

# NEURAL VARIATIONAL SPARSE TOPIC MODEL

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Effectively inferring discriminative and coherent latent topics of short texts is a critical task for many real world applications. Nevertheless, the task has been proven to be a great challenge for traditional topic models due to the data sparsity problem induced by the characteristics of short texts. Moreover, the complex inference algorithm also become a bottleneck for these traditional models to rapidly explore variations. In this paper, we propose a novel model called Neural Variational Sparse Topic Model (NVSTM) based on a sparsity-enhanced topic model named Sparse Topical Coding (STC). In the model, the auxiliary word embeddings are utilized to improve the generation of representations. The Variational Autoencoder (VAE) approach is applied to inference the model efficiently, which makes the model easy to explore extensions for its black-box inference process. Experimental results on Web Snippets, 20Newsgroups, BBC and Biomedical datasets show the effectiveness and efficiency of the model.

## 1 INTRODUCTION

With the great popularity of social networks and Q&A networks, short texts have been the prevalent information format on the Internet. Uncovering latent topics from huge volume of short texts is fundamental to many real world applications such as emergencies detection (Sakaki et al., 2010), user interest modeling (Sasaki et al., 2014), and automatic query-reply (Peng et al., 2016). However, short texts are characteristic of short document length, a very large vocabulary, a broad range of topics, and snarled noise, leading to much sparse word co-occurrence information. Thus, the task has been proven to be a great challenge to traditional topic models. Moreover, the complex inference algorithm also become a bottleneck for these traditional models to rapidly explore variations.

To address the aforementioned issue, there are many previous works introducing new techniques such as word embeddings and neural variational inference to topic models. Word embeddings are the low-dimensional real-valued vectors for words. It have proven to be effective at capturing syntactic and semantic information of words. Recently, many works have tried to incorporate word embeddings into topic models to enrich topic modeling (Das et al., 2015; Hu & Tsujii, 2016; Xun et al., 2017). Yet these models general rely on computationally expensive inference procedures like Markov Chain Monte Carlo, which makes them hard to rapidly explore extensions. Even minor changes to model assumptions requires a re-deduction of the inference algorithms, which is mathematic challenging and time consuming. With the advent of deep neural networks, the neural variational inference has emerged as a powerful approach to unsupervised learning of complicated distributions (Kingma & Welling, 2013; Rezende et al., 2014; Mnih & Gregor, 2014). It approximates the posterior of a generative model with a variational distribution parameterized by a neural network, which allows back-propagation based function approximations in generative models. The variational autoencoder (VAE) (Kingma & Welling, 2013), one of the most popular deep generative models, has shown great promise in modeling complicated data. Motivated by the promising potential of VAE in building generative models with black-box inference process, there are many works devoting to inference topic models with VAE (Srivastava & Sutton, 2017; Miao et al., 2017; Card et al., 2017). However, these methods yield the same poor performance in short texts as LDA.

Based on the analysis above, we propose a Neural Variational Sparse Topic Model (NVSTM) based on a sparsity-enhanced topic model STC for short texts. The model is parameterized with neural networks and trained with VAE. It still follows the probabilistic characteristics of STC. Thus, the model inherits the advantages of both sparse topic models and deep neural networks. Additionally, we exploit the auxiliary word embeddings to improve the generation of short text representations.

To summarize, the main contributions of this paper are as follows:

1. We propose a novel Neural Variational Sparse Topic Model (NVSTM) to learn sparse representations of short texts. The VAE is utilized to inference the model effectively.
2. The general word semantic information is introduced to improve the sparse representations of short texts via word embeddings.
3. We conduct experiments on four datasets. Experimental results demonstrate our model’s superiority in topic coherence and text classification accuracy.

The rest of this paper is organized as follows. First, we reviews related work. Then, we present the details of the proposed NVSTM, followed by the experimental results. Finally, we draw our conclusions.

## 2 RELATED WORK

**Topic models.** Traditional topic models and their extensions (Archambeau et al., 2015; Blei et al., 2003; Mcauliffe & Blei, 2008) have been widely applied to many tasks such as information retrieval, document classification and so on. These models work well on long texts which have abundant word co-occurrence information for learning, but get stuck in short texts. There have been many efforts to address the data sparsity problem of short texts. To achieve sparse representations in the document-topic and topic-term distributions, Williamson et al. (2010) introduced a Spike and Slab prior to model the sparsity in finite and infinite latent topic structures of text. Similarly, Lin et al. (2014) proposed a dual-sparse topic model that addresses the sparsity in both the topic mixtures and the word usage. These models are inspired by the effect of the variation of the Dirichlet prior on the probabilistic topic models. There are also some non-probabilistic sparse topic models aiming at extracting focused topics and words by imposing various sparsity constraints. Heiler & Schnörr (2006) formalized topic modeling as a problem of minimizing loss function regularized by lasso. Subsequently, Zhu & Xing (2011) presented sparse topical coding (STC) by utilizing the Laplacian prior to directly control the sparsity of inferred representations. However, over complicated inference procedure of these sparse topic models has limited their applications and extensions.

