



This repository

Search

[Pull requests](#) [Issues](#) [Gist](#)[nltk / nltk](#)[Watch](#)

334

[Star](#)

3,764

[Fork](#)

1,147

[Code](#)[Issues](#) 289[Pull requests](#) 8[Projects](#) 0[Wiki](#)[Pulse](#)[Graphs](#)Branch: [model](#)[nltk / nltk / test / model.doctest](#)[Find file](#)[Copy path](#)**Copper-Head** Unigram counts now also keep track of words in sentence-initial conte...

f5b7c47 17 days ago

1 contributor

332 lines (246 sloc) 10.5 KB

[Raw](#)[Blame](#)[History](#)

```
1  .. Copyright (C) 2001-2015 NLTK Project
2  .. For license information, see LICENSE.TXT
3
4  .. -*- coding: utf-8 -*-
5
6  =====
7  Counting Ngrams
8  =====
9
10 Building a Vocabulary
11 =====
12
13     >>> from nltk.corpus import gutenberg
14     >>> sents = gutenberg.sents("burgess-busterbrown.txt")
15     >>> test_sents = sents[3:5]
16     >>> test_words = [w for s in test_sents for w in s]
17     >>> test_words[:5]
18     ['Buster', 'Bear', 'yawned', 'as', 'he']
19
20     >>> from nltk.model import build_vocabulary
21
22 The first argument to the `build_vocabulary` function is a cutoff value.
23
24     >>> vocab = build_vocabulary(2, test_words)
25
26 This function returns an instance of NgramModelVocabulary, its characteristics
27 are described below.
28 Firstly, tokens with counts greater than or equal to the cutoff value will
29 be considered part of the vocabulary.
30
31     >>> vocab['the']
32     3
33     >>> 'the' in vocab
34     True
35     >>> vocab['he']
36     2
37     >>> 'he' in vocab
38     True
39
40 Tokens with frequency counts less than the cutoff value will be considered not
41 part of the vocabulary even though their entries in the count dictionary are
42 preserved.
43
44     >>> vocab['Buster']
45     1
46     >>> 'Buster' in vocab
47     False
48     >>> vocab['aliens']
49     0
50     >>> 'aliens' in vocab
51     False
52
53 Keeping the count entries for seen words allows us to change the cutoff value
54 without having to recalculate the counts.
55
56     >>> vocab.cutoff = 1
```

```

57     >>> "Buster" in vocab
58     True
59     >>> "aliens" in vocab
60     False
61
62     The cutoff value influences not only membership checking but also the result of
63     getting the size of the vocabulary using the built-in `len`.
64     Note that while the number of keys in the vocab dictionary stays the same, passing
65     it to `len` yields different numbers depending on the cutoff
66
67     >>> len(vocab.keys())
68     37
69     >>> len(vocab)
70     38
71     >>> vocab.cutoff = 2
72     >>> len(vocab)
73     8
74
75     Note also that we add 1 to the size of the vocabulary, so even when cutoff=1
76     `len(vocab) > len(vocab.keys())`. This is done because in many language modeling
77     algorithms the size of the vocabulary is used for normalizing scores and because
78     it needs to account for the unknown token label.
79     This is definitely not an uncontroversial design choice, one that may be reverted
80     based on discussion and usage.
81
82     Vocabulary from multiple sources
83     ~~~~~
84
85     The `build_vocabulary` function can handle multiple text arguments, i.e.
86
87     >>> test_words2 = sents[5]
88     >>> test_words3 = sents[6]
89     >>> vocab2 = build_vocabulary(2, test_words, test_words2, test_words3)
90     >>> len(vocab2)
91     14
92
93
94     Counting Ngrams
95     -----
96
97     >>> from nltk.model import count_ngrams
98
99     The first argument to `count_ngrams` is the highest ngram order to consider, the
100    second is an NgramModelVocabulary instance, followed by one (or more) texts.
101    A text is a sequence of sentences, eg. a list of lists of strings.
102
103    >>> bigram_counts = count_ngrams(2, vocab, sents[3:7])
104
105    This returns an instance of an NgramCounter class which provides an interface
106    to the ngram counts.
107
108    For all non-unigram ngrams (order > 1) the counts are stored in the `ngrams`
109    attribute indexed by ngram order.
110
111    >>> bigram_counts.ngrams[2]
112    <ConditionalFreqDist with 9 conditions>
113
114    The keys of this ConditionalFreqDist are the contexts preceding a word.
115
116    >>> sorted(bigram_counts.ngrams[2].conditions()) # doctest: +NORMALIZE_WHITESPACE
117    [('.',), ('<UNK>',), ('<s>',), ('and',),
118     ('he',), ('his',), ('the',), ('to',), ('yawned',)]
119
120    Each context has a FreqDist of tokens found following it in the text.
121
122    >>> bigram_counts.ngrams[2][('.',)]
123    FreqDist({'</s>': 4})
124
125    Accessing unigram counts works a little differently because they don't have any
126    preceding contexts. They are namely stored as a simple FreqDist in the `unigrams`
127    attribute.
128
129    >>> print(bigram_counts.unigrams)
130    <FreqDist with 10 samples and 121 outcomes>

