10 NOV 24 | DAY - 71 | Natural Language Processing

# #100DAYSOFDATA SCIENCE

PYTHON | SQL | STATISTICS | MACHINE LEARNING |

# Stopwords

**What Are Stopwords?**

Stopwords are common words in a language, such as "and," "the," "is," "in," and "at," that often appear frequently in text data but add little meaning or context. In NLP, removing stopwords is a standard preprocessing step to help models focus on the main content of the text, making data analysis more efficient and improving the performance of algorithms by eliminating these "filler" words.

**Key Aspects of Stopwords:**

1. **Efficiency**: By removing frequent but uninformative words, text data becomes easier and faster to process.
2. **Focus**: Filtering out stopwords enables models to focus on more meaningful words that contribute to the context and insights.
3. **Use Cases**: Stopwords removal is crucial in tasks like sentiment analysis, search engines, and text summarization, where high-frequency words can dilute the meaningful content.

**Common Libraries for Stopword Removal:**

1. **NLTK**: Provides a built-in list of stopwords for many languages, allowing you to customize and expand the list for specific needs.
2. **SpaCy**: Includes a customizable set of stopwords that can be easily adjusted by adding or removing words as per project requirements.

**Sample List of Common Stopwords in English:**

- **Domain-Specific Stopwords**: Some NLP projects may define custom stopwords, such as "data," "science," or "technology," depending on the analysis requirements.

```
English Stopwords:
['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', 'again', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', '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"]
------------------------------------------------------------
```

