# 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 20 years old or younger
- People older than 20
- 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 by Fig. 1, which is created by Python packages `Pandas` and `Matplotlib`. 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) Zodiac: The distribution over zodiac is reasonable even.
4) 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.
5) Category: the most frequent category is unknown, which is trivial, while the second frequent one is student, far more than other categories.

### 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 Fig. 2, and the details of each step are described in the following subsections.

*1) Data Cleaning:* Before applying any text mining techniques, it is important to do basic data cleaning to improve data quality. In this step, a few operations for preprocessing will be carried out based on the observations of the dataset, with details as follows.

- **Problem**: At some place there is no whitespace between a punctuation and the word following it, which causes wrong tokenisation. Specifically, the punctuation might be tokenised with the following word as one token.
  **Solution**: Add whitespace after a punctuation if a word immediately follows it.
- **Problem**: Two or more consecutive quote symbols may cause wrong tokenisation.
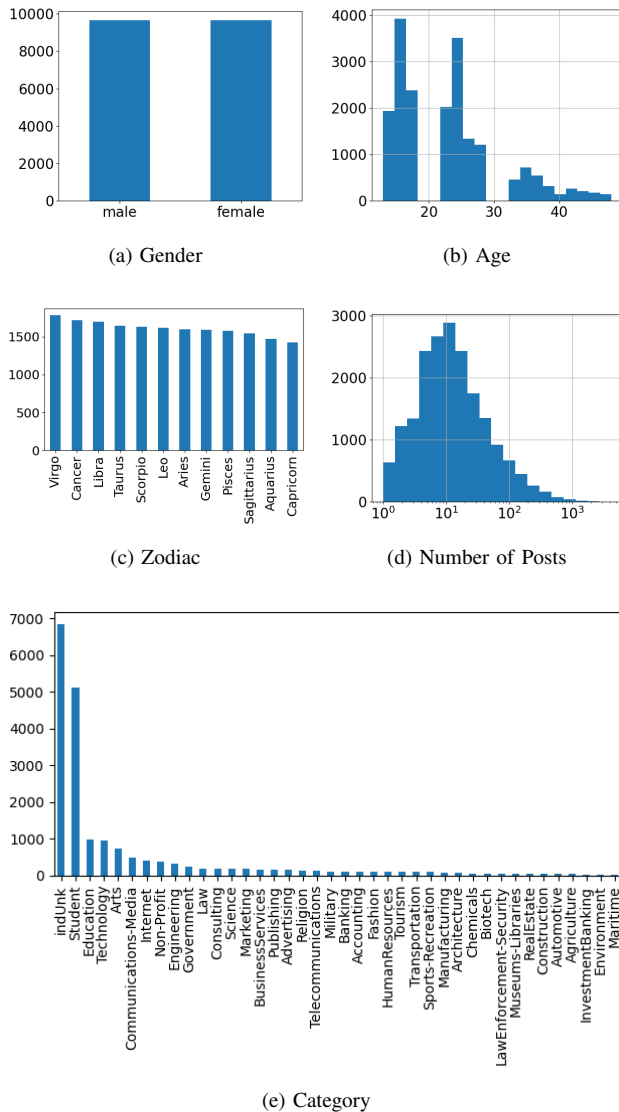  **Solution**: Replace two or more quotes as a double quote.

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

- **Problem**: The unicode quote may affect tokenisation and stopwords matching.
  **Solution**: Replace unicode quote by ASCII quote.
- **Problem**: The unicode quote may affect tokenisation and stopwords matching.
  **Solution**: Replace unicode quote by ASCII quote.
- **Problem**: Characters that are usually not part of normal English text may disturb tokenisation and POS tagging.
  **Solution**: Remove invalid characters such as "*","#", and so on.
- **Problem**: Sometimes people repeat a certain letter in a word for emphasis, but it will result in wrong words and also increase the vocabulary size.
  **Solution**: No English word has more than two consecutive appearances of the same letter, so three or more repetition of a letter is squeezed into two.

- **Problem**:
  **Solution**:

These operations are implemented by regex matching and substitution, or simple text replacing. To use regex, the Python's `re` package are imported.

