# Context Reinforced Neural Topic Modeling over Short Texts

**Jiachun Feng[1], Zusheng Zhang[1], Cheng Ding[1],**
**Yanghui Rao[1]***, **Haoran Xie[2]**

[1]School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China
[2]Department of Computing and Decision Sciences, Lingnan University, Hong Kong
{fengjch5, zhangzsh3, dingch6}@mail2.sysu.edu.cn,
raoyangh@mail.sysu.edu.cn, hrxie2@gmail.com

## Abstract

As one of the prevalent topic mining tools, neural topic modeling has attracted a lot of interests for the advantages of high efficiency in training and strong generalisation abilities. However, due to the lack of context in each short text, the existing neural topic models may suffer from feature sparsity on such documents. To alleviate this issue, we propose a Context Reinforced Neural Topic Model (CRNTM), whose characteristics can be summarized as follows. Firstly, by assuming that each short text covers only a few salient topics, CRNTM infers the topic for each word in a narrow range. Secondly, our model exploits pre-trained word embeddings by treating topics as multivariate Gaussian distributions or Gaussian mixture distributions in the embedding space. Extensive experiments on two benchmark datasets validate the effectiveness of the proposed model on both topic discovery and text classification.

## 1 Introduction

Mining topics from texts is significant for various applications of natural language processing, e.g., text classification, sentiment analysis, and recommender systems. As one of the most popular approaches for discovering latent topics, topic modeling (Blei et al., 2003; Yin and Wang, 2014) is capable of producing interpretable results. Generally, the dominant methods for parameter estimation in topic models are variational inference (Blei et al., 2003) and Gibbs sampling (Griffiths and Steyvers, 2004), both of which, however, require complex re-derivation when there is any minor changes to the model structure. Moreover, with the growth of data scale, the generative process is getting tricky and expensive, which leads to mathematically arduous derivation and high computational cost in training. These limitations make it difficult to extend the models to new variations flexibly.

With the development of deep learning, variational auto-encoder (VAE) (Kingma and Welling, 2014) has provided another promising solution for topic modeling. Benefiting from the flexibility of neural networks, the VAE framework is competent to learn complicated non-linear distributions and is convenient to be applied to various tasks. Furthermore, by using the back-propagation for optimization, VAE is highly efficient in training when compared with the models based on variational inference or Gibbs sampling. Considering the above advantages, several models built on VAE have been proposed, such as neural variational document model (NVDM) (Miao et al., 2016), neural variation latent Dirichlet allocation (NVLDA) (Srivastava and Sutton, 2017), Gaussian softmax model (GSM) (Miao et al., 2017), Dirichlet variational auto-encoder (DVAE) (Burkhardt and Kramer, 2019), and neural variational correlated topic modeling (NVCTM) (Liu et al., 2019). Although the VAE-based models reduce the computational cost impressively, they still suffer from the feature sparsity problem in short texts. In this case, the number of word occurrences in each text is relatively small, while the vocabulary corresponding to the corpus is large and the range of topics is broad.

To alleviate the above issue, many Bayesian approaches specific to short texts have been proposed (Yan et al., 2013; Lin et al., 2014; Li et al., 2016). Nonetheless, the above models all resort to Gibbs sampling or variational inference and hence incur the problems as mentioned before. In recent years, models built on VAE are also introduced for short texts, such as Graph-based inference network for the biterm topic model (GraphBTM) (Zhu et al., 2018) and neural sparsemax topic model (NSMTM) (Lin et al., 2019). However, learning context information is still challenging in these models due to significant word non-overlap in short texts. Relatedness information between word pairs may not be fully captured owing to the lack of

---

*The corresponding author.

word-overlap between such short messages.

In this paper, we propose a VAE-based topic model for short texts, where the context information for each text is effectively enhanced. Firstly, as can be observed, a short text generally covers only a subset of topics due to the limited text length. Therefore, we propose to filter irrelevant topics by setting a *topic controller* for each topic, encouraging each short text to focus on some salient topics. Through this way, the topic inference range is narrowed down and thus the topic sparsity can be achieved indirectly. Secondly, we incorporate pre-trained word embeddings into our model to explicitly enrich the context information. Specifically, we model each topic by a multivariate Gaussian distribution or a Gaussian mixture distribution in the embedding space, through which the relatedness of synonymous word pairs can be effectively inferred regardless of word non-overlap in short texts. In this way, our model can discover more interpretable topics than other topic models. We name the proposed model as Context Reinforced Neural Topic Model (CRNTM) and conclude the main contributions of our work as follows:

- We assume that each short text only focuses on a few salient topics. By setting a *topic controller* for each topic to filter irrelevant topics, CRNTM narrows down the topic inference space and achieves topic sparsity indirectly.

- Pre-trained word embeddings are incorporated to explicitly enrich the limited context information for each short message. By treating topic distributions over words as multivariate Gaussian distributions or Gaussian mixture distributions in the embedding space, CRNTM can produce more interpretable topics.

The rest of this paper is organized as follows. We discuss relevant research work in Section 2, and detail our proposed model in Section 3. Experimental settings and results are presented in Section 4. Finally, we draw the conclusion in Section 5.

