# Neural Topic Modeling with Deep Mutual Information Estimation

Kang Xu[a], Xiaoqiu Lu[b], Yuan-fang Li[c], Tongtong Wu[b], Guilin Qi[b,*], Ning Ye[a], Dong Wang[d], Zheng Zhou[d]

[a]*School of Computer Science, Nanjing University of Posts and Telecommunications, China*
[b]*School of Computer Science and Engineering, Southeast University,China*
[c]*Faculty of Information Technology, Monash University, Australia*
[d]*Department of Earth System Science, Tsinghua University, Beijing 100084, China*

## Abstract

The emerging neural topic models make topic modeling more easily adaptable and extendable in unsupervised text mining. However, the existing neural topic models is difficult to retain representative information of the documents within the learnt topic representation. In this paper, we propose a neural topic model which incorporates deep mutual information estimation, i.e., Neural Topic Modeling with Deep Mutual Information Estimation(NTM-DMIE). NTM-DMIE is a neural network method for topic learning which maximizes the mutual information between the input documents and their latent topic representation. To learn robust topic representation, we incorporate the discriminator to discriminate negative examples and positive examples via adversarial learning. Moreover, we use both global and local mutual information to preserve the rich information of the input documents in the topic representation. We evaluate NTM-DMIE on several metrics, including accuracy of text clustering, with topic representation, topic uniqueness and topic coherence. Compared to the existing methods, the experimental results show that NTM-DMIE can outperform in all the metrics on the four datasets.

*Keywords:* Neural Topic Modeling, Deep Mutual Information, Topic Discovery,

*This is to indicate the corresponding author.

*Email addresses:* kxu@njupt.edu.cn (Kang Xu), luxq0823@163.com (Xiaoqiu Lu), yuanfang.li@monash.edu (Yuan-fang Li), wutong8023@seu.edu.cn (Tongtong Wu), gqi@seu.edu.cn (Guilin Qi), yening@njupt.edu.cn (Ning Ye), wangdong19@mails.tsinghua.edu.cn (Dong Wang), zhouz@mail.tsinghua.edu.cn (Zheng Zhou)

## 1. Introduction

Topic models aim at discovering latent *semantic topics* from a corpus of text documents and have been widely employed in information retrieval and related fields. The field of topic modeling has shifted away from "Bag-of-Words" representations such as Latent Dirichlet Allocation (LDA) [1] to neural networks based methods [2, 3, 4, 5], which achieve state-of-the-art performance.

Srivastava and Sutton [6] use an Autocoder-based topic model which constructs a Laplace approximation to the Dirichlet prior, and the proposed ProdLDA uses product of experts to learn topic representations. Similarly, a number of VAE-based neural topic models have also been proposed [2, 3, 7, 8]. Different from the aforementioned neural topic model based on Gaussian distributions, Esmaeili et al [4] use a neural topic model (VALTA) with Gumbel-Softmax[9] to simulate discrete distributions. The advantage of using Gumbel-Softmax is that it promotes sparsity and leads to more disentangled representations, i.e. topics [10]. Another type of widely used neural topic models are based on Generative Adversarial Networks (GAN) [11]. Wang et al [12] propose the Adversarial neural Topic Model (ATM) that is based on adversarial training. Moreover, they also propose the Bidirectional Adversarial Topic (BAT) method [5] that models topics with the Dirichlet prior and builds a two-way transformation between the document-topic distribution and the document-word distribution via bidirectional adversarial training.

A main limitation of the existing neural topic models is that they only try to constrain the learned topic representations without regarding the useful information conveyed by the input documents. The useful information is the representative information to distinguish the text from others. For example, in the text which are composed of words, the topics of the text and its words can convey the useful representative information which is informative in text clustering and classification. Here, a *good* representation is one that can retain as much useful information of the input text as possible [11]. In topic modeling, the amount of useful information in learned topic representations

is important for the tasks, i.e., topic distribution learning and topic word mining. To alleviate this problem, a simple and effective way is to train a representation learning network to maximize the mutual information (MI) [13], between the input documents and their latent topic representations. Since mutual information can characterize both the relevance and the redundancy between random variables, it can effectively model the association of different variables. However, mutual information is difficult to estimate, especially in high-dimensional and continuous settings.

In this paper, we introduce deep mutual information estimation [14] to topic modeling. Our method, Neural Topic Modeling with Deep Mutual Information Estimation (NTM-DMIE), effectively estimates and maximizes mutual information between high-dimensional input (document) and output (topic) pairs with deep neural networks. To best utilize the rich information contained in documents in learning topic representations, we regard documents as *global* information and words contained in documents as *local* information. Globally, we maximize MI between the documents attached with negative examples and the learned topic representations. Locally, NTM-DMIE maximizes the average MI between topic representations with document words to further improve the representation quality.

The main contributions of our work are summarized as follows:

- We propose a novel neural topic modeling technique with deep mutual information estimation to better utilize document information. To the best of our knowledge, this is the first work to incorporate deep mutual information estimation into topic learning to improve the quality of topic representations and topic mining tasks.

- We propose to use global and local mutual information maximization to preserve the rich information contained in documents for learning their latent topic representations.

- Extensive experimental results on four benchmark datasets show that our NTM-DMIE model outperforms recent, strong baseline methods.

3

## 2. Related Work

Our work is mainly related to neural topic modeling and deep mutual information estimation. We briefly discuss their recent progress.

### 2.1. Neural Topic Modeling

Recently, neural networks[15, 11] have been employed in topic modeling, which are more effective and efficient at approximating the hidden, complex variables in the topic models. Based on Variational Auto-encoder (VAE), Miao et al [2] proposed the Neural Variational Document Model (NVDM), which builds a deep neural network conditioned on text to approximate the intractable distributions over the latent variables. Moreover, they [3] further proposed the Gaussian Softmax topic model (GSM), parameterized with neural networks. NVDM was extended for generalizing topic models to model with covariates, interactions, and customized regularizers [7]. Card et al [16] developed a supervised neural topic model (SCHOLAR) which models metadata as a covariate or a predicted variable. ProdLDA [6] and Variational Aspect-based Latent Topic Allocation (VALTA) [4] also embed relationships between documents, topics, and words in differentiable functions. Moreover, Neural Topic Model (NTM) [17] and Variational Topic Model with Reinforcement Learning(VTMRL) [8] both incorporated topic coherence into topic modeling.

Some neural topic models were designed based on Generative Adversarial Network(GAN), Wang et al [12, 5, 5] proposed the Adversarial neural Topic Model (ATM) based on adversarial training.Gupta et al [18] proposed a neural autoregressive topic model, DocNADE, to exploit the full context information around words in a document in a language modeling fashion.

