

Generative Topic Embedding: a Continuous Representation of Documents

Shaohua Li^{1,2} Tat-Seng Chua¹ Jun Zhu³ Chunyan Miao²
shaohua@gmail.com dcscts@nus.edu.sg dcszj@tsinghua.edu.cn ascymiao@ntu.edu.sg

1. School of Computing, National University of Singapore

2. Joint NTU-UBC Research Centre of Excellence in Active Living for the Elderly (LILY)

3. Department of Computer Science and Technology, Tsinghua University

Abstract

Word embedding maps words into a low-dimensional continuous embedding space by exploiting the local word collocation patterns in a small context window. On the other hand, topic modeling maps documents onto a low-dimensional topic space, by utilizing the global word collocation patterns in the same document. These two types of patterns are complementary. In this paper, we propose a generative topic embedding model to combine the two types of patterns. In our model, topics are represented by embedding vectors, and are shared across documents. The probability of each word is influenced by both its local context and its topic. A variational inference method yields the topic embeddings as well as the topic mixing proportions for each document. Jointly they represent the document in a low-dimensional continuous space. In two document classification tasks, our method performs better than eight existing methods, with fewer features. In addition, we illustrate with an example that our method can generate coherent topics even based on only one document.

1 Introduction

Representing documents as fixed-length feature vectors is important for many document processing algorithms. Traditionally documents are represented as a bag-of-words (BOW) vectors. However, this simple representation suffers from being high-dimensional and highly sparse, and loses semantic relatedness across the vector dimensions.

Word Embedding methods have been demonstrated to be an effective way to represent words

as continuous vectors in a low-dimensional embedding space (Bengio et al., 2003; Mikolov et al., 2013; Pennington et al., 2014; Levy et al., 2015). The learned embedding for a word encodes its semantic/syntactic relatedness with other words, by utilizing local word collocation patterns. In each method, one core component is the *embedding link function*, which predicts a word's distribution given its context words, parameterized by their embeddings.

When it comes to documents, we wish to find a method to encode their overall semantics. Given the embeddings of each word in a document, we can imagine the document as a “bag-of-vectors”. Related words in the document point in similar directions, forming *semantic clusters*. The centroid of a semantic cluster corresponds to the most representative embedding of this cluster of words, referred to as the *semantic centroids*. We could use these semantic centroids and the number of words around them to represent a document.

In addition, for a set of documents in a particular domain, some semantic clusters may appear in many documents. By learning collocation patterns across the documents, the derived semantic centroids could be more topical and less noisy.

Topic Models, represented by Latent Dirichlet Allocation (LDA) (Blei et al., 2003), are able to group words into topics according to their collocation patterns across documents. When the corpus is large enough, such patterns reflect their semantic relatedness, hence topic models can discover coherent topics. The probability of a word is governed by its latent topic, which is modeled as a categorical distribution in LDA. Typically, only a small number of topics are present in each document, and only a small number of words have high probability in each topic. This intuition motivated Blei et al. (2003) to regularize the topic distributions with Dirichlet priors.

Semantic centroids have the same nature as topics in LDA, except that the former exist in the embedding space. This similarity drives us to seek the common semantic centroids with a model similar to LDA. We extend a generative word embedding model PSDVec (Li et al., 2015), by incorporating topics into it. The new model is named TopicVec. In TopicVec, an embedding link function models the word distribution in a topic, in place of the categorical distribution in LDA. The advantage of the link function is that the semantic relatedness is already encoded as the *cosine* distance in the embedding space. Similar to LDA, we regularize the topic distributions with Dirichlet priors. A variational inference algorithm is derived. The learning process derives *topic embeddings* in the same embedding space of words. These topic embeddings aim to approximate the underlying semantic centroids.

To evaluate how well TopicVec represents documents, we performed two document classification tasks against eight existing topic modeling or document representation methods. Two setups of TopicVec outperformed all other methods on two tasks, respectively, with fewer features. In addition, we demonstrate that TopicVec can derive coherent topics based only on *one* document, which is not possible for topic models.

The source code of our implementation is available at <https://github.com/askerlee/topicvec>.

2 Related Work

Li et al. (2015) proposed a generative word embedding method PSDVec, which is the precursor of TopicVec. PSDVec assumes that the conditional distribution of a word given its context words can be factorized approximately into independent log-bilinear terms. In addition, the word embeddings and regression residuals are regularized by Gaussian priors, reducing their chance of overfitting. The model inference is approached by an efficient Eigendecomposition and blockwise-regression method (Li et al., 2016b). TopicVec differs from PSDVec in that in the conditional distribution of a word, it is not only influenced by its context words, but also by a topic, which is an embedding vector indexed by a latent variable drawn from a Dirichlet-Multinomial distribution.

Hinton and Salakhutdinov (2009) proposed to model topics as a certain number of binary hidden variables, which interact with all words in the doc-

ument through weighted connections. Larochelle and Lauly (2012) assigned each word a unique topic vector, which is a summarization of the context of the current word.

Huang et al. (2012) proposed to incorporate global (document-level) semantic information to help the learning of word embeddings. The global embedding is simply a weighted average of the embeddings of words in the document.

