

# 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

People's concerns and opinions are important reference for innovations of new products or services. However, accomplishing such task by humans is expensive, time-consuming and difficult to scale. As a response, a number of individuals and organisations are leveraging text mining technologies to mining meaningful information from large volume of text such as news media [1]. Among a variety of studies and applications, topic modelling is an important method to extract hot topics which reflects public attention and opinion from massive texts [1]–[3]. However, effective method of extracting useful information from text on the Internet remains an open challenge [3].

Evaluation of topics mined from text is another challenge, mostly due to the lack of ground truth because topic modelling is an unsupervised learning task [4].

The goal of this project is to mine most popular topics that people were discussing from blog posts by utilising various text mining algorithms and tools. Specifically, we will find two most popular topics for each group in the following demographics:

- Males
- Females
- People younger than or equal to 20 years old
- People older than 20 years old
- Everyone

The remainder of this article is organised as follows. In section II related works on topic mining and evaluation will be reviewed. The methodology of topic mining are detailed in section III, while the results, analysis and evaluations are presented in section IV. The works of this article are summarised in section V and open issues and future works are discussed in section VI.

## II. RELATED WORK

Jacobi, Atteveldt, and Welbers [1] conducted an in-depth study of how to apply topic modelling technologies on analysis of qualitative data in academic research.

## III. RESEARCH DESIGN

In this section the solution will be described in detail. First an overview of the dataset is given, and then the algorithm of topic mining is detailed.

### A. Data Description

The dataset contains 19,320 files in XML format, each containing articles of one person posted generally between 2001 and 2004. Metadata of the bloggers includes gender, age, category, and zodiac. In addition, the number of posts for each person are also counted. The result is summarised in figure 1. From this figure we can acquire some basic statistics of the dataset, including

- 1) Gender: data samples are quite evenly distributed over both genders.
- 2) Age: most bloggers are younger than 30, almost of them under 20. On the other hand, there are two gaps around 20 and 30 which may implies some missing data points in the dataset
- 3) Category: the most frequent category is unknown, which is trivial, while the second frequent one is student, far more than other categories.
- 4) Zodiac: The distribution over zodiac is reasonable even.
- 5) Number of posts: most bloggers published less than 100 posts, while the peak appears at 10, which implies people are most likely to write around 10 posts.

### B. Topic Mining Algorithm

The general idea for mining popular topics used in this project is to find the most significant “things” mentioned in the overall dataset, as well as the closely related information.

The overall architecture of the algorithm is shown as figure 2.

- 1) *Data Cleaning:*
- 2) *Tokenization:*
- 3) *POS Tagging:*
- 4) *NER:*
- 5) *Stopwords Removal:*
- 6) *Stemming and Lemmatization:*
- 7) *Counting and TF-IDF:*

## IV. RESULTS, ANALYSIS, AND EVALUATION

## V. CONCLUSION

## VI. OPEN ISSUES AND FUTURE WORKS

## REFERENCES

- [1] C. Jacobi, W. v. Atteveldt, and K. Welbers, “Quantitative analysis of large amounts of journalistic texts using topic modelling,” *Digital Journalism*, vol. 4, no. 1, pp. 89–106, Jan. 2, 2016, ISSN: 2167-0811. DOI: 10.1080/21670811.2015.1093271.

