

# Short Text Topic Modeling Techniques, Applications, and Performance: A Survey

Jipeng Qiang, Zhenyu Qian, Yun Li, Yunhao Yuan, and Xindong Wu, *Fellow, IEEE*,

**Abstract**—Analyzing short texts infers discriminative and coherent latent topics that is a critical and fundamental task since many real-world applications require semantic understanding of short texts. Traditional long text topic modeling algorithms (e.g., PLSA and LDA) based on word co-occurrences cannot solve this problem very well since only very limited word co-occurrence information is available in short texts. Therefore, short text topic modeling has already attracted much attention from the machine learning research community in recent years, which aims at overcoming the problem of sparseness in short texts. In this survey, we conduct a comprehensive review of various short text topic modeling techniques proposed in the literature. We present three categories of methods based on Dirichlet multinomial mixture, global word co-occurrences, and self-aggregation, with example of representative approaches in each category and analysis of their performance on various tasks. We develop the first comprehensive open-source library, called STTM, for use in Java that integrates all surveyed algorithms within a unified interface, benchmark datasets, to facilitate the expansion of new methods in this research field. Finally, we evaluate these state-of-the-art methods on many real-world datasets and compare their performance against one another and versus long text topic modeling algorithm.

**Index Terms**—Topic modeling, Short text, Sparseness, Short text topic modeling.



## 1 INTRODUCTION

Short texts have become an important information source including news headlines, status updates, web page snippets, tweets, question/answer pairs, etc. Short text analysis has been attracting increasing attention in recent years due to the ubiquity of short text in the real-world [1]–[3]. Effective and efficient models infer the latent topics from short texts, which can help discover the latent semantic structures that occur in a collection of documents. Short text topic modeling algorithms are always applied into many tasks such as topic detection [4], classification [5], comment summarization [6], user interest profiling [7].

Traditional topic modeling algorithms such as probabilistic latent semantic analysis (PLSA) [8] and latent Dirichlet allocation (LDA) [9] are widely adopted for discovering latent semantic structure from text corpus without requiring any prior annotations or labeling of the documents. In these algorithms, each document may be viewed as a mixture of various topics and each topic is characterized by a distribution over all the words. Statistical techniques (e.g., Variational methods and Gibbs sampling) are then employed to infer the latent topic distribution of each document and the word distribution of each topic using higher-order word co-occurrence patterns [10]. These algorithms and their variants have had a major impact on numerous applied fields in modeling text collections news articles, research papers, and blogs [11]–[13]. However, traditional topic models experience large performance degradation over short texts due to the lack of word co-occurrence information in each short text [1], [14]. Therefore, short text topic modeling has already attracted much

attention from the machine learning research community in recent years, which aims at overcoming the problem of sparseness in short texts.

Earlier works [15], [16] still used traditional topic models for short texts, but exploited external knowledge or metadata to bring in additional useful word co-occurrences across short texts, and therefore may boost the performance of topic models. For example, Phan et al. [16] first learned latent topics from Wikipedia, and then inferred topics from short texts. Weng et al. [7] and Mehrotra et al. [17] aggregated tweets for pseudo-document using hashtags and the same user respectively. The problem lies in that auxiliary information or metadata is not always available or just too costly for deployment. These studies suggest that topic models specifically designed for general short texts are imperative. This survey will provide a taxonomy that captures the existing short text topic modeling algorithms and their application domains.

News aggregation websites often rely on news headlines to cluster different source news about the same event. In Table 1, we show an event about artificial intelligence reported on March 1, 2018. As presented, all these short texts were reported about the same event. From these short texts, we can found these following characteristics. (1) Obviously, each short text lacks enough word co-occurrence information. (2) Due to a few words in each text, most texts are probably generated by only one topic (e.g., text 1, text2, text 3). (3) Statistical information of words among texts cannot fully capture words that are semantically related but rarely co-occur. For example, President Trump of text 1 and White House of text 2 are highly semantically related, and AI is short for Artificial Intelligence. (4) The single-topic assumption may be too strong for some short texts. For example, text 3 is probably associated with a small number of topics (e.g., one to three topics). Considering these characteristics, existing short text topic modeling algorithms were proposed by trying to solve one or two of these characteristics. Here, we divide the short text topic modeling algorithms basically into the following three major

- J. Qiang, Z. Qian, Y. Li and Y. Yuan are with the Department of Computer Science, Yangzhou, Jiangsu, P. R. China, 225127.  
E-mail: {jpqiang, liyun, yhyuan}@yzu.edu.cn
- X. Wu is with Department of Computer Science, University of Louisiana at Lafayette, Louisiana, USA.  
E-mail: xwu@louisiana.edu

TABLE 1  
An event about artificial intelligence was reported by different news media on March 1, 2018.

Number	Media	Headline
1	Lawfare	President Trump's Executive Order on Artificial Intelligence
2	Nextgov	White Houses Race to Maintain AI Dominance Misses Opportunity
3	Forbes	Artificial Intelligence Regulation may be Impossible
4	CognitiveWorld	Pop Culture, AI and Ethics

categories.

**(1) Dirichlet multinomial mixture (DMM) based methods:**

A simple and effective model, Dirichlet Multinomial Mixture model, has been adopted to infer latent topics in short texts [18], [19]. DMM follows the simple assumption that each text is sampled from only one latent topic. Considering the characteristics (1) and (2) in short texts, this assumption is reasonable and suitable for short texts compared to the complex assumption adopted by LDA that each text is modeled over a set of topics [20], [21]. Nigam et al. [22] proposed an EM-based algorithm for Dirichlet Multinomial Mixture (DMM) model. Except for the basic expectation maximization (EM), a number of inference methods have been used to estimate the parameters including variation inference and Gibbs sampling. For example, Yu et al. [23] proposed the DMAFP model based on variational inference algorithm [24]. Yin et al. [18] proposed a collapsed Gibbs sampling algorithm for DMM. Other variations based on DMM [25]–[27] were proposed for improving the performance. The above models based on DMM ignore the characteristic (3). Therefore, many models by incorporating word embeddings into DMM were proposed [28], [29], because word embeddings learned from millions of external documents contain semantic information of words [30]. Not only word co-occurrence words belong to one topic, but words with high similarity have high probability belonging to one topic, which can effectively solve the data sparsity issue. To highlight the characteristic (4), a Poisson-based DMM model (PDMM) was proposed that allows each short text is sampled by a limited number of topics [31]. Accordingly, Li et al. [31] proposed a new model by directly extending the PDMM model using word embeddings.

**(2) Global word co-occurrences based methods:** Considering the characteristic (1), some models try to use the rich global word co-occurrence patterns for inferring latent topics [14], [32]. Due to the adequacy of global word co-occurrences, the sparsity of short texts is mitigated for these models. According to the utilizing strategies of global word co-occurrences, this type of models can be divided into two types. 1) The first type directly uses the global word co-occurrences to infer latent topics. Biterm topic modeling (BTM) [14] posits that the two words in a biterm share the same topic drawn from a mixture of topics over the whole corpus. Some models extend the Biterm Topic Modeling (BTM) by incorporating the burstiness of biterms as prior knowledge [21] or distinguishing background words from topical words [33]. 2) The second type first constructs word co-occurrence network using global word co-occurrences and then infers latent topics from this network, where each word correspond to one node and the weight of each edge stands for the empirical co-occurrence probability of the connected two words [32], [34].

**(3) Self-aggregation based methods:** Self-aggregation based methods are proposed to perform topic modeling and text self-aggregation during topic inference simultaneously. Short texts are merged into long pseudo-documents before topic inference that can help improve word co-occurrence information. Different from

the aforementioned aggregation strategies [7], [17], this type of methods SATM [20] and PTM [35] posit that each short text is sampled from a long pseudo-document unobserved in current text collection, and infer latent topics from long pseudo-documents, without depending on auxiliary information or metadata. Considering the characteristic (3), Qiang et al. [36] and Bicalho et al. [37] merged short texts into long pseudo-documents using word embeddings.

## 1.1 Our contributions

This survey has the following three-pronged contribution:

(1) We propose a taxonomy of algorithms for short text topic modeling and explain their differences. We define three different tasks, i.e., application domains of short text topic modeling techniques. We illustrate the evolution of the topic, the challenges it faces, and future possible research directions.

(2) To facilitate the expansion of new methods in this field, we develop the first comprehensive open-source JAVA library, called STTM, which not only includes all short text topic modeling algorithms discussed in this survey with a uniform easy-to-use programming interface but also includes a great number of designed modules for the evaluation and application of short text topic modeling algorithms. STTM is open-sourced at <https://github.com/qiang2100/STTM>.

(3) We finally provide a detailed analysis of short text topic modeling techniques and discuss their performance on various applications. For each model, we analyze their results through comprehensive comparative evaluation on some common datasets.

## 1.2 Organization of the survey

The rest of this survey is organized as follows. In Section 2, we introduce the task of short text topic modeling. Section 3 proposes a taxonomy of short text topic modeling algorithms and provides a description of representative approaches in each category. The list of applications for which researchers have used the short text topic modeling algorithms is provided in Section 4. Section 5 presents our Java library for short text topic modeling algorithms. In the next two sections, we describe the experimental setup (Section 6) and evaluate the discussed models (Section 7). Finally, we draw our conclusions and discuss potential future research directions in Section 8.

## 2 DEFINITIONS

In this section, we formally define the problem of short text topic modeling.

Given a short text corpus  $\mathcal{D}$  of  $N$  documents, with a vocabulary  $\mathcal{W}$  of size  $V$ , and  $K$  pre-defined latent topics. One document  $d$  is represented as  $(w_{d,1}, w_{d,2}, \dots, w_{d,n_d})$  in  $\mathcal{D}$  including  $n_d$  words.

A topic  $\phi$  in a given collection  $\mathcal{D}$  is defined as a multinomial distribution over the vocabulary  $\mathcal{W}$ , i.e.,  $\{p(w|\phi)\}_{w \in \mathcal{W}}$ .

TABLE 2  
The notations of symbols used in the paper

$D, N$	Documents and number of documents in the corpus
$W, V$	The vocabulary and number of words in the vocabulary
$K$	Number of pre-defined latent topics
$\bar{l}$	Average length of each document in $D$
$n_k$	Number of words associated with topic $k$
$m_k$	Number of documents associated with topic $k$
$n_k^w$	Number of word $w$ associated with topic $k$ in $\vec{d}$
$n_d$	Number of words in document $d$
$n_d^w$	Number of word $w$ in document $d$
$n_d^k$	Number of word associated with topic $k$ in document $d$
$n_{k,d}^w$	Number of word $w$ associated with topic $k$ in document $d$
$P$	Long pseudo-document set generated by models
$\phi$	Topic distribution
$\theta$	Document-topic distribution
$z$	Topic indicator
$U$	Number of dimensions in word embeddings
$\zeta$	Time cost of considering GPU model
$\varsigma$	Maximum number of topics allowable in a short text
$c$	Size of sliding window

The topic representation of a document  $d$ ,  $\theta_d$ , is defined as a multinomial distribution over  $K$  topics, i.e.,  $\{p(\phi_k|\theta_d)\}_{k=1,\dots,K}$ . The general task of topic modeling aims to find  $K$  salient topics  $\phi_{k=1,\dots,K}$  from  $D$  and to find the topic representation of each document  $\theta_{d=1,\dots,N}$ .

Most classical probabilistic topic models adopt the Dirichlet prior for both the topics and the topic representation of documents, which are first used in LDA [9], which is  $\phi_k \sim \text{Dirichlet}(\beta)$  and  $\theta_d \sim \text{Dirichlet}(\alpha)$ . In practice, the Dirichlet prior smooths the topic mixture in individual documents and the word distribution of each topic, which alleviates the overfitting problem of probabilistic latent semantic analysis (PLSA) [8], especially when the number of topics and the size of vocabulary increase. Therefore, all of existing short text topic modeling algorithms adopt Dirichlet distribution as prior distribution.

Given a short text corpus  $D$  with a vocabulary of size  $V$ , and the predefined number of topics  $K$ , the major tasks of short text topic modeling can be defined as to:

- (1). Learn the word representation of topics  $\phi$ ;
- (2). Learn the sparse topic representation of documents  $\theta$ .

All the notations used in this paper are summarized in Table 2.

### 3 ALGORITHMIC APPROACHES: A TAXONOMY

In the past decade, there has been much work to discover latent topics from short texts using traditional topic modeling algorithms by incorporating external knowledge or metadata. More recently, researchers focused on proposing new short text topic modeling algorithms. In the following, we present historical context about the research progress in this domain, then propose a taxonomy of short text topic modeling techniques including: (1) Dirichlet Multinomial Mixture (DMM) based methods, (2) Global word co-occurrence based methods, and (3) Self-aggregation based methods.

#### 3.1 Short Text Topic Modeling Research Context and Evolution

Traditional topic modeling algorithms such as probabilistic latent semantic analysis (PLSA) [8] and latent Dirichlet allocation (LDA) [9] are widely adopted for discovering latent semantic

structure from text corpus by capturing word co-occurrence pattern at the document level. Hence, more word co-occurrences would bring in more reliable and better topic inference. Due to the lack of word co-occurrence information in each short text, traditional topic models have a large performance degradation over short texts. Earlier works focus on exploiting external knowledge to help enhance the topic inference of short texts. For example, Phan et al. [16] adopted the learned latent topics from Wikipedia to help infer the topic structure of short texts. Similarly, Jin et al. [15] searched auxiliary long texts for short texts to infer latent topics of short texts for clustering. A large regular text corpus of high quality is required by these models, which bring in big limitation for these models.

