

Popular Topic Mining from Blog Text

Stone Fang (Student ID: 19049042)
Computers and Information Sciences
Auckland University of Technology
Auckland, New Zealand
fkn7060@autuni.ac.nz

Abstract—Topic mining is an important way that provides insights of text data about opinions and preference of people. In this project, a complete topic mining solution is conducted, extracting two most popular topics among 19,320 blog posts grouped by the demographic of authors. Topics are generated from ranking objects that are mentioned in the dataset, and two ranking methods are used, one using word-count and the other based on TF-IDF. The result shows most topics are meaningful and coherent, but the topic from two methods are very different. The quality of mined topics are analysed and further discussions are presented. Meanwhile, some topics have relatively lower quality and the causes are analysed. In addition, with open issues and challenges pointed out, future works and possible improvements are provided.

Index Terms—Topic mining, TF-IDF, Named entity recognition, Topic coherence

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 ?? 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. LITERATURE REVIEW

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. This work also included a case study of analysis on nuclear technology presented in New York Times from 1942. Waila, Singh, and Singh [2] conducted an analysis on blog text with combination of topic modelling, NER, and sentiment analysis. Key themes are extracted from the dataset, and topic perplexity is calculated for evaluation. In addition, important entities such as person, place, and organisation are also recognised and displayed by word cloud. Guo, Zhang, Tan, *et al.* [3] proposed a Frequent Pattern stream (FP-stream) mining algorithm for hot topic discovery on Twitter. The author argues that traditional clustering-based topic detection algorithms are not suitable for the short, sparse, and fast-spreading Twitter data. Experiments were carried out and the result showed the hot topics and the trend of change over time.

Boyd-Graber, Mimno, and Newman [4] provides a summary of topic evaluation methods, which are divided into three categories: human evaluation, diagnostic metrics, and coherent metrics. The first one needs human effort so it is expensive and time-consuming, while the other two can be calculated by computer without human interference.

Human Evaluation requires human involvement in the evaluation task. One method in this category is accomplished by word intrusion task. Specifically, a person will be presented by a list of words and is asked to find an intruder in the meaning of not belonging to others. The words list are constructed by first selecting highly possible words from a topic, and then randomly choose one word with low probability in the same topic but high probability in a different topic. If the intruders are easily to be identified, then the topic is more likely coherent [4].

Diagnostic Metrics only compute statistics of topics without requirements of external knowledge source. Some methods in this category are [4]:

- **Topic Size:** measured by the sum of numbers of tokens belonging to a certain topic. Generally speaking, small topic size means low quality.

- Word Length: average length of N most dominant words in a topic. The usefulness of this metric is corpus dependent.
- Corpus Distribution Distance: A probability distribution can be derived from a topic over the vocabulary, and further normalised by global word count in the whole dataset. The distance between different topics reflects how much these topics are separated.

Coherence Metrics is a type of methods which automatically compute score of topic coherence, and their accuracy is close to human performance. The basic idea is measuring how a pair of words from top N dominant words are associated [4]. It is formalised as

$$TC-f(\mathbf{w}) = \sum_{i < j} f(w_i, w_j), i, j \in \{1 \dots N\}$$

where $\mathbf{w} = \{w_1, w_2, \dots, w_N\}$ is the list of N most dominant words, and f is the scoring function of association between two words. A typical value of N is 10. There are a variety of ways to compute f , such as counting the co-occurrence of two words, or counting the number of documents containing both words. Two popular implementations of f are Pointwise Mutual Information (PMI) and Log Conditional Probability (LCP) [4].

$$PMI(w_i, w_j) = \log \frac{P(w_i, w_j)}{P(w_i)P(w_j)}$$

$$LCP(w_i, w_j) = \log \frac{P(w_i, w_j)}{P(w_j)}$$

Fang, Macdonald, Ounis, *et al.* [5] proposed a metric based on Word Embeddings (WE). Similar to Latent Semantic Analysis (LSA) which represents each word as a real number vector obtained by Singular Value Decomposition (SVD), the WE metric represents each word by a WE vector, and the association scoring function is defined by cosine similarity of two vectors:

$$f(w_i, w_j) = \text{cosine}(V(w_i), V(w_j))$$

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) Gender and age group: since we divide the dataset by the author's age, it is beneficial to inspect some data distribution over age groups. Therefor, we plot the histogram for gender distribution for each age group and the data amounts for both genders in either age group are shown to be evenly distributed.
- 4) Zodiac: The distribution over zodiac is reasonable even.
- 5) Category: the most frequent category is unknown, which is trivial, while the second frequent one is student, far more than other categories.
- 6) 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.

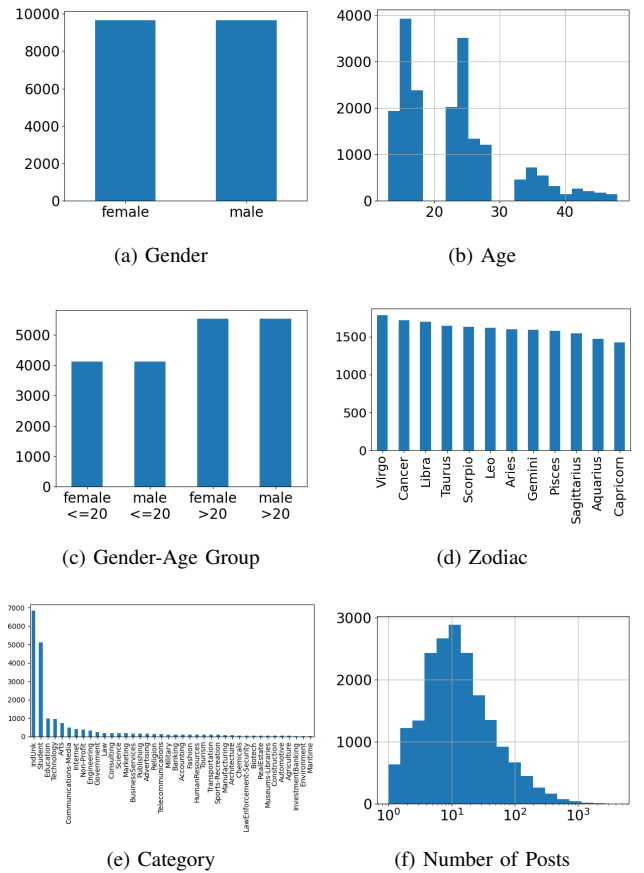


Fig. 1. Data Overview. Histogram over (a) gender (b) age (c) gender and age group (d) zodiac (e) category (f) number of posts

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 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.

