

NLP PROJECT REPORT

Submitted by

Team 'SAAR'

Rupesh Maheshwari - 20UCS165

Aayush Ashokbhai Sheth - 20UCS004

Ajinkya Eknath Kadam - 20DCS001

Sankalp Jain - 20UCS172

in partial fulfillment for the award of the degree of

B.Tech. in CSE



Github Repository Link -: <https://github.com/jinks2882/NLP-project>

ACKNOWLEDGEMENT

We would like to express our special thanks to our project guide Dr. Sakthi Balan who gave us the golden opportunity to do this wonderful project on the topic **Text Analysis [NLP Project Round -1]** which also helped us in doing a lot of research and it was a great learning experience.

DATE : 28 OCT 2022

Rupesh Maheshwari	(20UCS165)
Aayush Ashokbhai Sheth	(20UCS004)
Ajinkya Eknath Kadam	(20DCS001)
Sankalp Jain	(20UCS172)

TABLE OF CONTENTS

TITLE

	ABSTRACT	1
1.0	LITERATURE REVIEW	2
2.0	INTRODUCTION	3
3.0	PYTHON LIBRARIES	4
4.0	METHODOLOGY	5
5.0	CONCLUSION	15
6.0	REFERENCES	16

ABSTRACT

Text Analytics is a very important aspect in the field of natural language processing and in this project, we worked on text preprocessing, PoS tagging, and various other operations the book ***Understanding Cryptography by Christopher Paar and Jan Pelzl***. Working on this project on this book was particularly interesting because the book is written in a very simple manner which give easy understanding of the book. For all the operations performed, we have used NLTK (natural language Tool Kit) to perform all the preprocessing, tokenization, and tagging. The project also described the frequency distribution and Word Cloud of the book and helped us understand some fields of text mining. The graphs that are plotted in the report also say a lot about the input text which is derived from the book and writing style of the author, the words that he used frequently, main terms, definitive words etc. This project can also be used in understanding vocabulary in a certain text file, frequency of words, difficulty in the text file.

1. LITERATURE REVIEW

Many researchers worked on NLP, building tools and systems which make NLP what it is today. Tools like Sentiment Analyzer, Parts of Speech (POS) Taggers, Chunking, Named Entity Recognition (NER), Emotion detection, Semantic Role Labelling made NLP a good topic for research.

Related Work :

- Sentiment analyzer (Jeonghee et al.,2003) [26] works by extracting sentiments about a given topic.
- Parts of speech taggers for the languages like European languages, research is being done on making parts of speech taggers for other languages like Arabic, Sanskrit (Namrata Tapswi, Suresh Jain , 2012) [27], Hindi (Pradipta Ranjan Ray et al., 2003) [28], etc. It can efficiently tag and classify words as nouns, adjectives, verbs, etc.
- The Sanskrit part of speech tagger specifically uses the treebank technique.
- Arabic uses the Support Vector Machine (SVM) (Mona Diab et al.,2004) [29] approach to automatically tokenize, tag parts of speech, and annotate base phrases in Arabic text.

2. INTRODUCTION

In this Project we imported a book in text format in order to perform text analysis using NLP techniques on it. We tokenized and lemmatized the imported text file, analyzed the frequency distribution and performed PoS tagging on the text file, further we are going to visualize the data which is going to be a good learning experience in the field of NLP.

We have chosen book related to our course Computer Security. This book is among the top downloaded books in field of cryptography.

- ***Understanding Cryptography written by Christopher Paar and Jen Pelzl***

For the tasks given in this project we have used NLTK which is a suite of libraries and programs for symbolic and statistical natural language processing for English written in the Python programming language. And we have imported some modules and functions in order to perform different activities such as preprocessing, tokenizing, removing stop words etc. And after performing such operations on the text file we have plotted frequency distribution for text file and created word clouds

3. Python Libraries and Modules Used

Nltk.tokenizer package : Tokenizer divides strings into a list of substrings.

Nltk.stem package : Interfaces used to remove morphological affixes from words, leaving only the word stem.

Nltk.probability : A probability distribution specifies how likely it is that an experiment will have any given output.

