

Neural Topic Modeling with Bidirectional Adversarial Training

Rui Wang[†] Xuemeng Hu[†] Deyu Zhou^{†*} Yulan He[§]
Yuxuan Xiong[†] Chenchen Ye[†] Haiyang Xu[‡]

[†]School of Computer Science and Engineering, Key Laboratory of Computer Network and Information Integration, Ministry of Education, Southeast University, China

[§]Department of Computer Science, University of Warwick, UK

[‡]AI Labs - Didi Chuxing Co., Ltd. - Beijing, China

{rui_wang, xuemenghu, d.zhou, yuxuanxiong, chenchenye}@seu.edu.cn,
yulan.he@warwick.ac.uk, xuhaiyangsnow@didiglobal.com

Abstract

Recent years have witnessed a surge of interests of using neural topic models for automatic topic extraction from text, since they avoid the complicated mathematical derivations for model inference as in traditional topic models such as Latent Dirichlet Allocation (LDA). However, these models either typically assume improper prior (e.g. Gaussian or Logistic Normal) over latent topic space or could not infer topic distribution for a given document. To address these limitations, we propose a neural topic modeling approach, called Bidirectional Adversarial Topic (BAT) model, which represents the first attempt of applying bidirectional adversarial training for neural topic modeling. The proposed BAT builds a two-way projection between the document-topic distribution and the document-word distribution. It uses a generator to capture the semantic patterns from texts and an encoder for topic inference. Furthermore, to incorporate word relatedness information, the Bidirectional Adversarial Topic model with Gaussian (Gaussian-BAT) is extended from BAT. To verify the effectiveness of BAT and Gaussian-BAT, three benchmark corpora are used in our experiments. The experimental results show that BAT and Gaussian-BAT obtain more coherent topics, outperforming several competitive baselines. Moreover, when performing text clustering based on the extracted topics, our models outperform all the baselines, with more significant improvements achieved by Gaussian-BAT where an increase of near 6% is observed in accuracy.

1 Introduction

Topic models have been extensively explored in the Natural Language Processing (NLP) community for unsupervised knowledge discovery. Latent Dirichlet Allocation (LDA) (Blei et al., 2003), the

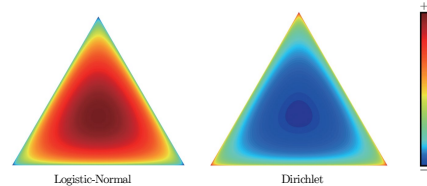


Figure 1: Illustrated probability simplex with Logistic-Normal distribution and Dirichlet distribution.

most popular topic model, has been extended (Lin and He, 2009; Zhou et al., 2014; Cheng et al., 2014) for various extraction tasks. Due to the difficulty of exact inference, most LDA variants require approximate inference methods, such as mean-field methods and collapsed Gibbs sampling. However, these approximate approaches have the drawback that small changes to the modeling assumptions result in a re-derivation of the inference algorithm, which can be mathematically arduous.

One possible way in addressing this limitation is through neural topic models which employ black-box inference mechanism with neural networks. Inspired by variational autoencoder (VAE) (Kingma and Welling, 2013), Srivastava and Sutton (2017) used the Logistic-Normal prior to mimic the simplex in latent topic space and proposed the Neural Variational LDA (NVLDA). Moreover, they replaced the word-level mixture in NVLDA with a weighted product of experts and proposed the ProdLDA (Srivastava and Sutton, 2017) to further enhance the topic quality.

Although Srivastava and Sutton (2017) used the Logistic-Normal distribution to approximate the Dirichlet distribution, they are not exactly the same. An illustration of these two distributions is shown in Figure 1 in which the Logistic-Normal distribution does not exhibit multiple peaks at the vertices of the simplex as that in the Dirichlet distribution and as such, it is less capable to capture

*corresponding author

the multi-modality which is crucial in topic modeling (Wallach et al., 2009). To deal with the limitation, Wang et al. (2019a) proposed the Adversarial-neural Topic Model (ATM) based on adversarial training, it uses a generator network to capture the semantic patterns lying behind the documents. However, given a document, ATM is not able to infer the document-topic distribution which is useful for downstream applications, such as text clustering. Moreover, ATM take the bag-of-words assumption and do not utilize any word relatedness information captured in word embeddings which have been proved to be crucial for better performance in many NLP tasks (Liu et al., 2018; Lei et al., 2018).

To address these limitations, we model topics with Dirichlet prior and propose a novel Bidirectional Adversarial Topic model (BAT) based on bidirectional adversarial training. The proposed BAT employs a generator network to learn the projection function from randomly-sampled document-topic distribution to document-word distribution. Moreover, an encoder network is used to learn the inverse projection, transforming a document-word distribution into a document-topic distribution. Different from traditional models that often resort to analytic approximations, BAT employs a discriminator which aims to discriminate between real distribution pair and fake distribution pair, thereby helps the networks (generator and encoder) to learn the two-way projections better. During the adversarial training phase, the supervision signal provided by the discriminator will guide the generator to construct a more realistic document and thus better capture the semantic patterns in text. Meanwhile, the encoder network is also guided to generate a more reasonable topic distribution conditioned on specific document-word distributions. Finally, to incorporate the word relatedness information captured by word embeddings, we extend the BAT by modeling each topic with a multivariate Gaussian in the generator and propose the Bidirectional Adversarial Topic model with Gaussian (Gaussian-BAT).

The main contributions of the paper are:

- We propose a novel Bidirectional Adversarial Topic (BAT) model, which is, to our best knowledge, the first attempt of using bidirectional adversarial training in neural topic modeling;
- We extend BAT to incorporate the word re-

latedness information into the modeling process and propose the Bidirectional Adversarial Topic model with Gaussian (Gaussian-BAT);

- Experimental results on three public datasets show that BAT and Gaussian-BAT outperform the state-of-the-art approaches in terms of topic coherence measures. The effectiveness of BAT and Gaussian-BAT is further verified in text clustering.