Example:

```
- input:  I brought...stuff...like clothes
- output: I brought... stuff... like clothes
```

- **Problem**: Two or more consecutive quote symbols may cause wrong tokenisation.

Solution: Replace two or more quotes as a double quote.

- **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.

Example:

```
- input:  going to be in deeeeeeep
- output: going to be in deep
```

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, which can be extracted by matching of pattern “the + NOUN”. However, due to limited time, this technique is not used in the experiment and might be considered in the future.

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 stopwords 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 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 [6]. For example, stemming against the third-person singular form “flies” returns “fli”, while lemmatisation returns “fly”. In this project, these two methods

¹<https://gist.github.com/sebleier/554280>

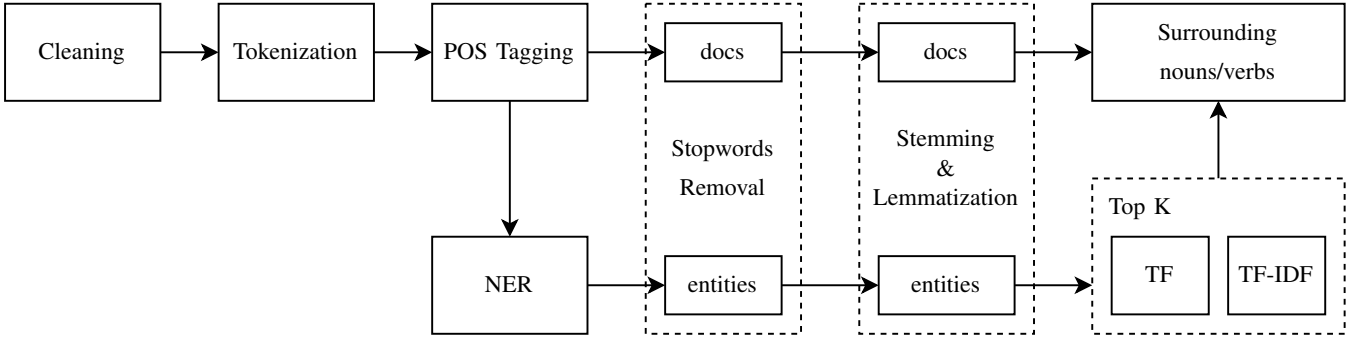


Fig. 2. Overall Architecture of Topic Mining Algorithm

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 equation [7]

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

$\text{TF}(t_i, d_j)$ 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 . In order to avoid large IDF value for some terms only appearing in a few documents, penalty is given to small DF values by adding a smooth constant C_1 to DF

$$\overline{\text{TF-IDF}}(t_i, d_j) = \text{TF}(t_i, d_j) \log \frac{N}{\text{DF}(t_i) + C_1}$$

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. Similarly, penalty is also given to terms having small DF values during the averaging process because terms appears in too few documents are hardly to be defined as popular. It is implemented by another smoothing constant C_2 which is not necessarily equal to C_1 . The effect of the penalty will be further discussed in section IV-B2.

$$\text{score}(t_i) = \frac{1}{\text{DF}(t_i) + C_2} \sum_{d_j} \overline{\text{TF-IDF}}(t_i, d_j)$$

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

C. Evaluation Method

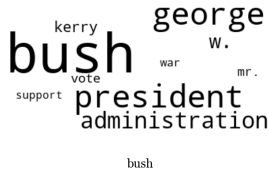
Evaluation of topics is challenging due to its nature of unsupervised learning. There are several topic evaluation approaches, including human evaluation, diagnostic metrics, topic coherence, evaluation on classification corpus. Human evaluation requires manual works so it is costly and slow. Evaluation on classification corpus is a method using labelled text classification corpus and checking the consistency between topic words and class labels. Although we have categories in the metadata for each document, a large portion of categories are labelled as “unknown”, so this approach is not suitable for this project. Therefore, diagnostic and coherence metrics are chosen to evaluate the result of the methodology. Details will be provided in section IV

IV. RESULT, EVALUATION AND ANALYSIS

This section will show the results generated from the methodology described above, and provides evaluation and discussion of how good the topics are. Two hottest topics are extracted for each demographic, and each topic contains 10 words. In the TF-IDF-based method, penalty parameters are set to $C_1 = C_2 = 100$. Due to the large volume of the original dataset, the experiments were conducted with 5,000 and 10,000 documents randomly sampled out of the total 19,320 ones, and the results demonstrated consistency. Therefore, in the following part of this section, only results from 10,000 samples are presented. The complete implementation is listed in appendix A, including two files, one for topic mining algorithm and the other for result evaluation and visualisation.

A. Result and Evaluation

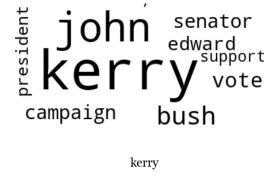
The results are displayed as word cloud generated by `wordcloud` and `matplotlib` package. Fig. 3 shows the topics mined by word count while Fig. 4 shows that by TF-IDF.



(a) Male



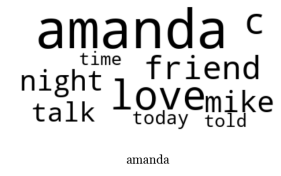
(b) Female



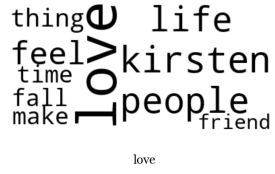
(a) Male



(b) Female



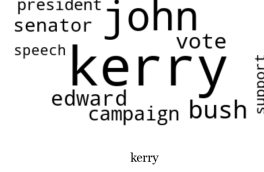
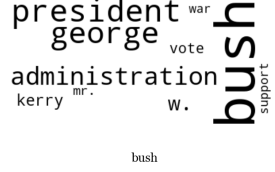
(c) 20 or younger



(c) 20 or younger



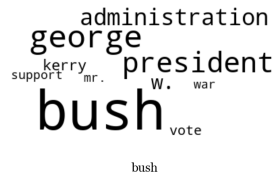
(d) Over 20



(d) Over 20



(e) Everyone



(e) Everyone



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

For evaluation of the quality of topics, we computed two diagnostic metrics, topic size and word length, and two coherence metrics, LSA and Word Embedding. The LSA vectors are provided by [8]. The pre-trained GloVe [9] word vectors with 100 dimension size are used as word embeddings. The metrics are summarised in table I and table II for two methods, respectively.

B. Analysis

In this section we will first conduct in-depth analysis on topics from both methods, and then compare these two results and discuss the pros and cons for each method.

