

WEB SCRAPED NLP TEXTUAL ANALYSIS

Abstract

This project aims to automate the extraction and analysis of textual data from online articles using Python programming language. Utilising web scraping tools such as BeautifulSoup, systematic retrieval and saving of article content were achieved. Subsequently, text analysis was conducted to compute various linguistic and readability metrics, including polarity, subjectivity, fog index, and word complexity. The final output conformed to a predefined structure for consistent reporting. This comprehensive pipeline provides valuable insights into content sentiment and quality, enabling enhanced content evaluation and categorization. The project demonstrates the effective application of Natural Language Processing (NLP) for automated data extraction, processing, and sentiment analysis from web-based textual sources.

Keywords

Web Scraping, NLP, Sentiment analysis, analysis of readability, Data analysis, Text Analysis, Python programming, Fog Index, Data Extraction, Automated Content Evaluation

Introduction

In today's data-driven world, the abundance of unstructured information, particularly in the form of online articles and blogs, presents both an opportunity and a challenge. While these textual resources are rich in knowledge and insights, extracting meaningful data and conducting comprehensive analysis from them requires sophisticated tools and techniques. This project aims to address this challenge by developing a Python-based solution to automate the extraction and analysis of online article content, thus enabling deeper insights through Natural Language Processing (NLP).

The project simulates a real-world application of text mining by combining web scraping and NLP to evaluate a set of articles based on predefined metrics. The primary objective is to extract article titles and their content from a list of URLs, clean the data, and perform a series of analytical operations that help quantify the content's sentiment, complexity, and readability. These operations are not only important for academic exercises but are also widely used in real-life applications like content moderation, SEO optimisation, media monitoring, and opinion mining.

The solution approach involves multiple components of text processing. Initially, URLs are scraped using Python libraries such as BeautifulSoup and Requests, with care taken to avoid extracting irrelevant parts like website headers or footers. After successful data extraction, textual analysis is performed using custom logic and NLP tools such as NLTK, with key metrics such as positive and negative sentiment scores, polarity, subjectivity, and personal pronoun usage being calculated.

Beyond sentiment analysis, this project delves into linguistic complexity by calculating features such as the FOG Index, average sentence length, and syllable count per word. These indicators reflect how difficult or easy the content is to understand. Such insights are particularly useful in fields like education technology, journalism, and marketing, where understanding content readability can shape

audience engagement strategies.

Furthermore, the implementation is structured to comply with a strict output format as specified in the assignment instructions. The final results are compiled into a structured CSV file that presents all computed variables in a tabular format. The Python code was modularized for reusability and clarity, while also providing the flexibility to be adapted for future extensions such as keyword extraction, named entity recognition, or multilingual analysis.

This project not only strengthened core programming skills but also offered practical exposure to working with real-world data and using open-source tools to tackle common challenges in data extraction and processing. Overall, this assignment exemplifies how the integration of web scraping and NLP can transform unstructured data into valuable, structured information suitable for further analysis and decision-making.

Requirements and Dependencies

1. A computer having a python interpreter or an IDE to execute the python code.
2. Following python libraries should be installed:

- a) numpy
- b) pandas
- c) BeautifulSoup
- d) warnings
- e) nltk
- f) requests
- g) re
- h) chardet
- i) openpyxl

To install any of these libraries use the code: “pip install library_name”.

3. Store the input file given along the assignment with the path: “Question/Input.xlsx”.
4. Form a folder with name “Extracted Data” to store the data extracted from the URL in the dataset.
5. Store the stopwords attached to the assignment with the path: “Question/StopWords/StopWords_file_name”
6. Store the positive and negative words with the path: “Question/MasterDictionary/file_name”.

Tech Contribution

The successful execution of this project depends on the effective use of multiple technologies and Python libraries that contribute to each phase of the workflow, from data extraction to advanced text analysis. Here’s a detailed list of the prominent technologies and their functions:

1. **Python Programming Language:** Python, known for its simplicity, extensive library support, and suitability for data-centric applications, was the foundation of this project. Its modular syntax and robust community-supported packages made it ideal for web scraping, text processing, and data analysis workflows.

