**web_scrape_clean_html_test.py**

```python
1  #
2  #  Course: COSC 4P02
3  #  Assignment: Group Project
4  #  Group: 9
5  #  Version: 1.0
6  #  Date: March 2024
7  #
8  import pytest # Import the pytest module for automated testing.
9  import requests # Import the requests module to make HTTP requests.
10 from langdetect import detect # Import the detect method from the langdetect module.
11 from bs4 import BeautifulSoup # Import the BeautifulSoup class from the bs4 module.
12 from web_scrape_clean_html import get_clean_content, clean_text, scrape_and_clean # Import
   the get_clean_content, clean_text, and scrape_and_clean methods from the custom
   web_scrape_clean_html.
13
14 def test_clean_text():
15     #
16     # Test Case 1: Testing the clean_text method of the web scraper independently. This will
   confirm base functionality of the clean_text method.
17     # Execution: python -m pytest web_scrape_clean_html_test.py -k "test_clean_text" -s -v #
   Only use -s to view the contents of the text in the test.
18     # The method will replace newline and carriage return characters with spaces, replace
   quotes with pipe characters, and reduce multiple spaces to a single space.
19     # Expected Result: Pass. The cleaned text should match the expected output.
20     #
21     input_text = "Test line 1.\nTest            line 2.\rTest line              3.
   \"quotes\"." # Sample text to be cleaned.
22     expected_output = "Test line 1. Test line 2. Test line 3. |quotes|." # Expected output of
   the cleaned text.
23     cleaned_text = clean_text(input_text)  # Call the clean_text method on the URL to be
   tested.
24     assert isinstance(cleaned_text, str)  # Confirm that the cleaned text is a string.
25     assert cleaned_text == expected_output  # Confirm that the cleaned text matches the
   expected output.
26     print(input_text) # Print the input text. This line will only show if the -s is included
   when running the test.
27     print(cleaned_text)  # Print the cleaned text. This line will only show if the -s is
   included when running the test.
28
29 def test_scrape_and_clean():
30     #
31     #  Test Case 2: Testing the scrape_and_clean method with a valid URL independently. For
   this test, not a Wikipedia page, but a CNN page.
32     #  Execution: python -m pytest web_scrape_clean_html_test.py -k "test_scrape_and_clean" -
   s -v # Only use -s to view the contents of the URL in the test.
33     #  This method will test the scrape_and_clean function on a valid URL. The function
   should scrape the page, clean the content,
34     #  and return the page content.
35     #  Expected Result: Pass. The content scraped from the URL is a valid string with only
   permitted characters.
```

```python
36       #
37       url = "https://www.cnn.com/2025/03/10/business/usmca-tariff-delay-trump/index.html" # URL
   to be tested. This is a valid URL.
38       content = scrape_and_clean(url) # Call the scrape_and_clean method on the URL to be
   tested.
39       assert isinstance(content, str) # Confirm that the content scraped from the URL is a
   string.
40       assert len(content) > 50  # Confirm that the content scraped from the URL is more than 50
   characters long. This ensures that the full content of the URL is not empty.
41       assert '\n' not in content # Confirm that the content scraped from the URL does not
   contain newline characters.
42       assert '\r' not in content # Confirm that the content scraped from the URL does not
   contain carriage return characters.
43       assert '"' not in content #Confirm that the content scraped from the URL does not contain
   quotes.
44       print(content) # Print the content scraped from the URL. This line will only show if the
   -s is included when running the test.
45
46   def test_ValidURL1():
47       #
48       #  Test Case 3: Testing the get_clean_content method with a valid URL. For this test,
   Wikipedia.
49       #  Execution: python -m pytest web_scrape_clean_html_test.py -k "test_ValidURL1" -s -v #
   Only use -s to view the contents of the URL in the test.
50       #  This method will test the get_clean_content method of the web scraper with a valid
   Wikipedia URL. The URL will be scraped and
51       #  the content will be cleaned. The content will be checked to ensure that it is a string
   and that it is not empty.
52       #  Expected Result: Pass. The content scraped from the URL is a valid string and is
   significantly long.
53       #
54       #
55       url = "https://en.wikipedia.org/wiki/Agile_software_development" # URL to be tested. This
   is a valid URL.
56       content = get_clean_content(url) # Call the get_clean_content method on the URL to be
   tested.
57       assert isinstance(content, str) # Confirm that the content scraped from the URL is a
   string.
58       assert len(content) > 50  # Confirm that the content scraped from the URL is more than 50
   characters long. This ensures that the full content of the URL is not empty.
59       print(content) # Print the content scraped from the URL. This line will only show if the
   -s is included when running the test.
60
61   def test_ValidURL2():
62       #
63       #  Test Case 4: Testing the get_clean_content method with a valid URL. For this test, CBS
   Sports.
64       #  Execution: python -m pytest web_scrape_clean_html_test.py -k "test_ValidURL2" -s -v #
   Only use -s to view the contents of the URL in the test.
65       #  This method will test the get_clean_content method of the web scraper with a valid CBS
   Sports URL. The URL will be scraped and
66       #  the content will be cleaned. The content will be checked to ensure that it is a string
   and that it is not empty.
```

```python
67        #  Expected Result: Pass. The content scraped from the URL is a valid string and is
      significantly long.
68        #
69        #
70        url = "https://www.cbssports.com/soccer/news/champions-league-bold-predictions-liverpool-
      will-rely-on-alisson-to-save-them-pedri-shines-for-barcelona/" # URL to be tested. This is a
      valid URL.
71        content = get_clean_content(url) # Call the get_clean_content method on the URL to be
      tested.