**Topic Models with Word Embeddings.** Since word embeddings can capture the semantic meanings of words via low-dimensional real-valued vectors, there have been a large number of works on topic models that incorporate word embeddings to improve topic modeling. (Das et al., 2015) proposed a new technique for topic modeling by treating the document as a collection of word embeddings and topics itself as multivariate Gaussian distributions in the embedding space. However, the assumption that topics are unimodal in the embedding space is not appropriate, since topically related words can occur distantly from each other in the embedding space. Therefore, (Hu & Tsujii, 2016) proposed latent concept topic model (LCTM), which modeled a topic as a distribution of concepts, where each concept defined another distribution of word vectors. (Nguyen et al., 2015) proposed Latent Feature Topic Modeling (LFTM), which extended LDA to incorporate word embeddings as latent features. Lately, (Xun et al., 2017) proposed a novel correlated topic model using word embeddings, which is enable to exploit the additional word-level correlation information in word embeddings and directly model topic correlation in the continuous word embedding space. However, these models also have trouble to rapidly explore extensions.

**Neural Variational Inference for topic models.** Neural variational inference is capable of approximating the posterior of a generative model with a variational distribution parameterized by a neural network (Kingma & Welling, 2013; Rezende et al., 2014; Mnih & Gregor, 2014). The variational autoencoder (VAE), as one of the most popular neural variational inference approach, has shown great promise in building generative models with black-box inference process (Kingma & Welling, 2013). To break the bottleneck of over complicated inference procedure in topic models, there are many efforts devoting to inference topic models with VAE. Srivastava & Sutton (2017) presents auto-encoding variational Bayes (AEVB) based inference method for latent Dirichlet allocation (LDA), tackling the problems caused by the Dirichlet prior and component collapsing in AEVB. Miao et al. (2017) presents alternative neural approaches in topic modeling by providing parameterized distributions over topics. It allows training the topic model via back-propagation under the framework of neural variational inference. Card et al. (2017) combines certain motivating ideas behind variations on topic models with modern techniques for variational inference to produce

a flexible framework for topic modeling that allows for rapid exploration of different models. Nevertheless, aforementioned works are based on traditional LDA, thus bypass the sparsity problem of short texts.

Drawing inspiration from the above analysis, we propose a novel neural variational sparse topic model NVSTM based on VAE for short texts, which combines the merits of neural networks and sparsity-enhanced topic models.

### 3 NEURAL VARIATIONAL SPARSE TOPIC MODEL

In this section, we start from describing Sparse Topical Coding (STC). Based on it, we further propose Neural Variational Sparse Topic Model (NVSTM). Later, we focus on the discussion of the inference process for NVSTM.

#### 3.1 SPARSE TOPICAL CODING

In standard STC, each document and each word is represented as a low-dimensional code in topic space. Based on the topic dictionary  $\beta$  with  $K$  topic bases sampled from a uniform distribution, the generative process is described as follows:

1. Sample the document code  $\theta$  from a prior  $p(\theta) \sim Laplace(\lambda_1)$ .
2. For each observed word  $n$ :
  - (a) Sample the word code  $s_n$  from a conditional distribution  $p(s_n|\theta) \sim supergaussian(\theta, \lambda_2)$ .
  - (b) Sample the observed word count  $w_n$  from a distribution  $p(w_n|s_n^T \beta_n) \sim Poisson(s_n^T \beta_n)$

STC reconstructs each observed word count from a linear combination of a set of topic bases, where the word code is utilized as the coefficient vector. According to the above generative process, we have the joint distribution:

$$p(\theta, s, w|\beta) = \prod_d p(\theta_d) \prod_n p(s_n|\theta) p(w_n|s_n, \beta) \quad (1)$$

To simplify the calculation, the document code can be collapsed and later obtained via an aggregation of the individual word codes of all its terms. Although STC has closed form coordinate descent equations for parameters  $(\theta, s, \beta)$ , it is inflexible for its complex inference process.

#### 3.2 GENERATIVE PROCESS FOR NVSTM

To address the aforementioned issue, we introduce black box inference methods into STC. We present NVSTM based on VAE via the reparameterization trick and introduces word embeddings. According to (Bai et al., 2013), we collapse the document code and obtain it via :

$\theta_d = \frac{\sum_{n=1}^{N_d} \sum_{k=1}^K s_{d,nk} \beta_{kn}}{\sum_{n=1}^{N_d} \sum_{k=1}^K s_{d,nk}}$  to simplify the model structure. Analogous to the generative process in STC, our model follows the generative story below for each document  $d$ :

1. For each word  $n$  in document  $d$ :
  - (a) Sample a latent variable word code  $s_n \sim U(-0.5, 0.5)$ .
  - (b) Sample the observed word count  $w_n$  from  $p(w_n|s_n^T \beta_n) \sim Poisson(s_n^T \beta_n)$

The graphical representation of NVSTM is depicted in Figure 1. To yield sparse word codes, we choose the Laplace distribution as the variational distribution of word codes, in which the base distribution  $p(\varepsilon)$  is uniform distribution  $U(-0.5, 0.5)$  according to the reparameterisation method. Therefore, in the above generative process, each word code vector is generated from the uniform prior distribution. The observed word count is sampled from Poisson distribution. Different from traditional STC, we replace the uniform distribution of the topic dictionary with a topic dictionary neural network. In the topic dictionary neural network, we introduce the word semantic information

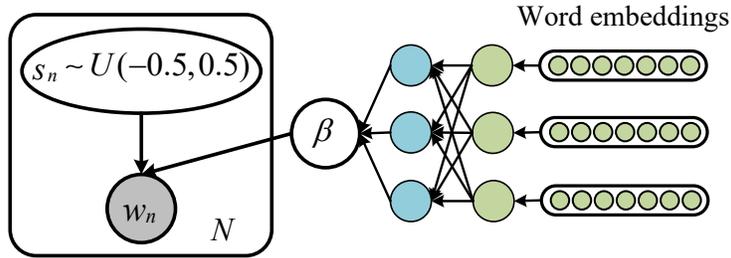


Figure 1: The graphical model of NVSTM.

via word embeddings to enrich the feature space for short texts. The topic dictionary neural network is comprised of two layers:

**Word embedding layer** ( $E \in \mathbb{R}^{N \times 300}$ ): Supposing the word number of the vocabulary is  $N$ , this layer devotes to transform each word to a distributed embedding representation. Here, we adopt the pre-trained embeddings by GloVe based on a large Wikipedia dataset<sup>1</sup>. Given a word embedding matrix  $E$ , we map each word to a 300-dimensional embedding vector, which can capture subtle semantic relationships between words.