With the rapid development of topic-aware related work, topic models were designed with many other machine learning methods. Zhao X et al [19] proposed the Variational Auto-Encoder Topic Model (VAETM) by combining word vector representation and entity vector representation to address the limitations for mining high-quality topics from short texts. Panwar M et al [20] proposed the Topic Attention Networks for Neural Topic Modeling(TAN-NTM), which processed document as a sequence of tokens through an LSTM whose contextual outputs are attended in a topic-aware manner.

Bahrainian S A et al [21] proposed a new light-weight Self-Supervised Neural Topic Model (SNTM) that learns a rich context by learning a topic representation jointly from three co-occurring words and a document that the triplet originates from. Jin Y et al [22] proposed a variational autoencoder (VAE) NTM model that jointly reconstructs the sentence and document word counts using combinations of bag-of-words (BoW) topical embeddings and pre-trained semantic embeddings. Zhao H et al [23] proposed to learn the topic distribution of a document by directly minimising its OT distance to the document's word distributions. Ma Z et al [24] proposed a novel topic model named Semantic-based Bidirectional Adversarial Neural Topic Model (SNTM), which introduces semantic information into Bidirectional Generative Adversarial Networks (BiGAN) by adding the word embedding and BiLSTM-Attention mechanism. Wang Y et al [25] developed a novel neural topic model, namely Layer-Assisted Neural Topic Model (LANTM), to enhance the topic represen- tation encoding by not only using text contents, but also the assisted network links. Yang Y et al [26] proposed TopNet, to leverage the recent advances in neural topic modeling to obtain high-quality skeleton words to complement the short input. Gupta P et al [27] proposed a neural topic mod- eling framework using multi-view embedding spaces: pretrained topic-embeddings, and pretrained word-embeddings (context-insensitive from Glove and context-sensitive from BERT models) jointly from one or many sources to improve topic quality and bet- ter deal with polysemy.

Despite the continual research of neural topic modeling, existing works do not yet fully exploit the useful information contained in documents. Our method is the first neural topic model that employs deep mutual information estimation [28]. It differs from the aforementioned neural topic models in the following main ways. (1) Un- like GSM, NVDM, and ProdLDA with a Gaussian or a logistic prior, we approximate the discrete topic assignment from a continuous distribution with Gumbel-Softmax[9], which can approximate categorical samples and whose parameters can be easily com- puted via the reparameterization trick. (2) Taking advantage of deep mutual informa- tion's capability of modeling the non-linear statistical dependence of the documents and the latent topics, our model estimates and maximizes mutual information between documents and topics to learn better representations of documents and topics.

*2.2. Deep Mutual Information Estimation*

One core objective of topic modeling is to learn useful topic representations. Similarly, deep mutual information estimation also aims to train an encoder for representation learning to maximize the mutual information (MI) [13] between its input and output. To the best of our knowledge, there is no work using Deep Mutual Information Estimation on topic modeling, hence we briefly survey recent research work on deep mutual information estimation. Belghazi et al [14] proposed MINE, a method to compute mutual information with neural network. Hjelm et al [28] proposed Deep InfoMax (DIM) to estimate and maximize the mutual information between input data, with global and local information, and its high-level representation with adversarial learning [11]. Yang et al [29] proposed a dual autoencoder network with mutual information estimation to learn the robust and discriminative latent representations for deep spectral clustering. Guo et al [30] proposed a method to learn disentangled representations, which incorporates deep mutual information estimation into the objective of cross-modal retrieval. Sanchez et al [31] proposed a method to learn the disentangled representations of images via deep mutual information estimation. Bachman et al [32] proposed a method, Augmented Multiscale Deep InfoMax (AMDIM), for representation learning based on maximizing mutual information between features from multiple views of a shared context and the latent representation. Qian et al [33] proposed VAE-MINE, which incorporates mutual information estimation (MINE) into variational autoencoder (VAE), to learn the latent representation. Zhou et al [34] proposed Text Matching with Deep Info Max (TIM), which maximizes the mutual information between the input and output with global and local information, for text matching. Benefiting from these work, we want to incorporate deep mutual information estimation into topic modeling to learn the topic representation more specific to given documents or sentences.

## 3. Neural Topic Model with Deep Mutual Information Estimation

The overall framework of our Neural Topic Modeling with Deep Mutual Information Estimation (NTM-DMIE) can be seen in Figure 1. Our framework consists
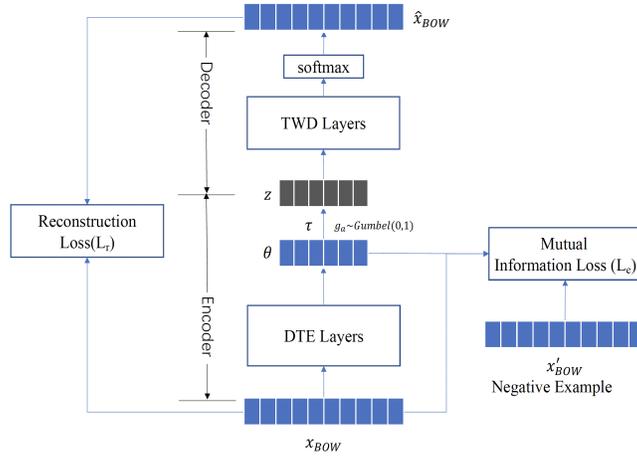
Figure 1: The overall framework of NTM-DMIE

of two main components, i.e., Document-Topic Encoder and Topic-Word Decoder. (1) The Document-Topic Encoder, which learns robust latent topic representations of documents with the documents themselves and their negative examples, simulates the document-topic distribution in LDA. To preserve the rich information of documents, we introduce mutual information estimation between the documents and their topic representations in the encoder. (2) The Topic-Word Decoder, which embeds the latent topic representations into document words, learns the topic-word distribution as in LDA. These two components are jointly optimized in a unified framework. The generative process of the common topic modeling is given as follows:

- Draw a topic distribution $\theta \sim Dirichlet(\alpha)$ ,

- For each word in the document, draw $w_n \sim Multinomial(\sigma(\beta\theta))$.

Where $\alpha$ is the parameter of the Dirichlet distribution, $w_n$ is the $n$-th word in the document, $\beta$ is the topic word distribution, $\sigma$ is the softmax function, and $\theta$ and $\beta$ are the parameters of the document-topic and topic-word distributions respectively, which are computed with neural network in our work.

