

# SFWRENG 4NL3 Assignment 1

Sumanya Gulati

20 January 2025

# Contents

<b>1</b>	<b>Data</b>	<b>2</b>
<b>2</b>	<b>Methodology</b>	<b>2</b>
2.1	Approach . . . . .	2
2.2	Available Options . . . . .	2
<b>3</b>	<b>Sample Output</b>	<b>3</b>
<b>4</b>	<b>Discussion</b>	<b>7</b>
4.1	Findings . . . . .	7
4.2	Reflection . . . . .	7

# List of Tables

# List of Figures

1	Top 25 words with their count. . . . .	4
2	Last 25 words with their count. . . . .	5
3	Bar plot showing the frequencies of top 45 words. . . . .	6
4	Last 25 words with their count. . . . .	7

# 1 Data

`myfile.txt` is a text file that contains inaugural speeches by United States Presidents' from George Washington to Barack Obama. The dataset was acquired from the Project Gutenberg collection and was downloaded as a Plain Text UTF-8 encoded file.

Prior to tokenization, the dataset was manually pre-processed which included removing the title page, table of contents, disclaimers and such information from the beginning and the ending of the text file to ensure consistency and remove any characters unrelated to the actual content.

I chose this dataset in light of today's inaugural speech of the 47th U.S. President. I also wanted to be able to use the script from this assignment to gather insights about how inaugural addresses have changed over centuries and analyze whether the increased use of words like 'immigration' and 'economy' is a recent trend over the past few decades.

I also want to analyze if the use of words propogating sentiments of nationalism and unity increases in years following major national of global events like the Vietnam war, 9/11, the 2008 recession, COVID-19 pandemic and more. If you would be interested in viewing the results of these analyses, feel free to shoot me an email and I would be glad to share them (assuming, I am actually able to get it done soon).

## 2 Methodology

This section outlines the approach employed in writing this script and explains available normalization options while describing their use and a rationale for using them.

### 2.1 Approach

The approach I used to write this script involved breaking it down into steps and implementing one required step at a time, tetsing it and then moving onto the next step. I started by writing a skeletal script that reads a file and implements command-line argument parsing. After that, I wrote the `tokenize_text()` function to pattern-match the input text file and break it into tokens, storing the result as a list.

The next step involved implementing the *lowercase* and *remove\_stopwords* options, followed by implementing the *lemmatization* and *stemming* options, the code snippets for which were copied from the lecture slides. The `visualize()` function was then implemented to produce first, a DataFrame (basically, a nicer table to view the output) and then a bar plot. A trial and error approach was used to figure out if either or both axes should be set as a log scale and a log-log line plot was generated to compare my result with that of datasets obeying Zipf's law as shown in this Wikipedia article.

Generative AI was used to resolve a logical bug and write line 12 of the script. A search query weas employed using ChatGPT with the following prompt: *"I am trying to tokenize an input text file using regex pattern matching but some of my tokens in the list are preceeded by the letters, ufeff. Could this be random or is there a logical error in my code?"*. As outlined in the course outline, the carbon footprint of this query is 4.32g of CO<sub>2</sub>.

### 2.2 Available Options

The following options, in any combination, can be used by the user during the text normalization steps:

1. `--lowercase`: Change letter case
2. `--stem`: Apply stemming
3. `--lemmatize`: Apply lemmatization

4. `--remove-stopwords`: Remove commonly used stopwords in the English language
5. `--remove-numbers`: Remove numbers

Beyond the defined requirements, I added the option to remove numbers. Since the text file contains speeches spanning centuries, it seemed like a fair assumption (backed by a brief look at the raw data) that the text file contains the use of digits to write out for example, the date or the year of the speech. I believe this is a useful option as it removes the numbers added as a part of the description, producing a more consistent dataset that focuses purely on the content of the speeches by excluding the remaining descriptive text.

I recognize that this option may cause unwanted issues by removing numbers that are a part of the actual speech or not remove numbers that are spelled out as opposed to being written as digits. Configuring this option to not cause these issues seems like a more complex task that would require additional time and effort.

### 3 Sample Output

Using the `pandas` library and the `DataFrame` structure, a table has been created to sort the list of unique tokens in decreasing order along with their frequency/count.

Figure 1 displays a table showing the top 25 words and their frequency based on the program's output.