2. **BeautifulSoup:** A Python library specifically designed for web scraping, parsed HTML content and extracted article titles and body content from URLs by traversing the HTML structure and excluding extraneous components.
3. **Requests:** The Requests library transmitted HTTP requests to URLs and retrieved corresponding HTML pages, which were then processed by BeautifulSoup.
4. **Regular Expressions (re):** Python's re module was used for text cleaning and preprocessing, removing punctuation, filtering HTML tags, identifying personal pronouns, and extracting syllables for complexity analysis.
5. **NLTK (Natural Language Toolkit):** NLTK was instrumental in linguistic and syntactic analysis, tokenizing words and sentences, removing stopwords, identifying complex words, and calculating metrics like average sentence length, word count, and FOG index, which are indicators of readability.
6. **Pandas:** Pandas was used for data analysis and processing. Pandas was used extensively for data manipulation and organisation. It read the initial dataset (Input.xlsx), stored intermediate and final results, including scores and metrics, in dataframes, and created the final output file in CSV format. NumPy processed numeric arrays and applied formulas for metrics like polarity and subjectivity scores. Predefined stopword files and sentiment lexicons were integrated to filter irrelevant terms and score sentiment, contributing to accurate analysis.

Metrics used

- **Positive Score:** This score is calculated by assigning the value of +1 for each word if found in the Positive Dictionary and then adding up all the values.
- **Negative Score:** This score is calculated by assigning the value of -1 for each word if found in the Negative Dictionary and then adding up all the values. We multiply the score with -1 so that the score is a positive number.
- **Polarity Score:** This is the score that determines if a given text is positive or negative in nature. It is calculated by using the formula:

$$\text{Polarity Score} = (\text{Positive Score} - \text{Negative Score}) / ((\text{Positive Score} + \text{Negative Score}) + 0.000001)$$
Range is from -1 to +1
- **Subjectivity Score:** This is the score that determines if a given text is objective or subjective. It is calculated by using the formula:

$$\text{Subjectivity Score} = (\text{Positive Score} + \text{Negative Score}) / ((\text{Total Words after cleaning}) + 0.000001)$$
Range is from 0 to +1
- **Average Sentence Length** = the number of words / the number of sentences
- **Percentage of Complex words** = the number of complex words / the number of words
- **Fog Index** = $0.4 * (\text{Average Sentence Length} + \text{Percentage of Complex words})$
- **Average Number of Words Per Sentence** = the total number of words / the total number of sentences
- **Complex words** are words in the text that contain more than two syllables.
- **number of Syllables** in each word of the text is the count of vowels present in each word. We also handle some exceptions like words ending with "es", "ed" by not counting them as a syllable.

- To calculate **Personal Pronouns** mentioned in the text, we use regex to find the counts of the words - “I,” “we,” “my,” “ours,” and “us”. Special care is taken so that the country name US is not included in the list.
- **Average Word Length** is calculated by the formula:
Sum of the total number of characters in each word/Total number of words

Organization

1. Requirement Analysis

- Reviewed the assignment objective to extract article content from URLs and compute text-based metrics.
- Understood the variables and output structure from the provided documentation.

2. Environment Setup

- Installed necessary Python libraries (beautifulsoup4, requests, nltk, pandas, etc.).
- Created folder structure for input files, stopwords, dictionary files, and output storage.

3. Article Extraction

- Fetched HTML content using requests.
- Parsed article titles and bodies using BeautifulSoup, excluding irrelevant elements.
- Saved cleaned articles as .txt files named by URL_ID.

4. Text Preprocessing

- Cleaned text using regex and removed stopwords using NLTK and custom files.
- Lowercased text for uniformity and consistency.

5. Sentiment Analysis

- Loaded sentiment word lists.
- Calculated positive, negative, polarity, and subjectivity scores.

6. Readability & Complexity Metrics

- Used nltk for tokenization.
- Computed metrics like FOG index, average sentence length, complex word count, and syllables per word.

7. Additional Text Features

- Counted personal pronouns using regex.
- Calculated word count and average word length.

8. Output Compilation

- Structured results as per Output Data Structure.xlsx.
- Exported final data to OUTPUT.csv.

9. Final Submission

- Packaged Python script, output file, and documentation for submission as per instructions.