Le and Mikolov (2014) proposed Paragraph Vector. It assumes each piece of text has a latent paragraph vector, which influences the distributions of all words in this text, in the same way as a latent word. It can be viewed as a special case of TopicVec, with the topic number set to 1. Typically, however, a document consists of multiple semantic centroids, and the limitation of only one topic may lead to underfitting.

Nguyen et al. (2015) proposed Latent Feature Topic Modeling (LFTM), which extends LDA to incorporate word embeddings as latent features. The topic is modeled as a mixture of the conventional categorical distribution and an embedding link function. The coupling between these two components makes the inference difficult. They designed a Gibbs sampler for model inference. Their implementation¹ is slow and infeasible when applied to a large corpus.

Liu et al. (2015) proposed Topical Word Embedding (TWE), which combines word embedding with LDA in a simple and effective way. They train word embeddings and a topic model separately on the same corpus, and then average the embeddings of words in the same topic to get the embedding of this topic. The topic embedding is concatenated with the word embedding to form the topical word embedding of a word. In the end, the topical word embeddings of all words in a document are averaged to be the embedding of the document. This method performs well on our two classification tasks. Weaknesses of TWE include: 1) the way to combine the results of word embedding and LDA lacks statistical foundations; 2) the LDA module requires a large corpus to derive semantically coherent topics.

Das et al. (2015) proposed Gaussian LDA. It uses pre-trained word embeddings. It assumes that words in a topic are random samples from a multivariate Gaussian distribution with the topic embedding as the mean. Hence the probability that a

¹<https://github.com/datquocnguyen/LFTM/>

Name	Description
\mathcal{S}	Vocabulary $\{s_1, \dots, s_W\}$
\mathbf{V}	Embedding matrix $(\mathbf{v}_{s_1}, \dots, \mathbf{v}_{s_W})$
\mathcal{D}	Document set $\{d_1, \dots, d_M\}$
\mathbf{v}_{s_i}	Embedding of word s_i
$a_{s_i s_j}, \mathbf{A}$	Bigram residuals
$\mathbf{t}_{ik}, \mathbf{T}_i$	Topic embeddings in doc d_i
r_{ik}, \mathbf{r}_i	Topic residuals in doc d_i
z_{ij}	Topic assignment of the j -th word j in doc d_i
ϕ_i	Mixing proportions of topics in doc d_i

Table 1: Table of notations

word belongs to a topic is determined by the Euclidean distance between the word embedding and the topic embedding. This assumption might be improper as the Euclidean distance is not an optimal measure of semantic relatedness between two embeddings².

3 Notations and Definitions

Throughout this paper, we use uppercase bold letters such as \mathcal{S}, \mathbf{V} to denote a matrix or set, lowercase bold letters such as \mathbf{v}_{w_i} to denote a vector, a normal uppercase letter such as N, W to denote a scalar constant, and a normal lowercase letter as s_i, w_i to denote a scalar variable.

Table 1 lists the notations in this paper.

In a document, a sequence of words is referred to as a *text window*, denoted by w_i, \dots, w_{i+l} , or $w_i:w_{i+l}$. A text window of chosen size c before a word w_i defines the *context* of w_i as w_{i-c}, \dots, w_{i-1} . Here w_i is referred to as the *focus word*. Each context word w_{i-j} and the focus word w_i comprise a bigram w_{i-j}, w_i .

We assume each word in a document is semantically similar to a *topic embedding*. Topic embeddings reside in the same N -dimensional space as word embeddings. When it is clear from context, topic embeddings are often referred to as *topics*. Each document has K candidate topics, arranged in the matrix form $\mathbf{T}_i = (\mathbf{t}_{i1} \dots \mathbf{t}_{iK})$, referred to as the *topic matrix*. Specifically, we fix $\mathbf{t}_{i1} = \mathbf{0}$, referring to it as the *null topic*.

In a document d_i , each word w_{ij} is assigned to a topic indexed by $z_{ij} \in \{1, \dots, K\}$. Geometrically this means the embedding $\mathbf{v}_{w_{ij}}$ tends to align

²Almost all modern word embedding methods adopt the exponentiated cosine similarity as the link function, hence the cosine similarity may be assumed to be a better estimate of the semantic relatedness between embeddings derived from these methods.

with the direction of $\mathbf{t}_{i,z_{ij}}$. Each topic \mathbf{t}_{ik} has a document-specific prior probability to be assigned to a word, denoted as $\phi_{ik} = P(k|d_i)$. The vector $\phi_i = (\phi_{i1}, \dots, \phi_{iK})$ is referred to as the *mixing proportions* of these topics in document d_i .

4 Link Function of Topic Embedding

In this section, we formulate the distribution of a word given its context words and topic, in the form of a link function.

The core of most word embedding methods is a *link function* that connects the embeddings of a focus word and its context words, to define the distribution of the focus word. Li et al. (2015) proposed the following link function:

$$P(w_c | w_0 : w_{c-1}) \approx P(w_c) \exp \left(\sum_{l=0}^{c-1} \mathbf{v}_{w_c}^\top \mathbf{v}_{w_l} + \sum_{l=0}^{c-1} a_{w_l w_c} \right) \quad (1)$$

Here $a_{w_l w_c}$ is referred as the bigram residual, indicating the non-linear part not captured by $\mathbf{v}_{w_c}^\top \mathbf{v}_{w_l}$. It is essentially the logarithm of the normalizing constant of a softmax term. Some literature, e.g. (Pennington et al., 2014), refers to such a term as a bias term.

