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', 'اللتين', 'اللتان', 'اللاتي', 'اللائي', 'الذين', 'الذي', 'التي', 'إلا', 'ألا', 'أكثر', 'أقل', 'أف', 'إذن', 'إذما', 'إذا', 'إذ
'أنت', 'أنتما', 'أنتم', 'أنت', 'إليكن', 'إليكم', 'إليك', 'إلى', 'اللواتي', 'اللذين
'بع', 'بس', 'بخ', 'إيه', 'أينما', 'أين', 'أي', 'آي', 'أوه', 'أولئك', 'أولاء', 'أو', 'آه', 'آها', 'أنى', 'أنه', 'إنه، 'إنما', 'إن، 'إنا', 'إن، 'أن، 'إما', 'أما', 'أما', 'أم، 'إليكما
'تل', 'بيد', 'بين', 'بي', 'بهما', 'بهم', 'بها', 'بمن، 'بماذا', 'بما', 'بنا، 'بلى', 'بل، 'بكن', 'بكما', 'بكم', 'بك، 'بعض
'ذان', 'ذات، 'ذا، 'إذا، 'دون، 'خلا، 'حينما', 'حيثما', 'حيث', 'حتى، 'حبذا', 'جير، 'ثمة، 'ثم، 'حاشا', 'تينك، 'تين، 'تي، 'ته، 'تلكما', 'تلكم
'على', 'عسى', 'شتان، 'عدا، 'فيها، 'فيم، 'فيما، 'في، 'فمن، 'فلا، 'فإن، 'فإذا', 'غير، 'عند، 'عن، 'عما', 'عليه', 'عليك
'لدى', 'لاسيما', 'لا، 'كيفما', 'كيت، 'كي، 'كما، 'كم، 'كليهما', 'كلتا، 'كلاهما', 'كلا، 'كل، 'كذلك، 'كذا، 'كأى
'لها، 'له، 'لنا، 'لن، 'لما، 'لكيلا', 'لكي، 'لكنما', 'لكن، 'لكم، 'لك، 'لعل، 'لسنا', 'لست، 'لستما', 'لستم
'هذه، 'هذان', 'هذا، 'هاهنا', 'هاك، 'هاتين', 'هاته، 'هاتان', 'هاتا', 'ها، 'نعم، 'نحو، 'نحن، 'مهما', 'مه، 'منها', 'منذ، 'منه، 'من
'وإذ، 'والذين', 'والذي', 'هيهات', 'هيت، 'هيا، 'هي، 'هؤلاء، 'هو، 'هلا، 'هم، 'هما، 'هن، 'هنا، 'هناك، 'هنالك، 'هل، 'هكذا', 'هذين
'آذار، 'كانون، 'أوت، 'جويلية', 'جوان، 'ماي، 'أفريل، 'فيفري', 'جانفي', 'ديسمبر', 'نوفمبر', 'أكتوبر، 'سبتمبر، 'أغسطس، 'يوليو
'سن، 'هللة، 'مليم، 'قرش، 'جنيه، 'ليرة، 'درهم، 'ريال، 'دينار، 'دولار، 'تشرين، 'أيلول، 'آب، 'تموز، 'حزيران، 'أيار، 'نيسان، 'أار
'إحد، 'اثني، 'اثنا، 'أحد، 'عشرة، 'تسعة، 'ثمانية، 'سبعة، 'ستة، 'خمسة، 'أربعة، 'ثلاثة، 'اثنان، 'واحد، 'شيكل، 'يوان، 'بن، 'يورو، 'تيم
'ثان، 'ثان، 'أول، 'جمعة، 'خميس، 'أربعاء، 'ثلاثاء، 'اثنين، 'أحد، 'سبت، 'ثمان، 'عشر، 'تسع، 'سبع، 'ست، 'خمس، 'أربع، 'ثلاث، 'ي
'خ، 'جيم، 'ثاء، 'تاء، 'باء، 'الف، 'ة، 'أ، 'ن، 'ة، 'ث، 'ت، 'ب، 'أ، 'ك، 'ق، 'غ، 'ع، 'ظ، 'طا، 'ص، 'ش، 'ر، 'ز، 'س، 'رابع، 'خامس، 'سادس، 'سابع، 'تاسع، 'ثامن، 'ثالث، 'م
'يا، 'واو، 'هاء، 'نون، 'ميم، 'لام، 'كاف، 'قاف، 'فاء، 'غين، 'عين، 'ظاء، 'طاء، 'ضاد، 'صاد، 'شين، 'سين، 'ه، 'كن، 'ك، 'نا، 'ي، 'همزة، 'ء
'قده، 'قذه، 'قذان، 'قذي، 'إياكما', 'إياكم', 'إياك، 'إياها، 'إياهما', 'إياهم، 'إياه، 'إيانا', 'إياي، 'إياكن، 'أولالك، 'تان، 'تانك
