

```

1 # Appendix C5 - feature_extractor.py
2
3 from nltk.tokenize import word_tokenize
4 from nltk.corpus import stopwords
5 from nltk.stem import PorterStemmer
6 import nltk
7
8 class FeatureExtractor():
9     def __init__(self, documents, n_features=5000):
10         self.documents = documents
11         self.stemmer = PorterStemmer()
12         self.vocabulary = self.top_words(n_features, self.freq_dist(self.make_vocabulary(
13 )))
14
15     def tokenize(self, document=None):
16         if document:
17             documents = [document]
18         else:
19             documents = self.documents
20
21         return [token for doc in documents for token in word_tokenize(doc.content)]
22
23     def process(self, vocabulary):
24         ADDITIONAL_STOP_WORDS = {'january', 'please', 'https', 'email', 'detail', 'email',
25 , 'send', 'if', 'december', 'october', 'kb', 'february', 'within', 'november', 'may', 'plea
26 se', '.mb', 'what', 'pdf', 'june', 'mach', 'good', 'august', 'september', 'html', 'july', '
27 beta', 'document', 'eg', 'published', 'april'}
28         stop_words = set(stopwords.words("english"))
29
30         processed_words = []
31         for word in vocabulary:
32             # select only words shorter than 20 char
33             if len(word) < 20:
34                 word = word.lower()
35                 # do not select stopwords
36                 if word not in (stop_words | ADDITIONAL_STOP_WORDS):
37                     # stem words
38                     word = self.stemmer.stem(word)
39                     # do not select words shorter than 2 characters
40                     if word.isalpha():
41                         if len(word) > 1:
42                             processed_words.append(word)
43                     else:
44                         processed_words.append(word)
45         return processed_words
46
47     def make_vocabulary(self, document=None):
48         if document:
49             vocabulary = self.tokenize(document)
50         else:
51             vocabulary = self.tokenize()
52
53         vocabulary = self.process(vocabulary)
54         return vocabulary
55
56     def bag_of_words(self, document):
57         doc_words = set(self.make_vocabulary(document))
58         bag_of_words = {}
59
60         for word in self.vocabulary:
61             bag_of_words[word] = (word in doc_words)
62
63         return bag_of_words
64
65     def freq_dist(self, vocabulary):
66         return nltk.FreqDist(vocabulary)
67
68

```

```
64 def top_words(self, n_features, freq_dist):
65     return list(freq_dist.keys())[:n_features]
66
```