2 Related work

Our work is related to two lines of research, which are adversarial training and neural topic modeling.

2.1 Adversarial Training

Adversarial training, first employed in Generative Adversarial Network (GAN) (Goodfellow et al., 2014), has been extensively studied from both theoretical and practical perspectives.

Theoretically, Arjovsky (2017) and Gulrajani (2017) proposed the Wasserstein GAN which employed the Wasserstein distance between data distribution and generated distribution as the training objective. To address the limitation that most GANs (Goodfellow et al., 2014; Radford et al., 2015) could not project data into a latent space, Bidirectional Generative Adversarial Nets (BiGAN) (Donahue et al., 2016) and Adversarially Learned Inference (ALI) (Dumoulin et al., 2016) were proposed.

Adversarial training has also been extensively used for text generation. For example, SeqGAN (Yu et al., 2017) incorporated a policy gradient strategy for text generation. RankGAN (Lin et al., 2017) ranked a collection of human-written sentences to capture the language structure for improving the quality of text generation. To avoid mode collapse when dealing with discrete data, MaskGAN (Fedus et al., 2018) used an actor-critic conditional GAN to fill in missing text conditioned on the context.

2.2 Neural Topic Modeling

To overcome the challenging exact inference of topic models based on directed graph, a replicated softmax model (RSM), based on the Restricted Boltzmann Machines was proposed in (Hinton and Salakhutdinov, 2009). Inspired by VAE, Miao et al. (2016) used the multivariate Gaussian as the prior distribution of latent space and proposed the

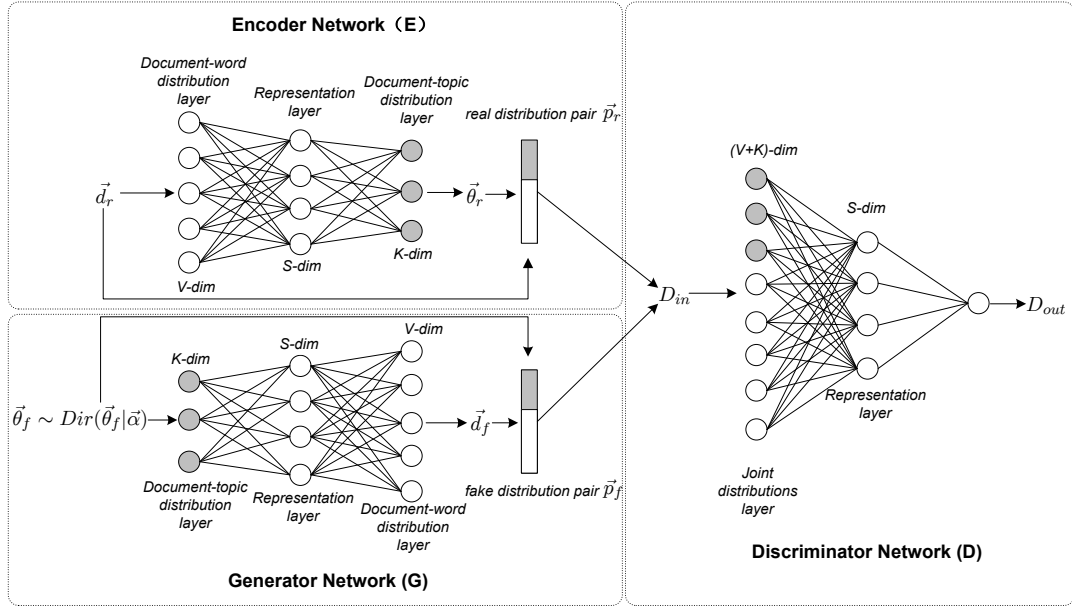


Figure 2: The framework of the Bidirectional Adversarial Topic (BAT) model.

Neural Variational Document Model (NVDM) for text modeling. To model topic properly, the Gaussian Softmax Model (GSM) (Miao et al., 2017) which constructs the topic distribution using a Gaussian distribution followed by a softmax transformation was proposed based on the NVDM. Likewise, to deal with the inappropriate Gaussian prior of NVDM, Srivastava and Sutton (2017) proposed the NVLDA which approximates the Dirichlet prior using a Logistic-Normal distribution. Recently, the Adversarial-neural Topic Model (ATM) (Wang et al., 2019a) is proposed based on adversarial training, it models topics with Dirichlet prior which is able to capture the multi-modality compared with logistic-normal prior and obtains better topics. Besides, the Adversarial-neural Event (AEM) (Wang et al., 2019b) model is also proposed for open event extraction by representing each event as an entity distribution, a location distribution, a keyword distribution and a date distribution.

Despite the extensive exploration of this research field, scarce work has been done to incorporate Dirichlet prior, word embeddings and bidirectional adversarial training into neural topic modeling. In this paper, we propose two novel topic modeling approaches, called BAT and Gaussian-BAT, which are different from existing approaches in the following aspects: (1) Unlike NVDM, GSM, NVLDA and ProLDA which model latent topic with Gaussian or logistic-normal prior, BAT and Gaussian-BAT explicitly employ Dirichlet prior to

model topics; (2) Unlike ATM which could not infer topic distribution of a given document, BAT and Gaussian-BAT uses a encoder to generate the topic distribution corresponding to the document; (3) Unlike neural topic models that only utilize word co-occurrence information, Gaussian-BAT models topic with multivariate Gaussian and incorporates the word relatedness into modeling process.

