

Extracting Topics with Simultaneous Word Co-occurrence and Semantic Correlation Graphs: Neural Topic Modeling for Short Texts

Yiming Wang^{1,2}, Ximing Li^{1,2} ^{*}†, Xiaotang Zhou³, Jihong Ouyang^{1,2*}

¹College of Computer Science and Technology, Jilin University, China

²Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, China

³School of Computer Science and Engineering, Changchun University of Technology, China

{yimingw17, liximing86}@gmail.com

zhouxiaotang@ccut.edu.cn, ouyj@jlu.edu.cn

Abstract

Short text nowadays has become a more fashionable form of text data, *e.g.*, Twitter posts, news titles, and product reviews. Extracting semantic topics from short texts plays a significant role in a wide spectrum of NLP applications, and neural topic modeling is now a major tool to achieve it. Motivated by learning more coherent and semantic topics, in this paper we develop a novel neural topic model named Dual Word Graph Topic Model (DWGTM), which extracts topics from simultaneous word co-occurrence and semantic correlation graphs. To be specific, we learn word features from the global word co-occurrence graph, so as to ingest rich word co-occurrence information; we then generate text features with word features, and feed them into an encoder network to get topic proportions per-text; finally, we reconstruct texts and word co-occurrence graph with topical distributions and word features, respectively. Besides, to capture semantics of words, we also apply word features to reconstruct a word semantic correlation graph computed by pre-trained word embeddings. Upon those ideas, we formulate DWGTM in an auto-encoding paradigm and efficiently train it with the spirit of neural variational inference. Empirical results validate that DWGTM can generate more semantically coherent topics than baseline topic models.

1 Introduction

The topic modeling family targets at learning latent topic representations from text document collections (Blei, 2012). During the past decades, it has been extensively applied in many tasks of natural language processing, *e.g.*, sentiment analysis (Lin and He, 2009), summarization (Ma et al., 2012) and classification (Zeng et al., 2018), to name just a few. Conventional topic models such as Latent

Dirichlet Allocation (LDA) (Blei et al., 2003) are often inferred by approximate inference methods, *e.g.*, mean-field variational inference (Jordan et al., 1999) and collapsed Gibbs sampling (Griffiths and Steyvers, 2004), which require model-specific derivations. The recent inference method with neural networks, such as Variational Auto-Encoder (VAE) (Kingma and Welling, 2014), works in a black-box manner, providing a more generic and flexible solution to topic models beyond traditional approximate inference methods. Broadly speaking, the models inferred with neural networks are referred to as **neural topic models**, and they have been recently drawn much more attention from the natural language processing community (Zhu et al., 2018; Burkhardt and Kramer, 2019; Dieng et al., 2020; Wu et al., 2020).

Unfortunately, whether for conventional or neural topic models, they tend to perform poorly on **short text**, a more fashionable and significant form of text data, *e.g.*, Twitter posts, news titles, and product reviews. The main reason is that short texts lack document-level word co-occurrences, known as the sparsity problem, which hinders models to capture coherent word patterns. Many conventional topic models have been developed to handle short texts. For example, given very few words per-text, Dirichlet Multinomial Mixture (DMM) (Nigam et al., 2000; Yin and Wang, 2014) constrains that each text covers a signal topic. Biterm Topic Model (BTM) (Yan et al., 2013; Cheng et al., 2014) directly learns topics from corpus-level word co-occurrence patterns. Recently, there are also few attempts of neural topic models aiming to address the sparsity problem of short texts. GraphBTM (Zhu et al., 2018) extracts topics from word graphs of randomly drawn mini-corpus. Negative sampling and Quantization Topic Model (NQTM) (Wu et al., 2020) applies a topic distribution quantization method to pursue peakier topic proportions of texts. As reported in (Wu et al., 2020), those neu-

* Corresponding Author

† Contributing equally with the first author.

ral topic models can empirically induce coherent topics from short texts.

Motivated by learning more coherent and semantic topics, in this paper we develop a novel neural topic model for short texts, namely **Dual Word Graph Topic Model (DWGTM)**. As the name suggests, in DWGTM we apply two word graphs, including the **word co-occurrence graph** constructed by aggregating word co-occurrence patterns of each text to alleviate the sparsity problem, and the **word semantic correlation graph** generated by using the pre-trained word embeddings to capture the semantic information of words. Specifically, we formulate DWGTM in an auto-encoding paradigm with four main components: (1) We encode the word co-occurrence graph as word features by applying a Graph Convolutional Network (GCN) module (Kipf and Welling, 2016a). (2) For each text, we construct its feature with corresponding word features, and encode it as topic proportion. (3) Reconstruct texts with topical distributions. (4) Reconstruct the two word graphs with word features. With the word semantic correlation graph, DWGTM can output topics that are associated with the semantic information of words. Besides, we propose a novel topic quality metric to measure the semantic coherence of learned topics, namely **Topical Semantics Coherence (TSC)**. We conduct extensive experiments to evaluate DWGTM, and empirical results indicate that DWGTM can learn more semantically coherent topics than existing baseline models.

In a nutshell, the major contributions of this paper are listed below:

- We propose a novel neural topic model **DWGTM** for short texts, extracting topics from simultaneous word co-occurrence and semantic correlation graphs.
- We propose a novel topic quality metric called **TSC**, which measures the semantic coherence of learned topics.
- On three benchmark datasets of short texts, **DWGTM** empirically outputs more semantically coherent topics than strong baseline models.

2 Related Work

In this section, we briefly review related works on conventional topic modeling of short texts and neural topic modeling.