1) *Method of Word Count*: Most topics are meaningful and coherent. For instance, the topic “bush” is about U.S. president George W. Bush and relative concepts, such as his rival John Kerry in 2004 U.S. president election and the Iraq War began in 2003. The word length is relatively high and both LSA and GloVe coherences are also high. However, the topic size is low and the reason could be that it is specific on politics. Another example is “god” which appears in multiple demographics. Though the word length is low, the topic size and coherence are high. We can also see that the keywords in this topic are meaningful and relevant, such as “love”, “life”, “people”, “bless”, and so forth. On the other hand, the topics “love” and

TABLE I
METRICS FOR WORD COUNT

Demographic	Topic	Diagnostic		Coherence	
		Topic Size	Word Length	LSA	GloVe
male	bush	68k	5.7	0.189	0.541
	american	174k	5.4	0.154	0.424
female	god	792k	4.3	0.193	0.502
	good	1,223k	4.4	0.137	0.637
20 or younger	god	816k	4.3	0.199	0.508
	love	949k	4.8	0.177	0.485
over 20	god	897k	4.0	0.201	0.492
	bush	68k	5.7	0.189	0.541
everyone	god	792k	4.0	0.209	0.460
	bush	68k	5.7	0.189	0.541

TABLE II
METRICS FOR TF-IDF

Demographic	Topic	Diagnostic		Coherence	
		Topic Size	Word Length	LSA	GloVe
male	kerry	173k	5.5	0.214	0.473
	gmail	83k	5.5	0.087	0.346
female	josh	798k	4.6	0.119	0.294
	amanda	734k	4.3	0.151	0.435
20 or younger	singapore	391k	6.2	0.116	0.416
	sarah	633k	5.2	0.145	0.416
over 20	kerry	66k	6.0	0.203	0.523
	india	372k	5.5	0.156	0.379
everyone	singapore	587k	5.9	0.120	0.466
	kyle	771k	4.6	0.115	0.479

“good” are too generic, thus have lower coherence scores.

2) *Method of TF-IDF*: The topics ranked by TF-IDF are different from word count. A topic relevant to the result from word count is “kerry”, which we can tell from the keywords represents U.S. Senator John Kerry and also the Democratic nominee for presidency competing with George W. Bush. The coherence of this topic is also high. Several other topics are about artists, singers or movie stars, but the qualities vary. For example, the topic “sarah” contains Sarah McLachlan. By contrast, the topic “josh” actually contains two people, Josh Groban and Josh Hartnett, so it should be two topics. Accordingly, the coherence for topic “josh” is very low. However, some topics are hard to tell what they represent, such as “amanda” and “kyle”, which are probably just common names in English world.

Two other meaningful topics are “singapore” and “india”, which are likely to be created by users from these two countries. An evidence of the significant number of Singaporean users is that the word “haiz” as a common saying of “sigh” in Singaporean English came out as a high ranked topic candidate before it was filtered out as a stopword. The “india” topic is more informative, which is telling us about the important job

outsourcing in India which is consistent to our knowledge. However, the coherence metrics are not significantly high, which could be distracted by not-so-relevant words such as “today”, “time”, “day”. The “india” topic has less distracting words and higher LSA coherence score, but the GloVe score is low, which is a possible issue of WE coherence metrics.

The topic “gmail” is a different case. From human’s point of view, it is a meaningful topic representing the launch of beta release of Google’s mail service at 1 April, 2004. However, both coherence scores are quite low, which may implies some drawbacks in the metrics definition.

3) *Comparison*: As can be seen from above results and analysis, the topics mined by two methods are highly different and share little in common. The difference comes from difference natures of these two methods.

The word-count-based method takes the “thing” mentioned the most times as the most popular one. This approach is simple and straightforward, and consistent to humans’ concept of popularity. The more widely spreaded across the dataset an object is, the more likely it is selected as the hottest topic. That is why we got very general topics such as “god”, “good”, “love”, as well as “bush” during the period of U.S. president election.

On the other hand, the TF-IDF-based method gives larger score on less frequent terms. However, if some term is mentioned too few times, it cannot be a hot topic. Therefore, penalty is given to such terms to make them lower ranked, leading to a trade-off between specialty and frequency and the choice of parameter representing the penalty is more a trick. Experimental shows that if the penalty is too high, the result from TF-IDF is quite similar to that from word count; but if it is too small, result is more likely to be some concept only mentioned in a few document. Generally speaking, this approach prefers topics that are frequently mentioned in a certain subgroup of people, for example, pop and movie stars, and people from a specific country such as Singapore and India.

Each method has its own advantages and disadvantages, but some topics from word-count-based method may be less valuable because they are too general. For instance, the topic “good” or “love” hardly provides any value to marketing or product decision-making because it is something that people are always talking about. By contrast, though the result from TF-IDF-based method may be only popular in a small group of people, it does provide insights into the opinion of that group. However, some topic from the first method is also valuable, such as the topic “bush” capturing people and events related to then U.S. president George W. Bush.

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. In order to rank topics from a candidate list formed by objects, two methods are used: word-count-based and TF-IDF-based. The results are displayed by word cloud

and metrics are computed and summarised. The results are also discussed, compared and evaluated in detail, and the rationale is given. Further and in-depth analysis 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.

An important concern is the quality of NER which has big influence on the final results. The NER result is used for topic candidate generating and wrong entities will cause wrong topics. The experiment showed that the quality of NER is moderate and needs improvement. To do that we can use other NER tools such as `spaCy` and combine the results for better performance.

Data cleaning is also an important stage in the whole algorithm. The data from internet is quite noisy. Though we employed several pre-processing to clean it, there are still many noises. For example, the advertisement in blog posts is a great interferential factor because the same advertisement appears in the text very frequently thus looks like a popular topic.

Spell checking is another means that can be used as an improvement because blog text as a kind of informal text often has many spelling errors which makes it more difficulty to do POS tagging and NER correctly. There are a few spell checking packages in Python, but they require very long time to finish large dataset, so they are not included in the current methodology. Therefore, faster and more scalable spell checking algorithms is another direction of future work.