3 Methodology

Our proposed neural topic models are based on bidirectional adversarial training (Donahue et al., 2016) and aim to learn the two-way non-linear projection between two high-dimensional distributions. In this section, we first introduce the Bidirectional Adversarial Topic (BAT) model that only employs the word co-occurrence information. Then, built on BAT, we model topics with multivariate Gaussian in the generator of BAT and propose the Bidirectional Adversarial Topic model with Gaussian (Gaussian-BAT), which naturally incorporates word relatedness information captured in word embeddings into modeling process.

3.1 Bidirectional Adversarial Topic model

As depicted in Figure 2, the proposed BAT consists of three components: (1) The *Encoder E* takes the V -dimensional document representation \vec{d}_r sampled from text corpus C as input and transforms it into the corresponding K -dimensional topic distribution $\vec{\theta}_r$; (2) The *Generator G* takes a random

topic distribution $\vec{\theta}_f$ drawn from a Dirichlet prior as input and generates a V -dimensional fake word distribution \vec{d}_f ; (3) The *Discriminator* D takes the real distribution pair $\vec{p}_r = [\vec{\theta}_r; \vec{d}_r]$ and fake distribution pair $\vec{p}_f = [\vec{\theta}_f; \vec{d}_f]$ as input and discriminates the real distribution pairs from the fake ones. The outputs of the discriminator are used as supervision signals to learn E , G and D during adversarial training. In what follows, we describe each component in more details.

3.1.1 Encoder Network

The encoder learns a mapping function to transform document-word distribution to document-topic distribution. As shown in the top-left panel of Figure 2, it contains a V -dimensional document-word distribution layer, an S -dimensional representation layer and a K -dimensional document-topic distribution layer, where V and K denote vocabulary size and topic number respectively.

More concretely, for each document d in text corpus, E takes the document representation \vec{d}_r as input, where \vec{d}_r is the representation weighted by TF-IDF, and it is calculated by:

$$tf_{i,d} = \frac{n_{i,d}}{\sum_v n_{v,d}}, \quad idf_i = \log \frac{|C|}{|C_i|}$$

$$tf-idf_{i,d} = tf_{i,d} * idf_i, \quad d_r^i = \frac{tf-idf_{i,d}}{\sum_v tf-idf_{v,d}}$$

where $n_{i,d}$ denotes the number of i -th word appeared in document d , $|C|$ represents the number of documents in the corpus, and $|C_i|$ means the number of documents that contain i -th word in the corpus. Thus, each document could be represented as a V -dimensional multinomial distribution and the i -th dimension denotes the semantic consistency between i -th word and the document.

With \vec{d}_r as input, E firstly projects it into an S -dimensional semantic space through the representation layer as follows:

$$\vec{h}_s^e = \text{BN}(W_s^e \vec{d}_r + \vec{b}_s^e) \quad (1)$$

$$\vec{o}_s^e = \max(\vec{h}_s^e, leak * \vec{h}_s^e) \quad (2)$$

where $W_s^e \in \mathbb{R}^{S \times V}$ and \vec{b}_s^e are weight matrix and bias term of the representation layer, \vec{h}_s^e is the state vector normalized by batch normalization $\text{BN}(\cdot)$, $leak$ denotes the parameter of LeakyReLU activation and \vec{o}_s^e represents the output of representation layer.

Then, the encoder transforms \vec{o}_s^e into a K -dimensional topic space based on the equation below:

$$\vec{\theta}_r = \text{softmax}(W_t^e \vec{o}_s^e + \vec{b}_t^e) \quad (3)$$

where $W_t^e \in \mathbb{R}^{K \times S}$ is the weight matrix of topic distribution layer, \vec{b}_t^e represents the bias term, $\vec{\theta}_r$ denotes the corresponding topic distribution of the input \vec{d}_r and the k -th ($k \in \{1, 2, \dots, K\}$) dimension θ_r^k represents the proportion of k -th topic in document d .

3.1.2 Generator network

The generator G is shown in the bottom-left panel of Figure 2. Contrary to encoder, it provides an inverse projection from document-topic distribution to document-word distribution and contains a K -dimensional document-topic layer, an S -dimensional representation layer and a V -dimensional document-word distribution layer.

As pointed out in (Wallach et al., 2009), the choice of Dirichlet prior over topic distribution is important to obtain interpretable topics. Thus, BAT employs the Dirichlet prior parameterized with $\vec{\alpha}$ to mimic the multi-variate simplex over topic distribution $\vec{\theta}_f$. It can be drawn randomly based on the equation below:

$$p(\vec{\theta}_f | \vec{\alpha}) = \text{Dir}(\vec{\theta}_f | \vec{\alpha}) \triangleq \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^K [\theta_f^k]^{\alpha_k - 1} \quad (4)$$

where $\vec{\alpha}$ is the K -dimensional hyper-parameter of Dirichlet prior, K is the topic number that should be set in BAT, $\theta_f^k \in [0, 1]$, follows the constrain that $\sum_{k=1}^K \theta_f^k = 1$, represents the proportion of the k -th topic in the document, and normalization term $\Delta(\vec{\alpha})$ is defined as $\frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$.

To learn the transformation from document-topic distribution to document-word distribution, G firstly projects $\vec{\theta}_f$ into an S -dimensional representation space based on equations:

$$\vec{h}_s^g = \text{BN}(W_s^g \vec{\theta}_f + \vec{b}_s^g) \quad (5)$$

$$\vec{o}_s^g = \max(\vec{h}_s^g, leak * \vec{h}_s^g) \quad (6)$$

where $W_s^g \in \mathbb{R}^{S \times K}$ is weight matrix of the representation layer, \vec{b}_s^g represents bias term, \vec{h}_s^g is the state vector normalized by batch normalization, Eq. 6 represents the LeakyReLU activation parameterized with $leak$, and \vec{o}_s^g is the output of the representation layer.