Let $x_{BOW} \in V^{\mathbb{Z}_{\geq 0}}$ denote an input document in the bag-of-words representation, where $V$ is the vocabulary and $\mathbb{Z}_{\geq 0}$ denotes non-negative integers. $x'_{BOW}$ denotes

an negative example for learning the discriminative topic representation of $x_{BOW}$ and $\hat{x}_{BOW}$ is the document reconstructed by the autoencoder. $\hat{\theta}$ denotes the topic distribution of the documents and $\theta$ is the distribution after normalization. $z$ denotes their corresponding topic assignment based on the topic distribution $\theta$, and $K$ is the topic number and the dimension of $\theta$, $\hat{\theta}$ and $z$. Moreover, $\tau$, the temperature parameter of Gumbel-Softmax, and $g_a$, the base distribution, are used for sampling Gumbel-Softmax in the reparameterization trick. We use $x$ and $x_{BOW}$ interchangeably when there is no ambiguity.

Our goal is to train a document encoder with deep mutual estimation [14] to learn robust, discriminative topic representations for a given document together with its negative examples.

### 3.1. Document-Topic Encoder

Learning the topic representation of documents is the core part of neural topic modeling and a good topic representation can improve the quality of learning the topic word distribution. Mutual information, as shown in Eq (1), can model the essential correlation between two objects. We use it to measure the association between documents $x \in X$ and latent topics $z \in Z$ to learn robust topic representations, where $X = \{x_1, ..., x_n\}$ denote the input documents, $Z = \{z_1, ..., z_n\}$ denote their corresponding latent topic representations.

$$
\begin{aligned}
I(X, Z) &= \int \int p(z|x)p(x) \log \frac{p(z|x)}{p(z)} dx dz \\
&= KL(p(z|x)p(x)||p(z)p(x))
\end{aligned}
\tag{1}
$$

In Eq 1, $p(x)$ is the distribution of documents, $p(z|x)$ is the distribution of the latent topic and the marginal distribution of topic $p(z)$ is computed by $p(z) = \int p(z|x)p(x)dx$. The objective of our encoder is to maximize the mutual information as Eq (2). To make the learned topic more representative to the documents, the marginal distribution of latent topic $p(z)$ must obey the prior distribution of Dirichlet distribution $q(z)$.

$$
p(z|x) = \max_{W_E}\{I(X, Z)\}
\tag{2}
$$

where $W_E$ is the parameter of the encoder $E$.

Based on Eq (2), the objective of the encoder can be summarized as Eq (3).

$$\hat{p}(z|x) = \min_{W_E}\{ - \beta I(X, Z)$$
$$+ \gamma \int \int p(z|x)p(x) \log \frac{p(z|x)}{q(z)} dx dz\} \tag{3}$$

The first term of Eq (3) is the mutual information defined in Eq (1), and the second term is the KL divergence of the posterior $p(z|x)$ and the prior $q(z)$, which is beneficial to make the latent topic space more regular. Here, $\gamma$ is the smoothing parameter of mutual information and KL divergence. To resolve the problem of unbounded KL divergence in the calculation of mutual information (Eq (1)), the Jensen-Shannon (JS) divergence is employed for mutual information estimation. The objective thus can be rewritten as Eq (4).

$$\hat{p}(z|x) = \min_{W_E}\{ - \beta JS(p(z|x)p(x), p(z)p(x))$$
$$+ \gamma E_{x \sim p(x)}[KL(p(z|x)||q(z))]\} \tag{4}$$

Inspired by Nowozin et al [35], we also use the variational estimation of JS divergence with adversarial learning. The first term of Eq (4) can thus be optimized as Eq (5).

$$\max_{T}\{\mathbb{E}_{x,z \sim p(z|x)p(x)}[\log(\sigma(T(x, z)))]$$
$$+\mathbb{E}_{x,z \sim p(z)p(x)}[\log(1 - \sigma(T(x, z)))]\} \tag{5}$$

where $T(x, z)$ is given in Eq (6) and $\sigma(T(x, z))$ is a discriminator.

$$T(x, z) = \log \frac{2p(z|x)p(x)}{p(z|x)p(x) + p(z)p(x)} \tag{6}$$

To learn robust topic representations, the encoder uses a discriminator to differentiate the original document and its negative examples to measure the topic representation of the original document. In Eq (5), $\sigma(T(x, z))$ is used for discriminating the original document $x$ and its negative examples $\{x'\}$, where the negative examples are selected from the corpus with a deterministic strategy (which we will describe below). Eq (5) only considers the global mutual information: that between the whole document $x$ and
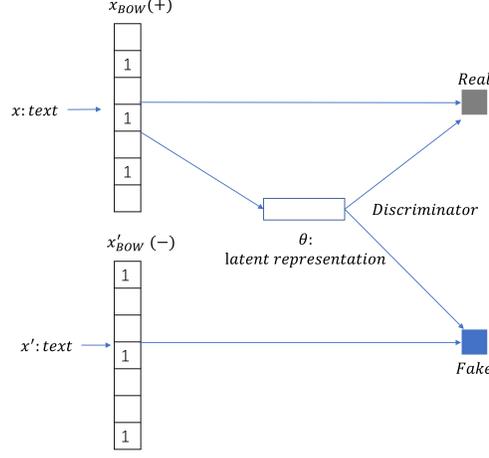
Figure 2: Global information: document and its negative examples. The latent topic representation is concatenated with the global "Bag of Words" representation. A 1 × 1 convolutional discriminator is used to score the 'real' document and its topic representation, while 'fake' is the randomly selected document with the learnt topic representation.

its topic representation $z$, as shown in Figure 2. The words in the document, i.e. local information, also play an important role in learning the topic representation of the document. Hence, we introduce the local mutual information to model the association between document words and the topic representation as shown in Figure 3. Similar to the global mutual information, negative examples are also used in the local mutual information. The loss function of the encoder is given as Eq (7), where $|x|$ is the number of words in the document, $x_i$ represents the $i$-th word in the document and $q(z)$ is the symmetrical Dirichlet distribution as the prior distribution.

$$
\begin{aligned}
L_e = &-\beta(\mathbb{E}_{x,z\sim p(z|x)p(x)}[\log(\sigma(T(x,z)))] \\
&+\mathbb{E}_{x,z\sim p(z)p(x)}[\log(1-\sigma(T(x,z)))]) \\
&-\frac{\beta}{|x|}\sum_i(\mathbb{E}_{x,z\sim p(z|x)p(x)}[\log(\sigma(T(x_i,z)))] \\
&+\mathbb{E}_{x,z\sim p(z)p(x)}[\log(1-\sigma(T(x_i,z)))]) \\
&+\gamma\mathbb{E}_{x\sim p(x)}[KL(p(z|x)||q(z))]
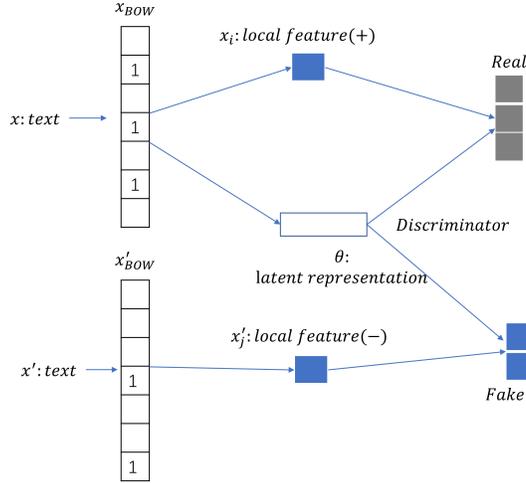\end{aligned}
\tag{7}
$$

10

Figure 3: Local information: document word and its negative examples. The latent topic representation is concatenated with the local word information. 'real' corresponds to the pair with word in the given document and the topic representation of the document, while 'fake' corresponds to the the pair with the randomly selected word and the given topic representation.

