

ENTEVYUV 11.0

WEB SCRAPING AND TEXT ANALYSIS OF NEWS ARTICLES

Data Cleansing, EDA, NLP

PRESENTED BY
Vinodhini Rajamanickam

Index

1 Problem Statement

2 Project aim

3 Tools Used

4 Approach

5 Exploratory Data Analysis
(EDA) Insights

6 Conclusion

Problem Statement

Developing an Automated Solution Utilizing Web Scraping Tools for Efficient Extraction and Comprehensive Text Analysis of Data from Designated URLs to Uncover Valuable Insights

Project Aim

Enhancing Practical Skills: Exploring Web Content for Insights. This project hones web scraping, text analysis, and variable computation abilities. Gain valuable expertise in data science, natural language processing, and content analysis by deciphering content traits, sentiment, and themes from textual data

Tools and Technologies



PROGRAMMING LANGUAGE

Python

DATA LOADING AND MANIPULATION

Pandas

DATA VISUALIZATION

Matplotlib

Seaborn

DATA CLEANING

re

WEB SCRAPING

BeautifulSoup

request

NATURAL LANGUAGE PROCESSING

Nltk

Approach

- Collected data by scraping news articles from various URLs using the BeautifulSoup library in Python.
- Utilized BeautifulSoup to extract text content from web articles and stored them in text files.
- Tokenized article text to break it into individual words.
- Converted all text to lowercase to standardize text data.
- Removed special characters, numbers, and punctuation using regular expressions.
- Eliminated common stopwords to reduce noise in the text data.
- Applied stemming to reduce words to their base form (e.g., running → run).
- Calculated sentiment scores based on positive and negative word frequencies.
- Defined and calculated additional variables like sentiment scores, average sentence length, complex word percentage, etc.
- Converted the processed data from text files to an Excel file for structured storage.
- Utilized Pandas to load the Excel data into a DataFrame, facilitating exploratory analysis.
- Prepared a clean and preprocessed dataset for subsequent analysis.
- Ensured that the dataset was ready for exploratory analysis and modeling.

Benefits:

- Enhanced the dataset with calculated features for deeper insights.
- Enabled seamless data transformation and manipulation using DataFrame operations.

Challenges:

- Ensured accuracy in sentiment scoring and feature calculations.
- Maintained consistency in data preprocessing across different articles.

Outcome:

- Generated a well-structured dataset with enriched features, ready for detailed analysis.
- Prepared the data for visualization and uncovering patterns in sentiment and article characteristics.

Insights

Question 13: What is the average complexity (Complex_Word_Count) for articles with high and low subjectivity?

```
In [73]: 1 # Calculate average complexity for high and low subjectivity articles
2 avg_complexity_high_subjectivity = data[data["Subjectivity_Score"] > 0.5]["Complex_Word_Count"].mean()
3 avg_complexity_low_subjectivity = data[data["Subjectivity_Score"] <= 0.5]["Complex_Word_Count"].mean()
4
5 print(f"Average Complexity for High Subjectivity Articles: {avg_complexity_high_subjectivity:.2f}")
6 print(f"Average Complexity for Low Subjectivity Articles: {avg_complexity_low_subjectivity:.2f}")
7
```

Average Complexity for High Subjectivity Articles: 961.22
Average Complexity for Low Subjectivity Articles: 1023.75

```
In [33]: 1 # Sentiment Analysis Insights
2 positive_sentiment = data[data["Positive_Score"] > data["Negative_Score"]]
3 negative_sentiment = data[data["Negative_Score"] > data["Positive_Score"]]
4
5 print("Number of articles with strong positive sentiment:", len(positive_sentiment))
6 print("Number of articles with strong negative sentiment:", len(negative_sentiment))
```

Number of articles with strong positive sentiment: 65
Number of articles with strong negative sentiment: 45

Question 19: How many articles have a higher average word length and a high positive score?

```
In [79]: 1 # Count articles with high average word length and high positive score
2 high_word_length_high_positive = data[(data["Average_Word_Length"] > 5) & (data["Positive_Score"] > data["Negative_Score"])]
3
4 num_articles = len(high_word_length_high_positive)
5 print(f"Number of articles with high average word length and high positive score: {num_articles}")
6
```

Number of articles with high average word length and high positive score: 38

Question 17: Are articles with longer sentences more complex?