Proposed Workflow of the Project

Step 1: Import necessary libraries and the dataset

- a) Import necessary libraries like numpy, pandas, BeautifulSoup, warnings, nltk, requests, re and openpyxl.
- b) Import the input dataset given using pandas library.

Step 2: Data Extraction

- a) Using BeautifulSoup library, construct a function to extract article text from one URL by studying any webpage given in the dataset.
- b) Add the functionality of removing the HTML tags present in the text using regular expressions.
- c) Loop the function created through the URL column and save the text in a text file with 'URL_UD' as its name. Also append it to the dataframe created while importing the dataset.

Step 3: Remove the stopwords of sentiment analysis

- a) Determine the encoding of files given for stopwords.
- b) Construct a function to extract stopwords from the files given and save them in a set. Also, lowercase the words for uniformity.
- c) Construct a function to remove these stopwords from each article and apply it on the column containing article text.

Step 4: Performing sentiment analysis

- a) Determine the encoding of files given for positive and negative words. Construct a function to extract these words from the files given and save them in a set. Also, lowercase the words for uniformity.
- b) Remove URLs present and punctuation marks present in the text without stopwords.
- c) Construct a function for calculating positive and negative score from a text and then pass article text without stopwords to the function to get scores for each row.
- d) After calculating total words after cleaning of stopwords, find the polarity score and subjectivity score with the help of formulas given.

Step 5: Analysis of readability and complex word count:

- a) Using nltk library, tokenize the article text and calculate total number of sentences and words in the article text. Also calculate average sentence length and average number of words per sentence using these.
- b) Create a function to calculate total number of complex words keeping in mind the rules present in assignment guidelines for a complex word by not considering words ending with 'es' or 'ed'.
- c) After passing article text through this function, find value of complex word count along with percentage of complex words using formula given.
- d) Also calculate, Fog Index using these variables.

Step 6: Word Count

- a) Using nltk library, tokenize the text and remove the stopwords present in article text using nltk stopwords package.
- b) Also remove punctuation from the text, and calculate word count for the article text.

Step 6: Syllable Count per word

- a) Using rules given in the assignment, create a function for constructing a list that gives number of syllables for each word in a sentence.
- b) Apply the function through the article text column. Get the list for each article in the dataframe.

Step 7: Count Personal pronouns

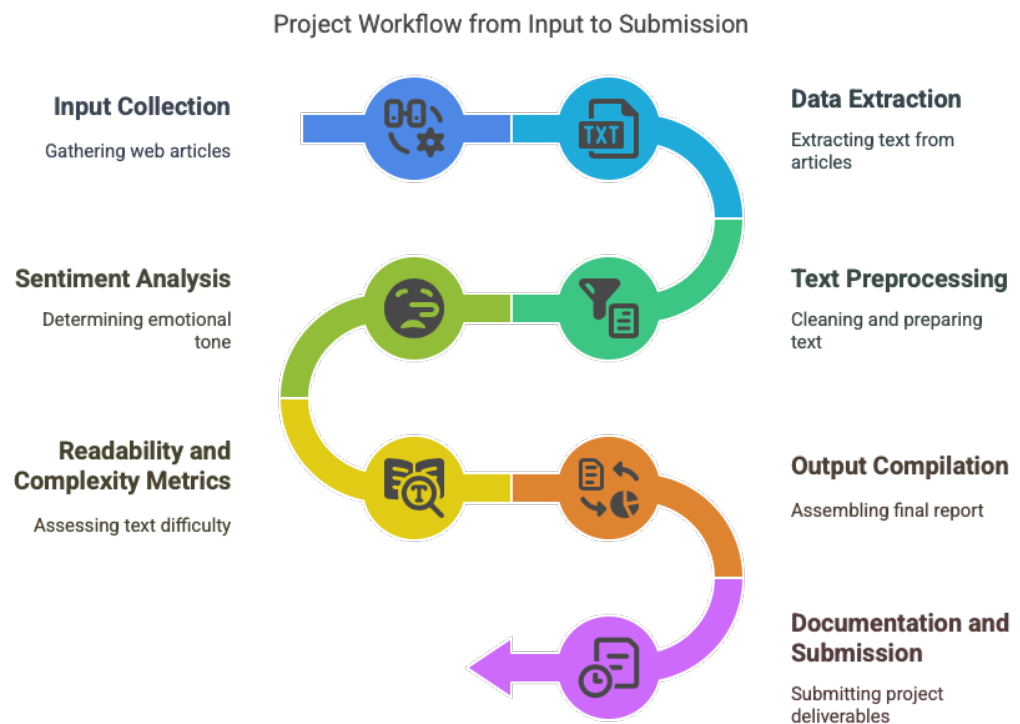
- a) Remove the word 'US' for taking special care so that it doesn't occur in the list. Lowercase the text as well.
- b) Using regular expressions count the number of personal pronouns.