In Eq (7), $q(z)$ is the standard prior Dirichlet distribution and $p(x)$ is given by the corpus. The core part, $p(z|x)$, is the encoder of Fig 1, and $p(z)$ can be computed via $\int p(z|x)p(x)dx$. To model $p(z|x)$, note that the input $x$ is represented as $x_{BOW}$. Then, a feedforward neural network with two hidden layers (FC Layers) is utilized to embed $x_{BOW}$ into a hidden vector $\hat{\theta}$. To solve the problem of *posterior collapse* in VAE [36], a batch normalization layer [10] is added and the hidden vector $\hat{\theta}$ with normalization is transformed as $\theta$, which is the latent representation of the document $x$, and the parameter of the discrete distribution, i.e. $p(z|x)$. In the case of the Dirichlet distribution, we use the concrete distribution, a relaxation of discrete distribution via a Gumbel-Softmax [4], for sampling via the reparameterization trick.

### 3.2. Topic-Word Decoder

The decoder of NTM-DMIE, the TWD layers of Fig 1, is a linear transformation layer that maps $z$ for document $x$ to the predicated probability of words $\hat{x}$, i.e., the reconstructed document. The reconstruction loss mainly depends on two parts: the

11

distribution of latent topics and the generative performance of the decoder network. Our goal is to compute the word distribution of each topic $\tau$ via the decoder network. The reconstructed document can be obtained by Eq (8), where $\tau \in \mathbb{R}^{|V| \times K}$. Each column of $\tau$, $\tau_k$, represents the word distribution of the topic $k$.

$$p(\hat{x}|\tau, z) = softmax(\tau z + b) \tag{8}$$

The reconstruction loss is defined over the reconstructed document $\hat{x}$ and the original document $x$ as below.

$$L_r = ||\hat{x} - x||_F^2 \tag{9}$$

*3.3. Model Training*

The entire NTM-DMIE model is trained in an end-to-end manner. The overall loss function $L$ is a weighted sum of the mutual information loss and the reconstruction loss.

$$L = \mu * L_r + (1 - \mu) * L_e \tag{10}$$

## 4. Experimental Setup

*4.1. Dataset*

We evaluate the performance of NTM-DMIE on four public datasets, including two labeled datasets and two unlabeled datasets. **20 Newsgroups**[1]**.** is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups. Documents in the 20 newsgroups collections have class labels, so the dataset is used for evaluating the performance of text classification. **AG_News** is a collection of more than one million news articles which include four different groups: World, Sports, Business and Sci/Tec. News articles have been gathered from more than 2,000 news sources by ComeToMyHead in more than 1 year of activity. We selected its subset, which contains 12,760 articles, for the evaluation. **NY times** is a collection of news articles published between 1987 and 2007, and contains a wide range of topics,

---

[1]http://qwone.com/~jason/20Newsgroups/

such as sports, politics, education, etc. **Wikitext-103** is a collection of more than 100 million sentences, all extracted from Wikipedia's Good and Featured articles. It has been widely used in language modeling.

Table 1: Statistics of the four datasets.

| Dataset | Num of Docs | Average_Size | Num of Labels |
|---|---|---|---|
| 20NG | 19,999 | 220.5 | 20 |
| AG_News | 12,760 | 26.3 | 4 |
| NYTimes | 242,798 | 7.29 | - |
| Wikitext-103 | 307,807 | 176.2 | - |

We conducted the following common preprocessing steps: conversion into lowercase, word tokenization, lemmatization, and removal of stop words special characters. After preprocessing, the statistics of the four datasets are summarized in Table 1.

*4.2. Baselines*

We compared our NTM-DMIE model with the following state-of-the-art methods:

- **ProdLDA** [6] is an Autocoder-based topic model that constructs a Laplace approximation to the Dirichlet prior.

- **GSM** [3] is a Gaussian Softmax topic model parameterized with neural networks.

- **NTM** [17] is a neural topic model which incorporates a topic coherence objective.

- **Scholar** [16] is a supervised neural topic model that allows for metadata to appear as either a covariate or a predicted variable in the model structure.

- **Gaussian-BAT** [5] models topics with the Dirichlet prior and builds a two-way transformation between document-topic distribution and document-word distribution via bidirectional adversarial training.

### 4.3. Implementation Details

The hyperparameters of ProdLDA, GSM, NTM, Scholar and Gaussian-BAT are set according to the best hyperparameters reported in their original papers. Topic number was set as 10, 20 and 50 for all the four datasets to test the ability of text clustering and the quality of topics in our model compared with the baselines. During model training, we used the Adam optimizer with a learning rate of $10E - 4$ on the all datasets. For the weight of the loss function of the encoder given in Eq 7, it was set as $\beta = \gamma = 1$. For $\mu$ of the overall loss function in Eq (10), it was empirically set as $\mu = 0.4$. All the experiments were conducted for 100 epochs with batch size 128. And all models are implemented by PyTorch with a single Nvidia GTX 1080Ti graphic card, running for four times. In our model, negative examples are used for learning robust topic representations of the document. In our experiment, we use different strategies for selecting negative examples: random selection and similarity-based selection. Hence, our model have two variants, NTM-DMIE (random) and NTM-DMIE (similarity).

### 4.4. Evaluation Metrics

We compare model performance on *topic coherence* and *topic uniqueness* to evaluate the quality of topics. We also perform text clustering to measure the reconstruction ability of latent features, for which we use accuracy.

**Topic Coherence (NPMI)** [37] Topic coherence indicates that the words in a topic should be as coherent as possible. For this we use the widely-used metric Normalized Pointwise Mutual Information (NPMI), which assumes coherent words should co-occur within a certain distance. Given the top $M$ topic words ordered by their probabilities, the NPMI score of the topic can be calculated as follows:

$$NPMI = \frac{1}{M} \sum_{w_i, w_j} \frac{\log \frac{p(w_i, w_j) + \epsilon}{p(w_i)p(w_j)}}{-\log(p(w_i, w_j) + \epsilon)} \tag{11}$$

where $p(w_i)$ is the probability of word $w_i$, $p(w_i, w_j)$ is the co-occurrence probability of $w_i$, $w_j$ within a window in the reference corpus and $\epsilon$ is used to avoid division by zero.

