# Required Dependencies

1. **NLTK**: For tokenizing, stopword removal, and text processing.
2. **Pandas**: For handling input and output data in tabular format.
3. **OpenPyXL**: For reading and writing Excel files.
4. **BeautifulSoup4**: For web scraping and extracting text from HTML.
5. **Requests**: For fetching webpages.
6. **Textstat**: For calculating readability metrics.

---

## Install Dependencies

Run the following command in your terminal to install all dependencies:

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

---

## Ensure Proper Setup in the Script

Add these imports at the beginning of your script:

```
import os
import re
import nltk
import pandas as pd
from nltk.tokenize import word_tokenize, sent_tokenize
from nltk.corpus import stopwords
from bs4 import BeautifulSoup
import requests
from textstat import textstat
```

## Download NLTK Resources

Ensure required NLTK resources are downloaded in your script:

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

---

## File Structure for Execution

Make sure your project directory contains the following files:

```
/project-directory
  ├── text_analysis.py        # Your Python script
  ├── input.xlsx              # Input file with URL ID and URL
columns
  ├── output.xlsx             # Output file (generated after script
execution)
  ├── MasterDictionary
  │   ├── positive-words.txt  # Positive word dictionary
  │   └── negative-words.txt  # Negative word dictionary
```

---

## Running the Script

**Navigate to the Script Directory**:

```
cd /path/to/project-directory
```

**Run the Script**:
```
python text_analysis.py
```

**Check Output**: The output will be saved in `output.xlsx` with columns ordered as **URL ID**, **URL**, and computed metrics.