*2) Tokenisation:* Tokenisation is usually the first step of all text mining pipelines, which includes sentence and word tokenisation. Sentence tokenisation is to split the whole text into sentences, while word tokenisation splits a sentence into word or tokens. In this project we use `nltk` package to do such task. This package provides two functions `sent_tokenize()` and `word_tokenize()` for both tokenisation. A document is first tokenised into sentences, and then each sentence is tokenised into words. Finally, a document is represented as list of lists, as each sentence is a list of words.

*3) POS Tagging:* Part-Of-Speech (POS) tagging is the second step following tokenisation. In this step, each word is assigned by a POS tag. `nltk` provides a handy function `pos_tag()` to do this task. This function works on sentence level, and maps each word into a tuple which is the pair of word and POS tag.

*4) Entity Extraction:* In this project, a topic is defined as a "thing" or "object". Therefore, in order to find the topics, we need to find all "things" or "objects" first. There are a few options to do this task, among which two methods will be employed by this project: Named Entity Recognition (NER) and parsing.

*a) NER:* Named entities are ideal candidates of topics as they denote real-world objects. `nltk` provides a function `ne_chunk()` to extract entities from sentences. The input of the function should be a list of tokens with POS tags, which is another reason why POS tagging should be done in previous step. The return value of this function is a list of chunks, each of which is basically a list and may contain a `label` attribute if it is a recognised entity. The entity type can be acquired by `label()` and the entity itself should be acquired by joining all the elements of the chunk.

*b) Parsing:* Another way to extract objects is parsing by pre-defined patterns. For example, it is reasonable to treat definite nouns as objects according to the grammar.

*5) Stopwords Removal:* Stopwords are most common words which carries no significant meanings. Removing stopwords can reduce the size of data to be proceeded as well as increase the result accuracy. `nltk` provides an out-of-the-box stopwords collection, but the experiment shows that some common words carrying no meaning are not included in the list. In order to expand the stopword list, more words are collected from website [1].

Stopwords removal is conducted after POS tagging and entity extraction because these two steps are sequence model, which means their performance rely on word order. If stopwords are removed before them, we will get sentences which do not comply with English grammar. In addition, stopwords removal are carried out on tagged documents as well as entities
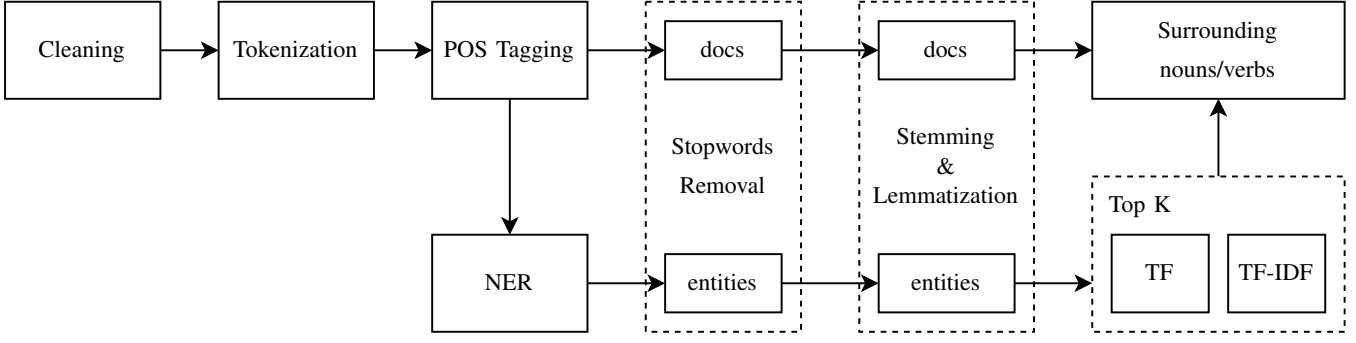
---

[1]https://gist.github.com/sebleier/554280

Fig. 2. Algorithm

extracted. Theoretically, stopwords cannot be entity, but errors will happen in any POS tagging and NER model. Therefore, trying to remove stopwords can reduce the error introduced in previous steps.

*6) Stemming and Lemmatisation:* Stemming and lemmatisation are both techniques for text normalisation, that is, convert an inflected word into its root form. However, stemming and lemmatisation work in different way. Stemming removes suffix or prefix from a word, returning a word stem which is not necessarily a word. On the other hand, lemmatisation always looks for the lemma from word variations with morphological analysis. For example, stemming against the third-person singular form "flies" returns "fli", while lemmatisation returns "fly". In this project, these two methods are combined together to reach the maximum extent of word normalisation.