## 2   Related Work

### 2.1   Neural Topic Modeling

With the development of deep learning, models built on neural networks have been proposed to discover latent topics, and most of them are based on VAE. In this vein, NVDM (Miao et al., 2016) is a neural variational framework for generative modeling on texts. It consists of an inference network and a multinomial softmax generative module. The inference network is used to estimate continuous hidden variables, which can represent the semantic content of documents, while the generative module aims to reconstruct the documents from the latent topic distributions. GSM (Miao et al., 2017) constructs the topic distributiona explicitly with a softmax function applied to the projection of the Gaussian random vector. ProdLDA (Srivastava and Sutton, 2017) replaces the mixture model in latent Dirichlet allocation (LDA) with a product of experts for better topic modeling. NVLDA (Srivastava and Sutton, 2017) approximates the Dirichlet distribution by using Laplace approximation. DVAE (Burkhardt and Kramer, 2019) decouples the properties of sparsity and smoothness by rewriting the Dirichlet parameter vector into a product of a sparse binary vector and a smoothness vector. NVCTM (Liu et al., 2019) enhances the capability of capturing the correlations among topics by reshaping topic distributions.

### 2.2   Short Text Topic Discovery

Topic models (Blei et al., 2003) provide a valuable solution for implicit semantic mining and understanding over documents. However, the feature sparsity problem arises for topic models when applied to short texts (Zhao et al., 2011), because such corpora are lack of word co-occurrences at the document level.

To overcome this limitation, the external documents were first introduced to enrich the contextual information in short texts (Sahami and Heilman, 2006; Jin et al., 2011; Phan et al., 2008). Unfortunately, it requires the external documents to be semantically close to the original corpus. Some approaches tackle the task by aggregating short texts into lengthy pseudo-documents and then applying a well established topic model. For this category of methods, short texts can be aggregated by utilizing the side information, e.g., user characteristics tags (Feng et al., 2020), user ID (Zhao et al., 2011), and timestamp (Diao et al., 2012). Another alternative methods directly modify the prior of Bayesian models to enrich word co-occurrences, so as to remedy the feature sparsity problem. For instance, the Biterm Topic Model (BTM) (Cheng et al., 2014), which models the global word co-occurrences at the corpus level, could lengthen short texts by con-

verting documents into biterm sets.

While the above methods are developed based on Bayesian models, some neural network based approaches have been introduced for short texts. Zhu et al. (2018) proposed a graph-based inference network named GraphBTM for accelerating the above BTM. This model sampled a fixed number of texts as a training instance to overcome the feature sparsity issue. Lin et al. (2019) proposed a neural model which is called NSMTM by providing sparse posterior distributions over topics based on the Gaussian sparsemax construction. Gupta et al. (2019) designed a neural autoregressive topic model named iDocNADE in a language modeling fashion. They also incorporated word embeddings as fixed prior in the model to introduce complementary information. However, the above approach does not model topic distributions explicitly.

# 3 Model Description

In this section, we describe our context reinforced neural topic model (CRNTM) in details. The overall architecture is illustrated in Figure 1, which consists of two major modules: an inference network for learning latent topics, and a Gaussian decoder for reconstructing documents.
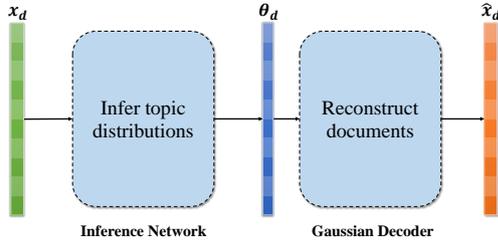


Figure 1: Structure of our CRNTM.

## 3.1 Problem Definition

Given a corpus with $D$ short texts, we denote the corresponding vocabulary as $W = \{w_1, w_2, ..., w_V\}$, with $V$ being the vocabulary size. Following (Miao et al., 2016), each document is processed into a bag-of-words (BOW) representation, i.e., $x_d = [x_{d,1}, x_{d,2}, ..., x_{d,V}]$, where $x_{d,i}$ denotes the number of times for word $w_i$ appearing in document $d$.

In the inference network, we use $\theta_d \in \mathbb{R}^K$ to denote the topic distribution of document $d$ and use $z_k \in \{z_1, z_2, ..., z_K\}$ to denote the topic as-

signment for an observed word, where $K$ denotes the number of topics inherent in the given corpus. Specifically, $\theta_d$ is drawn from the Gaussian distribution $\mathcal{N}(\mu_d, \Sigma_d)$, where both $\mu_d$ and $\Sigma_d$ are prior parameters. Furthermore, we set a *topic controller* $\lambda_{d,k} \in [0, 1]$ for each topic $z_k$ in document $d$: the topic will be kept when $\lambda_{d,k} = 1$, or it will be filtered out when $\lambda_{d,k} = 0$. For document $d$, the *topic controller* $\lambda_d = \{\lambda_{d,k}\}_{k=1}^K$ is drawn from a Beta distribution, i.e., $\lambda_d \sim \text{Beta}(\alpha_d, \beta_d)$, where $\alpha_d$ and $\beta_d$ are the prior parameters of $\lambda_d$.

For the Gaussian decoder, we denote the word embedding matrix corresponding to the vocabulary as $WE \in \mathbb{R}^{V \times r}$, where $r$ indicates the dimension of word embeddings. Moreover, the embedding of word $w_i$ is represented as $WE_i$. We use a matrix $TW \in \mathbb{R}^{K \times V}$ to denote the probabilities of words conditioned to topics, in which, $TW_{(k,i)}$ represents the conditional probability of word $w_i$ over topic $z_k$. In this study, $TW_{(k,i)}$ is drawn from a multivariate Gaussian distribution $\mathcal{N}(\mu_k, \Sigma_k)$ or a Gaussian mixture distribution, where $\mu_k \in \mathbb{R}^r$ and $\Sigma_k \in \mathbb{R}^{r \times r}$ are learnable parameters.