```

```

131     >>> expected_unigram_counts = {'.': 4,
132     ... '</s>': 4,
133     ... '<s>': 4,
134     ... '<UNK>': 80,
135     ... 'and': 5,
136     ... 'he': 6,
137     ... 'his': 5,
138     ... 'the': 6,
139     ... 'to': 4,
140     ... 'yawned': 3}
141     >>> bigram_counts.unigrams == expected_unigram_counts
142     True
143
144     Tweaking the unknown label
145     ~~~~~
146
147     While counting ngrams all tokens that do not occur frequently enough (read: are
148     not "in" the NgramModelVocabulary instance) are replaced by a special "unknown word"
149     label. By convention this tends to be something like "<UNK>".
150     This is the default for the NgramCounter and it gets returned by
151     the `check_against_vocab` method for words that don't pass the cutoff.
152
153     >>> bigram_counts.check_against_vocab('the')
154     'the'
155     >>> bigram_counts.check_against_vocab('aliens')
156     '<UNK>'
157
158     This can be changed by passing a different value for the `unk_label` argument.
159
160     >>> diff_unk = count_ngrams(2, vocab, sents[3:7], unk_label="UNKNOWN")
161     >>> diff_unk.check_against_vocab('aliens')
162     'UNKNOWN'
163
164
165     Changing the cutoff
166     ~~~~~
167
168     Sometimes we want to compare the effect of different cutoff values on the ngram
169     counts. We can do this specifying an `unk_cutoff` argument that overrides
170     the vocabulary's cutoff value.
171
172     >>> cutoff_1 = count_ngrams(2, vocab, unk_cutoff=1)
173     >>> cutoff_1.vocabulary.cutoff
174     1
175
176     Note that this does not affect the original vocabulary cutoff value.
177
178     >>> vocab.cutoff
179     2
180
181     Changing how ngrams are generated
182     ~~~~~
183
184     NgramCounter uses nltk's `util.ngrams` function to break a text up into ngrams.
185     The behavior of `util.ngrams` is controlled by passing it keyword arguments stored
186     in an attribute of NgramCounter.
187
188     >>> default_kwargs = {'left_pad_symbol': '<s>',
189     ... 'pad_left': True,
190     ... 'pad_right': True,
191     ... 'right_pad_symbol': '</s>'}
192     >>> bigram_counts.ngrams_kwargs == default_kwargs
193     True
194
195     The default values of these arguments have been chosen based on what I think are
196     the most common practices for ngram splitting. You can easily override the settings
197     by passing corresponding key=value pairs to `count_ngrams`
198
199     >>> no_padding = count_ngrams(2, vocab, sents[3:6], pad_left=False, pad_right=False)
200
201     You can always test your ngram generation by calling the `to_ngrams` method
202     with an example sentence.
203
204     >>> list(no_padding.to_ngrams(sents[4])) # doctest: +NORMALIZE_WHITESPACE

