

1. My approach to the solution

To address the challenge of analysing web content for sentiment and readability, I adopted a structured approach focusing on systematic data extraction, processing, and analysis. My solution entailed several key steps, leveraging Python's powerful libraries to automate and streamline the process.

Here is how I approached it:

Initial Setup and Data Collection

Library Imports: I began by importing essential libraries. `requests` and `BeautifulSoup` were chosen for fetching and parsing web content. `pandas` was used for data manipulation, while `os`, `re`, and `string` facilitated file operations, regex, and string processing. For natural language processing tasks, such as tokenization, stop word removal, and syllable counting, I utilized `nltk` and `syllapy`.

Text Extraction and Preprocessing

Extracting Text: The function `extract_text_from_url` was designed to request web pages and parse HTML content, extracting text from `<p>` tags. This process was essential for gathering the raw data needed for further analysis.

Saving and Organizing Text: The extracted text was saved into files using `save_text_to_file`, organizing the data for accessibility and processing.

Cleaning Text: I loaded a comprehensive set of stop words from various files within a specified directory, implementing `load_stop_words`. This was crucial for filtering out common but insignificant words, refining the data for analysis.

Sentiment and Readability Analysis

Sentiment Lexicons: I prepared positive and negative word lists to serve as the foundation for sentiment analysis, employing the `load_words` function to load these words into memory.

Calculating Sentiment Scores: Using `calculate_scores`, I quantified the sentiment by counting the occurrences of positive and negative words within the text, determining the overall polarity and subjectivity scores.

Readability Metrics: I utilized `nltk` and `syllapy` for detailed readability analysis, including syllable counting with `count_syllables`, and calculating various metrics like average sentence length and the Fog Index through `calculate_metrics`. These metrics offered insights into the complexity and accessibility of the text.

Aggregation and Reporting

Data Aggregation: After processing each text file, I aggregated the results into a structured format, enabling comprehensive analysis and interpretation.

Exporting Results: I meticulously organized the aggregated data in a pandas DataFrame, then reordered and exported this information to an Excel file, ensuring the results were easily accessible for review or further analysis.

Throughout this process, I emphasized clean, modular code and robust error handling to ensure reliability and scalability. My approach was methodical, starting from the foundational step of data collection and moving through to detailed text analysis, ultimately compiling the findings into a structured report. This process not only facilitated a thorough analysis of web content but also exemplified the power of Python for automating and simplifying complex data processing tasks.

In the provided code, three main directories are utilized to organize the workflow of extracting, processing, and saving text data. The purpose and use of these directories, as mentioned in the code, are as follows:

1. Input Directory (EXDT)

Purpose: This directory serves as the initial storage location for text files that are created from the content extracted from URLs. After fetching the HTML content of each URL listed in the input Excel file, the text content (specifically from paragraph tags) is saved into text files. Each file corresponds to one URL's content.

Use in Code: The script checks if this directory exists at the beginning of execution. If it doesn't exist, it creates the directory using `os.makedirs(input_directory)`.

Extracted text for each URL is saved into this directory by the `save_text_to_file` function, which creates a text file named after the URL's identifier (`url_id`) and writes the extracted text into it.

2. Output Directory (SWdata)

Purpose: This directory is intended for storing the final analysis results, particularly the Excel file named `Output.xlsx` that contains the sentiment and readability metrics calculated for each URL's text content.

Use in Code: Similar to the input directory, the script ensures this directory exists before attempting to save the output Excel file. If it doesn't exist, it's created using `os.makedirs(output_directory)`.

After calculating the sentiment scores and readability metrics for each text file, compiling these into a pandas DataFrame, and arranging the columns in the specified order, the DataFrame is saved as an Excel file (`Output.xlsx`) in this directory.

3. Stop Words Directory (StopWords)

Purpose: This directory contains text files with lists of stopwords. Stopwords are common words (like "the", "is", "at") that are often removed before processing text to ensure the analysis focuses on the more meaningful content.