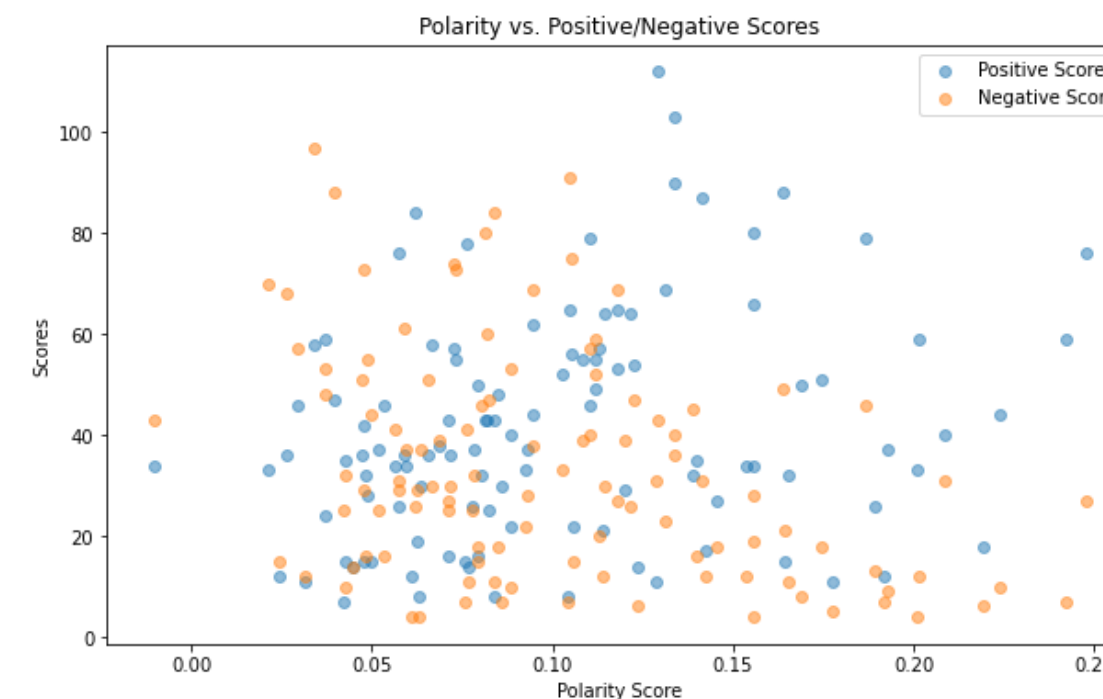
```
1 # Calculate correlation between Average_Sentence_Length and Complex_Word_Count
2 correlation = data["Average_Sentence_Length"].corr(data["Complex_Word_Count"])
3
4 print(f"Correlation between Average Sentence Length and Complex Word Count: {correlation:.2f}")
5
```

Correlation between Average Sentence Length and Complex Word Count: 1.00

In [35]:

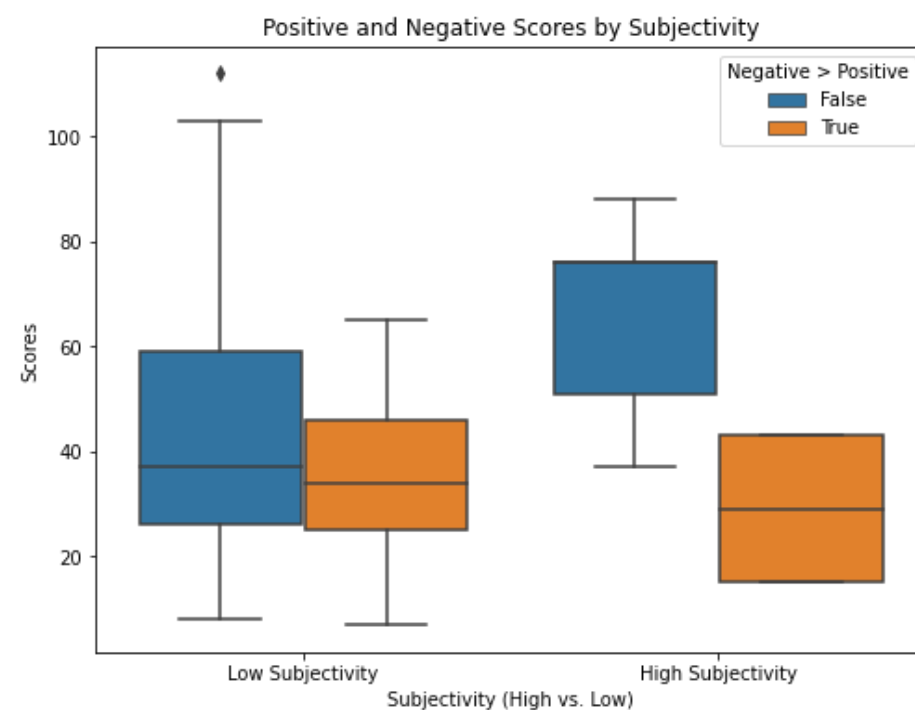
```
1 # Scatter plot of Polarity_Score vs. Positive_Score/Negative_Score
2 plt.figure(figsize=(10, 6))
3 plt.scatter(data["Polarity_Score"], data["Positive_Score"], label="Positive Score", alpha=0.5)
4 plt.scatter(data["Polarity_Score"], data["Negative_Score"], label="Negative Score", alpha=0.5)
5 plt.xlabel("Polarity Score")
6 plt.ylabel("Scores")
7 plt.title("Polarity vs. Positive/Negative Scores")
8 plt.legend()
9 plt.show()
10
```

The graph suggests that texts with higher polarity tend to have higher positive or negative scores, while texts with lower polarity tend to have lower scores. This means that texts that are more neutral or mixed in sentiment have lower scores than texts that are more clearly positive or negative.



Question 8: How does the distribution of positive and negative scores vary based on the subjectivity of articles?

```
1 # Box plot of Positive and Negative Scores by Subjectivity
2 plt.figure(figsize=(8, 6))
3 sns.boxplot(x=data["Subjectivity_Score"] > 0.5, y=data["Positive_Score"], hue=data["Negative_Score"] > data["Positive_Score"])
4 plt.xlabel("Subjectivity (High vs. Low)")
5 plt.ylabel("Scores")
6 plt.title("Positive and Negative Scores by Subjectivity")
7 plt.xticks([0, 1], ["Low Subjectivity", "High Subjectivity"])
8 plt.legend(title="Negative > Positive")
9 plt.show()
10
```



The graph shows that subjective texts have higher scores than objective texts, and that negative texts have higher scores than positive texts. This means that texts that are more emotional or negative have higher scores than texts that are more factual or positive.