## 3.2 Inference Network

The first component of CRNTM is the inference network, which is applied to infer the topic distributions for the input documents. The structure of our inference network is illustrated in Figure 2. Following the framework of VAE, CRNTM infers the parameters $\mu_d$ and $\Sigma_d$ via deep neural networks that are elaborately designed for the observed data. Being fed with the input document $x_d$, the inference network first outputs an encoded vector $\pi_d$. Then, $\pi_d$ is linearly transformed to obtain $\mu_d$ and $\Sigma_d$, which are used to parameterize the Gaussian prior $\mathcal{N}(\mu_d, \Sigma_d)$. The above process is described by $\pi_d = \text{MLP}_1(x_d)$, $\mu_d = l_1(\pi_d)$, $\log \sigma_d = l_2(\pi_d)$, and $\Sigma_d = \text{diag}(\sigma_d^2)$, where $\text{MLP}_1$ is a multilayer perceptron, $l_1(\cdot)$ and $l_2(\cdot)$ are linear transformations. Note that the diagonal elements $\sigma_d^2$ of covariance matrix $\Sigma_d$ are non-negative. The output of $l_2(\cdot)$ is regarded as the logarithmic form $\log \sigma_d$, which is a real number.

A Gaussian random vector $h_d$ is passed through a softmax function to parameterize the multinomial document-topic distribution $\theta_d'$. The process is defined as: $\epsilon_d \sim \mathcal{N}(0, I^2)$, $h_d = \mu_d + \epsilon_d * \sigma_d$, and $\theta_d' = \text{softmax}(W_\theta \cdot h_d + b_\theta)$, where $h_d$ is drawn from the Gaussian prior $\mathcal{N}(\mu_d, \Sigma_d)$ with reparameterization, allowing the parameters to be
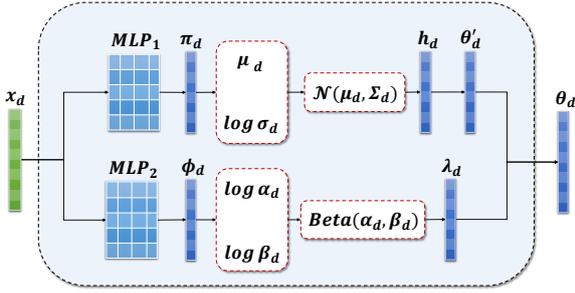
optimized by back-propagation.



Figure 2: Inference network.

Due to the limited text length, a short document only contains a few of words, resulting in the feature sparsity problem during the inference process. However, it can be observed that a short text generally focuses on a subset of topics. This inspires us to alleviate the above issue by narrowing down the scope for topic inference. Instead of letting the topic mixtures navigate freely in the simplex, CRNTM allows short texts to cover a narrow range of topics. This is achieved by setting a *topic controller* $\lambda_{d,k} \in [0,1]$ for each topic. Topic $z_k$ is focused when $\lambda_{d,k} = 1$, and it will be filtered out when $\lambda_{d,k} = 0$. The topic controllers are drawn from the Beta distribution, i.e., $\lambda_d \sim \text{Beta}(\alpha_d, \beta_d)$, as described in Section 3.1, which can guarantee that each component $\lambda_{d,k} \in \lambda_d$ is in the range of $[0,1]$. The parameters $\alpha_d$ and $\beta_d$ are inferred as follows:

$$\phi_d = \text{MLP}_2(x_d), \quad (1)$$
$$\log \alpha_d = l_3(\phi_d), \quad (2)$$
$$\log \beta_d = l_4(\phi_d), \quad (3)$$

where $\text{MLP}_2$ is a multilayer perceptron, $l_3(\cdot)$ and $l_4(\cdot)$ are linear transformations. The output of $l_3(\cdot)$ and $l_4(\cdot)$ are treated as the logarithmic form $\log \alpha_d$ and $\log \beta_d$, since both $\alpha_d$ and $\beta_d$ are non-negative.

The Beta sampling can not be differentiated directly, making it intractable to update model parameters through back-propagation. Therefore, we use the re-parameterization technique to obtain $\lambda_d$ by following (Naesseth et al., 2017). The sampling operation of $\text{Beta}(\alpha_d, \beta_d)$ can be decoupled into $\text{Gamma}(\alpha_d, 1)$ and $\text{Gamma}(\beta_d, 1)$, which is formulated by $\lambda_d = \frac{\lambda_{d,1}}{\lambda_{d,1}+\lambda_{d,2}}$, where $\lambda_{d,1} \sim$ $\text{Gamma}(\alpha_d, 1)$ and $\lambda_{d,2} \sim \text{Gamma}(\beta_d, 1)$. For the Gamma distribution $\text{Gamma}(\alpha, 1)$ with $\alpha > 1$, the re-parameterization can be accomplished by the reject sampling method:

$$\lambda_{d,1} = (\alpha_d - \frac{1}{3})(1 + \frac{\epsilon_d}{\sqrt{9\alpha_d - 3}})^3, \quad (4)$$

where $\epsilon_d \sim \mathcal{N}(0, I^2)$. On the other hand, the shape augmentation method is applied to convert $\alpha \leq 1$ to $\alpha > 1$ to increase the accept rate of each rejection sampler, which is formulated by $\lambda_{d,1} = \rho^{\frac{1}{\alpha_d}} \tilde{\lambda}_{d,1}$, where $\rho$ is drawn from a uniform distribution, i.e., $\rho \sim U[0,1]$, and $\tilde{\lambda}_{d,1} \sim$ $Gamma(\alpha + 1, 1)$ can be obtained according to Equation (4).