**Sample List of Common Stopwords in Arabic:**

```
Arabic Stopwords:
['اللذان', 'اللتيا', 'اللتان', 'اللائي', 'اللاتي', 'اللاي', 'الذين', 'الذي', 'التي', 'إلا', 'ألا', 'أكثر', 'أقل', 'أف', 'إذن', 'إذما', 'إذا', 'إذ
'أنت', 'أنتما', 'أنتم', 'أنت', 'أن', 'إنا', 'إن', 'إما', 'أما', 'أن', 'إليكن', 'إليكما', 'إليكم', 'إليك', 'إلى', 'اللواتي', 'اللذين
'بع', 'بس', 'بخ', 'إيه', 'أينما', 'أين', 'أي', 'أيها', 'أيا', 'آي', 'أوه', 'أولاء', 'أولئك', 'آه', 'آها', 'أو', 'أنى', 'أنت', 'إنه', 'إنما', 'إن
'تل', 'بيد', 'بين', 'بين', 'بهما', 'بهن', 'بها', 'بهم', 'بنا', 'بمن', 'بماذا', 'بما', 'بلى', 'بل', 'بكن', 'بكما', 'بكم', 'بك', 'بعض، 'بع
'ذان', 'ذا', 'ذات', 'ذاك', 'ذا', 'دون', 'خلا', 'حين', 'حيثما', 'حيث', 'حتى', 'حبذا', 'جعل', 'ثمة', 'ثم', 'تينك', 'تين', 'تي', 'ته', 'تلكما
'على', 'عل', 'عسى', 'شتان', 'سوى', 'سوف', 'ريث', 'ذينك', 'ذين', 'ذي', 'ذواتا', 'ذواتي', 'ذوا', 'ذو', 'ذه', 'ذلكن', 'ذلكما', 'ذلكم', 'ذلك', 'ذانك
'كأن', 'كأنما', 'كأي', 'كأين', 'قد', 'فيها', 'فيه', 'فيما', 'فيم', 'في', 'فمن', 'فلا', 'فإن', 'فإذا', 'غير', 'عند', 'عن', 'عما', 'عليه', 'عليك
'لدى', 'لاسيما', 'لا', 'كيتما', 'كيت', 'كيف', 'كما', 'كم', 'كليهما', 'كلتا', 'كلا', 'كلاهما', 'كل', 'كذلك', 'كذا', 'كي
'لها', 'له', 'لنا', 'لن', 'لما', 'لكيلا', 'لكي', 'لكنما', 'لكن', 'لكم', 'لكما', 'لك', 'لعل', 'لسنا', 'لست', 'لستم', 'لستن
'هذه', 'هذان', 'هذا', 'هاهنا', 'هاك', 'هاتين', 'هاتي', 'هاته', 'هاتان', 'ها', 'نعم', 'نحو', 'نحن', 'مهما', 'مه', 'منها', 'منذ', 'من', 'مع', 'مما', 'مذ', 'متى', 'ما، 'ماذا', 'مليار', 'ليسوا', 'ليستا', 'ليست', 'ليسا', 'ليس', 'ليت', 'لي', 'لن', 'لوما', 'لولا', 'لو', 'لهن', 'لهما', 'لهم
'هيا', 'هيت', 'هيهات', 'والذي', 'والذين', 'وإذ', 'وإذا', 'وإن', 'ولا', 'ولكن', 'وما', 'ومن', 'وهو', 'وا', 'ويك', 'وي، 'هنا', 'هناك', 'هنالك', 'هو', 'هؤلاء', 'هلا', 'هل', 'هكذا', 'هذي', 'هذين
'شباط', 'آذار', 'نيسان', 'أيار', 'حزيران', 'تموز', 'آب', 'أيلول', 'تشرين', 'دولار', 'دينار', 'ريال', 'درهم', 'جنيه', 'مليم', 'فلس', 'قرش، 'يوليو', 'أغسطس', 'سبتمبر', 'أكتوبر', 'نوفمبر', 'ديسمبر', 'جانفي', 'فيفري', 'مارس', 'أفريل', 'ماي', 'جوان', 'جويلية', 'أوت', 'كانون
'ثان', 'ثان', 'ثالث', 'رابع', 'خامس', 'سادس', 'سابع', 'ثامن', 'تاسع', 'عاشر', 'حادي', 'ثاني', 'أحد', 'اثنين', 'ثلاثاء', 'أربعاء', 'خميس', 'جمعة', 'أول', 'ثان
'إحد', 'اثنا', 'اثني', 'ثلاث', 'أربع', 'خمس', 'ست', 'سبع', 'ثمان', 'تسع', 'عشر', 'ثمان', 'أحد', 'سبت', 'اثنين، 'ثلاثة', 'أربعة', 'خمسة', 'ستة', 'سبعة', 'ثمانية', 'تسعة', 'عشرة، 'واحد', 'شيكل', 'يوان', 'بن', 'يورو', 'تيم
'خ', 'حاء', 'جيم', 'ثاء', 'تاء', 'باء', 'ألف', 'ة', 'أ', 'و', 'ن', 'م', 'ل', 'ك', 'ق', 'ف', 'غ', 'ع', 'ظ', 'طا، 'م', 'ص', 'ز', 'ر', 'ذ', 'د', 'ح', 'ج', 'ث', 'ت', 'ب', 'أ', 'ي', 'و', 'ها', 'ع، 'سادس', 'خامس', 'رابع
'يا', 'واو', 'هاء', 'نون', 'ميم', 'لام', 'كاف', 'قاف', 'فاء', 'غين', 'عين', 'ظاء', 'طاء', 'ضاد', 'صاد', 'شين', 'سين', 'زاي', 'راء', 'ذال', 'دال
'هذه', 'هذان', 'هذا', 'هاتين', 'هاتي', 'هاته', 'هاتان', 'هؤلاء', 'ذين', 'ذي', 'ذه', 'ثمة، 'تين', 'تي', 'ته', 'تانك', 'تان', 'إياي', 'إياكن', 'إياك', 'إياكما', 'إياكم', 'إياك', 'إياهن', 'إياها', 'إياهم', 'إياهما', 'أولالك
'آه', 'آم', 'آو', 'آمين', 'وا', 'فلان', 'بضع', 'كأيّن', 'كأيّ', 'ذيت', 'كأنّ', 'أيّان', 'أيّ', 'أنّى', 'أيّ', 'ذان، 'ذه', 'ذي', 'ذيّن', 'هؤلاء', 'ذين', 'ذي', 'ذه', 'ثمّة', 'تين', 'تي', 'ته', 'تانك، 'اللائي', 'اللاء، 'الألى
'حتّ', 'حيّ', 'حذار', 'حاي', 'تلك', 'بطآن', 'بتّ', 'بسّ', 'بج', 'أوّه', 'آيّ', 'أيّ', 'أفّ، 'أمّ', 'إلّ', 'إلَيك', 'إلَيك، 'إلَيك، 'إذ، 'أبو', 'أمامك', 'أمام', 'آهّ
'تر', 'تخذ', 'اتخذ', 'الفيّ', 'اتخذ', 'تفعلين', 'تفعلون', 'يفعلون', 'تفعلان', 'يفعلان', 'وقّ', 'وشكان', 'وراءك', 'واها', 'وا', 'هتّهات', 'شتّان', 'سرعان', 'رويدك، 'دونك
'ك', 'تعلّم', 'جعل', 'حجا', 'حبيب', 'خال', 'حسّ', 'خال', 'حسب', 'درى', 'رأى', 'زعم', 'كسا، 'أخبر', 'أرى', 'أعلم', 'عدّ', 'علم', 'طفق', 'ظنّ', 'صبر', 'صيّر', 'غادر', 'علم، 'ذهب', 'هبّ', 'وجد', 'ورد', 'وهب', 'أسكن
'كي', 'كلّا', 'كأنّ', 'في', 'علّ', 'علا', 'رّب', 'جلل', 'جير', 'ب', 'نتّا', 'أفعل به', 'حذّ', 'خبّر', 'أخبر', 'أيا', 'أعلم، 'أرى', 'إلى', 'أنّ', 'إنّ', 'إلّا', 'إذا', 'أجل', 'ء', 'آنّا، 'ألا', 'لكن', 'لعلّ', 'أطعم', 'أعطى', 'رزق', 'زود', 'سقى', 'كسا', 'نبّأ، 'قلم
'سبحان', 'حسّ', 'جميع', 'تلقاء', 'تجاه', 'ا', 'ه', 'و', 'لقا', 'ن', 'ت', 'ك', 'إلّا', 'وا', 'الّ', 'هلّا', 'نّ', 'مّ', 'الكنّ، 'م
'ثمانو', 'تسعملة', 'تمنملة', 'سبعملة', 'ستملة', 'خمسملة', 'أربعملة', 'ثلاثملة', 'ملتان', 'ملة', 'قو', 'حمو', 'أخو', 'أبو', 'معاذ', 'مثل', 'العمر', 'شبه، 'م
'ثمانو', 'سبعون', 'ستون', 'خمسون', 'أربعون', 'ثلاثون', 'تسعين', 'ستين', 'سبعين', 'ثمانين', 'تسعون', 'عشرون', 'تسعمائة', 'ثمانمائة', 'سبعمائة', 'ستمائة', 'خمسمائة', 'أربعمائة', 'ثلاثمائة', 'مائة، 'ا', 'أصلا', 'أهلا', 'أيضا', 'بؤسا', 'بعدا', 'حقا', 'تعسا', 'خلافا، 'حمدا', 'خاصة', 'دواليك', 'سحقا', 'سرا', 'سمعا', 'صبرا', 'صدقا', 'صراحة، 'لا سيم
'أ', 'آناء', 'آنفا', 'أمس', 'الآن', 'أمد', 'الآن', 'إصلا', 'إزاء', 'أبدا', 'معاذ', 'البتك', 'كثيرا', 'قاطبة', 'فضلا', 'فصلا', 'فرادى', 'غاليا', 'عيانا', 'عجبا', 'طرا
'أم', 'خلف', 'قبل', 'مرّة', 'لدن', 'لمّا', 'كلّما', 'قطّ', 'غدا', 'عوض', 'صحوة', 'مساء', 'صباح', 'صبح', 'حقا', 'ثمّ', 'ثمّة', 'نمّ', 'تارة', 'أثناء، 'أوّل
'ظ', 'صار', 'راح', 'جار', 'رجع', 'تحوّل', 'تبدّل', 'بات', 'انقلب', 'انفك', 'آض', 'أمسى', 'أصبح', 'أضحى', 'استحال', 'ارتدّ', 'شمال', 'يمين', 'تحت', 'فوق
'هبّ', 'كاد', 'كرب', 'قام', 'علق', 'طفق', 'شرع، 'جرى', 'أوشك', 'جعل', 'أنشأ', 'أقبل', 'انبرى', 'أقبل', 'اخلولق', 'أخذ', 'ابتدأ', 'مافتئ', 'مازال', 'مادام', 'ما برح', 'ما انفك', 'كان', 'غدا', 'عاد
-------------------------------------------
```

