

```

1 # Appendix C13 - test_content_importer.py
2
3 from klassify.src.tables import Document
4 from klassify.src.content_importer import ContentImporter
5 import os
6 import pytest
7
8 database_name = "test_klassify"2
9 if os.path.exists("%s.db" % database_name):
10     os.remove("%s.db" % database_name)
11
12 DOCUMENT = Document(
13     base_path = "/intelligent-machines",
14     title = "The Intelligent Machines",
15     html = open("test/fixtures/document_page.html", 'r').read())
16 STRING_PRESENT_IN_BOTH_HEADER_AND_FOOTER = "How government works"
17 STRING_PRESENT_IN_SCRIPT_TAG = "<![CDATA["
18 STRING_PRESENT_IN_TITLE = "HM Revenue & Customs"
19
20 def setup_module(module):
21     global IMPORTER
22     IMPORTER = ContentImporter(db_name="test_klassify")
23     IMPORTER.DBH.session.add(DOCUMENT)
24     IMPORTER.DBH.session.commit()
25 def teardown_module(module):
26     IMPORTER.DBH.session.close()
27     IMPORTER.DBH.destroy_db_if_present()
28
29 def test_cleaning_methods():
30     doc = IMPORTER.DBH.session.query(Document).first()
31     page = IMPORTER.parse_page(doc.html)
32
33     assert STRING_PRESENT_IN_BOTH_HEADER_AND_FOOTER in page.text
34     assert STRING_PRESENT_IN_SCRIPT_TAG in page.text
35     page = IMPORTER.remove_unwanted_tags(page)
36     assert STRING_PRESENT_IN_BOTH_HEADER_AND_FOOTER not in page.text
37     assert STRING_PRESENT_IN_SCRIPT_TAG not in page.text
38
39     assert STRING_PRESENT_IN_TITLE in page.text
40     page = IMPORTER.get_body(page)
41     assert STRING_PRESENT_IN_TITLE not in page.text
42
43     page_content = IMPORTER.extract_page_content(page)
44     page_content = IMPORTER.remove_non_relevant_content(page_content)
45     for phrase in IMPORTER.NON_RELEVANT_PHRASES:
46         assert phrase not in page_content
47
48     assert "2016" in page_content
49     page_content = IMPORTER.remove_punctuation_and_numbers(page_content)
50     assert "2016" not in page_content
51
52 def test_extract_content_single_method():
53     doc = IMPORTER.DBH.session.query(Document).first()
54
55     assert STRING_PRESENT_IN_BOTH_HEADER_AND_FOOTER in doc.html
56     assert STRING_PRESENT_IN_SCRIPT_TAG in doc.html
57
58     clean_content = IMPORTER.extract_content(doc)
59
60     assert STRING_PRESENT_IN_BOTH_HEADER_AND_FOOTER not in clean_content
61     assert STRING_PRESENT_IN_SCRIPT_TAG not in clean_content
62     for phrase in IMPORTER.NON_RELEVANT_PHRASES:
63         assert phrase not in clean_content
64

```