# DOCUMENTATION OF SENTIMENTAL ANALYSIS

The news article that is used for this project is from the link:
"https://astronomy.com/news/2023/03/two-potentially-habitable-planets-found-orbiting-distant-star".

## Two potentially habitable planets found orbiting distant star

*With more than 5,000 known exoplanets, astronomers are shifting their focus from discovering additional distant worlds to identifying which are good candidates for further study.*

By Joey Rodriguez, The Conversation, Andrew Vanderburg, The Conversation | Published: Wednesday, March 8, 2023

### THE DATA:

```
13    url = "https://astronomy.com/news/2023/03/two-potentially-habitable-planets-found-orbiting-distant-star"
14    res = requests.get(url)
15
16    s = BeautifulSoup(res.content, 'html.parser')
17    cont = s.get_text(strip=True)
18
19    cont = re.sub('[^a-zA-Z]', ' ', cont)
20    cont = cont.lower()
21    tk = nltk.word_tokenize(cont)
22    stopwords = nltk.corpus.stopwords.words('english')
23    tk = [tx for tx in tk if tx not in stopwords]
```

First the url is set and request is set in res. Then,

- **s = BeautifulSoup(res.content, 'html.parser')**: This line uses the BeautifulSoup library to parse the HTML content of the webpage that was retrieved using the requests library. The parsed content is stored in the variable s.

- **cont = s.get_text(strip=True):** This line extracts the text from the parsed HTML content stored in s using the get_text method of the BeautifulSoup object. The strip=True argument removes any leading or trailing white space from the text.

- **cont = re.sub('[^a-zA-Z]', ' ', cont):** This line uses the re library to substitute any character that is not a letter (i.e., anything that is not in the range a-z or A-Z) with a space. This effectively removes any punctuation, numbers, or other non-letter characters from the text.

- **cont = cont.lower():** This line converts all the remaining letters in cont to lowercase. This is often done to simplify text analysis by reducing the number of unique words that need to be considered.

- **tk = nltk.word_tokenize(cont):** This line uses the word_tokenize function from the nltk library to split the text into a list of individual words (tokens).

- **stopwords = nltk.corpus.stopwords.words('english'):** This line loads a list of common English stopwords from the nltk library. Stopwords are words that are considered to be uninformative or redundant for text analysis purposes, such as "the", "a", and "an".

- **tk = [tx for tx in tk if tx not in stopwords]:** This line creates a new list of tokens (tk) that excludes any tokens that appear in the stopwords list. This helps to further clean the text data by removing uninformative words.
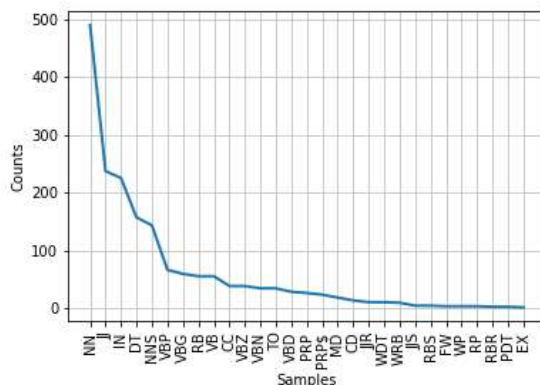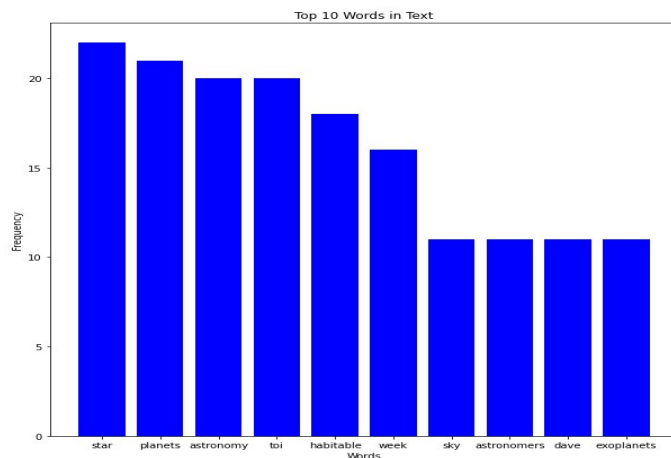
*EDA :*

```
26    # EDA
27
28    # bar chart for top 10 frequent words
29    fd = nltk.FreqDist(tk)
30    tw= fd.most_common(10)
31    lab = [w[0] for w in tw]
32    val = [w[1] for w in tw]
33    plt.figure(figsize=(10, 10))
34    plt.bar(lab, val, color='blue')
35    plt.title('Top 10 Words in Text')
36    plt.xlabel('Words')
37    plt.ylabel('Frequency')
38    plt.show()
39
40    # bigrams
41    bg = list(nltk.bigrams(tk))
42    gr = nx.Graph()
43    for bigram in bg:
44        if bigram[0] != bigram[1]:
45            gr.add_edge(bigram[0], bigram[1])
46    plt.figure(figsize=(50, 50))
47    pos = nx.spring_layout(gr, k=0.5, iterations=50)
48    nx.draw_networkx_nodes(gr, pos, node_color='lightblue', node_size=1000)
49    nx.draw_networkx_edges(gr, pos, edge_color='gray', width=1)
50    nx.draw_networkx_labels(gr, pos, font_size=15, font_family='sans-serif')
51    plt.title('Bigram Network Graph', fontsize=20)
52    plt.axis('off')
53    plt.show()
54
55    # wordcloud
56    wc = WordCloud().generate(cont)
57    plt.imshow(wc, interpolation='bilinear')
58    plt.axis("off")
59    plt.show()
60
61    # POS tag
62    tk1 = word_tokenize(cont)
63    pos_tags = nltk.pos_tag(tk1)
64    tag = nltk.FreqDist(tag for (word, tag) in pos_tags)
65    tag.plot()
66    plt.show()
67
68
```

