

Approach to the solution

1. Problem Understanding:

- The task involves scraping web articles, analyzing their content, and generating various metrics.
- This requires handling web scraping, text processing, and data analysis.

2. Data Acquisition:

- Implemented web scraping using `requests` and `BeautifulSoup` to extract article content from URLs.
- Input URLs are read from an Excel file.

3. Text Preprocessing:

- Developed functions to clean and normalize text (removing special characters, converting to lowercase).
- Implemented stop word removal using NLTK and custom stop word lists.

4. Analysis Components:

- Sentiment Analysis:
 - Used predefined positive and negative word lists.
 - Calculated positive score, negative score, polarity score, and subjectivity score.
- Readability Analysis:
 - Implemented metrics like average sentence length, percentage of complex words, and Fog Index.
- Text Statistics:
 - Counted syllables, personal pronouns, and calculated average word length.

5. Integration:

- Created a main `analyze_text` function that combines all analysis components.
- This function processes each article and returns a dictionary of all calculated metrics.

6. Data Processing Pipeline:

- Set up a loop to process each URL from the input Excel file.
- For each URL, the program scrapes the content, saves it to a file, analyzes it, and stores the results.

7. Output Generation:

- Results for all articles are compiled into a pandas DataFrame.
- The DataFrame is then exported to an Excel file for easy viewing and further analysis.

8. Error Handling:

- Implemented basic error handling to catch and report issues with processing individual URLs.

Required Dependencies

To run the web scraping and analysis script, you need to install the following Python packages:

1. pandas
2. requests
3. beautifulsoup4
4. nltk
5. openpyxl (for Excel file handling)

Use this command to install the dependencies

```
pip install pandas requests beautifulsoup4 nltk openpyxl
```

To download additional NLTK data run the command

```
import nltk  
nltk.download('punkt')  
nltk.download('stopwords')
```

Make sure you have the input data in the same directory as your script.