

Correlated Topic Model with Transformer Embeddings  
トランスフォーマーの埋め込みによる相関トピックモデル

by

Chun Wa Leung

梁 俊華

A Master Thesis

修士論文

Submitted to

the Graduate School of the University of Tokyo

on January 19, 2022

in Partial Fulfillment of the Requirements

for the Degree of Master of Information Science and

Technology

in Computer Science

Thesis Supervisor: Akihiko Takano 高野 明彦

Professor of Computer Science

## ABSTRACT

Topic modeling is one of the most common information retrieval task in natural language processing. In particular, Correlated topic model(CTM) is a topic model which captures the correlation between topics associated. However, such a classic statistical approach was not able to capture positional information from sequential input. At that point, traditional topic models may perform poorly in generating words from large number of topics. In this research, we introduce Correlated Topic Model with Transformer embeddings, a generative model where combine the advantage of using positional information of words and topic correlation. Specifically, transformer embedding maps topic words into latent space and further assign to its assigned topic. We attempted to add a covariance prior to the topic model, LKJ correlation prior to logistic-normal distribution, which aims to fit the correlation information from the data. In addition, we extended our model to handle time-series data integrated with Gaussian Process latent variable model(GPLVM), which also capturing temporal information from words occurrence of documents over time. The model was optimized using Stochastic Variational Inference (SVI), allows handling massive data sets with mini-batching. As compared to empirical results from experiments, our approach performs a better fit of the data than existing generative topic model and exhibit a better capability in obtaining high quality topics.

## Acknowledgements

I would like to express my thanks of gratitude to my guiding professor Akihiko Takano as well as professor Takeshi Abekawa me along the master thesis. It is my pleasure to receive many valuable suggestions on thesis writing which helped me a lot in doing research.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Applications . . . . .	2
1.3	Literature Review . . . . .	3
1.4	Objective and Outline . . . . .	3
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Problem definition . . . . .	5
2.2	Bag-of-word assumption . . . . .	5
2.3	LKJ Correlation Distribution . . . . .	6
2.4	Topic Models . . . . .	6
2.4.1	Correlated Topic Model . . . . .	6
2.4.2	Embedded Topic Model . . . . .	8
2.5	Representation learning . . . . .	8
2.5.1	Word Embedding . . . . .	8
2.5.2	Transformer . . . . .	8
2.6	Time Series model . . . . .	9
2.6.1	Gaussian Process (GP) . . . . .	9
2.6.2	Gaussian Process Latent Variable Model (GPLVM) . . . . .	10
2.6.3	Dynamic Topic Model . . . . .	11
2.7	Posterior Inference . . . . .	11
2.7.1	Variational Inference . . . . .	12
2.7.2	Stochastic Variational Inference . . . . .	12
2.7.3	Collapsing Parameters . . . . .	12
2.7.4	Autoencoding Variational Bayes (AEVB) . . . . .	12
2.7.5	Reparameterization trick . . . . .	13
2.8	Evaluation metrics . . . . .	13
2.8.1	Perplexity . . . . .	13
2.8.2	Topic Coherence . . . . .	13
2.8.3	Topic Diversity . . . . .	13
<b>3</b>	<b>Transformer embedding with Correlated Topic Model</b>	<b>14</b>
3.1	Introduction . . . . .	14
3.2	Related works . . . . .	15
3.3	Model description . . . . .	15
3.3.1	LKJ Correlation prior . . . . .	16
3.3.2	Transformer Embeddings . . . . .	16
3.3.3	Marginal Likelihood . . . . .	16
3.3.4	Joint Distribution . . . . .	17
3.4	Inference and Estimation . . . . .	17
3.4.1	Variational distribution . . . . .	17
3.4.2	Evidence Lower Bound(ELBO) . . . . .	17