**Sample List of Common Stopwords in French:**

```
French Stopwords:
['au', 'aux', 'avec', 'ce', 'ces', 'dans', 'de', 'des', 'du', 'elle', 'en', 'et', 'eux', 'il', 'ils', 'je', 'la', 'le', 'les', 'leur', 'lui', 'ma', 'mai
s', 'me', 'même', 'mes', 'moi', 'mon', 'ne', 'nos', 'notre', 'nous', 'on', 'ou', 'par', 'pas', 'pour', 'qu', 'que', 'qui', 'sa', 'se', 'ses', 'son', 'su
r', 'ta', 'te', 'tes', 'toi', 'ton', 'tu', 'un', 'une', 'vos', 'votre', 'vous', 'c', 'd', 'j', 'l', 'à', 'm', 'n', 's', 't', 'y', 'été', 'étée', 'étées',
'étés', 'étant', 'étante', 'étants', 'étantes', 'suis', 'es', 'est', 'sommes', 'êtes', 'sont', 'serai', 'seras', 'sera', 'serons', 'serez', 'seront', 'se
rais', 'serait', 'serions', 'seriez', 'seraient', 'étais', 'était', 'étions', 'étiez', 'étaient', 'fus', 'fut', 'fûmes', 'fûtes', 'furent', 'sois', 'soi
t', 'soyons', 'soyez', 'soient', 'fusse', 'fusses', 'fût', 'fussions', 'fussiez', 'fussent', 'ayant', 'ayante', 'ayantes', 'ayants', 'eu', 'eue', 'eues',
'eus', 'ai', 'as', 'avons', 'avez', 'ont', 'aurai', 'auras', 'aura', 'aurons', 'aurez', 'auront', 'aurais', 'aurait', 'aurions', 'auriez', 'auraient', 'a
vais', 'avait', 'avions', 'aviez', 'avaient', 'eut', 'eûmes', 'eûtes', 'eurent', 'aie', 'aies', 'ait', 'ayons', 'ayez', 'aient', 'eusse', 'eusses', 'eû
t', 'eussions', 'eussiez', 'eussent']
```