- Bar chart for top 10 frequent words: This graph helps in identifying the most frequent words in the text, which can give an idea about the main topics or themes present in the text.

- Bigram network graph: This graph shows the relationships between the bigrams (two-word phrases) in the text, which can give insights into the co-occurrence patterns and semantic relationships between words.
- Wordcloud: This graph provides a visual representation of the most frequent words in the text, with the size of each word indicating its frequency. It can help in identifying the most prominent themes or topics in the text at a glance.
- Part-of-speech (POS) tag plot: This graph shows the distribution of different parts of speech (nouns, verbs, adjectives, etc.) in the text, which can give insights into the syntactic structure and complexity of the text.

**GRAPHS:**









**SENTIMENTAL ANALYSIS :**

```
69    # Sentiment analyzer
70
71    #nltk.download('punkt')
72    #nltk.download('stopwords')
73    #nltk.download('punkt')
74    nltk.download('stopwords')
75    aly = SentimentIntensityAnalyzer()
76    sc = aly.polarity_scores(cont)
77    print('Positive:', sc['pos'])
78    print('Negative:', sc['neg'])
79    print('Neutral:', sc['neu'])
80    print('Compound:', sc['compound'])
```

```
81
82    # Visualize the sentiment
83
84    # Pie chart
85    plt.pie(sc.values(), labels=sc.keys(), colors=['red', 'green', 'blue'])
86    plt.title('Sentiment Analysis Pie Chart')
87    plt.show()
88
89
90    # Bar Plot
91    sentiment_scores = {
92        'Positive': sc['pos'],
93        'Negative': sc['neg'],
94        'Neutral': sc['neu']
95    }
96    labels = list(sentiment_scores.keys())
97    values = list(sentiment_scores.values())
98    plt.bar(labels, values, color=['green', 'red', 'blue'])
99    plt.title('Sentiment Analysis Bar Plot')
100   plt.xlabel('Sentiment')
101   plt.ylabel('Score')
102   plt.show()
```

- The code analyzes the sentiment of a given text using the VADER (Valence Aware Dictionary and sEntiment Reasoner) algorithm.
- The sentiment analyzer is created by importing the SentimentIntensityAnalyzer class from the nltk.sentiment.vader module.
- The analyzer object is used to compute the polarity scores of the given text by calling the polarity_scores() method on it. The method returns a dictionary of scores for positive, negative, neutral, and compound sentiment.
- The scores are then printed to the console to show the relative strength of positive, negative, and neutral sentiment in the text. The compound score is also printed, which is a single value that represents an overall sentiment score ranging from -1 (most negative) to +1 (most positive).
- The code then visualizes the sentiment using two different plots: a pie chart and a bar plot.
  For the pie chart, the sentiment scores are passed to the plt.pie() method to create a pie chart showing the relative strength of each sentiment type.
  For the bar plot, a dictionary of sentiment scores is created, and the scores are plotted using the plt.bar() method to create a bar chart showing the relative strength of each sentiment type.

**GRAPHS:**



Sentiment Analysis Pie Chart



Sentiment Analysis Bar Plot