Step 8: Average Word Length

- a) Calculate the total number of characters present in the article after tokenizing the text.

Step 9: Forming dataframe as given in the output data structure.

- a) Rename the required columns as given in the Output data structure.
- b) Create a new dataframe and add all the columns sequentially as given in output data structure.
- c) Save the dataframe created as a csv file as well using pandas library.



Results

The project successfully achieved its goal of extracting and analysing textual data from online articles using Python-based automation. The pipeline was developed to process a list of URLs, extract relevant text, clean and analyse it, and finally compile the results into a structured format.

The primary outcomes of this project are as follows:

1. Accurate Text Extraction

- Only the article title and content were extracted from each URL using BeautifulSoup and requests, effectively removing unwanted web elements such as navigation bars, ads, and footers.
- The extracted data was saved in separate .txt files named by URL_ID, ensuring traceability.

2. Sentiment and Subjectivity Metrics

- Each article was analysed for sentiment using a custom positive and negative word list.
- The project computed sentiment-based metrics including Positive Score, Negative Score, Polarity Score, and Subjectivity Score, enabling a deep understanding of article tone.

3. Readability and Linguistic Features

- The articles were evaluated for complexity using standard readability formulas like the FOG Index.
- Complex Word Count, Average Sentence Length, and Percentage of Complex Words were calculated to assess textual difficulty.

4. Word-Level and Structural Insights

- Word Count, Syllables per Word, Personal Pronoun Count, and Average Word Length were calculated using regex and NLP techniques.
- These metrics added granularity to the analysis and highlighted structural properties of the text.

5. Final Output Compilation

- All analytical metrics, along with the input metadata, were compiled into a clean and organised OUTPUT.csv file.

Resulting CSV file:

AutoSave

OUTPUT

Search (Cmd + Ctrl + U)

HomeInsertDrawPage LayoutFormulasDataReviewViewAutomate

CommentsShare

Paste

Aptos Narrow (Bod... | 12 | A⁺ A⁻ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Conclusion

This project successfully demonstrated the integration of web scraping and Natural Language Processing (NLP) to extract, clean, and analyse online article content. By automating the entire pipeline using Python, the system was able to fetch articles from a given set of URLs, process them efficiently, and generate valuable sentiment and readability metrics.

All predefined output variables were computed and structured into a standardised CSV file, matching the expected format. The implementation showcases a practical application of text analysis using libraries such as BeautifulSoup, NLTK, pandas, and python programming. The metrics derived—such as polarity score, FOG index, word count, and personal pronoun usage—offer meaningful insights into the sentiment and complexity of the articles.

The project helped reinforce foundational concepts in web automation, linguistic analysis, and data wrangling, while also developing problem-solving and software structuring skills. Overall, it provides a solid foundation for scaling to more advanced real-world applications.

- **Future ScopeReal-time Data Processing:**
Extend the system to operate in real time by integrating live RSS feeds or APIs for continuously updated article analysis.
- **Multilingual Support:**
Integrate language detection and translation modules to analyse articles in multiple languages.

- **Advanced Sentiment Models:**
Replace dictionary-based sentiment scoring with deep learning or transformer-based models (e.g., BERT) for more nuanced emotional understanding.
- **Named Entity Recognition (NER):**
Add NER to extract and classify people, places, organisations, and dates mentioned in the articles.
- **Visualisation Dashboard:**
Create a web-based or desktop dashboard using tools like Plotly, Dash, or Streamlit for dynamic visualisation of the results.
- **Content Categorization:**
Implement topic modelling techniques like LDA to group articles by subject matter or theme.
- **SEO Optimisation Analysis:**
Add features to assess content quality based on SEO metrics such as keyword density and readability grade level.