Since 2010, research on topic discovery from short texts has been shifted to merging short texts into long pseudo-documents using different aggregation strategies before adopting traditional topic modeling to infer the latent topics. For example, Weng et al. [7] merge all tweets of one user into a pseudo-document before using LDA. Other information includes hashtags, timestamps, and named entities have been tread as metadata to merging short texts [17], [19], [38]. However, helpful metadata may not be accessible in any domains, e.g., news headlines and search snippets. These studies suggest that topic models specifically designed for general short texts are crucial. This survey will provide a taxonomy that captures the existing strategies and these application domains.

#### 3.2 A Taxonomy of Short Text Topic Modeling Methods

We propose a taxonomy of short text topic modeling approaches. We categorize topic modeling approaches into three broad categories: (1) Dirichlet Multinomial Mixture (DMM) based, (2) Global Word co-occurrences based, and (3) Self-aggregation based. Below we describe the characteristics of each of these categories and present a summary of some representative methods for each category (cf. Table 3).

#### 3.3 Dirichlet Multinomial Mixture based Methods

Dirichlet Multinomial Mixture model (DMM) was first proposed by Nigam et al. [22] based on the assumption that each document is sampled by only one topic. The assumption is more fit for short texts than the assumption that each text is generated by multiple topics. Therefore, many models for short texts were proposed based on this simple assumption [19], [23], [25]. Yin et al. [18] proposed a DMM model based on collapse Gibbs sampling. Zhao et al. [19] proposed a Twitter-LDA model by assuming that one tweet is generated from one topic. Recently, more work incorporates word embeddings into DMM [29], [31].

##### 3.3.1 GSDMM

DMM respectively chooses Dirichlet distribution for topic-word distribution  $\phi$  and document-topic distribution  $\theta$  as prior distribution with parameter  $\alpha$  and  $\beta$ . DMM samples a topic  $z_d$  for the document  $d$  by Multinomial distribution  $\theta$ , and then generates all words in the document  $d$  from topic  $z_d$  by Multinomial distribution  $\phi_{z_d}$ . The graphical model of DMM is shown in Figure 1. The generative process for DMM is described as follows:

- (1). Sample a topic proportion  $\theta \sim \text{Dirichlet}(\alpha)$ .
- (2). For each topic  $k \in \{1, \dots, K\}$ :  
Draw a topic-word distribution  $\theta_k \sim \text{Dirichlet}(\beta)$ .
- (3). For each document  $d \in D$ :  
(b) Sample a topic  $z_d \sim \text{Multinomial}(\theta)$ .

TABLE 3  
List of short text topic modeling approaches

Category	Year	Published	Authors	Method	Time Complexity of One Iteration
DMM	2014	KDD [18]	J. Yin & et al.	GSDMM	$O(KNl)$
	2015	TACL [28]	D. Nguyen & et al.	LF-DMM	$O(O(2KN\bar{l} + KVU))$
	2016	SIGIR [29]	C. Li & et al.	GPU-DMM	$O(KN\bar{l} + N\bar{l}\zeta + KV)$
	2017	TOIS [31]	C. Li & et al.	GPU-PDMM	$O(N\bar{l}\sum_{i=1}^{c-1}C_K^i + N\bar{l}\zeta + KV)$
Global word co-occurrences	2013	WWW [14]	X. Chen & et al.	BTM	$O(KNlc)$
	2016	KAIS [32]	Y. Zuo & et al.	WNTM	$O(KNlc(c-1))$
Self-aggregation	2015	IJCAI	X. Quan & et al.	SATM	$O(N\bar{l}PK)$
	2016	KDD	Y. Zuo & et al.	PTM	$O(N\bar{l}(P+K))$

(c) For each word  $w \in \{w_{d,1}, \dots, w_{d,n_d}\}$ :  
Sample a word  $w \sim \text{Multinomial}(\phi_{z_d})$ .

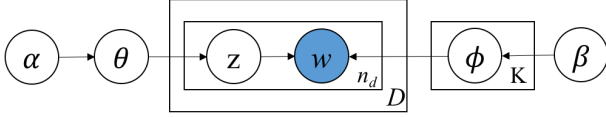


Fig. 1. Graphical Model of GSDMM.

Gibbs sampling algorithm for Dirichlet Multinomial Mixture model is denoted as GSDMM, which is based on the assumption that each text is sampled by a single topic [18]. Here, for better representation of latent topics, we represent a topic with the topic feature (CF) vector, which essentially is a big document combined with its documents.

The TF vector of a topic  $k$  is defined as a tuple  $\{n_k^w (w \in W), m_k, n_k\}$ , where  $n_k^w$  is the number of word  $w$  in topic  $k$ ,  $m_k$  is the number of documents in topic  $k$ , and  $n_k$  is the number of words in topic  $k$ .

The topic feature (TF) presents important addible and deletable properties, as described next.

(1) **Addible Property.** A document  $d$  can be efficiently added to topic  $k$  by updating its TF vector as follows.

$$\begin{aligned} n_k^w &= n_k^w + n_d^w \quad \text{for each word } w \text{ in } d \\ m_k &= m_k + 1; n_k = n_k + n_d \end{aligned}$$

(2) **Deletable Property.** A document  $d$  can be efficiently deleted from topic  $k$  by updating its TF vector as follows.

$$\begin{aligned} n_k^w &= n_k^w - n_d^w \quad \text{for each word } w \text{ in } d \\ m_k &= m_k - 1; n_k = n_k - n_d \end{aligned}$$

The hidden multinomial variable ( $z_d$ ) for document  $d$  is sampled based on collapsed Gibbs sampling, conditioned on a complete assignment of all other hidden variables. GSDMM uses the following conditional probability distribution to infer its topic,

$$p(z_d = k | \mathbf{Z}_{-d}, \mathbf{D}) \propto \frac{m_{k,-d} + \alpha}{N - 1 + K\alpha} \frac{\prod_{j=1}^{n_d} (n_{k,-d}^j + \beta + j - 1)}{\prod_{i=1}^{n_d} (n_{k,-d} + V\beta + i - 1)} \quad (1)$$

where  $\mathbf{Z}$  represents all topics of all documents, the subscript  $-d$  means document  $d$  is removed from its current topic feature (TF) vector, which is useful for the update learning process of GSDMM.

For each document, we first delete it from its current TF vector with the deletable property. Then, we reassign the document to a topic according to the probability of the document belonging to each of the  $K$  topics using Equation 1. After obtaining the topic of the document, we add it from its new TF vector with the addible property. Finally, the posterior distribution of each word belonging to each topic is calculated as the follows,

$$\phi_k^w = \frac{n_k^w + \beta}{n_k + V\beta} \quad (2)$$

### 3.3.2 LF-DMM

The graphical model of LF-DMM is shown in Figure 2. Based on the assumption that each text is sampled by a single topic, **LF-DMM generates the words by Dirichlet multinomial model or latent feature model.** Given two latent-feature vectors  $\tau$  associated with topic  $k$  and  $\omega$  associated with word  $w$ , latent feature model generates a word  $w$  using *softmax* function by the formula,

$$\sigma(w | \tau_k \omega^T) = \frac{e^{(\tau_k \cdot \omega_w)}}{\sum_{w' \in W} e^{(\tau_k \cdot \omega_{w'})}} \quad (3)$$

where  $\omega$  is pre-trained word vectors of all words  $W$ , and  $\omega_w$  is the word vector of word  $w$ .

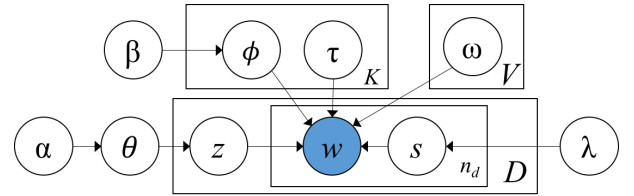


Fig. 2. Graphical Model of LF-DMM.

For each word  $w$  of document  $d$ , a binary indicator variable  $\mathbb{S}_{d,w}$  is sampled from a Bernoulli distribution to determine whether Dirichlet multinomial model or latent feature model will be used to generate  $w$ . The generative process is described as follows:

- (1). Sample a topic proportion  $\theta \sim \text{Dirichlet}(\alpha)$ .
- (2). For each topic  $k \in \{1, \dots, K\}$ :
  - (a) Draw a topic-word distribution  $\theta_k \sim \text{Dirichlet}(\beta)$ .
- (3). For each document  $d \in \mathbf{D}$ :
  - (a) Sample a topic  $z_d \sim \text{Multinomial}(\theta)$ .
  - (b) For each word  $w \in \{w_{d,1}, \dots, w_{d,n_d}\}$ :
    - (i) Sample a binary indicator variable  $\mathbb{S}_{d,w} \in \text{Bernoulli}(\lambda)$
    - (ii) Sample a word  $w \sim (1 - s_w)\text{Multinomial}(\phi_{z_d}) + s_w(\sigma(\tau_{z_d} \omega^T))$ .

Here, the hyper-parameter  $\lambda$  is the probability of a word being generated by latent feature model, and  $\mathbb{S}_{d,w}$  indicates whether Dirichlet multinomial model or latent feature model is applied to word  $w$  of document  $d$ . The topic feature (TF) in LF-DMM is similar with GSDMM, so we do not present the addible and deletable properties for LF-DMM.

Based on collapsed Gibbs sampling, LF-DMM uses the following conditional probability distribution to infer the topic of the document  $d$ ,

$$p(z_d = k | \mathbf{Z}_{-d}, \mathbf{D}, \boldsymbol{\tau}, \boldsymbol{\omega}) \propto (m_{k,-d} + \alpha) \prod_{w \in d} ((1 - \lambda) \frac{n_{k,-d}^w + \beta}{n_{k,-d} + V\beta} + \lambda \sigma(w | \boldsymbol{\tau}_k \boldsymbol{\omega}^T))^{n_d^w} \quad (4)$$

where  $n_d^w$  is the number of word  $w$  in document  $d$ .

The binary indicator variable  $\mathbb{S}_{d,w}$  for word  $w$  in document  $d$  conditional on  $z_d = k$  is inferred using the following distribution,

$$p(\mathbb{S}_{d,w} = s | z_d = k) \propto \begin{cases} (1 - \lambda) \frac{n_{k,-d}^w + \beta}{n_{k,-d} + V\beta} & \text{for } s = 0, \\ \lambda \sigma(w_i | \boldsymbol{\tau}_k \boldsymbol{\omega}^T) & \text{for } s = 1. \end{cases} \quad (5)$$

where the subscript  $-d$  means document  $d$  is removed from its current topic feature (TF) vector.

After each iteration, LF-DMM estimates the topic vectors using the following optimization function,

$$L_k = - \sum_{w \in W} F_k^w (\boldsymbol{\tau}_k \cdot \boldsymbol{\omega}_w - \log(\sum_{w' \in W} e^{\boldsymbol{\tau}_k \cdot \boldsymbol{\omega}_{w'}})) + \mu \|\boldsymbol{\tau}_k\|_2^2 \quad (6)$$

where  $F_k^w$  is the number of times word  $w$  generated from topic  $k$  by latent feature model. LF-DMM adopted L-BFGS<sup>1</sup> [40] to find the topic vector  $\boldsymbol{\tau}_k$  that minimizes  $L_k$ .

### 3.3.3 GPU-DMM

Based on DMM model, GPU-DMM [29] promotes the semantically related words under the same topic during the sampling process by the generalized Pólya urn (GPU) model [41]. When a ball of a particular color is sampled, a certain number of balls of similar colors are put back along with the original ball and a new ball of that color. In this case, sampling a word  $w$  in topic  $k$  not only increases the probability of  $w$  itself under topic  $k$ , but also increases the probability of the semantically similar words of word  $w$  under topic  $k$ .

Given pre-trained word embeddings, the semantic similarity between two words  $w_i$  and  $w_j$  is denoted by  $\cos(w_i, w_j)$  that are measured by cosine similarity. For all word pairs in vocabulary, if the semantic similarity score is higher than a pre-defined threshold  $\epsilon$ , the word pair is saved into a matrix  $\mathbb{M}$ , i.e.,  $\mathbb{M} = \{(w_i, w_j) | \cos(w_i, w_j) > \epsilon\}$ . Then, the promotion matrix  $\mathbb{A}$  with respect to each word pair is defined below,

$$\mathbb{A}_{w_i, w_j} = \begin{cases} 1 & w_i = w_j \\ \mu & w_j \in \mathbb{M}_{w_i} \text{ and } w_j \neq w_i \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

where  $\mathbb{M}_{w_i}$  is the row in  $\mathbb{M}$  corresponding to word  $w_i$  and  $\mu$  is the pre-defined promotion weight.

GPU-DMM and DMM share the same generative process and graphical representation but differ in the topic inference process that they use. Different from DMM and LF-DMM, GPU-DMM first samples a topic for a document, and then only reinforces only the semantically similar words if and only if a word has strong ties with the sampled topic. Therefore, a nonparametric probabilistic sampling process for word  $w$  in document  $d$  is as follows:

$$\mathbb{S}_{d,w} \sim \text{Bernoulli}(\lambda_{w,z_d}) \quad (8)$$

$$\lambda_{w,z_d} = \frac{p(z|w)}{p_{\max}(z'|w)} \quad (9)$$

$$p_{\max}(z'|w) = \max_k p(z = k|w)$$

$$p(z = k|w) = \frac{p(z = k)p(w|z = k)}{\sum_{i=1}^K p(z = i)p(w|z = i)} \quad (10)$$

where  $\mathbb{S}_{d,w}$  indicates whether GPU is applied to word  $w$  of document  $d$  given topic  $z_d$ . We can see that GPU model is more likely to be applied to  $w$  if word  $w$  is highly relate to topic  $z_d$ .