2.1 Topic Modeling for Short Texts

Short texts lack the document-level word co-occurrence information, making conventional topic models such as LDA (Blei et al., 2003) much less effective. To resolve this issue of short text, existing models mainly adopt the methodology of word co-occurrence enrichment (Yan et al., 2013; Yin and Wang, 2014; Quan et al., 2015; Zuo et al., 2016a,b; Li et al., 2016, 2018; Shi et al., 2018; Li et al., 2019a,b, 2020a). First, one straightforward way is to generate long pseudo-texts by adaptively aggregating short texts and then learn topics from them by applying LDA. Several representatives (Quan et al., 2015; Zuo et al., 2016a; Li et al., 2018) jointly estimate long pseudo-texts and topics, however, they are often time consuming as well as sensitive to the number of long pseudo-texts. Second, another mainstream is to extract more word co-occurrences at the corpus level. The BTM (Yan et al., 2013; Cheng et al., 2014) directly induces topics from all word co-occurrence patterns of the corpus. Semantics-assisted Non-negative Matrix Factorization (SeaNMF) (Shi et al., 2018) regards each word type as a pseudo-text consisting of the words that co-occur with it in the same short text, and learns topics with those auxiliary word type pseudo-texts. Additionally, other attempts (Li et al., 2016, 2019a) upgrade existing models, *e.g.*, DMM and BTM, by further leveraging auxiliary knowledge or techniques such as word semantic correlations measured by pre-trained word embeddings (Mikolov et al., 2013; Pennington et al., 2014). In contrast to aforementioned models, our DWGTM is built on the framework of neural variational inference with GCN (Kipf and Welling, 2016a), enabling to effectively extract topics with word co-occurrence patterns.

2.2 Neural Topic Modeling

Along the new research line of integrating VAE (Kingma and Welling, 2014), a number of neural topic models have been proposed. Generally, the basic idea of neural topic modeling is to apply neural networks as topic encoders to induce topic representations of texts, and reconstruct texts with topical distributions. Benefiting from the effectiveness and flexibility of neural networks in unsupervised representation learning, neural topic models can induce more significant topics from texts. Nowadays, the representatives include Neural Variational Document Model (NVDM) (Miao et al., 2016), Product

of expert LDA (ProdLDA) (Srivastava and Sutton, 2017), and Embedded Topic Model (ETM) (Ding et al., 2020), *etc.* Besides these “naive” neural variants of LDA, many other models have been investigated by applying (1) various neural modules to the topic encoder, *e.g.*, recurrent module (Rezaee and Ferraro, 2020), attention mechanism (Li et al., 2020b), and graphical connection (Zhu et al., 2018; Yang et al., 2020), and (2) new learning paradigms, *e.g.*, adversarial training (Wang et al., 2019), reinforcement learning (Gui et al., 2019), and lifelong learning (Gupta et al., 2020). However, despite their effectiveness on normal long texts, those models suffer from the sparsity problem of short texts (Zeng et al., 2018).

To our knowledge, there are only a few neural topic models for addressing the sparsity problem of short texts (Zeng et al., 2018; Zhu et al., 2018; Wu et al., 2020). Inspired by BTM (Yan et al., 2013), the GraphBTM method (Zhu et al., 2018) directly learns topics from the aggregated word co-occurrence patterns of randomly generated mini-corpus. The NQTM method (Wu et al., 2020) is based on the assumption that the peakier topic proportions of texts are more appropriate for modeling short texts as demonstrated in DMM (Yin and Wang, 2014). To achieve this, it applies a topic distribution quantization method, and meanwhile it adopts a negative sampling step to avoid repetitive topics. Orthogonal to those models, our DWGTM further employs the pre-trained word embeddings to capture the semantic information of words, so as to output more semantically coherent topics.

3 The Proposed DWGTM Model

In this section, we introduce the proposed **Dual Word Graph Topic Model (DWGTM)**. For convenience, the important notations used in this paper are summarized in Table 1.

3.1 Overview of DWGTM

The topic modeling family such as LDA (Blei et al., 2003) refers to the probabilistic model that describes the generative process of documents. Basically, it posits totally k topics $\phi_{1:k}$, each of which is a multinomial distribution over the vocabulary, and each document is represented by a topic proportion θ . Given a corpus \mathcal{D} consisting of n documents $\mathbf{x}_{1:n}$, the main goal of topic modeling is to estimate topics $\phi_{1:k}$ and topic proportions $\theta_{1:n}$ from \mathcal{D} . However, it is commonly intractable to

Table 1: Notation summary.

Notation	Description
n	number of texts
v	vocabulary size
k	number of topics
\mathbf{x}	word frequency vector of text
\mathcal{G}^c	word co-occurrence graph
\mathcal{G}^s	word semantic correlation graph
\mathbf{z}^w	word feature
\mathbf{z}^t	latent text feature
ϕ	topic distribution
θ	topic proportion of text
\mathbf{W}^c	learnable parameter of WCG-Encoder
\mathbf{W}^t	learnable parameter of TP-Encoder

accurately estimate $\{\phi, \theta\}$ for short texts, due to the lack of document-level word co-occurrences, known as the sparsity problem.

To effectively handle short texts, we propose a novel neural topic model named **DWGTM** by not only leveraging the corpus-level word co-occurrence information to address the sparsity problem (Yan et al., 2013), but also capturing word semantic correlations measured by pre-trained word embeddings (Mikolov et al., 2013; Pennington et al., 2014). Specifically, as depicted in Fig.1, DWGTM consists of the four components in the auto-encoding manner. (1) **WCG-Encoder**: We construct a global word co-occurrence graph \mathcal{G}^c , and then encode \mathcal{G}^c as word features $\mathbf{z}_{1:v}^w$, where v denotes the vocabulary size. (2) **TP-Encoder**: We construct latent text features $\mathbf{z}_{1:n}^t$ by using $\mathbf{z}_{1:v}^w$, and then encode $\mathbf{z}_{1:n}^t$ as topic proportions $\theta_{1:n}$. (3) **Text-Decoder**: We reconstruct the texts $\mathbf{x}_{1:n}$ with $\theta_{1:n}$ and topics $\phi_{1:k}$. (4) **DualWG-Decoder**: We reconstruct \mathcal{G}^c with $\mathbf{z}_{1:v}^w$. Meanwhile, to further capture semantic information of words, we construct a word semantic correlation graph \mathcal{G}^s by using pre-trained word embeddings, and reconstruct \mathcal{G}^s with also $\mathbf{z}_{1:v}^w$. In the following part, we introduce each component of DWGTM in more details.

