```python
1  # Appendix C3 - doc_importer.py
2
3  import math
4  from .tables import Subtopic, Document
5  from .db_handler import DBHandler
6  import requests
7  import sqlalchemy
8  import time
9
10 class DocumentImporter(object):
11     def __init__(self, db_name="klassify"):
12         self.ROOT_URL = "https://www.gov.uk/api/search.json?reject_specialist_sectors=_MI
SSING"
13         self.PAGE_URL = "https://www.gov.uk/api/search.json?reject_specialist_sectors=_MI
SSING&count=1000&start="
14         self.DBH = DBHandler(db_name, echo=False)
15
16     def api_response(self, url):
17         time.sleep(0.15)
18         return requests.get(url).json()
19
20     def total_documents(self, document_data):
21         self.document_count = document_data["total"]
22         return self.document_count
23
24     def pages(self, number_of_documents):
25         return math.ceil(number_of_documents / 1000)
26
27     def urls(self, number_of_pages):
28         urls = []
29         for i in range(number_of_pages):
30             item_count = i * 1000
31             url_with_pagination = self.PAGE_URL + str(item_count)
32             urls.append(url_with_pagination)
33         return urls
34
35     def associate_document_with_subtopics(self, document, subtopics):
36         # remove duplicates by converting topics to a set and then back to a list
37         subtopics = set(subtopics)
38         subtopics = list(subtopics)
39         document.subtopics = subtopics
40
41         return document
42
43     def make_document(self, document_data):
44         link = document_data["link"]
45         title = document_data["title"]
46         description = document_data["description"]
47         doc = Document(
48             web_url="https://www.gov.uk" + link,
49             description=description,
50             base_path=link,
51             title=title
52         )
53
54         return doc
55
56     def find_subtopics(self, document_data):
57         subtopics_data = document_data["specialist_sectors"]
58
59         subtopics = []
60         for subtopic_data in subtopics_data:
61             subtopic = self.DBH.session.query(Subtopic).filter_by(base_path=subtopic_data
['link']).first()
62             if subtopic: subtopics.append(subtopic)
63
64         return subtopics
```

```python
65
66      def run(self):
67          root_data = self.api_response(self.ROOT_URL)
68          number_of_documents = self.total_documents(root_data)
69          pages = self.pages(number_of_documents)
70          urls = self.urls(pages)
71
72          count = 0
73          duplicate_documents = []
74
75          for url in urls:
76              list_of_documents =  self.api_response(url)
77              documents_data = list_of_documents['results']
78              for document_data in documents_data:
79                  document = self.make_document(document_data)
80                  subtopics = self.find_subtopics(document_data)
81                  if subtopics:
82                      self.associate_document_with_subtopics(document, subtopics)
83                  try:
84                      self.DBH.session.add(document)
85                      self.DBH.session.commit()
86                  except sqlalchemy.exc.IntegrityError:
87                      duplicate_documents.append(document.base_path)
88                      self.DBH.session.rollback()
89                  except:
90                      self.DBH.session.rollback()
91                      raise
92                  if count % 250 == 0: print("Documents processed: %d/%d" % (count, self.do
cument_count))
93                  count = count + 1
94
95          self.DBH.session.close()
96
97          print("Documents with duplicates that have been ignored: %d" % len(duplicate_docu
ments))
98
```