The Topic feature vector of a topic  $k$  in GPU-DMM is defined as a tuple  $\{\tilde{n}_k^w(w \in W), m_k, \tilde{n}_k\}$ .

TF makes the same changes with GSDMM when no GPU is applied, namely  $\mathbb{S}_{d,w} = 0$ . Under  $\mathbb{S}_{d,w} = 1$ , the addible and deletable properties of topic feature (TF) in GPU-DMM are described below.

(1) **Addible Property.** A document  $d$  will be added into topic  $k$  by updating its TF vector as follows,

$$\begin{aligned} \tilde{n}_k &= \tilde{n}_k + n_d^{w_i} \cdot \mathbb{A}_{w_i, w_j} & \text{for each word } w_j \in \mathbb{M}_{w_i} \\ \tilde{n}_k^{w_j} &= \tilde{n}_k^{w_j} + n_d^w \cdot \mathbb{A}_{w_i, w_j} & \text{for each word } w_j \in \mathbb{M}_{w_i} \\ m_k &= m_k + 1 \end{aligned}$$

(2) **Deletable Property.** A document  $d$  will be deleted from topic  $k$  by updating its TF vector as follows,

$$\begin{aligned} \tilde{n}_k &= \tilde{n}_k - n_d^{w_i} \cdot \mathbb{A}_{w_i, w_j} & \text{for each word } w_j \in \mathbb{M}_{w_i} \\ \tilde{n}_k^{w_j} &= \tilde{n}_k^{w_j} - n_d^w \cdot \mathbb{A}_{w_i, w_j} & \text{for each word } w_j \in \mathbb{M}_{w_i} \\ m_k &= m_k - 1 \end{aligned}$$

Accordingly, based on Gibbs sampling, the conditional distribution to infer the topic for each document in Equation 1 is rewritten as follows:

$$p(z_d = k | \mathbf{Z}_{-d}, \mathbf{D}) \propto \frac{m_{k,-d} + \alpha}{N - 1 + K\alpha} \times \frac{\prod_{w \in d} \prod_{j=1}^{n_d^w} (\tilde{n}_{k,-d}^w + \beta + j - 1)}{\prod_{i=1}^{n_d} (\tilde{n}_{k,-d} + V\beta + i - 1)} \quad (11)$$

During each iteration, GPU-PDMM first delete it from its current TF vector with the deletable property. After obtaining the topic of the document, GPU-DMM first updates  $\mathbb{S}_{d,w}$  for GPU using Equation 8, and then updates TF vector for each word using the addible property. Finally, the posterior distribution in Equation 2 for GPU-DMM is rewritten as follows:

$$\phi_k^w = \frac{\tilde{n}_k^w + \beta}{\tilde{n}_k + V\beta} \quad (12)$$

1. LF-DMM used the implementation of the Mallet toolkit [39]

### 3.3.4 GPU-PDMM

Considering the single-topic assumption may be too strong for some short text corpus, Li et al. [31] first proposed Poisson-based Dirichlet Multinomial Mixture model (PDMM) that allows each document can be generated by one or more (but not too many) topics. Then PDMM can be extended as GPU-PDMM model by incorporating generalized Pólya urn (GPU) model during the sampling process.

In GPU-PDMM, each document is generated by  $t_d$  ( $0 < t_d \leq \varsigma$ ) topics, where  $\varsigma$  is the maximum number of topics allowable in a document. GPU-PDMM uses Poisson distribution to model  $t_d$ . The graphical model of GPU-PDMM is shown in Figure 3. The generative process of GPU-PDMM is described as follows.

- (1). Sample a topic proportion  $\theta \sim \text{Dirichlet}(\alpha)$ .
- (2). For each topic  $k \in \{1, \dots, K\}$ :
  - (a) Draw a topic-word distribution  $\theta_k \sim \text{Dirichlet}(\beta)$ .
- (3). For each document  $d \in \mathcal{D}$ :
  - (a) Sample a topic number  $t_d \sim \text{Poisson}(\lambda)$ .
  - (b) Sample  $t_d$  distinct topics  $\mathbf{Z}_d \sim \text{Multinomial}(\theta)$ .
  - (c) For each word  $w \in \{w_{d,1}, \dots, w_{d,n_d}\}$ :
    - (i) Uniformly sample a topic  $z_{d,w} \sim \mathbf{Z}_d$ .
    - (ii) Sample a word  $w \sim \text{Multinomial}(\phi_{z_{d,w}})$ .

Here  $t_d$  is sampled using Poisson distribution with parameter  $\lambda$ , and  $\mathbf{Z}_d$  is the topic set for document  $d$ .

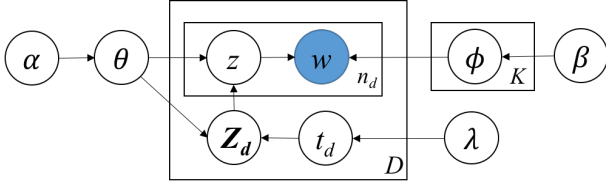


Fig. 3. Graphical Model of GPU-PDMM.

The topic feature (TF) vector of a topic  $k$  in GPU-PDMM is defined as a tuple  $\tilde{n}_k, \tilde{n}_k^w (w \in W), c_k, d_k, n_{k,d}, n_{k,d}^w$ , where  $c_k$  is the number of words associated with topic  $k$  and  $d_k$  represents the word set in topic  $k$ .

The addible and deletable properties of topic feature (TF) in GPU-PDMM are described below. The  $\tilde{n}_k$  and  $\tilde{n}_k^w$  of TF in GPU-PDMM makes the same changes with GPU-DMM. Here, we only describe other variables in TF.

(1) **Addible Property.** Suppose that word  $w$  in document  $d$  will be added to topic  $k$ , TF feature is updates as follows,

$$\begin{aligned} c_k &= c_k + 1 ; d_k = d_k + w \\ n_{k,d} &= n_{k,d} + 1 ; n_{k,d}^w = n_{k,d}^w + 1 \end{aligned}$$

(2) **Deletable Property.** Suppose that word  $d$  will be deleted from topic  $k$ . TF feature is updated as follows,

$$\begin{aligned} c_k &= c_k - 1 ; d_k = d_k - w \\ n_{k,d} &= n_{k,d} - 1 ; n_{k,d}^w = n_{k,d}^w - 1 \end{aligned}$$

The Gibbs sampling process of GPU-PDMM is similar to GPU-DMM, it updates the topic for word  $w$  in document  $d$  using the following equation,

$$p(z_{d,w} = k | z_{\neg(d,w)}, \mathbf{Z}_d, \mathbf{D}) \propto \frac{1}{t_d} \times \frac{\tilde{n}_{k,\neg(d,w)}^w + \beta}{\sum_w \tilde{n}_{k,\neg(d,w)}^w + V\beta} \quad (13)$$

Conditioned on all  $z_{d,w}$  in document  $d$ , GPU-PDMM samples each possible  $\mathbf{Z}_d$  as follows,

$$\begin{aligned} p(\mathbf{Z}_d | \mathbf{Z}_{\neg d}, \mathbf{D}) &\propto \frac{\lambda^{t_d}}{t_d^{n_d}} \\ &\times \frac{\prod_{k \in \mathbf{Z}_d} (c_{k,\neg d} + \alpha)}{\prod_{i=0}^{t_d-1} (\sum_k c_{k,\neg d} + K\alpha - i)} \\ &\times \prod_{k \in \mathbf{Z}_d} \frac{\prod_{w \in d_k} \prod_{i=0}^{n_{k,d}^w-1} (\tilde{n}_{k,\neg d}^w + n_{k,d}^w - i + \beta)}{\prod_{i=0}^{n_{k,d}-1} (\sum_w \tilde{n}_{k,\neg d}^w + n_{k,d} - i + V\beta)}. \end{aligned} \quad (14)$$

During each iteration, for each document  $d$ , GPU-PDMM first updates TF vector using Deletable Property and the topic for each word  $w$  in  $d$  using Equation 13. Then GPU-PDMM samples each possible  $\mathbf{Z}_d$  using Equation 14. Finally, GPU-PDMM sets all the values of  $z_{d,w}$  based on the updates  $\mathbf{Z}_d$ , updates  $\mathbb{S}_{d,w}$  for GPU using Equation 8, and then updates TF vector for each word using the addible property.

Here, due to the computational costs involved in sampling  $\mathbf{Z}_d$ , GPU-PDMM only samples the more relevant topics for each document. Specifically, GPU-PDMM infers the topic probability  $p(z|d)$  of each document  $d$  using the follows,

$$p(z = k | d) \propto \sum_{w \in d} p(z = k | w) p(w | d)$$

where  $p(w | d) = \frac{n_d^w}{n_d}$ . GPU-PDMM only chooses the top  $M$  topics for document  $d$  based on the probability  $p(z|d)$  to generate  $\mathbf{Z}_d$ , where  $\varsigma < M \leq K$ . The topic-word distribution can be calculated by Equation 12.

### 3.4 Global Word Co-occurrences based Methods

The closer the two words, the more relevance the two words. Utilizing this idea, global word co-occurrences based methods learn the latent topics from the global word co-occurrences obtained from the original corpus. This type of methods needs to set sliding window for extracting word co-occurrences. In general, if the average length of each document is larger than 10, they use sliding window and set the size of the sliding window as 10, else they can directly take each document as a sliding window.

#### 3.4.1 BTM

BTM [14] first generate biterns from the corpus  $\mathcal{D}$ , where any two words in a document is treated as a bitern. Suppose that the corpus contains  $n_b$  biterns  $\mathbf{B} = \{b_i\}_{i=1}^{n_b}$ , where  $b_i = (w_{i,1}, w_{i,2})$ . BTM infers topics over the biterns  $\mathcal{B}$ . The generative process of BTM is described as follows, and its graphical model is shown in Figure 4.

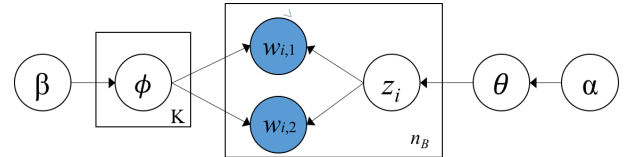


Fig. 4. Graphical Model of BTM.

- (1). Draw  $\theta \sim \text{Dirichlet}(\alpha)$ .
- (2). For each topic  $k \in [1, K]$ 
  - (a) draw  $\phi_k \sim \text{Dirichlet}(\beta)$ .



- (3). For each biterm  $b_i \in \mathbf{B}$ 
  - (a) draw  $z_i \sim \text{Multinomial}(\theta)$ ,
  - (b) draw  $w_{i,1}, w_{i,2} \sim \text{Multinomial}(\phi_{z_i})$ .

The TF vector of a topic  $k$  in BTM is defined as a tuple  $\{n_k(w \in W), n_k\}$ . The addible and deletable properties of topic feature (TF) in BTM are described below.

(1) **Addible Property.** A biterm  $b_i$  can be efficiently added to topic  $k$  by updating its TF vector as follows.

$$n_k^{w_{i,1}} = n_k^{w_{i,1}} + 1; n_k^{w_{i,2}} = n_k^{w_{i,2}} + 1; n_k = n_k + 1$$

(2) **Deletable Property.** A biterm  $b_i$  can be efficiently deleted from topic  $k$  by updating its TF vector as follows.

$$n_k^{w_{i,1}} = n_k^{w_{i,1}} - 1; n_k^{w_{i,2}} = n_k^{w_{i,2}} - 1; n_k = n_k - 1$$

Using the technique of collapsed Gibbs sampling, BTM samples the topic  $z_i$  of biterm  $b_i$  using the following conditional distribution,

$$p(z_i = k | \mathbf{Z}_{-i}, \mathbf{B}) \propto (n_{k,-i} + \alpha) \times \frac{(n_{k,-i}^{w_{i,1}} + \beta)(n_{k,-i}^{w_{i,2}} + \beta)}{(n_{k,-i} + V\beta + 1)(n_{k,-i} + V\beta)} \quad (15)$$

where  $\mathbf{Z}_{-i}$  denotes the topics for all biterms except the current biterm  $b_i$ , and  $n_k$  is the number of biterms assigned to topic  $k$ .

For each biterm, we first delete it from its current TF vector with the deletable property. Then, we reassign the biterm to a topic using Equation 4. Accordingly, we update the new TF vector with the addible property. After finishing the iterations, BTM estimates  $\phi$  and  $\theta$  as follows,

$$\phi_k^w = \frac{n_k^w + \beta}{n_k + V\beta}, \quad (16)$$

$$\theta_d^k = \sum_{i=1}^{n_d^b} p(z_i = k) \quad (17)$$

where  $\theta_d^k$  is the probability of topic  $k$  in document  $d$ , and  $n_d^b$  is the number of biterms in document  $d$ .

### 3.4.2 WNTM

WNTM [32] uses global word co-occurrence to construct word co-occurrence network, and learns the distribution over topics for each word from word co-occurrence network using LDA. WNTM first set the size of a sliding window, and the window is moving word by word. Suppose window size is set as 10 in the original paper, if one document has 15 words, it will have 16 windows in this document. As WNTM scanning word by word in one window, two distinct words in the window are regarded as co-occurrence. WNTM construct undirected word co-occurrence network, where each node of the word co-occurrence network represents one word and the weight of each edge is the number of co-occurrence of the two connected words. We can see that the number of nodes is the number of vocabulary  $V$ .