3.2 WCG-Encoder

Given a corpus \mathcal{D} , we first construct a word co-occurrence graph $\mathcal{G}^c = (\mathcal{V}, \mathcal{E}^c)$, where \mathcal{V} and \mathcal{E}^c denote the sets of word nodes and word co-occurrence edges, respectively. That is, the graph can be represented by the co-occurrence adjacency matrix $\mathbf{A}^c \in \mathbb{R}^{v \times v}$, where each element \mathbf{A}_{ij}^c denotes the count of words w_i and w_j co-occurring in the same text. The WCG-encoder targets at encoding \mathbf{A}^c as word features $\mathbf{Z}^w = [\mathbf{z}_1^w, \dots, \mathbf{z}_v^w]^\top$, so that more

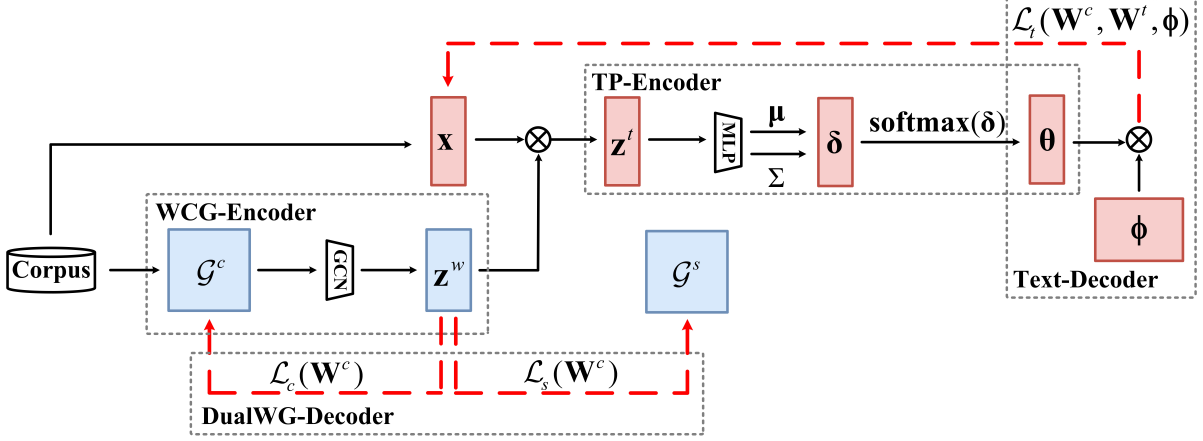


Figure 1: Overall model structure of DWGTM with four components, *i.e.*, **WCG-Encoder**, **TP-Encoder**, **Text-Decoder**, and **DualWG-Decoder**.

frequently co-occurring words tend to have more similar word features. This is achieved by applying a GCN module parameterized by \mathbf{W}^c :

$$\mathbf{Z}^w = f_{\text{GCN}}(\mathbf{A}^c; \mathbf{W}^c) \quad (1)$$

Following (Kipf and Welling, 2016a), each layer of the GCN module is formulated below:

$$\mathbf{z}_{(l)}^w = \psi(\tilde{\mathbf{A}}^c \mathbf{z}_{(l-1)}^w \mathbf{W}_{(l)}^c), \quad l = 1, \dots, l^c, \quad (2)$$

where l^c is the number of layers; $\mathbf{W}^c = \{\mathbf{W}_{(l)}^c\}_{l=1}^{l^c}$ are the learnable parameters; $\mathbf{z}_{(0)}^w$ is initialized by the identity matrix \mathbf{I}_v with the shape of v ; $\psi(\cdot)$ denotes the Tanh activation function; $\tilde{\mathbf{A}}^c = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A}^c + \mathbf{I}_v)\mathbf{D}^{-\frac{1}{2}}$ is the symmetrically normalized adjacency matrix; and \mathbf{D} denotes the degree matrix of $\mathbf{A}^c + \mathbf{I}_v$.

3.3 TP-Encoder

Naturally, the resulting word features \mathbf{Z}^w learned from \mathcal{G}^c are rich in global word co-occurrence information. Accordingly, we can use \mathbf{Z}^w to generate latent text features $\mathbf{z}_{1:n}^t$, enabling to alleviate the sparsity problem of short texts. For each short text, the latent text feature can be easily obtained by aggregating its corresponding word features, formulated below:

$$\mathbf{z}_d^t = (\mathbf{Z}^w)^\top \frac{\mathbf{x}_d}{|\mathbf{x}_d|}, \quad d = 1, \dots, n, \quad (3)$$

where \mathbf{x}_d and $|\mathbf{x}_d|$ denote the word frequency vector of the d^{th} document and its total number of word tokens, respectively.