**Topic dictionary layer** ( $\beta \in \mathbb{R}^{N \times K}$ ): This layer aims at converting  $E$  to a topic dictionary similar to the one in STC.

$$\beta(n) = \text{relu}(E \times W) \quad (2)$$

where  $W \in \mathbb{R}^{300 \times K}$  is the weight matrix between the word embedding layer and the topic dictionary layer. To conform to the framework of STC, we make a simplex projection among the output of topic dictionary neural network. We normalize each column of the dictionary via the simplex projection as follow:

$$\beta_{.k} = \text{project}(\beta_{.k}), \forall k \quad (3)$$

The simplex projection is the same as the sparsemax activation function in (Martins & Astudillo, 2016), which declares how the Jacobian of the projection can be efficiently computed, providing the theoretical base of its employment in a neural network trained with backpropagation. After the simplex projection, each column of the topic dictionary is promised to be sparse, non-negative and united.

Based the above generative process, the traditional variational inference for the model is to minimize the follow optimization problem, which is a lower bound to the marginal log likelihood:

$$L(\gamma|\beta) = D_{KL}[q(s|\gamma)||p(s|w, \beta)] - \log p(w|\beta) \quad (4)$$

where  $q(s|\gamma)$  is approximate variational posterior, and  $\gamma$  is the variational parameter.

### 3.3 VARIATIONAL AUTOENCODER FOR NVSTM

In this paper, we employ the VAE to carry out neural variational inference for our model. Variational Autoencoder (VAE) is one of the most popular deep generative network. It is a black-box variational method which bridges the conceptual and language gap of neural networks and probability generative models. From neural network perspective, a variational autoencoder consists of an encoder network, a decoder network, and a loss function. In our model, the encoder network is to parametrize the approximate posterior  $q_\theta(s|w)$ , which takes input as the observed word count to output the latent variable  $s$  with the variational parameters  $\theta$ . To derive the sparse word representation, we choose a Laplace variational distribution  $q_{\theta_n}(s_n|w_n) = \text{Laplace}(s_n; 0, b_n(w_n))$ , and define the encoder network as a feedforward neural network  $b_n(w_n) = f(w_n, \theta_n)$ . The decoder network outputs the observed data  $w$  with given  $s$  and the generative parameters  $\phi$ , which is denoted as  $p_\phi(w|s, \beta)$ . According to STC, we choose a Poisson distribution  $p_{\phi_n}(w_n|s_n) = \text{Poisson}(w_n; (s_n \beta_n))$ , and define the decoder network as a feedforward neural network ( $w_n(s_n) = f(s_n, \beta_n), \beta_n = f(E_n, \phi_n)$ ). Based on VAE, we rewrite the ELBO as:

$$L(\theta, \phi|\beta) = -D_{KL}[q_\theta(s|w)||p(s)] + E_{q_\theta(s|w)}(\log p_\phi(w|s, \beta)) \quad (5)$$

<sup>1</sup><http://nlp.stanford.edu/projects/glove/>

The first term is a regularizer that constraints the Kullback-Leibler divergence between the encoder’s distribution and the prior of the latent variables. The second term is the reconstruction loss, which encourages the decoder to reconstruct the data in minimum cost.

### 3.4 THE REPARAMETERIZATION TRICK AND OPTIMIZING

We devote to differentiate and optimize the lower bound above with stochastic gradient decent (SGD). However, the gradient of the lower bound is tricky since the error is unable to back propagate through a random drawn variable  $s$ , which is a non-continuous and has no gradient. Similar to the standard VAE, we make a differentiable transformation, called reparameterization trick. We approximate  $s$  with an auxiliary noise variable  $\varepsilon \sim U(-0.5, 0.5)$ :

$$s_n \sim Laplace(0, b_n) \rightarrow s_n = -b_n \text{sign}(\varepsilon) \ln(1 - 2|\varepsilon|), \varepsilon \sim U(-0.5, 0.5) \tag{6}$$

Through reparametrization, we can take  $s$  as a function with the parameter  $b$  deriving from the encoder network. It allows the reconstruction error to flow through the whole network. Figure 2 presents the complete VAE inference process for NVSTM. Moreover, in order to achieve interpretable word codes as in STC, we constrain  $s$  to be non-negative, activation function on the output  $s$  of encoder. After apply the reparameterization trick to the variational lower bound, we can yield

$$L(\Theta) = \sum_{d=1}^D \sum_{n=1}^N (1 + \log 2b_{d,n}) + E_{\varepsilon \sim U(-0.5, 0.5)} \sum_{d=1}^D \sum_{n=1}^N \left( \sum_{k=1}^K s_{d,nk} \beta_{nk} - w_{d,n} \ln \left( \sum_{k=1}^K s_{d,nk} \beta_{nk} \right) \right) \tag{7}$$

where  $s_n = -b_n \text{sign}(\varepsilon) \ln(1 - 2|\varepsilon|), \varepsilon \sim U(-0.5, 0.5)$ , and  $\Theta$  represents the set of all the model. As explained above, the decoding term  $\log p(w_{d,n} | s_{d,nk}, \beta_{nk})$  is the Poisson distribution, and  $\beta$  is generated by a topic dictionary neural network. After the differentiable transformation, the variation objective function can be computed in closed form and efficiently solved with SGD. The detailed algorithm is shown in Algorithm 1.