(1) is based on the assumption that the conditional distribution $P(w_c | w_0 : w_{c-1})$ can be factorized approximately into independent log-bilinear terms, each corresponding to a context word. This approximation leads to an efficient and effective word embedding algorithm *PSDVec* (Li et al., 2015). We follow this assumption, and propose to incorporate the topic of w_c in a way like a latent word. In particular, in addition to the context words, the corresponding embedding \mathbf{t}_{ik} is included as a new log-bilinear term that influences the distribution of w_c . Hence we obtain the following extended link function:

$$P(w_c | w_0:w_{c-1}, z_c, d_i) \approx P(w_c) \cdot \exp \left(\sum_{l=0}^{c-1} \mathbf{v}_{w_c}^\top \mathbf{v}_{w_l} + \mathbf{t}_{z_c} + \sum_{l=0}^{c-1} a_{w_l w_c} + r_{z_c} \right), \quad (2)$$

where d_i is the current document, and r_{z_c} is the logarithm of the normalizing constant, named the *topic residual*. Note that the topic embeddings \mathbf{t}_{z_c} may be specific to d_i . For simplicity of notation, we drop the document index in \mathbf{t}_{z_c} . To restrict the impact of topics and avoid overfitting, we constrain the magnitudes of all topic embeddings, so that they are always within a hyperball of radius γ .

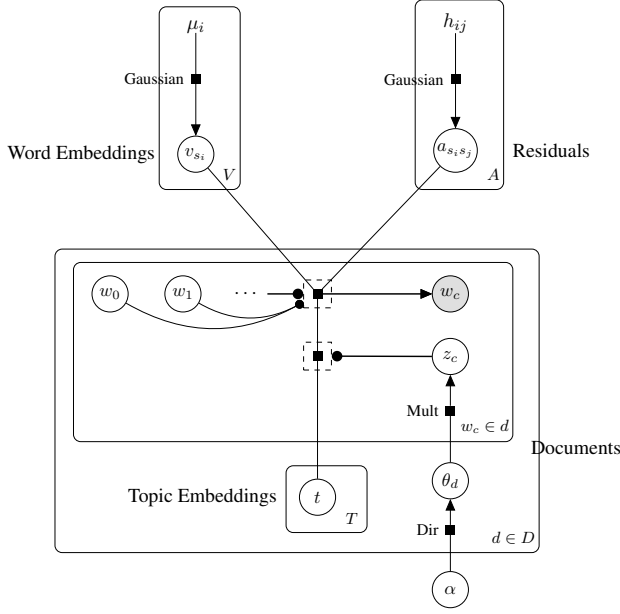


Figure 1: Graphical representation of TopicVec.

It is infeasible to compute the exact value of the topic residual r_k . We approximate it by the context size $c = 0$. Then (2) becomes:

$$P(w_c | k, d_i) = P(w_c) \exp \mathbf{v}_{w_c}^\top \mathbf{t}_k + r_k \quad (3)$$

It is required that $\sum_{w_c \in \mathcal{S}} P(w_c | k) = 1$ to make (3) a distribution. It follows that

$$r_k = -\log \sum_{s_j \in \mathcal{S}} P(s_j) \exp \{ \mathbf{v}_{s_j}^\top \mathbf{t}_k \} \quad (4)$$

(4) can be expressed in the matrix form:

$$\mathbf{r} = -\log(\mathbf{u} \exp \{ \mathbf{V}^\top \mathbf{T} \}), \quad (5)$$

where \mathbf{u} is the row vector of unigram probabilities.

5 Generative Process and Likelihood

The generative process of words in documents can be regarded as a hybrid of LDA and PSDVec. Analogous to PSDVec, the word embedding \mathbf{v}_{s_i} and residual $a_{s_i s_j}$ are drawn from respective Gaussians. For the sake of clarity, we ignore their generation steps, and focus on the topic embeddings. The remaining generative process is as follows:

1. For the k -th topic, draw a topic embedding uniformly from a hyperball of radius γ , i.e. $\mathbf{t}_k \sim \text{Unif}(B_\gamma)$;
2. For each document d_i :
 - (a) Draw the mixing proportions ϕ_i from the Dirichlet prior $\text{Dir}(\alpha)$;

(b) For the j -th word:

- i. Draw topic assignment z_{ij} from the categorical distribution $\text{Cat}(\phi_i)$;
- ii. Draw word w_{ij} from \mathcal{S} according to $P(w_{ij} | w_{i,j-c}:w_{i,j-1}, z_{ij}, d_i)$.

The above generative process is presented in plate notation in Figure (1).