**Types of Stopwords and Their Impact:**
1. **High-Frequency Words**: Words that occur frequently across documents but add minimal meaning. Removing these helps reduce text noise.
2. **Domain-Specific Stopwords**: Some stopwords may be tailored for specific industries or contexts, such as excluding "data" in data science projects, which can improve focus on more relevant terms.

**When to Use Stopword Removal:**
- **Ideal For**: Tasks where meaningful keywords are essential, like **topic modeling** or **document classification**.
- **Avoid For**: Applications where sentence structure is important, like **text generation** or **language translation**, where stopwords contribute to the natural flow.

**Selecting Stopwords:** Choosing the right stopwords list depends on the NLP project's context. For instance, legal documents might require a custom list, while social media texts might include common slang and abbreviations.

# ing-the-noise-for-better-analysis

November 10, 2024

```
[87]: corpus = '''Natural language processing is an exciting area of artificial␣
       ↪intelligence that focuses on enabling computers to understand and respond to␣
       ↪human language. By applying techniques like tokenization, stemming, and␣
       ↪lemmatization, NLP systems can break down sentences into their core␣
       ↪components, allowing computers to process language in a way that's closer to␣
       ↪human understanding. However, not all words contribute equally to meaning,␣
       ↪so stopwords such as 'is,' 'and,' 'the,' and 'by' are often filtered out to␣
       ↪improve processing efficiency. As the technology advances, applications of␣
       ↪NLP continue to expand, helping us with tasks ranging from simple text␣
       ↪summarization to complex sentiment analysis.'''

       print(corpus)
```