The TP-Encoder aims at encoding $\mathbf{z}_{1:n}^t$ as topic proportions $\theta_{1:n}$. Inspired by (Miao et al., 2016; Dieng et al., 2020), we apply the VAE-like paradigm

with logistic-normal prior distribution. Specifically, suppose that for each short text the topic proportion is drawn from a logistic-normal prior as follows:

$$\delta_d \sim \mathcal{N}(\mu_0, \Sigma_0); \quad \theta_d = \text{softmax}(\delta_d), \quad d = 1, \dots, n, \quad (4)$$

where δ can be regarded as the unnormalized topic proportion; and $\mathcal{N}(\mu_0, \Sigma_0)$ denotes a Gaussian prior probability. We apply a fully-connected module, *a.k.a.*, variational inference network (Dieng et al., 2020), which ingests each latent text feature \mathbf{z}_d^t and outputs the mean μ_d and covariance Σ_d of the unnormalized topic proportion δ_d , formulated below:

$$\mathbf{H}_{(l)} = \rho(\mathbf{W}_{(l)}^t \mathbf{H}_{(l-1)}), \quad l = 1, \dots, l^t. \quad (5)$$

$$\mu_d = \mathbf{W}_\mu^t \cdot \mathbf{H}_{(l^t)} \quad (6)$$

$$\Sigma_d = \mathbf{W}_\Sigma^t \cdot \mathbf{H}_{(l^t)}. \quad (7)$$

where l^t denotes the number of layers; $\mathbf{W}^t = \{\{\mathbf{W}_{(l)}^t\}_{l=1}^{l^t}, \mathbf{W}_\mu^t, \mathbf{W}_\Sigma^t\}$ are the learnable parameters; $\mathbf{H}_{(0)}$ is initialized by \mathbf{z}_d^t ; and $\rho(\cdot)$ denotes the Tanh activation function. We then compute the topic proportion θ_d by leveraging the reparameterization trick (Kingma and Welling, 2014):

$$\theta_d = \text{softmax}(\mu_d + \Sigma_d \odot \epsilon), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k) \quad (8)$$

where \odot denotes element-wise product; and ϵ is a sample drawn from the Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I}_k)$. Due to the space limitation, we omit background descriptions of this VAE-like paradigm and reparameterization, and refer the readers to more details in (Kingma and Welling, 2014; Mnih and Gregor, 2014; Rezende et al., 2014; Miao et al., 2016).

Remark. Strictly speaking, the reparameterization trick (*i.e.*, Eq.8) is really meant for forming the Monte Carlo approximation of the variational objective (Kingma and Welling, 2014), and it should even be described in the decoding process. We kindly emphasize that we introduce Eq.8 as a step of encoding θ for the sake of a more intuitive expression for the TP-Encoder.

3.4 Text-Decoder

In this component, we reconstruct the texts $\mathbf{x}_{1:n}$ with topic proportions $\theta_{1:n}$ and topics $\phi_{1:k}$. Following the spirit of VAE derivation (Kingma and Welling, 2014), the reconstruction loss of $\mathbf{x}_{1:n}$ consists of a log marginal likelihood term and a KL-divergence regularizer as follows:

$$\mathcal{L}_t(\mathbf{W}^c, \mathbf{W}^t, \phi) = -\log p(\mathbf{x}) + \mathcal{R}_{\text{KL}}, \quad (9)$$

We adhere to the generative assumption of LDA-like models, therefore the marginal likelihood term of texts can be formulated below:

$$p(\mathbf{x}) = \prod_{d=1}^n \prod_{i \in \mathbf{x}_d} \sum_{t=1}^k \theta_{dt} \phi_{ti}, \quad (10)$$

where $\theta_{i:n}$ are computed by Eq.8. Second, the KL-divergence regularizer admits a closed-form expression as follows:

$$\mathcal{R}_{\text{KL}} = -\frac{1}{2} \sum_{d=1}^n (1 + \log |\Sigma_d| - \boldsymbol{\mu}_d^\top \boldsymbol{\mu}_d - \text{Tr}(\Sigma_d)), \quad (11)$$

where $\text{Tr}(\cdot)$ denotes the trace of a matrix.

3.5 DualWG-Decoder

As its name suggests, the aim of DualWG-Decoder is two-fold: applying the word features $\mathbf{z}_{1:v}^w$ to reconstruct the word co-occurrence graph \mathcal{G}^c and also an auxiliary word semantic correlation graph \mathcal{G}^s .

Reconstruction of \mathcal{G}^c . Following (Kipf and Welling, 2016b), we apply an inner product decoder with word features. Accordingly, the reconstruction loss is formulated as follows:

$$\mathcal{L}_c(\mathbf{W}^c) = - \sum_{\{i,j\} \in \mathcal{E}^c} \mathbf{A}_{ij}^c \log \sigma \left((\mathbf{z}_i^w)^\top \mathbf{z}_j^w \right), \quad (12)$$

where $\sigma(\cdot)$ denotes the Sigmoid function.

Reconstruction of \mathcal{G}^s . Besides extracting topics by applying the word co-occurrence statistics, We expect to take the semantic information of words into consideration, so as to generate more semantically coherent topics (Li et al., 2016, 2019a). To achieve this, we construct a word semantic correlation graph $\mathcal{G}^s = (\mathcal{V}, \mathcal{E}^s)$, where \mathcal{E}^s denotes the set of word semantic correlation edges. Let $\mathbf{A}^s \in \mathbb{R}^{v \times v}$ be the corresponding adjacency matrix, where each element \mathbf{A}_{ij}^s reflects the cosine similarity between pre-trained GloVe embeddings¹ of words w_i and w_j . To be specific, it is formulated as follows:

$$\mathbf{A}_{ij}^s = \begin{cases} \gamma_{ij}, & \text{if } \gamma_{ij} > \gamma^* \\ 0, & \text{otherwise} \end{cases}, \quad (13)$$

where $\gamma_{ij} = \cos(\mathbf{g}_i, \mathbf{g}_j)$ denotes the cosine similarity; the notation \mathbf{g} specifies the pre-trained GloVe word embedding; and γ^* is a word semantic correlation threshold.