Question 15: Do articles with higher subjectivity tend to have higher polarity scores?

```
1 # Calculate average polarity for high and low subjectivity articles
2 avg_polarity_high_subjectivity = data[data["Subjectivity_Score"] > 0.5]["Polarity_Score"].mean()
3 avg_polarity_low_subjectivity = data[data["Subjectivity_Score"] <= 0.5]["Polarity_Score"].mean()
4
5 print(f"Average Polarity for High Subjectivity Articles: {avg_polarity_high_subjectivity:.2f}")
6 print(f"Average Polarity for Low Subjectivity Articles: {avg_polarity_low_subjectivity:.2f}")
7
```

Average Polarity for High Subjectivity Articles: 0.13

Average Polarity for Low Subjectivity Articles: 0.10

Question 1: What is the average positive and negative sentiment across all articles?

```
1 # Calculate average positive and negative sentiment
2 average_positive_score = data["Positive_Score"].mean()
3 average_negative_score = data["Negative_Score"].mean()
4
5 print(f"Average Positive Score: {average_positive_score:.2f}")
6 print(f"Average Negative Score: {average_negative_score:.2f}")
7
```

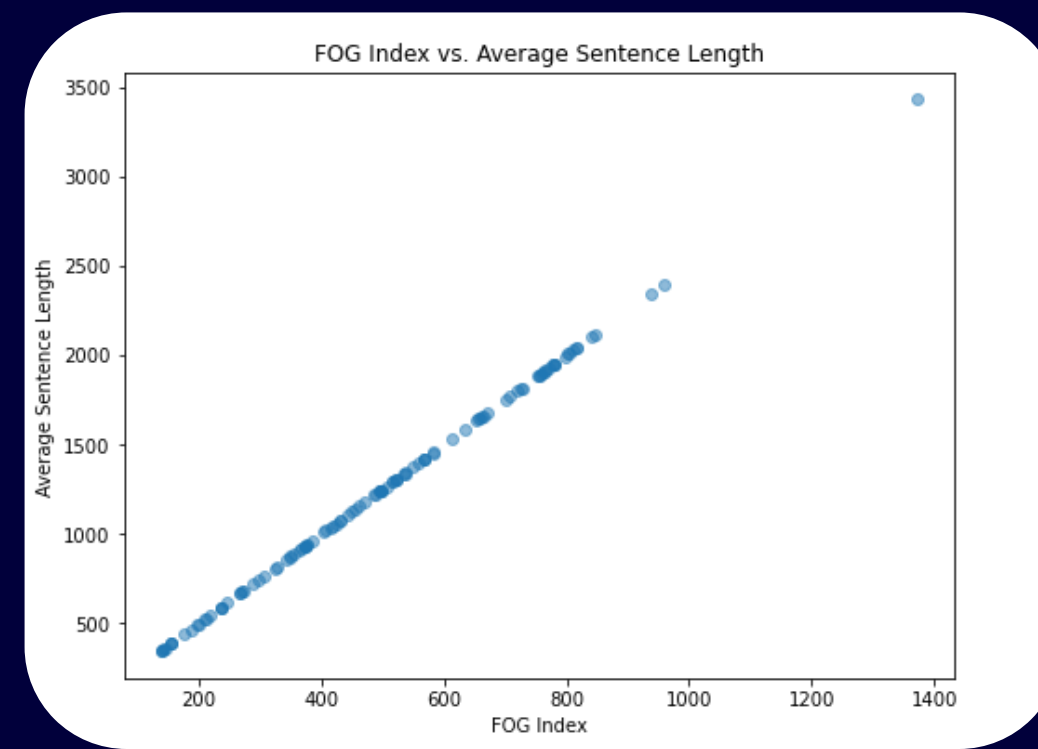
Average Positive Score: 41.09

Average Negative Score: 33.31

Question 4: Is there a correlation between average sentence length and complexity (FOG Index)?

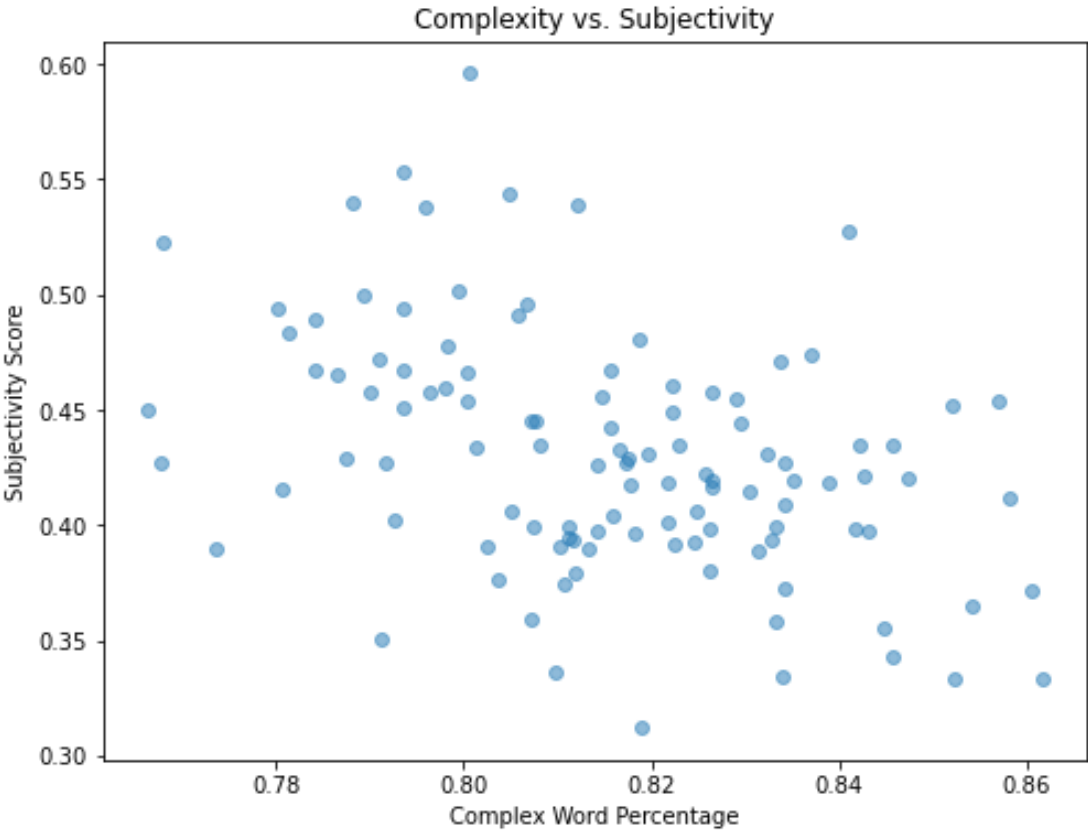
```
1 # Calculate correlation between Average_Sentence_Length and FOG_Index
2 correlation = data["Average_Sentence_Length"].corr(data["FOG_Index"])
3
4 print(f"Correlation between Average Sentence Length and FOG Index: {correlation:.2f}")
5
```

