



[Zolzaya Luvsandorj](#)

Jun 20, 2020

<https://towardsdatascience.com/simple-wordcloud-in-python-2ae54a9f58e5>

Simple word cloud in Python

💡 **Word cloud** is a technique for visualising frequent words in a text where the size of the words represents their frequency.

One easy way to make a word cloud is to search 'word cloud' on Google to find one of those free websites that generate a word cloud. You can possibly customise how it looks like. Quick and easy!

What's more exciting is that you can build one yourself in Python 🌀. This post will show how to create a word cloud like the example below.



word cloud from section 3. Fancier word cloud 🌀

0. Python setup 🛠️

I assume the reader (👁️ yes, you!) has access to and is familiar with Python including installing packages, defining functions and other basic tasks. If you are new to Python, [this](#) is a good place to get started.

I have used and tested the scripts in Python 3.7.1 in Jupyter Notebook.

Let's make sure you have the following libraries installed before we get started:

- To create a word cloud:** *wordcloud*
- To import an image:** *pillow* (will later import is as *PIL*)
- To scrape text from Wikipedia:** *wikipedia*

Last package is optional, you can instead load up or create your own text data without having to pull text via web scraping.

1. Data 📄

As our sample text, we will use scraped text from a Wikipedia page on 'Web scraping'.

```
# Import packages
import wikipedia
import re# Specify the title of the Wikipedia page
wiki = wikipedia.page('Web scraping')# Extract the plain text content of the
page
text = wiki.content# Clean text
text = re.sub(r'==.*?==+', '', text)
text = text.replace('\n', '')
```

⇒ I have explained what this script does in a [separate post](#) on scraping. In short, this script will pull out the plain text content in the paragraphs and assign it to *text* string.

2. Word cloud ☁️

Firstly, let's prepare a function that plots our word cloud:

```
# Import package
import matplotlib.pyplot as plt# Define a function to plot word cloud
def plot_cloud(wordcloud):
    # Set figure size
    plt.figure(figsize=(40, 30))
    # Display image
    plt.imshow(wordcloud)
    # No axis details
    plt.axis("off");
```

Secondly, let's create our first word cloud and plot it:

```
# Import package
from wordcloud import WordCloud, STOPWORDS# Generate word cloud
wordcloud = WordCloud(width= 3000, height = 2000, random_state=1,
background_color='salmon', colormap='Pastell', collocations=False, stopwords =
STOPWORDS).generate(text)# Plot
plot_cloud(wordcloud)
```



Ta-da! We just built a word cloud! Here are some notes regarding the arguments for *WordCloud* function:

- **width/height:** You can change the word cloud dimension to your preferred width and height with these.
- **random_state:** If you don't this set this to a number of your choice, you are likely to get a slightly different word cloud every time you run the same script on the same input data. By setting this parameter, you ensure reproducibility of the exact same word cloud. You could play around with random numbers until you find the one that results in the word cloud you like.
- **background_colour:** 'white' and 'black' are common background colours. If you would like to explore more colours, [this](#) may come in handy. Please note that some colours may not work. Hope you will find something you fancy.
- **colormap:** With this argument, you can set up the colour theme that the words are displayed in. There are many beautiful [Matplotlib colormaps](#) to choose from. Some of my favourites are 'rainbow', 'seismic', 'Pastell' and 'Pastel2'.

■ **collocations:** Set this to *False* to ensure that the word cloud doesn't appear as if it contains any duplicate words. Otherwise, you may see 'web', 'scraping' and 'web scraping' as a collocation in the word cloud, giving an impression that words have been duplicated.

■ **stopwords:** Stopwords are common words which provide little to no value to the meaning of the text. 'We', 'are' and 'the' are examples of stopwords. I have explained stopwords in more detail [here](#) (scroll to 'STEP3. REMOVE STOPWORDS' section).

Here, we used *STOPWORDS* from the *wordcloud* package. To see the set of stopwords, use `print(STOPWORDS)` and to add custom stopwords to this set, use this template `STOPWORDS.update(['word1', 'word2'])`, replacing *word1* and *word2* with your custom stopwords before generating a word cloud.

There are other arguments that you can also customise. Check out the [documentation](#) for more information.

Let's generate another word cloud with a different *background_colour* and *colormap* 🌀. You could play with different combinations until you find the one that you like. I find the following combination quite nice:

```
# Generate wordcloud
wordcloud = WordCloud(width = 3000, height = 2000, random_state=1,
background_color='black', colormap='Set2', collocations=False, stopwords =
STOPWORDS).generate(text)# Plot
plot_cloud(wordcloud)
```


■ [Two simple ways to scrape text from Wikipedia in Python](#)

(Below lists a series of posts on Introduction to NLP)

■ [Part 1: Preprocessing text in Python](#)

■ [Part 2: Difference between lemmatisation and stemming](#)

■ [Part 3: TF-IDF explained](#)

■ [Part 4: Supervised text classification model in Python](#)

■ [Part 5A: Unsupervised topic model in Python \(sklearn\)](#)

■ [Part 5B: Unsupervised topic model in Python \(gensim\)](#)