

Branch: **develop** ▾**gensim** / **gensim** / **corpora** / **wikicorpus.py**

Find file

Copy path

 **menshikh-iv** Add tox and pytest to gensim, integration with Travis and Appveyor. Fix

8766edc 25 days ago

18 contributors 

Executable File 398 lines (322 sloc) 16.3 KB

```
1  #!/usr/bin/env python
2  # -*- coding: utf-8 -*-
3  #
4  # Copyright (C) 2010 Radim Rehurek <radimrehurek@seznam.cz>
5  # Copyright (C) 2012 Lars Buitinck <larsmans@gmail.com>
6  # Licensed under the GNU LGPL v2.1 - http://www.gnu.org/licenses/lgpl.html
7
8
9  """
10 Construct a corpus from a Wikipedia (or other MediaWiki-based) database dump.
11
12 If you have the `pattern` package installed, this module will use a fancy
13 lemmatization to get a lemma of each token (instead of plain alphabetic
14 tokenizer). The package is available at https://github.com/clips/pattern .
15
16 See scripts/process_wiki.py for a canned (example) script based on this
17 module.
18 """
19
20
21 import bz2
22 import logging
```

```
23 import multiprocessing
24 import re
25 import signal
26 from xml.etree.cElementTree import \
27     iterparse # LXML isn't faster, so let's go with the built-in solution
28
29 from gensim import utils
30 # cannot import whole gensim.corpora, because that imports wikicorpus...
31 from gensim.corpora.dictionary import Dictionary
32 from gensim.corpora.textcorpus import TextCorpus
33
34 logger = logging.getLogger(__name__)
35
36 # ignore articles shorter than ARTICLE_MIN_WORDS characters (after full preprocessing)
37 ARTICLE_MIN_WORDS = 50
38
39 # default thresholds for lengths of individual tokens
40 TOKEN_MIN_LEN = 2
41 TOKEN_MAX_LEN = 15
42
43
44 RE_P0 = re.compile(r'<!--.*?-->', re.DOTALL | re.UNICODE) # comments
45 RE_P1 = re.compile(r'<ref(> ].*?)(</ref>|/>)', re.DOTALL | re.UNICODE) # footnotes
46 RE_P2 = re.compile(r'(\n\\[([a-z][a-z][\w-]*:[^\s\\]+\\])+$', re.UNICODE) # links to languages
47 RE_P3 = re.compile(r'{{{([^\}]+)}}}', re.DOTALL | re.UNICODE) # template
48 RE_P4 = re.compile(r'{{{([^\}]+)}}}', re.DOTALL | re.UNICODE) # template
49 RE_P5 = re.compile(r'\\((\\w+):\\w\\/(.*?))(( .*?))|()\\]', re.UNICODE) # remove URL, keep description
50 RE_P6 = re.compile(r'\\([([^\]]*)\\|([^\]]*)\\|\\]', re.DOTALL | re.UNICODE) # simplify links, keep description
51 RE_P7 = re.compile(r'\\n\\[\\[i\\]mage(.*?)(\\|.*?)*\\|(.*)\\|\\]', re.UNICODE) # keep description of images
52 RE_P8 = re.compile(r'\\n\\[\\[f\\]ile(.*?)(\\|.*?)*\\|(.*)\\|\\]', re.UNICODE) # keep description of files
53 RE_P9 = re.compile(r'<nowiki(> ].*?)(</nowiki>|/>)', re.DOTALL | re.UNICODE) # outside links
54 RE_P10 = re.compile(r'<math(> ].*?)(</math>|/>)', re.DOTALL | re.UNICODE) # math content
55 RE_P11 = re.compile(r'<(.*?)>', re.DOTALL | re.UNICODE) # all other tags
56 RE_P12 = re.compile(r'\\n((\\|)|\\|\\-)|\\|\\|)(.*?)(?=\\n)', re.UNICODE) # table formatting
57 RE_P13 = re.compile(r'\\n(\\|\\|!)(.*?\\|)*\\|\\|', re.UNICODE) # table cell formatting
```

```
58 RE_P14 = re.compile(r'\\\[Category:[^]]*\]\\', re.UNICODE) # categories
59 # Remove File and Image template
60 RE_P15 = re.compile(r'\\\[([fF]ile:|[iI]mage)[^]]*(\\\[\\])', re.UNICODE)
61
62 # MediaWiki namespaces (https://www.mediawiki.org/wiki/Manual:Namespace) that
63 # ought to be ignored
64 IGNORED_NAMESPACES = [
65     'Wikipedia', 'Category', 'File', 'Portal', 'Template',
66     'MediaWiki', 'User', 'Help', 'Book', 'Draft', 'WikiProject',
67     'Special', 'Talk'
68 ]
69
70
71 def filter_wiki(raw):
72     """
73     Filter out wiki mark-up from `raw`, leaving only text. `raw` is either unicode
74     or utf-8 encoded string.
75     """
76     # parsing of the wiki markup is not perfect, but sufficient for our purposes
77     # contributions to improving this code are welcome :)
78     text = utils.to_unicode(raw, 'utf8', errors='ignore')
79     text = utils.decode_htmlentities(text) # '&nbsp;' --> '\xa0'
80     return remove_markup(text)
81
82
83 def remove_markup(text):
84     text = re.sub(RE_P2, '', text) # remove the last list (=languages)
85     # the wiki markup is recursive (markup inside markup etc)
86     # instead of writing a recursive grammar, here we deal with that by removing
87     # markup in a loop, starting with inner-most expressions and working outwards,
88     # for as long as something changes.
89     text = remove_template(text)
90     text = remove_file(text)
91     iters = 0
92     while True:
```

```
93     old, iters = text, iters + 1