5.1 Likelihood Function

Given the embeddings \mathbf{V} , the bigram residuals \mathbf{A} , the topics \mathbf{T}_i and the hyperparameter α , the complete-data likelihood of a single document d_i is:

$$\begin{aligned} & p(d_i, \mathbf{Z}_i, \phi_i | \alpha, \mathbf{V}, \mathbf{A}, \mathbf{T}_i) \\ &= p(\phi_i | \alpha) p(\mathbf{Z}_i | \phi_i) p(d_i | \mathbf{V}, \mathbf{A}, \mathbf{T}_i, \mathbf{Z}_i) \\ &= \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{j=1}^K \phi_{ij}^{\alpha_j - 1} \cdot \prod_{j=1}^{L_i} \phi_{i,z_{ij}} P(w_{ij}) \\ & \quad \cdot \exp \sum_{l=j-c}^{j-1} \mathbf{v}_{w_{il}}^\top \mathbf{v}_{w_{ij}} + \mathbf{t}_{z_{ij}} \\ & \quad + \sum_{l=j-c}^{j-1} a_{w_{il}w_{ij}} + r_{i,z_{ij}}, \end{aligned} \quad (6)$$

where $\mathbf{Z}_i = (z_{i1}, \dots, z_{iL_i})$, and $\Gamma(\cdot)$ is the Gamma function.

Let $\mathbf{Z}, \mathbf{T}, \phi$ denote the collection of all the document-specific $\{\mathbf{Z}_i\}_{i=1}^M, \{\mathbf{T}_i\}_{i=1}^M, \{\phi_i\}_{i=1}^M$, respectively. Then the complete-data likelihood of the whole corpus is:

$$\begin{aligned} & p(\mathbf{D}, \mathbf{A}, \mathbf{V}, \mathbf{Z}, \mathbf{T}, \phi | \alpha, \gamma, \mu) \\ &= \prod_{i=1}^M P(\mathbf{v}_{s_i}; \mu_i) \prod_{i,j=1}^{W,W} P(a_{s_i s_j}; f(h_{ij})) \prod_k \text{Unif}(B_\gamma) \\ & \quad \cdot \prod_{i=1}^M \{ p(\phi_i | \alpha) p(\mathbf{Z}_i | \phi_i) p(d_i | \mathbf{V}, \mathbf{A}, \mathbf{T}_i, \mathbf{Z}_i) \} \\ &= \frac{1}{\mathcal{Z}(\mathbf{H}, \mu) U_\gamma^K} \exp \left\{ - \sum_{i,j=1}^{W,W} f(h_{i,j}) a_{s_i s_j}^2 - \sum_{i=1}^M \mu_i \|\mathbf{v}_{s_i}\|^2 \right\} \\ & \quad \cdot \prod_{i=1}^M \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{j=1}^K \phi_{ij}^{\alpha_j - 1} \cdot \prod_{j=1}^{L_i} \phi_{i,z_{ij}} P(w_{ij}) \\ & \quad \cdot \exp \sum_{l=j-c}^{j-1} \mathbf{v}_{w_{il}}^\top \mathbf{v}_{w_{ij}} + \mathbf{t}_{z_{ij}} + \sum_{l=j-c}^{j-1} a_{w_{il}w_{ij}} + r_{i,z_{ij}}, \end{aligned} \quad (7)$$

where $P(\mathbf{v}_{s_i}; \mu_i)$ and $P(a_{s_i s_j}; f(h_{ij}))$ are the two Gaussian priors as defined in (Li et al., 2015).

Following the convention in (Li et al., 2015), h_{ij} , \mathbf{H} are empirical bigram probabilities, $\boldsymbol{\mu}$ are the embedding magnitude penalty coefficients, and $\mathcal{Z}(\mathbf{H}, \boldsymbol{\mu})$ is the normalizing constant for word embeddings. U_γ is the volume of the hyperball of radius γ .

Taking the logarithm of both sides, we obtain

$$\begin{aligned} & \log p(\mathbf{D}, \mathbf{A}, \mathbf{V}, \mathbf{Z}, \mathbf{T}, \phi | \boldsymbol{\alpha}, \gamma, \boldsymbol{\mu}) \\ = & C_0 - \log \mathcal{Z}(\mathbf{H}, \boldsymbol{\mu}) - \|\mathbf{A}\|_{f(\mathbf{H})}^2 - \sum_{i=1}^W \mu_i \|\mathbf{v}_{s_i}\|^2 \\ & + \sum_{i=1}^M \sum_{k=1}^K \log \phi_{ik} (m_{ik} + \alpha_k - 1) + \sum_{i=1}^{L_i} r_{i, z_{ij}} \\ & + \mathbf{v}_{w_{ij}}^\top \mathbf{v}_{w_{il}} + \mathbf{t}_{z_{ij}} + a_{w_{il}w_{ij}} \quad , \quad (8) \end{aligned}$$

where $m_{ik} = \sum_{j=1}^{L_i} \delta(z_{ij} = k)$ counts the number of words assigned with the k -th topic in d_i , $C_0 = M \log \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} + \sum_{i,j=1}^{M, L_i} \log P(w_{ij}) - K \log U_\gamma$ is constant given the hyperparameters.

6 Variational Inference Algorithm

6.1 Learning Objective and Process

Given the hyperparameters $\boldsymbol{\alpha}, \gamma, \boldsymbol{\mu}$, the learning objective is to find the embeddings \mathbf{V} , the topics \mathbf{T} , and the word-topic and document-topic distributions $p(\mathbf{Z}_i, \phi_i | d_i, \mathbf{A}, \mathbf{V}, \mathbf{T})$. Here the hyperparameters $\boldsymbol{\alpha}, \gamma, \boldsymbol{\mu}$ are kept constant, and we make them implicit in the distribution notations.