Then, WNTM generates one pseudo-document  $l$  for each vertex  $v$  which is consisted of the adjacent vertices of this vertex in word network. The occur times of this adjacent vertex in  $l$  is determined by the weight of the edge. The number of words in  $l$

is the degree of the vertex  $v$  and the number of pseudo-documents  $P$  is the number of vertices.

After obtaining pseudo-documents  $P$ , WNTM adopts LDA to learn latent topics from pseudo-documents. Therefore, the topic feature (TF) in LF-DMM is same with LDA. For each word  $w$  in  $l$ , WNTM infers its topic using the following conditional distribution,

$$p(z_{l,w} = k | \mathbf{Z}_{-(l,w)}, P, \alpha, \beta) \propto (n_{l,-(l,w)}^k + \alpha) \frac{n_{k,-(l,w)}^w + \beta}{n_{k,-(l,w)} + V\beta} \quad (18)$$

where  $n_l^k$  is the number of topic  $k$  belonging to pseudo-document  $l$ , and  $-(l, w)$  means word  $w$  is removed from its pseudo-document  $l$ .

Because each pseudo-document is each word's adjacent word-list, the document-topic distribution learned from pseudo-document is the topic-word distribution in WNTM. Suppose pseudo-document  $l$  is generated from word  $w$ , the topic-word distribution of  $w$  is calculated using the following Equation,

$$\phi_k^w = \frac{n_l^k + \alpha}{n_l + K\alpha} \quad (19)$$

where  $n_l$  is the number of words in  $l$ .

Given topic-word distribution, the document-word distribution  $\theta_d$  can be calculated as,

$$\theta_d^k = \sum_{i=1}^{n_d} \phi_k^{w_{d,i}} p(w_{d,i} | d)$$

$$p(w_{d,i} | d) = \frac{n_d^{w_{d,i}}}{n_d}$$

where  $n_d^{w_{d,i}}$  is the number of word  $w_{d,i}$  in document  $d$ .

## 3.5 Self-aggregation based Methods

Self-aggregation based methods alleviate the problem of sparseness by merging short texts into long pseudo-documents  $P$  before inferring the latent topics [15], [17], [36]. The previous self-aggregation based methods first merged short texts, and then applied topic models. Recently, SATM and PTM simultaneously integrate clustering and topic modeling in one iteration. In general, the number of pseudo-documents  $|P|$  is significantly less than the number of short texts, namely  $|P| \ll N$ .

### 3.5.1 SATM

Self-aggregation based topic modeling (SATM) [20] supposes that each short text is sampled from an unobserved long pseudo-document, and infers latent topics from pseudo-documents using standard topic modeling. The Gibbs sampling process in SATM can be described in two indispensable steps.

The first step calculates the probability of the occurrence of a pseudo-document  $l$  in  $P$  conditioned on short document  $d$  in short corpus, which is estimated using the mixture of unigrams model [22],

w.d,i

$$p(l|d) = \frac{p(l) \prod_{i=1}^V (\frac{n_l^{w_{d,i}}}{n_l})^{n_d^{w_{d,i}}}}{\sum_{m=1}^{|P|} p(m) \prod_{i=1}^V (\frac{n_m^{w_{d,i}}}{n_m})^{n_d^{w_{d,i}}}} \quad (20)$$

where  $p(l) = \frac{n_l}{N}$  represents the probability of pseudo-document  $p_l$ ,  $n_l^{w_{d,i}}$  is the number of word  $w_{d,i}$  in pseudo-document  $p_l$ , and  $n_l$  is the number of words in  $p_l$ .

The second step estimates draws a pair of pseudo-document label  $l_{d,w}$  and topic label  $z_{d,w}$  jointly for word  $w$  in document  $d$ , which is similar with standard topic modeling (author-topic modeling) [42].

The addible and deletable properties of pseudo-document and topic feature (PTF) in SATM are described below.

(1) **Addible Property.** A word  $w$  can be efficiently added into pseudo-document  $l$  and topic  $k$  by updating its TPF vector as follows.

$$\begin{aligned} n_l^w &= n_l^w + 1; n_l^k = n_l^k + 1; n_l = n_l + 1 \\ n_k^w &= n_k^w + 1; n_k = n_k + 1 \end{aligned}$$

(2) **Deletable Property.** A word  $w$  can be efficiently deleted from pseudo-document  $l$  and topic  $k$  by updating its PTF vector as follows.

$$\begin{aligned} n_l^w &= n_l^w - 1; n_l^k = n_l^k - 1; n_l = n_l - 1 \\ n_k^w &= n_k^w - 1; n_k = n_k - 1 \end{aligned}$$

The pair of pseudo-document label  $l_{d,w}$  and topic label  $z_{d,w}$  jointly for word  $w$  in document  $d$  can be calculated by,

$$\begin{aligned} p(l_{d,w} = l, z_{d,w} = k | \mathbf{Z}_{-(d,w)}, P_{-(d,w)}) &\propto \\ p(l|d) \times \frac{n_l^k + \alpha}{n_l + K\alpha} \cdot \frac{n_k^w + \beta}{n_k + V\beta} \end{aligned} \quad (21)$$

where  $n_l^k$  is the number of words in pseudo-document  $l$  belonging to topic  $k$ .

After finishing the iterations, SATM estimates  $\phi$  and  $\theta$  as follows,

$$\phi_k^w = \frac{n_k^w + \beta}{n_k + V\beta}, \quad (22)$$

$$\theta_d^k = \prod_{i=1}^{n_d} \phi_k^{w_{d,i}} \quad (23)$$

### 3.5.2 PTM

The pseudo-document-based topic modeling (PTM) [35] supposes that each short text is sampled from one long pseudo-document  $p_l$ , and then infers the latent topics from long pseudo-documents  $P$ . A multinomial distribution  $\varphi$  is used to model the distribution of short texts over pseudo-documents. The graphical model of PTM is shown in Figure 5. The generative process of PTM is described as follows,

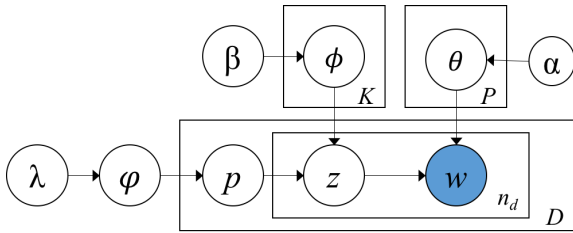


Fig. 5. Graphical Model of PTM.

(1). Sample  $\varphi \sim \text{Dir}(\lambda)$

- (2). For each topic  $k \in [1, K]$ 
  - (a) draw  $\phi_k \sim \text{Dirichlet}(\beta)$ .
- (3). For each pseudo-document  $l$ 
  - (a) sample  $\theta_l \sim \text{Dir}(\alpha)$
- (4). For each document  $d \in \mathbf{D}$ :
  - (a) Sample a pseudo-document  $l \sim \text{Multinomial}(\varphi)$ .
  - (b) For each word  $w \in \{w_{d,1}, \dots, w_{d,n_d}\}$  in  $d$ :
    - (i) Sample a topic  $z \sim \text{Multinomial}(\theta_l)$ .
    - (ii) Sample a word  $w \sim \text{Multinomial}(\phi_z)$ .

The addible and deletable properties of pseudo-document and topic feature (PTF) in PTM are described below.

(1) **Addible Property.** A document  $d$  can be efficiently added into pseudo-document  $l$  by updating its PTF vector as follows.

$$\begin{aligned} n_l^k &= n_l^k + 1 \quad \text{for } z_{d,w} = k \text{ in } d \\ m_l &= m_l + 1; n_l = n_l + n_d \end{aligned}$$

(2) **Deletable Property.** A document  $d$  can be efficiently deleted from pseudo-document  $l$  by updating its PF vector as follows.

$$\begin{aligned} n_l^k &= n_l^k - 1 \quad \text{for } z_{d,w} = k \text{ in } d \\ m_l &= m_l - 1; n_l = n_l - n_d \end{aligned}$$

Integrating out  $\theta, \phi$  and  $\varphi$ , the pseudo-document assignment  $l$  for short text  $d$  based on collapsed Gibbs sampling can be estimated as follows,

$$\begin{aligned} p(l_d = l | \overrightarrow{P_{-d}}, \mathbf{D}) &\propto \\ \frac{m_{l,-d}}{N - 1 + \lambda|P|} \frac{\prod_{k \in d} \prod_{j=1}^{n_d} (n_{l,-d}^k + \alpha + j - 1)}{\prod_{i=1}^{n_d} (n_{l,-d} + K\alpha + i - 1)} \end{aligned} \quad (24)$$

where  $m_l$  is the number of short texts associated with pseudo-document  $l$ ,  $n_l^k$  is the number of words associated with topic  $k$  in pseudo-document  $l$ .

After obtaining the pseudo-document for each short text, PTM samples the topic assignment for each word  $w$  in document  $d$ . That is,

$$p(z_{d,w} = k | \mathbf{Z}_{-(d,w)}, \mathbf{D}) \propto (n_l^k + \alpha) \frac{n_k^w + \beta}{n_k + V\beta} \quad (25)$$

where  $n_l^k$  is the number of words associated with topic  $k$  in pseudo-document  $l$ .

The document-word distribution  $\theta_d$  can be calculated as,

$$\theta_d^k = \frac{n_d^k + \alpha}{n_d + K\alpha} \quad (26)$$

## 4 APPLICATIONS

With the emerging of social media, topic models have been used for social media content analysis, such as content characterizing and recommendation, text classification, event tracking. However, although the corpus is composed of short texts, some previous work directly applied traditional topic models for topic discovery, since no specific short text topic models were proposed at that time. Therefore, it brings a new chance for short text topic modeling to improve the performance of these tasks.



#### 4.1 Content characterizing and recommendation

Microblogging sites are used as publishing platforms to create and consume content from sets of users with overlapping and disparate interests, which results in many contents are useless for users. These work [19], [43] have been devoted to content analysis of Twitter. Ramage et al. [43] used topic models to discover latent topics from the tweets that can be roughly categorized into four types: substance topics about events and ideas, social topics recognizing language used toward a social end, status topics denoting personal updates, and style topics that embody broader trends in language usage. Next, they characterize selected Twitter users along these learned dimensions for providing interpretable summaries or characterizations of users tweet streams. Zhao et al. [19] performed content analysis on tweets using topic modeling to discover the difference between Twitter and traditional medium.

Content analysis is crucial for content recommendation for microblogging users [44], [45]. Phelan et al. [46] identified emerging topics of interest from Twitter information using topic modeling, and recommended news by matching emerging topics and recent news coverage in an RSS feed. Chen et al. [47] also studied content recommendation based on Twitter for better capture users' attention by exploring three separate dimensions in designing such a recommender: content sources, topic interest models for users, and social voting.

#### 4.2 Text Classification

Topic models for text classification are mainly from the following two aspects. The first one is topics discovered from external large-scale data corpora are added into short text documents as external features. For example, Phan et al. [16] built a classifier on both a set of labeled training data and a set of latent topics discovered from a large-scale data collection. Chen et al. [48] integrated multi-granularity hidden topics discovered from short texts and produced discriminative features for short text classification. Vo et al. [49] explored more external large-scale data collections which contain not only Wikipedia but also LNCS and DBLP for discovering latent topics.

The other one is that topic models are used to obtain a low dimensional representation of each text, and then classify text using classification methods [50], [51]. Compared with traditional statistical methods, the representation using topic models can get a compact, dense and lower dimensional vector in which each dimension of the vector usually represents a specific semantic meaning (e.g., a topic) [22]. Dai et al. [50] used the topic information from training data to extend representation for short text. Recent topic modeling methods on text representation have explicitly evaluated their models on this task. They showed that a low dimensional representation for each text suffices to capture the semantic information.

#### 4.3 Event Tracking

Nowadays, a large volume of text data is generated from the social communities, such as blogs, tweets, and comments. The important task of event tracking is to observe and track the popular events or topics that evolve over time [52]–[54]. Lin et al. [52] proposed a novel topic modeling that models the popularity of events over time, taking into consideration the burstiness of user interest, information diffusion in the network structure, and the evolution of latent topics. Lau et al. [55] designed a novel topic modeling for event detecting, whose model has an in-built update

mechanism based on time slices by implementing a dynamic vocabulary.

### 5 A JAVA LIBRARY FOR SHORT TEXT TOPIC MODELING

We released an open-source Java library, STTM (Short Text Topic Modeling)<sup>2</sup>, which is the first comprehensive open-source library, which not only includes the state-of-the-art algorithms with a uniform easy-to-use programming interface but also includes a great number of designed modules for the evaluation and application of short text topic modeling algorithms. The design of STTM follows three basic principles. (1) Preferring integration of existing algorithms rather than implementing them. If the original implementations are open, we always attempt to integrate the original codes rather than implement them. The work that we have done is to consolidate the input/output file formats and package these different approaches into some newly designed java classes with a uniform easy-to-use member functions. (2) Including traditional topic modeling algorithms for long texts. The classical topic modeling algorithm (LDA [9] and its variation LF-LDA [28]) are integrated, which is easy for users to the comparison of long text topic modeling algorithms and short text topic modeling algorithms. (3) Extendibility. Because short text topic modeling is an emerging research field, many topics have not been studied yet. For incorporating future work easily, we try to make the class structures as extendable as possible when designing the core modules of STTM.