Then, to project $\vec{\theta}_s^g$ into word distribution \vec{d}_f , a subnet contains a linear layer and a softmax layer is used and the transformation follows:

$$\vec{d}_f = \text{softmax}(W_w^g \vec{\theta}_s^g + \vec{b}_w^g) \quad (7)$$

where $W_w^g \in \mathbb{R}^{V \times S}$ and \vec{b}_w^g are weight matrix and bias of word distribution layer, \vec{d}_f is the word distribution correspond to $\vec{\theta}_f$. For each $v \in \{1, 2, \dots, V\}$, the v -th dimension d_f^v is the probability of the v -th word in fake document \vec{d}_f .

3.1.3 Discriminator network

The discriminator D is constituted by three layers (a $V + K$ -dimensional joint distribution layer, an S -dimensional representation layer and an output layer) as shown in the right panel of Figure 2. It employs real distribution pair \vec{p}_r and fake distribution pair \vec{p}_f as input and then outputs D_{out} to identify the input sources (fake or real). Concretely, a higher value of D_{out} represents that D is more prone to predict the input as real and vice versa.

3.2 BAT with Gaussian (Gaussian-BAT)

In BAT, the generator models topics based on the bag-of-words assumption as in most other neural topic models. To incorporate the word relatedness information captured in word embeddings (Mikolov et al., 2013a,b; Pennington et al., 2014; Joulin et al., 2017; Athiwaratkun et al., 2018) into the inference process, we modify the generator of BAT and propose Gaussian-BAT, in which G models each topic with a multivariate Gaussian as shown in Figure 3.

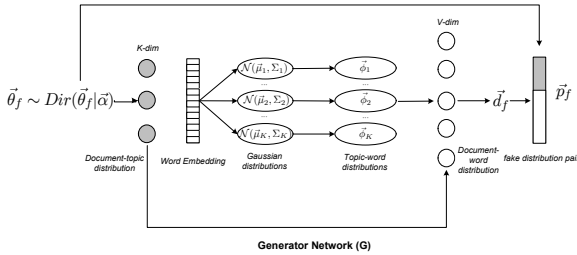


Figure 3: The generator of Gaussian-BAT.

Concretely, Gaussian-BAT employs the multivariate Gaussian $\mathcal{N}(\vec{\mu}_k, \Sigma_k)$ to model the k -th topic. Here, $\vec{\mu}_k$ and Σ_k are trainable parameters, they represent mean and covariance matrix respectively. Following its probability density, for each word $v \in \{1, 2, \dots, V\}$, the probability in the k -th

topic $\phi_{k,v}$ is calculated by:

$$p(\vec{e}_v | \text{topic} = k) = \mathcal{N}(\vec{e}_v; \vec{\mu}_k, \Sigma_k) = \frac{\exp(-\frac{1}{2}(\vec{e}_v - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{e}_v - \vec{\mu}_k))}{\sqrt{(2\pi)^{D_e} |\Sigma_k|}} \quad (8)$$

$$\phi_{k,v} = \frac{p(\vec{e}_v | \text{topic} = k)}{\sum_{v=1}^V p(\vec{e}_v | \text{topic} = k)} \quad (9)$$

where \vec{e}_v means the word embedding of v -th word, V is the vocabulary size, $|\Sigma_k| = \det \Sigma_k$ is the determinant of covariance matrix Σ_k , D_e is the dimension of word embeddings, $p(\vec{e}_v | \text{topic} = k)$ is the probability calculated by density, and $\vec{\phi}_k$ is the normalized word distribution of k -th topic. With randomly sampled topic distribution $\vec{\theta}_f$ and the calculated topic-word distributions $\{\vec{\phi}_1, \vec{\phi}_2, \dots, \vec{\phi}_K\}$, the fake word distribution \vec{d}_f corresponding to $\vec{\theta}_f$ can be obtained by:

$$\vec{d}_f = \sum_{k=1}^K \vec{\phi}_k * \theta_k \quad (10)$$

where θ_k is the topic proportion of the k -th topic. Then, $\vec{\theta}_f$ and \vec{d}_f are concatenated to form the fake distribution pair \vec{p}_f as shown in Figure 3. And encoder and discriminator of Gaussian-BAT are same as BAT, shown as Figure 2. In our experiments, the pre-trained 300-dimensional *Glove* (Pennington et al., 2014) embedding is used.

3.3 Objective and Training Procedure

In Figure 2, the real distribution pair $\vec{p}_r = [\vec{\theta}_r; \vec{d}_r]$ and the fake distribution pair $\vec{p}_f = [\vec{\theta}_f; \vec{d}_f]$ can be viewed as random samples drawn from two $(K + V)$ -dimensional joint distributions \mathbb{P}_r and \mathbb{P}_f , each of them comprising of a K -dimensional Dirichlet distribution and a V -dimensional Dirichlet distribution. The training objective of BAT and Gaussian-BAT is to make the generated joint distribution \mathbb{P}_f close to the real joint distribution \mathbb{P}_r as much as possible. In this way, a two-way projection between document-topic distribution and document-word distribution could be built by the learned encoder and generator.

To measure the distance between \mathbb{P}_r and \mathbb{P}_f , we use the Wasserstein-distance as the optimization objective, since it was shown to be more effective compared to Jensen-Shannon divergence (Arjovsky et al., 2017):

$$Loss = \mathbb{E}_{\vec{p}_f \sim \mathbb{P}_f} [D(\vec{p}_f)] - \mathbb{E}_{\vec{p}_r \sim \mathbb{P}_r} [D(\vec{p}_r)] \quad (11)$$

where $D(\cdot)$ represents the output signal of the discriminator. A higher value denotes that the discriminator is more prone to consider the input as a real distribution pair and vice versa. In addition, we use weight clipping which was proposed to ensure the Lipschitz continuity (Arjovsky et al., 2017) of D .