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 improvements," *Handbook of Mixed Membership Models and Their Applications*, p. 30,
- [5] A. Fang, C. Macdonald, I. Ounis, and P. Habel, "Using word embedding to evaluate the coherence of topics from twitter data," in *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval - SIGIR '16*, Pisa, Italy: ACM Press, 2016, pp. 1057–1060, ISBN: 978-1-4503-4069-4. DOI: 10.1145/2911451.2914729.
- [6] H. Jabeen. (Oct. 23, 2018). Stemming and lemmatization in python, DataCamp Community, [Online]. Available: <https://www.datacamp.com/community/tutorials/stemming-lemmatization-python> (visited on 06/15/2020).
- [7] W. Scott. (May 21, 2019). TF-IDF for document ranking from scratch in python on real world dataset, Medium, [Online]. Available: <https://towardsdatascience.com/tf-idf-for-document-ranking-from-scratch-in-python-on-real-world-dataset-796d339a4089> (visited on 06/23/2020).
- [8] V. Rus, M. Lintean, R. Banjade, N. Niraula, and D. Stefanescu, "SEMILAR: The semantic similarity toolkit," *Proceedings of the 51st annual meeting of the association for computational linguistics: system demonstrations*, pp. 163–168,
- [9] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–1543. [Online]. Available: <http://www.aclweb.org/anthology/D14-1162>.

APPENDIX

SOURCE CODE IN PYTHON

A. Code for topic mining

```
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 import nltk
24 from nltk.corpus import stopwords
25 from nltk.corpus import wordnet
26 from nltk.stem import PorterStemmer, LancasterStemmer, WordNetLemmatizer
27
28 # 10000 documents are sampled from the whole dataset due to limited computing capacity.
29 _DEBUG = False
30 NUM_SAMPLES = 10000
31
32
33 STOPWORDS = set(stopwords.words("english"))
34 # Add more stopwords manually
35 with open('stopwords1.txt') as f:
36     STOPWORDS.update(w.strip().lower() for w in f)
37 STOPWORDS.update(['i\'m', 'dont', '\t', '\m', '\s', '\re', '\ve',
38     'haha', 'hah', 'wow', 'hehe', 'heh',
39     'ah', 'ahh', 'hm', 'hmm', 'urllink', 'ok', 'hey', 'yay', 'yeah'])
40
41 #####
42 #                               Utility functions                               #
43 #####
44
45 def len2d(iter2d):
46     return sum(len(d) for d in iter2d)
47
48 def list2d(iter2d):
49     return [[x for x in inner] for inner in iter2d]
50
51 def flatten2d(list2d):
52     return itertools.chain.from_iterable(list2d)
53
54 def flatten3d(list3d):
55     return itertools.chain.from_iterable(flatten2d(list3d))
56
57 def mapbar(f, seq, desc):
58     for e in tqdm(seq, desc):
59         yield f(e)
60
61 def map2d(f, docs):
62     '''Map a list of lists to a new list of lists by applying f to the elements
63     of the inner lists. Usually works on the sentence-level
64     '''
65
66     with tqdm(total=len2d(docs)) as pbar:
67         def _helper(sent):
68             pbar.update(1)
69             return f(sent)
```



```

70         return [list(map(_helper, doc)) for doc in docs]
71
72
73 def map3d(f, docs):
74     '''Map a list of lists of lists to a new list of lists of lists
75     by applying f to the elements of the inner lists.
76     Usually works on the token-level.
77     '''
78
79     with tqdm(total=len2d(docs)) as pbar:
80         def _helper(sent):
81             pbar.update(1)
82             return [f(word) for word in sent]
83
84         return [list(map(_helper, doc)) for doc in docs]
85
86 def foreach3d(f, docs):
87     with tqdm(total=len2d(docs)) as pbar:
88         for doc in docs:
89             for sent in doc:
90                 for word in sent:
91                     f(word)
92                 pbar.update(1)
93
94 def foreach2d(f, docs):
95     with tqdm(total=len2d(docs)) as pbar:
96         for doc in docs:
97             for sent in doc:
98                 f(sent)
99             pbar.update(1)
100
101 def filter3d(f, docs):
102     ret = []
103     with tqdm(total=len2d(docs)) as pbar:
104         def _helper_doc(doc):
105             for sent in doc:
106                 pbar.update(1)
107                 out = [word for word in sent if f(word)]
108                 if len(out) > 0:
109                     yield out
110
111         for doc in docs:
112             out = list(_helper_doc(doc))
113             ret.append(out)
114     return ret
115
116 def load_pkl(fpath):
117     print('load dataset from cached pickle file ' + fpath)
118     with open(fpath, 'rb') as f:
119         dataset = pickle.load(f)
120     return dataset
121
122 def save_pkl(obj, fpath):
123     if _DEBUG:
124         fpath = fpath.replace('.pkl', '-debug.pkl')
125     with open(fpath, 'wb') as f:
126         print('save dataset to pickle file ' + fpath)
127         pickle.dump(obj, f)
128
129 def save_json(obj, fpath, indent=2):
130     with open(fpath, 'w', encoding="utf8") as f:
131         print('save dataset to json file ' + fpath)
132         json.dump(obj, f, indent=indent)
133
134 #####
135 #                               Codes for data reading & transformation                               #
136 #####
137
138 Record = namedtuple('Record', ['meta', 'posts'])
139 Post = namedtuple('Post', ['date', 'text'])
140 MetaData = namedtuple('MetaData', ['id', 'gender', 'age', 'category', 'zodiac'])
141
142 def parse_meta_data(meta_data_str):
143     arr = meta_data_str.lower().strip().split('.')

```

```

144     return MetaData(arr[0], arr[1], int(arr[2]), arr[3], arr[4])
145
146 def read_blog_file(fpath):
147     '''Read one XML file and extract texts
148     '''
149
150     try:
151         with open(fpath, encoding='utf-8', errors='ignore') as f:
152             soup = BeautifulSoup(f.read(), "xml")
153             blog = soup.Blog
154     except:
155         print('Error: invalid xml file {}'.format(fpath))
156         return []
157
158     posts = []
159     state = 'date'
160     for c in blog.find_all(recursive=False):
161         if c.name != state:
162             # The <date> and <text> tags should appear alternatively
163             print('Warning: inconsistent format in file {}'.format(fpath))
164         if state == 'date':
165             try:
166                 date_str = c.text.strip()
167                 date = date_str
168             except ValueError:
169                 print('Warning: invalid date {} in file {}' \
170                       .format(c.text, fpath))
171             state = 'post'
172         else:
173             text = c.text.strip()
174             state = 'date'
175             posts.append(Post(date, text))
176     posts.sort(key=lambda p: p.date)
177     return posts
178
179 def read_blogs(path, force=False, cache_file='blogs.pkl'):
180
181     dataset = read_blogs_xml(path)
182
183     # save to pickle file for fast loading next time
184     if cache_file is not None:
185         save_pkl(dataset, cache_file)
186
187     return dataset
188
189 def read_blogs_xml(path):
190     '''Read all blog files from the data directory
191     '''
192
193     print('reading all data files from directory {} ...'.format(path))
194     dataset = []
195
196     if _DEBUG: # use small files for fast debugging
197         files = [os.path.join(path, fname) for fname in ['3998465.male.17.indUnk.Gemini.xml',
198                                                         '3949642.male.25.indUnk.Leo.xml', '3924311.male.27.HumanResources.Gemini.xml']]
199     elif NUM_SAMPLES is None:
200         files = glob(os.path.join(path, '*'))
201     else:
202         files = random.sample(list(glob(os.path.join(path, '*'))), NUM_SAMPLES)
203
204     for fpath in tqdm(files):
205         fname = os.path.basename(fpath)
206         meta_data = parse_meta_data(fname)
207         posts = read_blog_file(fpath)
208         rec = Record(meta_data, posts)
209         dataset.append(rec)
210     return dataset
211
212 #####
213 #                               Codes for topic mining                               #
214 #####
215
216 punct_re = re.compile(r'([\.\!?,;:;])?(?=[a-zA-Z])') # add space between a punctuation and a word
217 # replace two or more consecutive single quotes to a double quote

