

```

1 # Appendix C4 - content_importer.py
2
3 from .db_handler import DBHandler
4 from .tables import Document
5 from bs4 import BeautifulSoup
6 import requests
7 import time
8
9 # Future implementation: Tuning features by adding Documents' to their content. Maybe with
  a multiplier.
10 class ContentImporter(object):
11     def __init__(self, db_name="klassify"):
12         self.DBH = DBHandler(db_name, echo=False)
13         self.ROOT_URL = "https://www.gov.uk"
14         self.NON_RELEVANT_PHRASES = [
15             "Skip to main content",
16             "Find out more about cookies"
17             "GOV.UK uses cookies to make the site simpler",
18             "Is there anything wrong with this page",
19             "Last updated",
20             "Other ways to apply",
21             "Before you start",
22             "Elsewhere on the web",
23             "Find out about call charges",
24             "find out more about beta services",
25             "Return to top ↑",
26             "Find out more about cookies",
27             "GOV.UK",
28             "Don't include personal or financial information",
29             "Help us improve",
30             "This file may not be suitable for users of assistive technology"
31             "If you use assistive technology and need a version of this document in a mor
e accessible format",
32             "tell us what format you need It will help us if you say what assistive techn
ology you use",
33             "Request a different format",
34             "What you were doing",
35             "What went wrong",
36             "uses cookies to make the site simpler."
37         ]
38
39     def parse_page(self, page):
40         soup = BeautifulSoup(page, 'html.parser')
41         return soup
42
43     def extract_page_content(self, page):
44         return page.text
45
46     # Iterate through each Document in the database, get their URL on the site and
47     # query it to obtain their HTML and eventually store it.
48     def import_documents_html(self):
49         documents = self.DBH.session.query(Document).all()
50
51         count = 0
52         for doc in documents:
53             if doc.html == None:
54                 time.sleep(0.75)
55                 doc.html = requests.get(doc.web_url).text
56                 self.DBH.session.commit()
57                 count += 1
58             if count % 250 == 0: print("Documents processed: %d/%d" %(count, len(document
s)))
59
60     # Iterate through the Documents' HTML, parse it and store it.
61     def extract_documents_content(self):
62         documents = self.DBH.session.query(Document).all()
63

```

```

64         count = 0
65         for doc in documents:
66             doc.content = self.extract_content(doc)
67             self.DBH.session.commit()
68             count += 1
69             if count % 250 == 0: print("Documents processed: %d/%d" %(count, len(document
s)))
70
71     def extract_content(self, document):
72         page = self.parse_page(document.html)
73         page = self.remove_unwanted_tags(page)
74         page = self.get_body(page)
75
76         page_content = self.extract_page_content(page)
77         page_content = self.remove_non_relevant_content(page_content)
78         page_content = self.remove_punctuation_and_numbers(page_content)
79         return page_content
80
81     def get_body(self, page):
82         return page.body
83
84     # Discard anything inside footer, header and scripts
85     def remove_unwanted_tags(self, page):
86         for tag in page.find_all(['footer', 'script', 'header']):
87             tag.replace_with('')
88
89         return page
90
91     def remove_non_relevant_content(self, page):
92         for phrase in self.NON_RELEVANT_PHRASES:
93             page = page.replace(phrase, "")
94         return page
95
96     def remove_punctuation_and_numbers(self, page):
97         punctuation = [ '\\', '>', '_', '`', '{', '}', '*', '[',
98                         '^', '+', '!', '(', ':', ';', '"', "'",
99                         '<', '|', '"', '?', '=', '}', '&', '/',
100                        '$', ')', '~', '#', '%', ',']
101
102         page = ''.join(ch for ch in page if ch not in punctuation)
103         page = ''.join([i for i in page if not i.isdigit()])
104         return page

```