During the inference process, CRNTM determines whether topic $z_k$ is kept according to $\lambda_{d,k}$. By filtering out some certain topics, the short texts are allowed to focus on a few specific topics, and thus the feature sparsity problem can be alleviated. Finally, the topic distribution of document $x_d$ is obtained by $\theta_d = \theta'_d * \lambda_d$.

### 3.3 Gaussian Decoder

Context information is important for topic mining (Gupta et al., 2019). Words that appear together frequently are more likely to belong to the same topic, which implies that closer words in the embedding space are more likely to reflect the same topic. In Bayesian models, word embeddings that are trained on a large corpus have shown to effectively bring auxiliary context information for short texts (Li et al., 2016). Considering this advantage, we propose to introduce word embeddings into the decoder named Gaussian decoder. To our best knowledge, this is the first work of incorporating pre-trained word embeddings into the decoder of VAE to enhance the ability of capturing context information. The basic structure of our Gaussian decoder is shown in Figure 3(a).

The Gaussian decoder is applied to decode the topic distribution $\theta_d$, based of which a new document can be reconstructed. Concretely, the decoder employs the multivariate Gaussian distribution $\mathcal{N}(\mu_k, \Sigma_k)$ to model the $k$-th topic in the embedding space. Since the elements of the diagonal matrix $\Sigma_k$ are non-negative, a transformation similar to the one applied to $\Sigma_d$ is used here, as follows: $\Sigma_k = \text{diag}(\sigma_k^2)$.

By incorporating pre-trained word embeddings, the probability of word $w_i$ conditioned on topic $z_k$, i.e., $TW_{(k,i)}$, can be formulated by $TW_{(k,i)} = \frac{\exp(g(WE_i))}{(2\pi)^{r/2}|\Sigma_k|^{1/2}}$, where $r$ is the word embedding dimension, and $g(WE_i) = -\frac{1}{2}(WE_i -$

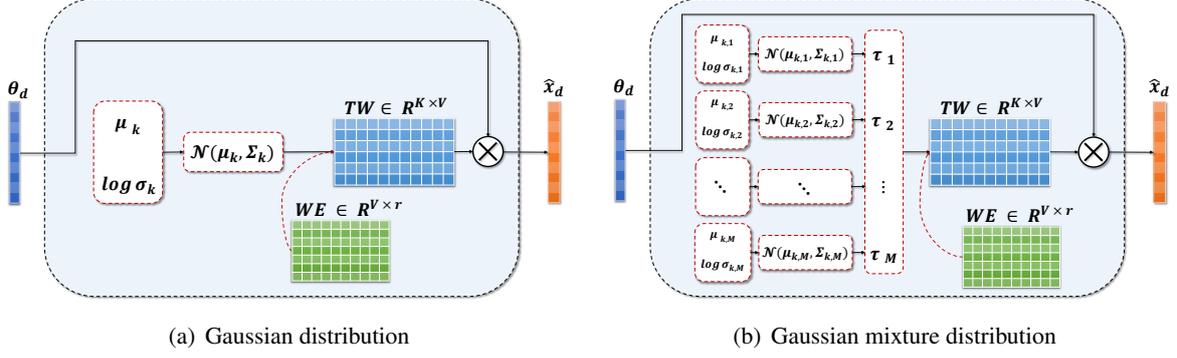(a) Gaussian distribution      (b) Gaussian mixture distribution

Figure 3: Gaussian decoder.

$\mu_k)^T \Sigma_k^{-1}(WE_i - \mu_k)$. It is worth noting that the parameters $\mu_k$ and $\Sigma_k$ can be regarded as the topic centroid and topic concentration in the embedding space. According to the properties of Gaussian distribution, words that are closer to the topic centroid have higher probabilities. Meanwhile, words in each topic are related to the context information implied by word embeddings. Therefore, CRNTM can enrich context information via pre-trained word embeddings to address the feature sparsity problem. Finally, we estimate the conditional probability $p(w_{d,i}|\theta_d, \lambda_d)$ by $p(w_{d,i}|\theta_d, \lambda_d) = \sum_k \theta_{d,k} \cdot TW_{(k,i)}$.

The above method can be easily adjusted by assuming that $TW_{(k,i)}$ obeys a Gaussian mixture distribution, as shown in Figure 3(b). In this case, $TW_{(k,i)} = \sum_{m=1}^M \tau_m \frac{\exp(g_m(WE_i))}{(2\pi)^{r/2}|\Sigma_{k,m}|^{1/2}}$, where $M$ is number of Gaussian components, $\tau_m$ is the coefficient of Gaussian distributions, and $g_m(WE_i) = -\frac{1}{2}(WE_i - \mu_{k,m})^T \Sigma_{k,m}^{-1}(WE_i - \mu_{k,m})$.