Algorithm 1 Training procedure for BAT and Gaussian-BAT

Input: $K, c, n_d, m, \alpha_1, \beta_1, \beta_2$

Output: The trained encoder E and generator G .

```

1: Initialize  $D, E$  and  $G$  with  $\omega_d, \omega_e$  and  $\omega_g$ 
2: while  $\omega_e$  and  $\omega_g$  have not converged do
3:   for  $t = 1, \dots, n_d$  do
4:     for  $j = 1, \dots, m$  do
5:       Sample  $\vec{d}_r \sim \mathbb{P}_r^d$ ,
6:       Sample a random  $\vec{\theta}_f \sim Dir(\vec{\theta}_f | \vec{\alpha})$ 
7:        $\vec{d}_f \leftarrow G(\vec{\theta}_f), \vec{\theta}_r \leftarrow E(\vec{d}_r)$ 
8:        $\vec{p}_r = [\vec{\theta}_r; \vec{d}_r], \vec{p}_f = [\vec{\theta}_f; \vec{d}_f]$ 
9:        $L^{(j)} = D(\vec{p}_f) - D(\vec{p}_r)$ 
10:    end for
11:     $\omega_d \leftarrow Adam(\nabla_{\omega_d} \frac{1}{m} \sum_{j=1}^m L^{(j)}, \omega_d, p_a)$ 
12:     $\omega_d \leftarrow clip(\omega_d, -c, c)$ 
13:  end for
14:   $\omega_g \leftarrow Adam(\nabla_{\omega_g} \frac{-1}{m} \sum_{j=1}^m D(\vec{p}_f^j), \omega_g, p_a)$ 
15:   $\omega_e \leftarrow Adam(\nabla_{\omega_e} \frac{1}{m} \sum_{j=1}^m D(\vec{p}_r^j), \omega_e, p_a)$ 
16: end while

```

The training procedure of BAT and Gaussian-BAT is given in Algorithm. 1. Here, c is the clipping parameter, n_d represents the number of discriminator iterations per generator iteration, m is the batch size, α_1 is the learning rate, β_1 and β_2 are hyper-parameters of Adam (Kingma and Ba, 2014), and p_a represents $\{\alpha_1, \beta_1, \beta_2\}$. In our experiments, we set the $n_d = 5, m = 64, \alpha_1 = 1e-4, c = 0.01, \beta_1 = 0.5$ and $\beta_2 = 0.999$.

3.4 Topic Generation and Cluster Inference

After model training, learned G and E will build a two-way projection between document-topic distribution and document-word distribution. Thus, G and E could be used for topic generation and cluster inference.

To generate the word distribution of each topic, we use $\vec{t}_{S(k)}$, a K -dimensional vector, as the one-hot encoding of the k -th topic. For example, $\vec{t}_{S_2} = [0, 1, 0, 0, 0, 0]^T$ in a six topic setting. And the word

distribution of the k -th topic is obtained by:

$$\vec{\phi}_k = G(\vec{t}_{S(k)}) \quad (12)$$

Likewise, given the document representation \vec{d}_r , topic distribution $\vec{\theta}_r$ obtained by BAT/Gaussian-BAT could be used for cluster inference based on:

$$\vec{\theta}_r = E(\vec{d}_r); \quad c_r = \arg \max \vec{\theta}_r \quad (13)$$

where c_r denotes the inferred cluster of \vec{d}_r .

4 Experiments

In this section, we first present the experimental setup which includes the datasets used and the baselines, followed by the experimental results.

4.1 Experimental Setup

We evaluate BAT and Gaussian-BAT on three datasets for topic extraction and text clustering, 20Newsgroups¹, Grolier² and NYTimes³. Details are summarized below:

20Newsgroups (Lang, 1995) is a collection of approximately 20,000 newsgroup articles, partitioned evenly across 20 different newsgroups.

Grolier is built from Grolier Multimedia Encyclopedia, which covers almost all the fields in the world. *NYTimes* is a collection of news articles published between 1987 and 2007, and contains a wide range of topics, such as sports, politics, education, etc.

We use the full datasets of 20Newsgroups¹ and Grolier². For the NYTimes dataset, we randomly select 100,000 articles and remove the low frequency words. The final statistics are shown in Table 1:

Dataset	#Doc (Train)	#Doc (Test)	#Words
20Newsgroups	11,259	7,488	1,995
Grolier	29,762	-	15,276
NYtimes	99,992	-	12,604

Table 1: The statistics of datasets.

We choose the following models as baselines: **LDA** (Blei et al., 2003) extracts topics based on word co-occurrence patterns from documents. We implement LDA following the parameter setting suggested in (Griffiths and Steyvers, 2004). **NVDM** (Miao et al., 2016) is an unsupervised text modeling approach based on VAE. We use the original implementation of the paper⁴.

