

```

1 # Appendix C6 - doc_operator.py
2
3 from .db_handler import DBHandler
4 from .tables import Topic, Subtopic, Document
5 from .feature_extractor import FeatureExtractor
6 import random
7
8 class DocumentOperator():
9     def __init__(self, db_name="klassify", n=3, min_docs=None, max_docs=None, n_features=N
one):
10         self.DBH = DBHandler(db_name=db_name, echo=False)
11         self.topics = self.pick_random_topics(n, min_docs)
12         self.max_docs = max_docs
13         self.topic_labels = [topic.title for topic in self.topics]
14         self.docs_with_labels = self.docs_with_labels()
15         self.featuresets = []
16         self.processor = FeatureExtractor([doc for doc, cat in self.docs_with_labels], n_
features)
17
18     def pick_random_topics(self, n, min_docs):
19         topics = self.DBH.session.query(Topic).all()
20         if min_docs:
21             topics = [topic for topic in topics if len(topic.documents()) > min_docs]
22             random.shuffle(topics)
23             topics = topics[:n]
24         return topics
25
26     def find_random_doc_by_title(self, title):
27         topic = self.DBH.session.query(Topic).filter(Topic.title == title).first()
28         subtopic = random.choice(topic.subtopics)
29         return random.choice(subtopic.documents)
30
31     def random_document(self):
32         all_topics = self.DBH.session.query(Topic).all()
33         topic = random.choice(all_topics)
34         subtopic = random.choice(topic.subtopics)
35         doc = random.choice(subtopic.documents)
36         bag_of_words = self.baggify_document(doc)
37         return doc, bag_of_words
38
39     def docs_with_labels(self):
40         docs_with_filtered_labels = []
41
42         for topic in self.topics:
43             docs_with_labels = topic.documents_with_labels()
44
45             if self.max_docs:
46                 random.shuffle(docs_with_labels)
47                 docs_with_labels = docs_with_labels[:self.max_docs]
48
49             for doc, doc_labels in docs_with_labels:
50                 filtered_labels = []
51                 for label in doc_labels:
52                     # Filter out labels that are not the selected topics
53                     if label in self.topic_labels:
54                         filtered_labels.append(label)
55                 docs_with_filtered_labels.append([doc, filtered_labels])
56
57         return docs_with_filtered_labels
58
59     def build_feature_sets(self):
60         document_set_with_category = self.docs_with_labels
61         random.shuffle(document_set_with_category)
62
63         count = 0
64         for (document, category) in document_set_with_category:
65             count = count + 1

```

```
66         self.featuresets.append([self.baggify_document(document), category])
67
68     def baggify_document(self, doc):
69         return self.processor.bag_of_words(doc)
70
```