaking view."

# Remove stopwords using Gensim
new_filtered_text = remove_stopwords(new_text)

print("Original Text:", new_text)
print("Text after Stopword Removal:", new_filtered_text)
```

## Output:

Original Text: The majestic mountains provide a breathtaking view.

Text after Stopword Removal: The majestic mountains provide breathtaking view.

The provided Python code utilizes [Gensim's](#) `remove_stopwords` function to preprocess a sample text. In this specific example, the original text is "The

efficiently eliminates common English stopwords, resulting in a filtered version of the text, which is then printed alongside the original text.

## Removing stop words with SkLearn

### Python3

```
from sklearn.feature_extraction.text import CountVectorizer
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

# Another sample text
new_text = "The quick brown fox jumps over the lazy dog."

# Tokenize the new text using NLTK
new_words = word_tokenize(new_text)

# Remove stopwords using NLTK
new_filtered_words = [
    word for word in new_words if word.lower() not in stopwords.words('english')]

# Join the filtered words to form a clean text
new_clean_text = ' '.join(new_filtered_words)

print("Original Text:", new_text)
print("Text after Stopword Removal:", new_clean_text)
```

### Output:

```
Original Text: The quick brown fox jumps over the lazy dog.
Text after Stopword Removal: quick brown fox jumps lazy dog .
```

The provided Python code combines [scikit-learn](#) and NLTK for stopword removal and text processing. First, the sample text, “The quick brown fox jumps over the lazy dog,” is tokenized into words using NLTK’s [word\\_tokenize](#) function. Subsequently, common English stopwords are removed by iterating through the tokenized words and checking their absence in the NLTK stopwords set. The final step involves joining the non-stopword tokens to

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#)



[CountVectorizer](#) could be utilized for further text analysis, such as creating a bag-of-words representation.

Also Check:

- [Natural Language Processing.\(NLP\) Tutorial](#)
- [Python | Lemmatization with NLTK](#)
- [Introduction to Stemming](#)

## Frequently Asked Questions

### 1. What are stopwords in NLP?

*Stopwords are common words in a language that are often filtered out during NLP tasks because they are considered to carry little meaning or contribute minimally to the overall understanding of a text. Examples include “the,” “is,” “and,” “in,” etc.*

### 2. Why do we remove stopwords in NLP?

*Stopwords are removed in NLP to focus on the more meaningful and informative words in a text. This is often done to reduce noise, improve efficiency in processing, and highlight keywords that carry the essential meaning of the text.*

### 3. Can the list of stopwords vary between NLP libraries?

*Yes, the list of stopwords can vary between NLP libraries. Different libraries, such as NLTK, spaCy, Gensim, and scikit-learn, may provide their own sets of stopwords. Users can also customize the list based on their specific needs.*

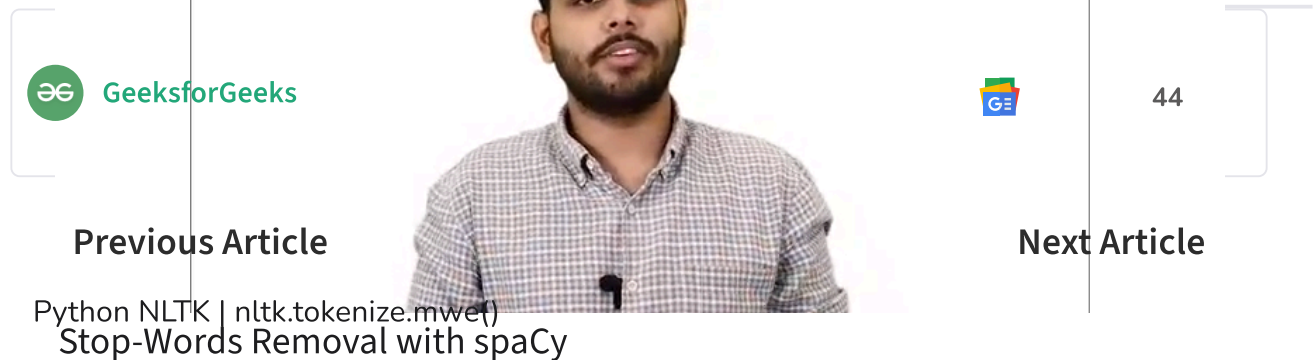
*Stopwords should be retained when they contribute to the context, coherence, or specific linguistic features of the text. Tasks like machine translation, text summarization, and certain sentiment analysis scenarios may benefit from retaining stopwords.*

## 5. What are some common English stopwords?

*Common English stopwords include “the,” “and,” “is,” “in,” “for,” “where,” “when,” “to,” “at,” etc. These words are frequently used but are often excluded during text analysis.*