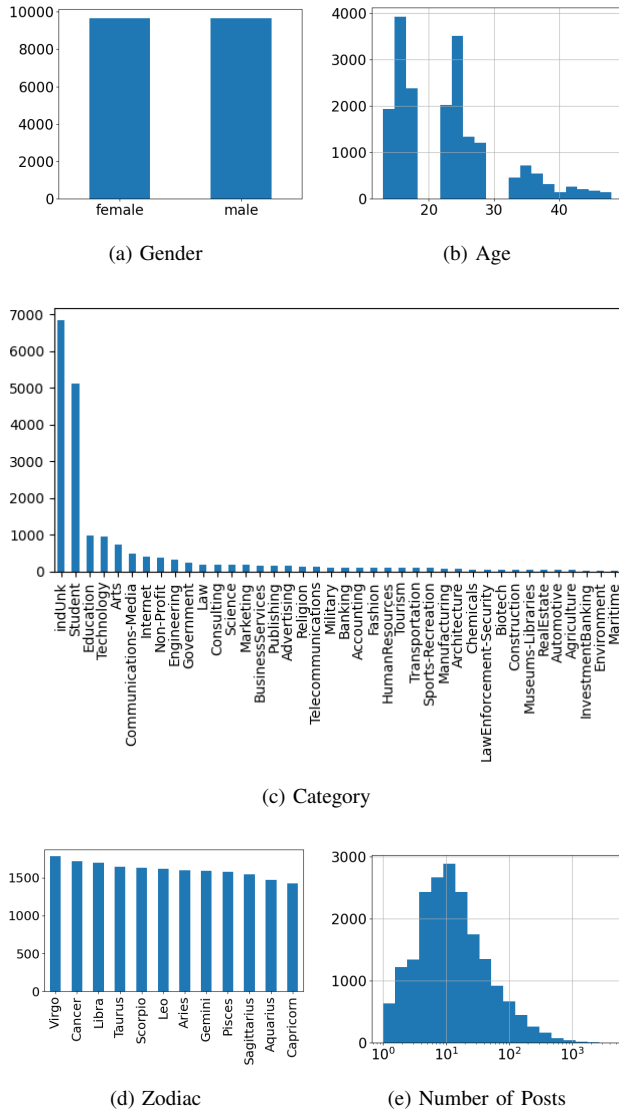


Fig. 1. Data Overview. Histogram over (a) gender (b) age (c) category (d) zodiac (e) number of posts

- [2] P. Waila, V. K. Singh, and M. K. Singh, "Blog text analysis using topic modeling, named entity recognition and sentiment classifier combine," in *2013 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, Mysore: IEEE, Aug. 2013, pp. 1166–1171, ISBN: 978-1-4673-6217-7 978-1-4799-2432-5 978-1-4799-2659-6. DOI: 10.1109/ICACCI.2013.6637342.
- [3] J. Guo, P. Zhang, JianlongTan, and L. Guo, "Mining hot topics from twitter streams," *Procedia Computer Science*, vol. 9, pp. 2008–2011, 2012, ISSN: 18770509. DOI: 10.1016/j.procs.2012.04.224.
- [4] J. Boyd-Graber, D. Mimno, and D. Newman, "Care and feeding of topic models: Problems, diagnostics, and improvements," *Handbook of Mixed Membership Models and Their Applications*, p. 30,