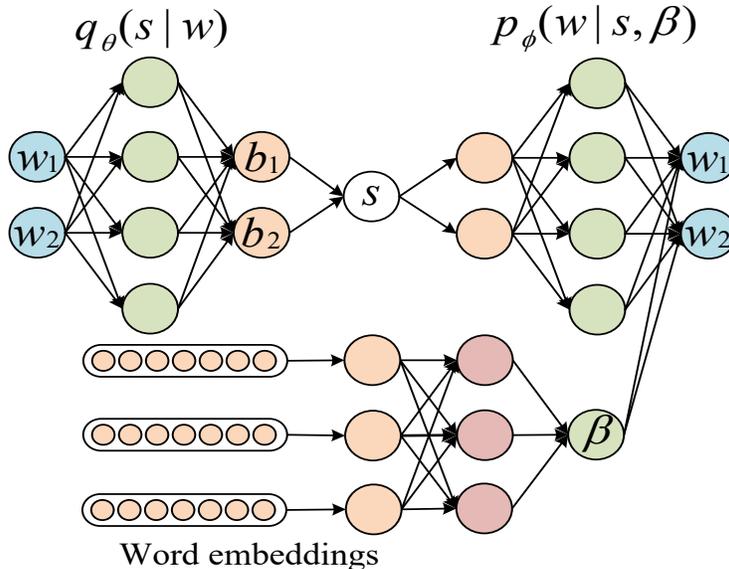


Figure 2: The VAE inference for NVSTM, the VAE is implemented as a feedforward neural network.

## 4 EXPERIMENTS

### 4.1 DATA AND SETTING

To evaluate the performance of our model, we present a series of experiments below. The objectives of the experiments include: (1) the qualitative evaluations: classification accuracy of documents and

**Algorithm 1** Training Algorithm for NVSTM

---

**Input:** initialize  $\theta, \phi, W$

- 1: **repeat**
- 2:  $w^M \leftarrow$  Random mini-batch of  $M$  word counts from full datasets
- 3:  $\varepsilon \leftarrow$  Random samples from noise distribution  $p(\varepsilon)$
- 4:  $g \leftarrow \nabla_{\theta, \phi, W} L(\theta, \phi; w^M, \varepsilon)$
- 5:  $\theta, \phi, W \leftarrow$  Update parameters using SGD
- 6: **until** convergence

---

sparse ratio of latent representations; (2) the qualitative inspection: the quality of extracted topics and document representations. Our evaluation is based on the four datasets:

- **20Newsgroups:** The classic 20 newsgroups dataset, which is comprised of 18775 news-group articles with 20 categories, and contains 60698 unique words <sup>2</sup>.
- **Web Snippet:** The web snippet dataset, which includes 12340 Web search snippets in 8 categories. We remove the words with fewer than 3 characters or whose document frequency less than 3 in the dataset. After the preprocessing, it contains 5581 unique words.<sup>3</sup>
- **BBC:** It consists of 2225 BBC news articles from 2004-2005 with 5 classes. We only use the title and headline of each article. We remove the stop words and the words whose document frequency less than 3 in the dataset <sup>4</sup>.
- **Biomedical:** It consists of 20000 paper titles from 20 different MeSH in BioASQ’s official website. We convert letters into lower case and remove the words whose document frequency less than 3 in the dataset. After preprocessing, there are 19989 documents with 20 classes <sup>5</sup>.

Table 1: Statistics on the two datasets. Label: the number of ground truth labels or categories; Docs: the total number of documents; Words: average number of words per document.

Dataset	Label	Docs	Words	Vocabulary
20Newsgroups	20	18775	135	60698
Web Snippet	8	12265	10.72	5581
BBC	5	2225	11.97	2453
Biomedical	20	19989	7.95	6887

Statistics on the four datasets after preprocessing is reported in Table 1.

We compare our model with five topic models:

- **LDA** (Blei et al., 2003). A classical probabilistic topic model. We use the open source LDA implemented by collapsed Gibbs sampling <sup>6</sup>. We use the default settings with iteration number  $n = 2000$ , the Dirichlet parameter for distribution over topics  $\alpha = 0.1$  and the Dirichlet parameter for distribution over words  $\eta = 0.01$ .
- **STC** (Zhu & Xing, 2011). A sparsity-enhanced topic model which has been proven to perform better than many existing models. We adopt the implementation of STC released by its authors <sup>7</sup>. We set the regularization constants as  $\lambda = 0.2$ ,  $\rho = 0.001$  and the maximum number of iterations of hierarchical sparse coding, dictionary learning as 100.
- **NTM** (Cao et al., 2015). A recently proposed neural network based topic model, which has been reported to outperform the Replicated Softmax model <sup>8</sup>. In NTM, the learn-

<sup>2</sup><http://www.qwone.com/jason/20Newsgroups/>

<sup>3</sup><http://jwebpro.sourceforge.net/data-web-snippets.tar.gz>

<sup>4</sup><http://mlg.ucd.ie/datasets/bbc.html>

<sup>5</sup><http://participants-area.bioasq.org/>

<sup>6</sup><https://pypi.python.org/pypi/lda>

<sup>7</sup><http://bigml.cs.tsinghua.edu.cn/jun/stc.shtml/>

<sup>8</sup><https://github.com/elbamos/NeuralTopicModels>

ing rate is 0.01 and the regularization factor is 0.001. During the pre-training procedure for all weight matrices, they are initialized with a uniform distribution in interval  $[-4\sqrt{6/(n_{\text{visible}}+n_{\text{hidden}})}, 4\sqrt{6/(n_{\text{visible}}+n_{\text{hidden}})}]$ , where  $n_{\text{visible}}=784$  and  $n_{\text{hidden}}=500$ .