Are you passionate about data and looking to make one giant leap into your career? Our [Data Science Course](#) will help you change your game and, most importantly, allow students, professionals, and working adults to tide over into the data science immersion. Master state-of-the-art methodologies, powerful tools, and industry best practices, hands-on projects, and real-world applications. Become the executive head of industries related to **Data Analysis**, **Machine Learning**, and **Data Visualization** with these growing skills. Ready to Transform Your Future? ***Enroll Now to Be a Data Science Expert!***

---



## Similar Reads

### Removing stop words with NLTK in Python

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...

9 min read

### Python NLTK | nltk.tokenizer.word\_tokenize()

With the help of nltk.tokenize.word\_tokenize() method, we are able to extract the tokens from string of characters by using tokenize.word\_tokenize() method. It actually returns the syllables from a single word. A...

1 min read

### Correcting Words using NLTK in Python

nltk stands for Natural Language Toolkit and is a powerful suite consisting of libraries and programs that can be used for statistical natural language processing. The libraries can implement tokenization, classification,...

4 min read

### Stop Word Removal In R

In Natural Language Processing, different words carry different amounts of information. We want to select only those words which are meaningful to the machine learning models. By analyzing the text, some words...

9 min read

### How To Remove Nltk From Python

In Python, NLTK, or Natural Language Toolkit, is a powerful library that is used for human language data. This

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#)

## Tokenize text using NLTK in python

To run the below python program, (NLTK) natural language toolkit has to be installed in your system. The NLTK module is a massive tool kit, aimed at helping you with the entire Natural Language Processing (NLP)...

3 min read

---

## Python NLTK | tokenize.regexp()

With the help of NLTK tokenize.regexp() module, we are able to extract the tokens from string by using regular expression with RegexpTokenizer() method. Syntax : tokenize.RegexpTokenizer() Return : Return arra...

1 min read

---

## How to remove punctuations in NLTK

Natural Language Processing (NLP) involves the manipulation and analysis of natural language text by machines. One essential step in preprocessing text data for NLP tasks is removing punctuations. In this articl...

4 min read

---

## Python - Compute the frequency of words after removing stop words and stemming

In this article we are going to tokenize sentence, paragraph, and webpage contents using the NLTK toolkit in the python environment then we will remove stop words and apply stemming on the contents of sentences,...

8 min read

---

## Python NLTK | nltk.TweetTokenizer()

With the help of NLTK nltk.TweetTokenizer() method, we are able to convert the stream of words into small tokens so that we can analyse the audio stream with the help of nltk.TweetTokenizer() method. Syntax :...

1 min read

---

## Python NLTK | nltk.WhitespaceTokenizer

With the help of nltk.tokenize.WhitespaceTokenizer() method, we are able to extract the tokens from string of words or sentences without whitespaces, new line and tabs by using tokenize.WhitespaceTokenizer() metho...

1 min read

---

## Python NLTK | nltk.tokenize.SpaceTokenizer()

With the help of nltk.tokenize.SpaceTokenizer() method, we are able to extract the tokens from string of words on the basis of space between them by using tokenize.SpaceTokenizer() method. Syntax :...

1 min read

---

## Python NLTK | nltk.tokenize.SExprTokenizer()

With the help of nltk.tokenize.SExprTokenizer() method, we are able to extract the tokens from string of characters or numbers by using tokenize.SExprTokenizer() method. It actually looking for proper brackets to

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#)