### 3.4 Optimization Objective

The optimization objective of CRNTM is $\mathcal{L} = \log p(D)$, where $\log p(D)$ is the likelihood of observed samples. According to the assumption, there is $\log p(D) = \sum_d \log p(d)$. Since the true distributions of documents are unknown, variational inference is used here to convert the optimization to its evidence lower bound (ELBO), that is, $\log p(d) \geq \mathcal{L}(d)$. According to the variational inference method, $\mathcal{L}(d)$ is derived as follows:

$$
\begin{aligned}
\mathcal{L}(d) &= \iint q(\theta_d, \lambda_d|x_d)[-\log q(\theta_d, \lambda_d|x_d) \\
&\quad + \log p(x_d, \theta_d, \lambda_d)]\mathrm{d}\theta_d \mathrm{d}\lambda_d \\
&= E_{q(\theta_d|x_d)q(\lambda_d|x_d)}[\log p(x_d|\theta_d, \lambda_d)] \\
&\quad - D_{KL}[q(\theta_d|x_d) \parallel p(\theta_d)] \\
&\quad - D_{KL}[q(\lambda_d|x_d) \parallel p(\lambda_d)], \quad (5)
\end{aligned}
$$

where $E_{q(\theta_d|x_d)q(\lambda_d|x_d)}[\log p(x_d|\theta_d, \lambda_d)]$ is often regarded as the reconstruction loss. $p(x_d|\theta_d, \lambda_d) = \prod_{i=1}^{n_d} p(w_{d,i}|\theta_d, \lambda_d)$. $p(\theta_d)$ is the prior distribution of $\theta_d$, $q(\theta_d|x_d)$ is the variational approximation of $p(\theta_d)$, $p(\lambda_d)$ is the prior distribution of $\lambda_d$, $q(\lambda_d|x_d)$ is the variational approximation of $p(\lambda_d)$, and $E_{q(\theta_d|x_d)}(\cdot)$ is approximated by sampling of $\theta_d \sim q(\theta_d|x_d)$. For $\theta_d$, we assume that the true prior $p(\theta_d)$ is a normal Gaussian distribution $\mathcal{N}(0, I)$ by following (Kingma and Welling, 2014; Miao et al., 2016; Liu et al., 2019). Therefore, the KL divergence term $D_{KL}[q(\theta_d|x_d) \parallel p(\theta_d)]$ can be derived by $D_{KL}[q(\theta_d|x_d) \parallel p(\theta_d)] = \frac{1}{2}(-n + \mu_d^2 - \log|\Sigma_d| + |\Sigma_d|)$.

Similarly, we take $Beta(\alpha', \beta')$ as the true prior of $p(\lambda_d)$, and the KL divergence term $D_{KL}[q(\lambda_d|x_d) \parallel p(\lambda_d)]$ can be computed by $D_{KL}[q(\lambda_d|x_d) \parallel p(\lambda_d)] = \ln \frac{\Delta(\alpha', \beta')}{\Delta(\alpha_d, \beta_d)} - (\alpha' - \alpha_d)\psi(\alpha_d) - (\beta' - \beta_d)\psi(\beta_d) + (\alpha' - \alpha_d + \beta' - \beta_d)\psi(\alpha_d + \beta_d)$, where $\Delta(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$, $\Gamma(\cdot)$ is the Gamma function, and $\psi(\cdot)$ is the Digamma function. In our model, the *topic controller* acts as a switch to filter out irrelevant topics and keep related topics with higher probabilities. Since $\alpha$ and $\beta$ determine the shape of Beta distribution, we set both $\alpha'$ and $\beta'$ to 0.5, so that the probabilities are sharp in values of 0 and 1.

## 4 Experiments

In this section, we first introduce the experimental setting, and then evaluate the effectiveness of our model by a series of experiments.

### 4.1 Datasets

To compare the model performance on both topic mining and text classification, we employ 20News-

Table 1: The statistics of datasets.

| Dataset | Train | Test | $V$ | $AvgD$ | $L$ |
|---|---|---|---|---|---|
| 20NewsGroups | $11,314$ | $7,531$ | $2,000$ | $12.3$ | $20$ |
| Snippets | $10,060$ | $2,280$ | $5,000$ | $14.3$ | $8$ |

Groups[1] and Snippets[2] with document labels as our datasets. 20NewsGroups is a collection of short news messages, which contains $11,314$ training and $7,531$ testing samples. These short texts are grouped into 20 different categories. Snippets is collected from the results of web search transaction over 8 domain labels. The officially divided 10,060 and 2,280 search transaction documents are used for training and testing, respectively. For data preprocessing, we remove stopwords and take the most frequent $2,000$ words and $5,000$ words as vocabularies. The statistics of the processed corpora are shown in Table 1, where $AvgD$ and $L$ denote the averaged number of words for each document and the number of categories, respectively.

## 4.2 Baseline Methods

We use the following mainstream VAE based methods as baselines for evaluation: NVDM (Miao et al., 2016), NVLDA & ProdLDA (Srivastava and Sutton, 2017), GSM (Miao et al., 2017), TMN (Zeng et al., 2018), NVCTM (Liu et al., 2019), and DVAE (Burkhardt and Kramer, 2019). Among these methods, NVDM is one of the first neural document models, NVLDA, ProdLDA, and GSM are classical neural topic models. TMN consists of a neural topic model and a topic memory mechanism, which are trained in an end-to-end learning manner. NVCTM exploits the Centralized Transformation Flow (CTF) to capture the topic correlations by reshaping topic distributions. DVAE achieves a competitive topic coherence and a high log-likelihood by decoupling the properties of sparsity and smoothness in VAE-based topic models for short texts.

Note that iDocNADE (Gupta et al., 2019) and NSMTM (Lin et al., 2019) are not adopted for comparison, because the former does not model topic distributions explicitly while the training process of the latter is too sensitive to continue based on our implementation. Besides, since the ELBO is

typically used and necessary to evaluate the performance of VAE based methods (Miao et al., 2016, 2017), we do not use Bayesian models such as (Yan et al., 2013; Lin et al., 2014; Li et al., 2016) as baselines for fair comparison. Finally, GraphBTM (Zhu et al., 2018) which only models a mini-corpus is unsuitable to be evaluated in this study.