- **DocNADE** (Larochelle & Lauly, 2012). An unsupervised neural network topic model of documents and have shown that it is a competitive model both as a generative model and as a document representation learning algorithm <sup>9</sup>. In DocNADE, we choose the sigmoid activate function, the hidden size is 50, the learning rate is 0.01, the bath size is 64 and the max training number is 1000.
- **GaussianLDA** (Das et al., 2015). A new technique for topic modeling by treating the document as a collection of word embeddings and topics itself as multivariate Gaussian distributions in the embedding space <sup>10</sup>. We use default values for the parameters.

Our model is implemented in Python via TensorFlow. For four datasets, we utilize the pre-trained 300-dimensional word embeddings from Wikipedia by GloVe, which is fixed during training. For each out-of-vocabulary word, we sample a random vector from a normal distribution in interval  $[0, 1]$ . We adopted ADAM optimizer for weight updating with an initial learning rate of  $4e - 4$  for four dataset. All weight matrices are initialized with a uniform distribution in interval  $[0, 1e - 5]$ . In practice, we found that our model is stable with the size of hidden layer, and set it to 500.

#### 4.2 CLASSIFICATION ACCURACY

To evaluate the effectiveness of the representation of documents learned by NVSTM, we perform text classification tasks on web snippet, 20NG, BBC and Biomedical using the document codes learned by topic models as the feature representation in a multi-class SVM. For each method, after obtaining the document representations of the training and test sets, we trained a classifier on the training set using the scikit-learn library. We then evaluated its predictive performance on the test set. On web snippet, we utilize 80% documents for training and 20% for testing. On the 20NG dataset, we keep 60% documents for training and 40% for testing, which is the same configuration as in (Lin et al., 2014). For BBC and Biomedical dataset, we also keep 60% documents for training and 40% for testing. Table 2 and Table 3 report the classification accuracy under different methods with different settings on the number of topics among the four datasets. It clearly denotes that 1) In the four datasets, the NVSTM yields the highest accuracy. 2) In general, the neural network based NVSTM, NTM, DocNADE and GLDA generate better document representations than STC and LDA, demonstrating the representative advantage of neural networks in distributed word representations. 3) Sparse models NVSTM are superior to non-sparse models (DocNADE, NTM, GLDA and LDA) separately. It indicates that sparse topic models are more capable to extract topics from short documents.

Table 2: Classification accuracy of different models on Web snippet and 20NG, with different number of topic K settings.

Dataset	Snippet					20NG				
	50	75	100	125	150	50	100	150	200	250
LDA	0.682	0.615	0.592	0.583	0.573	0.545	0.615	0.607	0.613	0.623
STC	0.678	0.686	0.699	0.724	0.701	0.602	0.631	0.647	0.652	0.654
NTM	0.660	0.667	0.723	0.732	0.747	0.623	0.627	0.641	0.632	0.667
DocNADE	0.656	0.656	0.645	0.646	0.647	0.682	0.670	0.646	0.583	0.573
GLDA	0.669	0.689	0.675	0.670	0.623	0.367	0.438	0.465	0.496	0.526
NVSTM	0.742	0.808	0.799	0.805	0.818	0.654	0.671	0.672	0.683	0.691

#### 4.3 CHARACTERISTICS OF CODE REPRESENTATION

In this part, we quantitatively investigate the word codes and documents codes learned by our model.

<sup>9</sup><https://github.com/huashiyiqike/TMBP/tree/master/DocNADE>

<sup>10</sup>[https://github.com/rajarshd/Gaussian\\_LDA](https://github.com/rajarshd/Gaussian_LDA)

Table 3: Classification accuracy of different models on BBC and Biomedical, with different number of topic K settings.

Dataset	BBC					Biomedical				
	k	20	30	40	50	60	50	100	150	200
LDA	0.784	0.774	0.796	0.762	0.758	0.536	0.534	0.547	0.534	0.541
STC	0.602	0.593	0.599	0.634	0.604	0.351	0.405	0.439	0.464	0.494
NTM	0.639	0.727	0.710	0.699	0.654	0.533	0.545	0.573	0.627	0.657
DocNADE	0.793	0.839	0.832	0.834	0.819	0.597	0.588	0.588	0.583	0.582
GLDA	0.609	0.566	0.573	0.564	0.567	0.482	0.515	0.497	0.483	0.513
NVSTM	0.783	0.835	0.833	0.836	0.813	0.567	0.623	0.645	0.671	0.664

**Word code:** We compute the average word code as  $\bar{s}_n = \frac{1}{D_n} \sum_{d \in D_n} s_{d,n}$  over all documents that word  $n$  appears in. Table 4 shows the average word codes of some representative words learned by NVSTM and LDA in 8 categories of web snippet. For each category, we also present the topics learned by NVSTM in Table 5. We list top-9 words according to their probabilities under each topic. In Table 4, the results illustrate that the codes discovered by NVSTM are apparently much sparser than those discovered by LDA. It tends to focus on narrow spectrum of topics and obtains discriminative and sparse representations of word. In contrast, LDA generates word codes with many non-zeros due to the data sparsity, leading to a confused topic distribution. Besides, in NVSTM, it is clear that each non-zero element in the word codes represents the topical meaning of words in corresponding position. The weights of these elements express their relationship with the topics. Noticed that there are words (e.g. candidates) have only a small range of topical meanings, indicating a narrow usage of those terms. While other words (e.g. hockey and marketing) tend to have a broad spectrum of topical meanings, denoting a general usage of those terms.