## Python NLTK | nltk.tokenize.TabTokenizer()

With the help of nltk.tokenize.TabTokenizer() method, we are able to extract the tokens from string of words on the basis of tabs between them by using tokenize.TabTokenizer() method. Syntax : tokenize.TabTokenizer(...

1 min read

---

## Python NLTK | tokenize.WordPunctTokenizer()

With the help of nltk.tokenize.WordPunctTokenizer()() method, we are able to extract the tokens from string of words or sentences in the form of Alphabetic and Non-Alphabetic character by using...

1 min read

---

## NLP | Expanding and Removing Chunks with RegEx

RegexpParser or RegexpChunkRule.fromstring() doesn't support all the RegexpChunkRule classes. So, we need to create them manually. This article focusses on 3 of such classes : ExpandRightRule: It adds chunk...

2 min read

---

## Top K Nearest Words using Edit Distances in NLP

We can find the Top K nearest matching words to the given query input word by using the concept of Edit/Levenshtein distance. If any word is the same word as the query input(word) then their Edit distance wou...

3 min read

---

## Python NLTK | nltk.tokenize.StanfordTokenizer()

With the help of nltk.tokenize.StanfordTokenizer() method, we are able to extract the tokens from string of characters or numbers by using tokenize.StanfordTokenizer() method. It follows stanford standard for...

1 min read

---

## Python NLTK | nltk.tokenize.mwe()

With the help of NLTK nltk.tokenize.mwe() method, we can tokenize the audio stream into multi\_word expression token which helps to bind the tokens with underscore by using nltk.tokenize.mwe() method....

1 min read

---

### Article Tags :

[AI-ML-DS](#)[NLP](#)[AI-ML-DS With Python](#)[Python-nltk](#)

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#)

Tower, Sector- 136, Noida, Uttar Pradesh  
(201305) | Registered Address:- K 061,  
Tower K, Gulshan Vivante Apartment,  
Sector 137, Noida, Gautam Buddh  
Nagar, Uttar Pradesh, 201305



## Company

About Us  
Legal  
In Media  
Contact Us  
Advertise with us  
GFG Corporate Solution  
Placement Training Program  
GeeksforGeeks Community

## DSA

Data Structures  
Algorithms  
DSA for Beginners  
Basic DSA Problems  
DSA Roadmap  
Top 100 DSA Interview Problems  
DSA Roadmap by Sandeep Jain  
All Cheat Sheets

## Web Technologies

HTML  
CSS  
JavaScript  
TypeScript  
ReactJS  
NextJS  
Bootstrap  
Web Design

## Computer Science

Operating Systems

## Languages

Python  
Java  
C++  
PHP  
GoLang  
SQL  
R Language  
Android Tutorial  
Tutorials Archive

## Data Science & ML

Data Science With Python  
Data Science For Beginner  
Machine Learning  
ML Maths  
Data Visualisation  
Pandas  
NumPy  
NLP  
Deep Learning

## Python Tutorial

Python Programming Examples  
Python Projects  
Python Tkinter  
Web Scraping  
OpenCV Tutorial  
Python Interview Question  
Django

## DevOps

Cit

We use cookies to ensure you have the best browsing experience on our website. By using our site, you acknowledge that you have read and understood our [Cookie Policy](#) & [Privacy Policy](#).

Software Engineering  
Digital Logic Design  
Engineering Maths  
Software Development  
Software Testing

Docker  
Kubernetes  
Azure  
GCP  
DevOps Roadmap

### System Design

High Level Design  
Low Level Design  
UML Diagrams  
Interview Guide  
Design Patterns  
OOAD  
System Design Bootcamp  
Interview Questions

### Interview Preparation

Competitive Programming  
Top DS or Algo for CP  
Company-Wise Recruitment Process  
Company-Wise Preparation  
Aptitude Preparation  
Puzzles

### School Subjects

Mathematics  
Physics  
Chemistry  
Biology  
Social Science  
English Grammar  
Commerce  
World GK

### GeeksforGeeks Videos

DSA  
Python  
Java  
C++  
Web Development  
Data Science  
CS Subjects

@GeeksforGeeks, Sanchhaya Education Private Limited, All rights reserved