

```

1 # Appendix C16 - test_feature_extractor.py
2
3 from klassify.src.feature_extractor import FeatureExtractor
4 from klassify.src.tables import Document
5
6 initial_document_1 = Document(title="Test title 1",
7                               base_path="/test-1",
8                               content="This is a test document - one")
9 initial_document_2 = Document(title="Test title 2",
10                               base_path="/test-2",
11                               content="This is a test document - two")
12 initial_document_3 = Document(title="Test title 3",
13                               base_path="/test-3",
14                               content="This is a test document - three")
15
16 EXTRACTOR = FeatureExtractor([
17     initial_document_1,
18     initial_document_2,
19     initial_document_3,
20 ])
21
22 new_document = Document(title="Self assessment deadlines 3",
23                          base_path="/self-assessment-3",
24                          html="<strong>PAY NOW 3</strong>",
25                          content="This has a different content - four")
26
27 def test_tokenize():
28     tokenized_content = EXTRACTOR.tokenize(initial_document_1)
29     assert tokenized_content == ['This', 'is', 'a', 'test', 'document', "-", 'one']
30
31 def test_make_vocabulary():
32     # without document
33     assert EXTRACTOR.make_vocabulary() == ['test', 'one', 'test', 'two', 'test', 'three']
34     # with document
35     assert EXTRACTOR.make_vocabulary(new_document) == ['differ', 'content', 'four']
36
37 def test_bag_of_words():
38     # This is built against the vocabulary.
39     # The vocabulary is the sum of all the different terms in all the documents provided
40     # at instantiation.
41     assert EXTRACTOR.bag_of_words(initial_document_3) == {'one': False, 'test': True, 'three': True, 'two': False}
42     assert EXTRACTOR.bag_of_words(new_document) == {'one': False, 'test': False, 'three': False, 'two': False}
43
44 def test_process():
45     # What is bein discarded: Single letter words, Stop words, Long words
46     # Additionally, remaining words will be stemmed.
47     document_with_unfiltered_content = Document(title="Test", base_path="/test",
48         content=" within https .mb , a b c reallylongwordthatshouldbefilteredout cloudy regular words should be stemmed in this process")
49
50     tokenized_content = EXTRACTOR.tokenize(document_with_unfiltered_content)
51
52     assert EXTRACTOR.process(tokenized_content) == ['cloudi', 'regular', 'word', 'stem', 'process']
53

```