`nltk` provides various stemming algorithms such as `PorterStemmer` and `LancasterStemmer`, and one lemmatisation algorithm `WordNetLemmatizer`. In the code we use `WordNetLemmatizer` followed by `PorterStemmer`.

*7) Word Count and TF-IDF:* After all "objects" have been extracted and normalised, the next step is to find most popular ones as the most dominant topics. Popularity can be defined in various ways, and in this project two approaches are used: word count and Term Frequency-Inverse Document Frequency (TF-IDF). In the first method, we simply count the appearances of each entity and get the most two frequent ones. In the second method, we calculate the TF-IDF value of each entity word, following the definition

$$\text{TF-IDF}(t_i, d_j) = \text{TF}(t_i, d_j) \times \text{IDF}(t_i)$$
$$= \text{TF}(t_i, d_j) \log \frac{N}{\text{DF}(t_i)}$$

$\text{TF}(t_i, d_j)$ is the Term Frequency of term $t_i$ in document $d_j$, which is computed by count of $t_i$ in $d_j$ divided by the total number of terms in $d_j$. $\text{DF}(t_i)$ is the Document Frequency, which is the number of documents that contains $t_i$. As we can see here, TF-IDF is a term-document-wise number so a term has different TF-IDF values in different documents. In order to rank all terms over the whole dataset, TF-IDF values of a term are averaged over all documents as the score of that term.

$$\text{score}(t_i) = \underset{d_j}{\text{avg}} \ \text{TF-IDF}(t_i, d_j)$$

Two different methods might return different results, which will be compared and analysed in section IV.

### C. Evaluation

Evaluation of topics is challenging due to its nature of unsupervised learning. Among existing metrics, xx is chosen to evaluate the result of the methodology.

## IV. RESULT, ANALYSIS, AND EVALUATION

### A. Result

The results are presented by word cloud. Fig. 3 shows the topics mined by word count while Fig. 4 shows that by TF-IDF.

### B. Analysis

### C. Evaluation

## V. CONCLUSION

This project has designed and implemented a complete solution to mine most popular topics from blogs. A variety of text mining technologies are employed and combined together to reach the goal. The results are compared and evaluated. Further and in-depth discussion is also provided.

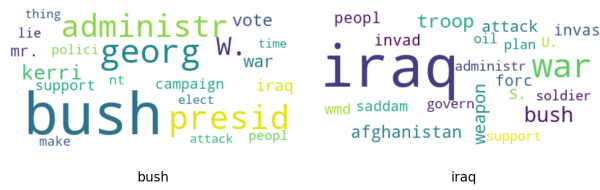## VI. OPEN ISSUES AND FUTURE WORKS

There are still a few open issues remaining in the solution which can be improved by future work or changed if re-do this project.

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn gdwgd fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

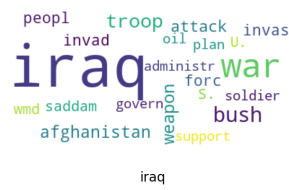fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

Fig. 3. Topics mined by word count. (a) male (b) female (c) 20 or younger (d) over 20 (e) everyone



Fig. 4. Topics mined by TF-IDF. (a) male (b) female (c) 20 or younger (d) over 20 (e) everyone

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

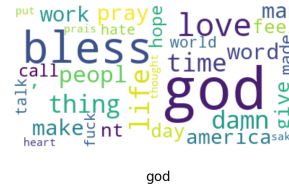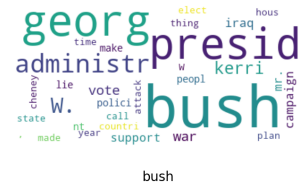fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt

d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn fake wtdwt dwt wd twdt w d dw w

wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

fake wtdwt dwt wd twdt w d dw w wdt wd tw dtdw wd dwt d tagfdygafdgkdwlrn

## 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.

[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, J. Tan, 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 improvementes," *Handbook of Mixed Membership Models and Their Applications*, p. 30,