```

```

218 # e.g. '' -> "      ''' -> "
219 quotes_re = re.compile(r"['\"]{2,}")
220 def preprocess(text):
221     '''Text pre-processing, removing invalid characters by regex substitution
222     '''
223
224     out = punct_re.sub(r'\1 ', text)
225     out = quotes_re.sub(r'""', out)
226     out = remove_invalid(out)
227     return out
228
229 leading_quote_re = re.compile(r'[\'\\".\~*%^%#|!|-]+([a-zA-Z].*)')
230 def clean_word(word):
231     if word in ("ve", "re", "s", "t", "ll", "m", "d", "", "'"):
232         return word
233     word = leading_quote_re.sub(r'\1', word)
234     return word.strip()
235
236 def tokenise(dataset):
237     '''
238     consider all the blogs from one person as a document
239
240     Returns
241     -----
242     docs: list of list of list
243           a list of documents, each of which is a list of sentences,
244           each of which is a list of words.
245     '''
246
247     print('tokenising the text dataset...')
248     docs = []
249     with tqdm(total=sum(len(rec.posts) for rec in dataset)) as pbar:
250         for rec in dataset:
251             doc = []
252             for post in rec.posts:
253                 for sent_str in nltk.sent_tokenize(post.text):
254                     sent_str = preprocess(sent_str)
255                     sent = [clean_word(w) for w in nltk.word_tokenize(sent_str)]
256                     sent = [w for w in sent if w != '']
257                     doc.append(sent)
258                     pbar.update(1)
259             docs.append(doc)
260
261     return docs
262
263 def calc_vocab(docs):
264     '''Calculate the vocabulary (set of distinct words) from a collection
265     of documents.
266     '''
267
268     print('calculating the vocabulary...')
269     vocab = set()
270
271     def _helper(sent):
272         vocab.update(sent)
273
274     foreach2d(_helper, docs)
275     return sorted(vocab)
276
277 def calc_pos_tags(docs):
278     print('POS tagging...')
279     def _f(sent):
280         try:
281             return nltk.pos_tag(sent)
282         except IndexError:
283             print('error sentence: {}'.format(sent))
284             raise
285     tagged_docs = map2d(_f, docs)
286     return tagged_docs
287
288 pattern = re.compile(r'([^\.])(\1{2,})')
289 pattern_ellipse = re.compile(r'\.{4,}')
290 invalid_chars = re.compile(r'[*\^#]')
291 def remove_invalid(text):

```

```

292 '''Basic cleaning of words, including:
293
294 1. rip off characters repeated more than twice as English words have a max
295    of two repeated characters.
296 2. remove characters which are not part of English words
297 '''
298
299 text = invalid_chars.sub(' ', text)
300 text = pattern.sub(r'\1\1', text)
301 text = pattern_ellipse.sub('...', text)
302 return text.strip()
303
304 def remove_invalid_all(docs):
305     print('reduce lengthily repeated characters...')
306     return filter3d(lambda w: len(w) > 0, map3d(remove_invalid, docs))
307
308
309 def correct_spelling(word):
310     spell = SpellChecker()
311     if not wordnet.synsets(word) and not word in STOPWORDS:
312         return spell.correction(word)
313     else:
314         return word
315
316 def correct_spelling_all(docs):
317     print('running spelling correction...')
318     return map3d(correct_spelling, docs)
319
320 def remove_stopwords(docs):
321     print('removing stopwords...')
322     return filter3d(lambda wp: wp[0].lower() not in STOPWORDS, docs)
323
324 lemmatizer = WordNetLemmatizer()
325 porter = PorterStemmer()
326 lancaster = LancasterStemmer()
327 def stem_word(word):
328     '''Do stemming followed by lemmatisation for maximum normalisation
329     '''
330
331     return porter.stem(lemmatizer.lemmatize(word))
332
333 def do_stemming(docs):
334     print('stemming or lemmatising words...')
335     return map3d(lambda wp: (stem_word(wp[0]), wp[1]), docs)
336
337 def calc_ne_all(docs):
338     print('extracting named entities...')
339     def _calc_ne(sent):
340         ne = []
341         for chunk in nltk.ne_chunk(sent):
342             if hasattr(chunk, 'label'):
343                 ne.append((' '.join(c[0] for c in chunk), chunk.label()))
344         return ne
345     return map2d(_calc_ne, docs)
346
347
348 def calc_df(docs):
349     '''Calculate the Document Frequency (DF) for each token
350     '''
351
352     df = defaultdict(lambda: 0)
353     for doc in docs:
354         for w in set(doc):
355             df[w] += 1
356     return df
357
358 def calc_tfidf(docs):
359     '''The original TF-IDF is a document-wise score. This function will
360     calculate the average TF-IDF on whole dataset as an overall scoring.
361     '''
362
363     tf_idf = defaultdict(lambda: 0)
364     df = calc_df(docs)
365     num_docs = len(docs)