Figure 6 shows the hierarchical architecture of STTM. STTM supports the entire knowledge discovery procedure including analysis, inference, evaluation, application for classification and clustering. In the data layer, STTM is able to read a text file, in which each line represents one document. Here, a document is a sequence of words/tokens separated by whitespace characters. If we need to evaluate the algorithm, we also need to read a gold file, in which each line is the class label of one document. STTM provides implementations of DMM [18], LF-DMM [28], GPU-DMM [29], GPU-PDMM [31], BTM [14], WNTM [32], SATM [20], and PTM [35]. For each model, we not only provide how to train a model on existing corpus but also give how to infer topics on a new/unseen corpus using a pre-trained topic model. In addition, STTM presents three aspects of how to evaluate the performance of the algorithms (i.e., topic coherence, clustering, and classification). For topic coherence, we use the point-wise mutual information (PMI) to measure the coherence of topics [35]. Short text topic modeling algorithms are widely used for clustering and classification. In the clustering module, STTM provides two measures (NMI and Purity) [21]. Based on the latent semantic representations learned by short text topic modeling, accuracy metric is chosen in classifications [14].

### 6 EXPERIMENTAL SETUP

In this section, we specify the parameter setting of the introduced short text topic models, dataset and evaluation metrics we used. All of these are implemented in our library STTM. The experiments were performed on a Ubuntu 18.04(bionic) system with 6 cores, Intel Xeon E5645 CPU and 12288 KB cache.

For all models in comparison, we used the recommended setting by the authors and set the number of iterations as 2,000

2. <https://github.com/qiang2100/STTM>

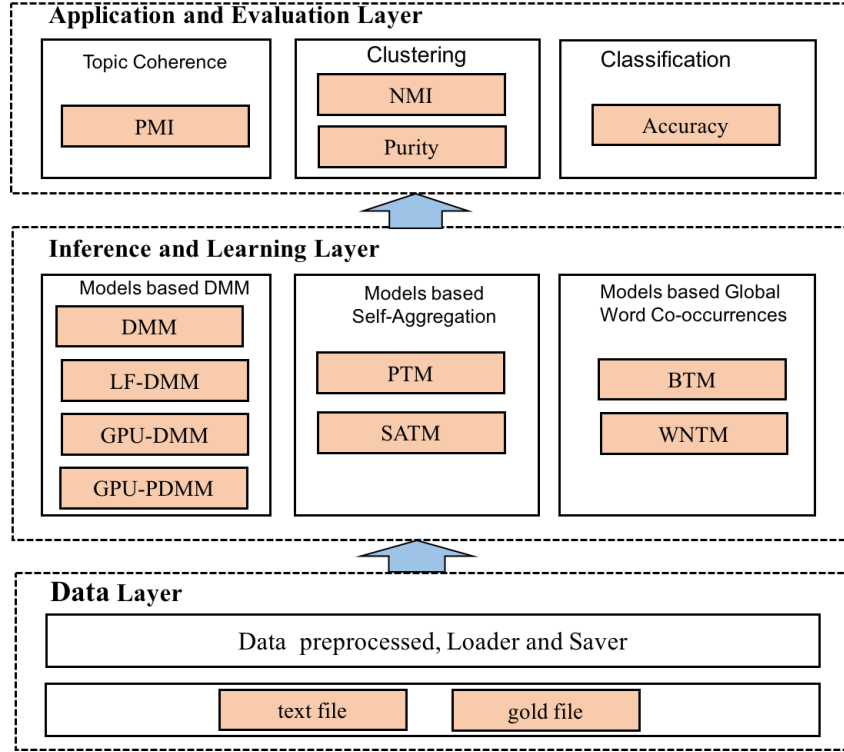


Fig. 6. The architecture of STTM

unless explicitly specified elsewhere. The word embeddings of LF-DMM, GPU-DMM, and GPU-PDMM are trained by Glove [56]. In this paper, we used the pre-trained word embeddings "glove.6B.200d.txt", where the dimension of the vector is 200.

### 6.1 Parameter Setting

**LDA:** LDA is the most popular and classic topic modeling. We choose it as a baseline to the comparison. The hyper-parameters of LDA are set as  $\alpha = 0.05$  and  $\beta = 0.01$  that are proved in the paper (BTM). The authors tuned parameters via grid search on the smallest collection to get the best performance.

**GSDMM:** We set  $k = 300$ ,  $\alpha = 0.1$  and  $\beta = 0.1$  declared in the paper(GSDMM).

**LF-DMM:** We set  $\lambda = 0.6$ ,  $\alpha = 0.1$  and  $\beta = 0.01$  shown in their paper. We set the iterations for baseline models as 1,500 and ran the further iterations 500 times.

**GPU-DMM:** We use the hyper-parameter settings provided by the authors,  $\alpha = 50/k$ ,  $\beta = 0.01$ . The number of iterations is 1,000 in the paper.

**GPU-PDMM:** All the settings are same as the model GPU-DMM. We set  $\lambda = 1.5$ ,  $\zeta = 1.5$ , and  $M=10$  declared in the original paper.

**BTM:** The parameters  $\alpha = 50/K$  and  $\beta = 0.01$  are used, the model gets optimal performance. Each document is treated as one window.

**WNTM:** We set  $\alpha = 0.1$  and  $\beta = 0.1$  used in the original paper. The window size is set as 10 words.

**SATM:** We set the number of pseudo-numbers as 300, and the hyper-parameters  $\alpha = 50/k$ ,  $\beta = 0.1$ . The number of iterations is set as 1,000.

**PTM:** The hyper-parameters are  $\alpha = 0.1$  and  $\beta = 0.01$ . We also set the number of pseudo-document as 1000.

### 6.2 Datasets

To show the effects and differences of the above nine models, we select the following six datasets to verify the models. After preprocessing these datasets, we present the key information of the datasets that are summarized in Table 4, where  $K$  corresponds to the number of topics per dataset,  $N$  represents the number of documents in each dataset,  $Len$  shows the average length and maximum length of each document, and  $V$  indicates the size of the vocabulary.

TABLE 4  
The basic information of the corpus

<i>Dataset</i>	<i>K</i>	<i>N</i>	<i>Len</i>	<i>V</i>
SearchSnippets	8	12,295	14.4/37	5,547
StackOverflow	20	16,407	5.03/17	2,638
Biomedicine	20	19,448	7.44/28	4498
Tweet	89	2,472	8.55/20	5,096
GoogleNews	152	11,109	6.23/14	8,110
PascalFlickr	20	4,834	5.37/19	3,431

**SearchSnippets:** Given the predefined phrases of 8 different domains, this dataset was chosen from the results of web search transaction. The 8 domains are Business, Computers, Culture-Arts, Education-Science, Engineering, Health, Politics-Society, and Sports, respectively.

**StackOverflow:** The dataset is released on Kaggle.com. The raw dataset contains 3,370,528 samples from July 31st, 2012 to August 14, 2012. Here, the dataset randomly selects 20,000 question titles from 20 different tags.

**Biomedicine:** Biomedicine makes use of the challenge data delivered on BioASQ's official website.

**Tweet:** In the 2011 and 2012 microblog tracks at Text REtrieval Conference (TREC), there are 109 queries for using. After

removing the queries with none highly-relevant tweets, Tweet dataset includes 89 clusters and totally 2,472 tweets.

**GoogleNews:** In the Google news site, the news articles are divided into clusters automatically. GoogleNews dataset is downloaded from Google news site on November 27, 2013, and crawled the titles and snippets of 11,109 news articles belonging to 152 clusters.

**PascalFlickr:** PascalFlickr dataset are a set of captions [57], which is used as evaluation for short text clustering [58].

### 6.3 Evaluation Metrics

It is still an open problem about how to evaluate short text topic models. A lot of metrics have been proposed for measuring the coherence of topics in texts [59], [60]. Although some metrics tend to be reasonable for long texts, they can be problematic for short texts [20]. Most conventional metrics (e.g., perplexity) try to estimate the likelihood of held-out testing data based on parameters inferred from training data. However, this likelihood is not necessarily a good indicator of the quality of the extracted topics [61]. To provide a good evaluation, we evaluate all models from many aspects using different metrics,

**Classification Evaluation:** Each document can be represented using document-topic distribution  $p(z|d)$ . Therefore, we can evaluate the performance of topic modeling using text classification. Here, we choose accuracy as a metric for classification. Higher accuracy means the learned topics are more discriminative and representative. We use a linear kernel Support Vector Machine (SVM) classifier in LIBLINEAR<sup>3</sup> with the default parameter settings. The accuracy of classification is computed through fivefold cross-validation on all datasets.

**Cluster Evaluation (Purity and NMI):** By choosing the maximum of topic probability for each document, we can get the cluster label for each text. Then, we can compare the cluster label and the golden label using metric Purity and NMI [18], [24].

**Topic Coherence:** Computing topic coherence, additional dataset (Wikipedia) as a single meta-document is needed to score word pairs using term co-occurrence in the paper (Automatic Evaluation of Topic Coherence). Here, we calculate the point-wise mutual information (PMI) of each word pair, estimated from the entire corpus of over one million English Wikipedia articles [31]. Using a sliding window of 10 words to identify co-occurrence, we computed the PMI of all a given word pair. The Wikipedia corpus can be downloaded here<sup>4</sup>. Then, we can transfer the dataset from HTML to text using the code in the STTM package. Finally, due to the large size, we only choose 1,000,000 sentences from it.

## 7 EXPERIMENTS AND ANALYSIS

In this section, we conduct experiments to evaluate the performance of the nine models. We run each model 20 times on each dataset and report the mean and standard deviation.

### 7.1 Classification Accuracy

Classification accuracy is used to evaluate document-topic distribution. We represent each document with its document-topic distribution and employ text classification method to assess. For DMM based methods, we use  $p(z_d = k)$  to represent each document. For other models, we adopt the giving equation  $\theta_d^k$ .

The classification accuracy on six datasets using nine models is shown in Figure 7. We observe that although the performance of methods is dataset dependent, DMM based methods which utilize word embeddings outperform others, especially on Tweet and GoogleNews datasets. This is because GoogleNews and Tweet are general (not domain-specific) datasets and word embeddings used in this paper are trained in general datasets. If we try to use these models (LF-DMM, GPU-DMM, and GPU-PDMM) on domain-specific datasets, we can further improve the performance by re-training word embeddings on domain-specific datasets.

We also observe that self-aggregation based methods are unable to achieve high accuracy, especially the SATM method. The performance of self-aggregation based methods is affected by generating long pseudo-documents. Without any auxiliary information or metadata, the error of this step of generating pseudo-documents will be amplified in the next step.

In conclusion, these models based on the simple assumption (BTM and GSDMM) always outperform than LDA in all datasets, which indicate that two words or all words in one document are very likely to from one topic. Here we can see that the performance of other models (LF-DMM, GPU-DMM, GPU-PDMM, WNTM) are highly data set dependent. For example, WNTM achieves good performance on Tweet, GoogleNews and StackOverflow, but performs poorly on other data sets. GPU-PDMM achieves the best performance on all data sets, except SearchSnippets.

### 7.2 Clustering

Another important application of short text topic modeling is short text clustering. For each document, we choose the maximum value from its topic distribution as the cluster label. We report the mean value of each modeling on all datasets in the last column. The best results for each dataset using each metric are highlighted in bold.

Table 5 illustrates the results using cluster metrics. We can see that all models outperform long text topic modeling (LDA), except SATM. Here, similar to the conclusions in classification, we can observe that the performance of approaches is highly data set dependent. WNTM achieves the best performance on several datasets but performs poorly on PascalFlickr. GPU-PDMM performs very well on all datasets except SearchSnippets.

For self-aggregation based methods, PTM performs better than SATM. For global word-occurrences based methods, two methods are very different from each other. WNTM performs better than BTM on Tweet and StackOverflow, and BTM achieves good performance on GoogleNews and PascalFlickr. For DMM based methods, GSDMM without incorporating word embeddings outperforms other methods on Biomedicine and SearchSnippets.