72        assert isinstance(content, str) # Confirm that the content scraped from the URL is a
      string.
73        assert len(content) > 50  # Confirm that the content scraped from the URL is more than 50
      characters long. This ensures that the full content of the URL is not empty.
74        print(content) # Print the content scraped from the URL. This line will only show if the
      -s is included when running the test.
75
76  def test_InvalidURL():
77        #
78        #  Test Case 5: Testing the get_clean_content method with an invalid URL. For this test,
      www.cosc4p02group9fakeurl.com.
79        #  Execution: python -m pytest web_scrape_clean_html_test.py -k "test_InvalidURL" -s -v #
      Only use -s to view the printed error message of the test.
80        #  This method will test the get_clean_content method of the web scraper with a custom
      invalid URL. The web scraper returns "Error:" if the URL is invalid,
81        #  so the test will check for this string in the content. If the string is found, the
      test will pass, the URL was invalid.
82        #  Expected Result: Pass. "Error:"" will exist in the content because the URL is invalid.
83        #
84        url = "http://www.cosc4p02group9fakeurl.com/" # URL to be tested. This is an invalid URL.
85        content = get_clean_content(url) # Call the get_clean_content method on the URL to be
      tested.
86        assert "Error:" in content # Confirm that the content scraped from the URL contains the
      string "Error:". This indicates that the scraper found no URL, hence it is invalid.
87        if "Error:" in content: # If the content contains the string "Error:", the URL is
      invalid.
88            print(f"URL '{url}' is not valid.") # Print that the URL is not valid.
89
90  def test_NoMainContent():
91        #
92        # Test Case 6: Testing the get_clean_content method with a valid URL but no main content.
      For this test, https://www.example.com/.
93        # Execution: python -m pytest web_scrape_clean_html_test.py -k "test_NoMainContent" -s -v
      # Only use -s to view the printed message of the test.
94        # This method will test the get_clean_content method on a valid URL that doesn't have any
      main content. The page may be blank, or there may be no main content. The web scraper returns
95        # "Main content not found." if the main content is not found, so the test will check for
      this string in the content. If the string is found, the test will pass, the URL had no main
      content.
96        # Expected Result: Pass. "Main content not found." will exist in the content because the
      URL has no main content.
97        #
```

```python
 98     url = "https://www.example.com/"  # URL to be tested. This is a valid URL without main
        content.
 99     content = get_clean_content(url)  # Call the get_clean_content method on the URL to be
        tested.
100     assert isinstance(content, str) # Confirm that the content scraped from the URL is a
        string.
101     assert "Main content not found." in content  # Confirm that the content scraped from the
        URL contains the string "Main content not found." This indicates that the scraper found no
        main content.
102     if "Main content not found." in content:
103         print(f"URL '{url}' has no main content.") # Print that the URL has no main content.
104
105 def test_NonEnglishContent():
106     #
107     # Test Case 7: Testing the get_clean_content method with a valid URL that has content not
        in English.
108     # Execution: python -m pytest web_scrape_clean_html_test.py -k "test_NonEnglishContent" -
        s -v # Only use -s to view the English contents of the URL in the test.
109     # This method will test the get_clean_content method on a page with content in a language
        other than English.
110     # The web scraper should not include non-English lines and return only English text if
        present.
111     # Expected Result: Pass. The content should not contain any lines that contain other
        languages.
112     #
113     url = "https://en.wikipedia.org/wiki/French_Wikipedia"  # URL to be tested. This is a
        valid URL, with a French word in the content.
114     content = get_clean_content(url)  # Call the get_clean_content method on the URL to be
        tested.
115     assert isinstance(content, str) # Confirm that the content scraped from the URL is a
        string.
116     assert "Wikipédia en français" not in content  # Confirm that the French phrase
        "Wikipédia en français" should be filtered out. It should not be in content.
117     assert "This edition was started on 23 March 2001, two months after the official creation
        of Wikipedia." in content # Confirm the English lines are still present in content. This is
        the first full English line of the page.
118     print(content) # Print the content scraped from the URL. This line will only show if the
        -s is included when running the test.
119
120 def test_TagsInContent():
121     #
122     # Test Case 8: Testing the get_clean_content method with a valid URL that contains HTML
        tags embedded in the content. Citations, Footnotes, etc.
123     # Execution: python -m pytest web_scrape_clean_html_test.py -k "test_TagsInContent" -s -v
        # Only use -s to view the contents of the URL in the test.
124     # This method will test the get_clean_content method on a page with HTML tags within the
        content. The web scraper should remove the HTML tags and return only readable text.
125     # Most Wikipedia articles contain many HTML tags, so this test could be performed on
        other Wikipedia articles to showcase the removal of embedded tags. Other sites also have
        tags, but Wikipedia is a good example.
126     # Expected Result: Pass. The content should only contain clean text.
127     #
```

```python
128        url = "https://en.wikipedia.org/wiki/Python_(programming_language)"  # URL to be tested.
    This is a valid URL, with many tags and styles.
129        content = get_clean_content(url)  # Call the get_clean_content method on the URL to be
    tested.
130        assert isinstance(content, str) # Confirm that the content scraped from the URL is a
    string.
131        assert "<" not in content  # Confirm there are no tags in the content. The content should
    not contain any HTML tags. 1/2
132        assert ">" not in content  # Confirm there are no tags in the content. The content should
    not contain any HTML tags. 2/2
133        assert "Python" in content  # Confirm there are English words remaining in the content.
    The content should contain the word "Python".
134        print(content) # Print the content scraped from the URL. This line will only show if the
    -s is included when running the test.
135
136
```