Correlation between Average Sentence Length and FOG Index: 1.00

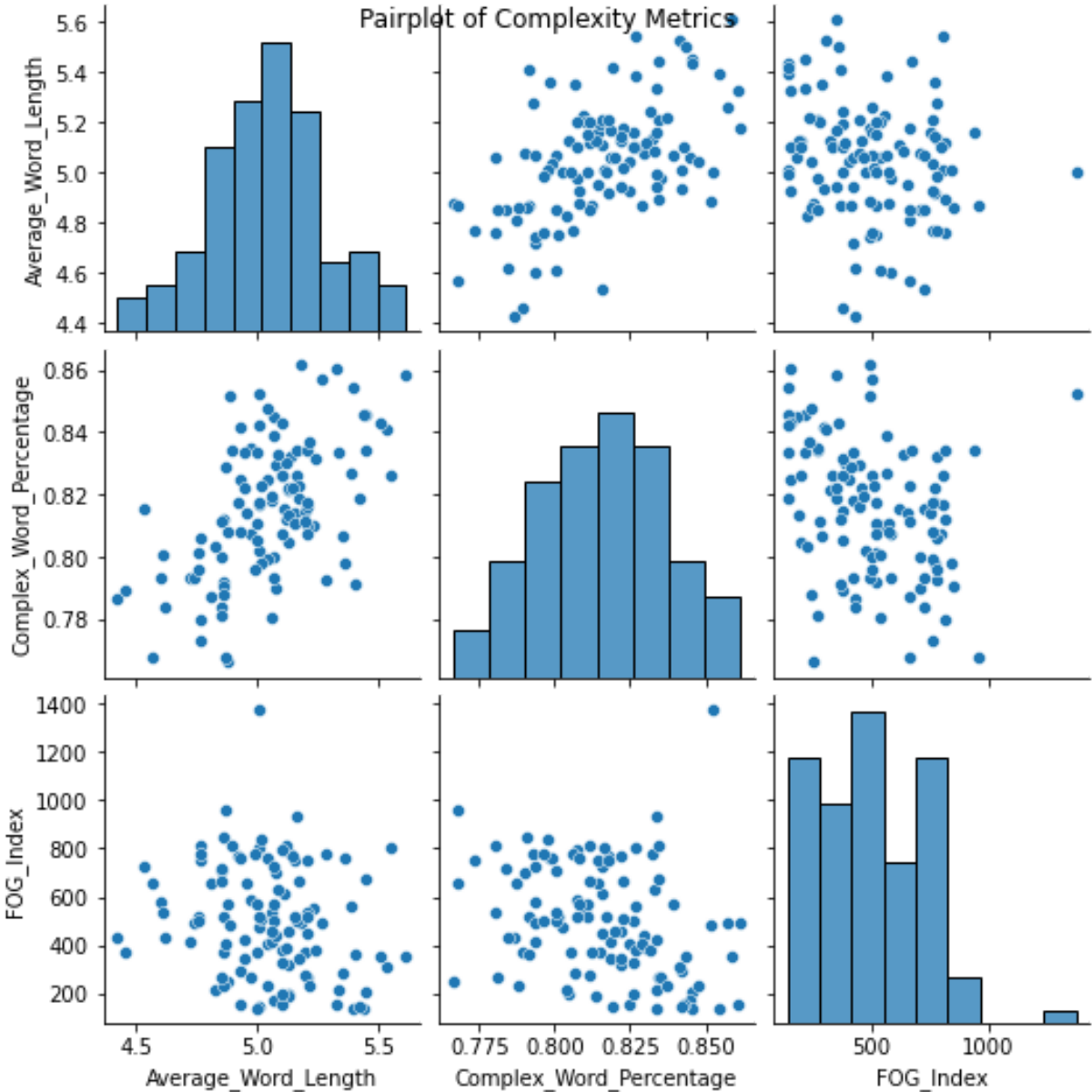


Question 7: Are articles with higher complexity more subjective?

```
1 # Scatter plot of Complex_Word_Percentage vs. Subjectivity_Score
2 plt.figure(figsize=(8, 6))
3 plt.scatter(data["Complex_Word_Percentage"], data["Subjectivity_Score"], alpha=0.5)
4 plt.xlabel("Complex Word Percentage")
5 plt.ylabel("Subjectivity Score")
6 plt.title("Complexity vs. Subjectivity")
7 plt.show()
8
```



The graph shows that texts with more hard words tend to have more opinion or emotion, and texts with fewer hard words tend to have less opinion or emotion.