### 7.3 Topic Coherence

Topic coherence is used to evaluate the quality of topic-word distribution. Here, we only choose the top 10 words for each topic based on the word probability. The results are shown in Figure 8. DMM based methods achieve the best performance on all datasets. LF-DMM has the best performance on four datasets (Biomedicine, GoogleNews, SearchSnippets, and Tweet), GPU-DMM has the best performance on StackOverflow, and GPU-PDMM achieves the best on PascalFlickr. It means that incorporating word embeddings into DMM can help to alleviate the sparseness. Two methods based on global word co-occurrences perform very well and achieve a similar result on each dataset,

3. <https://liblinear.bwaldvogel.de/>

4. <https://dumps.wikimedia.org/enwiki/>

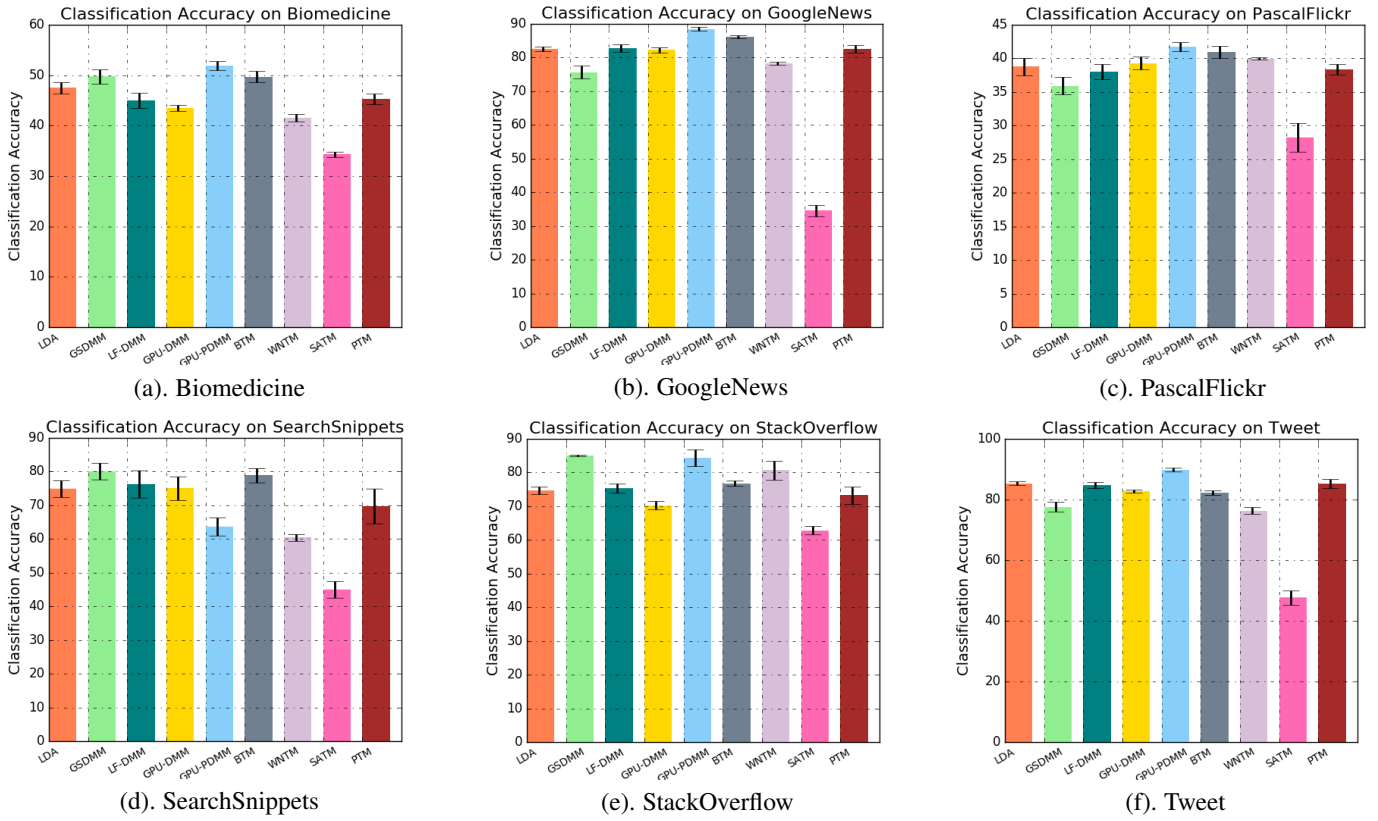


Fig. 7. Average classification accuracy of all models on six datasets.

which indicates that the adequacy of global word co-occurrences can mitigate the sparsity of short texts. Similar to the above results using other metrics, self-aggregation based methods perform very poorly.

We also present the qualitative evaluation of latent topics. Here, we choose SearchSnippets dataset as an example, since it only contains eight topics that are Health, Politics-Society (politics), Engineering (engine.), Culture-Arts (culture), Sports, Computers, Business, and Education-Science (education). Table 6 shows the eight topics learned by the nine models. Each topic is visualized by the top ten words. Words that are noisy and lack of representativeness are highlighted in bold.

From Table 6, we observe that LF-DMM can achieve a similar conclusion with topic coherence, which can learn more coherent topics with fewer noisy and meaningless words. GPU-DMM and GPU-PDMM can not discriminate the topic 'Engineering'. SATM remains the worst method in all short text topic models, which cannot discriminate three topics 'Engineering', 'Politics-Society' can 'Culture. Except LDA, PTM, WNTM, and SATM, other models can identify at least seven topics from all eight topics.

#### 7.4 Influence of the number of iterations

In this subsection, we try to investigate the influence of the number of iterations to the performance of all models using NMI metric. Since all models have converged when the number of iterations reaches 2000, we vary the number of iterations from 2 to 2024.

The results are shown in Figure 9. We can see that models based DMM can converge fast to the optimal solutions and almost get stable within 30 iterations. Models based global word co-

occurrences get stable within 60 iterations. Models based self-aggregation has the slowest convergence speed and the lowest iterative performance.

#### 7.5 Efficiency

In this part, we compare the efficiency of various short text topic models. Here, we choose the largest dataset "Biomedicine" from all datasets to do the experiments.

The average runtime of the initiation and per iteration for each model are reported in Table 7. Among all models evaluated, LDA and DMM are the most efficient methods as expected. GPU-DMM is slightly slower than DMM and LDA, due to a similar Gibbs sampling process with GSDMM. LF-DMM and GPU-PDMM take much more time than GPU-DMM, because GPU-PDMM spends more time for the computational costs involved in sampling  $\mathbf{Z}_d$  and LF-DMM need much time for optimizing the topic vectors. We can see that GPU-PDMM is the slowest modeling compared with other models.

Global word co-occurrences based methods are much slower than GSDMM, LDA and GPU-DMM and faster than the rest models. This is expected since they extend the number of words by extracting word co-occurrences. For self-aggregation based methods, the time is affected by the number of pseudo-documents. PTM is much faster than SATM but much slower than global word co-occurrences based methods. In addition, the models by incorporating word embeddings (GPU-DMM, LF-DMM, and GPU-PDMM) have the slowest time for the initiation due to the computational cost for the similarity between words.

TABLE 5  
Purity and NMI value of all models on six datasets.

Model		Biome dicine	Google News	Pascal Flickr	Search Snippets	Stack Overflow	Tweet	Mean Value
LDA	Purity	0.456 $\pm$ 0.011	0.793 $\pm$ 0.005	0.376 $\pm$ 0.013	0.740 $\pm$ 0.013	0.562 $\pm$ 0.013	0.821 $\pm$ 0.006	0.625 $\pm$ 0.013
	NMI	0.356 $\pm$ 0.004	0.825 $\pm$ 0.002	0.321 $\pm$ 0.006	0.517 $\pm$ 0.025	0.425 $\pm$ 0.006	0.805 $\pm$ 0.004	0.542 $\pm$ 0.008
GSDMM	Purity	<b>0.494</b> $\pm$ <b>0.011</b>	0.754 $\pm$ 0.014	0.360 $\pm$ 0.012	<b>0.801</b> $\pm$ <b>0.024</b>	0.713 $\pm$ 0.002	0.785 $\pm$ 0.011	0.650 $\pm$ 0.013
	NMI	<b>0.396</b> $\pm$ <b>0.006</b>	0.851 $\pm$ 0.004	0.317 $\pm$ 0.005	<b>0.608</b> $\pm$ <b>0.023</b>	0.593 $\pm$ 0.002	0.801 $\pm$ 0.007	<b>0.590</b> $\pm$ <b>0.001</b>
LF-DMM	Purity	0.421 $\pm$ 0.019	0.828 $\pm$ 0.009	0.381 $\pm$ 0.009	0.762 $\pm$ 0.042	0.518 $\pm$ 0.0217	0.856 $\pm$ 0.009	0.630 $\pm$ 0.018
	NMI	0.348 $\pm$ 0.005	0.875 $\pm$ 0.005	0.365 $\pm$ 0.007	0.579 $\pm$ 0.026	0.443 $\pm$ 0.007	0.843 $\pm$ 0.006	0.578 $\pm$ 0.009
GPU-DMM	Purity	0.433 $\pm$ 0.008	0.818 $\pm$ 0.005	<b>0.395</b> $\pm$ <b>0.010</b>	0.751 $\pm$ 0.035	0.511 $\pm$ 0.013	0.830 $\pm$ 0.006	0.623 $\pm$ 0.013
	NMI	0.366 $\pm$ 0.006	0.852 $\pm$ 0.002	<b>0.370</b> $\pm$ <b>0.004</b>	0.561 $\pm$ 0.026	0.429 $\pm$ 0.003	0.810 $\pm$ 0.006	0.565 $\pm$ 0.008
GPU-PDMM	Purity	0.481 $\pm$ 0.011	<b>0.860</b> $\pm$ <b>0.002</b>	0.368 $\pm$ 0.008	0.537 $\pm$ 0.030	0.702 $\pm$ 0.032	<b>0.869</b> $\pm$ <b>0.005</b>	0.636 $\pm$ 0.015
	NMI	0.381 $\pm$ 0.005	0.871 $\pm$ 0.001	0.322 $\pm$ 0.003	0.341 $\pm$ 0.014	0.607 $\pm$ 0.013	0.830 $\pm$ 0.003	0.559 $\pm$ 0.007
BTM	Purity	0.458 $\pm$ 0.012	0.849 $\pm$ 0.005	0.392 $\pm$ 0.011	0.765 $\pm$ 0.032	0.537 $\pm$ 0.019	0.814 $\pm$ 0.008	0.636 $\pm$ 0.014
	NMI	0.380 $\pm$ 0.004	0.875 $\pm$ 0.003	0.368 $\pm$ 0.006	0.566 $\pm$ 0.027	0.456 $\pm$ 0.008	0.808 $\pm$ 0.005	0.575 $\pm$ 0.009
WNTM	Purity	0.472 $\pm$ 0.009	0.837 $\pm$ 0.007	0.324 $\pm$ 0.005	0.712 $\pm$ 0.016	<b>0.750</b> $\pm$ <b>0.026</b>	0.856 $\pm$ 0.012	<b>0.658</b> $\pm$ <b>0.013</b>
	NMI	0.369 $\pm$ 0.004	<b>0.876</b> $\pm$ <b>0.004</b>	0.295 $\pm$ 0.003	0.464 $\pm$ 0.011	<b>0.659</b> $\pm$ <b>0.006</b>	<b>0.850</b> $\pm$ <b>0.009</b>	0.585 $\pm$ 0.006
SATM	Purity	0.384 $\pm$ 0.007	0.654 $\pm$ 0.008	0.237 $\pm$ 0.059	0.459 $\pm$ 0.055	0.505 $\pm$ 0.019	0.392 $\pm$ 0.011	0.438 $\pm$ 0.027
	NMI	0.27 $\pm$ 0.001	0.76 $\pm$ 0.005	0.186 $\pm$ 0.049	0.205 $\pm$ 0.036	0.366 $\pm$ 0.011	0.507 $\pm$ 0.006	0.382 $\pm$ 0.018
PTM	Purity	0.425 $\pm$ 0.012	0.807 $\pm$ 0.010	0.359 $\pm$ 0.012	0.674 $\pm$ 0.057	0.481 $\pm$ 0.034	0.839 $\pm$ 0.007	0.597 $\pm$ 0.022
	NMI	0.353 $\pm$ 0.003	0.866 $\pm$ 0.005	0.336 $\pm$ 0.010	0.457 $\pm$ 0.045	0.442 $\pm$ 0.016	0.846 $\pm$ 0.006	0.550 $\pm$ 0.014

TABLE 7  
The average runtime of initiation and per iteration of each model on Biomedicine (in milliseconds).

Model	Initiation time	Per iteration time
LDA	77	41.50
GSDMM	46	48.15
LF-DMM	3329	2243.03
GPU-DMM	13610	53.83
GPU-PDMM	13037	9685.44
BTM	320	160.43
WNTM	192	220.12
SATM	41	2015.03
PTM	126	818.32

## 8 CONCLUSION AND FUTURE WORK

The review of short text topic modeling (STTM) techniques covered three broad categories of methods: DMM based, global word co-occurrences based, and self-aggregation based. We studied the structure and properties preserved by various topic modeling algorithms and characterized the challenges faced by short text topic modeling techniques in general as well as each category of approaches. We presented various applications of STTM including content characterizing and recommendation, text classification, and event tracking. We provided an open-source Java library, named STTM, which is consisted of short text topic modeling approaches surveyed and evaluation tasks including classification,

clustering, and topic coherence. Finally, we evaluated the surveyed approaches to these evaluation tasks using six publicly available real datasets and compared their strengths and weaknesses.

Short text topic modeling is an emerging field in machine learning, and there are many promising research directions: (1) Visualization: as shown in the survey, we display topics by listing the most frequent words of each topic (see Figure 6). This new ways of labeling the topics may be more reasonable by either choosing different words or displaying the chosen words differently [62], [63]. How to display a document using topic models is also a difficult problem? For each document, topic modeling provides useful information about the structure of the document. Binding with topic labels, this structure can help to identify the most interesting parts of the document. (2) Evaluation: Useful evaluation metrics for topic modeling algorithms have never been solved [10]. Topic coherence cannot distinguish the differences between topics. In addition, existing metrics only evaluate one part of topic modeling algorithms. One open direction for topic modeling is to develop new evaluation metrics that match how the methods are used. (3) Model checking. From the experimental results on this paper, each method has different performance on different datasets. When dealing with a new corpus or a new task, we cannot decide which topic modeling algorithms should I use. How can I decide which of the many modeling assumptions are suitable for my goals? New computational answers to these questions would be a significant contribution to topic modeling.

TABLE 6  
The top ten words of each topic by each model on SearchSnippets dataset.