**Topic Uniqueness (TU)** [38] measures the diversity of a set of topics, and can be used to determine how distinguished the topics are from each other. Given the top $M$

words for each of the $K$ topics, TU for topic $k = 1, \ldots, K$ can be defined as follows:

$$TU(k) = \frac{1}{M} \sum_{i=1}^{M} \frac{1}{cnt(i, k)} \tag{12}$$

where $cnt(i, k)$ is the occurrence count of the $i$-th top word in topic $k$ in the top words across all topics. The range of TU value is between $1/K$ and $1$. A higher TU value indicates that fewer words are repeated across topics, thus the produced topics are more diverse.

**Clustering Accuracy (ACC)** [5] measures the effectiveness of learned topics on document clustering, in which the learned topic distributions are used as features for clustering. Model performance is evaluated by accuracy (ACC) as follows:

$$ACC = \max_{mapping} \frac{\sum_{i=1}^{N_t} \mathbb{1}(l_i = mapping(c_i))}{N_t} \tag{13}$$

where $N_t$ is the number of documents in the test set, $\mathbb{1}()$ is the indicator function, $l_i$ is the ground-truth label of the $i$-th document, $c_i$ is the cluster assignment of the $i$-th document, and $mapping$ ranges over all possible one-to-one mappings between labels and clusters. A higher ACC score means that the model is more likely to capture features that are representative of the given corpus.

## 5. Results and Analysis

In this section we present a comprehensive empirical evaluation on our proposed method. Our experiment evaluate the metrics, topic coherence for evaluating the quality of the learnt topics, topic uniqueness for the diversity of the topics, and the text clustering for evaluating the representative ability of the topic representations. Finally, the qualitative analysis of the learnt topics are given below. The source code of our experiment is given in the link https://github.com/Asuper-code/NTM-DMIE.

### 5.1. Topic Coherence and Topic Uniqueness

Table 2 presents the results on topic coherence (measured by NPMI) and topic uniqueness with the number of topics set to 20. Results with the number of topics set to 10 , 50 and 100 can be found in Figure 4, Figure 5 and Figure 6, where our

model exhibits similar superiority over the compared state-of-the-art models. Among Figure 4 , 5 and 6, we also give the result of NMPI and TU of the two variants, NTM-DMIE(random) and NTM-DMIE(similarity). And the result show that NTM-DMIE(similarity) performed better than NTM-DMIE(random). The detailed analysis of the experiment about selection of negative examples will be discussed in section 5.4.

Table 2: Topic quality evaluation for 20 topics. The numbers in each cell are NPMI/TU, showing the 95% confidence interval. Both NPMI and TU values are the higher the better.

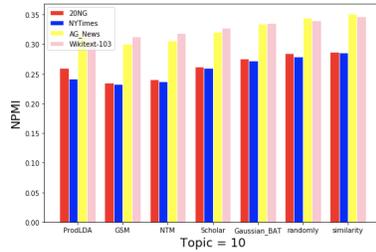| Method | 20NG | | NYTimes | | AG_News | | Wikitext-103 | |
|---|---|---|---|---|---|---|---|---|
| | NPMI | TU | NPMI | TU | NPMI | TU | NPMI | TU |
| ProdLDA | 0.267±0.002 | 0.58±0.01 | 0.319±0.001 | 0.67±0.03 | 0.247±0.003 | 0.65±0.02 | 0.325±0.001 | 0.69±0.02 |
| GSM | 0.243±0.001 | 0.65±0.02 | 0.303±0.002 | 0.79±0.02 | 0.239±0.003 | 0.70±0.03 | 0.317±0.001 | 0.71±0.02 |
| NTM | 0.252±0.003 | 0.62±0.02 | 0.310±0.003 | 0.88±0.01 | 0.245±0.001 | 0.67±0.01 | 0.320±0.003 | 0.79±0.02 |
| Scholar | 0.273±0.004 | 0.73±0.01 | 0.328±0.001 | 0.89±0.01 | 0.265±0.002 | 0.78±0.02 | 0.333±0.002 | 0.85±0.01 |
| Gaussian-BAT | 0.285±0.001 | 0.85±0.01 | 0.344±0.002 | 0.95±0.01 | 0.274±0.001 | 0.86±0.01 | 0.338±0.003 | 0.93±0.02 |
| NTM-DMIE(random) | 0.294±0.001 | 0.88±0.03 | 0.350±0.002 | 0.94±0.01 | 0.281±0.004 | 0.90±0.02 | 0.341±0.003 | 0.91±0.01 |
| NTM-DMIE(similarity) | **0.298**±0.002 | **0.93**±0.01 | **0.357**±0.003 | **0.96**±0.01 | **0.285**±0.002 | **0.94**±0.02 | **0.347**±0.001 | **0.94**±0.02 |

In Table 2, in terms of topic coherence, NTM-DMIE achieves the highest NPMI scores on all four datasets. Moreover, the NPMI of our model is substantially higher than all the baselines. This result demonstrates that our model can obtain more coherent topic words than the state-of-the-art baselines.

As for topic uniqueness, NTM-DMIE also achieves the highest scores than all baselines on all the four datasets. This result indicates that our model can obtain topics with less repetition better than the baseline models.
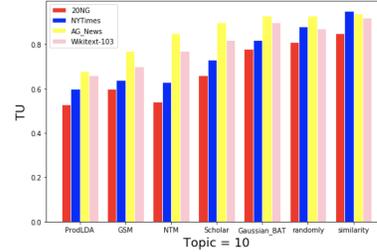
## 5.2. Ablation Study

We conduct an ablation study to examine the effectiveness of the local information and global information components in our framework. Experiments are conducted on all four datasets with the topic number set to 20.

In Table 3 and Table 4 (for labeled and unlabeled datasets respectively), the full NTM-DMIE model outperforms both variants significantly. It attests to the effectiveness of both local and global mutual information in our framework. On NPMI, we can see that NTM-DMIE with local information only performs better than with global
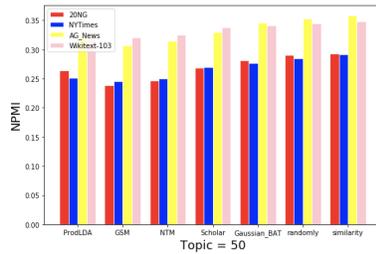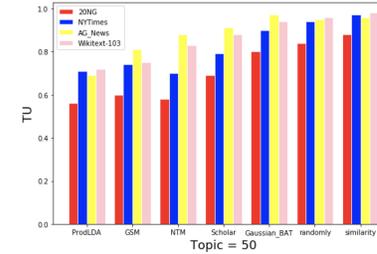
(a) NPMI results.



