

Popular Topic Mining from Blogs

Stone Fang (Student ID: 19049045)
Computers and Information Sciences
Auckland University of Technology
Auckland, New Zealand
fnk7060@autuni.ac.nz

I. OVERVIEW

A. *Objective*

II. RELATED WORK

III. SOLUTION

IV. ACCURACY INSURANCE

V. CONCLUSION, OPEN ISSUES AND FUTURE WORK

APPENDIX
SOURCE CODE IN PYTHON

```
1 #!/usr/bin/env python
2 # -*- coding: utf-8 -*-
3
4 import sys
5 import os.path
6 from glob import glob
7 from tqdm import tqdm
8 import pickle
9
10 import itertools
11 from collections import namedtuple, Counter
12 import pandas as pd
13 from bs4 import BeautifulSoup
14
15 import nltk
16 from nltk.corpus import stopwords
17 from sklearn.feature_extraction.text import TfidfVectorizer
18
19 _DEBUG = True
20
21 sw = stopwords.words("english")
22
23
24 #####
25 #                               Utility functions                               #
26 #####
27
28 def len2d(iter2d):
29     return sum(len(d) for d in iter2d)
30
31 def list2d(iter2d):
32     return [[x for x in inner] for inner in iter2d]
33
34 def flatten2d(list2d):
35     return itertools.chain.from_iterable(list2d)
36
37 def flatten3d(list3d):
38     return itertools.chain.from_iterable(flatten2d(list3d))
39
40 def mapbar(f, seq, desc):
41     for e in tqdm(seq, desc):
42         yield f(e)
43
44 def map2d(f, docs):
45     with tqdm(total=len2d(docs)) as pbar:
46         def _helper(sent):
47             pbar.update(1)
48             return f(sent)
49
50         for doc in docs:
51             yield map(_helper, doc)
52
53 def map3d(f, docs):
54     with tqdm(total=len2d(docs)) as pbar:
55         def _helper(sent):
56             pbar.update(1)
57             return [f(word) for word in sent]
58
59         for doc in docs:
60             yield map(_helper, doc)
61
62 def foreach3d(f, docs):
63     with tqdm(total=len2d(docs)) as pbar:
64         for doc in docs:
```

```

65         for sent in doc:
66             for word in sent:
67                 f(word)
68             pbar.update(1)
69
70 def foreach2d(f, docs):
71     with tqdm(total=len2d(docs)) as pbar:
72         for doc in docs:
73             for sent in doc:
74                 f(sent)
75             pbar.update(1)
76
77 #####
78 #           Codes for data reading & transformation           #
79 #####
80
81 Record = namedtuple('Record', ['meta', 'posts'])
82 Post = namedtuple('Post', ['date', 'text'])
83 MetaData = namedtuple('MetaData', ['id', 'gender', 'age', 'category', 'zodiac'])
84
85 def parse_meta_data(meta_data_str):
86     arr = meta_data_str.strip().split('.')
87     return MetaData(arr[0], arr[1], int(arr[2]), arr[3], arr[4])
88
89 def read_blog_file(fpath):
90     try:
91         with open(fpath, encoding='utf-8', errors='ignore') as f:
92             soup = BeautifulSoup(f.read(), "xml")
93             blog = soup.Blog
94     except ParseError:
95         print('Error: invalid xml file {}'.format(fpath))
96         raise
97     return []
98
99     posts = []
100     state = 'date'
101     for c in blog.find_all(recursive=False):
102         if c.name != state:
103             print('Warning: inconsistent format in file {}'.format(fpath))
104         if state == 'date':
105             try:
106                 date_str = c.text.strip()
107                 date = date_str
108             except ValueError:
109                 print('Warning: invalid date {} in file {}'.format(c.text, fpath) \
110                       .format(c.text, fpath))
111             state = 'post'
112         else:
113             text = c.text.strip()
114             state = 'date'
115             posts.append(Post(date, text))
116     posts.sort(key=lambda p: p.date)
117     return posts
118
119 def read_blogs(path, force=False, cache_file='blogs.pkl'):
120     if not force and cache_file is not None and os.path.exists(cache_file):
121         print('load dataset from cached pickle file ' + cache_file)
122         with open(cache_file, 'rb') as f:
123             dataset = pickle.load(f)
124         return dataset
125
126     dataset = read_blogs_xml(path)
127
128     # save to pickle file for fast loading next time
129     if cache_file is not None:
130         with open(cache_file, 'wb') as f:
131             print('save dataset to pickle file ' + cache_file)