LDA				GSDMM				LF-DMM			
Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)	Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)	Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)
health <b>information</b> cancer <b>gov</b> medical <b>news</b> research disease healthy nutrition	<b>wikipedia</b> <b>encyclopedia</b> <b>wiki</b> <b>political</b> <b>culture</b> <b>democracy</b> <b>system</b> <b>party</b> <b>republic</b> <b>philosophy</b>	car engine electrical <b>com</b> products digital home motor energy calorie	music movie <b>com</b> <b>news</b> film movies yahoo art video arts	health <b>information</b> cancer gov medical <b>news</b> healthy disease nutrition hiv	political culture culture democracy party war republic <b>information</b> government	car engine <b>com</b> electrical <b>wikipedia</b> system wheels <b>olympic</b> digital <b>trade</b>	movie <b>com</b> art film fashion <b>motor</b> <b>wikipedia</b> books arts movies	health cancer disease healthy medical drug treatment physical food care	culture party democratic war political democracy congress presidential communist philosophy	motor engine wheels electronics electric cars models <b>phone</b> <b>graduation</b> <b>fashion</b>	film music art <b>com</b> fashion movie books arts rock band
Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)	Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)	Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)
com <b>news</b> sports football games <b>amazon</b> game soccer world tennis	computer business web software <b>com</b> <b>news</b> <b>market</b> <b>stock</b> internet programming	<b>gov</b> <b>business</b> <b>information</b> <b>school</b> <b>trade</b> <b>edu</b> <b>research</b> <b>home</b> <b>law</b> <b>economic</b>	edu science theory <b>journal</b> theoretical physics computer <b>information</b> university	<b>news</b> <b>music</b> <b>com</b> sports football games <b>movie</b> game <b>wikipedia</b> tennis	computer software web programming <b>wikipedia</b> memory <b>com</b> intel internet data	business market <b>news</b> <b>information</b> stock gov <b>com</b> finance services <b>home</b> <b>computer</b>	research edu science theory <b>information</b> school university journal physics <b>computer</b>	sports football games <b>news</b> league <b>com</b> hockey game soccer golf	computer intel software device linux digital network hardware web computers	business financial bank economic trade <b>news</b> market services law stock	research edu graduate resources science university school faculty center national
GPU-DMM				GPU-PDMM				BTM			
Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)	Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)	Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)
health <b>information</b> cancer medical gov <b>news</b> healthy disease nutrition hiv	political culture democracy <b>wikipedia</b> party system <b>information</b> government <b>news</b> gov	<b>theory</b> <b>theoretical</b> <b>physics</b> <b>wikipedia</b> <b>edu</b> <b>information</b> <b>science</b> <b>research</b> <b>amazon</b> <b>com</b>	movie music <b>com</b> film <b>news</b> movies <b>wikipedia</b> art fashion <b>amazon</b>	health gov cancer medical disease healthy nutrition physical hiv diet	<b>wikipedia</b> <b>encyclopedia</b> <b>wiki</b> political system democracy party government gov war	<b>com</b> <b>news</b> <b>information</b> <b>home</b> <b>online</b> <b>world</b> <b>web</b> <b>music</b> <b>index</b> <b>amazon</b>	culture art history <b>car</b> arts <b>imdb</b> museum <b>income</b> literature	health <b>information</b> gov cancer medical <b>news</b> research disease healthy nutrition	political culture democracy <b>encyclopedia</b> party <b>wiki</b> system government war	car engine <b>intel</b> <b>com</b> digital motor wheels products automatic	movie music <b>com</b> film <b>news</b> movies books art video
Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)	Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)	Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)
sports <b>news</b> football <b>com</b> games soccer game <b>wikipedia</b> tennis world	computer web software programming <b>com</b> <b>wikipedia</b> memory intel linux digital	business market <b>news</b> trade stock <b>information</b> <b>com</b> services finance <b>home</b>	research edu science school journal university computer <b>information</b> department graduate	<b>movie</b> sports games football <b>yahoo</b> game <b>video</b> soccer film <b>movies</b>	computer software programming systems memory <b>engine</b> intel design electrical security	business trade management economic law international gov products jobs bank	research edu science theory school university journal theoretical physics department	sports <b>news</b> football <b>com</b> games game soccer match world tennis	computer software web programming <b>com</b> internet memory data wikipedia linux	business <b>news</b> market <b>information</b> trade stock services <b>home</b> gov finance	edu research science theory <b>information</b> university journal school theoretical physics
WNTM				SATM				PTM			
Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)	Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)	Topic1 (health)	Topic2 (politics)	Topic3 (engine.)	Topic4 (culture)
health <b>information</b> cancer hiv healthy nutrition disease medical news diet	political <b>wikipedia</b> culture <b>encyclopedia</b> democracy party government war world	<b>music</b> <b>engine</b> car <b>rock</b> <b>com</b> <b>motor</b> <b>reviews</b> <b>pop</b> <b>band</b> <b>wikipedia</b>	movie <b>com</b> <b>amazon</b> film art <b>books</b> <b>fashion</b> <b>online</b> <b>movies</b> <b>video</b>	health <b>information</b> <b>news</b> research gov medical <b>party</b> <b>home</b> disease healthy	<b>wikipedia</b> <b>research</b> <b>trade</b> <b>information</b> <b>wiki</b> <b>gov</b> <b>international</b> <b>journal</b> <b>programming</b> <b>business</b>	<b>culture</b> <b>amazon</b> <b>wikipedia</b> <b>democracy</b> <b>books</b> <b>com</b> <b>political</b> <b>history</b> <b>edu</b> <b>encyclopedia</b>	movie <b>news</b> <b>com</b> film <b>wikipedia</b> movies <b>reviews</b> <b>online</b> <b>digital</b> <b>articles</b>	health <b>information</b> cancer medical gov <b>news</b> disease healthy nutrition hiv	<b>wikipedia</b> <b>encyclopedia</b> <b>wiki</b> <b>political</b> <b>democracy</b> <b>system</b> <b>party</b> <b>war</b> <b>government</b> <b>house</b>	car engine <b>theory</b> <b>com</b> <b>theoretical</b> <b>physics</b> <b>books</b> <b>car</b> <b>engine</b> <b>models</b> <b>electrical</b>	music movie <b>com</b> <b>news</b> film movies video reviews <b>intel</b> <b>imdb</b>
Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)	Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)	Topic5 (sport)	Topic6 (computer)	Topic7 (business)	Topic8 (education)
sports football <b>news</b> games <b>com</b> soccer game tennis match league	computer web programming intel memory internet systems <b>com</b> data wikipedia	business market trade stock <b>news</b> <b>information</b> jobs finance <b>home</b> tax	research science edu school theory <b>information</b> journal university theoretical physics	sports games <b>com</b> <b>news</b> football game world soccer <b>online</b> <b>wikipedia</b>	system <b>com</b> web computer <b>information</b> <b>car</b> wikipedia memory device <b>engine</b>	business <b>news</b> <b>com</b> market yahoo stock internet services financial <b>information</b>	edu science research school university <b>program</b> <b>fashion</b> <b>department</b> <b>home</b>	<b>news</b> sports <b>com</b> football games game soccer <b>culture</b> world tennis	computer edu software web <b>school</b> research programming <b>university</b> <b>information</b> systems	business <b>news</b> <b>information</b> market services <b>com</b> trade stock <b>home</b> gov	research science edu journal <b>art</b> resources culture <b>information</b> directory library



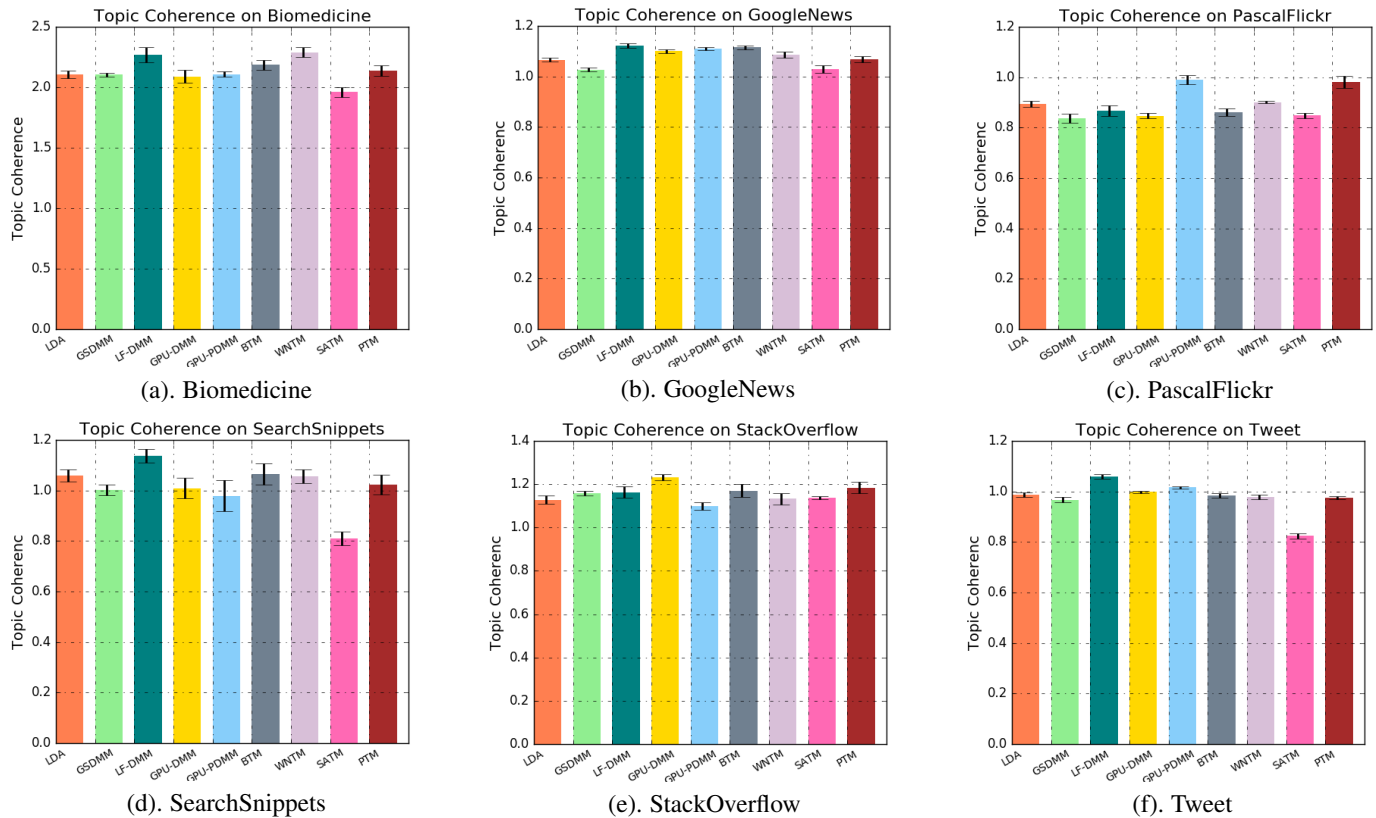


Fig. 8. Topic Coherence of all models on six datasets.

## ACKNOWLEDGEMENT

This research is partially supported by the National Natural Science Foundation of China under grant 61703362, and the Natural Science Foundation of Jiangsu Province of China under grant BK20170513.

## REFERENCES

- [1] T. Lin, W. Tian, Q. Mei, C. Hong, The dual-sparse topic model: mining focused topics and focused terms in short text, in: International Conference on World Wide Web, 2014.
- [2] J. Qiang, P. Chen, W. Ding, T. Wang, F. Xie, X. Wu, Topic discovery from heterogeneous texts, in: Tools with Artificial Intelligence (ICTAI), 2016 IEEE 28th International Conference on, IEEE, 2016, pp. 196–203.
- [3] T. Shi, K. Kang, J. Choo, C. K. Reddy, Short-text topic modeling via non-negative matrix factorization enriched with local word-context correlations, in: Proceedings of the 2018 World Wide Web Conference on World Wide Web, International World Wide Web Conferences Steering Committee, 2018, pp. 1105–1114.
- [4] X. Wang, C. Zhai, X. Hu, R. Sproat, Mining correlated bursty topic patterns from coordinated text streams, in: Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining, 2007, pp. 784–793.
- [5] B. Sriram, D. Fuhry, E. Demir, H. Ferhatosmanoglu, M. Demirbas, Short text classification in twitter to improve information filtering, in: Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval, 2010, pp. 841–842.
- [6] Z. Ma, A. Sun, Q. Yuan, G. Cong, Topic-driven reader comments summarization, in: Proceedings of the 21st ACM international conference on Information and knowledge management, 2012, pp. 265–274.
- [7] J. Weng, E.-P. Lim, J. Jiang, Q. He, Twitterrank: finding topic-sensitive influential twitterers, in: Proceedings of the third ACM international conference on Web search and data mining, 2010, pp. 261–270.
- [8] T. Hofmann, Probabilistic latent semantic indexing, in: Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval, ACM, 1999, pp. 50–57.
- [9] D. M. Blei, A. Y. Ng, M. I. Jordan, Latent dirichlet allocation, the Journal of machine Learning research 3 (2003) 993–1022.
- [10] D. M. Blei, Probabilistic topic models, Communications of the ACM 55 (4) (2012) 77–84.
- [11] M. D. Hoffman, D. M. Blei, F. Bach, Online learning for latent dirichlet allocation, in: International Conference on Neural Information Processing Systems, 2010.
- [12] P. Xie, E. P. Xing, Integrating document clustering and topic modeling, UAI.
- [13] P. Xie, D. Yang, E. P. Xing, Incorporating word correlation knowledge into topic modeling, in: NACACL, 2015.
- [14] X. Cheng, X. Yan, Y. Lan, J. Guo, Btm: Topic modeling over short texts, Knowledge and Data Engineering, IEEE Transactions on 26 (12) (2014) 2928–2941.
- [15] O. Jin, N. N. Liu, K. Zhao, Y. Yu, Q. Yang, Transferring topical knowledge from auxiliary long texts for short text clustering, in: Proceedings of the 20th ACM international conference on Information and knowledge management, ACM, 2011, pp. 775–784.
- [16] X.-H. Phan, L.-M. Nguyen, S. Horiguchi, Learning to classify short and sparse text & web with hidden topics from large-scale data collections, in: Proceedings of the 17th international conference on World Wide Web, ACM, 2008, pp. 91–100.
- [17] R. Mehrotra, S. Sanner, W. Buntine, L. Xie, Improving lda topic models for microblogs via tweet pooling and automatic labeling, in: Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval, ACM, 2013, pp. 889–892.
- [18] J. Yin, J. Wang, A dirichlet multinomial mixture model-based approach for short text clustering, in: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, 2014, pp. 233–242.
- [19] W. X. Zhao, J. Jiang, J. Weng, J. He, E.-P. Lim, H. Yan, X. Li, Comparing twitter and traditional media using topic models, in: Advances in Information Retrieval, 2011, pp. 338–349.
- [20] X. Quan, C. Kit, Y. Ge, S. J. Pan, Short and sparse text topic modeling via self-aggregation, in: Proceedings of the 24th International Conference on Artificial Intelligence, 2015, pp. 2270–2276.
- [21] X. Yan, J. Guo, Y. Lan, J. Xu, X. Cheng, A probabilistic model for bursty topic discovery in microblogs., in: AAAI, 2015, pp. 353–359.