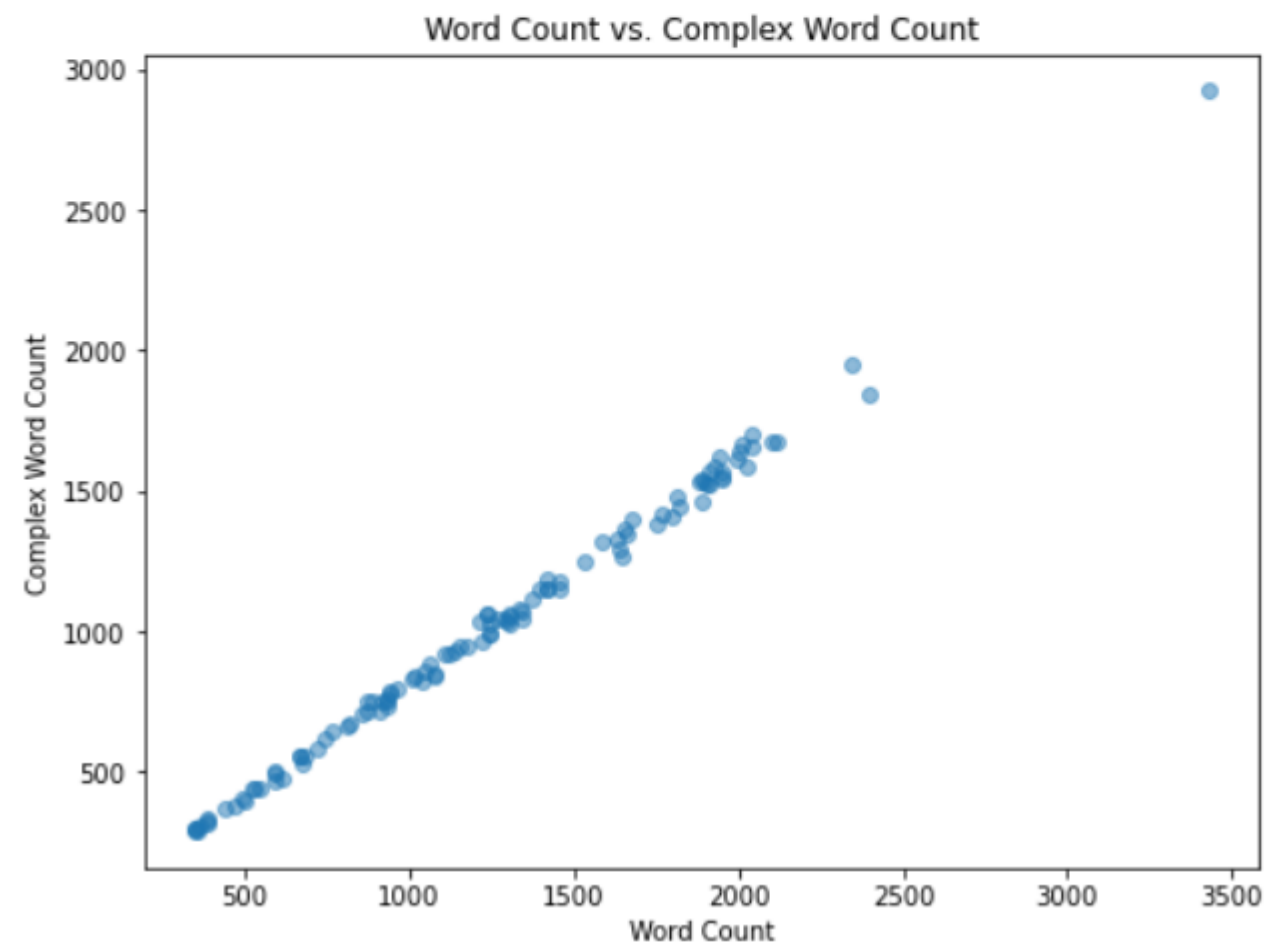


This image is a scatter plot matrix of complexity metrics. Complexity metrics are numbers that tell how hard a text is to read or understand. The matrix shows how different complexity metrics are related to each other. The plots on the diagonal show how many texts have different values of the same complexity metric.