# APPENDIX
## SOURCE CODE IN PYTHON

### A. Code for topic mining

```python
#!/usr/bin/env python
# -*- coding: utf-8 -*-

import sys
import os.path
from glob import glob
from tqdm import tqdm
import pickle
import json
from datetime import date
import pprint
pp = pprint.PrettyPrinter(indent=2)

import random
import itertools
from collections import namedtuple, Counter, OrderedDict, defaultdict
import heapq
from operator import itemgetter
import re
from bs4 import BeautifulSoup
import numpy as np

from spellchecker import SpellChecker
import nltk
from nltk.corpus import stopwords
from nltk.corpus import wordnet
from nltk.stem import PorterStemmer, LancasterStemmer, WordNetLemmatizer

NUM_SAMPLES = None
_DEBUG = False

_DEBUG = True

STOPWORDS = set(stopwords.words("english"))
# Add more stopwords manually
with open('stopwords1.txt') as f:
    STOPWORDS.update(w.strip().lower() for w in f)
STOPWORDS.update(['i\'m', 'dont', '\'t', '\'m', '\'s', '\'re', '\'ve',
    'haha', 'hah', 'wow', 'hehe', 'heh',
    'ah', 'ahh', 'hm', 'hmm', 'urllink', 'ok', 'hey', 'yay', 'yeah'])

##############################################################################
#                          Utility functions                                 #
##############################################################################

def len2d(iter2d):
    return sum(len(d) for d in iter2d)

def list2d(iter2d):
    return [[x for x in inner] for inner in iter2d]

def flatten2d(list2d):
    return itertools.chain.from_iterable(list2d)

def flatten3d(list3d):
    return itertools.chain.from_iterable(flatten2d(list3d))

def mapbar(f, seq, desc):
    for e in tqdm(seq, desc):
        yield f(e)

def map2d(f, docs):
```

```python
    with tqdm(total=len2d(docs)) as pbar:
        def _helper(sent):
            pbar.update(1)
            return f(sent)

        return [list(map(_helper, doc)) for doc in docs]

def map3d(f, docs):
    with tqdm(total=len2d(docs)) as pbar:
        def _helper(sent):
            pbar.update(1)
            return [f(word) for word in sent]

        return [list(map(_helper, doc)) for doc in docs]

def foreach3d(f, docs):
    with tqdm(total=len2d(docs)) as pbar:
        for doc in docs:
            for sent in doc:
                for word in sent:
                    f(word)
                pbar.update(1)

def foreach2d(f, docs):
    with tqdm(total=len2d(docs)) as pbar:
        for doc in docs:
            for sent in doc:
                f(sent)
                pbar.update(1)

def filter3d(f, docs):
    ret = []
    with tqdm(total=len2d(docs)) as pbar:
        def _helper_doc(doc):
            for sent in doc:
                pbar.update(1)
                out = [word for word in sent if f(word)]
                if len(out) > 0:
                    yield out

        for doc in docs:
            out = list(_helper_doc(doc))
            ret.append(out)
    return ret

def load_pkl(fpath):
    print('load dataset from cached pickle file ' + fpath)
    with open(fpath, 'rb') as f:
        dataset = pickle.load(f)
    return dataset

def save_pkl(obj, fpath):
    with open(fpath, 'wb') as f:
        print('save dataset to pickle file ' + fpath)
        pickle.dump(obj, f)

def save_json(obj, fpath, indent=2):
    with open(fpath, 'w', encoding="utf8") as f:
        print('save dataset to json file ' + fpath)
        json.dump(obj, f, indent=indent)

##############################################################################
#                Codes for data reading & transformation                    #
##############################################################################

Record = namedtuple('Record', ['meta', 'posts'])
Post = namedtuple('Post', ['date', 'text'])
```

```python
MetaData = namedtuple('MetaData', ['id', 'gender', 'age', 'category', 'zodiac'])

def parse_meta_data(meta_data_str):
    arr = meta_data_str.lower().strip().split('.')
    return MetaData(arr[0], arr[1], int(arr[2]), arr[3], arr[4])

def read_blog_file(fpath):
    try:
        with open(fpath, encoding='utf-8', errors='ignore') as f:
            soup = BeautifulSoup(f.read(), "xml")
        blog = soup.Blog
    except ParseError:
        print('Error: invalid xml file {}'.format(fpath))
        raise
        return []

    posts = []
    state = 'date'
    for c in blog.find_all(recursive=False):
        if c.name != state:
            print('Warning: inconsistent format in file {}'.format(fpath))
        if state == 'date':
            try:
                date_str = c.text.strip()
                date = date_str
            except ValueError:
                print('Warning: invalid date {} in file {}' \
                        .format(c.text, fpath))
            state = 'post'
        else:
            text = c.text.strip()
            state = 'date'
            posts.append(Post(date, text))
    posts.sort(key=lambda p: p.date)
    return posts

def read_blogs(path, force=False, cache_file='blogs.pkl'):
    if not force and cache_file is not None and os.path.exists(cache_file):
        return load_pkl(cache_file)

    dataset = read_blogs_xml(path)

    # save to pickle file for fast loading next time
    if cache_file is not None:
        save_pkl(dataset, cache_file)

    return dataset

def read_blogs_xml(path):
    print('reading all data files from directory {} ...'.format(path))
    dataset = []

    if _DEBUG:  # use small files for fast debugging
        files = [os.path.join(path, fname) for fname in ['3998465.male.17.indUnk.Gemini.xml',
            '3949642.male.25.indUnk.Leo.xml', '3924311.male.27.HumanResources.Gemini.xml']]
        files = random.sample(list(glob(os.path.join(path, '*'))), 100)
    elif NUM_SAMPLES is None:
        files = glob(os.path.join(path, '*'))
    else:
        files = random.sample(list(glob(os.path.join(path, '*'))), NUM_SAMPLES)

    for fpath in  tqdm(files):
        fname = os.path.basename(fpath)
        meta_data = parse_meta_data(fname)
        posts = read_blog_file(fpath)
        rec = Record(meta_data, posts)
        dataset.append(rec)
```

```python
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       print(text)
209       out = punct_re.sub(r'\1 ', text)
210       print(out)
211       out = quotes_re.sub(r'"', out)
212       print(out)
213       print(out)
214       out = remove_invalid(out)
215       print(out)
216       return out
217
218   leading_quote_re = re.compile(r'[\'\.~=*&^%#!|\-]+([a-zA-Z].*)')
219   def clean_word(word):
220       if word in ("'ve", "'re", "'s", "'t", "'ll", "'m", "'d", "'", "''"):
221           return word
222       word = leading_quote_re.sub(r'\1', word)
223       return word.strip()
224
225   def tokenise(dataset):
226       '''
227       consider all the blogs from one person as a document
228
229       Returns
230       ---------
231       docs: list of list of list
232           a list of documents, each of which is a list of sentences,
233           each of which is a list of words.
234       '''
235
236       print('tokenising the text dataset...')
237       docs = []
238       with tqdm(total=sum(len(rec.posts) for rec in dataset)) as pbar:
239           for rec in dataset:
240               doc = []
241               for post in rec.posts:
242                   for sent_str in nltk.sent_tokenize(post.text):
243                       sent_str = preprocess(sent_str)
244                       sent = [clean_word(w) for w in nltk.word_tokenize(sent_str)]
245                       sent = [w for w in sent if w != '']
246                       doc.append(sent)
247                   pbar.update(1)
248               docs.append(doc)
249
250       return docs
251
252   def calc_vocab(docs):
253       '''Calculate the vocabulary (set of distinct words) from a collection
254        of documents.
255       '''
256
257       print('calculating the vocabulary...')
258       vocab = set()
259
260       def _helper(sent):
261           vocab.update(sent)
262
263       foreach2d(_helper, docs)
```

```python
264          return sorted(vocab)
265
266  def calc_pos_tags(docs):
267          print('POS tagging...')
268          def _f(sent):
269                  try:
270                          return nltk.pos_tag(sent)
271                  except IndexError:
272                          print('error sentence: {}'.format(sent))
273                          raise
274          tagged_docs = map2d(_f, docs)
275          return tagged_docs
276
277  pattern = re.compile(r'([^\.])\1{2,}')
278  pattern_ellipse = re.compile(r'\.{4,}')
279  invalid_chars = re.compile(r'[*\^#]')
280  def remove_invalid(text):
281          '''Basic cleaning of words, including:
282
283            1. rip off characters repeated more than twice as English words have a max