(b) TU results.

Figure 4: NPMI and TU Performance when Topic=10 on the four datasets.
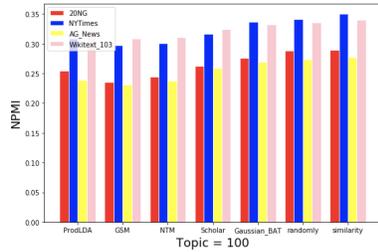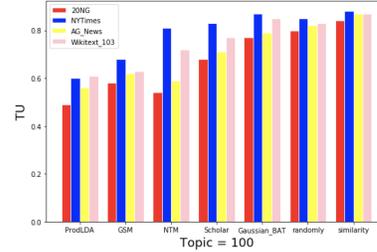


(a) NPMI results.



(b) TU results.

Figure 5: NPMI and TU Performance when Topic=50 on the four datasets.

information only, which can be attributed to the fact that local information helps the model capture more specific and high-quality features than global information. A similar observation can be made for TU, where local information performs better except a slightly worse TU score on the NYTimes dataset.

In the Table 3 and Table 4, it show the result of the variants of our model, which are the one without the global information, the one without local information, and the one without both the global information and local information. Based on the result, we can find that the NPMI and TU performance of the model obviously decline when our model remove the global information or the local information.

(a) NPMI results.



(b) TU results.

Figure 6: NPMI and TU Performance when Topic=100 on the four datasets.

Table 3: Ablation study of labeled datasets.

| Method | 20NG | | AG_News | |
|---|---|---|---|---|
| | NPMI | TU | NPMI | TU |
| NTM-DMIE | **0.298** | **0.93** | **0.285** | **0.94** |
| w/ gloal information only | 0.243 | 0.79 | 0.241 | 0.76 |
| w/ local information only | 0.252 | 0.82 | 0.247 | 0.79 |
| w/ both | 0.223 | 0.66 | 0.225 | 0.61 |

Table 4: Ablation study of unlabeled datasets.

| Method | NYTimes | | Wikitext-103 | |
|---|---|---|---|---|
| | NPMI | TU | NPMI | TU |
| NTM-DMIE | **0.357** | **0.96** | **0.347** | **0.94** |
| w/ gloal information only | 0.319 | 0.88 | 0.323 | 0.80 |
| w/ local information only | 0.321 | 0.89 | 0.329 | 0.84 |
| w/ both | 0.300 | 0.71 | 0.307 | 0.63 |

## 5.3. Text Clustering Performance

We evaluate the text clustering performance of the topic models learned by each method on the labeled datasets 20NG and AG_News, and Table 5 presents the overall

18

results for all the models where the topic number is set to 20, which is set based on the number of classes in the labeled datasets.

In Table 5, features captured by NTM-DMIE obtain the highest accuracy on text clustering on both datasets. We attribute this result to the mutual information framework of NTM-DMIE, which helps reconstruct the input and improves the quality of topics captured by the model.

Table 5: Text clustering performance of different methods on the 20Newsgroups dataset, where topic number is set to be 20. Higher value indicates better performance.

| Model | 20NG | AG_News |
|---|---|---|
| ProdLDA | 32.5% | 78.2% |
| GSM | 31.7% | 80.1% |
| NTM | 33.9% | 82.3% |
| Scholar | 35.7% | 82.9% |
| Gaussian_BAT | 39.9% | 84.5% |
| NTM-DMIE (random) | 43.5% | 85.1% |
| NTM-DMIE (similarity) | **45.2%** | **85.9%** |

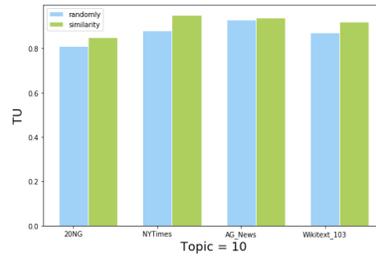### 5.4. Selection Strategies of Negative Examples

We employ a discriminator in the document-topic encoder to distinguish a document from its negative examples for learning robust topic representations of the document. Here we examine the effect of different strategies for selecting negative examples: random selection and similarity-based selection.

Random selection, as the name suggests, randomly chooses several negative examples for a given example. In contrast, the similarity-based selection strategy chooses the most dissimilar documents as negative examples. The results of NPMI and TU on all four datasets are shown in Table 2, with the topic number set to 20.

Table 2 shows that NTM-DMIE performs better with the similarity-based selection strategy than with the random strategy. We can also see the trend in Figure 7, Figure 8, Figure 9 and Figure 10. This makes intuitive sense as a randomly chosen "negative" example may not actually be sufficiently different from the given document, thus providing the discriminator with noisy training signals.
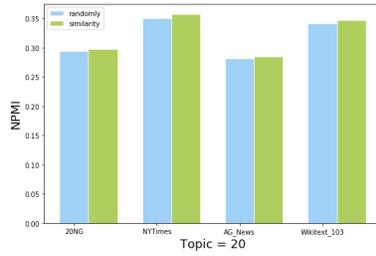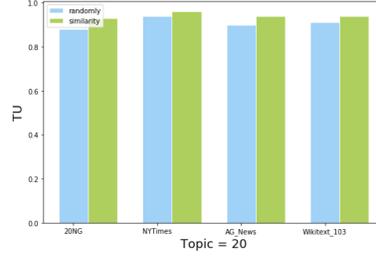
(a) Topic = 10_NPMI.

(b) Topic = 10_TU.

Figure 7: NPMI and TU with different ways of negative sample choosing when Topic = 10.
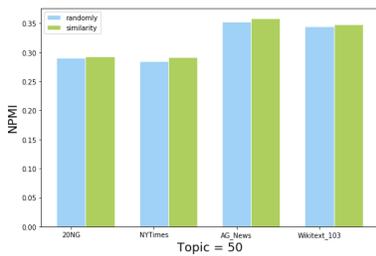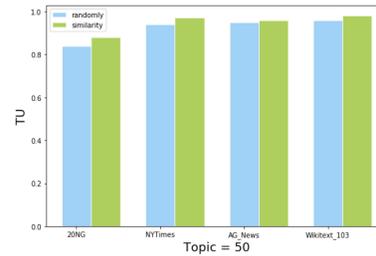


(a) Topic = 20_NPMI.

(b) Topic = 20_TU.

