

Reading: Stopwords and WordCloud

Estimated time: 15 Min

Objectives:

- Understanding stopwords
- Exploring the significance of stopwords in text analysis
- Familiarizing with the NLTK package
- Exploring Wordcloud Generation
- Analyzing the presence of stopwords in Wordclouds

What are Stopwords?

Stopwords are common words in natural language, such as articles, conjunctions, and prepositions. They are often filtered out during text processing because they lack significant meaning on their own.

Examples of stopwords include:

- **Articles:** "the", "a", "an"
- **Conjunctions:** "and", "but", "or"
- **Prepositions:** "in", "on", "at"
- **Commonly used verbs:** "is", "are", "was", "were"
- **Pronouns:** "I", "you", "he", "she", "it", "we", "they"

Significance of stopwords

The significance of stopwords in text analysis can be elucidated through several key points:

- **Noise reduction:** Stopwords often appear frequently in text but carry little semantic value. By filtering out these words, the overall noise in the text data is reduced. This noise reduction helps to focus the analysis on the more meaningful content of the text.
- **Improved computational efficiency:** Removing stopwords can significantly reduce the computational complexity of text processing tasks. Since stopwords occur frequently and are present in almost every text document, their removal decreases the amount of data that needs to be processed. This leads to faster execution of text analysis algorithms.
- **Enhanced accuracy:** Stopwords can sometimes dominate the frequency distribution of terms within a document but contribute little to the overall meaning. By excluding them from analysis, the accuracy of tasks such as text classification, sentiment analysis, and topic modeling is improved. This is because the focus shifts to more meaningful terms that better represent the content of the text.
- **Facilitation of meaningful insights:** The removal of stopwords enables analysts to extract more meaningful insights from the text data. By concentrating on content-bearing words, analysts can uncover trends, patterns, and relationships that are more relevant to the subject matter being analyzed.
- **Improved natural language processing (NLP):** Stopword removal is a common preprocessing step in NLP tasks. By eliminating these irrelevant words, NLP algorithms can better understand and interpret the semantics of the text. This leads to more accurate results in tasks such as named entity recognition, part-of-speech tagging, and semantic analysis.
- **Language independence:** Stopwords exist across different languages, although the specific words may vary. Therefore, the concept of stopwords and their significance extends to multilingual text analysis. Regardless of the language being analyzed, removing stopwords remains a crucial step in text preprocessing for various NLP tasks.
- **Customization and flexibility:** While there are standard lists of stopwords available for different languages, analysts can also customize the list of stopwords based on the specific context or domain of the text data. This flexibility allows for tailored preprocessing approaches that better suit the needs of the analysis.

What is NLTK?

Natural Language Toolkit (NLTK) is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to over 50 corpora and lexical resources such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning.

Key points about NLTK:

1. **Comprehensive toolkit:** NLTK offers a wide range of functionalities for natural language processing tasks, making it a comprehensive toolkit for researchers, developers, and educators.
2. **Open-source:** NLTK is an open-source library, which means it is freely available for anyone to download, use, and modify according to their requirements.
3. **Community support:** Being one of the most widely used NLP libraries, NLTK has a large community of users and contributors. This ensures continuous development, improvement, and support for the library.
4. **Ease of use:** NLTK provides simple and intuitive interfaces for various NLP tasks, making it accessible even to users with limited programming experience.
5. **Rich resource collection:** NLTK comes with a vast collection of corpora, lexical resources, and pre-trained models, facilitating research and experimentation in natural language processing.

Exploring Wordcloud generation

Word cloud generation is a visualization technique used to represent textual data in a visually appealing and intuitive manner. In a word cloud, words from a given text are displayed in different sizes and colors, with more frequent words appearing larger and more prominent while less frequent words are smaller and less prominent. Word

clouds are commonly used to identify patterns, trends, and key themes within a corpus of text data.