**Document code:** Here, each document code is calculated as  $\theta_d = \frac{\sum_{n=1}^{N_d} s_{d,nk} \beta_{kn}}{\sum_{n=1}^{N_d} \sum_{k=1}^K s_{d,nk} \beta_{kn}}$  (Bai et al., 2013). To demonstrate the quality of the learned representations by our model, we produce a t-SNE projection with for the document codes of the four datasets learned by our model in Figure 3. For Web Snippet, we sample 10% of the whole document codes. For 20newsgroups and Biomedical, we sample 30% of the whole document codes. As for BBC, we present the whole document codes. It is obvious to see that all documents are clustered into distinct categories, which is equal to the ground truth number of categories in the four datasets. It proves the semantic effectiveness of the documents codes learned by our model.

## 5 CONCLUSION

We propose a neural sparsity-enhanced topic model NVSTM, which is the first effort in introducing effective VAE inference algorithm to STC as far as we know. We take advantage of VAE to simplify the inference process, which require no model-specific algorithm derivations. With the employing of word embeddings and neural network framework, NVSTM is able to generate clearer and semantic-enriched representations for short texts. The evaluation results demonstrate the effectiveness and efficiency of our model. Future work can include extending our model with other deep generative models, such as generative adversarial network (GAN).

## REFERENCES

- Cedric Archambeau, Balaji Lakshminarayanan, and Guillaume Bouchard. Latent ibp compound dirichlet allocation. *IEEE transactions on pattern analysis and machine intelligence*, 37(2):321–333, 2015.
- Lu Bai, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. Group sparse topical coding: from code to topic. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pp. 315–324. ACM, 2013.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.

Table 4: The word codes of representative words for different categories discovered by NVSTM and LDA.

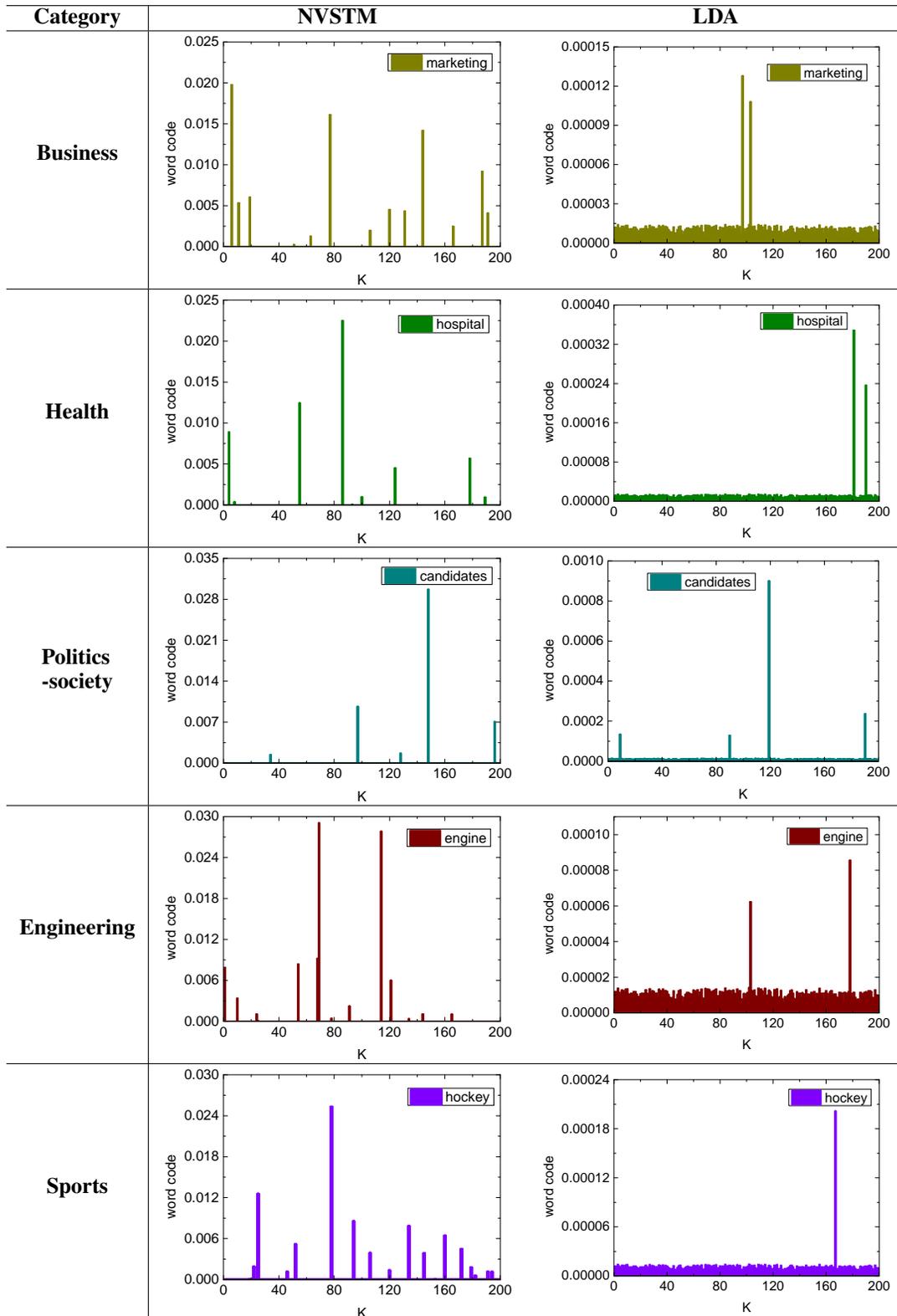


Table 5: The topics discovered by NVSTM.