```

```

205     [('Once', 'more'), ('more', 'he'), ('he', 'yawned'), ('yawned', ','),
206     (',', 'and'), ('and', 'slowly'), ('slowly', 'got'), ('got', 'to'),
207     ('to', 'his'), ('his', 'feet'), ('feet', 'and'), ('and', 'shook'),
208     ('shook', 'himself'), ('himself', '.')]
209
210 Compare this to the output of the default counter.
211
212 >>> list(bigram_counts.to_ngrams(sents[4])) # doctest: +NORMALIZE_WHITESPACE
213     [('<s>', 'Once'), ('Once', 'more'), ('more', 'he'), ('he', 'yawned'), ('yawned', ','),
214     (',', 'and'), ('and', 'slowly'), ('slowly', 'got'), ('got', 'to'),
215     ('to', 'his'), ('his', 'feet'), ('feet', 'and'), ('and', 'shook'),
216     ('shook', 'himself'), ('himself', '.'), (',', '</s>')]
217
218
219 Multiple (or no) sources
220 ~~~~~
221
222 Just like `build_vocab`, `count_ngrams` can handle multiple source text arguments.
223
224 >>> multiple_texts_counter = count_ngrams(2, vocab, sents[3:6], sents[6:9])
225
226 Moreover, it is also possible to specify no texts whatsoever, which simply creates
227 an empty NgramCounter object so that it can be trained later.
228
229 >>> no_initial_texts = count_ngrams(2, vocab)
230
231 In this case though, it's probably more reader-friendly to use the NgramCounter
232 class directly:
233
234 >>> from nltk.model.counter import NgramCounter
235 >>> no_initial_texts = NgramCounter(2, vocab)
236
237 If you choose to do this, you can call the `train_counts` method and pass it some text
238 to populate your counter instance with some numbers.
239
240 >>> no_initial_texts.train_counts(sents[3:7])
241 >>> expected_unigram_counts = {'.': 4,
242 ... '</s>': 4,
243 ... '<s>': 4,
244 ... '<UNK>': 80,
245 ... 'and': 5,
246 ... 'he': 6,
247 ... 'his': 5,
248 ... 'the': 6,
249 ... 'to': 4,
250 ... 'yawned': 3}
251 >>> no_initial_texts.unigrams == expected_unigram_counts
252 True
253
254
255 =====
256 From Counts to Scores
257 =====
258
259 Example: MLE
260 -----
261
262 Once we counted up the ngrams it is trivial to turn them into a proper model
263 with scores (probabilities): just pass your NgramCounter object to an ngram model class.
264 The module currently has classes for only the basic ngram model types:
265 MLE, Lidstone, and Laplace.
266 Here's an example of how to use an MLE estimator with the counts we've defined.
267
268 >>> from nltk.model import MLNgramModel
269 >>> bigram_model = MLNgramModel(bigram_counts)
270
271 We can access the `bigram_counts` object from the `ngram_counter` attribute:
272
273 >>> bigram_model.ngram_counter == bigram_counts
274 True
275
276 The most important method of all ngram model classes is `score`, which computes
277 the score (aka probability) of a word given some preceding context.
278 In our example we are working with a bigram model, so the context consists of only

```

```
279 one word.
280
281 >>> ex_score = bigram_model.score("yawned", ["he"])
282
283 Note that an MLE score is no more than the relative frequency of the word given
284 its context (it just looks much nicer!).
285
286 >>> bigram_counts.ngrams[2][('he',)].freq('yawned') == ex_score
287 True
288
289 Also note that both lists and tuples are supported as `context` arguments.
290
291 >>> bigram_model.score("yawned", ("he",)) == ex_score
292 True
293
294 In addition to `score` all ngram model classes provide the following methods:
295 - logscore
296 - entropy
297 - perplexity
298 Their usage is described in the following section.
299
300
301 Custom NgramModel Classes
302 -----
303
304 In case you need an ngram model estimator that's not currently in the module,
305 it's quite easy to define your own class for that, you simply subclass the
306 `BaseNgramModel` class. The following section describes.
307
308 >>> from nltk.model import BaseNgramModel
309 >>> base_model = BaseNgramModel(bigram_counts)
310
311 Since it is intended only as a base class for real estimators, the `BaseNgramModel`
312 doesn't have an interesting `score` method, it always returns 0.5.
313
314 >>> base_model.score("test", ('abc',))
315 0.5
316 >>> base_model.score("abc", ("test",))
317 0.5
318
319 Children classes are expected to override this method and provide their own
320 logic for how to compute a score.
321 The advantage of having a base ngram model class is that there are plenty of methods
322 that do not depend of `score` in their implementations which are all defined in
323 the `BaseNgramModel` class and which you get "for free" by inheriting from it.
324 Here are examples of these methods:
325
326 >>> base_model.logscore("abc", ("test",))
327 -1.0
328 >>> base_model.entropy(sents[4])
329 1.0
330 >>> base_model.perplexity(sents[4])
331 2.0
```