Figure 8: NPMI and TU with different ways of negative sample choosing when Topic = 20.
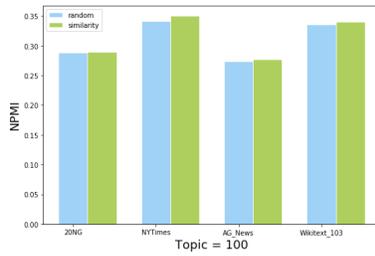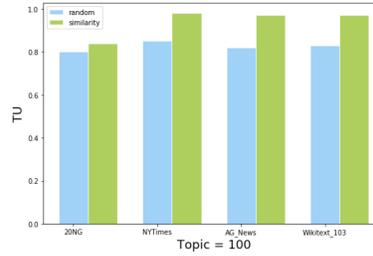


(a) Topic = 50_NPMI.

(b) Topic = 50_TU.

Figure 9: NPMI and TU with different ways of negative sample choosing when Topic = 50.

(a) Topic = 100_NPMI.

(b) Topic = 100_TU.

Figure 10: NPMI and TU with different ways of negative sample choosing when Topic = 100.

## 5.5. Qualitative Analysis

In order to more closely examine the accuracy of the topics and the corresponding keywords captured from each topic model, we compare NTM-DMIE with all the baseline models on the 20NG and NYTimes datasets, with the topic number set to 10.

Table 6: Topic-word example of Gaussain_BAT and NTM-DMIE on 20Newsgroups dataset with topic = 10

| | Topic:Politics and Military |
| --- | --- |
| ProdLDA | gun people law problem government think kill death say drug |
| GSM | gun government tax pay law weapon rights people firearm president |
| NTM | people gun government law rights article drug death clinton kill |
| Scholar | gun article run people law government say death problem weapon |
| Gaussian_BAT | gun problem people use government law article drug write kill |
| NTM-DMIE | gun government tax pay law weapon rights people firearm president |

| | Topic:Transportation and daily life |
| --- | --- |
| ProdLDA | car bike ride get dod run write go put like |
| GSM | car bike problem use dod work ride driver set get |
| NTM | car bike problem get use ride dod driver work drive |
| Scholar | car bike drive dod ride get engine driver use buy |
| Gaussian_BAT | car bike ball use dod work ride driver set get |
| NTM-DMIE | car bike drive ride engine dod buy sell like bmw |

| | Topic:Security |
| --- | --- |
| ProdLDA | key chip use problem encryption bit system work clipper government |
| GSM | key chip encryption use run clipper system government problem bit |
| NTM | key chip encryption clipper get work government use system bit |
| Scholar | key use chip problem encryption work system set machine bit |
| Gaussian_BAT | key chip use problem encryption bit system work clipper government |
| NTM-DMIE | key chip clipper encryption escrow nsa government system secure need |

| | Topic:Communication |
| --- | --- |
| ProdLDA | drive card controller disk monitor thanks port pc driver system |
| GSM | card drive controller disk monitor use port controller driver window |
| NTM | drive card use problem disk work sale offer monitor machine |
| Scholar | drive card disk windows use run sale problem monitor pc |
| Gaussian_BAT | drive card use disk monitor sale controller reply printer window |
| NTM-DMIE | card thanks please use window advance email port reply display |

Table 7: Topic-word example of Gaussain_BAT and NTM-DMIE on NYtimes dataset with topic = 10

| Topic:Business | |
|---|---|
| ProdLDA | company money percent pay state cost bill industry buy chief |
| GSM | company executive business chief sell buy price pay share |
| NTM | company president executive business sell market sale chief share buy |
| Scholar | company percent market price business sell sale industry buy executive |
| Gaussian_BAT | company computer system technology program number service price product information |
| NTM-DMIE | company percent market price business sell sale pay buy industry |

| Topic:Literature and Art | |
|---|---|
| ProdLDA | music play art film book world write performance audience present |
| GSM | music art play book film write world television performance life |
| NTM | music program art book director write production company feature performance |
| Scholar | play music film performance movie audience young write book character |
| Gaussian_BAT | play music write book life world art film character director |
| NTM-DMIE | music play art film movie book director write performance audience |

| Topic:Politics and Law | |
|---|---|
| ProdLDA | case police charge official court law judge yesterday officer state |
| GSM | case police charge officer law man life feel official rule |
| NTM | law charge case court police judge recieve graduate official rule |
| Scholar | state law issue official vote public case court member judge |
| Gaussian_BAT | case police charge official court law worker judge public officer |
| NTM-DMIE | case law official court state charge issue judge rule police |

| Topic:People and Life | |
|---|---|
| ProdLDA | man life feel thing tell ask woman friend son student |
| GSM | life man woman young thing write son love tell friend |
| NTM | man police woman life death son daughter family child kill |
| Scholar | man life father mother family student child graduate director school |
| Gaussian_BAT | man life woman father mother child young family feel friend |
| NTM-DMIE | shcool child father family mrs son mother student graduate daughter |

Table 6 shows that nearly all the models can extract the four topics, namely Politics and Military; Transportation and daily life; Security; and Communication. However, NTM-DMIE can mine more keywords that are more closely related to the specific topics than the other models, while other baselines may mine some words that do not

belong to the corresponding topic. For example, in topic Politics and Military, Gaussain_BAT captures **"article"** and **"write"**, which are not so closely related to politics or military while NTM-DMIE can capture words like **"firearm"** and **"president"** which are closely related to the topic. Table 7 also shows that nearly all the models can extract the four topics, namely Business; Literature and Art; Politics and Law; and People and Life. However, similarly, NTM-DMIE performs better than the other five baselines.

## 6. Conclusion

In this paper, we have proposed a framework to incorporate deep mutual information into neural topic modeling. Our framework maximizes the mutual information between the input documents and their latent topic representations. We capture mutual information on the global and local levels to preserve the rich information of the documents and words into their topic representations. A discriminator is also employed to discriminate a document from its negative examples for learning robust topic representations. Experiments on four public datasets show that our model outperforms state-of-the-art neural topic models on the metrics of topic coherence and topic uniqueness. A further experiment on text clustering demonstrates the quality of the learned topics in downstream tasks. In future work, we will investigate self-supervised learning approaches to extend our model.

## 7. Acknowledgements

## References

[1] D. M. Blei, A. Y. Ng, M. I. Jordan, Latent dirichlet allocation, the Journal of machine Learning research 3 (Jan) (2003) 993–1022.

[2] Y. Miao, L. Yu, P. Blunsom, Neural variational inference for text processing, in: Proceedings of the 33nd International Conference on Machine Learning, Vol. 48, JMLR.org, 2016, pp. 1727–1736.

[3] Y. Miao, E. Grefenstette, P. Blunsom, Discovering discrete latent topics with neural variational inference, in: Proceedings of the 34th International Conference on Machine Learning, Vol. 70, PMLR, 2017, pp. 2410–2419.