284               of two repeated characters.
285            2. remove characters which are not part of English words
286          '''
287
288          print(text)
289          text = invalid_chars.sub(' ', text)
290          print(text)
291          text = pattern.sub(r'\1\1', text)
292          print(text)
293          text = pattern_ellipse.sub('...', text)
294          print(text)
295          return text.strip()
296
297  def remove_invalid_all(docs):
298          print('reduce lengthily repreated characters...')
299          return filter3d(lambda w: len(w) > 0, map3d(remove_invalid, docs))
300
301  spell = SpellChecker()
302
303  def correct_spelling(word):
304          if not wordnet.synsets(word) and not word in STOPWORDS:
305                  return spell.correction(word)
306          else:
307                  return word
308
309  def correct_spelling_all(docs):
310          print('running spelling correction...')
311          return map3d(correct_spelling, docs)
312
313  def remove_stopwords(docs):
314          print('removing stopwords...')
315          return filter3d(lambda wp: wp[0].lower() not in STOPWORDS, docs)
316
317  lemmatizer = WordNetLemmatizer()
318  porter = PorterStemmer()
319  lancaster = LancasterStemmer()
320  def stem_word(word):
321          return lemmatizer.lemmatize(word)
322
323  def do_stemming(docs):
324          print('stemming or lemmatising words...')
325          return map3d(lambda wp: (stem_word(wp[0]), wp[1]), docs)
326
327  def calc_ne_all(docs):
328          print('extracting named entities...')
329          def _calc_ne(sent):
330                  ne = []
```

```python
331            for chunk in nltk.ne_chunk(sent):
332                if hasattr(chunk, 'label'):
333                    ne.append((' '.join(c[0] for c in chunk), chunk.label()))
334            return ne
335        return map2d(_calc_ne, docs)
336
337
338    def calc_df(docs):
339        df = defaultdict(lambda: 0)
340        for doc in docs:
341            for w in set(doc):
342                df[w] += 1
343        return df
344
345    def calc_tfidf(docs):
346        '''The original TF-IDF is a document-wise score. This function will
347        calculate the average TF-IDF on whole dataset as an overall scoring.
348        '''
349        tf_idf = defaultdict(lambda: 0)
350        df = calc_df(docs)
351        num_docs = len(docs)
352        for doc in docs:
353            counter = Counter(doc)
354            num_words = len(doc)
355            for token in set(doc):
356                tf = counter[token] / num_words
357                df_i = df[token]
358                idf = np.log(num_docs / df_i)
359                tf_idf[token] += tf * idf
360
361        for token in tf_idf:
362            tf_idf[token] /= num_docs
363
364        return tf_idf
365
366    def get_top_topics(named_entities, n=5, method='tf'):
367        print('calculating most popular topics by ' + method + '...')
368        if method == 'tf':
369            ranks = nltk.FreqDist(w for w, t in flatten3d(named_entities))
370            print(ranks.most_common(50))
371            ranks = dict(ranks)
372        elif method == 'tfidf':
373            ranks = calc_tfidf([[w for w, t in flatten2d(doc)] for doc in named_entities])
374        ranks = [(k, v) for k, v in ranks.items()]
375        print('n largest:', heapq.nlargest(200, ranks, key=itemgetter(1)))
376        topics = heapq.nlargest(n, ranks, key=itemgetter(1))
377        print('topics: ', topics)
378        return topics
379
380    def get_surroundings(words, docs, n=4):
381        '''expand the topic to be 2 verb/noun before and 2 verb/noun after the topic
382        '''
383
384        print('get surrounding 2 nouns/verbs for words {}'.format(words))
385
386        sur = {}
387        for w, c in words:
388            sur[w] = Counter()
389
390        # POS tags list for searching verbs/nouns
391
392        def _helper(sent):
393            sent_w = [w for w, p in sent]
394            for w, c in words:
395                try:
396                    idx = sent_w.index(w)
397                except ValueError:
```

```
398                    continue
399
400                after = 0
401                vicinity = [sent[i] for i in [idx-2, idx-1, idx+1, idx+2]
402                            if i >= 0 and i < len(sent)]
403                for (wi, pi) in vicinity:
404                    if pi.startswith('N') or pi.startswith('V'):
405                        sur[w][wi] += 1
406
407        foreach2d(_helper, docs)
408        ret = []
409        for w, c in words:
410            ret.append({'topic': w, 'score': c, 'keywords': sur[w].most_common(n)})
411        return ret
412
413  def calc_intermediate_data(dataset):
414        docs = tokenise(dataset)
415        vocab = calc_vocab(docs)
416        print('Size of vocabulary: {}'.format(len(vocab)))
417        print(vocab[1:2000:2])
418        print(vocab[1:100000:100])
419
420
421
422        tagged_docs = calc_pos_tags(docs)
423        docs = vocab = None
424
425        named_entities = calc_ne_all(tagged_docs)
426
427        # Remove stopwords after POS tagging and NER finished
428        tagged_docs = remove_stopwords(tagged_docs)
429        named_entities = remove_stopwords(named_entities)
430
431        tagged_docs = do_stemming(tagged_docs)
432        named_entities = do_stemming(named_entities)
433        return tagged_docs, named_entities
434
435  def mine_topics(dataset, intermediate_data, group='all'):
436        print('-' * 80)
437        print('mining most popular topics for group ' + group)
438        print('-' * 80)
439        tagged_docs, named_entities = intermediate_data
440
441        if group != 'all':
442            if group == 'male' or group == 'female':
443                idx = [i for i, rec in enumerate(dataset) if rec.meta.gender == group]
444            elif group == '<=20':
445                idx = [i for i, rec in enumerate(dataset) if rec.meta.age <= 20]
446            elif group == '>20':
447                idx = [i for i, rec in enumerate(dataset) if rec.meta.age > 20]
448            else:
449                raise NotImplementedError()
450            tagged_docs = [tagged_docs[i] for i in idx]
451            named_entities = [named_entities[i] for i in idx]
452
453        print('selected docs: {}, {}'.format(len(tagged_docs), len(named_entities)))
454
455        ret = {}
456        num_keywords = 200
457        print('-------------- result from TFIDF ------------------')
458        topics = get_top_topics(named_entities, n=50, method='tfidf')
459        keywords = get_surroundings(topics, tagged_docs, n=num_keywords)
460        ret['tfidf'] = keywords
461
462        print('-------------- result from TF ------------------')
463        topics = get_top_topics(named_entities, n=50, method='tf')
464        keywords = get_surroundings(topics, tagged_docs, n=num_keywords)
```

```
465      ret['tf'] = keywords
466      return ret
467
468  def main_intermediate():
469      if not _DEBUG and NUM_SAMPLES is None:
470          dataset = read_blogs('blogs')
471      else:
472          dataset = read_blogs('blogs', cache_file=None)
473
474      intermediate_data = calc_intermediate_data(dataset)
475      save_pkl(intermediate_data, 'intermediate_data.pkl')
476      return dataset, intermediate_data
477
478  def main_mine_topics(dataset=None, intermediate_data=None):
479      if dataset is None:
480          dataset = load_pkl('blogs.pkl')
481      if intermediate_data is None:
482          intermediate_data = load_pkl('intermediate_data.pkl')
483
484      topics = {}
485      topics['male'] = mine_topics(dataset, intermediate_data, group='male')
486      topics['female'] = mine_topics(dataset, intermediate_data, group='female')
487      topics['less_or_20'] = mine_topics(dataset, intermediate_data, group='<=20')
488      topics['over_20'] = mine_topics(dataset, intermediate_data, group='>20')
489      topics['all'] = mine_topics(dataset, intermediate_data, group='all')
490      if _DEBUG:
491          suffix = 'debug'
492      else:
493          suffix = date.today().strftime('%Y%m%d')
494          if NUM_SAMPLES > 0:
495              suffix += '-' + str(NUM_SAMPLES)
496
497      save_json(topics, 'topics-{}.json'.format(suffix))
498
499  def main():
500      if len(sys.argv) <= 1:
501          phases = [1, 2]
502      else:
503          phases = [int(i) for i in sys.argv[1].split(',')]
504
505      dataset = intermediate_data = None
506      for ph in phases:
507          if ph == 1:
508              dataset, intermediate_data = main_intermediate()
509          elif ph == 2:
510              main_mine_topics(dataset, intermediate_data)
511
512  if __name__ == '__main__':
513      main()
```