¹<http://qwone.com/~jason/20Newsgroups/>

²<https://cs.nyu.edu/~roweis/data/>

³<http://archive.ics.uci.edu/ml/datasets/Bag+of+Words>

⁴<https://github.com/ysmiao/nvdm>

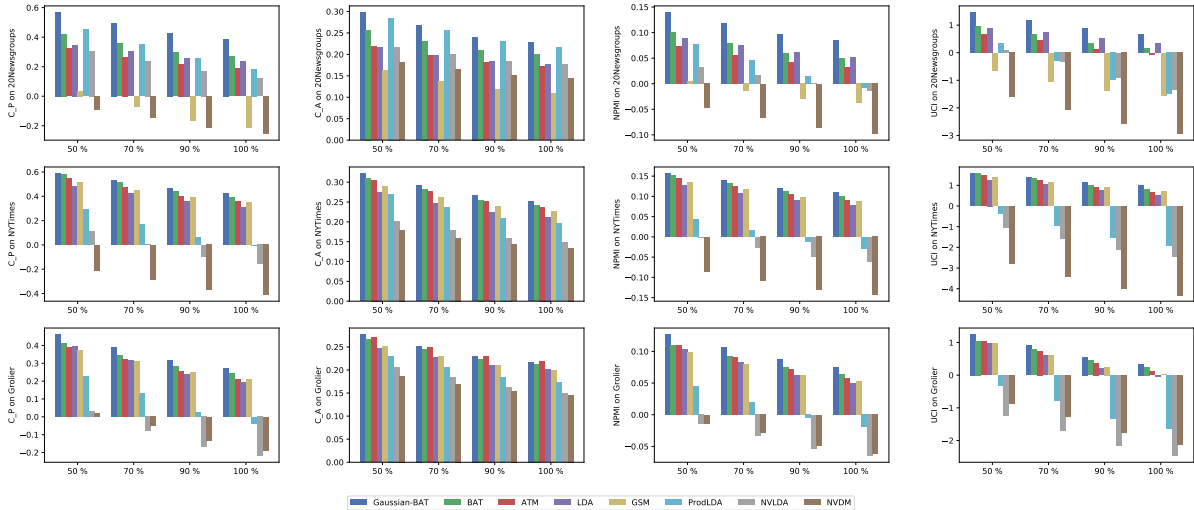


Figure 4: The comparison of average topic coherence vs. different topic proportion on three datasets.

GSM(Miao et al., 2017) is an enhanced topic model based on NVDM, we use the original implementation in our experiments⁵.

NVLDA (Srivastava and Sutton, 2017), also built on VAE but with the logistic-normal prior. We use the implementation provided by the author⁶.

ProdLDA (Srivastava and Sutton, 2017), is a variant of NVLDA, in which the distribution over individual words is a product of experts. The original implementation is used.

ATM (Wang et al., 2019a), is a neural topic modeling approach based on adversarial training, we implement the ATM following the parameter setting suggested in the original paper.

4.2 Topic Coherence Evaluation

Topic models are typically evaluated with the likelihood of held-out documents and topic coherence. However, Chang et al. (2009) showed that a higher likelihood of held-out documents does not correspond to human judgment of topic coherence. Thus, we follow (Röder et al., 2015) and employ four topic coherence metrics (C_P, C_A, NPMI and UCI) to evaluate the topics generated by various models. In all experiments, each topic is represented by the top 10 words according to the topic-word probabilities, and all the topic coherence values are calculated using the Palmetto library⁷.

We firstly make a comparison of topic coherence vs. different topic proportions. Experiments are

Dataset	Model	C_P	C_A	NPMI	UCI
20Newsgroups	NVDM	-0.2558	0.1286	-0.0984	-2.9496
	GSM	-0.2318	0.1067	-0.0400	-1.6083
	NVLDA	0.1205	0.1763	-0.0207	-1.3466
	ProdLDA	0.1858	0.2155	-0.0083	-1.5044
	LDA	0.2361	0.1769	0.0523	0.3399
	ATM	0.1914	0.1720	0.0207	-0.3871
	BAT	0.2597	0.1976	0.0472	0.0969
	Gaussian-BAT	0.3758	0.2251	0.0819	0.5925
Grolier	NVDM	-0.1877	0.1456	-0.0619	-2.1149
	GSM	0.1974	0.1966	0.0491	-0.0410
	NVLDA	-0.2205	0.1504	-0.0653	-2.4797
	ProdLDA	-0.0374	0.1733	-0.0193	-1.6398
	LDA	0.1908	0.2009	0.0497	-0.0503
	ATM	0.2105	0.2188	0.0582	0.1051
	BAT	0.2312	0.2108	0.0608	0.1709
	Gaussian-BAT	0.2606	0.2142	0.0724	0.2836
NYtimes	NVDM	-0.4130	0.1341	-0.1437	-4.3072
	GSM	0.3426	0.2232	0.0848	0.6224
	NVLDA	-0.1575	0.1482	-0.0614	-2.4208
	ProdLDA	-0.0034	0.1963	-0.0282	-1.9173
	LDA	0.3083	0.2127	0.0772	0.5165
	ATM	0.3568	0.2375	0.0899	0.6582
	BAT	0.3749	0.2355	0.0951	0.7073
	Gaussian-BAT	0.4163	0.2479	0.1079	0.9215

Table 2: Average topic coherence on three datasets with five topic settings [20, 30, 50, 75, 100].

conducted on the datasets with five topic number settings [20, 30, 50, 75, 100]. We calculate the average topic coherence values among topics whose coherence values are ranked at the top 50%, 70%, 90%, 100% positions. For example, to calculate the average C_P value of BAT @90%, we first compute the average C_P coherence with the selected topics whose C_P values are ranked at the top 90% for each topic number setting, and then average the five coherence values with each corresponding to a particular topic number setting.

The detailed comparison is shown in Figure 4. It can be observed that BAT outperforms the baselines on all the coherence metrics for NYTimes datasets. For Grolier dataset, BAT outperforms all the baselines on C_P, NPMI and UCI metrics, but

⁵<https://github.com/linkstrife/NVDM-GSM>

⁶[https://github.com/akashgit/autoencoding vi for topic models](https://github.com/akashgit/autoencoding_vi_for_topic_models)