[4] B. Esmaeili, H. Huang, B. Wallace, J.-W. van de Meent, Structured neural topic models for reviews, in: The 22nd International Conference on Artificial Intelligence and Statistics, Vol. 89, 2019, pp. 3429–3439.

[5] R. Wang, X. Hu, D. Zhou, Y. He, Y. Xiong, C. Ye, H. Xu, Neural topic modeling with bidirectional adversarial training, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, 2020, pp. 340–350.

[6] A. Srivastava, C. A. Sutton, Autoencoding variational inference for topic models, in: 5th International Conference on Learning Representations, OpenReview.net, 2017.

[7] D. Card, C. Tan, N. A. Smith, A neural framework for generalized topic models, CoRR abs/1705.09296.

[8] L. Gui, J. Leng, G. Pergola, Y. Zhou, R. Xu, Y. He, Neural topic model with reinforcement learning, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, Association for Computational Linguistics, 2019, pp. 3476–3481.

[9] E. Jang, S. Gu, B. Poole, Categorical reparameterization with gumbel-softmax, in: 5th International Conference on Learning Representations, OpenReview.net, 2017.

[10] S. Burkhardt, S. Kramer, Decoupling sparsity and smoothness in the dirichlet variational autoencoder topic model, Journal of Machine Learning Research 20 (2019) 131:1–131:27.

[11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Advances in Neural Information Processing Systems 27, 2014, pp. 2672–2680.

[12] R. Wang, D. Zhou, Y. He, ATM: adversarial-neural topic model, Information Processing & Management 56 (6).

[13] P. Viola, W. M. Wells III, Alignment by maximization of mutual information, International journal of computer vision 24 (2) (1997) 137–154.

[14] M. I. Belghazi, A. Baratin, S. Rajeswar, S. Ozair, Y. Bengio, R. D. Hjelm, A. C. Courville, Mutual information neural estimation, in: Proceedings of the 35th International Conference on Machine Learning, Vol. 80, PMLR, 2018, pp. 530–539.

[15] D. P. Kingma, M. Welling, Auto-encoding variational bayes, in: 2nd International Conference on Learning Representations, Conference Track Proceedings, 2014.

[16] D. Card, C. Tan, N. A. Smith, Neural models for documents with metadata, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, 2018, pp. 2031–2040.

[17] R. Ding, R. Nallapati, B. Xiang, Coherence-aware neural topic modeling, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, 2018, pp. 830–836.

[18] P. Gupta, Y. Chaudhary, F. Buettner, H. Schütze, Document informed neural autoregressive topic models with distributional prior, in: The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI Press, 2019, pp. 6505–6512.

[19] X. Zhao, D. Wang, Z. Zhao, W. Liu, C. Lu, F. Zhuang, A neural topic model with word vectors and entity vectors for short texts, Information Processing & Management 58 (2) (2021) 102455.

[20] M. Panwar, S. Shailabh, M. Aggarwal, B. Krishnamurthy, Tan-ntm: Topic attention networks for neural topic modeling, arXiv preprint arXiv:2012.01524.

[21] S. A. Bahrainian, M. Jaggi, C. Eickhoff, Self-supervised neural topic modeling, in: Findings of the Association for Computational Linguistics: EMNLP 2021, 2021, pp. 3341–3350.

[22] Y. Jin, H. Zhao, M. Liu, L. Du, W. Buntine, Neural attention-aware hierarchical topic model, arXiv preprint arXiv:2110.07161.

[23] H. Zhao, D. Phung, V. Huynh, T. Le, W. Buntine, Neural topic model via optimal transport, in: International Conference on Learning Representations, 2020.

[24] Z. Ma, J. Lu, J. Feng, Y. Zhang, W. Wu, Semantic-based bidirectional adversarial neural topic model, in: 2021 IEEE 33rd International Conference on Tools with Artificial Intelligence (ICTAI), IEEE, 2021, pp. 376–380.

[25] Y. Wang, X. Li, J. Ouyang, Layer-assisted neural topic modeling over document networks, in: International Joint Conference on Artificial Intelligence, 2021, pp. 3148–3154.

[26] Y. Yang, B. Pan, D. Cai, H. Sun, Topnet: Learning from neural topic model to generate long stories, in: Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, 2021, pp. 1997–2005.

[27] P. Gupta, Y. Chaudhary, H. Schütze, Multi-source neural topic modeling in multi-view embedding spaces, arXiv preprint arXiv:2104.08551.

[28] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, Y. Bengio, Learning deep representations by mutual information estimation and maximization, in: 7th International Conference on Learning Representations, OpenReview.net, 2019.

[29] X. Yang, C. Deng, F. Zheng, J. Yan, W. Liu, Deep spectral clustering using dual autoencoder network, in: IEEE Conference on Computer Vision and Pattern Recognition, Computer Vision Foundation / IEEE, 2019, pp. 4066–4075.

[30] W. Guo, H. Huang, X. Kong, R. He, Learning disentangled representation for cross-modal retrieval with deep mutual information estimation, in: Proceedings of the 27th ACM International Conference on Multimedia, ACM, 2019, pp. 1712–1720.

[31] E. H. Sanchez, M. Serrurier, M. Ortner, Learning disentangled representations via mutual information estimation, Computer Vision - ECCV 2020 - 16th European Conference 12367 (2020) 205–221.

[32] P. Bachman, R. D. Hjelm, W. Buchwalter, Learning representations by maximizing mutual information across views, in: Advances in Neural Information Processing Systems 32, 2019, pp. 15509–15519.

[33] D. Qian, W. K. Cheung, Enhancing variational autoencoders with mutual information neural estimation for text generation, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, Association for Computational Linguistics, 2019, pp. 4045–4055.

[34] X. Zhou, C. Li, J. Bu, C. Yao, K. Shi, Z. Yu, Z. Yu, Matching text with deep mutual information estimation, CoRR abs/2003.11521.

[35] S. Nowozin, B. Cseke, R. Tomioka, f-gan: Training generative neural samplers using variational divergence minimization, in: Advances in Neural Information Processing Systems 29, 2016, pp. 271–279.

[36] Q. Zhu, W. Bi, X. Liu, X. Ma, X. Li, D. Wu, A batch normalized inference network keeps the KL vanishing away, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, 2020, pp. 2636–2649.

[37] N. Aletras, M. Stevenson, Evaluating topic coherence using distributional seman-
tics, in: Proceedings of the 10th International Conference on Computational Se-
mantics, The Association for Computer Linguistics, 2013, pp. 13–22.

[38] F. Nan, R. Ding, R. Nallapati, B. Xiang, Topic modeling with wasserstein autoen-
coders, in: Proceedings of the 57th Conference of the Association for Compu-
tational Linguistics, Association for Computational Linguistics, 2019, pp. 6345–
6381.