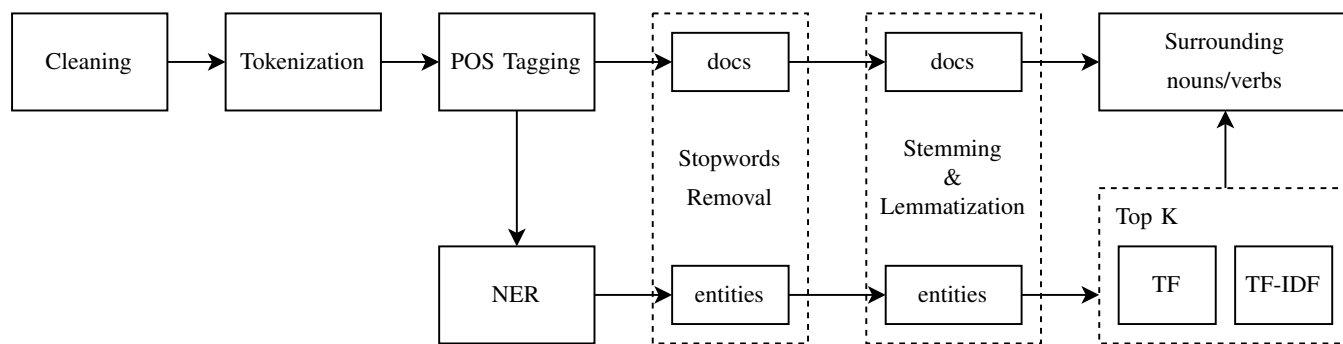


Fig. 2. Algorithm

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 import json
10 from datetime import date
11 import pprint
12 pp = pprint.PrettyPrinter(indent=2)
13
14 import random
15 import itertools
16 from collections import namedtuple, Counter, OrderedDict, defaultdict
17 import heapq
18 from operator import itemgetter
19 import re
20 from bs4 import BeautifulSoup
21 import numpy as np
22
23 from spellchecker import SpellChecker
24 import nltk
25 from nltk.corpus import stopwords
26 from nltk.corpus import wordnet
27 from nltk.stem import PorterStemmer, LancasterStemmer, WordNetLemmatizer
28
29 NUM_SAMPLES = None
30 _DEBUG = False
31
32 _DEBUG = True
33
34 STOPWORDS = set(stopwords.words("english"))
35 # Add more stopwords manually
36 with open('stopwords1.txt') as f:
37     STOPWORDS.update(w.strip().lower() for w in f)
38 STOPWORDS.update(['i\'m', 'dont', '\t', '\m', '\s', '\re', '\ve',
39                  'haha', 'hah', 'wow', 'hehe', 'heh',
40                  'ah', 'ahh', 'hm', 'hmm', 'urllink', 'ok', 'hey', 'yay', 'yeah'])
41
42 #####
43 #                               Utility functions                               #
44 #####
45
46 def len2d(iter2d):
47     return sum(len(d) for d in iter2d)
48
49 def list2d(iter2d):
50     return [[x for x in inner] for inner in iter2d]
51
52 def flatten2d(list2d):
53     return itertools.chain.from_iterable(list2d)
54
55 def flatten3d(list3d):
56     return itertools.chain.from_iterable(flatten2d(list3d))
57
58 def mapbar(f, seq, desc):
59     for e in tqdm(seq, desc):
60         yield f(e)
61
62 def map2d(f, docs):
63     with tqdm(total=len2d(docs)) as pbar:
64         def _helper(sent):
```

```

65         pbar.update(1)
66         return f(sent)
67
68     return [list(map(_helper, doc)) for doc in docs]
69
70 def map3d(f, docs):
71     with tqdm(total=len2d(docs)) as pbar:
72         def _helper(sent):
73             pbar.update(1)
74             return [f(word) for word in sent]
75
76     return [list(map(_helper, doc)) for doc in docs]
77
78 def foreach3d(f, docs):
79     with tqdm(total=len2d(docs)) as pbar:
80         for doc in docs:
81             for sent in doc:
82                 for word in sent:
83                     f(word)
84             pbar.update(1)
85
86 def foreach2d(f, docs):
87     with tqdm(total=len2d(docs)) as pbar:
88         for doc in docs:
89             for sent in doc:
90                 f(sent)
91             pbar.update(1)
92
93 def filter3d(f, docs):
94     ret = []
95     with tqdm(total=len2d(docs)) as pbar:
96         def _helper_doc(doc):
97             for sent in doc:
98                 pbar.update(1)
99                 out = [word for word in sent if f(word)]
100                 if len(out) > 0:
101                     yield out
102
103     for doc in docs:
104         out = list(_helper_doc(doc))
105         ret.append(out)
106     return ret
107
108 def load_pkl(fpath):
109     print('load dataset from cached pickle file ' + fpath)
110     with open(fpath, 'rb') as f:
111         dataset = pickle.load(f)
112     return dataset
113
114 def save_pkl(obj, fpath):
115     with open(fpath, 'wb') as f:
116         print('save dataset to pickle file ' + fpath)
117         pickle.dump(obj, f)
118
119 def save_json(obj, fpath, indent=2):
120     with open(fpath, 'w', encoding="utf8") as f:
121         print('save dataset to json file ' + fpath)
122         json.dump(obj, f, indent=indent)
123
124 #####
125 #           Codes for data reading & transformation           #
126 #####
127
128 Record = namedtuple('Record', ['meta', 'posts'])
129 Post = namedtuple('Post', ['date', 'text'])
130 MetaData = namedtuple('MetaData', ['id', 'gender', 'age', 'category', 'zodiac'])
131

```