## 4.3 Experimental Settings

In our experiments, the publicly available codes of NVDM[3], NVLDA & ProdLDA[4], TMN[5], and DVAE[6] are directly used. The baselines of GSM and NVCTM are implemented by us based on the code of NVDM, where the length of CTF in NVCTM is set to 10 according to the preliminary experiments. For our model, we use the widely adopted pre-trained word embedding from Glove (Pennington et al., 2014), and the embedding size is 300. All the models are trained alternatively by Adam optimizer with a learning rate of $1e^{-5}$ and a batch size of 64. In the task of topic discovery, $perplexity = \exp\{-\frac{1}{D}\sum_{d=1}^{D}\frac{1}{N_d}\sum_{i=1}^{N_d}\log p(w_{d,i})\}$ is used to evaluate the generalization performance of models on the testing set, where $D$ is the number of documents, $N_d$ is the number of words in document $d$ and $p(w_{d,i})$ is the log-likelihood of model on word $w_i$ in document $d$. To evaluate the quality of discovered topics, we also use the normalized pointwise mutual information (NPMI) (Lau et al., 2014) as the metric. The averaged values of NPMI on the top 5, 10, and 15 words for all topics is computed as the final results. For the task of text classification, we use the topic vector of each document generated by convergent models as the input of a classifier. MLPClassifier from scikit-learn[7] is chosen as the classifier in this study and accuracy is used as the metric. For each task, the topic numbers are set to 25, 50, and 100. We denote our models with Gaussian distribution and Gaussian mixture distribution in the decoder as CRNTM_GD and CRNTM_GMD, respectively. Unless explicitly specified, the number of Gaussian components is set to 25 for CRNTM_GMD. The source code, detailed parameter settings, and complementary

---

Table 2: Perplexity results of different models on both datasets, where the best scores are boldfaced.

| Model | 20NewsGroups | | | Snippets | | |
|---|---|---|---|---|---|---|
| | 25 | 50 | 100 | 25 | 50 | 100 |
| NVDM | 802 | 855 | 871 | 5144 | 5180 | 5328 |
| NVLDA | 1046 | 1252 | 1153 | 5336 | 5496 | 5374 |
| ProdLDA | 1106 | 1073 | 1035 | 5312 | 5379 | 5348 |
| GSM | 949 | 922 | 943 | 5237 | 5295 | 5434 |
| TMN | 1159 | 1136 | 1128 | **3177** | **3197** | **3236** |
| NVCTM | 758 | 738 | 744 | 5090 | 5121 | 5136 |
| DVAE | 1095 | 1066 | 1075 | 5090 | 5121 | 5136 |
| CRNTM_GD | 698 | 706 | 680 | 4822 | 4872 | 4861 |
| CRNTM_GMD | **574** | **586** | **590** | 4608 | 4695 | 4602 |

Table 3: Topic coherence results of different models on both datasets, where the best scores are boldfaced.

| Model | 20NewsGroups | | | Snippets | | |
|---|---|---|---|---|---|---|
| | 25 | 50 | 100 | 25 | 50 | 100 |
| NVDM | 0.041 | 0.061 | 0.053 | 0.068 | 0.067 | 0.069 |
| NVLDA | 0.065 | 0.062 | 0.061 | 0.042 | 0.045 | 0.041 |
| ProdLDA | 0.064 | 0.062 | 0.065 | 0.046 | 0.051 | 0.045 |
| GSM | 0.080 | 0.076 | 0.065 | 0.068 | 0.061 | 0.065 |
| TMN | 0.031 | 0.051 | 0.042 | 0.043 | 0.025 | 0.029 |
| NVCTM | 0.022 | 0.017 | 0.014 | 0.052 | 0.051 | 0.055 |
| DVAE | 0.065 | 0.075 | 0.069 | 0.039 | 0.052 | 0.040 |
| CRNTM_GD | 0.065 | 0.077 | 0.069 | 0.075 | 0.076 | 0.074 |
| CRNTM_GMD | **0.088** | **0.081** | **0.079** | **0.082** | **0.084** | **0.085** |

results of our models can be found at Github[8].

## 4.4 Comparison with Baselines

Table 2 presents the test document perplexities of all models, from which we can observe that our CRNTM_GD and CRNTM_GMD achieve the best results in most cases. Specifically, TMN performs the best on Snippets. The reason may be that TMN is basically a supervised model for text classification and that the supervision from labels can help mining topics on a corpus consisting of less formal texts (i.e., Snippets). We also report the results of topic coherence in Table 3. It can be observed that both of our models are significantly better than all the baselines, which shows that they are able to discover more meaningful and interpretable topics. The performance comparisons for text classification are shown in Table 4. We can find that CRNTM_GD and CRNTM_GMD obtain competitive performances when compared with the benchmark methods, which validates the effectiveness of our models on generating representative vectors for short text classification.
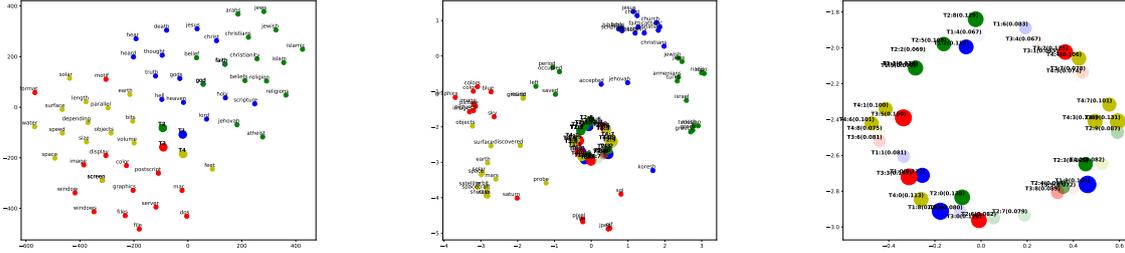
---