Use in Code: The script loads stopwords from text files located in this directory at the start of its execution. This is done through the `load_stop_words` function, which reads each file in the directory, adding the words to a set of stopwords.

These loaded stopwords are then used to filter out common words from the extracted text content of each URL, as part of the text cleaning process before calculating sentiment scores and readability metrics.

Each of these directories plays a crucial role in the structured workflow of the script, from managing raw text data extracted from URLs to storing the results of the analysis in an easily accessible format.

2. How to run the main.py file to generate output

To run the Python program and generate the output, we follow the following steps:

Step 1: Environment Setup

Ensure Python is installed on the system along with the required libraries: requests, BeautifulSoup4, pandas, nltk, syllapy, and openpyxl. If these libraries are not installed, they can be installed via pip: **pip install requests BeautifulSoup4 pandas nltk syllapy openpyxl**

Step 2: Prepare Input Data

The program requires an Excel file named Input.xlsx containing URLs to be processed. This file should have at least two columns: URL_ID (a unique identifier for each URL) and URL (the webpage address).

Step 3: Directory Structure

Ensure the existence of three directories as specified in the script: EXDT, SWdata, and StopWords.

EXDT: Stores raw text extracted from URLs.

SWdata: Destination for the final output, Output.xlsx.

StopWords: Contains text files with lists of stopwords to be filtered out during analysis.

If these directories do not exist, create them in the same location as the script.

Step 4: Master Dictionary

The sentiment analysis relies on lists of positive and negative words located in the MasterDictionary directory. Verify the presence of positive-words.txt and negative-words.txt files in this directory.

Step 5: Execute the Script

With the environment set up and input data prepared, execute the script from a terminal or command prompt by navigating to the directory containing the script and running: **python main.py**

Step 6: Verify the Output

Upon successful execution, the SWdata directory should contain the Output.xlsx file. This file will include the analysis results, with each row corresponding to a URL from the input file and columns detailing metrics like sentiment scores and readability metrics.

Troubleshooting: If issues arise during execution, consult the console output for error messages. Common problems may include missing libraries (resolve by installing through pip), incorrect file paths or directory names, or issues with internet connectivity that may affect URL content fetching.

3. Include all dependencies required

The Python program provided requires several external libraries and dependencies to function correctly. Here's a comprehensive list of these dependencies, along with a brief description of their purpose within the context of the program:

1. **requests:** Used for making HTTP requests to fetch web page content from URLs.
2. **BeautifulSoup (bs4):** A library for parsing HTML and XML documents. It's used to extract text from the HTML content fetched from URLs.
3. **pandas:** A powerful data manipulation and analysis library. In this program, it's used for reading the input Excel file containing URLs and saving the analysis results into an Excel file.
4. **os:** A module providing a way of using operating system-dependent functionality. It's used for directory operations, such as checking if a directory exists and creating directories.
5. **re (Regular Expression):** This module is used for searching text for strings matching a regular expression. In the program, it's utilized for identifying personal pronouns within the text as part of readability analysis.

6. **string**: This module contains common string operations, including constants for string processing. The program uses it for removing punctuation from text.
7. **nltk (Natural Language Toolkit)**: A leading platform for building Python programs to work with human language data. It's used for tokenizing the text into sentences and words, removing stopwords, and leveraging the CMU Pronouncing Dictionary for syllable counting.
8. **syllapy**: A library for counting the number of syllables in English words. It complements the CMU Pronouncing Dictionary from nltk for syllable counting, especially for words not covered by the dictionary.
9. **openpyxl**: A library for reading and writing Excel 2010 xlsx/xlsm/xltx/xltn files. It's required by pandas to write DataFrame objects directly to Excel files.

Make sure to have an internet connection to download these packages from the Python Package Index (PyPI). Additionally, some nltk resources need to be downloaded using nltk's download interface, specifically:

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

These commands can be run within the Python script itself or in an interactive Python session to ensure the necessary nltk datasets and tokenizers are available for use.