```

```

366 total_doc_len = sum(len(doc) for doc in docs)
367 total_df = sum(i for t, i in df.items())
368 num_words = len(set(flatten2d(docs)))
369 avg_doc_len = total_doc_len / num_docs
370 avg_count = total_doc_len / num_words
371 avg_df = total_df / num_words
372 print('num docs', num_docs, 'num words', num_words, len(df), 'avg_len', avg_doc_len, \
373       'avg word count:', avg_count, 'avg DF', avg_df)
374 min_df = 100
375 print('give penalty for DF less than the minimum of', min_df)
376 for doc in docs:
377     counter = Counter(doc)
378     num_words = len(doc)
379     for token in set(doc):
380         tf = (counter[token] + avg_count) / (num_words + avg_doc_len)
381         df_i = df[token]
382         idf = np.log(num_docs / (df_i + 100))
383         tf_idf[token] += tf * idf
384
385 for token in tf_idf:
386     # Smooth the average to avoid tokens with large TF-IDF value but
387     # only appeared in a few documents
388     tf_idf[token] /= (df[token] + min_df)
389
390 return tf_idf
391
392 def get_top_topics(named_entities, n=5, method='tf'):
393     '''Get most dominant objects as popular topics by either word count or TF-IDF'''
394
395
396     print('calculating most popular topics by ' + method + '...')
397     if method == 'tf':
398         ranks = nltk.FreqDist(w for w, t in flatten3d(named_entities))
399         print(ranks.most_common(50))
400         ranks = dict(ranks)
401     elif method == 'tfidf':
402         ranks = calc_tfidf([[w for w, t in flatten2d(doc)] for doc in named_entities])
403     ranks = [(k, v) for k, v in ranks.items()]
404     topics = heapq.nlargest(n, ranks, key=itemgetter(1))
405     return topics
406
407 def get_surroundings(words, docs, n=4):
408     '''expand the topic to be 2 verb/noun before and 2 verb/noun after the topic'''
409
410
411     print('get surrounding 2 nouns/verbs for words {}'.format(words))
412
413     sur = {}
414     for w, c in words:
415         sur[w] = Counter()
416
417     # POS tags list for searching verbs/nouns
418
419     def _helper(sent):
420         sent_w = [w for w, p in sent]
421         for w, c in words:
422             try:
423                 idx = sent_w.index(w)
424             except ValueError:
425                 continue
426
427             after = 0
428             vicinity = [sent[i] for i in [idx-2, idx-1, idx+1, idx+2]
429                          if i >= 0 and i < len(sent)]
430             for (wi, pi) in vicinity:
431                 if wi != w and (pi.startswith('N') or pi.startswith('V')):
432                     sur[w][wi] += 1
433
434     foreach2d(_helper, docs)
435     ret = []
436     for w, c in words:
437         ret.append({'topic': w, 'score': c, 'keywords': sur[w].most_common(n)})
438     return ret
439

```

```

440 def calc_intermediate_data(dataset):
441     docs = tokenise(dataset)
442     save_pkl(docs, 'tokenised_docs.pkl')
443     vocab = calc_vocab(docs)
444     print('Size of vocabulary: {}'.format(len(vocab)))
445
446
447     tagged_docs = calc_pos_tags(docs)
448     docs = vocab = None
449
450     named_entities = calc_ne_all(tagged_docs)
451
452     # Remove stopwords after POS tagging and NER finished
453     tagged_docs = remove_stopwords(tagged_docs)
454     named_entities = remove_stopwords(named_entities)
455
456     tagged_docs = do_stemming(tagged_docs)
457     named_entities = do_stemming(named_entities)
458     return tagged_docs, named_entities
459
460 def mine_topics(dataset, intermediate_data, group='all'):
461     print('-' * 80)
462     print('mining most popular topics for group ' + group)
463     print('-' * 80)
464     tagged_docs, named_entities = intermediate_data
465
466     if group != 'all':
467         if group == 'male' or group == 'female':
468             idx = [i for i, rec in enumerate(dataset) if rec.meta.gender == group]
469         elif group == '<=20':
470             idx = [i for i, rec in enumerate(dataset) if rec.meta.age <= 20]
471         elif group == '>20':
472             idx = [i for i, rec in enumerate(dataset) if rec.meta.age > 20]
473         else:
474             raise NotImplementedError()
475         tagged_docs = [tagged_docs[i] for i in idx]
476         named_entities = [named_entities[i] for i in idx]
477
478     entity_set = set(w for w, t in flatten3d(named_entities))
479     print('selected docs: {}, {}'.format(len(tagged_docs), len(named_entities)))
480
481     entity_words = named_entities
482
483     ret = {}
484     num_keywords = 200
485     print('----- result from TFIDF -----')
486     topics = get_top_topics(entity_words, n=50, method='tfidf')
487     keywords = get_surroundings(topics, tagged_docs, n=num_keywords)
488     ret['tfidf'] = keywords
489
490     print('----- result from TF -----')
491     topics = get_top_topics(entity_words, n=50, method='tf')
492     keywords = get_surroundings(topics, tagged_docs, n=num_keywords)
493     ret['tf'] = keywords
494     return ret
495
496 def main_intermediate():
497     if not _DEBUG:
498         dataset = read_blogs('blogs')
499     else:
500         dataset = read_blogs('blogs', cache_file=None)
501
502     intermediate_data = calc_intermediate_data(dataset)
503     save_pkl(intermediate_data, 'intermediate_data.pkl')
504     return dataset, intermediate_data
505
506 def main_mine_topics(dataset=None, intermediate_data=None):
507     if dataset is None:
508         dataset = load_pkl('blogs.pkl')
509     if intermediate_data is None:
510         intermediate_data = load_pkl('intermediate_data.pkl')
511
512     def _post_filter(word):
513

```

```

514         return len(word[0]) > 0 and word[0] not in ('lol', 'fuck', 'Im',
515             'urllink quizilla', 'urllink hello', 'haiz',
516             ',', '\', '@', ';', '.', '\\', '/', '!', '?')
517
518 docs, entities = intermediate_data
519 docs = filter3d(_post_filter, docs)
520 entities = filter3d(_post_filter, entities)
521 intermediate_data = docs, entities
522
523 topics = {}
524 topics['male'] = mine_topics(dataset, intermediate_data, group='male')
525 topics['female'] = mine_topics(dataset, intermediate_data, group='female')
526 topics['less_or_20'] = mine_topics(dataset, intermediate_data, group='<=20')
527 topics['over_20'] = mine_topics(dataset, intermediate_data, group='>20')
528 topics['all'] = mine_topics(dataset, intermediate_data, group='all')
529 if _DEBUG:
530     suffix = 'debug'
531 else:
532     suffix = date.today().strftime('%Y%m%d')
533     if NUM_SAMPLES > 0:
534         suffix += '-' + str(NUM_SAMPLES)
535
536 save_json(topics, 'topics-{}.json'.format(suffix))
537
538 def main():
539     if len(sys.argv) <= 1:
540         phases = [1, 2]
541     else:
542         phases = [int(i) for i in sys.argv[1].split(',')]
543
544     dataset = intermediate_data = None
545     for ph in phases:
546         if ph == 1:
547             dataset, intermediate_data = main_intermediate()
548         elif ph == 2:
549             main_mine_topics(dataset, intermediate_data)
550
551 if __name__ == '__main__':
552     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  from glob import glob
7  from tqdm import tqdm
8  from collections import Counter, defaultdict
9  import numpy as np
10 from numpy.linalg import norm
11 import pandas as pd
12 from as2 import load_pkl, save_pkl, Record, MetaData, Post, \
13     stem_word, flatten3d, foreach3d, lemmatizer
14 import matplotlib.pyplot as plt
15 from wordcloud import WordCloud
16
17 MAX_FONT_SIZE = 80
18
19 def show_summary(dataset):
20     '''This function describes the summary of dataset or human inspection.
21     It's not necessary for the mining process.
22
23     Parameters
24     -----
25     dataset : list of Record
26         The blog dataset
27     '''
28
29     df = pd.DataFrame([d.meta for d in dataset])
30     df['blog_count'] = [len(d.posts) for d in dataset]
31     df['char_count'] = [sum(len(p.text) for p in d.posts) for d in dataset]
32