3.4.3	Optimization step . . . . .	19
3.5	Results & Evaluations . . . . .	20
3.5.1	Experiment Testing . . . . .	20
3.5.2	Dataset . . . . .	20
3.5.3	Models . . . . .	21
3.5.4	Algorithmic Settings . . . . .	21
3.5.5	Quantitative Result . . . . .	22
3.5.6	Training . . . . .	23
3.5.7	Qualitative Result . . . . .	23
3.5.8	Visualization . . . . .	23
3.5.9	Discussion . . . . .	26
<b>4</b>	<b>Times Series Topic Retrieval with TECTM</b>	<b>28</b>
4.1	Introduction . . . . .	28
4.2	Related works . . . . .	29
4.3	Model description . . . . .	29
4.3.1	Generative Model . . . . .	30
4.3.2	Joint Distribution . . . . .	30
4.4	Inference and Estimation . . . . .	32
4.4.1	Variational Distribution . . . . .	32
4.4.2	Evidence lower bound (ELBO) . . . . .	32
4.5	Experiment and results . . . . .	36
4.5.1	Experiment settings . . . . .	36
4.5.2	Results . . . . .	37
4.5.3	Discussion . . . . .	39
<b>5</b>	<b>Conclusions</b>	<b>42</b>
	<b>References</b>	<b>43</b>

## List of Figures

2.1	Graphical representation for CTM . . . . .	7
3.1	Graphical model for TECTM . . . . .	15
3.2	20Newsgroups (k=20) . . . . .	24
3.3	20Newsgroups (k=50) . . . . .	24
3.4	t-SNE Visualization for our model #Topic:20 . . . . .	25
3.5	t-SNE visualization for ETM #Topic:20 . . . . .	25
3.6	Visualization #Topics:50 . . . . .	27
4.1	t-SNE visualization for documents labeled by topic(UN debates) .	38
4.2	Word trend for top-6 topics (UN debates) . . . . .	39
4.3	Topic trend for top-6 topics (UN debate) . . . . .	39
4.4	Word trend for top-6 topics (NeurIPS dataset) . . . . .	40
4.5	Topic trend for top-6 topics (NeurIPS dataset) . . . . .	41

## List of Tables

3.1	Result for Reuters-21578 dataset . . . . .	22
3.2	Result for 20Newsgroups dataset . . . . .	22
3.3	Top-9 words for each topic from 5 topics selected . . . . .	26
4.1	Result on UN debates dataset, k=30 . . . . .	37
4.2	Result on NeurIPS dataset (1987-2019), k=30 . . . . .	37
4.3	Word trend in topic reinforcement learning (5 years interval) . . .	40

## List of Algorithms

1	Generative Process for CTM . . . . .	7
2	Generative Process for ETM . . . . .	8
3	Generative Process for DTM . . . . .	11
4	Generative Process for TECTM . . . . .	16
5	Training on TECTM . . . . .	20
6	Generative Process for DTECTM . . . . .	30
7	Training on DTECTM . . . . .	35



# Chapter 1

## Introduction


Today, myriads of terabyte data is one of the crucial challenge for scientific researcher. Information Retrieval become increasingly important for building useful information from massive data set. Specifically, topic modeling is one of the most popular technique for extracting key point ideas and exploring documents. Specifically, correlated topic models (CTM) make use of the correlation between topics and deliver a better result. In the research, we would like to explore the application one of the topic modeling techniques and try to improve their performance. In section 1.1, we provide a brief introduction to the existing algorithm and covers the background of it. Then, section 1.2 will discuss the current application and state-of-art improvement on topic modeling advances. Following that, section 1.3 will elaborate the conduct the research and the general direction. And section 1.4 will explain the way the proposal model to be evaluated and compare to the existing model.

### 1.1 Motivation

Topic modeling is one of the most exciting domain in Information Retrieval (IR). It can be extended to accomplish versatile range of IR and data mining tasks. For instance, one of the topic model: Latent Dirichlet Allocation (LDA), was proposed and examined its capability on extracting latent topics and output keywords suggestion for each topic. As result, LDA has been implemented into various area of applications. However, There are several problems that LDA could not handle well.

**Computational complexity** Generally, LDA acquires to compute the posterior distribution for inference, which is relatively expensive to obtain an exact solution. In the same way, its variant, correlated topic model (CTM) requires to calculate the covariance matrix specifically, which makes it not feasible come into practical application.

**Statistical Laws** In particular, LDA does not take empirical statistical laws observed in text into account. For example, LDA's prior does not dependent on Zipf's Law or Heap's Law, which may not collaborate well with natural document text data. Similarly, Moody[32] proposed lda2vec, which exploit the meta-information of each document and evaluate their correlation between documents.

**Correlation information** Correlation information can be useful to identify topics. For instance, ho and soccer are correlated but uncorrelated other

topic like space and religion. Such intuition could help topic model to exploit those information.