Natural language processing is an exciting area of artificial intelligence that
focuses on enabling computers to understand and respond to human language. By
applying techniques like tokenization, stemming, and lemmatization, NLP systems
can break down sentences into their core components, allowing computers to
process language in a way that's closer to human understanding. However, not all
words contribute equally to meaning, so stopwords such as 'is,' 'and,' 'the,'
and 'by' are often filtered out to improve processing efficiency. As the
technology advances, applications of NLP continue to expand, helping us with
tasks ranging from simple text summarization to complex sentiment analysis.

```
[89]: from nltk.corpus import stopwords
```

```
[91]: import nltk
       nltk.download("stopwords")
```

```
[nltk_data] Downloading package stopwords to
[nltk_data]     C:\Users\Zahid.Shaikh\AppData\Roaming\nltk_data…
[nltk_data]   Package stopwords is already up-to-date!
```

```
[91]: True
```

```
[93]: print("English Stopwords: ")
       print(stopwords.words('english'))
       print('-' * 153)
       print("Arabic Stopwords: ")
```

```
print(stopwords.words('arabic'))
print('-' * 153)
print("French Stopwords: ")
print(stopwords.words('french'))
```

English Stopwords:
['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',
'again', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why',
'how', 'all', 'any', 'both', '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"]
--------------------------------------------------------------------------------
-------------------------------------------------------------------------
Arabic Stopwords:
[' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
' ', ' ', ' ', ' ', ' ', ' ', ' ', ' ',
```

```
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',
'      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      ',  '      '[
--------------------------------------------------------------------------------
--------------------------------------------------------------------------------
French Stopwords:
['au', 'aux', 'avec', 'ce', 'ces', 'dans', 'de', 'des', 'du', 'elle', 'en',
'et', 'eux', 'il', 'ils', 'je', 'la', 'le', 'les', 'leur', 'lui', 'ma', 'mais',
'me', 'même', 'mes', 'moi', 'mon', 'ne', 'nos', 'notre', 'nous', 'on', 'ou',
'par', 'pas', 'pour', 'qu', 'que', 'qui', 'sa', 'se', 'ses', 'son', 'sur', 'ta',
'te', 'tes', 'toi', 'ton', 'tu', 'un', 'une', 'vos', 'votre', 'vous', 'c', 'd',
'j', 'l', 'à', 'm', 'n', 's', 't', 'y', 'été', 'étée', 'étées', 'étés', 'étant',
'étante', 'étants', 'étantes', 'suis', 'es', 'est', 'sommes', 'êtes', 'sont',
'serai', 'seras', 'sera', 'serons', 'serez', 'seront', 'serais', 'serait',
'serions', 'seriez', 'seraient', 'étais', 'était', 'étions', 'étiez', 'étaient',
'fus', 'fut', 'fûmes', 'fûtes', 'furent', 'sois', 'soit', 'soyons', 'soyez',
'soient', 'fusse', 'fusses', 'fût', 'fussions', 'fussiez', 'fussent', 'ayant',
'ayante', 'ayantes', 'ayants', 'eu', 'eue', 'eues', 'eus', 'ai', 'as', 'avons',
'avez', 'ont', 'aurai', 'auras', 'aura', 'aurons', 'aurez', 'auront', 'aurais',
'aurait', 'aurions', 'auriez', 'auraient', 'avais', 'avait', 'avions', 'aviez',
'avaient', 'eut', 'eûmes', 'eûtes', 'eurent', 'aie', 'aies', 'ait', 'ayons',
'ayez', 'aient', 'eusse', 'eusses', 'eût', 'eussions', 'eussiez', 'eussent']
```

```python
from nltk.stem import PorterStemmer
stemmer = PorterStemmer()
sentences = nltk.sent_tokenize(corpus)
print(type(sentences))
print("Before Porter Stemmer:\n", sentences)
print('-' * 153)
for i in range(len(sentences)):
    words = nltk.word_tokenize(sentences[i])
    words = [stemmer.stem(word) for word in words if word not in set(stopwords.
 ↪words('english'))]
    sentences[i] = ' '.join(words)
print("After Porter Stemmer:\n", sentences)
```