Nltk.corpus : The modules in this package provide functions that can be used to read corpus files in a variety of formats.

Wordcloud package : Provides modules to create wordcloud in python.

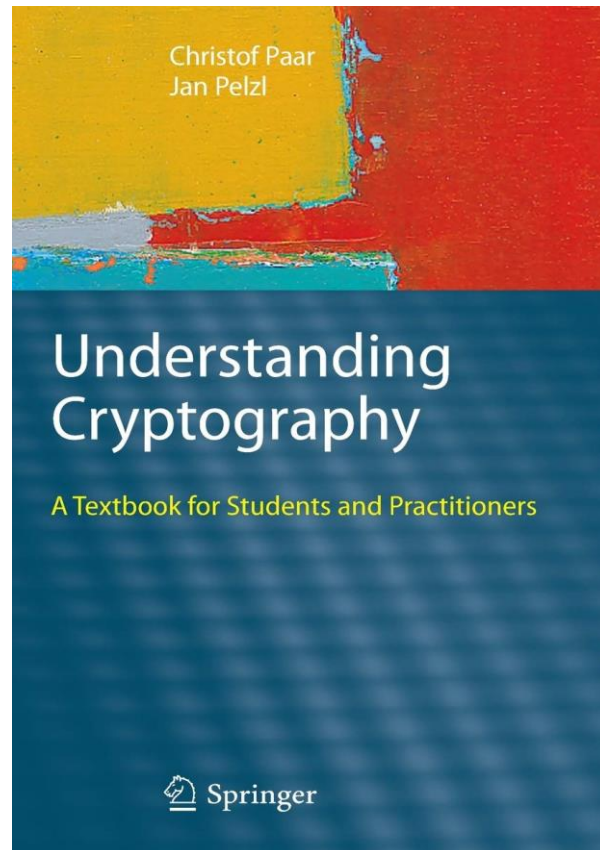
Collection modules : Provide different types of container data types.

4. METHODOLOGY

Downloading Books :

- We have downloaded the book *Understanding Cryptography written by Cristopher Paar and Jen Palzl* in plain text format for text processing from
(<https://swarm.cs.pub.ro/~mbarbulescu/cripto/Understanding%20Cryptography%20by%20Christof%20Paar%20.pdf>)

The downloaded file is .txt file



Cover of the book selected for text analysis

Importing the text

In this step we created a function (`txt_file_to_string`) to read the text imported from the book and convert it into string for processing in python. This function takes as input the path of the file to be read and returns the content of the file in string format.

```
Understanding Gp 7elcolele-lolly Ronee tee eto ees Understanding Cryptography Christof Paar - Jan Pelzl Understandi
ng Cryptography A Textbook for Students and Practitioners Foreword by Bart Preneel D) Springer Prof. Dr.-Ing. Christof
Paar Chair for Embedded Security Department of Electrical Engineering and Information Sciences Ruhr-Universitat Bochum
44780 Bochum Germany cpaar@crypto.rub.de ISBN 978-3-642-04 100-6 DOI 10.1007/978-3-642-04101-3 Springer Heidelberg
Dordrecht London New York ACM Computing Classification (1998): E.3, K.4.4, K.6.5. Library of Congress Control Number: 200
9940447 © Springer- Verlag Berlin Heidelberg 2010 This work is subject to copyright. All rights are reserved, whether th
e whole or part of the material is Dr.-Ing. Jan Pelzl escript GmbH – Embedded Security Zentrum fiir IT-Sicherheit Lise-
Meitner-Allee 4 44801 Bochum Germany jpelzl@escript.com e-ISBN 978-3-642-04101-3 concerned, specifically the rights o
f translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other w
ay, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions o
f the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained
from Springer. Violations are liable to prosecution under the German Copyright Law. The use of general descriptive names,
registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that
such names are exempt from the relevant protective laws and regulations and therefore free for general use. Cover desig
n: KuenkelLopka GmbH Printed on acid-free paper Springer is part of Springer Science+Business Media (www.springer.com)
To Flora, Maja, Noah and Sarah as well as to Karl, Greta and Nele While writing this book we noticed that for some reas
on the names of our spouses and children are limited to five letters. As far as we know, this has no cryptographic releva
nce. Foreword Academic research in cryptology started in the mid-1970s; today it is a mature re- search discipline with
an established professional organization (IACR, International Association for Cryptologic Research), thousands of researc
```

Text before Preprocessing

Text Pre-Processing and Tokenization :

1. Removing prefix and suffix to narrow down to text from eBook - The book from the website had additional prefix and suffix in its .txt file, apart from the contents of the eBook.
2. Lowercase - We converted our string to lowercase using the string function `string.lower()`.
3. Expansion of some Contractions - We expanded some generic contractions. For example : can't to can not, all instances of 'll to will etc. This is not very accurate and will lead to some incorrect expansion since disambiguation to the right expansion is not deterministic but this will be correct for most cases.

4. Removal Of Punctuations - We removed all the punctuation using regular expressions in two steps by first replacing everything other than word and whitespace characters with empty string and then replacing _ (underscore, which is considered part of word in python) by empty string.
5. Removing unnecessary repeated words
6. Replacing one or more continuous white space characters with single space to make the string evenly spaced.
7. Replacing numbers from integer to word form.
8. We tokenized the string into single words using word_tokenize() function imported from NLTK and stored them. Then we lemmatized these lists as we plan to analyze word frequencies later, therefore reducing the words to their lemma form will be suitable.

start of the project xi table of contents one introduction to cryptography and data security eleven overview of cryptology and this book twelve symmetric cryptography zero c eee eee eee eee one hundred and twenty-one basics zero ccc cece eens one hundred and twenty-two simple symmetric encryption the substitution cipher thirteen cryptanalysis zero zero eee eee e ee one hundred and thirty-one general thoughts on breaking cryptosystems one hundred and thirty-two how many key bits are enough fourteen modular arithmetic and more historical ciphers one hundred and forty-one modular arithmetic two eee eee e ee one hundred and forty-two integer rings000 zero ccc cece eens one hundred and forty-three shift cipher or caesar ciphe r twenty thousand and five one hundred and forty-four affine cipher zero 0c cee ees fifteen discussion and further readin g zero c cece eee eee sixteen lessons learned zero problems twenty-two cence eee eee two stream ciphers zero zero cece tw enty-one introduction zero e een eee two hundred and eleven stream ciphers vs block ciphers twenty-two thousand, two hund red and forty-five two hundred and twelve encryption and decryption with stream ciphers twenty-two random numbers and an unbreakable stream cipher two hundred and twenty-one random number generators0000000004 two hundred and twenty-two the on etime pad zero eee eee eee two hundred and twenty-three towards practical stream ciphers twenty-four twenty-three shift r egisterbased stream ciphers two hundred thousand and five two hundred and thirty-one linear feedback shift registers lfsr two hundred and thirty-two knownplaintext attack against single lfsrs two thousand, three hundred and thirty trivium twen ty-two n eee eee ee twenty-four discussion and further reading zero eee ee eee eleven xiii xiv table of contents twenty-f ive lessons learned zero c cece eee eee fifty problems zero neces fifty-two the data encryption standard des and alternat ives fifty-five thirty-one introduction to des zero o eee eee fifty-six three hundred and eleven confusion and diffusion zero 000s fifty-seven thirty-two overview of the des algorithm zero zero cee eee eee fifty-eight thirty-three internal st ructure of des two hundred and twenty eee ee sixty-one three hundred and thirty-one initial and final permutation20 sixty

Text after Preprocessing

```
Out[5]: ['start',  
         'of',  
         'the',  
         'project',  
         'xi',  
         'table',  
         'of',  
         'contents',  
         'one',  
         'introduction',  
         'to',  
         'cryptography',  
         'and',  
         'data',  
         'security',  
         'eleven',  
         'overview',  
         'of',  
         'cryptology',  
         '']
```

Tokenization before removing stop words

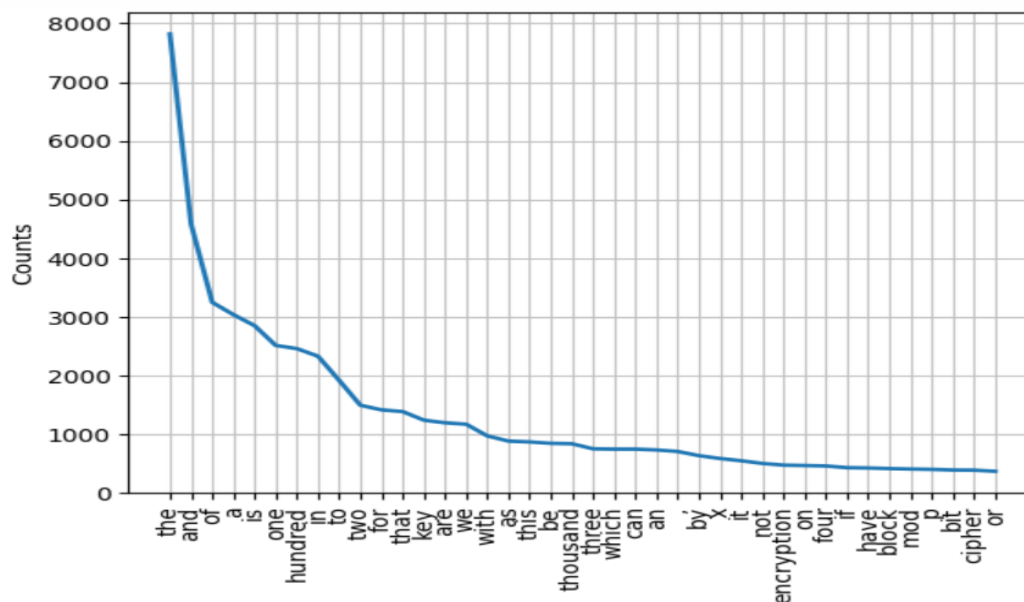
Plotting Frequency distribution of tokens :

- After tokenizing the text, we imported the function `FreqDist()` from the module `nltk.probability` which is helpful in probability calculations, where frequency distribution counts the number of times that each outcome of an experiment occurs. We stored the frequency distribution of the two strings in the `FreqDist` object by passing the tokenized and lemmatized lists to the above function.

```
FreqDist({'the': 7820, 'and': 4580, 'of': 3247, 'a': 3037, 'is': 2850, 'one': 2512, 'hundred': 2458, 'in': 2328, 'to': 1917, 'two': 1493, ...})
```

Frequency before removing stop words

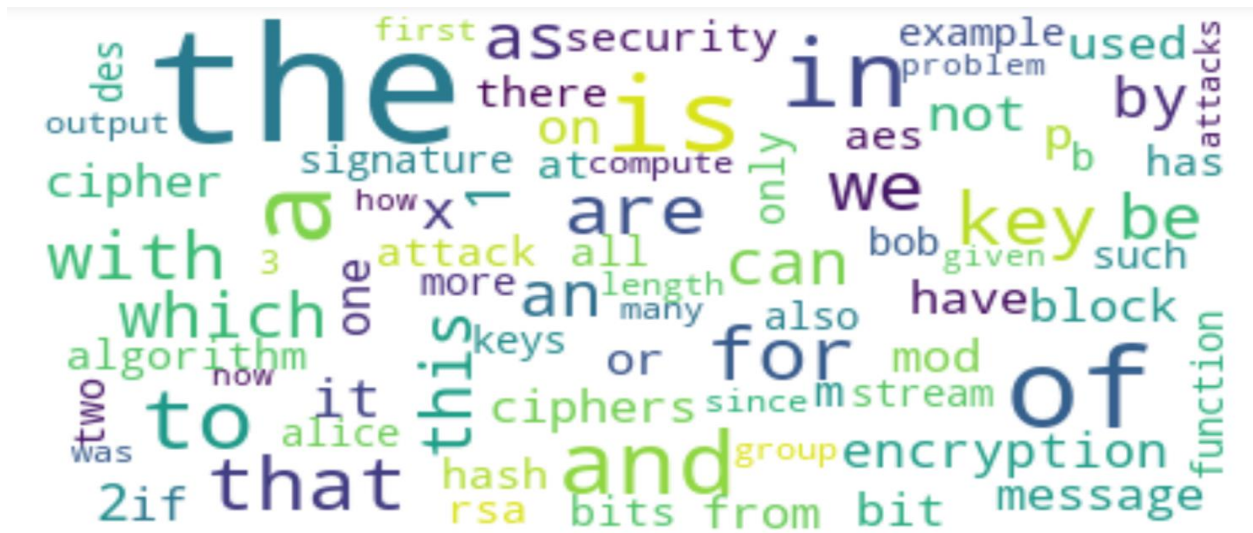
- For plotting the graph of frequency distribution we imported the function `figure()` from the module `matplotlib.pyplot` which is used to create a figure object. The whole figure is regarded as the figure object. Next we plotted frequency distribution graph:



Frequency distribution of top 40 words (with stop words)

Creating Word Cloud

- For creating word cloud we imported Counter from the module Collections. Then we imported wordcloud from the module wordcloud. Then we used functions plt.figure() , plt.axis(), plt.show() for the visualization of wordcloud. We made the word cloud for the top 80 most frequently used words which is attached below :



Word cloud (with stop words)

Inference :

- It is giving the visual representation of the most frequent words in the book Understanding Cryptography

- We observe that words like the, and, of, a are among the words that appear bigger and their frequency is also high in the previous plot.

- Infact, the overall word cloud is heavily dominated by stop words.

- Thus we can infer that the stop words occur in high frequency in the text.

Removing stop words and creating word cloud again:

- We remove the stop words from the text. We used the nltk.corpus stopwords from English and removed them from the text. We updated the frequency distribution of words in the texts and repeated the above process to obtain a frequency distribution, tokenization, frequency and Word cloud without stop words which is attached below :



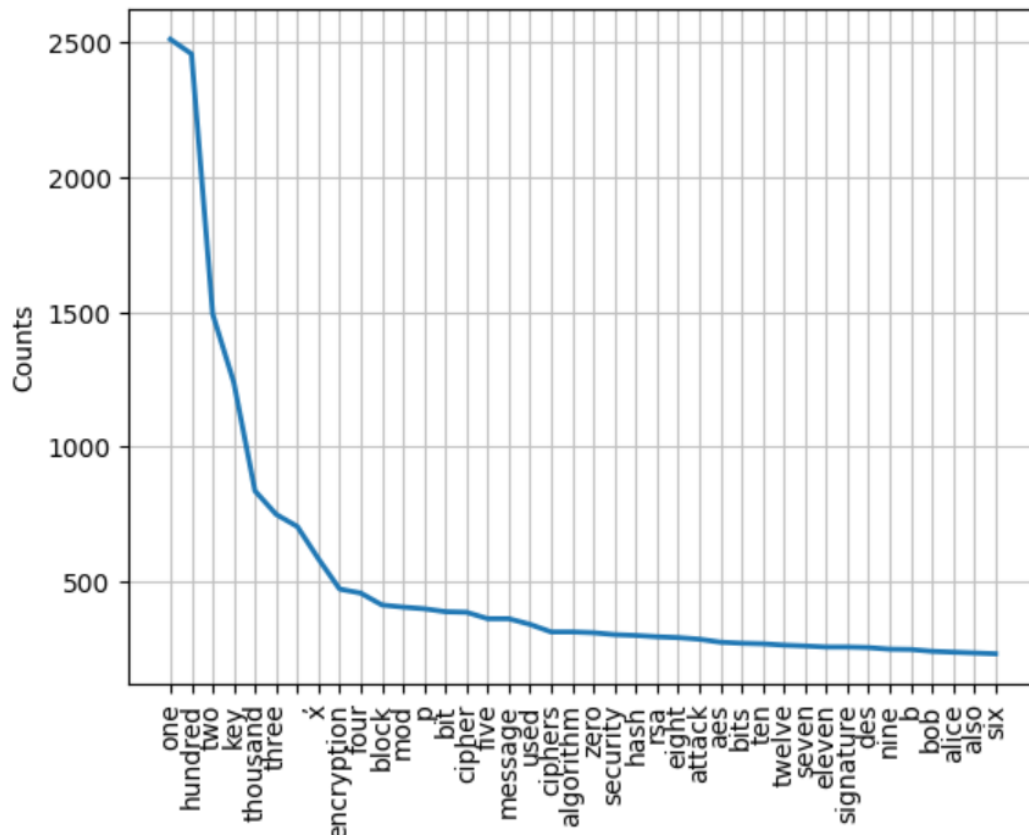
Word cloud (without stop words)

```
Out[11]: ['start',
           'project',
           'xi',
           'table',
           'contents',
           'one',
           'introduction',
           'cryptography',
           'data',
           'security',
           'eleven',
           'overview',
           'cryptology',
           'book',
           'twelve',
           'symmetric',
           'cryptography',
           'zero',
           'c',
```

Tokenization after removing stop words

```
FreqDist({'one': 2512, 'hundred': 2458, 'two': 1493, 'key': 1239, 'thousand': 837, 'three': 750, ',': 705, 'x': 585, 'encryp  
tion': 473, 'four': 458, ...})
```

Frequency after removing stop words



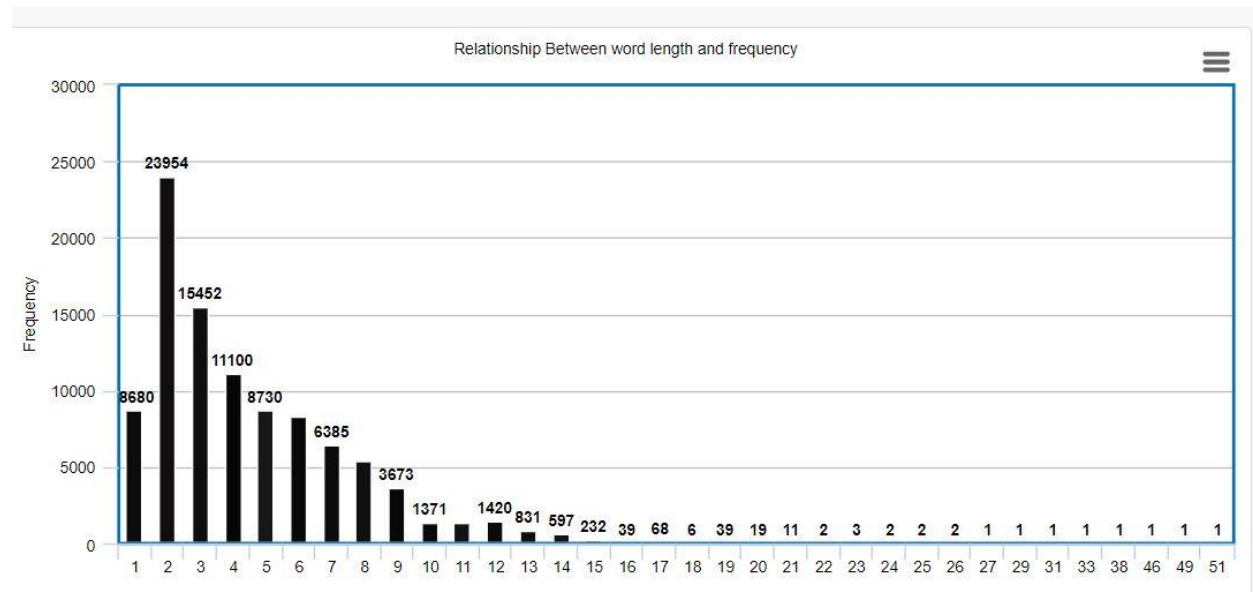
Frequency distribution of top 40 words(without stop words)

Inference :

- It is giving the visual representation of the most frequent words in the book Understanding Cryptography after removal of stop words.
- We observe that words like one, hundred, two and key are now amongst the words that appear bigger as their frequency is higher relatively after removal of stop words.