```

132 def parse_meta_data(meta_data_str):
133     arr = meta_data_str.lower().strip().split('.')
134     return MetaData(arr[0], arr[1], int(arr[2]), arr[3], arr[4])
135
136 def read_blog_file(fpath):
137     try:
138         with open(fpath, encoding='utf-8', errors='ignore') as f:
139             soup = BeautifulSoup(f.read(), "xml")
140             blog = soup.Blog
141     except ParseError:
142         print('Error: invalid xml file {}'.format(fpath))
143         raise
144     return []
145
146     posts = []
147     state = 'date'
148     for c in blog.find_all(recursive=False):
149         if c.name != state:
150             print('Warning: inconsistent format in file {}'.format(fpath))
151         if state == 'date':
152             try:
153                 date_str = c.text.strip()
154                 date = date_str
155             except ValueError:
156                 print('Warning: invalid date {} in file {}' \
157                       .format(c.text, fpath))
158             state = 'post'
159         else:
160             text = c.text.strip()
161             state = 'date'
162             posts.append(Post(date, text))
163     posts.sort(key=lambda p: p.date)
164     return posts
165
166 def read_blogs(path, force=False, cache_file='blogs.pkl'):
167     if not force and cache_file is not None and os.path.exists(cache_file):
168         return load_pkl(cache_file)
169
170     dataset = read_blogs_xml(path)
171
172     # save to pickle file for fast loading next time
173     if cache_file is not None:
174         save_pkl(dataset, cache_file)
175
176     return dataset
177
178 def read_blogs_xml(path):
179     print('reading all data files from directory {} ...'.format(path))
180     dataset = []
181
182     if _DEBUG: # use small files for fast debugging
183         files = [os.path.join(path, fname) for fname in ['3998465.male.17.indUnk.Gemini.xml',
184                                                         '3949642.male.25.indUnk.Leo.xml', '3924311.male.27.HumanResources.Gemini.xml']]
185         files = random.sample(list(glob(os.path.join(path, '*'))), 100)
186     elif NUM_SAMPLES is None:
187         files = glob(os.path.join(path, '*'))
188     else:
189         files = random.sample(list(glob(os.path.join(path, '*'))), NUM_SAMPLES)
190
191     for fpath in tqdm(files):
192         fname = os.path.basename(fpath)
193         meta_data = parse_meta_data(fname)
194         posts = read_blog_file(fpath)
195         rec = Record(meta_data, posts)
196         dataset.append(rec)
197     return dataset
198

```

```

199 #####
200 #           Codes for topic mining           #
201 #####
202
203 punct_re = re.compile(r'([\.\!?,;:;])(?=[a-zA-Z])') # add space between a punctuation and a word
204 # replace two or more consecutive single quotes to a double quote
205 #   e.g. ''' -> "      ''' -> "
206 quotes_re = re.compile(r"['\"]{2,}")
207 def preprocess(text):
208     out = punct_re.sub(r'\1 ', text)
209     out = quotes_re.sub(r'""', out)
210     out = remove_invalid(out)
211     return out
212
213 leading_quote_re = re.compile(r'[\'\\".\~*^%#|!|-]+([a-zA-Z].*)')
214 def clean_word(word):
215     if word in ("'ve", "'re", "'s", "'t", "'ll", "'m", "'d", "'", "'"):
216         return word
217     word = leading_quote_re.sub(r'\1', word)
218     return word.strip()
219
220 def tokenise(dataset):
221     '''
222     consider all the blogs from one person as a document
223
224     Returns
225     -----
226     docs: list of list of list
227         a list of documents, each of which is a list of sentences,
228         each of which is a list of words.
229     '''
230
231     print('tokenising the text dataset...')
232     docs = []
233     with tqdm(total=sum(len(rec.posts) for rec in dataset)) as pbar:
234         for rec in dataset:
235             doc = []
236             for post in rec.posts:
237                 for sent_str in nltk.sent_tokenize(post.text):
238                     sent_str = preprocess(sent_str)
239                     sent = [clean_word(w) for w in nltk.word_tokenize(sent_str)]
240                     sent = [w for w in sent if w != '']
241                     doc.append(sent)
242             pbar.update(1)
243             docs.append(doc)
244
245     return docs
246
247 def calc_vocab(docs):
248     '''Calculate the vocabulary (set of distinct words) from a collection
249     of documents.
250     '''
251
252     print('calculating the vocabulary...')
253     vocab = set()
254
255     def _helper(sent):
256         vocab.update(sent)
257
258     foreach2d(_helper, docs)
259     return sorted(vocab)
260
261 def calc_pos_tags(docs):
262     print('POS tagging...')
263     def _f(sent):
264         try:
265             return nltk.pos_tag(sent)

```