```

```

33 print(df.describe(include='all'))
34 print('{} possible values for "gender": {}'.format(
35     len(df.gender.unique()), ', '.join(sorted(df.gender.unique()))))
36 print('{} possible values for category: {}'.format(
37     len(df.category.unique()), ', '.join(sorted(df.category.unique()))))
38 print('{} possible values for zodiac: {}'.format(
39     len(df.zodiac.unique()), ', '.join(sorted(df.zodiac.unique()))))
40
41 plt.rcParams.update({'font.size': 20})
42 df['gender'].value_counts().plot(kind='bar')
43 plt.xticks(rotation=0)
44 plt.gcf().tight_layout()
45 plt.savefig('img/show-gender.png')
46
47 plt.rcParams.update({'font.size': 10})
48 plt.clf()
49 df['category'].value_counts().plot(kind='bar')
50 plt.gcf().tight_layout()
51 plt.savefig('img/show-category.png')
52
53 plt.rcParams.update({'font.size': 18})
54 plt.clf()
55 df['zodiac'].value_counts().plot(kind='bar')
56 plt.xticks(rotation=90)
57 plt.gcf().tight_layout()
58 plt.savefig('img/show-zodiac.png')
59
60 plt.rcParams.update({'font.size': 20})
61 plt.clf()
62 age = df['age']
63 df['age'].hist(bins=20)
64 plt.gcf().tight_layout()
65 plt.savefig('img/show-age.png')
66
67 plt.clf()
68 cnt = df['blog_count']
69 logbins = np.logspace(np.log10(cnt.min()), np.log10(cnt.max()), 20)
70 cnt.hist(bins=logbins)
71 plt.xscale('log')
72 plt.gcf().tight_layout()
73 plt.savefig('img/show-blog-count.png')
74
75 plt.clf()
76 cnt = df['char_count']
77 logbins = np.logspace(np.log10(cnt.min()), np.log10(cnt.max()), 20)
78 cnt.hist(bins=logbins)
79 plt.xscale('log')
80 plt.gcf().tight_layout()
81 plt.savefig('img/show-char-count.png')
82
83 plt.clf()
84 df['gender_age'] = [g + '\n' + ('<=20' if a <= 20 else '>20') \
85     for (g, a) in zip(df['gender'], df['age'])]
86 df['gender_age'].value_counts()[[2, 3, 1, 0]].plot(kind='bar')
87 plt.xticks(rotation=0)
88 plt.gcf().tight_layout()
89 plt.savefig('img/show-gender-age.png')
90
91 def calc_stem_map():
92     '''Map word stem back to the most representative word so we can display valid
93     English words in the word cloud, and also for calculating the coherence score'''
94
95
96     print('building map from stem to words ...')
97     docs = load_pkl('tokenised_docs.pkl')
98
99     stem2word = defaultdict(lambda *_: Counter())
100     def _helper(w):
101         s = stem_word(w)
102         stem2word[s][lemmatizer.lemmatize(w.lower())] += 1
103
104     print('calculating map...')
105     foreach3d(_helper, docs)
106

```



```

107 out = {}
108 for k, cnt in stem2word.items():
109     out[k] = cnt.most_common(10)
110 save_pkl(out, 'stem2word.pkl')
111 return out
112
113 # Colour functions for word cloud
114 def color_black(word, *args, **kwargs):
115     return '#000000'
116
117 def grey_color_func(word, font_size, position, orientation, random_state=None, **kwargs):
118     return 'hsl(0, 0%, {:d}%)'.format((MAX_FONT_SIZE - font_size) // (MAX_FONT_SIZE * 1))
119
120 STEM2WORD = None
121 def eval_topics(fpath, method='tf', top_k=2, num_words_in_topic=10):
122     '''Evaluate topics by:
123     1. plotting the word cloud
124     2. calculating the diagnostic and coherence metrics
125     '''
126
127     with open(fpath, encoding='utf8') as f:
128         result = json.load(f)
129
130     global STEM2WORD
131     if STEM2WORD is None:
132         STEM2WORD = load_pkl('stem2word.pkl')
133
134     def _2w(w):
135         if w in STEM2WORD:
136             return STEM2WORD[w][0][0]
137         else:
138             return w
139
140     topics_formatted = {}
141     for group, topics2 in result.items():
142         topics = topics2[method]
143         topics_formatted[group] = []
144         for i, topic in enumerate(topics[:top_k]):
145             topic_name = _2w(topic['topic'])
146             words = {}
147             words.update((_2w(kw[0]), kw[1]) for kw in topic['keywords'][: (num_words_in_topic-1)])
148             if method == 'tf':
149                 words[topic_name] = topic['score']
150             else:
151                 try:
152                     words[topic_name] = topic['keywords'][0][1] * 2 # fake frequency for display
153                 except IndexError:
154                     words[topic_name] = 1
155             topics_formatted[group].append((topic_name, words))
156
157     print(topics_formatted)
158     plot_topics(topics_formatted, method=method)
159     return calc_coherence_all(topics_formatted, method=method)
160
161 LSA = None
162 def load_lsa():
163     global LSA
164     if LSA is None:
165         print('loading LSA model...')
166         try:
167             LSA = load_pkl('lsa.pkl')
168         except:
169             print('failed')
170             print('loading LSA model...')
171             with open('similar/LSA-MODELS/LSA-MODEL-TASA-LEMMATIZED-DIM300/vocab.txt') as f:
172                 vocab = [x.strip() for x in f]
173             print('vocab size:', len(vocab))
174
175             with open('similar/LSA-MODELS/LSA-MODEL-TASA-LEMMATIZED-DIM300/lsaModel.txt') as f:
176                 vec = [np.array([float(x) for x in line.split()]) for line in f]
177             print('vector size:', len(vec), len(vec[0]))
178             LSA = {w: v for w, v in zip(vocab, vec)}
179             save_pkl(LSA, 'lsa.pkl')
180     return LSA