Category	Topic
<b>Business</b>	T6: marketing parascope development business sustainable partnerships movieactors developing partnership T63: finance loans equity loan mortgage financing banking investment mortgages T67: investing ratneshwar investments investment investors invest equity niddk income T133: products source product quality premium csail content manufacture socialsciences T144: development sciserv ecommerce developing innovation developers business marketing projects T176: trade trading markets commodities commodity stocks market parascope currencies
<b>Computers</b>	T38: processor microprocessor processors llnl signonsandiego cpu microprocessors intel cores T108: memory laptop computer computers processor nutritionsource laptops intel disk T112: firefox mozilla netscape macintosh linux windows adobe verizon zdnet T118: systems system control security controls remote automatic monitoring automation T121: msn yahoo firefox aol gmail java algorithm algorithms signonsandiego T159: quantum computing space nasa cpu computational computers astrophysics physics
<b>culture -arts- entertainment</b>	T3: ocos parascope space socialsciences living world academyawards planet intradoc T5: film films indie filmmaker filmmakers movie comedy screening filmmaking T10: sound audio voice acoustic recordings recording listening bass song T16: photography poetry poems prose poet writing getthejob poem photographer T58: sculptor painter artist sculpture sculptures paintings artists artwork surrealist T177: art sculpture socialsciences sculptures painter paintings sculptor painting pcguide
<b>education - science</b>	T41: mathematics physics maths professors students undergraduate science teachers acidod T59: undergraduate degree student undergraduates faculty students acts particles mathematics T82: teaching school mathforum english teacher mathematics education college schools T102: lecture book lectures papers essay journal seminar conference books T109: topics mathforum essays lectures articles journals emedicinehealth literature syllabus T147: science scientific research journals published theories sciences publications articles
<b>Engineering</b>	T7: machine machines software printer llnl pcguide converter freeware kurose T69: engine engines fuel diesel cylinder gearbox piston motor petrol T90: factory inc steel searchsmb chrome ford socialsciences wiewless ltd T100: device cancertopics devices modem cable died wireless semiconductor connection T114: gasoline diesel fuel petrol engines engine emissions gas combustion T141: salary cancertopics engineer nutritionsource engineering engineers jobs job talent
<b>Health</b>	T55: therapists medical physicians therapy nurses pediatric therapist clinics doctors T56: calories nutritional vitamins calorie diet carbohydrates fats nutrition foods T86: hospital webobjects home nutritionsource clinic homes nursing emedicinehealth center T124: disease cancer lung flu cancers infection influenza arthritis infections T142: vitamins foods herbal diet alcohol supplements vitamin oils nutritional T156: drugs drug medications medication tobacco smoking alcohol prescription cancer
<b>Politics -society</b>	T27: culture democracy capitalism politics socialism society socialist ideology democratic T103: secretary party fhfb highfat parliamentary managementhelp parties elections election T128: candidates candidate election elections presidential ballot electoral vote voting T148: scandal election president senator minister senate presidency elections chairman T155: participatory debates debate democracy activism democratic pluralism grassroots encourage T168: party parties election elections electoral parliamentary parliament presidential political
<b>Engineering</b>	T7: machine machines software printer llnl pcguide converter freeware kurose T69: engine engines fuel diesel cylinder gearbox piston motor petrol T90: factory inc steel searchsmb chrome ford socialsciences wiewless ltd T100: device cancertopics devices modem cable died wireless semiconductor connection T114: gasoline diesel fuel petrol engines engine emissions gas combustion T141: salary cancertopics engineer nutritionsource engineering engineers jobs job talent
<b>Sports</b>	T72: player socialsciences players essortment game games straight button bottom T79: hockey basketball football soccer volleyball sports baseball tennis sport T95: resort village resorts town mountain peaceful hotel valley coastal T119: tennis volleyball emnlp swimming tournaments cricket hockey softball skiing T125: poker tournaments games tournament game betting gaming competitions football T127: football league soccer rugby championship basketball championships playoffs leagues