We reconstruct \mathcal{G}^s by encouraging the resulting word features to capture word semantic correlations. Accordingly, the reconstruction loss of \mathcal{G}^s can be formulated below:

$$\mathcal{L}_s(\mathbf{W}^c) = \sum_{\{i,j\} \in \mathcal{E}^s} \|\cos(\mathbf{z}_i^w, \mathbf{z}_j^w) - \gamma_{ij}\|_2^2, \quad (14)$$

where $\|\cdot\|_2$ denotes the ℓ_2 norm.

3.6 Full Objective of DWGTM

We now outline the full objective of DWGTM. Except the reconstruction losses of $\mathbf{x}_{1:n}$, \mathcal{G}^c , and \mathcal{G}^s , we also incorporate the following entropy regularization term to encourage peakier topic proportions:

$$\mathcal{R}_E = - \sum_{d=1}^n \sum_{t=1}^k \theta_{dt} \log \theta_{dt} \quad (15)$$

Finally, we can reach the full objective with respect to the learnable parameters $\{\mathbf{W}^c, \mathbf{W}^t, \phi\}$ as follows:

$$\mathcal{L}(\mathbf{W}^c, \mathbf{W}^t, \phi) = \mathcal{L}_t(\mathbf{W}^c, \mathbf{W}^t, \phi) + \lambda_1 \mathcal{L}_c(\mathbf{W}^c) + \lambda_2 \mathcal{L}_s(\mathbf{W}^c) + \lambda_3 \mathcal{R}_E, \quad (16)$$

where λ_1 , λ_2 , and λ_3 are the scale parameters.

4 Experiment

Datasets. In the experiments, we select three benchmark short text datasets: *Trec*,² *Google*

¹ <https://nlp.stanford.edu/projects/glove/>

² <http://cogcomp.cs.illinois.edu/Data/QA/QC>

Table 2: Statistics of short text datasets after preprocessing. n : number of short texts. v : vocabulary size. \hat{n} : average document length. l : number of categories.

Dataset	n	v	\hat{n}	l
Trec	4,198	989	3.2	6
GoogleNews	9,284	659	3.7	152
YahooAnswer	22,937	1297	3.8	10

News,³ and *YahooAnswer*.⁴ For all datasets, we remove digits and words with term frequencies less than 20. Stop words and non-english words are filtered out by NLTK.⁵ For clarity, the statistics of those datasets are listed in Table 2.

Baseline Topic Models. We select 8 existing baselines, including 4 conventional topic models and 4 neural topic models. Following their original papers, the important implementation details of all baselines are described below.

- **LDA**⁶ (Blei et al., 2003): The model is inferred by variational inference, and the Dirichlet priors for topic proportions and topic distributions are set to 0.1 and 0.01, respectively.
- **DMM**⁷ (Yin and Wang, 2014): The two Dirichlet priors are set as $50/k$ and 0.01, respectively.
- **BTM**⁸ (Yan et al., 2013): The two Dirichlet priors are set as 0.01 and 0.001, respectively.
- **Generalized Pólya Urn DMM (GPUDDMM)**⁹ (Li et al., 2016): The two Dirichlet priors are set as $50/k$ and 0.01, respectively; and the similarity threshold is set as 0.8.
- **NVDM**¹⁰ (Miao et al., 2016): The model applies a 2-layer MLP encoder with 500 hidden neurons.
- **ProdLDA**¹¹ (Srivastava and Sutton, 2017): The model applies a 3-layer MLP encoder with 100 hidden neurons.

³<https://news.google.com/>

⁴<https://answers.yahoo.com>

⁵<https://nltk.org>

⁶<https://github.com/blei-lab/lda-c>

⁷<https://github.com/jackyin12/GSDMM>

⁸<https://github.com/xiaohuiyan/BTM>

⁹<https://github.com/NobodyWHU/GPUDDMM>

¹⁰<https://github.com/ysmiao/nvdm>

¹¹https://github.com/akashgit/autoencoding_vi_for_topic_models

- **GraphBTM**¹² (Zhu et al., 2018): The model applies a 3-layer GCN encoder with 100 hidden neurons and samples 3 documents as a mini-corpus.
- **NQTM**¹³ (Wu et al., 2020): The model applies a 3-layer MLP encoder with 100 hidden neurons and the word sample size for negative sampling is set as 20.

For DWGTM, we apply a 2-layer GCN WCG-Encoder and a 3-layer MLP TP-Encoder, where the hidden neurons of both encoders are set as 100-300-400-300- k . To avoid posterior collapsing, we adopt 0.4 dropout, batch normalization, and a shallower 1-layer Text-Decoder. The threshold γ^* is set to 0.6 for *Trec*, and 0.8 for *GoogleNews* and *YahooAnswer*. Scale parameters are set as $\lambda_1 = 0.1, \lambda_2 = 0.1, \lambda_3 = 1$. The number of epochs is 900 and mini-batch size is 200. To construct \mathcal{G}^s , we employ the 300-dimensional GloVe¹⁴ embeddings trained on Wikipedia2014 and Gigaword5. For fair comparisons, the baselines requiring word embeddings use the same GloVe embeddings.

Evaluation Metrics. To evaluate the topic quality, we adopt two metrics: Topic Coherence (TC) and Topical Semantics Coherence (TSC).

First, TC is the most popular topic quality metric that measures the co-occurrence statistics between top- m words of topics. Here, we compute the TC scores with the public TC project of *Palmetto*,¹⁵ where, especially, the setting of C_V is applied. Second, we propose a novel metric named TSC to measure the semantic coherence of topics. Analogy to TC, we suppose that higher similarities between top- m words of topics imply better semantic coherence for topics. Accordingly, TSC can be defined as follows:

$$\text{TSC} = \frac{2}{km^2} \sum_{t=1}^k \sum_{(w_i, w_j) \in \Omega_t} \frac{\cos(e_{w_i}, e_{w_j}) + 1}{2}, \quad (17)$$

where Ω_t is the top- m words of the t^{th} topic; and e_{w_i} and e_{w_j} denote the pre-trained word embeddings of w_i and w_j , respectively.