94     text = re.sub(RE_P0, '', text) # remove comments
95     text = re.sub(RE_P1, '', text) # remove footnotes
96     text = re.sub(RE_P9, '', text) # remove outside links
97     text = re.sub(RE_P10, '', text) # remove math content
98     text = re.sub(RE_P11, '', text) # remove all remaining tags
99     text = re.sub(RE_P14, '', text) # remove categories
100    text = re.sub(RE_P5, '\\3', text) # remove urls, keep description
101    text = re.sub(RE_P6, '\\2', text) # simplify links, keep description only
102    # remove table markup
103    text = text.replace('||', '\\n|') # each table cell on a separate line
104    text = re.sub(RE_P12, '\\n', text) # remove formatting lines
105    text = re.sub(RE_P13, '\\n\\3', text) # leave only cell content
106    # remove empty mark-up
107    text = text.replace('[]', '')
108    # stop if nothing changed between two iterations or after a fixed number of iterations
109    if old == text or iters > 2:
110        break
111
112    # the following is needed to make the tokenizer see '[[socialist]]s' as a single word 'socialists'
113    # TODO is this really desirable?
114    text = text.replace('[', '').replace(']', '') # promote all remaining markup to plain text
115    return text
116
117
118 def remove_template(s):
119     """Remove template wikimedia markup.
120
121     Return a copy of `s` with all the wikimedia markup template removed. See
122     http://meta.wikimedia.org/wiki/Help:Template for wikimedia templates
123     details.
124
125     Note: Since template can be nested, it is difficult remove them using
126     regular expressions.
127     """
```

```
128
129 # Find the start and end position of each template by finding the opening
130 # '{{' and closing '}}'
131 n_open, n_close = 0, 0
132 starts, ends = [], []
133 in_template = False
134 prev_c = None
135 for i, c in enumerate(iter(s)):
136     if not in_template:
137         if c == '{' and c == prev_c:
138             starts.append(i - 1)
139             in_template = True
140             n_open = 1
141     if in_template:
142         if c == '{':
143             n_open += 1
144         elif c == '}':
145             n_close += 1
146         if n_open == n_close:
147             ends.append(i)
148             in_template = False
149             n_open, n_close = 0, 0
150     prev_c = c
151
152 # Remove all the templates
153 return ''.join([s[end + 1:start] for start, end in zip(starts + [None], [-1] + ends)])
154
155
156 def remove_file(s):
157     """Remove the 'File:' and 'Image:' markup, keeping the file caption.
158
159     Return a copy of `s` with all the 'File:' and 'Image:' markup replaced by
160     their corresponding captions. See http://www.mediawiki.org/wiki/Help:Images
161     for the markup details.
162     """
```

```
163 # The regex RE_P15 match a File: or Image: markup
164 for match in re.finditer(RE_P15, s):
165     m = match.group(0)
166     caption = m[:-2].split('|')[-1]
167     s = s.replace(m, caption, 1)
168 return s
169
170
171 def tokenize(content, token_min_len=TOKEN_MIN_LEN, token_max_len=TOKEN_MAX_LEN, lower=True):
172     """
173     Tokenize a piece of text from wikipedia. The input string `content` is assumed
174     to be mark-up free (see `filter_wiki()`).
175
176     Set `token_min_len`, `token_max_len` as character length (not bytes!) thresholds for individual tokens.
177
178     Return list of tokens as utf8 bytestrings.
179     """
180     # TODO maybe ignore tokens with non-latin characters? (no chinese, arabic, russian etc.)
181     return [
182         utils.to_unicode(token) for token in utils.tokenize(content, lower=lower, errors='ignore')
183         if token_min_len <= len(token) <= token_max_len and not token.startswith('_')
184     ]
185
186
187 def get_namespace(tag):
188     """Returns the namespace of tag."""
189     m = re.match("^(.{*?})", tag)
190     namespace = m.group(1) if m else ""
191     if not namespace.startswith("http://www.mediawiki.org/xml/export-"):
192         raise ValueError("%s not recognized as MediaWiki dump namespace" % namespace)
193     return namespace
194
195
196 _get_namespace = get_namespace
197
```

```
198
199 def extract_pages(f, filter_namespaces=False):
200     """
201     Extract pages from a Mediawiki database dump = open file-like object `f`.
202
203     Return an iterable over (str, str, str) which generates (title, content, pageid) triplets.
204
205     """
206     elems = (elem for _, elem in iterparse(f, events=("end",)))