However, the coupling between \mathbf{A}, \mathbf{V} and $\mathbf{T}, \mathbf{Z}, \phi$ makes it inefficient to optimize them simultaneously. To get around this difficulty, we learn word embeddings and topic embeddings separately. Specifically, the learning process is divided into two stages:

1. In the first stage, considering that the topics have a relatively small impact to word distributions and the impact might be ‘‘averaged out’’ across different documents, we simplify the model by ignoring topics temporarily. Then the model falls back to the original PSDVec. The optimal solution $\mathbf{V}^*, \mathbf{A}^*$ is obtained accordingly;
2. In the second stage, we treat $\mathbf{V}^*, \mathbf{A}^*$ as constant, plug it into the likelihood function, and find the corresponding optimal $\mathbf{T}^*, p(\mathbf{Z}, \phi | \mathbf{D}, \mathbf{A}^*, \mathbf{V}^*, \mathbf{T}^*)$ of the full model.

As in LDA, this posterior is analytically intractable, and we use a simpler variational distribution $q(\mathbf{Z}, \phi)$ to approximate it.

6.2 Mean-Field Approximation and Variational GEM Algorithm

In this stage, we fix $\mathbf{V} = \mathbf{V}^*, \mathbf{A} = \mathbf{A}^*$, and seek the optimal $\mathbf{T}^*, p(\mathbf{Z}, \phi | \mathbf{D}, \mathbf{A}^*, \mathbf{V}^*, \mathbf{T}^*)$. As $\mathbf{V}^*, \mathbf{A}^*$ are constant, we also make them implicit in the following expressions.

For an arbitrary variational distribution $q(\mathbf{Z}, \phi)$, the following equalities hold

$$\begin{aligned} & E_q \log \frac{p(\mathbf{D}, \mathbf{Z}, \phi | \mathbf{T})}{q(\mathbf{Z}, \phi)} \\ = & E_q [\log p(\mathbf{D}, \mathbf{Z}, \phi | \mathbf{T})] + \mathcal{H}(q) \\ = & \log p(\mathbf{D} | \mathbf{T}) - \text{KL}(q || p), \quad (9) \end{aligned}$$

where $p = p(\mathbf{Z}, \phi | \mathbf{D}, \mathbf{T})$, $\mathcal{H}(q)$ is the entropy of q . This implies

$$\begin{aligned} & \text{KL}(q || p) \\ = & \log p(\mathbf{D} | \mathbf{T}) - E_q [\log p(\mathbf{D}, \mathbf{Z}, \phi | \mathbf{T})] + \mathcal{H}(q) \\ = & \log p(\mathbf{D} | \mathbf{T}) - \mathcal{L}(q, \mathbf{T}). \quad (10) \end{aligned}$$

In (10), $E_q [\log p(\mathbf{D}, \mathbf{Z}, \phi | \mathbf{T})] + \mathcal{H}(q)$ is usually referred to as the *variational free energy* $\mathcal{L}(q, \mathbf{T})$, which is a lower bound of $\log p(\mathbf{D} | \mathbf{T})$. Directly maximizing $\log p(\mathbf{D} | \mathbf{T})$ w.r.t. \mathbf{T} is intractable due to the hidden variables \mathbf{Z}, ϕ , so we maximize its lower bound $\mathcal{L}(q, \mathbf{T})$ instead. We adopt a mean-field approximation of the true posterior as the variational distribution, and use a variational algorithm to find q^*, \mathbf{T}^* maximizing $\mathcal{L}(q, \mathbf{T})$.

The following variational distribution is used:

$$\begin{aligned} & q(\mathbf{Z}, \phi; \boldsymbol{\pi}, \boldsymbol{\theta}) = q(\phi; \boldsymbol{\theta}) q(\mathbf{Z}; \boldsymbol{\pi}) \\ = & \prod_{i=1}^M \text{Dir}(\phi_i; \boldsymbol{\theta}_i) \prod_{j=1}^{L_i} \text{Cat}(z_{ij}; \boldsymbol{\pi}_{ij}) \quad . \quad (11) \end{aligned}$$

We can obtain (Li et al., 2016a)

$$\begin{aligned} & \mathcal{L}(q, \mathbf{T}) \\ = & \sum_{i=1}^M \sum_{k=1}^K \sum_{j=1}^{L_i} \pi_{ij}^k + \alpha_k - 1 \quad \psi(\theta_{ik}) - \psi(\theta_{i0}) \\ & + \text{Tr}(\mathbf{T}_i^\top \mathbf{v}_{w_{ij}} \boldsymbol{\pi}_{ij}^\top) + \mathbf{r}_i^\top \boldsymbol{\pi}_{ij} \\ & + \mathcal{H}(q) + C_1, \quad (12) \end{aligned}$$

where \mathbf{T}_i is the topic matrix of the i -th document, and \mathbf{r}_i is the vector constructed by concatenating all the topic residuals r_{ik} . $C_1 = C_0 - \log \mathcal{Z}(\mathbf{H}, \boldsymbol{\mu}) - \|\mathbf{A}\|_F^2(\mathbf{H}) - \sum_{i=1}^W \mu_i \|\mathbf{v}_{s_i}\|^2 + \sum_{i,j=1}^{M,L_i} \mathbf{v}_{w_{ij}}^\top \sum_{k=j-c}^{j-1} \mathbf{v}_{w_{ik}} + \sum_{k=j-c}^{j-1} a_{w_{ik}w_{ij}}$ is constant.

