

# C1M6L2\_Final\_Project\_V3

July 28, 2020

## 1 Final Project - Word Cloud

For this project, you'll create a "word cloud" from a text by writing a script. This script needs to process the text, remove punctuation, ignore case and words that do not contain all alphabets, count the frequencies, and ignore uninteresting or irrelevant words. A dictionary is the output of the `calculate_frequencies` function. The `wordcloud` module will then generate the image from your dictionary.

For the input text of your script, you will need to provide a file that contains text only. For the text itself, you can copy and paste the contents of a website you like. Or you can use a site like [Project Gutenberg](#) to find books that are available online. You could see what word clouds you can get from famous books, like a Shakespeare play or a novel by Jane Austen. Save this as a .txt file somewhere on your computer. Now you will need to upload your input file here so that your script will be able to process it. To do the upload, you will need an uploader widget. Run the following cell to perform all the installs and imports for your word cloud script and uploader widget. It may take a minute for all of this to run and there will be a lot of output messages. But, be patient. Once you get the following final line of output, the code is done executing. Then you can continue on with the rest of the instructions for this notebook. **Enabling notebook extension fileupload/extension... - Validating: OK**

```
In [1]: # Here are all the installs and imports you will need for
        # your word cloud script and uploader widget

        # !pip install wordcloud
        # !pip install fileupload
        # !pip install ipywidgets
        # !jupyter nbextension install --py --user fileupload
        # !jupyter nbextension enable --py fileupload

import wordcloud
import numpy as np
from matplotlib import pyplot as plt
from IPython.display import display
import fileupload
import io
import sys
```

Whew! That was a lot. All of the installs and imports for your word cloud script and uploader widget have been completed. **IMPORTANT!** If this was your first time running the above cell

containing the installs and imports, you will need save this notebook now. Then under the File menu above, select Close and Halt. When the notebook has completely shut down, reopen it. This is the only way the necessary changes will take affect. To upload your text file, run the following cell that contains all the code for a custom uploader widget. Once you run this cell, a “Browse” button should appear below it. Click this button and navigate the window to locate your saved text file.

In [2]: *# This is the uploader widget*

```
def _upload():

    _upload_widget = fileupload.FileUploadWidget()

    def _cb(change):
        global file_contents
        decoded = io.StringIO(change['owner'].data.decode('utf-8'))
        filename = change['owner'].filename
        print('Uploaded `{}` ( {:.2f} kB)'.format(
            filename, len(decoded.read()) / 2 **10))
        file_contents = decoded.getvalue()

    _upload_widget.observe(_cb, names='data')
    display(_upload_widget)

_upload_widget = FileUploadWidget(label='Browse', _dom_classes=('widget-item', 'btn-group'))
```

Uploaded `The Yellow Wallpaper.txt` (30.96 kB)

The uploader widget saved the contents of your uploaded file into a string object named *file\_contents* that your word cloud script can process. This was a lot of preliminary work, but you are now ready to begin your script.

Write a function in the cell below that iterates through the words in *file\_contents*, removes punctuation, and counts the frequency of each word. Oh, and be sure to make it ignore word case, words that do not contain all alphabets and boring words like “and” or “the”. Then use it in the *generate\_from\_frequencies* function to generate your very own word cloud! **Hint:** Try storing the results of your iteration in a dictionary before passing them into wordcloud via the *generate\_from\_frequencies* function.

```
In [3]: def calculate_frequencies(file_contents):
        # A list of punctuations and uninteresting words you can use to process text
        punctuations = '''!()-[]{};:'"\<>./?@#%&*_~'''
        uninteresting_words = ["the", "a", "to", "if", "is", "it", "of", "and", \
            "or", "an", "as", "i", "me", "my", "we", "our", "ours", "you", "your", \
            "yours", "he", "she", "him", "his", "her", "hers", "its", "they", "them", \
            "their", "what", "which", "who", "whom", "this", "that", "am", "are", \
```

```

"was", "were", "be", "been", "being", "have", "has", "had", "do", "does", \
"did", "but", "at", "by", "with", "from", "here", "when", "where", "how", \
"all", "any", "both", "each", "few", "more", "some", "such", "no", "nor", \
"too", "very", "can", "will", "just", "in", "on", "for", "not", "so"]

# LEARNER CODE START HERE
for x in punctuations:
    file_contents = file_contents.replace(x, " ")
wordcount = {}
words = file_contents.split()
for word in words:
    # Do nothing if the word is uninteresting
    if word.lower() in uninteresting_words:
        continue
    # Add the word to dictionary if it's a new word
    elif word.upper() not in wordcount:
        wordcount[word.upper()] = 0
    # Increase the count of word by 1
    wordcount[word.upper()] += 1

#wordcloud
cloud = wordcloud.WordCloud(width=3840, height=2160, background_color='white',
                             mode="RGBA", colormap='tab20')
cloud.generate_from_frequencies(wordcount)
return cloud.to_array()

```

If you have done everything correctly, your word cloud image should appear after running the cell below. Fingers crossed!

In [4]: *# Display your wordcloud image*

```

myimage = calculate_frequencies(file_contents)
plt.figure(figsize=(16, 9), dpi=300)
plt.imshow(myimage, interpolation = 'nearest')
plt.axis('off')
plt.show()

```