```

```

132         pickle.dump(dataset, f)
133
134     return dataset
135
136 def read_blogs_xml(path):
137     print('reading all data files from directory {} ...'.format(path))
138     dataset = []
139
140     if _DEBUG: # use small files for fast debugging
141         files = [os.path.join(path, fname) for fname in ['3998465.male.17.indUnk.Gemini.xml',
142                                                         '3949642.male.25.indUnk.Leo.xml', '3924311.male.27.HumanResources.Gemini.xml']]
143     else:
144         files = glob(os.path.join(path, '*'))
145
146     for fpath in tqdm(files):
147         fname = os.path.basename(fpath)
148         meta_data = parse_meta_data(fname)
149         posts = read_blog_file(fpath)
150         rec = Record(meta_data, posts)
151         dataset.append(rec)
152     return dataset
153
154 def show_summary(dataset):
155     '''This function describes the summary of dataset or human inspection.
156     It's not necessary for the mining process.
157
158     Parameters
159     -----
160     dataset : list of Record
161         The blog dataset
162     '''
163
164     df = pd.DataFrame([d.meta for d in dataset])
165     df['blog_count'] = [len(d.posts) for d in dataset]
166     print(df.describe(include='all'))
167     print('{} possible values for "gender": {}'.format(
168         len(df.gender.unique()), ', '.join(sorted(df.gender.unique()))))
169     print('{} possible values for category: {}'.format(
170         len(df.category.unique()), ', '.join(sorted(df.category.unique()))))
171     print('{} possible values for zodiac: {}'.format(
172         len(df.zodiac.unique()), ', '.join(sorted(df.zodiac.unique()))))
173
174     #####
175     #               Codes for topic mining               #
176     #####
177
178 def tokenise(dataset):
179     '''
180     consider all the blogs from one person as a document
181
182     Returns
183     -----
184     docs: list of list of list
185         a list of documents, each of which is a list of sentences,
186         each of which is a list of words.
187     '''
188
189     print('tokenising the text dataset...')
190     docs = []
191     vocab = set()
192     with tqdm(total=sum(len(rec.posts) for rec in dataset)) as pbar:
193         for rec in dataset:
194             doc = []
195             for post in rec.posts:
196                 for sent_str in nltk.sent_tokenize(post.text):
197                     sent = [w.lower() for w in nltk.word_tokenize(sent_str)]
198                     doc.append(sent)

```

```

199         vocab.update(sent)
200         pbar.update(1)
201         docs.append(doc)
202     return sorted(vocab), docs
203
204 def get_things(docs, n=5):
205     things = filter(lambda wp: wp[1] == 'NN', flatten3d(docs))
206     tf = nltk.FreqDist(things)
207     return tf.most_common(n)
208
209
210 def get_surroundings(words, docs, n=2):
211     '''expand the topic to be 2 verb/noun before and 2 verb/noun after the topic
212     '''
213
214     print('get surrounding {} nouns/verbs for words {}'.format(n, words))
215
216     sur = {w: Counter() for w in words}
217
218     target_pos_tags = ('NN', 'NNS', 'VB', 'VBP', 'VBD', 'VBN')
219
220     def _helper(sent):
221         for w in words:
222             try:
223                 idx = sent.index(w)
224                 except ValueError:
225                     continue
226
227                 after = 0
228                 for (wi, pi) in sent[(idx+1):]:
229                     if pi in target_pos_tags:
230                         sur[w][(wi, pi)] += 1
231                         after += 1
232                     if after == n:
233                         break
234
235                 before = 0
236                 for (wi, pi) in reversed(sent[:idx]):
237                     if pi in target_pos_tags:
238                         sur[w][(wi, pi)] += 1
239                         before += 1
240                     if before == n:
241                         break
242
243     foreach2d(_helper, docs)
244     return sur
245
246 def mine_topic_by_freq(dataset):
247     vocab, docs = tokenise(dataset)
248     print('Size of vocabulary: {}'.format(len(vocab)))
249     print(vocab[:100])
250
251     print('POS tagging...')
252     tagged_docs = list2d(map2d(lambda s: nltk.pos_tag(s), docs))
253
254     things = get_things(tagged_docs)
255     print('things: ', things)
256
257
258     thing_words = [(w, pos) for ((w, pos), c) in things]
259     keywords = get_surroundings(thing_words, tagged_docs, n=2)
260     print(keywords)
261
262
263 def main():
264     if _DEBUG:
265         dataset = read_blogs('blogs', cache_file=None)

```

```
266     else:
267         dataset = read_blogs('blogs')
268
269     mine_topic_by_freq(dataset)
270     return
271
272 if __name__ == '__main__':
273     main()
```