We proceed to optimize (12) with a Generalized Expectation-Maximization (GEM) algorithm w.r.t. q and \mathbf{T} as follows:

1. Initialize all the topics $\mathbf{T}_i = \mathbf{0}$, and correspondingly their residuals $\mathbf{r}_i = \mathbf{0}$;
2. Iterate over the following two steps until convergence. In the l -th step:
 - (a) Let the topics and residuals be $\mathbf{T} = \mathbf{T}^{(l-1)}$, $\mathbf{r} = \mathbf{r}^{(l-1)}$, find $q^{(l)}(\mathbf{Z}, \boldsymbol{\phi})$ that maximizes $\mathcal{L}(q, \mathbf{T}^{(l-1)})$. This is the Expectation step (E-step). In this step, $\log p(\mathbf{D}|\mathbf{T})$ is constant. Then the q that maximizes $\mathcal{L}(q, \mathbf{T}^{(l)})$ will minimize $\text{KL}(q||p)$, i.e. such a q is the closest variational distribution to p measured by KL-divergence;
 - (b) Given the variational distribution $q^{(l)}(\mathbf{Z}, \boldsymbol{\phi})$, find $\mathbf{T}^{(l)}, \mathbf{r}^{(l)}$ that improve $\mathcal{L}(q^{(l)}, \mathbf{T})$, using Gradient descent method. This is the generalized Maximization step (M-step). In this step, $\boldsymbol{\pi}, \boldsymbol{\theta}, \mathcal{H}(q)$ are constant.

6.2.1 Update Equations of $\boldsymbol{\pi}, \boldsymbol{\theta}$ in E-Step

In the E-step, $\mathbf{T} = \mathbf{T}^{(l-1)}$, $\mathbf{r} = \mathbf{r}^{(l-1)}$ are constant. Taking the derivative of $\mathcal{L}(q, \mathbf{T}^{(l-1)})$ w.r.t. π_{ij}^k and θ_{ik} , respectively, we can obtain the optimal solutions (Li et al., 2016a) at:

$$\pi_{ij}^k \propto \exp\{\psi(\theta_{ik}) + \mathbf{v}_{w_{ij}}^\top \mathbf{t}_{ik} + r_{ik}\}. \quad (13)$$

$$\theta_{ik} = \sum_{j=1}^{L_i} \pi_{ij}^k + \alpha_k. \quad (14)$$

6.2.2 Update Equation of \mathbf{T}_i in M-Step

In the Generalized M-step, $\boldsymbol{\pi} = \boldsymbol{\pi}^{(l)}$, $\boldsymbol{\theta} = \boldsymbol{\theta}^{(l)}$ are constant. For notational simplicity, we drop their superscripts (l).

To update \mathbf{T}_i , we first take the derivative of (12) w.r.t. \mathbf{T}_i , and then take the Gradient Descent method.

The derivative is obtained as (Li et al., 2016a):

$$\begin{aligned} & \frac{\partial \mathcal{L}(q^{(l)}, \mathbf{T})}{\partial \mathbf{T}_i} \\ &= \sum_{j=1}^{L_i} \mathbf{v}_{w_{ij}} \boldsymbol{\pi}_{ij}^\top + \sum_{k=1}^K \bar{m}_{ik} \frac{\partial r_{ik}}{\partial \mathbf{T}_i}, \end{aligned} \quad (15)$$

where $\bar{m}_{ik} = \sum_{j=1}^{L_i} \pi_{ij}^k = E[m_{ik}]$, the sum of the variational probabilities of each word being assigned to the k -th topic in the i -th document. $\frac{\partial r_{ik}}{\partial \mathbf{T}_i}$ is a gradient matrix, whose j -th column is $\frac{\partial r_{ik}}{\partial \mathbf{t}_{ij}}$.

Remind that $r_{ik} = -\log E_{P(s)}[\exp\{\mathbf{v}_s^\top \mathbf{t}_{ik}\}]$.

When $j \neq k$, it is easy to verify that $\frac{\partial r_{ik}}{\partial \mathbf{t}_{ij}} = \mathbf{0}$.

When $j = k$, we have

$$\begin{aligned} \frac{\partial r_{ik}}{\partial \mathbf{t}_{ik}} &= e^{-r_{ik}} \cdot E_{P(s)}[\exp\{\mathbf{v}_s^\top \mathbf{t}_{ik}\} \mathbf{v}_s] \\ &= e^{-r_{ik}} \cdot \sum_{s \in W} \exp\{\mathbf{v}_s^\top \mathbf{t}_{ik}\} P(s) \mathbf{v}_s \\ &= e^{-r_{ik}} \cdot \exp\{\mathbf{t}_{ik}^\top \mathbf{V}\} (\mathbf{u} \circ \mathbf{V}), \end{aligned} \quad (16)$$

where $\mathbf{u} \circ \mathbf{V}$ is to multiply each column of \mathbf{V} with \mathbf{u} element-by-element.

