

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

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
12 import pandas as pd
13 # from xml.etree import ElementTree
14 # from xml.etree.ElementTree import ParseError
15 from bs4 import BeautifulSoup
16 # from datetime import datetime
17
18 import nltk
19 from nltk.corpus import stopwords
20 sw = stopwords.words("english")
21
22
23 #####
24 #                               Utility functions                               #
25 #####
26
27 def flatten2d(list2d):
28     return itertools.chain.from_iterable(list2d)
29
30 def flatten3d(list3d):
31     return itertools.chain.from_iterable(flatten2d(list3d))
32
33 def map2d(f, docs):
34     for doc in docs:
35         yield [f(sent) for sent in doc]
36
37 def map3d(f, docs):
38     for doc in docs:
39         yield [[f(word) for word in sent] for sent in doc]
40
41 #####
42 #                               Codes for data reading & transformation          #
43 #####
44
45 Record = namedtuple('Record', ['meta', 'posts'])
46 Post = namedtuple('Post', ['date', 'text'])
47 MetaData = namedtuple('MetaData', ['id', 'gender', 'age', 'category', 'zodiac'])
48
49 def parse_meta_data(meta_data_str):
50     arr = meta_data_str.strip().split('.')
51     return MetaData(arr[0], arr[1], int(arr[2]), arr[3], arr[4])
52
53 # _parser = ElementTree.XMLParser(encoding="utf-8")
54 def read_blog_file(fpath):
55     try:
56         # tree = ElementTree.parse(fpath, parser=_parser)
57         with open(fpath, encoding='utf-8', errors='ignore') as f:
58             soup = BeautifulSoup(f.read(), "xml")
59             blog = soup.Blog
60     except ParseError:
61         print('Error: invalid xml file {}'.format(fpath))
62         raise
63     return []
64
```

```

65 posts = []
66 state = 'date'
67 for c in blog.find_all(recursive=False):
68     # print(c)
69     # print(c.text)
70     # check the <date> and <post> tags appear alternately
71     if c.name != state:
72         print('Warning: inconsistent format in file {}'.format(fpath))
73     if state == 'date':
74         try:
75             date_str = c.text.strip()
76             # date_str = date_str.replace('janvier', 'january') \
77             #     .replace('mars', 'march') \
78             #     .replace('avril', 'april') \
79             #     .replace('mai', 'may') \
80             #     .replace('juin', 'june') \
81             #     .replace('juillet', 'july')
82             # date = pd.to_datetime(date_str, format='%d,%B,%Y')
83             date = date_str
84         except ValueError:
85             print('Warning: invalid date {} in file {}'.format(c.text, fpath) \
86                   .format(c.text, fpath))
87             # date = datetime.fromtimestamp(0)
88             state = 'post'
89         else:
90             text = c.text.strip()
91             state = 'date'
92             posts.append(Post(date, text))
93             # print(c, c.text)
94 posts.sort(key=lambda p: p.date)
95 # print(posts)
96 # print(pd.DataFrame(posts))
97 # sys.exit()
98 return posts
99
100 def read_blogs(path, force=False, cache_file='blogs.pkl'):
101     if not force and os.path.exists(cache_file):
102         print('load dataset from cached pickle file ' + cache_file)
103         with open(cache_file, 'rb') as f:
104             dataset = pickle.load(f)
105         return dataset
106
107     dataset = read_blogs_xml(path)
108
109     # save to pickle file for fast loading next time
110     with open(cache_file, 'wb') as f:
111         print('save dataset to pickle file ' + cache_file)
112         pickle.dump(dataset, f)
113
114     return dataset
115
116 def read_blogs_xml(path):
117     print('reading all data files from directory {} ...'.format(path))
118     dataset = []
119     # for fpath in tqdm(glob(os.path.join(path, '*'))):
120     for fpath in list(glob(os.path.join(path, '*'))[:3]):
121         # print(fpath)
122         fname = os.path.basename(fpath)
123         meta_data = parse_meta_data(fname)
124         # print(meta_data)
125         posts = read_blog_file(fpath)
126         rec = Record(meta_data, posts)
127         dataset.append(rec)
128     return dataset
129
130 def show_summary(dataset):
131     '''This function describes the summary of dataset or human inspection.

```

```

132 It's not necessary for the mining process.
133
134 Parameters
135 -----
136 dataset : list of Record
137           The blog dataset
138 '''
139
140 df = pd.DataFrame([d.meta for d in dataset])
141 df['blog_count'] = [len(d.posts) for d in dataset]
142 # print(df)
143 print(df.describe(include='all'))
144 print('{} possible values for "gender": {}'.format(
145     len(df.gender.unique()), ', '.join(sorted(df.gender.unique()))))
146 # print('{} possible values for "{}": {}'.format(
147 #     len(df.age.unique()), ', '.join(sorted(df.age.unique()))))
148 print('{} possible values for category: {}'.format(
149     len(df.category.unique()), ', '.join(sorted(df.category.unique()))))
150 print('{} possible values for zodiac: {}'.format(
151     len(df.zodiac.unique()), ', '.join(sorted(df.zodiac.unique()))))
152
153 #####
154 # Codes for topic mining #
155 #####
156
157 def tokenise(dataset):
158     '''
159     consider all the blogs from one person as a document
160
161     Returns
162     -----
163     docs: list of list of list
164           a list of documents, each of which is a list of sentences,
165           each of which is a list of words.
166     '''
167
168     print('tokenising the text dataset...')
169     docs = []
170     for rec in dataset:
171         doc = []
172         for post in rec.posts:
173             # print(post)
174             for sent in nltk.sent_tokenize(post.text):
175                 doc.append([w.lower() for w in nltk.word_tokenize(sent)])
176             # print(doc)
177         docs.append(doc)
178         # print(doc)
179         # return(docs)
180     return docs
181
182 def get_things(docs):
183     things = filter(lambda wp: wp[1] == 'NN', flatten3d(docs))
184     tf = nltk.FreqDist(things)
185     print(tf.most_common(50))
186
187 #TODO: remove stop words
188 def count_word(dataset):
189     docs = tokenise(dataset)
190
191     print('POS tagging...')
192     tagged_docs = list(map2d(lambda s: nltk.pos_tag(s), docs))
193     print(tagged_docs[0][0])
194     # print(sorted(set(p for w, p in flatten3d(tagged_docs))))
195
196     things = get_things(tagged_docs)
197
198     print('counting word frequencies...')

```

```
199     tf = nltk.FreqDist(flatten3d(tagged_docs))
200     print(tf.most_common(50))
201
202 def main():
203     # read_blogs('blogs')
204     # read_blogs('blogs', force=True, cache_file='blogs-10.pkl')
205     # dataset = read_blogs('.', cache_file='blogs-10.pkl')
206     dataset = read_blogs('.', cache_file='blogs-3.pkl')
207
208     count_word(dataset)
209     return
210
211 if __name__ == '__main__':
212     main()
```