**Bag-of-word assumption** Typical topic model like Nonnegative Matrix Factorization (NMF) [26] and Latent Dirichlet Allocation (LDA) [9] do not consider positional information from the document set. This lead to the drawbacks of those models may not make good prediction on the topic words due to the limitations. For most of the NLP tasks, it is very common to let the model learning context by

**Transfer Learning** Due to the prevalence of deep neural network in recent years, Transform Learning has became a hot topic in research. The aim for transfer learning is to make find a way to reduce computational cost and improve re-usability of machine learning models. In the ascendant of powerful accelerator such as GPU and more memory, we are able to build more complex architecture and boost up computation time. Specifically, Transformer has been one of the most used NLP architecture. Number of variants have been built due to its success, such as, BERT[12], ROBERTA[28], and ELMO[35], etc.

## 1.2 Applications

Topic model are one of the crucial tasks in discovering hidden topic from document collections. The success of LDA make able it does not limit to topic modeling task. Graber[10] gave a verbose survey on topic model applications. Many tasks have been applied with the model, for examples below,

**Feature Extraction** For number of  $n$  topics, LDA can accomplish the task cluster them and extract a set of corpus with  $k$  terms which can represent each topic most and uniquely. Eren [15] uses LDA to analysis all literature related to COVID-19 and subdivided them into minor topics. As result, each subtopic were extracted with a set of keywords.

**Text Classification** Topic model can also treated to deal with classification task to identify unseen data. Kim [21] adapted the semi-supervised method with multi-co-training method to improve the overall classification performance. Moreover, the paper extended Word2Vec to Doc2Vec which maintain semantic relationship between two paragraphs. Doc2Vec transforms a paragraph into a  $d$ -dimensional vectors, which put documents with similar paragraph into near vector space.

**Recommender Systems** LDA often can be applied to recommendations. Xu[46] employed UIS-LDA (A User Recommendation based on Social Connections and Interests of Users in Uni-Directional Social Networks), which utilizes Generative Polya Urn (GPU) model and perform prediction for nearest user for the recommendations. Wang [44] implemented a LDA version which utilize the twitter datasets and recommend a serial of tourist location to user.

Moreover, in the growth of word embedding [30] enables an effective way to capture semantic meaning in language in a continuous vector space. Vocabularies that have similar meaning are close together by Euclidean distance.

### 1.3 Literature Review

These topic models take bag-of-word assumption and model each document as an admixture of latent topics, which are multinomial distributions over words. Nonnegative Matrix Factorization (NMF) [26] uses singular value decomposition to construct latent topic and topic-word distribution from the document, which consist of document-topic distribution matrix and topic word distribution matrix. Probabilistic Latent Semantic Analysis (PLSA) [20] is a probabilistic model assign every document in a single topic, and then assign word for every word position given the topic assignment. Similar to PLSA, Latent Dirichlet Allocation (LDA) [9] advanced PLSA from the topic assumption documents, which every document consist of admixture of topic distribution.

Some improvement exploit the correlation information between topics, which model the topic assignment with multivariate distribution to parameterize the relation between topics with mean and covariance.

Moreover, due to the success of LDA. there have been a numbers of topic models proposed on top of the LDA model. Dynamic Latent Dirichlet Allocation[5] was developed for continuous time data. Relational Latent Dirichlet Allocation [11] exploit the tuple information from the dataset and use it to inference the document set. Supervised Latent Dirichlet Allocation [29] includes labeled data which supposed to be helpful in several particular application areas such as movie review and sparse data prediction. Later LDA was extended to nonparametric version, hierarchical Latent Dirichlet Allocation (hLDA)[39], which follows a stochastic process called n-Chinese Restaurant Process (nCRP)[38]. hLDA maintain a hierarchical structure of topic instead of flat structure in LDA.

Amortized inference[22] are common in implementing to topic models, specifically, a neural network architecture with encoder-decoder are used into topic model structure for model inference. Srivastava[37] applied amortized variational inference to approximate the variational distribution of the model. Specifically, product of expert were used to collapse out the document-topic assignment parameter and simplify the inference process.

Some other attempts use graph techniques to model topic distributions. Yang[48] introduced new topic model with Graph neural network techniques. The paper introduced Graph Attention Topic Network (GATON) which hybridized the graph attention network (GAT) and amortized inference into application of topic modeling which supposed to reduce the require computation complexity. In past research, some considered n-gram to model the word pattern under a sentence structure which results a better prediction. Wallach [43] proposed a topic model using bi-gram information from the data set to yield a better performance in topic interpretability. Wang [45] extends the topic model to n-gram assumption with similar approach.