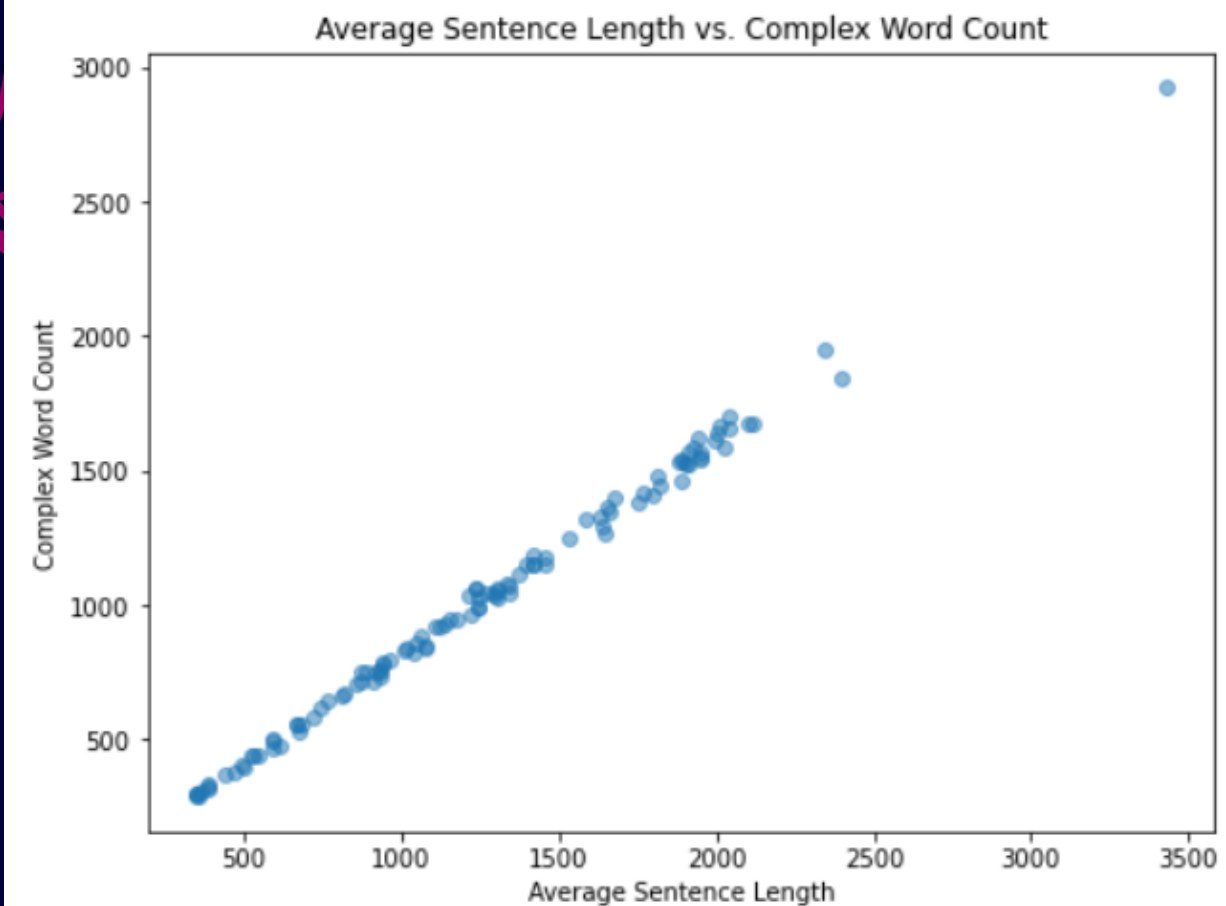
Question 11: Are longer articles generally more complex?

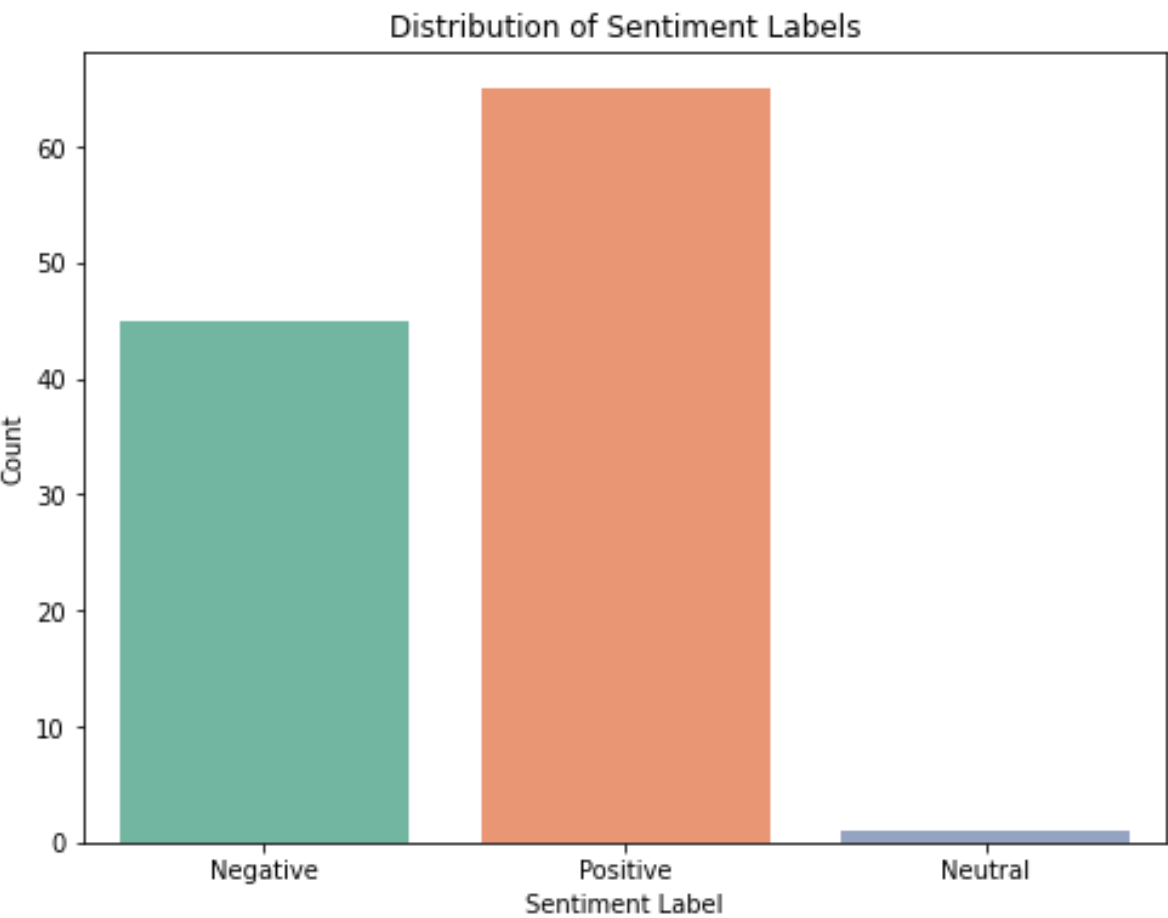
```
1 # Scatter plot of Word_Count vs. Complex_Word_Count
2 plt.figure(figsize=(8, 6))
3 plt.scatter(data["Word_Count"], data["Complex_Word_Count"], alpha=0.5)
4 plt.xlabel("Word Count")
5 plt.ylabel("Complex Word Count")
6 plt.title("Word Count vs. Complex Word Count")
7 plt.show()
8
```



Question 12: How does the average sentence length correlate with the complexity of articles?