*B. Code for analysis, evaluation and visualisation*

```
1   #!/usr/bin/env python
2   # -*- coding: utf-8 -*-
3
4   import sys
5   import json
6   import numpy as np
7   import pandas as pd
8   from as2 import load_pkl, Record, MetaData, Post
9   import matplotlib.pyplot as plt
10  from wordcloud import WordCloud
11
12
13  def show_summary(dataset):
14      '''This function describes the summary of dataset or human inspection.
15      It's not necessary for the mining process.
```

```
16
17      Parameters
18      --------------
19      dataset : list of Record
20          The blog dataset
21      '''
22
23      df = pd.DataFrame([d.meta for d in dataset])
24      df['blog_count'] = [len(d.posts) for d in dataset]
25      df['char_count'] = [sum(len(p.text) for p in d.posts) for d in dataset]
26
27      print(df.describe(include='all'))
28      print('{} possible values for "gender": {}'.format(
29              len(df.gender.unique()), ', '.join(sorted(df.gender.unique()))))
30      print('{} possible values for category: {}'.format(
31              len(df.category.unique()), ', '.join(sorted(df.category.unique()))))
32      print('{} possible values for zodiac: {}'.format(
33              len(df.zodiac.unique()), ', '.join(sorted(df.zodiac.unique()))))
34
35      plt.rcParams.update({'font.size': 20})
36      df['gender'].value_counts().plot(kind='bar')
37      plt.xticks(rotation=0)
38      plt.gcf().tight_layout()
39      plt.savefig('img/show-gender.png')
40
41      plt.rcParams.update({'font.size': 10})
42      plt.clf()
43      df['category'].value_counts().plot(kind='bar')
44      plt.gcf().tight_layout()
45      plt.savefig('img/show-category.png')
46
47      plt.rcParams.update({'font.size': 18})
48      plt.clf()
49      df['zodiac'].value_counts().plot(kind='bar')
50      plt.xticks(rotation=90)
51      plt.gcf().tight_layout()
52      plt.savefig('img/show-zodiac.png')
53
54      plt.rcParams.update({'font.size': 20})
55      plt.clf()
56      age = df['age']
57      df['age'].hist(bins=20)
58      plt.gcf().tight_layout()
59      plt.savefig('img/show-age.png')
60
61      plt.clf()
62      cnt = df['blog_count']
63      logbins = np.logspace(np.log10(cnt.min()),np.log10(cnt.max()), 20)
64      cnt.hist(bins=logbins)
65      plt.xscale('log')
66      plt.gcf().tight_layout()
67      plt.savefig('img/show-blog-count.png')
68
69      plt.clf()
70      cnt = df['char_count']
71      logbins = np.logspace(np.log10(cnt.min()),np.log10(cnt.max()), 20)
72      cnt.hist(bins=logbins)
73      plt.xscale('log')
74      plt.gcf().tight_layout()
75      plt.savefig('img/show-char-count.png')
76
77  def eval_topics(fpath, method='tf', top_k=2, num_words_in_topic=30):
78      with open(fpath, encoding='utf8') as f:
79          result = json.load(f)
80
81      for group, topics2 in result.items():
82          topics = topics2[method]
```

```python
        for i, topic in enumerate(topics[:top_k]):
            topic_name = topic['topic']
            words = {}
            words.update(tuple(kw) for kw in topic['keywords'][:num_words_in_topic+1])
            if method == 'tf':
                words[topic_name] = topic['score']
            else:
                words[topic_name] = topic['keywords'][0][1] * 2  # fake frequency for display

            print('topic: ', topic_name, 'number of keywords:', len(topic['keywords']))
            wc = WordCloud(background_color="white", max_font_size=80,
                    max_words=num_words_in_topic+1)
            wc.generate_from_frequencies(words)

            plt.clf()
            plt.imshow(wc, interpolation="bilinear")
            plt.axis("off")
            plt.title(topic_name, y=-0.25, fontsize=20)
            plt.gcf().tight_layout()
            fig_path = 'img/{}-{}-{}.png'.format(group, method, i+1, topic['topic'])
            print('drawing ' + fig_path)
            plt.savefig(fig_path)

def main():
    cmd = sys.argv[1]
    if cmd == 'show':
        show_summary(load_pkl('blogs.pkl'))
    elif cmd == 'eval':
        fpath = sys.argv[2]
        eval_topics(fpath, top_k=2)
        eval_topics(fpath, top_k=2, method='tfidf')

if __name__ == '__main__':
    main()
```