	<b>Tokens</b>	<b>Count</b>	<b>Rank</b>
<b>0</b>	the	9396	1
<b>1</b>	of	7047	2
<b>2</b>	and	5091	3
<b>3</b>	to	4384	4
<b>4</b>	a	3078	5
<b>5</b>	in	2553	6
<b>6</b>	our	1977	7
<b>7</b>	that	1743	8
<b>8</b>	it	1580	9
<b>9</b>	be	1473	10
<b>10</b>	is	1434	11
<b>11</b>	we	1191	12
<b>12</b>	for	1101	13
<b>13</b>	by	1051	14
<b>14</b>	have	1004	15
<b>15</b>	which	1003	16
<b>16</b>	not	940	17
<b>17</b>	with	898	18
<b>18</b>	will	869	19
<b>19</b>	I	834	20
<b>20</b>	are	795	21
<b>21</b>	all	773	22
<b>22</b>	their	726	23
<b>23</b>	this	713	24
<b>24</b>	The	623	25

Figure 1: Top 25 words with their count.

Similarly, 2 displays a table showing the last 25 words and their frequency based on the program's output.

	<b>Tokens</b>	<b>Count</b>	<b>Rank</b>
<b>8821</b>	Falls	1	8822
<b>8822</b>	Selma	1	8823
<b>8823</b>	Stonewall	1	8824
<b>8824</b>	sung	1	8825
<b>8825</b>	unsung	1	8826
<b>8826</b>	footprint	1	8827
<b>8827</b>	preacher	1	8828
<b>8828</b>	King	1	8829
<b>8829</b>	inextricably	1	8830
<b>8830</b>	daughter	1	8831
<b>8831</b>	gay	1	8832
<b>8832</b>	student	1	8833
<b>8833</b>	workforce	1	8834
<b>8834</b>	expelled	1	8835
<b>8835</b>	Detroit	1	8836
<b>8836</b>	Appalachia	1	8837
<b>8837</b>	lane	1	8838
<b>8838</b>	Newtown	1	8839
<b>8839</b>	contour	1	8840
<b>8840</b>	absolutism	1	8841
<b>8841</b>	reasoned	1	8842
<b>8842</b>	Philadelphia	1	8843
<b>8843</b>	realizes	1	8844
<b>8844</b>	END	1	8845
<b>8845</b>	P	1	8846

Figure 2: Last 25 words with their count.

Using the `matplotlib.pyplot` library, a bar graph depicted in figure 3 shows the top 45 words. It must be noted that the y-axis has been set to a log scale to better understand the generated output. Although a normal bar plot without a log scale did produce a readable graph, the difference between frequencies was not as apparent.