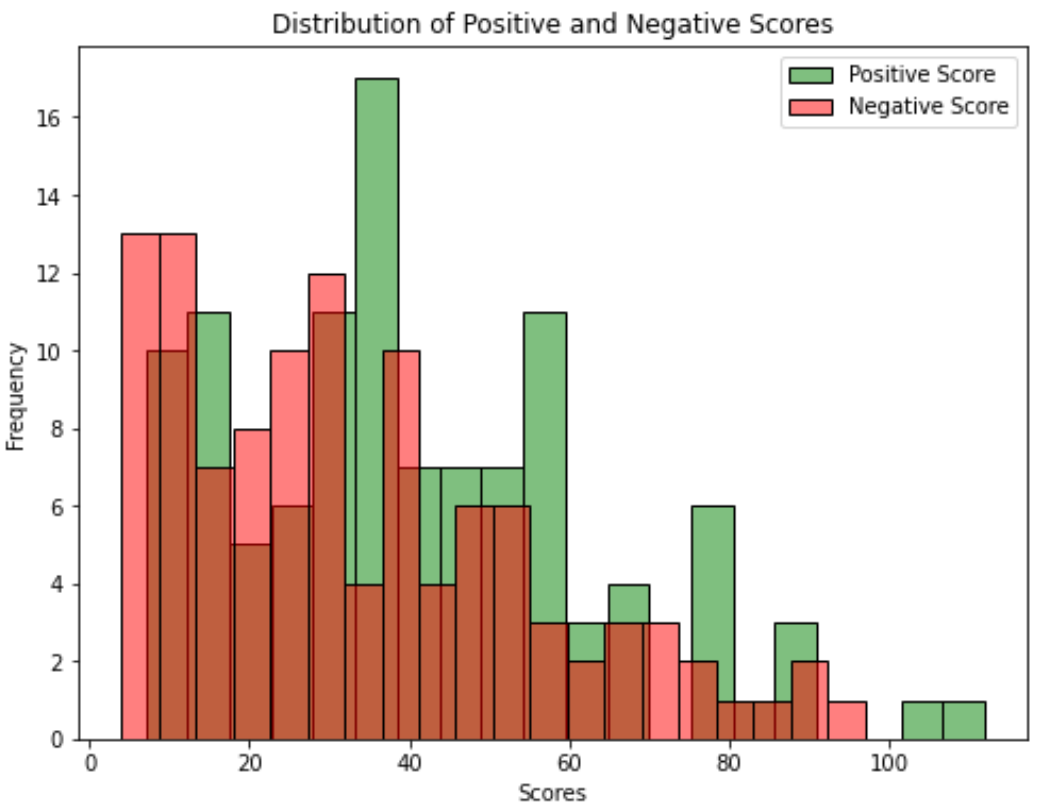
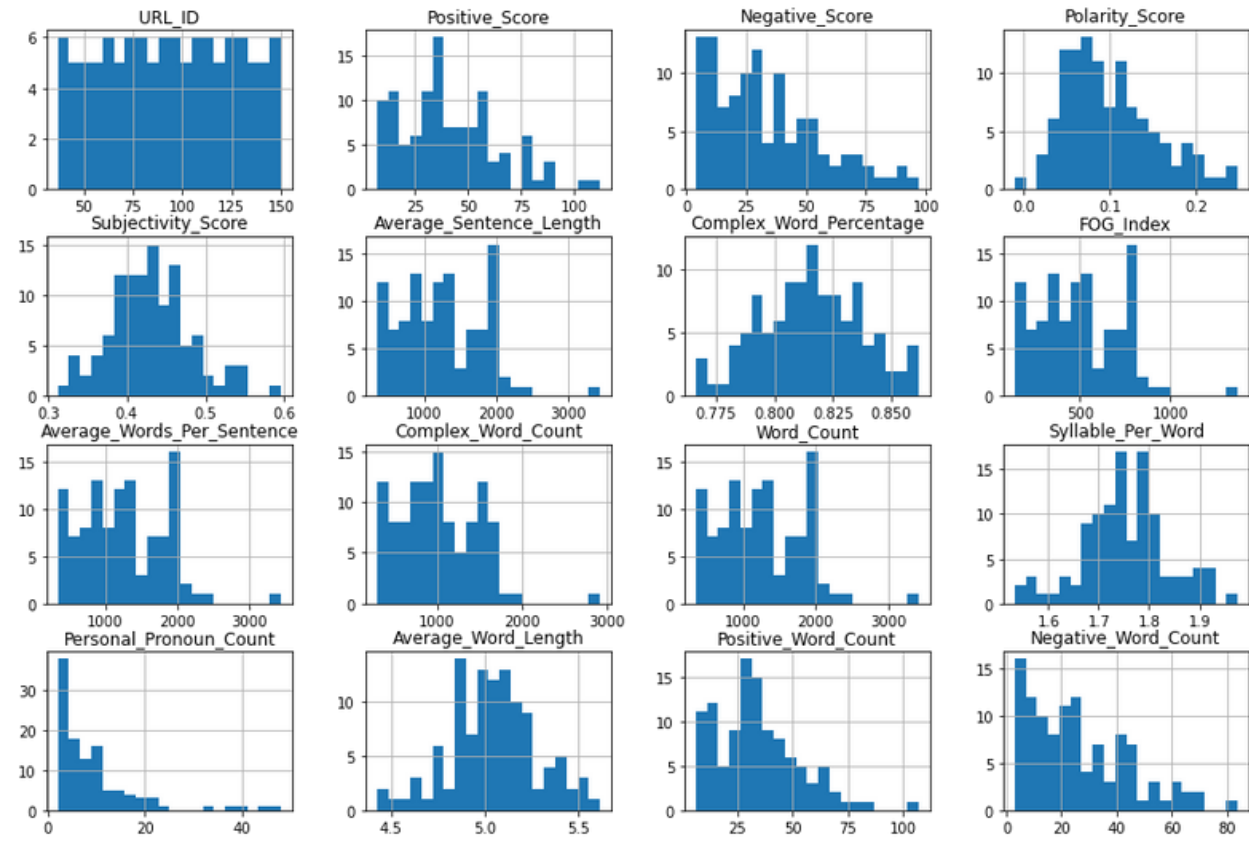
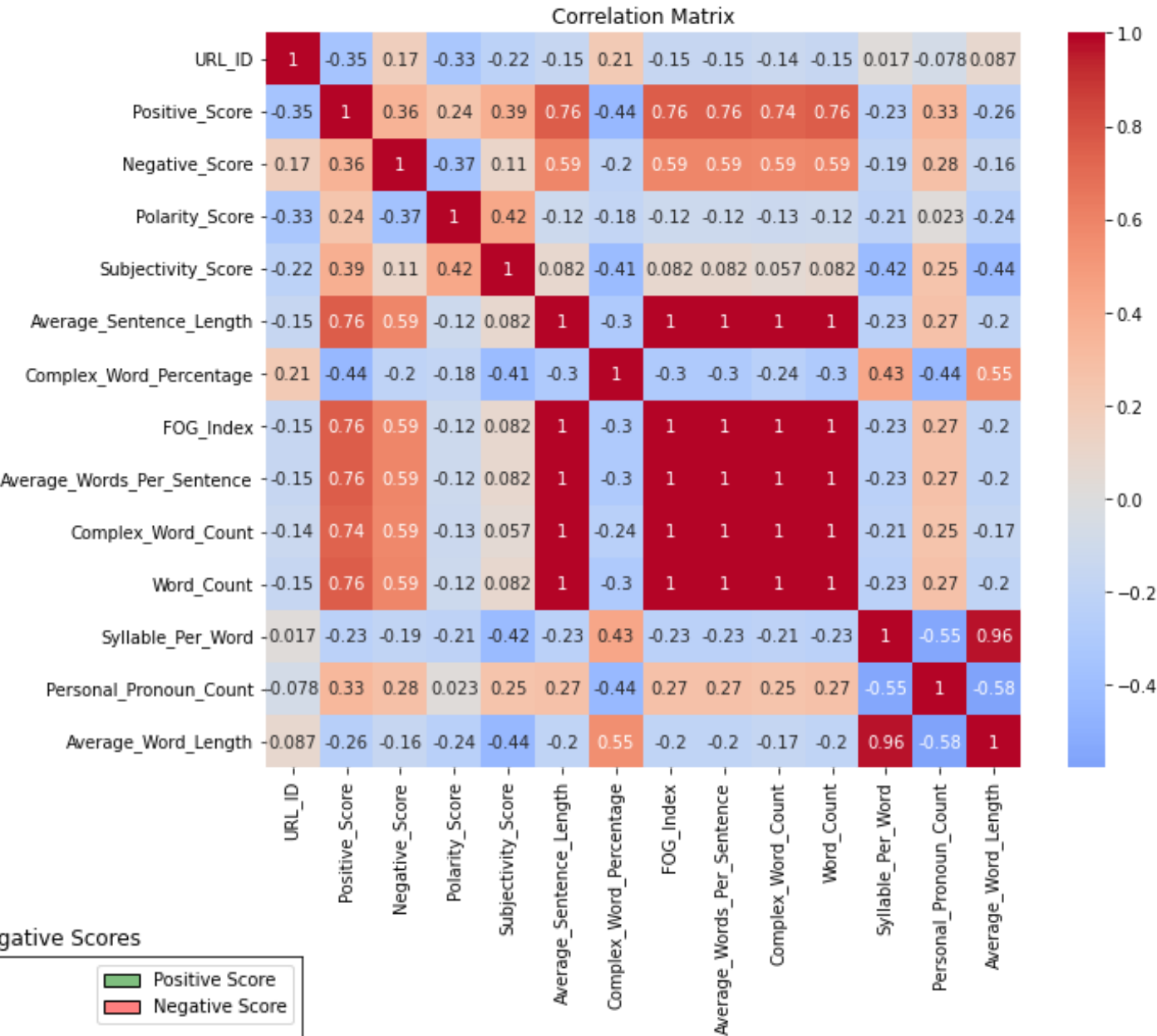
```
1 # Scatter plot of Average_Sentence_Length vs. Complex_Word_Count
2 plt.figure(figsize=(8, 6))
3 plt.scatter(data["Average_Sentence_Length"], data["Complex_Word_Count"], alpha=0.5)
4 plt.xlabel("Average Sentence Length")
5 plt.ylabel("Complex Word Count")
6 plt.title("Average Sentence Length vs. Complex Word Count")
7 plt.show()
8
```





Sentiment Label Counts:

Sentiment_Label	Count
Positive	65
Negative	45
Neutral	1



Thank you

NAME

Vinodhini Rajamanickam

GITHUB

<https://github.com/Vinodhini96>

COURSE AND BATCH

Data Science D50