207
208     # We can't rely on the namespace for database dumps, since it's changed
209     # it every time a small modification to the format is made. So, determine
210     # those from the first element we find, which will be part of the metadata,
211     # and construct element paths.
212     elem = next(elems)
213     namespace = get_namespace(elem.tag)
214     ns_mapping = {"ns": namespace}
215     page_tag = "{%(ns)s}page" % ns_mapping
216     text_path = "./{%(ns)s}revision/{%(ns)s}text" % ns_mapping
217     title_path = "./{%(ns)s}title" % ns_mapping
218     ns_path = "./{%(ns)s}ns" % ns_mapping
219     pageid_path = "./{%(ns)s}id" % ns_mapping
220
221     for elem in elems:
222         if elem.tag == page_tag:
223             title = elem.find(title_path).text
224             text = elem.find(text_path).text
225
226             if filter_namespaces:
227                 ns = elem.find(ns_path).text
228                 if ns not in filter_namespaces:
229                     text = None
230
231             pageid = elem.find(pageid_path).text
232             yield title, text or "", pageid      # empty page will yield None
```

```
233
234     # Prune the element tree, as per
235     # http://www.ibm.com/developerworks/xml/library/x-hiperfparse/
236     # except that we don't need to prune backlinks from the parent
237     # because we don't use LXML.
238     # We do this only for <page>s, since we need to inspect the
239     # ./revision/text element. The pages comprise the bulk of the
240     # file, so in practice we prune away enough.
241     elem.clear()
242
243
244 _extract_pages = extract_pages # for backward compatibility
245
246
247 def process_article(args, tokenizer_func=tokenize, token_min_len=TOKEN_MIN_LEN,
248                    token_max_len=TOKEN_MAX_LEN, lower=True):
249     """
250     Parse a wikipedia article, returning its content as a list of tokens
251     (utf8-encoded strings).
252
253     Set `tokenizer_func` (defaults to `tokenize`) parameter for languages like japanese or thai to perform better
254     tokenization. The `tokenizer_func` needs to take 4 parameters: (text, token_min_len, token_max_len, lower).
255     """
256     text, lemmatize, title, pageid = args
257     text = filter_wiki(text)
258     if lemmatize:
259         result = utils.lemmatize(text)
260     else:
261         result = tokenizer_func(text, token_min_len, token_max_len, lower)
262     return result, title, pageid
263
264
265 def init_to_ignore_interrupt():
266     """Should only be used when master is prepared to handle termination of child processes."""
267     signal.signal(signal.SIGINT, signal.SIG_IGN)
```



```
268
269
270 def _process_article(args):
271     """Should not be called explicitly. Use `process_article` instead."""
272
273     tokenizer_func, token_min_len, token_max_len, lower = args[-1]
274     args = args[:-1]
275
276     return process_article(
277         args, tokenizer_func=tokenizer_func, token_min_len=token_min_len,
278         token_max_len=token_max_len, lower=lower
279     )
280
281
282 class WikiCorpus(TextCorpus):
283     """
284     Treat a wikipedia articles dump (<LANG>wiki-<YYYYMMDD>-pages-articles.xml.bz2
285     or <LANG>wiki-latest-pages-articles.xml.bz2) as a (read-only) corpus.
286
287     The documents are extracted on-the-fly, so that the whole (massive) dump
288     can stay compressed on disk.
289
290     **Note:** "multistream" archives are *not* supported in Python 2 due to
291     `limitations in the core bz2 library
292     <https://docs.python.org/2/library/bz2.html#de-compression-of-files>`.
293
294     >>> wiki = WikiCorpus('enwiki-20100622-pages-articles.xml.bz2') # create word->word_id mapping, takes almost 8h
295     >>> MmCorpus.serialize('wiki_en_vocab200k.mm', wiki) # another 8h, creates a file in MatrixMarket format and mapping
296
297     """
298     def __init__(self, fname, processes=None, lemmatize=utils.has_pattern(), dictionary=None,
299                 filter_namespaces=('0',), tokenizer_func=tokenize, article_min_tokens=ARTICLE_MIN_WORDS,
300                 token_min_len=TOKEN_MIN_LEN, token_max_len=TOKEN_MAX_LEN, lower=True):
301         """
302         Initialize the corpus. Unless a dictionary is provided, this scans the
```

```
303     corpus once, to determine its vocabulary.
304
305     If `pattern` package is installed, use fancier shallow parsing to get
306     token lemmas. Otherwise, use simple regexp tokenization. You can override
