

1. Install Required Libraries: You'll need to install the following Python libraries:

pandas: To read data from Excel. beautifulsoup4: For web scraping. requests: To send HTTP requests. openpyxl: For reading the Excel file.

```
1 pip install pandas beautifulsoup4 requests openpyxl
2
```

```

➔ Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (2.2.2)
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.10/dist-packages (4.12.3)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (2.32.3)
Requirement already satisfied: openpyxl in /usr/local/lib/python3.10/dist-packages (3.1.5)
Requirement already satisfied: numpy>=1.22.4 in /usr/local/lib/python3.10/dist-packages (from pandas) (1.26.4)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas) (2024.2)
Requirement already satisfied: soupsieve>1.2 in /usr/local/lib/python3.10/dist-packages (from beautifulsoup4) (2.6)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests) (2024.8.30)
Requirement already satisfied: et-xmlfile in /usr/local/lib/python3.10/dist-packages (from openpyxl) (1.1.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)

```

```

1 ## 2. Read Input File (input.xlsx):
2 # Use the pandas library to read the URLs from the Excel file.

```

```

1 import pandas as pd
2
3 # Load the Excel file
4 df = pd.read_excel("/content/Input.xlsx")
5
6 # Replace 'URL_ID' and 'URL' with the actual names of the columns in your Excel file
7 # Based on your global variable `df`, the DataFrame has 2 columns, and we assume the columns have names 'URL_ID' and 'URL'
8 urls = df[['URL_ID', 'URL']].values

```

3. Web Scraping Using BeautifulSoup: For each URL, you'll extract the title and text from the page, ignoring any headers, footers, or other irrelevant content. Here's an example of how to do this using BeautifulSoup:

```

1 import requests
2 from bs4 import BeautifulSoup
3
4 def extract_article_text(url):
5     try:
6         response = requests.get(url)
7         if response.status_code == 200:
8             soup = BeautifulSoup(response.content, 'html.parser')
9
10            # Extract the title
11            title = soup.title.get_text() if soup.title else "No Title"
12

```

```
13     # Extract the main content, modify this based on the site's HTML structure
14     article_body = soup.find('article') # This may vary, adjust based on the actual site
15     if article_body:
16         text = article_body.get_text(separator=' ', strip=True)
17     else:
18         text = "No content found"
19
20     return title, text
21 else:
22     return "Failed to retrieve", "No content"
23 except Exception as e:
24     return "Error occurred", str(e)
25
```

4. Save Extracted Text to a File: After extracting the text, save it to a text file named after the URL\_ID.

```
1 def save_article_text(url_id, title, text):
2     filename = f"{url_id}.txt"
3     with open(filename, 'w', encoding='utf-8') as f:
4         f.write(f"Title: {title}\n\n")
5         f.write(text)
6
```

5. Iterate Through All URLs and Extract Data:

```
1 for url_id, url in urls:
2     title, text = extract_article_text(url)
3     save_article_text(url_id, title, text)
4     print(f"Article {url_id} processed and saved.")
5
```

```
➦ Article Netclan20241017 processed and saved.
Article Netclan20241018 processed and saved.
Article Netclan20241019 processed and saved.
Article Netclan20241020 processed and saved.
Article Netclan20241021 processed and saved.
Article Netclan20241022 processed and saved.
Article Netclan20241023 processed and saved.
Article Netclan20241024 processed and saved.
Article Netclan20241025 processed and saved.
Article Netclan20241026 processed and saved.
Article Netclan20241027 processed and saved.
Article Netclan20241028 processed and saved.
Article Netclan20241029 processed and saved.
Article Netclan20241030 processed and saved.
Article Netclan20241031 processed and saved.
Article Netclan20241032 processed and saved.
Article Netclan20241033 processed and saved.
Article Netclan20241034 processed and saved.
Article Netclan20241035 processed and saved.
Article Netclan20241036 processed and saved.
Article Netclan20241037 processed and saved.
```

```
Article Netclan20241038 processed and saved.
Article Netclan20241039 processed and saved.
Article Netclan20241040 processed and saved.
Article Netclan20241041 processed and saved.
Article Netclan20241042 processed and saved.
Article Netclan20241043 processed and saved.
Article Netclan20241044 processed and saved.
Article Netclan20241045 processed and saved.
Article Netclan20241046 processed and saved.
Article Netclan20241047 processed and saved.
Article Netclan20241048 processed and saved.
Article Netclan20241049 processed and saved.
Article Netclan20241050 processed and saved.
Article Netclan20241051 processed and saved.
Article Netclan20241052 processed and saved.
Article Netclan20241053 processed and saved.
Article Netclan20241054 processed and saved.
Article Netclan20241055 processed and saved.
Article Netclan20241056 processed and saved.
Article Netclan20241057 processed and saved.
Article Netclan20241058 processed and saved.
Article Netclan20241059 processed and saved.
Article Netclan20241060 processed and saved.
Article Netclan20241061 processed and saved.
Article Netclan20241062 processed and saved.
Article Netclan20241063 processed and saved.
Article Netclan20241064 processed and saved.
Article Netclan20241065 processed and saved.
Article Netclan20241066 processed and saved.
Article Netclan20241067 processed and saved.
Article Netclan20241068 processed and saved.
Article Netclan20241069 processed and saved.
Article Netclan20241070 processed and saved.
Article Netclan20241071 processed and saved.
Article Netclan20241072 processed and saved.
Article Netclan20241073 processed and saved.
```

```
1 !pip install textblob
2 from textblob import TextBlob
3 import pandas as pd
4
5 # Define necessary functions
6 def positive_score(text):
7     """Calculates the positive score of a text using TextBlob."""
8     analysis = TextBlob(text)
9     return analysis.sentiment.polarity # Assuming polarity represents positive sentiment
10
11 def negative_score(text):
12     """Calculates the negative score of a text using TextBlob."""
13     analysis = TextBlob(text)
14     return analysis.sentiment.subjectivity - analysis.sentiment.polarity # Example calculation
15
16 def polarity_score(text):
17     # Add your implementation
18     pass
19
20 def subjectivity_score(text):
```

```
21     # Add your implementation
22     pass
23
24 def avg_sentence_length(text):
25     # Add your implementation
26     pass
27
28 def percentage_complex_words(text):
29     # Add your implementation
30     pass
31
32 def fog_index(text):
33     # Add your implementation
34     pass
35
36 def avg_words_per_sentence(text):
37     # Add your implementation
38     pass
39
40 def complex_word_count(text):
41     # Add your implementation
42     pass
43
44 def word_count(text):
45     # Add your implementation
46     pass
47
48 def syllables_per_word(text):
49     # Add your implementation
50     pass
51
52 def personal_pronouns_count(text):
53     # Add your implementation
54     pass
55
56 def avg_word_length(text):
57     # Add your implementation
58     pass
59
60 output_data = []
61
62 for url_id in df['URL_ID']:
63     with open(f'{url_id}.txt', 'r', encoding='utf-8') as f:
64         text = f.read()
65
66     analysis = {
67         'URL_ID': url_id,
68         'POSITIVE SCORE': positive_score(text),
69         'NEGATIVE SCORE': negative_score(text),
70         'POLARITY SCORE': polarity_score(text),
71         'SUBJECTIVITY SCORE': subjectivity_score(text),
72         'AVG SENTENCE LENGTH': avg_sentence_length(text),
73         'PERCENTAGE OF COMPLEX WORDS': percentage_complex_words(text),
74         'FOG INDEX': fog_index(text),
75         'AVG NUMBER OF WORDS PER SENTENCE': avg_words_per_sentence(text),
```

```
76     'COMPLEX WORD COUNT': complex_word_count(text),
77     'WORD COUNT': word_count(text),
78     'SYLLABLE PER WORD': syllables_per_word(text),
79     'PERSONAL PRONOUNS': personal_pronouns_count(text),
80     'AVG WORD LENGTH': avg_word_length(text)
81 }
82
83     output_data.append(analysis)
84
85 # Save to Excel
86 output_df = pd.DataFrame(output_data)
87 output_df.to_excel("Output Data Structure.xlsx", index=False)
```

➞ Requirement already satisfied: textblob in /usr/local/lib/python3.10/dist-packages (0.17.1)  
Requirement already satisfied: nltk>=3.1 in /usr/local/lib/python3.10/dist-packages (from textblob) (3.8.1)  
Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from nltk>=3.1->textblob) (8.1.7)  
Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from nltk>=3.1->textblob) (1.4.2)  
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.10/dist-packages (from nltk>=3.1->textblob) (2024.9.11)  
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk>=3.1->textblob) (4.66.5)