⁷<https://github.com/dice-group/Palmetto>

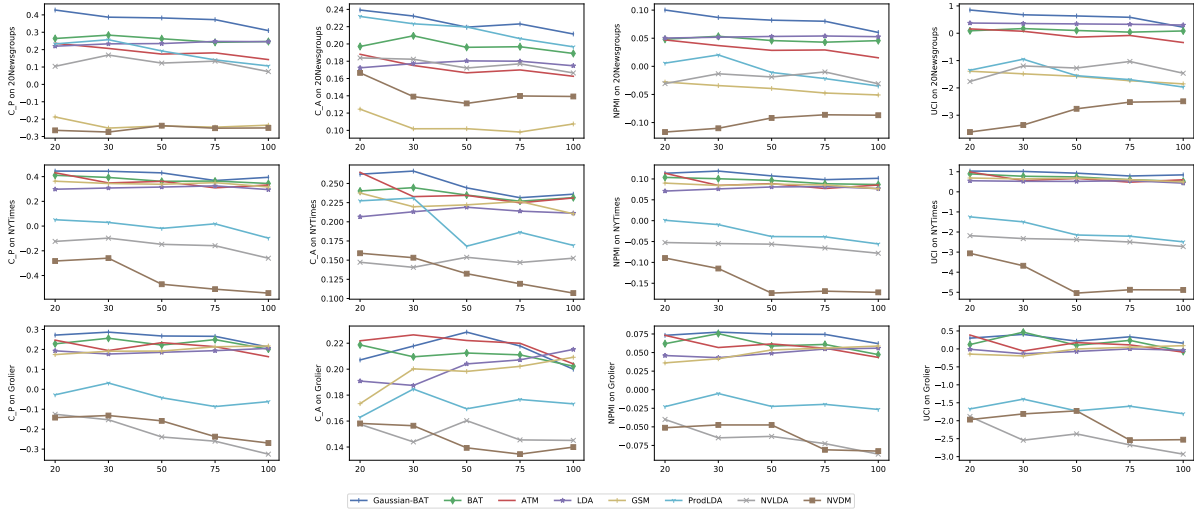


Figure 5: The comparison of average topic coherence vs. different topic number on 20Newsgroups, Grolier and NYTimes.

Model	Topics
Gaussian-BAT	voter campaign poll candidates democratic election republican vote presidential democrat song album music band rock pop sound singer jazz guitar film movie actor character movies director series actress young scenes flight airline passenger airlines aircraft shuttle airport pilot carrier planes
BAT	vote president voter campaign election democratic governor republican black candidates album band music rock song jazz guitar pop musician <i>record</i> film actor play acting role playing character father movie actress flight airline delay airlines plane pilot airport passenger carrier attendant
LDA	voter vote poll election campaign primary candidates republican race party music song band sound <i>record</i> artist album show musical rock film movie character play actor director movies <i>minutes</i> theater cast flight plane <i>ship crew</i> air pilot <i>hour boat</i> passenger airport
ATM	voter vote poll republican race <i>primary percent</i> election campaign democratic music song musical album jazz band <i>record recording</i> mp3 composer film movie actor director award movies character <i>theater production</i> play jet flight airline <i>hour</i> plane passenger <i>trip plan</i> travel pilot

Table 3: Topic examples extracted by models, italics means out-of-topic words. These topics correspond to ‘election’, ‘music’, ‘film’ and ‘airline’ respectively, and topic examples of other models are omitted due to poor quality.

gives slightly worse results compared to ATM on C_A. For 20Newsgroups dataset, BAT performs the best on C_P and NPMI, but gives slightly worse results compared to ProLDA on C_A, and LDA on UCI. By incorporating word embeddings through trainable Gaussian distribution, Gaussian-BAT outperforms all the baselines and BAT on four coherence metrics, often by a large margin, across all the three datasets except for Grolier dataset on C_A when considering 100% topics. This may be attribute to the following factors: (1) The Dirichlet prior employed in BAT and Gaussian-BAT could exhibit a multi-modal distribution in latent space and is more suitable for discovering semantic patterns from text; (2) ATM does not consider the relationship between topic distribution and word distribution since it only carry out adversarial training in word distribution space; (3) The incorpora-

tion of word embeddings in Gaussian-BAT helps generating more coherent topics.

We also compare the average topic coherence values (all topics taken into account) numerically to show the effectiveness of proposed BAT and Gaussian-BAT. The results of numerical topic coherence comparison are listed in Table 2 and each value is calculated by averaging the average topic coherences over five topic number settings. The best coherence value on each metric is highlighted in bold. It can be observed that Gaussian-BAT gives the best overall results across all metrics and on all the datasets except for Grolier dataset on C_A. To make the comparison of topics more intuitive, we provide four topic examples extracted by models in Table 3. It can be observed that the proposed BAT and Gaussian-BAT can generate more coherent topics.

Moreover, to explore how topic coherence varies with different topic numbers, we also provide the comparison of average topic coherence vs. different topic number on 20newsgroups, Grolier and NYTimes (all topics taken into account). The detailed comparison is shown in Figure 5. It could be observed that Gaussian-BAT outperforms the baselines with 20, 30, 50 and 75 topics except for Grolier dataset on C_A metric. However, when the topic number is set to 100, Gaussian-BAT performs slightly worse than LDA (e.g., UCI for 20NewsGroups and C_A for NYTimes). This may be caused by the increased model complexity due to the larger topic number settings. Likewise, BAT can achieve at least the second-best results among all the approaches in most cases for NYTimes dataset. For Grolier, BAT also performs the second-best except on C_A metric. However, for 20newsgroups, the results obtained by BAT are worse than ProLDA (C_A) and LDA (UCI) due to the limited training documents in the dataset, though it still largely outperforms other baselines.

4.3 Text Clustering

We further compare our proposed models with baselines on text clustering. Due to the lack of document label information in Grolier and NYTimes, we only use 20NewsGroups dataset in our experiments. The topic number is set to 20 (ground-truth categories) and the performance is evaluated by accuracy (ACC):

$$ACC = \max_{\text{map}} \frac{\sum_{i=1}^{N_t} \text{ind}(l_i = \text{map}(c_i))}{N_t} \quad (14)$$

where N_t is the number of documents in the test set, $\text{ind}(\cdot)$ is the indicator function, l_i is the ground-truth label of i -th document, c_i is the category assignment, and map ranges over all possible one-to-one mappings between labels and clusters. The optimal map function can be obtained by the Kuhn-Munkres algorithm (Kuhn, 1955). A larger accuracy value indicates a better text clustering results.