```
<class 'list'>
Before Porter Stemmer:
 ['Natural language processing is an exciting area of artificial intelligence
that focuses on enabling computers to understand and respond to human
language.', 'By applying techniques like tokenization, stemming, and
lemmatization, NLP systems can break down sentences into their core components,
```

allowing computers to process language in a way that's closer to human understanding.', "However, not all words contribute equally to meaning, so stopwords such as 'is,' 'and,' 'the,' and 'by' are often filtered out to improve processing efficiency.", 'As the technology advances, applications of NLP continue to expand, helping us with tasks ranging from simple text summarization to complex sentiment analysis.']
--------------------------------------------------------------------------------
--------------------------------------------------------------------------------
After Porter Stemmer:
 ['natur languag process excit area artifici intellig focus enabl comput understand respond human languag .', 'by appli techniqu like token , stem , lemmat , nlp system break sentenc core compon , allow comput process languag way ' closer human understand .', "howev , word contribut equal mean , stopword 'i , ' 'and , ' 'the , ' 'bi ' often filter improv process effici .", 'as technolog advanc , applic nlp continu expand , help us task rang simpl text summar complex sentiment analysi .']

```
[97]: from nltk.stem import LancasterStemmer
      stemmer = LancasterStemmer()
      sentences = nltk.sent_tokenize(corpus)
      print(type(sentences))
      print("Before Lancaster Stemmer:\n", sentences)
      print('-' * 153)
      for i in range(len(sentences)):
          words = nltk.word_tokenize(sentences[i])
          words = [stemmer.stem(word) for word in words if word not in set(stopwords.
       ↪words('english'))]
          sentences[i] = ' '.join(words)
      print("After Lancaster Stemmer:\n", sentences)
```

<class 'list'>
Before Lancaster Stemmer:
 ['Natural language processing is an exciting area of artificial intelligence that focuses on enabling computers to understand and respond to human language.', 'By applying techniques like tokenization, stemming, and lemmatization, NLP systems can break down sentences into their core components, allowing computers to process language in a way that's closer to human understanding.', "However, not all words contribute equally to meaning, so stopwords such as 'is,' 'and,' 'the,' and 'by' are often filtered out to improve processing efficiency.", 'As the technology advances, applications of NLP continue to expand, helping us with tasks ranging from simple text summarization to complex sentiment analysis.']
--------------------------------------------------------------------------------
--------------------------------------------------------------------------------
After Lancaster Stemmer:
 ['nat langu process excit are art intellig focus en comput understand respond hum langu .', 'by apply techn lik tok , stem , lem , nlp system break sent cor compon , allow comput process langu way ' clos hum understand .', "howev , word

```
contribut eq mean , stopword 'is , ' 'and , ' 'the , ' 'by ' oft filt improv
process efficy .", 'as technolog adv , apply nlp continu expand , help us task
rang simpl text summ complex senty analys .']
```

[99]:
```python
from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
sentences = nltk.sent_tokenize(corpus)
print(type(sentences))
print("Before Lemmatizer:\n", sentences)
print('-' * 153)
for i in range(len(sentences)):
    words = nltk.word_tokenize(sentences[i])
    words = [lemmatizer.lemmatize(word.lower(), pos='v') for word in words if␣
  ↪word not in set(stopwords.words('english'))]
    sentences[i] = ' '.join(words)
print("After Lemmatizer:\n", sentences)
```

```
<class 'list'>
Before Lemmatizer:
 ['Natural language processing is an exciting area of artificial intelligence
that focuses on enabling computers to understand and respond to human
language.', 'By applying techniques like tokenization, stemming, and
lemmatization, NLP systems can break down sentences into their core components,
allowing computers to process language in a way that's closer to human
understanding.', "However, not all words contribute equally to meaning, so
stopwords such as 'is,' 'and,' 'the,' and 'by' are often filtered out to improve
processing efficiency.", 'As the technology advances, applications of NLP
continue to expand, helping us with tasks ranging from simple text summarization
to complex sentiment analysis.']
---------------------------------------------------------------------------------
---------------------------------------------------------------------------

After Lemmatizer:
 ['natural language process excite area artificial intelligence focus enable
computers understand respond human language .', 'by apply techniques like
tokenization , stem , lemmatization , nlp systems break sentence core components
, allow computers process language way ' closer human understand .', "however ,
word contribute equally mean , stopwords 'is , ' 'and , ' 'the , ' 'by ' often
filter improve process efficiency .", 'as technology advance , applications nlp
continue expand , help us task range simple text summarization complex sentiment
analysis .']
```

Made with   by Zahid Salim Shaikh

[ ]: