

NLP PROJECT

ROUND 1 REPORT

TEAM NAME : Nlp_Learners

MOHIT JAIN (19UCS050)

AYUSH KABRA (19UCS113)

ARYAN VYAS (19UCS117)

Github repository: https://github.com/Aryan-Vyas/NlpLearners_NlpProject

Books used :

1. A town is drowning....by frederik pohl : T1
2. Nick Carter Stories....by Nick Carter : T2

TASKS

1. Import the text from two books, let's call it as T1 and T2.
2. Perform simple text pre-processing steps and tokenize the text T1 and T2.
3. Analyze the frequency distribution of tokens in T1 and T2 separately.
4. Create a Word Cloud of T1 and T2 using the token that you have got.
5. Remove the stopwords from T1 and T2 and then again create a word cloud.
6. Compare with word clouds before the removal of stopwords.
7. Evaluate the relationship between the word length and frequency for both T1 and T2.
8. Do PoS Tagging for both T1 and T2 using anyone of the four tagset studied in the class and Get the distribution of various tags

Python Libraries used in this project:

NLTK - Used for Tokenizing, Lemmatization and Removing Stopwords

Re - Used to remove URLs and Decontract Contractions in English Language

Wordcloud - Used to create WordClouds from Tokenized Data

Matplotlib - Used to Visualize our text data

DATA PREPROCESSING STEPS

1. Removing not needed text of Book

We will remove the documentation part of the text that is of no use to us.

2. Remove chapter name

We will remove chapter names as given in the task.

3. Convert all data to lowercase

We will convert all text data to lowercase.

4. Remove link

We will remove urls using regular expression

5. Remove Punctuations

6. Lemmatization

We do Lemmatization with the help of `WordNetLemmatizer()` function

Data Preparation

We apply all the functionalities we added above and Prepare our data for analysis

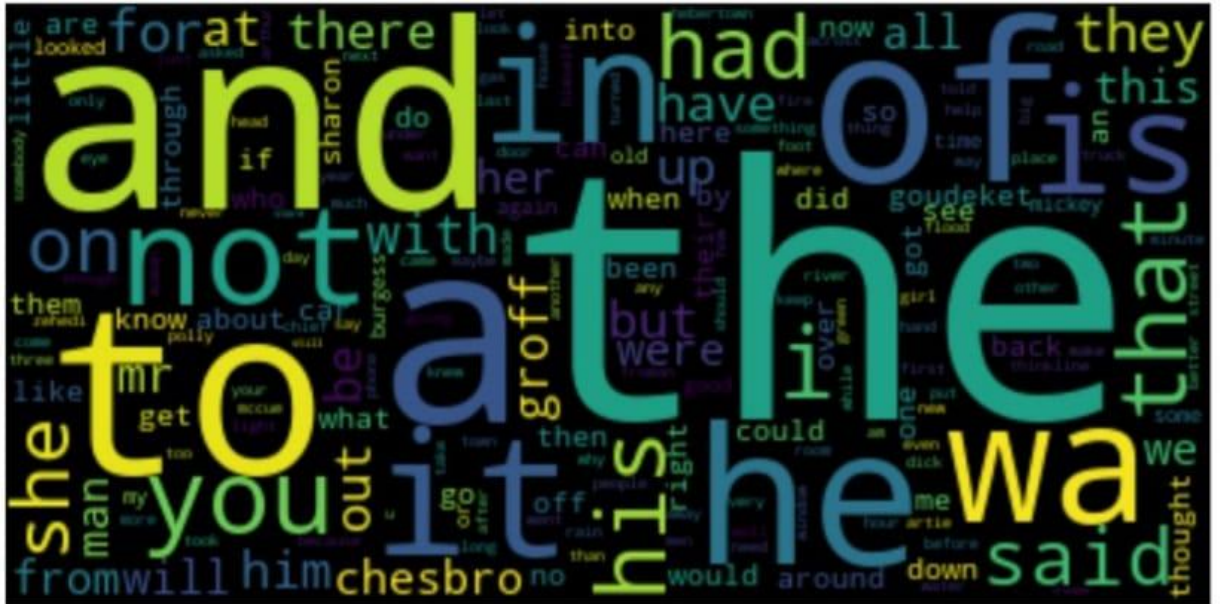
Problem Statements and Inferences

- Tokenization
- Analyze the frequency distribution of tokens in T1 and T2 separately
we use python library pandas First we tokenize the given data, and then we plot a bargraph of top 10 most frequent words

Generating a word cloud from T1 and T2

For this we use python library wordcloud and its function `WordCloud`.

For T1



For T2



Inferences

- Words like 'and', 'a', 'of', 'to' and 'the' are the most frequently used words in T1
- Words like 'and', 'the', 'he', 'to' are frequently used words in T2
- The words seen in wordcloud are called stop words

Generating new word clouds after removing stopwords

To remove stopwords, we use `STOPWORDS` function in `nltk`.

For T1



For T2

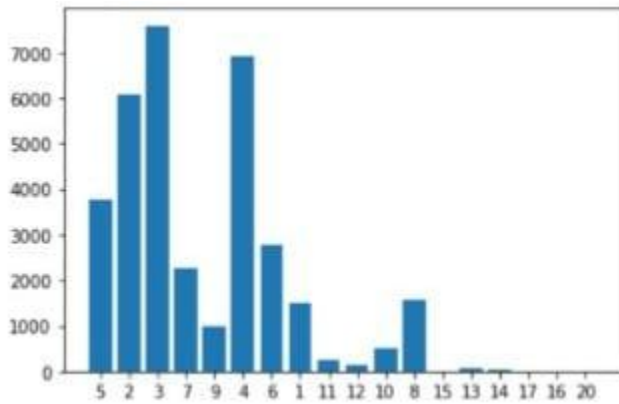


Inferences

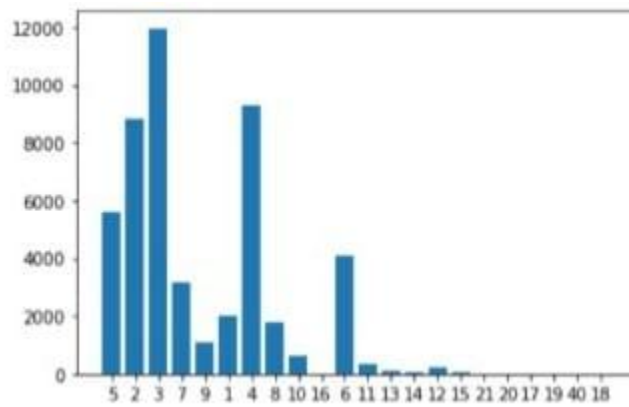
- Now the stopwords are removed
- New words like 'sharon', 'groff' in T1 and 'nick', 'carter', 'floyd', etc. can be seen in wordcloud which indicates the name of characters and other important stuffs around whom the book revolves.
- Therefore, after removing stopwords we get more meaning and understanding of the book.

Evaluate the relationship between the word length and frequency for both T1 and T2.

For T1



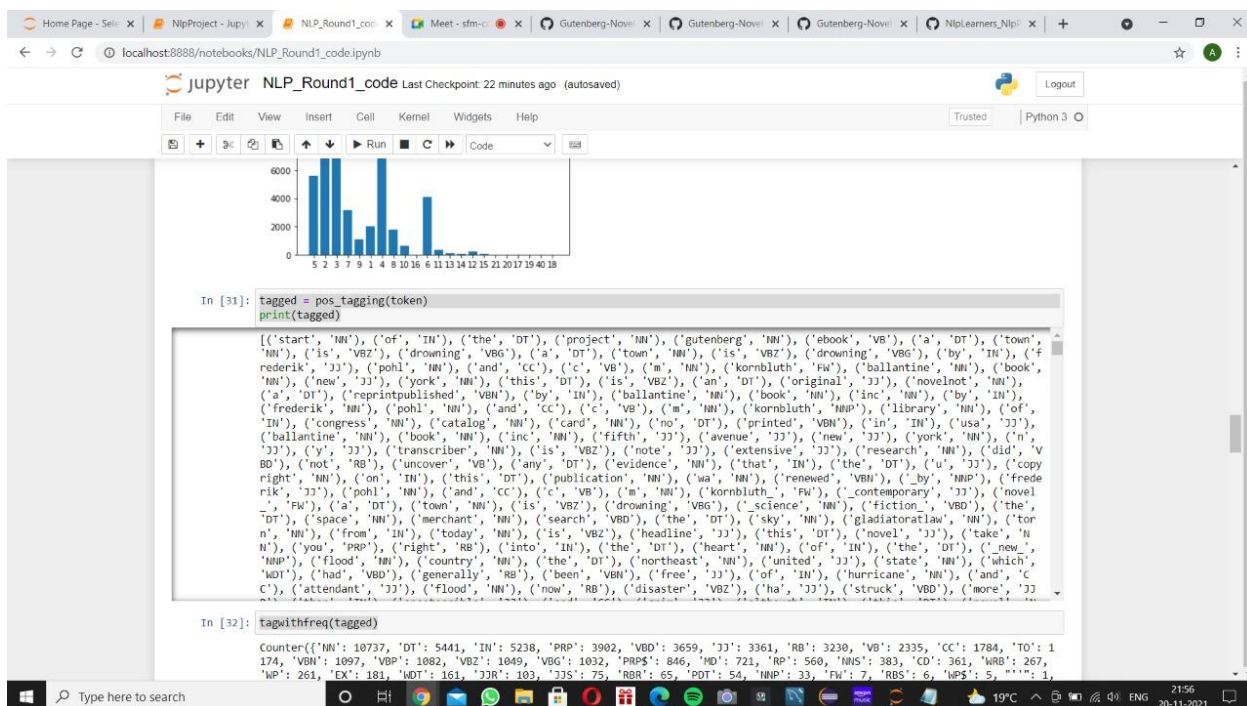
For T2



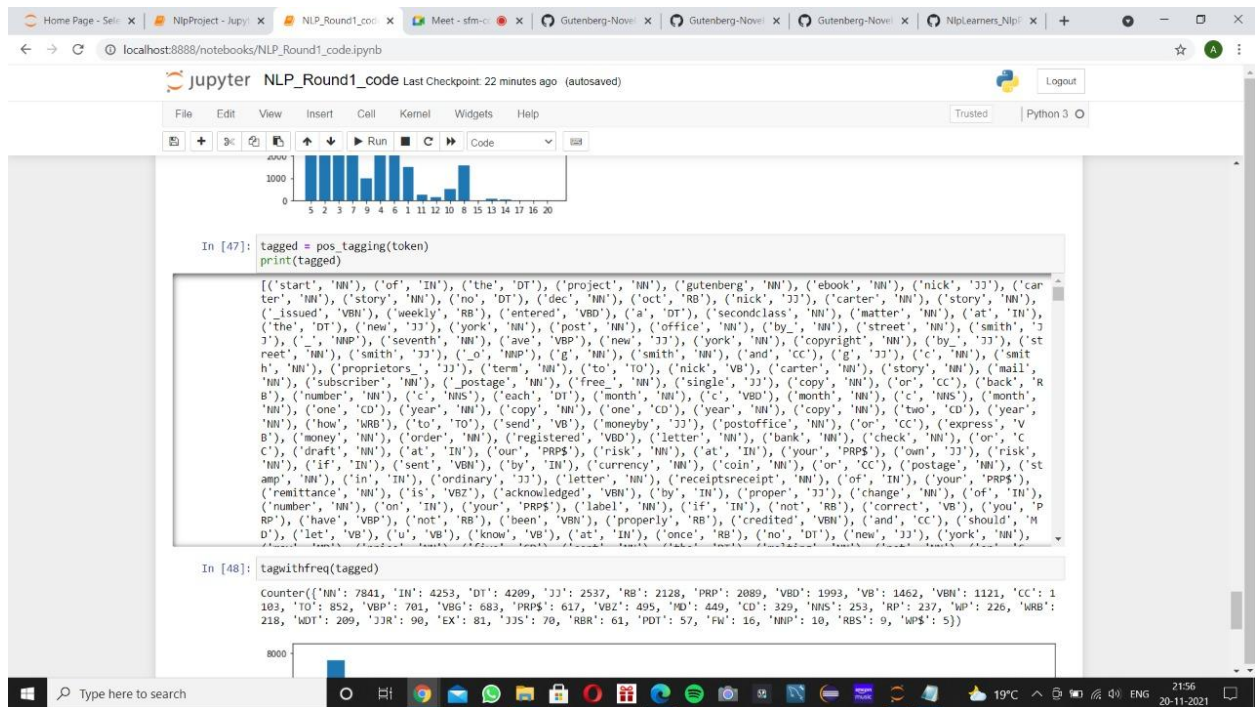
Performing POS Tagging

We will now perform the POS Tagging on T1 and T2 using inbuilt functions of nltk namely `post_tag()` which uses Penn Treebank tag set to perform POS tagging.

For T1



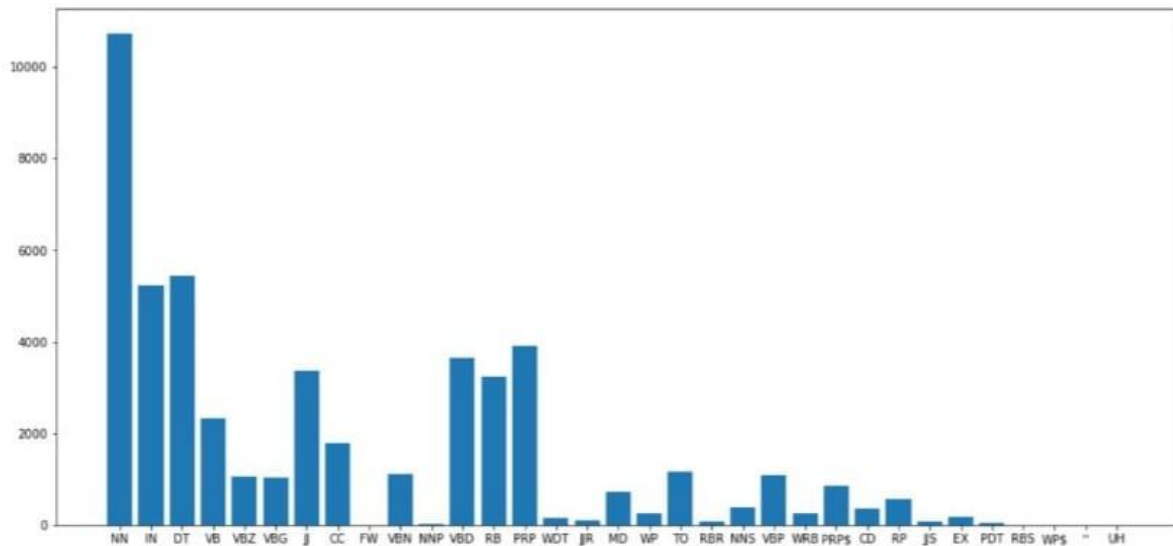
For T2



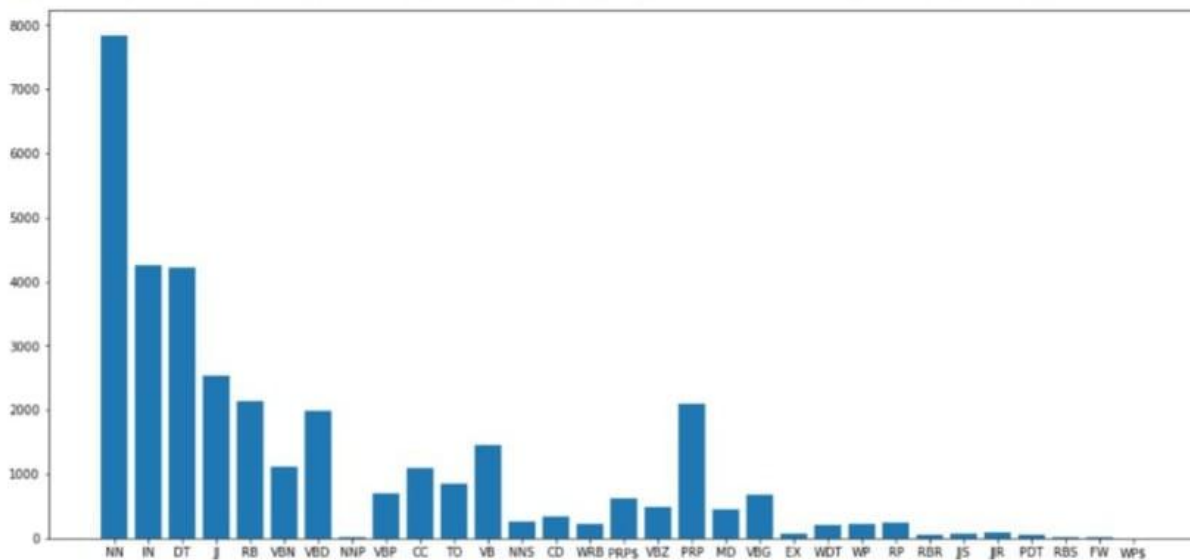
Frequency Distribution of Tags

Now, we plot the frequency distribution of tags after POS tagging on T1 and T2. For this we take the help of Counter() function from collections python library and FreqDist() function from nltk python library.

FOR T1



FOR T2



Inferences

From the above results we infer that the highest occurring tag is 'NN', and 'Determinant' Tags are on the lower frequency side. This is largely due to the removal of stopwords before POS Tagging.