Therefore $\frac{\partial r_{ik}}{\partial \mathbf{T}_i} = (\mathbf{0}, \dots, \frac{\partial r_{ik}}{\partial \mathbf{t}_{ik}}, \dots, \mathbf{0})$. Plugging it into (15), we obtain

$$\frac{\partial \mathcal{L}(q^{(l)}, \mathbf{T})}{\partial \mathbf{T}_i} = \sum_{j=1}^{L_i} \mathbf{v}_{w_{ij}} \boldsymbol{\pi}_{ij}^\top + (\bar{m}_{i1} \frac{\partial r_{i1}}{\partial \mathbf{t}_{i1}}, \dots, \bar{m}_{iK} \frac{\partial r_{iK}}{\partial \mathbf{t}_{iK}}).$$

We proceed to optimize \mathbf{T}_i with a gradient descent method:

$$\mathbf{T}_i^{(l)} = \mathbf{T}_i^{(l-1)} + \lambda(l, L_i) \frac{\partial \mathcal{L}(q^{(l)}, \mathbf{T})}{\partial \mathbf{T}_i},$$

where $\lambda(l, L_i) = \frac{L_0 \lambda_0}{l \cdot \max\{L_i, L_0\}}$ is the learning rate function, L_0 is a pre-specified document length threshold, and λ_0 is the initial learning rate. As the magnitude of $\frac{\partial \mathcal{L}(q^{(l)}, \mathbf{T})}{\partial \mathbf{T}_i}$ is approximately proportional to the document length L_i , to avoid the step size becoming too big on a long document, if $L_i > L_0$, we normalize it by L_i .

To satisfy the constraint that $\|\mathbf{t}_{ik}^{(l)}\| \leq \gamma$, when $\mathbf{t}_{ik}^{(l)} > \gamma$, we normalize it by $\gamma / \|\mathbf{t}_{ik}^{(l)}\|$.

After we obtain the new \mathbf{T} , we update $\mathbf{r}_i^{(m)}$ using (5).

Sometimes, especially in the initial few iterations, due to the excessively big step size of the gradient descent, $\mathcal{L}(q, \mathbf{T})$ may decrease after the update of \mathbf{T} . Nonetheless the general direction of $\mathcal{L}(q, \mathbf{T})$ is increasing.

6.3 Sharing of Topics across Documents

In principle we could use one set of topics across the whole corpus, or choose different topics for different subsets of documents. One could choose a way to best utilize cross-document information.

For instance, when the document category information is available, we could make the documents in each category share their respective set

of topics, so that M categories correspond to M sets of topics. In the learning algorithm, only the update of π_{ij}^k needs to be changed to cater for this situation: when the k -th topic is relevant to the document i , we update π_{ij}^k using (13); otherwise $\pi_{ij}^k = 0$.

An identifiability problem may arise when we split topic embeddings according to document subsets. In different topic groups, some highly similar redundant topics may be learned. If we project documents into the topic space, portions of documents in the same topic in different documents may be projected onto different dimensions of the topic space, and similar documents may eventually be projected into very different topic proportion vectors. In this situation, directly using the projected topic proportion vectors could cause problems in unsupervised tasks such as clustering. A simple solution to this problem would be to compute the pairwise similarities between topic embeddings, and consider these similarities when computing the similarity between two projected topic proportion vectors. Two similar documents will then still receive a high similarity score.

7 Experimental Results

To investigate the quality of document representation of our TopicVec model, we compared its performance against eight topic modeling or document representation methods in two document classification tasks. Moreover, to show the topic coherence of TopicVec on a single document, we present the top words in top topics learned on a news article.

7.1 Document Classification Evaluation

7.1.1 Experimental Setup

Compared Methods Two setups of TopicVec were evaluated:

- **TopicVec**: the topic proportions learned by TopicVec;
- **TV+WV**: the topic proportions, concatenated with the mean word embedding of the document (same as the MeanWV below).

We compare the performance of our methods against eight methods, including three topic modeling methods, three continuous document representation methods, and the conventional bag-of-words (**BOW**) method. The count vector of BOW is unweighted.

The topic modeling methods include:

- **LDA**: the vanilla LDA (Blei et al., 2003) in the gensim library³;
- **sLDA**: Supervised Topic Model⁴ (McAuliffe and Blei, 2008), which improves the predictive performance of LDA by modeling class

The same preprocessing steps were applied to all methods: words were lowercased; stop words and words out of the word embedding vocabulary (which means that they are extremely rare) were removed.

Experimental Settings TopicVec used the word embeddings trained using PSDVec on a March 2015 Wikipedia snapshot. It contains the most frequent 180,000 words. The dimensionality of word embeddings and topic embeddings was 500. The hyperparameters were $\alpha = (0.1, \dots, 0.1)$, $\gamma = 5$. For 20news and Reuters, we specified 15 and 12 topics in each category on the training set, respectively. The first topic in each category was always set to null. The learned topic embeddings were combined to form the whole topic set, where redundant null topics in different categories were removed, leaving us with 281 topics for 20News and 111 topics for Reuters. The initial learning rate was set to 0.1. After 100 GEM iterations on each dataset, the topic embeddings were obtained. Then the posterior document-topic distributions of the test sets were derived by performing one E-step given the topic embeddings trained on the training set.

LFTM includes two models: LF-LDA and LF-DMM. We chose the better performing LF-LDA to evaluate. TWE includes three models, and we chose the best performing TWE-1 to compare.