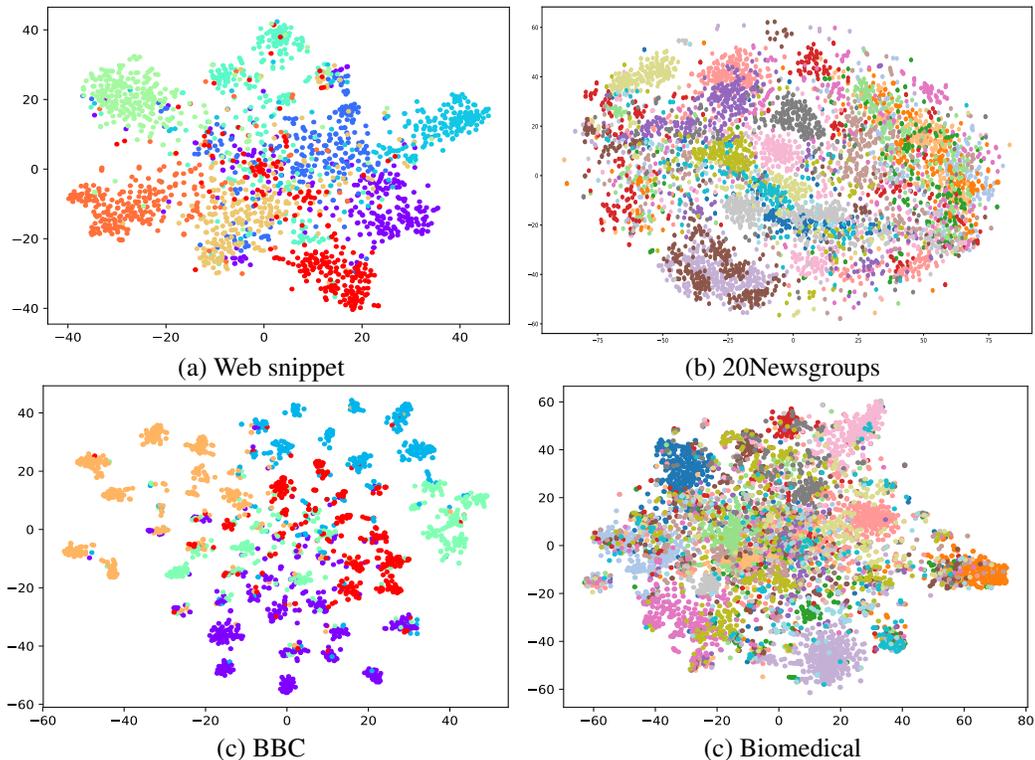


Figure 3: t-SNE projection of the estimated document codes from Web snippet, 20Newsgroups, BBC and Biomedical. The vectors are learned by the NVSTM model with 200 topics, each color represents one category from the different categories of the dataset.

Ziqiang Cao, Sujian Li, Yang Liu, Wenjie Li, and Heng Ji. A novel neural topic model and its supervised extension. In *AAAI*, pp. 2210–2216, 2015.

Dallas Card, Chenhao Tan, and Noah A Smith. A neural framework for generalized topic models. *arXiv preprint arXiv:1705.09296*, 2017.

Rajarshi Das, Manzil Zaheer, and Chris Dyer. Gaussian lda for topic models with word embeddings. In *ACL (1)*, pp. 795–804, 2015.

Matthias Heiler and Christoph Schnörr. Learning sparse representations by non-negative matrix factorization and sequential cone programming. *Journal of Machine Learning Research*, 7(Jul): 1385–1407, 2006.

Weihua Hu and Junichi Tsujii. A latent concept topic model for robust topic inference using word embeddings. In *The 54th Annual Meeting of the Association for Computational Linguistics*, pp. 380, 2016.

Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Hugo Larochelle and Stanislas Lauly. A neural autoregressive topic model. In *Advances in Neural Information Processing Systems*, pp. 2708–2716, 2012.

Tianyi Lin, Wentao Tian, Qiaozhu Mei, and Hong Cheng. The dual-sparse topic model: mining focused topics and focused terms in short text. In *Proceedings of the 23rd international conference on World wide web*, pp. 539–550. ACM, 2014.

Andre Martins and Ramon Astudillo. From softmax to sparsemax: A sparse model of attention and multi-label classification. In *International Conference on Machine Learning*, pp. 1614–1623, 2016.

- Jon D Mcauliffe and David M Blei. Supervised topic models. In *Advances in neural information processing systems*, pp. 121–128, 2008.
- Yishu Miao, Edward Grefenstette, and Phil Blunsom. Discovering discrete latent topics with neural variational inference. *arXiv preprint arXiv:1706.00359*, 2017.
- Andriy Mnih and Karol Gregor. Neural variational inference and learning in belief networks. *arXiv preprint arXiv:1402.0030*, 2014.
- Dat Quoc Nguyen, Richard Billingsley, Lan Du, and Mark Johnson. Improving topic models with latent feature word representations. *Transactions of the Association for Computational Linguistics*, 3:299–313, 2015.
- Min Peng, Binlong Gao, Jiahui Zhu, Jiajia Huang, Mengting Yuan, and Fei Li. High quality information extraction and query-oriented summarization for automatic query-reply in social network. *Expert Systems with Applications*, 44:92–101, 2016.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. Earthquake shakes twitter users: real-time event detection by social sensors. In *Proceedings of the 19th international conference on World wide web*, pp. 851–860. ACM, 2010.
- Kentaro Sasaki, Tomohiro Yoshikawa, and Takeshi Furuhashi. Online topic model for twitter considering dynamics of user interests and topic trends. In *EMNLP*, pp. 1977–1985, 2014.
- Akash Srivastava and Charles Sutton. Neural variational inference for topic models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- Sinead Williamson, Chong Wang, Katherine A Heller, and David M Blei. The ibp compound dirichlet process and its application to focused topic modeling. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 1151–1158, 2010.
- Guangxu Xun, Yaliang Li, Wayne Xin Zhao, Jing Gao, and Aidong Zhang. A correlated topic model using word embeddings. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. [doi: 10.24963/ijcai.2017/588], 2017.
- Jun Zhu and Eric P. Xing. Sparse topical coding. *CoRR*, abs/1202.3778, 2011.