Choosing and applying frequency and weights to create word clouds

Creating a word cloud involves determining the frequency of each word in the given text and assigning weights to these frequencies to visualize them effectively. Here's a breakdown of how frequency and weights are chosen and applied:

- **Tokenization:** The text is first broken down into individual words or tokens. Punctuation marks, numbers, and other non-alphabetic characters are typically removed at this stage.
- **Frequency calculation:** The frequency of each word is counted. This involves counting how many times each unique word appears in the text. For example, in the provided text, "acting" appears once, "movie" appears once, "phenomenal" appears once, and so on.
- **Weight assignment:** The frequency of each word is then assigned a weight. This weight determines the importance of the word in the word cloud. There are various methods to assign weights, but commonly used ones include:
 - **Linear scaling:** In this method, the frequency of each word is directly mapped to its size in the word cloud. Words with higher frequencies are assigned larger sizes, making them more prominent in the visualization.
 - **Logarithmic scaling:** This method involves taking the logarithm of the frequency of each word before mapping it to its size in the word cloud. Logarithmic scaling helps to reduce the influence of extremely high-frequency words, allowing less frequent but still important words to be more visible.
 - **TF-IDF (Term Frequency-Inverse Document Frequency):** TF-IDF considers both the frequency of a word in the current text (term frequency) and its rarity across all texts (inverse document frequency). Words with high TF-IDF scores are considered important in the current text but rare in the overall corpus, making them stand out in the word cloud.
- **Visualization:** Once weights are assigned, the words are arranged in the word cloud visualization. Words with higher weights (determined by frequency and weight assignment method) are typically displayed larger and more prominently, while words with lower weights are smaller and less prominent. The arrangement of words is often randomized or organized in a visually pleasing manner to create an aesthetic representation of the text.

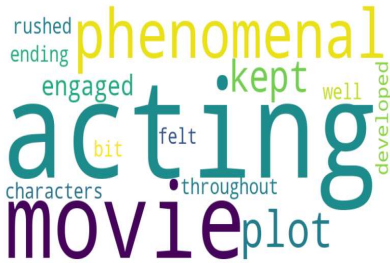
Types of word clouds:

1. **Frequency-based word clouds:** These word clouds display words based on their frequency of occurrence in the text, with more frequent words appearing larger and more prominent.

Let's say we have a text document containing the following sentences:

Text = "The acting in this movie was phenomenal! The plot kept me engaged throughout, and the characters were well-developed. However, the ending felt a bit rushed."

The frequency word cloud might look something like this:



The size of each word is determined solely by its frequency of occurrence in the text. The more frequently a word appears, the larger it appears in the word cloud.

2. **Weighted word clouds:** In weighted word clouds, users can assign custom weights or importance scores to individual words, influencing their size and prominence in the cloud.



The size of each word is determined not only by its frequency but also by some additional weighting factor. This weighting factor could represent significance, importance, or any other metric relevant to the context.

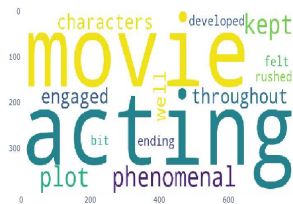
3. **Masked word clouds:** Masked word clouds use a specific image, known as a mask, to define the shape of the word cloud. Words are arranged within the boundaries of the mask image.



The size of each word is determined based on its frequency in the text. Words with higher frequencies are typically assigned larger sizes, making them more prominent in the visualization.

4. **Interactive word clouds:** Interactive word clouds allow users to interact with the visualization by hovering over words to view additional information or clicking on words to navigate to related content.

Interactive Word Cloud



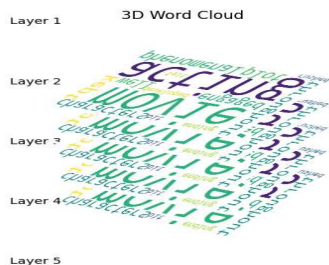
Word selection and representation in interactive word clouds can be influenced by various factors beyond just frequency. Users may have the ability to customize the visualization by specifying criteria such as relevance, sentiment, or metadata associated with the words. Additionally, words may be dynamically updated based on user interactions or changes in the underlying data.

5. **Dynamic word clouds:** Dynamic word clouds update in real time or dynamically change based on user input or changes in the underlying data. They can be useful for monitoring trends or displaying live data.



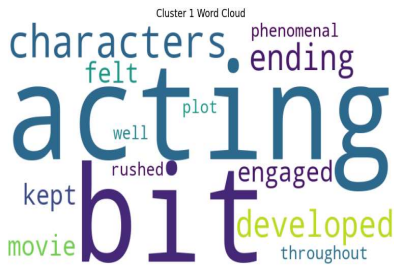
Words may be selected, resized, or repositioned based on various factors beyond just frequency. This could include relevance, sentiment analysis, or other contextual information. It's often used when there is a need to visualize changes or trends in textual data over time, such as social media trends, live feedback analysis, or streaming data analysis.