LDA, sLDA, LFTM and TWE used the specified 50 topics on Reuters, as this is the optimal topic number according to (Lu et al., 2011). On the larger 20news dataset, they used the specified 100 topics. Other hyperparameters of all compared methods were left at their default values.

GaussianLDA was specified 100 topics on 20news and 70 topics on Reuters. As each sampling iteration took over 2 hours, we only had time for 100 sampling iterations.

For each method, after obtaining the document representations of the training and test sets, we trained an ℓ -1 regularized linear SVM one-vs-all classifier on the training set using the scikit-learn library¹¹. We then evaluated its predictive performance on the test set.

Evaluation metrics Considering that the largest few categories dominate Reuters, we adopted macro-averaged precision, recall and F1 measures as the evaluation metrics, to avoid the average results being dominated by the performance of the

	20News			Reuters		
	Prec	Rec	F1	Prec	Rec	F1
BOW	69.1	68.5	68.6	92.5	90.3	91.1
LDA	61.9	61.4	60.3	76.1	74.3	74.8
sLDA	61.4	60.9	60.9	88.3	83.3	85.1
LFTM	63.5	64.8	63.7	84.6	86.3	84.9
MeanWV	70.4	70.3	70.1	92.0	89.6	90.5
Doc2Vec	56.3	56.6	55.4	84.4	50.0	58.5
TWE	69.5	69.3	68.8	91.0	89.1	89.9
GaussianLDA	30.9	26.5	22.7	46.2	31.5	35.3
TopicVec	71.4	71.3	71.2	91.8	92.0	91.7
TV+WV ¹	72.1	71.9	71.8	91.4	91.9	91.5

¹Combined features of TopicVec topic proportions and MeanWV.

Table 2: Performance on multi-class text classification. Best score is in boldface.

Avg. Features	BOW	MeanWV	TWE	TopicVec	TV+WV
20News	50381	500	800	281	781
Reuters	17989	500	800	111	611

Table 3: Number of features of the five best performing methods.

top categories.

Evaluation Results Table 2 presents the performance of the different methods on the two classification tasks. The highest scores were highlighted with boldface. It can be seen that TV+WV and TopicVec obtained the best performance on the two tasks, respectively. With only topic proportions as features, TopicVec performed slightly better than BOW, MeanWV and TWE, and significantly outperformed four other methods. The number of features it used was much lower than BOW, MeanWV and TWE (Table 3).

GaussianLDA performed considerably inferior to all other methods. After checking the generated topic embeddings manually, we found that the embeddings for different topics are highly similar to each other. Hence the posterior topic proportions were almost uniform and non-discriminative. In addition, on the two datasets, even the fastest Alias sampling in (Das et al., 2015) took over 2 hours for one iteration and 10 days for the whole 100 iterations. In contrast, our method finished the 100 EM iterations in 2 hours.

¹¹<http://scikit-learn.org/stable/modules/svm.html>

- and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 795–804, Beijing, China, July. Association for Computational Linguistics.
- Geoffrey E Hinton and Ruslan R Salakhutdinov. 2009. Replicated softmax: an undirected topic model. In *Advances in neural information processing systems*, pages 1607–1614.
- Eric H Huang, Richard Socher, Christopher D Manning, and Andrew Y Ng. 2012. Improving word representations via global context and multiple word prototypes. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 873–882. Association for Computational Linguistics.
- Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Q. Weinberger. 2015. From word embeddings to document distances. In David Blei and Francis Bach, editors, *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 957–966. JMLR Workshop and Conference Proceedings.
- Hugo Larochelle and Stanislas Lauly. 2012. A neural autoregressive topic model. In *Advances in Neural Information Processing Systems*, pages 2708–2716.
- Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1188–1196.
- Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics*, 3:211–225.
- Shaohua Li, Jun Zhu, and Chunyan Miao. 2015. A generative word embedding model and its low rank positive semidefinite solution. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1599–1609, Lisbon, Portugal, September. Association for Computational Linguistics.
- Shaohua Li, Tat-Seng Chua, Jun Zhu, and Chunyan Miao. 2016a. Generative topic embedding: a continuous representation of documents (extended version with proofs). Technical report. <https://github.com/askerlee/topicvec/blob/master/topicvec-ext.pdf>.
- Shaohua Li, Jun Zhu, and Chunyan Miao. 2016b. PSDVec: a toolbox for incremental and scalable word embedding. *To appear in Neurocomputing*.
- Yang Liu, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. 2015. Topical word embeddings. In *AAAI*, pages 2418–2424.
- Yue Lu, Qiaozhu Mei, and ChengXiang Zhai. 2011. Investigating task performance of probabilistic topic models: an empirical study of PLSA and LDA. *Information Retrieval*, 14(2):178–203.
- Jon D McAuliffe and David M Blei. 2008. Supervised topic models. In *Advances in neural information processing systems*, pages 121–128.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of NIPS 2013*, pages 3111–3119.
- Dat Quoc Nguyen, Richard Billingsley, Lan Du, and Mark Johnson. 2015. Improving topic models with latent feature word representations. *Transactions of the Association for Computational Linguistics*, 3:299–313.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. GloVe: Global vectors for word representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12.