Relationship between the word length and frequency :

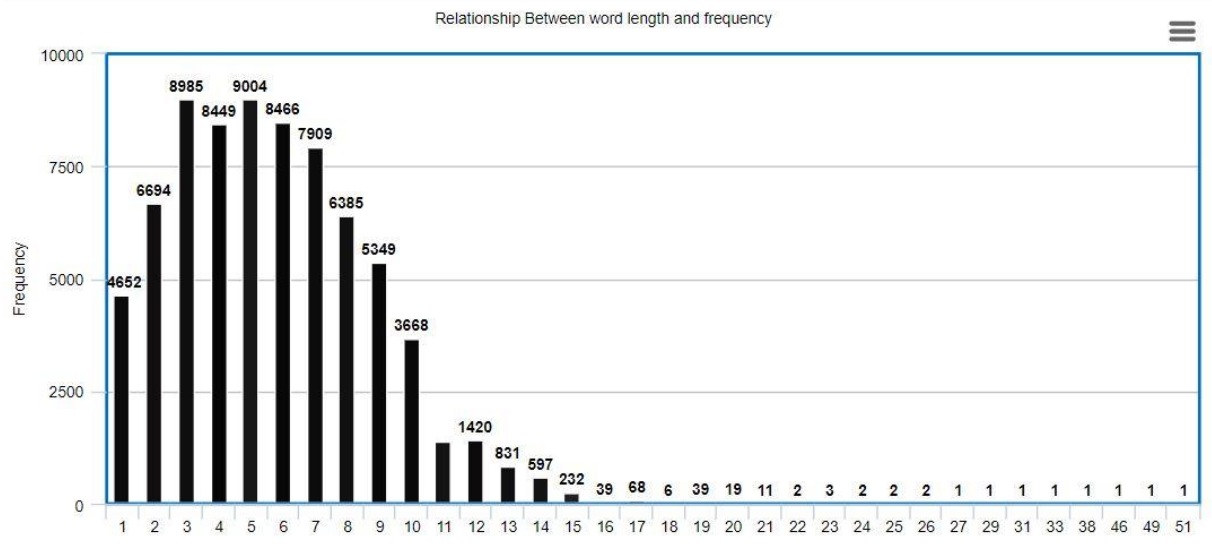
- We made an ordered dictionary to store the frequency of different word lengths in the text. The word lengths varied between 1 to 51. Next we plotted a bar chart showing the frequency of different word lengths. We did this for the book twice, once with the stop words and once after removal of stop words. These plots are attached below :



Relationship between the word length and frequency (with stop words)

Inference :

- We have plotted a bar chart for the words of different length and their frequency. But in this chart we have included stop words, so the words of small length have large frequencies.
- Words of length 2 have the highest frequency.
- We can infer that words of lengths between 2 and 4 (inclusive) form the majority of the text.



Relationship between the word length and frequency (without stop words)

Inference :

- We have plotted a bar chart for the words of different length and their frequency. But for this chart we have not included stop words.
- Word length 5 has the highest frequency.
- We can infer that words of lengths between 3 and 6 (inclusive) form the majority of the text (after stop word removal).
- Comparing the above range with the chart for stop words included, we can also infer that stop words mainly have length between 2 and 3 (inclusive).

5. CONCLUSION

Working on this project was a great learning experience for us in understanding the subject as well as team coordination. We all had surface-level knowledge about all the processes in text analytics but this project has helped us gain a better understanding of text processing using NLP techniques in python. Using python libraries and in-built toolkits, we came to a conclusion that this project highlights the basic understanding of text preprocessing, PoS tagging, Tokenization, etc. The Graphs, Bar Charts, Word clouds represented by using matplotlib helped us more in understanding the output and it is also beneficial for the visual representation of the data. Overall this was a great learning experience and it has encouraged us to explore more in the fields of NLP.

REFERENCES

- https://www.researchgate.net/publication/319164243_Natural_Language_Processing_State_of_The_Art_Current_Trends_and_Challenges
- <http://librarycarpentry.org/lc-tdm/index.html>
- <https://www.analyticsvidhya.com/blog/2021/06/text-preprocessing-in-nlp-with-python-codes/>
- <https://www.mygreatlearning.com/blog/nltk-tutorial-with-python/#3>
- <https://www.geeksforgeeks.org/text-analysis-in-python-3/>