Dataset	NVLDA	ProdLDA	LDA	BAT	G-BAT
20NG	33.31%	33.82%	35.36%	35.66%	41.25%

Table 4: Text clustering accuracy on 20NewsGroups (20NG). ‘G-BAT’ refers to ‘Gaussian-BAT’. The best result is highlighted in bold.

The comparison of text clustering results on 20NewsGroups is shown in Table 4. Due to the

poor performance of NVDM in topic coherence evaluation, its result is excluded here. Not surprisingly, NVLDA and ProdLDA perform worse than BAT and Gaussian-BAT that model topics with the Dirichlet prior. This might be caused by the fact that Logistic-Normal prior does not exhibit multiple peaks at the vertices of the simplex, as depicted in Figure 1. Compared with LDA, BAT achieves a comparable result in accuracy since both models have the same Dirichlet prior assumption over topics and only employ the word co-occurrence information. Gaussian-BAT outperforms the second best model, BAT, by nearly 6% in accuracy. This shows that the incorporation of word embeddings is important to improve the semantic coherence of topics and thus results in better consistency between cluster assignments and ground-truth labels.

5 Conclusion

In this paper, we have explored the use of bidirectional adversarial training in neural topic models and proposed two novel approaches: the Bidirectional Adversarial Topic (BAT) model and the Bidirectional Adversarial Topic model with Gaussian (Gaussian-BAT). BAT models topics with the Dirichlet prior and builds a two-way transformation between document-topic distribution and document-word distribution via bidirectional adversarial training. Gaussian-BAT extends from BAT by incorporating word embeddings into the modeling process, thereby naturally considers the word relatedness information captured in word embeddings. The experimental comparison on three widely used benchmark text corpus with the existing neural topic models shows that BAT and Gaussian-BAT achieve improved topic coherence results. In the future, we would like to devise a nonparametric neural topic model based on adversarial training. Besides, developing correlated topic models is another promising direction.

Acknowledgements

We would like to thank anonymous reviewers for their valuable comments and helpful suggestions. This work was funded by the National Key Research and Development Program of China(2017YFB1002801) and the National Natural Science Foundation of China (61772132). And YH is partially supported by EPSRC (grant no. EP/T017112/1).

References

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223.
- Ben Athiwaratkun, Andrew Wilson, and Anima Anandkumar. 2018. Probabilistic fasttext for multi-sense word embeddings. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1–11.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.
- Jonathan Chang, Sean Gerrish, Chong Wang, Jordan L Boyd-Graber, and David M Blei. 2009. Reading tea leaves: How humans interpret topic models. In *Advances in neural information processing systems*, pages 288–296.
- 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.
- Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. 2016. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*.
- Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. 2016. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*.
- William Fedus, Ian Goodfellow, and Andrew M Dai. 2018. Maskgan: better text generation via filling in the... *arXiv preprint arXiv:1801.07736*.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680.
- Thomas L Griffiths and Mark Steyvers. 2004. Finding scientific topics. *Proceedings of the National academy of Sciences*, 101(suppl 1):5228–5235.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. 2017. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777.
- Geoffrey E Hinton and Ruslan R Salakhutdinov. 2009. Replicated softmax: an undirected topic model. In *Advances in neural information processing systems*, pages 1607–1614.
- Armand Joulin, Edouard Grave, and Piotr Bojanowski Tomas Mikolov. 2017. Bag of tricks for efficient text classification. *EACL 2017*, page 427.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Harold W Kuhn. 1955. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2:83–97.
- Ken Lang. 1995. Newsweeder: Learning to filter netnews. In *Proceedings of the Twelfth International Conference on Machine Learning*, pages 331–339.
- Zeyang Lei, Yujiu Yang, and Min Yang. 2018. Saan: A sentiment-aware attention network for sentiment analysis. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1197–1200. ACM.
- Chenghua Lin and Yulan He. 2009. Joint sentiment/topic model for sentiment analysis. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pages 375–384. ACM.
- Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, and Ming-Ting Sun. 2017. Adversarial ranking for language generation. In *Advances in Neural Information Processing Systems*, pages 3155–3165.
- Qiao Liu, Haibin Zhang, Yifu Zeng, Ziqi Huang, and Zufeng Wu. 2018. Content attention model for aspect based sentiment analysis. In *Proceedings of the 2018 World Wide Web Conference*, pages 1023–1032. International World Wide Web Conferences Steering Committee.
- 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-Volume 70*, pages 2410–2419. JMLR. org.
- Yishu Miao, Lei Yu, and Phil Blunsom. 2016. Neural variational inference for text processing. In *International conference on machine learning*, pages 1727–1736.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.

- Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.
- Michael Röder, Andreas Both, and Alexander Hinneburg. 2015. Exploring the space of topic coherence measures. In *Proceedings of the eighth ACM international conference on Web search and data mining*, pages 399–408. ACM.
- Akash Srivastava and Charles Sutton. 2017. Autoencoding variational inference for topic models. *arXiv preprint arXiv:1703.01488*.
- Hanna M Wallach, David M Mimno, and Andrew McCallum. 2009. Rethinking lda: Why priors matter. In *Advances in neural information processing systems*, pages 1973–1981.
- Rui Wang, Deyu Zhou, and Yulan He. 2019a. Atm: Adversarial-neural topic model. *Information Processing & Management*, 56(6):102098.
- Rui Wang, Deyu Zhou, and Yulan He. 2019b. Open event extraction from online text using a generative adversarial network. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 282–291, Hong Kong, China.
- Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- Deyu Zhou, Liangyu Chen, and Yulan He. 2014. A simple bayesian modelling approach to event extraction from twitter. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 700–705.