307     this automatic logic by forcing the `lemmatize` parameter explicitly.
308     self.metadata if set to true will ensure that serialize will write out article titles to a pickle file.
309
310     Set `article_min_tokens` as a min threshold for article token count (defaults to 50). Any article below this is
311     ignored.
312
313     Set `tokenizer_func` (defaults to `tokenize`) with a custom function reference to control tokenization else use
314     the default regexp tokenization. Set this parameter for languages like japanese or thai to perform better
315     tokenization. The `tokenizer_func` needs to take 4 parameters: (text, token_min_len, token_max_len, lower). The
316     parameter values are as configured on the class instance by default.
317
318     Set `lower` to control if everything should be converted to lowercase or not (default True).
319
320     Set `token_min_len`, `token_max_len` as thresholds for token lengths that are returned (default to 2 and 15).
321
322     """
323     self.fname = fname
324     self.filter_namespaces = filter_namespaces
325     self.metadata = False
326     if processes is None:
327         processes = max(1, multiprocessing.cpu_count() - 1)
328     self.processes = processes
329     self.lemmatize = lemmatize
330     self.tokenizer_func = tokenizer_func
331     self.article_min_tokens = article_min_tokens
332     self.token_min_len = token_min_len
333     self.token_max_len = token_max_len
334     self.lower = lower
335
336     if dictionary is None:
337         self.dictionary = Dictionary(self.get_texts())
```

```
338         else:
339             self.dictionary = dictionary
340
341     def get_texts(self):
342         """
343         Iterate over the dump, returning text version of each article as a list
344         of tokens.
345
346         Only articles of sufficient length are returned (short articles & redirects
347         etc are ignored). This is control by `article_min_tokens` on the class instance.
348
349         Note that this iterates over the **texts**; if you want vectors, just use
350         the standard corpus interface instead of this function::
351
352         >>> for vec in wiki_corpus:
353             >>>     print(vec)
354         """
355
356         articles, articles_all = 0, 0
357         positions, positions_all = 0, 0
358
359         tokenization_params = (self.tokenizer_func, self.token_min_len, self.token_max_len, self.lower)
360         texts = \
361             ((text, self.lemmatize, title, pageid, tokenization_params)
362              for title, text, pageid
363              in extract_pages(bz2.BZ2File(self.fname), self.filter_namespaces))
364         pool = multiprocessing.Pool(self.processes, init_to_ignore_interrupt)
365
366         try:
367             # process the corpus in smaller chunks of docs, because multiprocessing.Pool
368             # is dumb and would load the entire input into RAM at once...
369             for group in utils.chunkize(texts, chunksize=10 * self.processes, maxsize=1):
370                 for tokens, title, pageid in pool.imap(_process_article, group):
371                     articles_all += 1
372                     positions_all += len(tokens)
```

```
373         # article redirects and short stubs are pruned here
374         if len(tokens) < self.article_min_tokens or \
375            any(title.startswith(ignore + ':') for ignore in IGNORED_NAMESPACES):
376             continue
377         articles += 1
378         positions += len(tokens)
379         if self.metadata:
380             yield (tokens, (pageid, title))
381         else:
382             yield tokens
383     except KeyboardInterrupt:
384         logger.warn(
385             "user terminated iteration over Wikipedia corpus after %i documents with %i positions "
386             "(total %i articles, %i positions before pruning articles shorter than %i words)",
387             articles, positions, articles_all, positions_all, ARTICLE_MIN_WORDS
388         )
389     else:
390         logger.info(
391             "finished iterating over Wikipedia corpus of %i documents with %i positions "
392             "(total %i articles, %i positions before pruning articles shorter than %i words)",
393             articles, positions, articles_all, positions_all, ARTICLE_MIN_WORDS
394         )
395         self.length = articles # cache corpus length
396     finally:
397         pool.terminate()
```