

# 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. METRICS

## APPENDIX

```

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 #           Codes for data reading & transformation           #
24 #####
25
26 Record = namedtuple('Record', ['meta', 'posts'])
27 Post = namedtuple('Post', ['date', 'text'])
28 MetaData = namedtuple('MetaData', ['id', 'gender', 'age', 'category', 'zodiac'])
29
30 def parse_meta_data(meta_data_str):
31     arr = meta_data_str.strip().split('.')
32     return MetaData(arr[0], arr[1], int(arr[2]), arr[3], arr[4])
33
34 # _parser = ElementTree.XMLParser(encoding="utf-8")
35 def read_blog_file(fpath):
36     try:
37         # tree = ElementTree.parse(fpath, parser=_parser)
38         with open(fpath, encoding='utf-8', errors='ignore') as f:
39             soup = BeautifulSoup(f.read(), "xml")
40             blog = soup.Blog
41     except ParseError:
42         print('Error: invalid xml file {}'.format(fpath))
43         raise
44     return []
45
46 posts = []
47 state = 'date'
48 for c in blog.find_all(recursive=False):
49     # print(c)
50     # print(c.text)
51     # check the <date> and <post> tags appear alternately
52     if c.name != state:
53         print('Warning: inconsistent format in file {}'.format(fpath))
54     if state == 'date':
55         try:
56             date_str = c.text.strip()
57             # date_str = date_str.replace('janvier', 'january') \
58             #     .replace('mars', 'march') \
59             #     .replace('avril', 'april') \
60             #     .replace('mai', 'may') \
61             #     .replace('juin', 'june') \
62             #     .replace('juillet', 'july')
63             # date = pd.to_datetime(date_str, format='%d,%B,%Y')
64             date = date_str
65         except ValueError:

```

```

66         print('Warning: invalid date {} in file {}'.format(c.text, fpath) \
67               .format(c.text, fpath))
68         date = datetime.fromtimestamp(0)
69         state = 'post'
70     else:
71         text = c.text.strip()
72         state = 'date'
73         posts.append(Post(date, text))
74         # print(c, c.text)
75     posts.sort(key=lambda p: p.date)
76     # print(posts)
77     # print(pd.DataFrame(posts))
78     # sys.exit()
79     return posts
80
81 def read_blogs(path, force=False, cache_file='blogs.pkl'):
82     if not force and os.path.exists(cache_file):
83         print('load dataset from cached pickle file ' + cache_file)
84         with open(cache_file, 'rb') as f:
85             dataset = pickle.load(f)
86         return dataset
87
88     dataset = read_blogs_xml(path)
89
90     # save to pickle file for fast loading next time
91     with open(cache_file, 'wb') as f:
92         print('save dataset to pickle file ' + cache_file)
93         pickle.dump(dataset, f)
94
95     return dataset
96
97 def read_blogs_xml(path):
98     print('reading all data files from directory {}'.format(path))
99     dataset = []
100     # for fpath in tqdm(glob(os.path.join(path, '*'))):
101     for fpath in list(glob(os.path.join(path, '*'))[:3]):
102         # print(fpath)
103         fname = os.path.basename(fpath)
104         meta_data = parse_meta_data(fname)
105         # print(meta_data)
106         posts = read_blog_file(fpath)
107         rec = Record(meta_data, posts)
108         dataset.append(rec)
109     return dataset
110
111 def show_summary(dataset):
112     '''This function describes the summary of dataset or human inspection.
113     It's not necessary for the mining process.
114
115     Parameters
116     -----
117     dataset : list of Record
118         The blog dataset
119     '''
120
121     df = pd.DataFrame([d.meta for d in dataset])
122     df['blog_count'] = [len(d.posts) for d in dataset]
123     # print(df)
124     print(df.describe(include='all'))
125     print('{} possible values for "gender": {}'.format(
126         len(df.gender.unique()), ', '.join(sorted(df.gender.unique()))))
127     # print('{} possible values for "{}": {}'.format(
128     #     len(df.age.unique()), ', '.join(sorted(df.age.unique()))))
129     print('{} possible values for category: {}'.format(
130         len(df.category.unique()), ', '.join(sorted(df.category.unique()))))
131     print('{} possible values for zodiac: {}'.format(
132         len(df.zodiac.unique()), ', '.join(sorted(df.zodiac.unique()))))

```

```

133 #####
134 # Codes for topic mining #
135 #####
136 #####
137
138 def tokenise(dataset):
139     """
140     consider all the blogs from one person as a document
141
142     Returns
143     -----
144     docs: list of list of list
145           a list of documents, each of which is a list of sentences,
146           each of which is a list of words.
147     """
148
149     print('tokenising the text dataset...')
150     docs = []
151     for rec in dataset:
152         doc = []
153         for post in rec.posts:
154             # print(post)
155             for sent in nltk.sent_tokenize(post.text):
156                 doc.append([w.lower() for w in nltk.word_tokenize(sent)])
157             # print(doc)
158         docs.append(doc)
159         # print(doc)
160         # return(docs)
161     return docs
162
163 def flatten2d(list2d):
164     return itertools.chain.from_iterable(list2d)
165
166 def flatten3d(list3d):
167     return itertools.chain.from_iterable(flatten2d(list3d))
168
169 def map2d(f, docs):
170     for doc in docs:
171         yield [f(sent) for sent in doc]
172
173 def map3d(f, docs):
174     for doc in docs:
175         yield [[f(word) for word in sent] for sent in doc]
176
177 def get_things(docs):
178     things = filter(lambda wp: wp[1] == 'NN', flatten3d(docs))
179     tf = nltk.FreqDist(things)
180     print(tf.most_common(50))
181
182 #TODO: remove stop words
183 def count_word(dataset):
184     docs = tokenise(dataset)
185
186     print('POS tagging...')
187     tagged_docs = list(map2d(lambda s: nltk.pos_tag(s), docs))
188     print(tagged_docs[0][0])
189     # print(sorted(set(p for w, p in flatten3d(tagged_docs))))
190
191     things = get_things(tagged_docs)
192
193     print('counting word frequencies...')
194     tf = nltk.FreqDist(flatten3d(tagged_docs))
195     print(tf.most_common(50))
196
197 def main():
198     # read_blogs('blogs')
199     # read_blogs('blogs', force=True, cache_file='blogs-10.pkl')

```

```
200     # dataset = read_blogs('.', cache_file='blogs-10.pkl')
201     dataset = read_blogs('.', cache_file='blogs-3.pkl')
202
203     count_word(dataset)
204     return
205
206 if __name__ == '__main__':
207     main()
```