¹²<https://github.com/valdersoul/GraphBTM>

¹³<https://github.com/BobXWu/NQTM>

¹⁴<https://nlp.stanford.edu/projects/glove/>

¹⁵<https://github.com/dice-group/Palmetto>

Table 3: Results (mean \pm std) of TC. The best scores are displayed in boldface.

Model	Trec		GoogleNews		YahooAnswer		AvgRank
	$k = 20$	$k = 50$	$k = 20$	$k = 50$	$k = 20$	$k = 50$	
LDA	.368 \pm .001	.361 \pm .002	.382 \pm .004	.381 \pm .002	.377 \pm .004	.377 \pm .003	6.50
DMM	.372 \pm .006	.379 \pm .006	.408 \pm .003	.393 \pm .004	.383 \pm .007	.373 \pm .004	4.00
BTM	.365 \pm .010	.370 \pm .001	.395 \pm .004	.385 \pm .002	.367 \pm .007	.367 \pm .008	7.00
GPUDMM	.377 \pm .005	.379 \pm .005	.403 \pm .007	.391 \pm .003	.374 \pm .005	.375 \pm .003	4.50
NVDM	.343 \pm .006	.349 \pm .006	.390 \pm .008	.395 \pm .006	.383 \pm .010	.377 \pm .008	5.67
ProdLDA	.372 \pm .022	.376 \pm .005	.369 \pm .010	.372 \pm .007	.396 \pm .014	.391 \pm .013	5.17
GraphBTM	.354 \pm .001	.340 \pm .001	.359 \pm .001	.383 \pm .001	.384 \pm .001	.371 \pm .001	7.50
NQTM	.397 \pm .006	.379 \pm .007	.416 \pm .004	.391 \pm .004	.404 \pm .015	.397 \pm .006	2.17
DWGTm	.402 \pm .006	.392 \pm .007	.419 \pm .008	.403 \pm .002	.406 \pm .001	.389 \pm .015	1.33

Table 4: Results (mean \pm std) of TSC. The best scores are displayed in boldface.

Model	Trec		GoogleNews		YahooAnswer		AvgRank
	$k = 20$	$k = 50$	$k = 20$	$k = 50$	$k = 20$	$k = 50$	
LDA	.396 \pm .003	.393 \pm .003	.386 \pm .004	.398 \pm .005	.485 \pm .003	.456 \pm .002	6.83
DMM	.442 \pm .002	.421 \pm .004	.406 \pm .004	.413 \pm .002	.509 \pm .004	.491 \pm .003	3.00
BTM	.418 \pm .004	.414 \pm .005	.407 \pm .006	.415 \pm .004	.522 \pm .008	.513 \pm .004	2.83
GPUDMM	.441 \pm .007	.420 \pm .003	.396 \pm .006	.415 \pm .003	.513 \pm .004	.499 \pm .003	3.00
NVDM	.400 \pm .005	.400 \pm .003	.387 \pm .006	.382 \pm .004	.460 \pm .011	.449 \pm .004	6.83
ProdLDA	.349 \pm .004	.348 \pm .003	.408 \pm .010	.423 \pm .005	.370 \pm .015	.374 \pm .017	6.17
GraphBTM	.348 \pm .001	.356 \pm .001	.363 \pm .001	.370 \pm .001	.399 \pm .001	.370 \pm .001	8.67
NQTM	.428 \pm .012	.417 \pm .003	.401 \pm .004	.411 \pm .002	.499 \pm .011	.479 \pm .006	4.50
DWGTm	.453 \pm .007	.418 \pm .005	.412 \pm .005	.406 \pm .004	.516 \pm .007	.479 \pm .040	2.83

Specially, we describe several details of metrics.

- (1) For both metrics, higher scores indicate better performance.
- (2) We fix m to 10 in all evaluations.
- (3) For fair comparisons, we employ the pre-trained word2vec¹⁶ embeddings to compute TSC, instead of the GloVe embeddings that have been used in some of comparing models.

4.1 Topic Quality Results

We independently run each comparing model 5 times, then report the average scores of TC and TSC in Tables 3 and 4. In terms of TC, it can be clearly seen that our DWGTm can achieve higher scores than baseline models in most cases. First, DWGTm outperforms the neural competitors GraphBTM and NQTM, which also focus on handling short texts. Second, the TC scores of DWGTm are higher than those conventional topic models in all settings, where the results demonstrate the GCN WCG-Encoder can better capture the word co-occurrence information from the corpora. In terms of TSC, our DWGTm gets competitive scores, and ranks the first averagely. Comparing with neural topic models, DWGTm can achieve higher scores in most cases, where more importantly it beats the most art NQTM. Surprisingly, conventional short text topic models, *e.g.*, DMM,

BTM, and GPUDMM, can achieve competitive TSC scores with DWGTm, and even perform better than NQTM. The possible reason is that those shallow models capture similar semantic information to word2vec, *i.e.*, the word embeddings used to compute TSC scores in the experiments. Specially, we kindly indicate that a potential problem of DWGTm is the TSC degradation with more topics compared to conventional topic models. We will further investigate this problem.