'آه، 'آوه، 'آمين، 'وا، 'فلان، 'بضع، 'كأين، 'كأي، 'ذيت، 'كأن، 'أيان، 'أنى، 'أي، 'أيها، 'الألى، 'الألاء، 'اللائي، 'ته، 'تي، 'تين، 'ثمة، 'ذان، 'ذي، 'ذه، 'قذي
'حم، 'حيّ، 'خذار، 'حاي، 'تلة، 'بطآن، 'بضع، 'بس، 'كأنّ، 'بجـ، 'كخ، 'عدس، 'إلتك، 'إليك، 'إليكن، 'أيّان، 'أيّ، 'أيّ، 'ألى، 'إبوـ، 'أمامك، 'أمامك، 'أوّه، 'صه، 'صمـ، 'طاق، 'طق، 'عدّ، 'ظنّ، 'صبر، 'زعم، 'رأى، 'دري، 'درى، 'خال، 'حسبـ، 'حسب، 'حبيبـ، 'حجا، 'جعل، 'تعلم، 'خ
'تر، 'تخذ، 'اتخذ، 'الفن، 'تفعلين، 'تفعلون، 'يفعلون، 'تفعلان، 'يفعلان، 'وقّ، 'وشّكان، 'وراءك، 'وراءك، 'مكانكن، 'مكا نكما، 'مكانكم، 'مكانك، 'مكانك، 'مكانك، 'رويدك، 'سرعان، 'شتان، 'هيّا، 'قتهات، 'واـ، 'واهاـ
'قلم، 'وهبـ، 'وردـ، 'وجد، 'ذهبـ، 'غادر، 'علمـ، 'خيّرـ، 'خذّ، 'حدّت، 'أينّا، 'أبا، 'أعلمـ، 'أرى، 'أخيرـ، 'كسا، 'سقيـ، 'زودـ، 'رزقـ، 'أعطـ، 'أطعم
'كي، 'كلّا، 'كأنّ، 'فـ، 'علّ، 'علـ، 'رتّ، 'نمّـ، 'جلـل، 'جير، 'بـ، 'إى، 'إلى، 'إنّ، 'أنّ، 'إلاّ، 'إمّا، 'إذاً، 'أجل، 'عزّ، 'لكنّ، 'لات، 'لـ
'سبحان، 'حس، 'جميع، 'تلقاء، 'تجاهـ، 'ثمّ، 'و، 'ه، 'ك، 'لقا، 'ن، 'ه، 'و، 'اـ، 'إلاّ، 'الّ، 'واـ، 'هلّا، 'نّ، 'مـ، 'لكنّ، 'لعلّ، 'لات
'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، 'ثمانين، 'ستين، 'سبعين، 'خمسين، 'أربعين، 'ثلاثين، 'عشرين، 'تسعمائة، 'ثمانمائة، 'سبعمائة، 'ستمائة، 'خمسمائة، 'أربعمائة، 'ثلاثمائة، 'الة
'صراحة، 'صدقا، 'صبرا، 'سمعا، 'سرا، 'دواليك، 'خاصة، 'خلافا، 'حمدا، 'حقا، 'بغتة، 'تعسا، 'قاطبة، 'قطا، 'أجمع، 'جميع، 'عامة، 'عين، 'نفس، 'لا سيم، 'أصلا، 'أهلا، 'أيضا، 'بؤسا، 'بعدا، 'ما، 'نيف، 'بضع، 'تسعين
'أ، 'آناء، 'آنفا، 'أمس، 'الآن، 'أمد، 'إزاء، 'أصلا، 'أبدا، 'معاذ، 'كثيرا، 'البتة، 'قطعا، 'قطبية، 'فصلا، 'فرادى، 'غاليا، 'عيانا، 'عجبا، 'طرا
'ام، 'خلف، 'قبل، 'مرّة، 'لدن، 'كلّما، 'قط، 'غداة، 'غدا، 'عوض، 'ضحوة، 'مساء، 'صباح، 'صباح، 'حقا، 'ثمّ، 'ثمّة، 'نمّ، 'تارة، 'أثان، 'أول، 'أمّ
'ظ، 'صار، 'خار، 'رجع، 'حار، 'تحوّل، 'تبدّل، 'بات، 'انقلب، 'أمسى، 'أضحى، 'أصبح، 'أضحى، 'استحال، 'ارتدّ، 'شمال، 'يمين، 'تحت، 'فوق
'حرى، 'هبّ، 'جعل، 'أوشك، 'أنشأ، 'أقبل، 'انبرى، 'أقبل، 'اخلولق، 'أخذ، 'ابتدأ، 'مافتئ، 'مازال، 'مادام، 'ما برح، 'ما انفك، 'كان، 'غدا، 'عاد
'هبّ، 'كاد، 'كرب، 'قام، 'علق، 'طفق، 'شرع']
------------------------------------------------
```

**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.