```

266         except IndexError:
267             print('error sentence: {}'.format(sent))
268             raise
269     tagged_docs = map2d(_f, docs)
270     return tagged_docs
271
272 pattern = re.compile(r'([\^\.])\1{2,}')
273 pattern_ellipse = re.compile(r'\.{4,}')
274 invalid_chars = re.compile(r'[*\^#]')
275 def remove_invalid(text):
276     '''Basic cleaning of words, including:
277
278     1. rip off characters repeated more than twice as English words have a max
279     of two repeated characters.
280     2. remove characters which are not part of English words
281     '''
282
283     text = invalid_chars.sub(' ', text)
284     text = pattern.sub(r'\1\1', text)
285     text = pattern_ellipse.sub('...', text)
286     return text.strip()
287
288 def remove_invalid_all(docs):
289     print('reduce lengthily repeated characters...')
290     return filter3d(lambda w: len(w) > 0, map3d(remove_invalid, docs))
291
292 spell = SpellChecker()
293
294 def correct_spelling(word):
295     if not wordnet.synsets(word) and not word in STOPWORDS:
296         return spell.correction(word)
297     else:
298         return word
299
300 def correct_spelling_all(docs):
301     print('running spelling correction...')
302     return map3d(correct_spelling, docs)
303
304 def remove_stopwords(docs):
305     print('removing stopwords...')
306     return filter3d(lambda wp: wp[0].lower() not in STOPWORDS, docs)
307
308 lemmatizer = WordNetLemmatizer()
309 porter = PorterStemmer()
310 lancaster = LancasterStemmer()
311 def stem_word(word):
312     return porter.stem(lemmatizer.lemmatize(word))
313
314 def do_stemming(docs):
315     print('stemming or lemmatising words...')
316     return map3d(lambda wp: (stem_word(wp[0]), wp[1]), docs)
317
318 def calc_ne_all(docs):
319     print('extracting named entities...')
320     def _calc_ne(sent):
321         ne = []
322         for chunk in nltk.ne_chunk(sent):
323             if hasattr(chunk, 'label'):
324                 ne.append((' '.join(c[0] for c in chunk), chunk.label()))
325         return ne
326     return map2d(_calc_ne, docs)
327
328
329 def calc_df(docs):
330     df = defaultdict(lambda: 0)
331     for doc in docs:
332         for w in set(doc):

```



```

333         df[w] += 1
334     return df
335
336 def calc_tfidf(docs):
337     '''The original TF-IDF is a document-wise score. This function will
338     calculate the average TF-IDF on whole dataset as an overall scoring.
339     '''
340     tf_idf = defaultdict(lambda: 0)
341     df = calc_df(docs)
342     num_docs = len(docs)
343     for doc in docs:
344         counter = Counter(doc)
345         num_words = len(doc)
346         for token in set(doc):
347             tf = counter[token] / num_words
348             df_i = df[token]
349             idf = np.log(num_docs / df_i)
350             tf_idf[token] += tf * idf
351
352     for token in tf_idf:
353         tf_idf[token] /= num_docs
354
355     return tf_idf
356
357 def get_top_topics(named_entities, n=5, method='tf'):
358     print('calculating most popular topics by ' + method + '...')
359     if method == 'tf':
360         ranks = nltk.FreqDist(w for w, t in flatten3d(named_entities))
361         print(ranks.most_common(50))
362         ranks = dict(ranks)
363     elif method == 'tfidf':
364         ranks = calc_tfidf([[w for w, t in flatten2d(doc)] for doc in named_entities])
365         ranks = [(k, v) for k, v in ranks.items()]
366         print('\n largest:', heapq.nlargest(200, ranks, key=itemgetter(1)))
367         topics = heapq.nlargest(n, ranks, key=itemgetter(1))
368         print('topics: ', topics)
369     return topics
370
371 def get_surroundings(words, docs, n=4):
372     '''expand the topic to be 2 verb/noun before and 2 verb/noun after the topic
373     '''
374
375     print('get surrounding 2 nouns/verbs for words {}'.format(words))
376
377     sur = {}
378     for w, c in words:
379         sur[w] = Counter()
380
381     # POS tags list for searching verbs/nouns
382
383     def _helper(sent):
384         sent_w = [w for w, p in sent]
385         for w, c in words:
386             try:
387                 idx = sent_w.index(w)
388                 except ValueError:
389                     continue
390
391                 after = 0
392                 vicinity = [sent[i] for i in [idx-2, idx-1, idx+1, idx+2]
393                             if i >= 0 and i < len(sent)]
394                 for (wi, pi) in vicinity:
395                     if pi.startswith('N') or pi.startswith('V'):
396                         sur[w][wi] += 1
397
398     foreach2d(_helper, docs)
399     ret = []

```