```

```

181
182 WIKI_PMI = None
183 def load_wiki_pmi():
184     global WIKI_PMI
185     if WIKI_PMI is None:
186         print('loading wiki PMI model...')
187         try:
188             WIKI_PMI = load_pkl('wiki-pmi.pkl')
189         except:
190             print('load from original files...')
191             WIKI_PMI = {}
192             for fname in tqdm(glob('semilar/wiki-pmi/*')):
193                 with open(fname) as f:
194                     next(f)
195                     next(f)
196                     next(f)
197                     next(f)
198                     for line in f:
199                         a, b, s = line.strip().split()
200                         s = float(s)
201                         WIKI_PMI[(a, b)] = s
202
203             save_pkl(WIKI_PMI, 'wiki-pmi.pkl')
204
205     return WIKI_PMI
206
207 GLOVE = None
208 def load_glove(ndim=100):
209     global GLOVE
210     if GLOVE is None:
211         print('loading glove embeddings...')
212         try:
213             GLOVE = load_pkl('glove{}.pkl'.format(ndim))
214         except:
215             print('failed')
216             GLOVE = {}
217             fname = 'embeddings/glove.6B.{}.txt'.format(ndim)
218             print('load from file', fname)
219             with open(fname) as f:
220                 for line in f:
221                     arr = line.strip().split()
222                     GLOVE[arr[0].strip()] = np.array([float(f) for f in arr[1:]])
223
224             save_pkl(GLOVE, 'glove{}.pkl'.format(ndim))
225
226     return GLOVE
227
228 def cosine_similarity(a, b):
229     return np.dot(a, b) / (norm(a) * norm(b))
230
231 def lsa_score(wi, wj):
232     lsa = load_lsa()
233     vi = lsa[wi]
234     vj = lsa[wj]
235     return cosine_similarity(vi, vj)
236
237 def pmi_score(wi, wj):
238     pmi = load_wiki_pmi()
239     p = (wi, wj)
240     if p in pmi:
241         return pmi[p]
242     else:
243         return pmi[wj, wi]
244
245 def glove_score(wi, wj):
246     lsa = load_glove()
247     vi = lsa[wi]
248     vj = lsa[wj]
249     return cosine_similarity(vi, vj)
250
251 def calc_coherence(words, f):
252     n = 0
253     score = 0.0
254     for i in range(len(words)):

```

```

255         for j in range(i+1, len(words)):
256             try:
257                 score += f(words[i], words[j])
258             except KeyError:
259                 print('warning: cannot find pairwise association: {} {}'.format(words[i], words[j]))
260                 continue
261             n += 1
262         return score / n
263
264
265 def calc_word_length(words):
266     return sum(len(w) for w in words) / len(words)
267
268 WORD_COUNT = None
269 def calc_topic_size(words):
270     global WORD_COUNT
271     if WORD_COUNT is None:
272         docs, _ = load_pkl('intermediate_data.pkl')
273         WORD_COUNT = Counter(w for w, t in flatten3d(docs))
274     return sum(WORD_COUNT[w] for w in words)
275
276 def calc_coherence_all(topics_all, method):
277     metrics = []
278     for group, topics in topics_all.items():
279         for i, (topic_name, words_freq) in enumerate(topics):
280             if len(words_freq) <= 1:
281                 continue
282             words = sorted(words_freq.keys())
283             print('topic:', topic_name, words)
284             lsa = calc_coherence(words, lsa_score)
285             print('lsa score', lsa)
286             we_glove = calc_coherence(words, glove_score)
287             print('we glove score', we_glove)
288             wl = calc_word_length(words)
289             print('avg word length', wl)
290             ts = calc_topic_size(words)
291             print('avg topic size', ts)
292             metrics.append((group, i, topic_name, words, lsa, we_glove, wl, ts))
293     return metrics
294
295 def print_metrics_as_table(metrics, fpath):
296     '''Plot metrics as \LaTeX table'''
297
298     with open(fpath, 'w') as f:
299         for group, i, topic, keywords, lsa, we_glove, wl, ts in metrics:
300             if i > 0:
301                 group = ' '
302             elif group == 'less_or_20':
303                 group = '20 or younger'
304             elif group == 'over_20':
305                 group = 'over 20'
306             elif group == 'all':
307                 group = 'everyone'
308             f.write(f'{group} & {topic} & {ts//1000:}k & {wl:.1f} & {lsa:.3f} & {we_glove:.3f} \\\n')
309             if i == 1: # two topics for each group
310                 f.write('\nhline\n')
311             else:
312                 f.write('\cline{2-6}\n')
313
314
315 def plot_topics(topics_json, method):
316     '''Plot word cloud for the keywords in each topic'''
317
318     for group, topics in topics_json.items():
319         for i, (topic_name, words) in enumerate(topics):
320             print('topic: ', topic_name, 'number of keywords:', len(words))
321             wc = WordCloud(background_color="white",
322                             max_font_size=80,
323                             max_words=len(words)+1,
324                             color_func=grey_color_func)
325             wc.generate_from_frequencies(words)
326
327             plt.clf()
328

```

```

329         plt.imshow(wc, interpolation="bilinear")
330         plt.axis("off")
331         plt.title(topic_name, y=-0.25, fontsize=20, fontname='Georgia')
332         plt.gcf().tight_layout()
333         fig_path = 'img/{}-{}-{}.png'.format(group, method, i+1, topic_name)
334         print('drawing ' + fig_path)
335         plt.savefig(fig_path)
336
337 def main():
338     cmd = sys.argv[1]
339     if cmd == 'show':
340         show_summary(load_pkl('blogs.pkl'))
341     elif cmd == 'eval':
342         fpath = sys.argv[2]
343         metrics_tf = eval_topics(fpath, top_k=2, method='tf')
344         metrics_tfidf = eval_topics(fpath, top_k=2, method='tfidf')
345         print('tf metrics', metrics_tf)
346         print_metrics_as_table(metrics_tf, 'metrics-tf.tex')
347         print('tfidf metrics', metrics_tfidf)
348         print_metrics_as_table(metrics_tfidf, 'metrics-tfidf.tex')
349     elif cmd == 'stem2word':
350         calc_stem_map()
351
352 if __name__ == '__main__':
353     main()

```