6. **3D word clouds:** In 3D word clouds, words are displayed in a three-dimensional space, adding depth and visual interest to the visualization.



3D word clouds represent words in three-dimensional space, often with varying sizes and positions to create depth and perspective. Words in a 3D word cloud are typically arranged in a three-dimensional grid or layout, with larger words appearing closer to the viewer and smaller words further away.

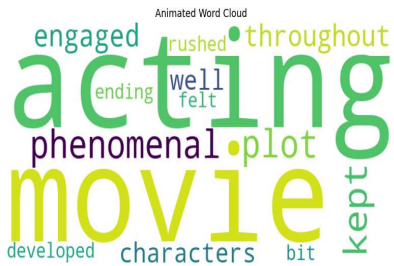
7. **Clustered word clouds:** Clustered word clouds group related words or terms together based on their semantic similarity or co-occurrence in the text, allowing for more structured and organized visualizations.



Clustered word clouds organize words into distinct groups or clusters based on their semantic similarity or other criteria. Each cluster may represent a different topic or theme.

Within each cluster, words are arranged in a traditional two-dimensional layout, with larger sizes indicating higher frequency or importance.

8. **Animated word clouds:** Animated word clouds feature dynamic movement or transitions, making them visually engaging and attention-grabbing.



Animated word clouds incorporate dynamic elements such as motion, transitions, or temporal changes to represent evolving data or trends.

The organization of words in an animated word cloud may vary depending on the specific implementation, but they often follow similar principles as static word clouds, with words arranged based on frequency or relevance.

Stopwords in Wordclouds

NLTK is a powerful library for natural language processing tasks in Python. It provides a simple way to access a list of stopwords for various languages.

Step 1: Install NLTK and WordCloud

If you haven't already installed NLTK and WordCloud, you can do so using pip:

```
pip install nltk
```

```
pip install wordcloud
```

Step 2: Import the required libraries

```
import nltk
from nltk.corpus import stopwords
from wordcloud import WordCloud
import matplotlib.pyplot as plt
```

Step 3: Download NLTK stopwords

```
nltk.download('stopwords')
```

Step 4: Define the text

```
text = "The acting in this movie was phenomenal! The plot kept me engaged throughout, and the characters were well-developed. However, the ending fel
```

Note: The text variable can be replaced with any text content, such as a product review, feedback, or any other textual data.

Step 5: Split the text into words

```
words = text.split()
```

The `split()` method is applied to the string variable `text`. This method is a built-in function in Python that is used to split a string into a list of substrings, based on a specified separator.

Step 6: Remove stopwords

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

Here's what's happening:

- **stopwords.words('english')**: This part of the code is a method call to `stopwords.words()` from the NLTK library. `stopwords.words()` is a method that returns a list of stopwords for a specified language. In this case, 'english' is passed as an argument to specify that we want stopwords for the English language.
- **set()**: This function is then used to convert the list of stopwords returned by `stopwords.words('english')` into a set. A set is an unordered collection of unique elements in Python. Using a set ensures that each stopword appears only once and eliminates any duplicates.
- **stop_words**: Finally, the resulting set of English stopwords is assigned to the variable `stop_words`.

Step 7: Add custom stopwords

```
custom_stopwords = {'movie', 'plot', 'characters', 'ending'} # Example custom stopwords  
stop_words.update(custom_stopwords)
```

Here's what's happening:

- **custom_stopwords = {'movie', 'plot', 'characters', 'ending'}:** This line creates a set called custom_stopwords containing four example stopwords specific to the context of the text analysis. Each stopword, such as 'movie', 'plot', 'characters', and 'ending', is enclosed within curly braces {} and separated by commas, representing individual elements of the set.
- **stop_words.update(custom_stopwords):** The update() method is called on the stop_words set to incorporate the custom stopwords contained in the custom_stopwords set. This method adds elements from custom_stopwords to the stop_words set, ensuring that the stop_words set now includes both the default NLTK stopwords for English as well as the custom stopwords defined in custom_stopwords.