4.2 Topic Visualization

For qualitative evaluations, we show the top-10 words of two selected topics about *politics* and *credit* across *YahooAnswer*. As presented in Table 5, we can observe that DWGTm can effectively learn informative word patterns from corpora, where the top topical words are exactly associated with *politics* and *credit*, being consistent in the judgment of human-beings to some extent. In contrast to baseline models, the topics learned by DWGTm seem more coherent, where some of baselines often generate several less informative words, *e.g.*, for the topic of *credit*, {"best", "bad", "long"} in LDA, {"old", "weight", "stomach"} in NVDM, and {"salt", "water", "ice"} in NQTM for the topic of *credit*. Besides, we also observe that the top-10 words of GPUDMM and DWGTm contain semantically related words. This implies

¹⁶<https://wikipedia2vec.github.io/wikipedia2vec/pretrained/>

Table 5: Visualization of the top-10 words of two example topics learned by comparing models

Model	Top-10 word list
LDA	many bush people president ok war come world end times best way take really long get credit good place bad
DMM	bush think president war people iraq us world george like get need find school job know credit go money want
BTM	bush people president think us war states united america george find get need credit know good money free online help
GPUDMM	bush think president war people world iraq george americans american get credit find much money need business pay home company
NVDM	united states president bush school want high really know george rid get old back mean weight credit yr stomach loan
ProdLDA	president marriage movie iraq every form get need remember cancer uick money development wont wants longer sports care base treat
GraphBTM	immigration music bike republicans three operation step named income god illegal install stand female affect true turn accept choice pay
NQTM	president united states george bush god war clinton iraq nuclear credit card debt salt bank water ice loan green interview
DWGTM	iraq bush democrats clinton george president america war us democracy money debt pay credit tax company taxes loan income paying

applying pre-trained word embeddings can effectively capture semantics for topic modeling.

4.3 Ablative Study

We conduct an ablative study to evaluate whether the two reconstruction losses of the DualWG-Decoder and also the entropy regularization of θ (*i.e.*, Eq.15) have positive effects on topic extraction. To achieve this, we examine three simplified versions that independently remove the loss of \mathcal{G}^c ($\lambda_1 = 0$), the loss of \mathcal{G}^s ($\lambda_2 = 0$), and the entropy regularization term ($\lambda_3 = 0$).

We show the topic quality results of different versions of DWGTM on *Trec* when $k = 20$. As shown in Table 6, we can observe that the full DWGTM method outperforms all three simplified versions, indicating that all three components have positive effects on topic extraction. Specifically, DWGTM w/o the losses of \mathcal{G}^c and \mathcal{G}^s (*i.e.*, $\lambda_1 = 0$ and $\lambda_2 = 0$) lead to TSC deficiency over 0.01, indicating that the two reconstruction processes in the DualWG-Decoder can help capturing the semantic information of words. Besides, the gain of DWGTM over the version without the entropy regularization (*i.e.*, $\lambda_3 = 0$) shows more significant validity. This coincides with the fact that the entropy regularization tends to compute peakier topic proportions, which are beneficial for extracting topics from short texts with extremely limited words.

Specially, we have evaluated different values of $\{\lambda_1, \lambda_2, \lambda_3\}$ and also the threshold γ^* from the

Table 6: Results of the ablative study.

Metric	DWGTM	$\lambda_1 = 0$	$\lambda_2 = 0$	$\lambda_3 = 0$
TC	0.402	0.397	0.395	0.393
TSC	0.453	0.439	0.440	0.447

range $\{0.1, 0.2, \dots, 1\}$ in the early experiments. The results show that $\{\lambda_1, \lambda_2\}$ and λ_3 perform better with smaller and larger values, respectively; and γ^* performs relatively stable with different values. Due to the space limitation, we omit the detailed results and will show them in the next version.

5 Conclusion

In this paper, we develop a novel neural topic model for short texts, called DWGTM. The proposed DWGTM model extracts topics by simultaneously applying the word co-occurrence graph and word semantic correlation graph. Specifically, it consists of four main components: (1) Encode the word co-occurrence graph as word features. (2) Generate text features with word features, and encode them as topic proportions. (3) Reconstruct the texts with topical distributions. (4) Reconstruct both graphs with word features. We also propose a novel metric to evaluate the semantic coherence of topics, called TSC. Empirically, the effectiveness of DWGTM was validated on three benchmark datasets of short texts. We show that the topics learned by DWGTM can simultaneously capture meaningful patterns and semantic correlations of words.

Acknowledgements

The research is supported by the National Natural Science Foundation of China (NSFC) (No.61876071) and Scientific and Technological Developing Scheme of Jilin Province (No.20180201003SF, No.20190701031GH) and Energy Administration of Jilin Province (No.3D516L921421) and Key Program of Science and Technology Research during the 13th Five-Year Plan Period, the Educational Department of Jilin Province of China (Grant No. JJKH20200677KJ).