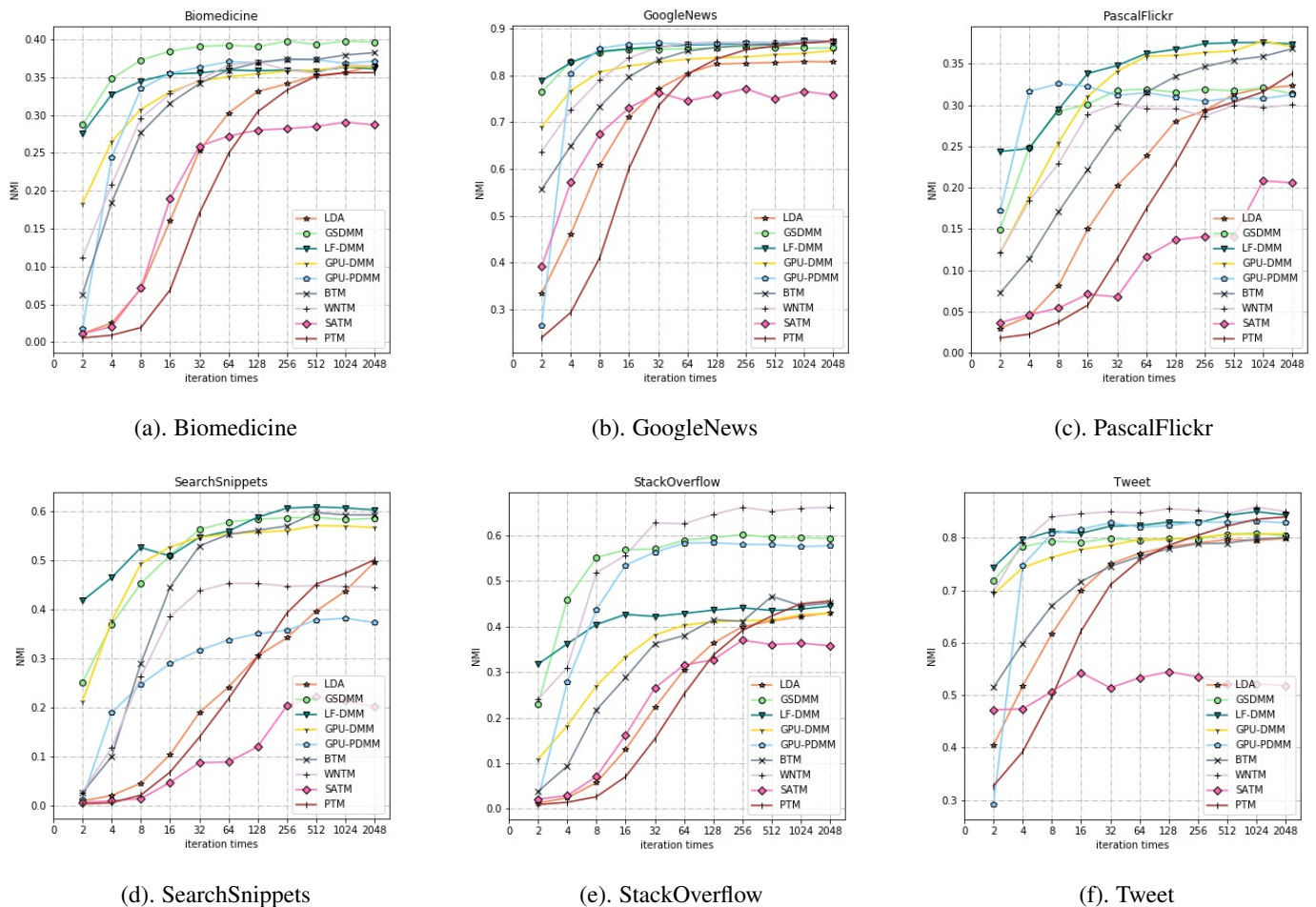


Fig. 9. NMI values with different number of iterations on every corpora.

- [22] K. Nigam, A. K. McCallum, S. Thrun, T. Mitchell, Text classification from labeled and unlabeled documents using em, *Machine learning* 39 (2-3) (2000) 103–134.
- [23] G. Yu, R. Huang, Z. Wang, Document clustering via dirichlet process mixture model with feature selection, in: *SIGKDD, ACM*, 2010, pp. 763–772.
- [24] R. Huang, G. Yu, Z. Wang, J. Zhang, L. Shi, Dirichlet process mixture model for document clustering with feature partition, *Knowledge and Data Engineering, IEEE Transactions on* 25 (8) (2013) 1748–1759.
- [25] J. Qiang, Y. Li, Y. Yuan, X. Wu, Short text clustering based on pitman-yor process mixture model, *Applied Intelligence* 48 (7) (2018) 1802–1812.
- [26] J. Yin, J. Wang, A text clustering algorithm using an online clustering scheme for initialization, in: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM*, 2016, pp. 1995–2004.
- [27] J. Yin, D. Chao, Z. Liu, W. Zhang, X. Yu, J. Wang, Model-based clustering of short text streams, in: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM*, 2018, pp. 2634–2642.
- [28] D. Q. Nguyen, R. Billingsley, L. Du, M. Johnson, Improving topic models with latent feature word representations, *Transactions of the Association for Computational Linguistics* 3 (2015) 299–313.
- [29] C. Li, H. Wang, Z. Zhang, A. Sun, Z. Ma, Topic modeling for short texts with auxiliary word embeddings, in: *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval, ACM*, 2016, pp. 165–174.
- [30] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: *Advances in neural information processing systems*, 2013, pp. 3111–3119.
- [31] C. Li, Y. Duan, H. Wang, Z. Zhang, A. Sun, Z. Ma, Enhancing topic modeling for short texts with auxiliary word embeddings, *ACM Transactions on Information Systems (TOIS)* 36 (2) (2017) 11.
- [32] Y. Zuo, J. Zhao, K. Xu, Word network topic model: a simple but general solution for short and imbalanced texts, *Knowledge and Information Systems* 48 (2) (2016) 379–398.
- [33] W. Chen, J. Wang, Y. Zhang, H. Yan, X. Li, User based aggregation for biern topic model, in: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, Vol. 2, 2015, pp. 489–494.
- [34] W. Wang, H. Zhou, K. He, J. E. Hopcroft, Learning latent topics from the word co-occurrence network, *National Conference of Theoretical Computer Science* (2017) 18–30.
- [35] Y. Zuo, J. Wu, H. Zhang, H. Lin, F. Wang, K. Xu, H. Xiong, Topic modeling of short texts: A pseudo-document view, in: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining, ACM*, 2016, pp. 2105–2114.
- [36] J. Qiang, P. Chen, T. Wang, X. Wu, Topic modeling over short texts by incorporating word embeddings, in: *Pacific-Asia Conference on Knowledge Discovery and Data Mining, Springer*, 2017, pp. 363–374.
- [37] P. V. Bicalho, M. Pita, G. Pedrosa, A. Lacerda, G. L. Pappa, A general framework to expand short text for topic modeling, *Information Sciences* 393 (2017) 66–81.
- [38] L. Hong, B. D. Davison, Empirical study of topic modeling in twitter, in: *Proceedings of the first workshop on social media analytics, ACM*, 2010, pp. 80–88.
- [39] A. K. McCallum, Mallet: A machine learning for language toolkit.
- [40] D. C. Liu, J. Nocedal, On the limited memory bfgs method for large scale optimization, *Mathematical programming* 45 (1-3) (1989) 503–528.
- [41] H. Mahmoud, *Pólya urn models*, Chapman and Hall/CRC, 2008.
- [42] M. Rosen-Zvi, T. L. Griffiths, M. Steyvers, P. Smyth, The author-topic model for authors and documents, *uncertainty in artificial intelligence* (2004) 487–494.

- [43] D. Ramage, S. Dumais, D. Liebling, Characterizing microblogs with topic models, in: Fourth International AAAI Conference on Weblogs and Social Media, 2010.
- [44] I. Guy, Social recommender systems, in: Recommender systems handbook, Springer, 2015, pp. 511–543.
- [45] X. Qian, H. Feng, G. Zhao, T. Mei, Personalized recommendation combining user interest and social circle, IEEE transactions on knowledge and data engineering 26 (7) (2014) 1763–1777.
- [46] O. Phelan, K. McCarthy, B. Smyth, Using twitter to recommend real-time topical news, in: Proceedings of the third ACM conference on Recommender systems, ACM, 2009, pp. 385–388.
- [47] J. Chen, R. Nairn, L. Nelson, M. Bernstein, E. Chi, Short and tweet: experiments on recommending content from information streams, in: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, ACM, 2010, pp. 1185–1194.
- [48] M. Chen, X. Jin, D. Shen, Short text classification improved by learning multi-granularity topics, in: IJCAI, 2011, pp. 1776–1781.
- [49] D.-T. Vo, C.-Y. Ock, Learning to classify short text from scientific documents using topic models with various types of knowledge, Expert Systems with Applications 42 (3) (2015) 1684–1698.
- [50] Z. Dai, A. Sun, X.-Y. Liu, Crest: Cluster-based representation enrichment for short text classification, in: Pacific-Asia Conference on Knowledge Discovery and Data Mining, Springer, 2013, pp. 256–267.
- [51] A. H. Razavi, D. Inkpen, Text representation using multi-level latent dirichlet allocation, in: Canadian Conference on Artificial Intelligence, Springer, 2014, pp. 215–226.
- [52] C. X. Lin, B. Zhao, Q. Mei, J. Han, Pet: a statistical model for popular events tracking in social communities, in: Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2010, pp. 929–938.
- [53] C. C. Aggarwal, K. Subbian, Event detection in social streams, in: Proceedings of the 2012 SIAM international conference on data mining, SIAM, 2012, pp. 624–635.
- [54] A. Ritter, O. Etzioni, S. Clark, et al., Open domain event extraction from twitter, in: Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2012, pp. 1104–1112.
- [55] J. H. Lau, N. Collier, T. Baldwin, On-line trend analysis with topic models: \# twitter trends detection topic model online, Proceedings of COLING 2012 (2012) 1519–1534.
- [56] J. Pennington, R. Socher, C. Manning, Glove: Global vectors for word representation, in: Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), 2014, pp. 1532–1543.
- [57] C. Rashtchian, P. Young, M. Hodosh, J. Hockenmaier, Collecting image annotations using amazon’s mechanical turk, in: Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk, Association for Computational Linguistics, 2010, pp. 139–147.
- [58] C. Finegan-Dollak, R. Coke, R. Zhang, X. Ye, D. Radev, Effects of creativity and cluster tightness on short text clustering performance, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Vol. 1, 2016, pp. 654–665.
- [59] D. Newman, J. H. Lau, K. Grieser, T. Baldwin, Automatic evaluation of topic coherence, in: Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, 2010, pp. 100–108.
- [60] D. Mimno, H. M. Wallach, E. Talley, M. Leenders, A. McCallum, Optimizing semantic coherence in topic models, in: Proceedings of the Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, 2011, pp. 262–272.
- [61] J. Chang, S. Gerrish, C. Wang, J. L. Boyd-Graber, D. M. Blei, Reading tea leaves: How humans interpret topic models, in: Advances in neural information processing systems, 2009, pp. 288–296.
- [62] J. Chuang, C. D. Manning, J. Heer, Termite: Visualization techniques for assessing textual topic models, in: Proceedings of the international working conference on advanced visual interfaces, ACM, 2012, pp. 74–77.
- [63] C. Sievert, K. Shirley, Ldavis: A method for visualizing and interpreting topics, in: Proceedings of the workshop on interactive language learning, visualization, and interfaces, 2014, pp. 63–70.