[8] https://github.com/Deloris-NLP/CRNTM

Table 4: Classification accuracies of different models on both datasets, where the best scores are boldfaced.

| Model | 20NewsGroups | | | Snippets | | |
|---|---|---|---|---|---|---|
| | 25 | 50 | 100 | 25 | 50 | 100 |
| NVDM | 0.64 | 0.64 | 0.67 | 0.15 | 0.17 | 0.16 |
| NVLDA | 0.40 | 0.45 | 0.42 | 0.12 | 0.13 | 0.13 |
| ProdLDA | 0.43 | 0.44 | 0.40 | 0.14 | 0.14 | 0.15 |
| GSM | 0.45 | 0.46 | 0.45 | 0.11 | 0.12 | 0.11 |
| TMN | 0.40 | 0.48 | 0.51 | 0.15 | 0.16 | 0.13 |
| NVCTM | 0.64 | 0.64 | 0.65 | **0.16** | **0.18** | **0.18** |
| DVAE | 0.32 | 0.37 | 0.34 | 0.08 | 0.09 | 0.06 |
| CRNTM_GD | 0.64 | **0.65** | **0.68** | 0.15 | 0.16 | 0.14 |
| CRNTM_GMD | **0.69** | **0.65** | 0.66 | **0.16** | 0.16 | 0.17 |

## 4.5 Evaluation on Gaussian Decoder via Topic Visualization

To investigate the quality of topics discovered by our models, we report top 15 words of 4 representative topics and visualize these topics by their embedding vectors using 20NewsGroups. Particularly, we extract $\mu_k$ of Gaussian distributions as the topic centroid and utilize t-SNE (van der Maaten, 2009) for visualization. Topic visualization of the results in CRNTM_GD is depicted in Figure 4(a). The points with different colors indicate different topics, and the centroid of topic $k$ is denoted as $Tk$. For the convenience of comparison, we manually annotate each topic by referring to the ground truth category. Accordingly, $T1$, $T2$, $T3$, and $T4$ in CRNTM_GD are annotated as "soc.religion.christian", "talk.politics", "comp.sys.ibm.pc.hardware", and "comp.graphics", respectively. We can see that all top words of the same topics are close to each other and to the corresponding topic centroids in the continuous vector space. This validates that our Gaussian decoder can effectively capture the context information via word embeddings in mining topics. We also present 4 topics generated by CRNTM_GMD whose semantics are similar to those in CRNTM_GD to verify the effectiveness of Gaussian mixture distributions. Topic visualization of the results in CRNTM_GMD is shown in Figure 4(b). For clarity, the coefficients of Gaussian components are indicated by different point sizes and shades of colour. The bigger the points are and the stronger the color is, the higher coefficient of the corresponding Gaussian components is. The topic centroids and their probabilities are detailed in Figure 4(c). We can observe that for topic $T3$ named as "comp.graphics", the main components such as $T3{:}5$, $T3{:}3$, and $T3{:}0$ are close to the cluster of red points, while $T2{:}8$ and $T2{:}5$ are close to sub-clusters of top words in

(a) Top 15 words in CRNTM_GD, and $T1, T2, T3, T4$ are topic centroids.

(b) Top 15 words in CRNTM_GMD. $Tk{:}m$ is the $m$th centroid of topic $k$.

(c) Topic centroids and their probabilities in CRNTM_GMD.

Figure 4: Characteristics of 4 representative topics generated by our models on 20NewsGroups.

$T2$.

To make a comprehensive comparison, we present the results of all models on generating topic "soc.religion.christian" in Table 5. It can be observed that our models can discover quite meaningful topics. We can also observe that the numbers of semantically irrelevant words of TMN and NVCTM are more than other models, which is consistent to the topic coherence results in Table 3.

Table 5: Top 10 words of manually labeled topic "soc.religion.christian" from all models on 20News-Groups, where irrelevant words are underlined.

| Model | Top words |
|---|---|
| NVDM | god sin scsi bible jesus rutgers homosexuality christian ide christians |
| NVLDA | god scsi sin drive jesus bible christian christians homosexuality love |
| ProdLDA | god christians jesus bible doctrine interpretation belief homosexuality christianity eternal |
| GSM | god jesus bible christ church people christian believe christians sin |
| TMN | sin myers eternal president mary god heaven christ doctor jobs |
| NVCTM | church catholic christians magnus scripture duke andrew turkey sex christianity |
| DVAE | jesus scripture christ bible doctrine sin christians god canon homosexuality |
| CRNTM_GD | jesus god christ heaven death holy truth gods faith lord |
| CRNTM_GMD | god christians bible christ jesus sin religion church lord doctrine |

### 4.6 Impact of Gaussian Mixture Numbers

We further study the impact of the number of Gaussian components. Table 6 presents the results of CRNTM_GMD on 20NewsGroups when varying component numbers under 25 topics. We can observe that CRNTM_GMD with more Gaussian com-