Step 8: Filter out stopwords

```
filtered_words = [word for word in words if word.lower() not in stop_words]
```

Note: You can choose stopwords from any language by replacing 'english' with the desired language code.

Here's a breakdown of what each part does:

- **filtered_words:** This is the variable that will store the filtered words after the list comprehension is applied.
- **[word for word in words if word.lower() not in stop_words]:** This is the list comprehension itself, which consists of several parts.
- **word:** This represents each individual word in the words list.
- **for word in words:** This is the iteration over each word in the words list.
- **if word.lower() not in stop_words:** This is a conditional statement that checks if the lowercase version of the current word (word.lower()) is not present in the stop_words set. Here's what happens within this conditional statement:
- **word.lower():** This converts the current word to lowercase. This ensures that case differences (for example, "The" vs. "the") don't affect the comparison.
- **not in stop_words:** This checks if the lowercase version of the current word is not present in the stop_words set.

Step 9: Join the filtered words back into a single string

```
filtered_text = ' '.join(filtered_words)
```

Here's what each part does:

- **filtered_text:** This is the variable that will store the resulting string after joining the words.
- **' '.join(filtered_words):** This is the join() method applied to the string ' '.
- **' ':** This is the separator that will be used to join the elements of the list. In this case, it's a single space, indicating that each word in the list will be separated by a space.
- **.join(filtered_words):** This is the join() method applied to the separator ' '. It joins the elements of the filtered_words list into a single string, with each element separated by the specified separator. Here, filtered_words is the list containing the words that passed the filtering process (i.e., the stopwords were removed). Joining these words with spaces between them reconstructs the original text but without the stopwords.

Step 10: Generate WordCloud

```
wordcloud = WordCloud(width=800, height=800,
                      background_color='white',
                      stopwords=stop_words,
                      min_font_size=10).generate(filtered_text)
```

Here, we are creating a WordCloud object named wordcloud with the following specifications:

- **Width** and **height** are set to 800 pixels each.
- The **background color** of the word cloud is set to white.
- The **stopwords** parameter is set to the `stop_words` set, indicating that stopwords should be excluded from the word cloud.
- The **min_font_size** parameter is set to 10, specifying the minimum font size for words in the word cloud.
- The **generate()** method is called on the `WordCloud` object, with **filtered_text** passed as an argument. This generates the word cloud based on the provided text after stopwords have been removed.

Step 11: Plot the WordCloud

```
plt.figure(figsize=(8, 8), facecolor=None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad=0)
plt.show()
```

Here we are utilizing the matplotlib to display the generated word cloud:

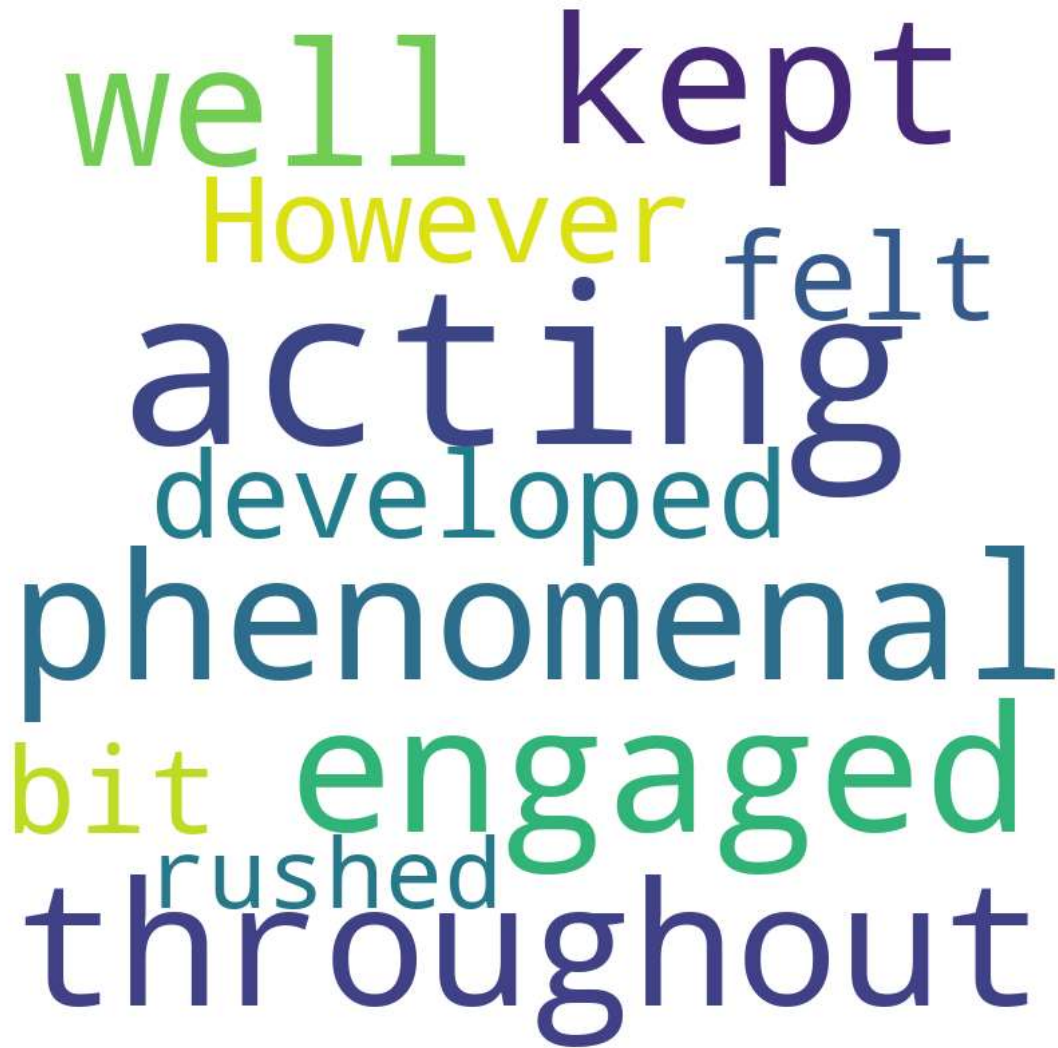
- **plt.figure(figsize=(8, 8), facecolor=None):** This line creates a new figure with a size of 8x8 inches. The `facecolor=None` argument sets the background color of the figure to be transparent.
- **plt.imshow(wordcloud):** This line displays the word cloud (`wordcloud`) on the previously created figure.
- **plt.axis("off"):** This line turns off the axis lines and labels in the plot, ensuring that only the word cloud is displayed without any extra elements.
- **plt.tight_layout(pad=0):** This line adjusts the padding of the layout to 0, ensuring that the word cloud fills the entire figure without any extra whitespace.
- **plt.show():** This line displays the plot, showing the word cloud visualization generated by Matplotlib.

Complete Code

Here is the complete code:

```
import nltk
from nltk.corpus import stopwords
from wordcloud import WordCloud
import matplotlib.pyplot as plt
# Download NLTK stopwords
nltk.download('stopwords')
# Define the text
text = "The acting in this movie was phenomenal! The plot kept me engaged throughout, and the characters were well-developed. However, the ending fel
# Split the text into words
words = text.split()
# Remove NLTK stopwords
stop_words = set(stopwords.words('english'))
# Add custom stopwords
custom_stopwords = {'movie', 'plot', 'characters', 'ending'} # Example custom stopwords
stop_words.update(custom_stopwords)
# Filter out stopwords
filtered_words = [word for word in words if word.lower() not in stop_words]
# Join the filtered words back into a single string
filtered_text = ' '.join(filtered_words)
# Generate WordCloud
wordcloud = WordCloud(width=800, height=800,
                      background_color='white',
                      stopwords=stop_words,
                      min_font_size=10).generate(filtered_text)
# Plot the WordCloud
plt.figure(figsize=(8, 8), facecolor=None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad=0)
plt.show()
```

Desired output



Analysis :

- **Positive Sentiment:** The most prominent words are all positive, including "phenomenal", "engaged", "well-developed", and "kept." This suggests that the overall sentiment of the text is positive.
- **Focus on Acting and Plot:** The largest words are "acting" and "plot", which suggests that these are the two aspects of the movie that the reviewer found most impressive.

Conclusion

Customizing stopwords based on the context of your text data empowers you to perform more accurate and relevant analyses. By integrating NLTK's functionality, you can efficiently preprocess your text data and extract meaningful insights.

Author(s)

[Akansha Yadav](#)



Skills Network