```

400     for w, c in words:
401         ret.append({'topic': w, 'score': c, 'keywords': sur[w].most_common(n)})
402     return ret
403
404 def calc_intermediate_data(dataset):
405     docs = tokenise(dataset)
406     vocab = calc_vocab(docs)
407     print('Size of vocabulary: {}'.format(len(vocab)))
408     print(vocab[1:2000:2])
409     print(vocab[1:100000:100])
410
411
412     tagged_docs = calc_pos_tags(docs)
413     docs = vocab = None
414
415     named_entities = calc_ne_all(tagged_docs)
416
417     # Remove stopwords after POS tagging and NER finished
418     tagged_docs = remove_stopwords(tagged_docs)
419     named_entities = remove_stopwords(named_entities)
420
421     tagged_docs = do_stemming(tagged_docs)
422     named_entities = do_stemming(named_entities)
423     return tagged_docs, named_entities
424
425
426 def mine_topics(dataset, intermediate_data, group='all'):
427     print('-' * 80)
428     print('mining most popular topics for group ' + group)
429     print('-' * 80)
430     tagged_docs, named_entities = intermediate_data
431
432     if group != 'all':
433         if group == 'male' or group == 'female':
434             idx = [i for i, rec in enumerate(dataset) if rec.meta.gender == group]
435         elif group == '<=20':
436             idx = [i for i, rec in enumerate(dataset) if rec.meta.age <= 20]
437         elif group == '>20':
438             idx = [i for i, rec in enumerate(dataset) if rec.meta.age > 20]
439         else:
440             raise NotImplementedError()
441         tagged_docs = [tagged_docs[i] for i in idx]
442         named_entities = [named_entities[i] for i in idx]
443
444     print('selected docs: {}, {}'.format(len(tagged_docs), len(named_entities)))
445
446     ret = {}
447     print('----- result from TFIDF -----')
448     topics = get_top_topics(named_entities, n=50, method='tfidf')
449     keywords = get_surroundings(topics, tagged_docs, n=200)
450     ret['tfidf'] = keywords
451
452     print('----- result from TF -----')
453     topics = get_top_topics(named_entities, n=50, method='tf')
454     keywords = get_surroundings(topics, tagged_docs, n=20)
455     ret['tf'] = keywords
456     return ret
457
458 def main_intermediate():
459     if not _DEBUG and NUM_SAMPLES is None:
460         dataset = read_blogs('blogs')
461     else:
462         dataset = read_blogs('blogs', cache_file=None)
463
464     intermediate_data = calc_intermediate_data(dataset)
465     save_pkl(intermediate_data, 'intermediate_data.pkl')
466     return dataset, intermediate_data

```

```

467
468 def main_mine_topics(dataset=None, intermediate_data=None):
469     if dataset is None:
470         dataset = load_pkl('blogs.pkl')
471     if intermediate_data is None:
472         intermediate_data = load_pkl('intermediate_data.pkl')
473
474     topics = {}
475     topics['male'] = mine_topics(dataset, intermediate_data, group='male')
476     topics['female'] = mine_topics(dataset, intermediate_data, group='female')
477     topics['no_more_than_20'] = mine_topics(dataset, intermediate_data, group='<=20')
478     topics['more_than_20'] = mine_topics(dataset, intermediate_data, group='>20')
479     topics['all'] = mine_topics(dataset, intermediate_data, group='all')
480     if _DEBUG:
481         suffix = 'debug'
482     else:
483         suffix = date.today().strftime('%Y%m%d')
484         if NUM_SAMPLES > 0:
485             suffix += '-' + str(NUM_SAMPLES)
486
487     save_json(topics, 'topics-{}.json'.format(suffix))
488
489 def main():
490     if len(sys.argv) <= 1:
491         phases = [1, 2]
492     else:
493         phases = [int(i) for i in sys.argv[1].split(',')]
494
495     dataset = intermediate_data = None
496     for ph in phases:
497         if ph == 1:
498             dataset, intermediate_data = main_intermediate()
499         elif ph == 2:
500             main_mine_topics(dataset, intermediate_data)
501
502 if __name__ == '__main__':
503     main()

```