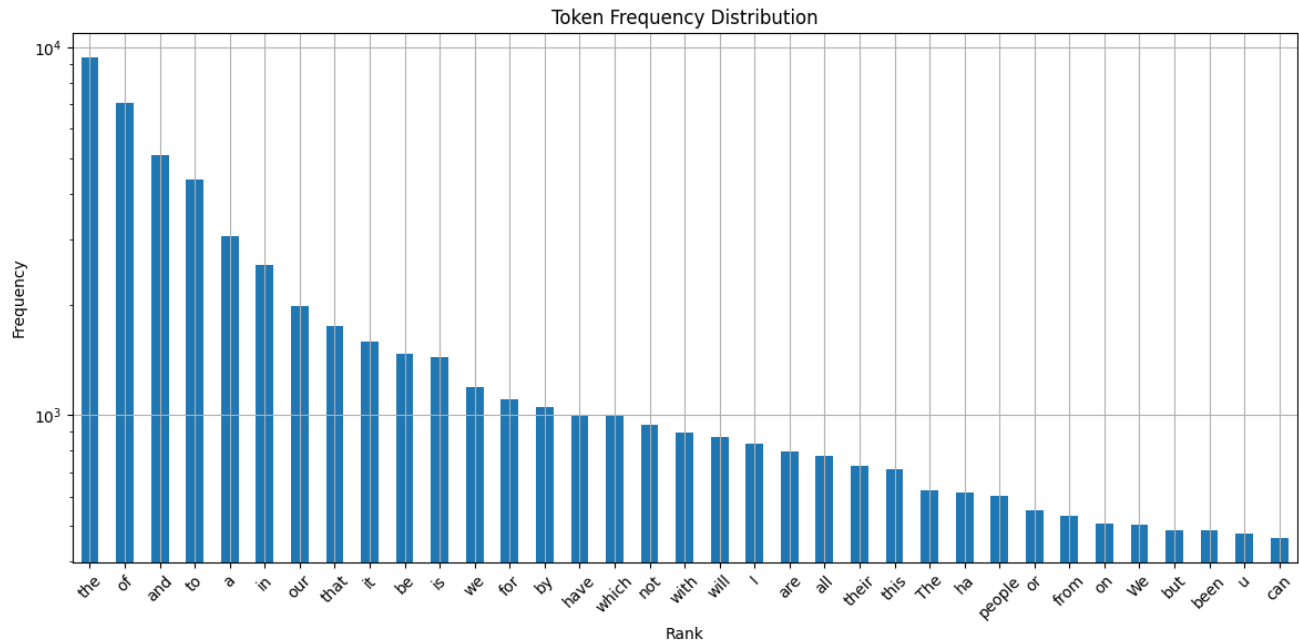


Figure 3: Bar plot showing the frequencies of top 45 words.

To better compare the output with Zipf's Law, a line graph as shown in figure 4 has been generated. It must be noted that this is a log-log plot.

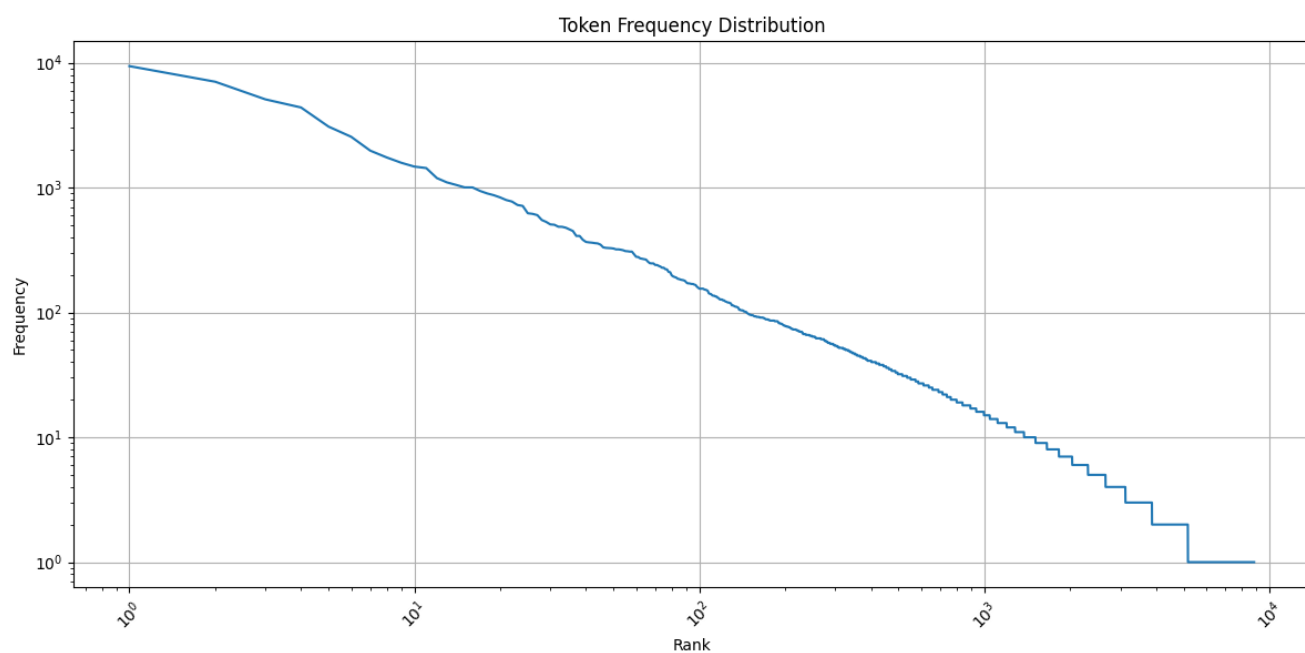


Figure 4: Last 25 words with their count.

## 4 Discussion

### 4.1 Findings

### 4.2 Reflection