Table 6: Performance of CRNTM_GMD with different Gaussian mixture numbers on 20NewsGroups, where the best results are boldfaced.

| M | Perplexity | Coherence | Accuracy |
|---|---|---|---|
| 5 | 634 | 0.060 | 0.65 |
| 10 | 616 | 0.071 | 0.66 |
| 15 | 597 | 0.081 | 0.68 |
| 20 | 588 | 0.084 | 0.68 |
| 25 | 574 | **0.088** | **0.69** |
| 30 | 574 | 0.080 | 0.66 |
| 35 | **571** | 0.081 | 0.66 |

ponents generally performs better than that with less ones, which demonstrates that a more sophisticated mixture possesses a stronger capacity of learning high quality topics. The best topic coherence and classification accuracy are obtained when the component number is set to 25, and a larger value may not further boost the model performance.

## 5 Conclusion

In this paper, we propose a Context Reinforced Neural Topic Model (CRNTM) to address the feature sparsity problem in short texts. By introducing a *topic controller* to the inference network, CRNTM infers the topic for each word in a narrow range. Besides, pre-trained word embeddings are incorporated with multivariate Gaussian distributions or Gaussian mixture distributions into our model to enrich the context information of short messages. To quantitatively validate the effectiveness of CRNTM, we conduct various experiments on two benchmark datasets in terms of perplexity, topic coherence, and text classification accuracy. The results indicate that the proposed model largely

improves the performance of topic modeling by enriching the context information effectively.

## References

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.

Qiming Diao, Jing Jiang, Feida Zhu, and Ee-Peng Lim. 2012. Finding bursty topics from microblogs. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*, pages 536–544.

Jiachun Feng, Yanghui Rao, Haoran Xie, Fu Lee Wang, and Qing Li. 2020. User group based emotion detection and topic discovery over short text. *World Wide Web*, 23:1553–1587.

Thomas L Griffiths and Mark Steyvers. 2004. Finding scientific topics. *Proceedings of the National Academy of Sciences of the United States of America*, 101:5228–5235.

Pankaj Gupta, Yatin Chaudhary, Florian Buettner, and Hinrich Schütze. 2019. Document informed neural autoregressive topic models with distributional prior. In *Proceedings of the 33rd Conference on Artificial Intelligence*, pages 6505–6512.

Ou Jin, Nathan Nan Liu, Kai Zhao, Yong Yu, and Qiang Yang. 2011. Transferring topical knowledge from auxiliary long texts for short text clustering. In *Proceedings of the 20th Conference on Information and Knowledge Management*, pages 775–784.

Diederik P. Kingma and Max Welling. 2014. Autoencoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations*.

Jey Han Lau, David Newman, and Timothy Baldwin. 2014. Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 530–539.

Chenliang Li, Haoran Wang, Zhiqian Zhang, Aixin Sun, and Zongyang Ma. 2016. 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*, pages 165–174.

Tianyi Lin, Zhiyue Hu, and Xin Guo. 2019. Sparsemax and relaxed wasserstein for topic sparsity. In *Proceedings of the 20th ACM International Conference on Web Search and Data Mining*, pages 141–149.

Tianyi Lin, Wentao Tian, Qiaozhu Mei, and Hong Cheng. 2014. The dual-sparse topic model: mining focused topics and focused terms in short text. In *Proceedings of the 23rd International World Wide Web Conference*, pages 539–550.

Luyang Liu, Heyan Huang, Yang Gao, Yongfeng Zhang, and Xiaochi Wei. 2019. Neural variational correlated topic modeling. In *Proceedings of the 28th International World Wide Web Conference*, pages 1142–1152.

Laurens van der Maaten. 2009. Learning a parametric embedding by preserving local structure. In *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics*, pages 384–391.

Yishu Miao, Edward Grefenstette, and Phil Blunsom. 2017. Discovering discrete latent topics with neural variational inference. In *Proceedings of the 34th International Conference on Machine Learning*, pages 2410–2419.

Yishu Miao, Lei Yu, and Phil Blunsom. 2016. Neural variational inference for text processing. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 1727–1736.

Christian A. Naesseth, Francisco J. R. Ruiz ands Scott W. Linderman, and David M. Blei. 2017. Reparameterization gradients through acceptance-rejection sampling algorithms. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, pages 489–498.

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1532–1543.

Xuan Hieu Phan, Minh Le Nguyen, and Susumu Horiguchi. 2008. 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*, pages 91–100.

Mehran Sahami and Timothy D. Heilman. 2006. A web-based kernel function for measuring the similarity of short text snippets. In *Proceedings of the 15th International Conference on World Wide Web*, pages 377–386.

Akash Srivastava and Charles A. Sutton. 2017. Autoencoding variational inference for topic models. In *Proceedings of the 5th International Conference on Learning Representations*.

Xiaohui Yan, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. 2013. A biterm topic model for short texts. In *Proceedings of the 22nd International World Wide Web Conference*, pages 1445–1456.

Jianhua Yin and Jianyong Wang. 2014. 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*, pages 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 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3120–3131.

Wayne Xin Zhao, Jing Jiang, Jianshu Weng, Jing He, Ee-Peng Lim, Hongfei Yan, and Xiaoming Li. 2011. Comparing twitter and traditional media using topic models. In *Proceedings of the 33rd European Conference on IR Research*, pages 338–349.

Qile Zhu, Zheng Feng, and Xiaolin Li. 2018. Graphbtm: Graph enhanced autoencoded variational inference for biterm topic model. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4663–4672.