References

- David M Blei. 2012. Probabilistic topic models. *Communications of The ACM*, 55(4):77–84.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022.
- Sophie Burkhardt and Stefan Kramer. 2019. Decoupling sparsity and smoothness in the Dirichlet variational autoencoder topic model. *Journal of Machine Learning Research*, 20(131):1–27.
- Xueqi Cheng, Xiaohui Yan, Yanyan Lan, and Jiafeng Guo. 2014. BTM: topic modeling over short texts. *IEEE Transactions on Knowledge and Data Engineering*, 26(12):2928–2941.
- Adji B Dieng, Francisco J R Ruiz, and David M Blei. 2020. Topic modeling in embedding spaces. *Transactions of the Association for Computational Linguistics*, 8:439–453.
- Thomas L Griffiths and Mark Steyvers. 2004. Finding scientific topics. *Proceedings of the National Academy of Sciences of the United States of America*, 101(suppl 1):5228–5235.
- Lin Gui, Jia Leng, Gabriele Pergola, Yu Zhou, Ruifeng Xu, and Yulan He. 2019. Neural topic model with reinforcement learning. In *Empirical Methods in Natural Language Processing*, pages 3478–3483.
- Pankaj Gupta, Yatin Chaudhary, Thomas Runkler, and Hinrich Schuetze. 2020. Neural topic modeling with continual lifelong learning. In *International Conference on Machine Learning*, pages 3907–3917. PMLR.
- Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. 1999. An introduction to variational methods for graphical models. *Machine Learning*, 37(2):105–161.
- Diederik P Kingma and Max Welling. 2014. Auto-encoding variational Bayes. In *International Conference on Learning Representations*.
- Thomas N. Kipf and Max Welling. 2016a. Variational graph auto-encoders. In *NIPS Workshop on Bayesian Deep Learning*.
- Thomas N. Kipf and Max Welling. 2016b. Variational graph auto-encoders. In *NIPS Workshop on Bayesian Deep Learning*.
- Changchun Li, Ximing Li, Jihong Ouyang, and Yiming Wang. 2020a. Semantics-assisted Wasserstein learning for topic and word embeddings. In *IEEE International Conference on Data Mining*, pages 292–301.
- Chenliang Li, Haoran Wang, Zhiqian Zhang, Aixin Sun, and Zongyang Ma. 2016. Topic modeling for short texts with auxiliary word embeddings. In *International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 165–174.
- Shuangyin Li, Yu Zhang, and Rong Pan. 2020b. Bi-directional recurrent attentional topic model. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 14(6):1–30.
- Ximing Li, Changchun Li, Jinjin Chi, and Jihong Ouyang. 2018. Short text topic modeling by exploring original documents. *Knowledge and Information Systems*, 56(2):443–462.
- Ximing Li, Ang Zhang, Changchun Li, Lantian Guo, Wenting Wang, and Jihong Ouyang. 2019a. Relational biterm topic model: Short text topic modeling using word embeddings. *The Computer Journal*, 62(3):359–372.
- Ximing Li, Jiaojiao Zhang, and Jihong Ouyang. 2019b. Dirichlet multinomial mixture with variational manifold regularization: topic modeling over short texts. In *AAAI Conference on Artificial Intelligence*, pages 7884–7891.
- Chenghua Lin and Yulan He. 2009. Joint sentiment/topic model for sentiment analysis. In *ACM international Conference on Information and Knowledge Management*, pages 375–384.
- Zongyang Ma, Aixin Sun, Quan Yuan, and Gao Cong. 2012. Topic-driven reader comments summarization. In *ACM international Conference on Information and Knowledge Management*, pages 265–274.
- Yishu Miao, Lei Yu, and Phil Blunsom. 2016. Neural variational inference for text processing. In *International Conference on Machine Learning*, pages 1727–1736.
- Tomas Mikolov, Kai Chen, Greg S Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. In *International Conference on Learning Representations*.
- Andriy Mnih and Karol Gregor. 2014. Neural variational inference and learning in belief networks. In *International Conference on Machine Learning*, pages 1791–1799.

- Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell. 2000. Text classification from labeled and unlabeled documents using EM. *Machine Learning*, 39(2-3):103–134.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. GloVe: global vectors for word representation. In *Empirical Methods in Natural Language Processing*, pages 1532–1543.
- Xiaojun Quan, Chunyu Kit, Yong Ge, and Sinno Jialin Pan. 2015. Short and sparse text topic modeling via self-aggregation. In *International Joint Conference on Artificial Intelligence*, pages 2270–2276.
- Mehdi Rezaee and Francis Ferraro. 2020. A discrete variational recurrent topic model without the reparametrization trick. *arXiv preprint arXiv:2010.12055*.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pages 1278–1286.
- Tian Shi, Kyeongpil Kang, Jaegul Choo, and Chandan K. Reddy. 2018. Short-text topic modeling via non-negative matrix factorization enriched with local word-context correlations. In *The Web Conference*, pages 1105–1114.
- Akash Srivastava and Charles A. Sutton. 2017. Autoencoding variational inference for topic models. In *International Conference on Learning Representations*.
- Rui Wang, Deyu Zhou, and Yulan He. 2019. Atm: Adversarial-neural topic model. *Information Processing & Management*, 56(6):102098.
- Xiaobao Wu, Chunping Li, Yan Zhu, and Yishu Miao. 2020. Short text topic modeling with topic distribution quantization and negative sampling decoder. In *Empirical Methods in Natural Language Processing*, pages 1772–1882.
- Xiaohui Yan, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. 2013. A biterm topic model for short text. In *International Conference on World Wide Web*, pages 1445–1456.
- Liang Yang, Fan Wu, Junhua Gu, Chuan Wang, Xiaochun Cao, Di Jin, and Yuanfang Guo. 2020. Graph attention topic modeling network. In *Proceedings of The Web Conference 2020*, pages 144–154.
- Jianhua Yin and Jianyong Wang. 2014. A Dirichlet multinomial mixture model-based approach for short text clustering. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 233–242.
- Jichuan Zeng, Jing Li, Yan Song, Cuiyun Gao, Michael R Lyu, and Irwin King. 2018. Topic memory networks for short text classification. In *Empirical Methods in Natural Language Processing*, pages 3120–3131.
- Qile Zhu, Zheng Feng, and Xiaolin Li. 2018. GraphBTM: Graph enhanced autoencoded variational inference for biterm topic model. In *Empirical Methods in Natural Language Processing*, pages 4663–4672.
- Yuan Zuo, Junjie Wu, Hui Zhang, Hao Lin, Fei Wang, Ke Xu, and Hui Xiong. 2016a. Topic modeling of short texts: A pseudo-document view. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2105–2114.
- Yuan Zuo, Jichang Zhao, and Ke Xu. 2016b. Word network topic model: a simple but general solution for short and imbalanced texts. *Knowledge and Information Systems*, 48(2):379–398.