### 1.4 Objective and Outline

In the thesis, we would like to construct a topic model that make use of the positional information we obtain from the document set. Moreover, we would exploit the usage for imposing a prior to model the covariance on the document-topic proportion. Particularly, we develop the Transformer Embedding Correlated Topic Model(TECTM), a model that combine word embedding and topic model together to make a better fit of the dataset. Moreover, we integrate the Transformer into embedding, such that we can also take assumption of word position and convert it into meaningful contextual embeddings. In its generative

process, the model uses the topic embedding to form a per-topic distribution over the vocabulary. Specifically, the **TMTE** uses a log-linear model that takes the inner product of the word embedding matrix and the topic embedding. With this form, the **TMTE** assigns high probability to a word  $v$  in topic  $k$  by measuring the agreement between the word's embedding and the topic's embedding. To evaluate our model, we applied the proposed model on *20Newsgroups* and *Reuter-21578* dataset. The experiment results demonstrate that our model is capable to obtain high quality topics than the state-of-the-art model. In chapter 4, we also extended the model to handle time-series information, the model extends the architecture based on chapter 3. Specifically, we built Dynamic Transformer Embedding Correlated Topic Model(DTECTM), a model on top of the one from chapter 3, that make use of Gaussian Process Latent Variable Model(GPLVM) to captures time series information. We put our model into a set of experiments with other instances to examine the effectiveness. The models are compared with *NIPS* dataset and *UN debates* dataset, which each of them consist of a time label that represents the year a specific document belongs to. Additionally, we also visualize the time-to-topic proportion that the model obtained to explore the topics evolve over time. The result shows that To give a outline of this thesis, chapter 2 will give a brief background to the problem description and the related knowledge, including the existing topic models, the related methodology in , and the evaluation metric in NLP domain such as perplexity, and topic coherence and topic diversity, which is specifically for topic model evaluations. Then chapter 3 will specify the methodology and explain the detail of our model, we compare the proposed model with LDA and the latest model. In 4, we explore the DTECTM model to train with time-series data. We will specify the model and the Lastly, chapter 5 will sum up the merit and limitation overall the research, and prospect the future works.

## Chapter 2

### Background

In this section, we give short description on topic model techniques and several key components and terminologies related to LDA model. In section 2.1, we formulate the problem description of topic modeling. In 2.5.1, we give description to word embedding. In section 2.4.1, we explain the algorithm of Correlated topic model(CTM). Then, section 2.3 will give the definition for LKJ correlation distribution. Section 2.4.2 will cover the Embedded Topic Model(ETM). Finally, section 2.5.2 will give introduction to transformer embedding that we are used in the model. Several topic modeling techniques include non-negative matrix factorization (NMF) [26], Latent Semantic Analysis [23], probabilistic Latent Semantic Analysis (pLSA) [20] and Latent Dirichlet Allocation (LDA)[9].

#### 2.1 Problem definition

Massive data sets in internet has made the task of understanding data by accessing them one-by-one became not humanly possible. The raise of topic modeling gives possibility to summarize a given set of document collections. To describe topic modeling intuitively, for a document collection  $d = \{1 \dots D\}$  and define a number of topic  $K$ , the model outputs topic-word distribution  $\beta \in \mathbb{R}^{K \times V}$ , where  $K$  is number of topics and  $V$  is number of vocabularies in document set. For  $\{\beta_{k,1}, \beta_{k,2}, \dots, \beta_{k,V}\}_{k=1}^K$ , each  $\{\beta_{k,v}\}_{k=1}^K$  resulting the expression power vocabulary  $v$  could represent in topic  $k$ . The higher value a vocabulary obtained in tuple  $\{\beta_{k,v}\}_{v=1}^V$ , the high representation power that the word is related to a latent topic. Practically, we capture top- $m$  words from topic model for each topic  $k$ . For a set of top- $m$  words obtained from topic  $k$  in descending order, defined as  $\{\beta_{k,v_m}\}_{m=1}^M$  where  $\beta_{k,v_1} \succeq \dots \succeq \beta_{k,v_m} \succeq \dots \succeq \beta_{k,v_M}$ .

#### 2.2 Bag-of-word assumption

Suppose we have a collection of documents  $\in \mathbb{R}$ . The document and vocabularies Specifically, Topic model is a generative model that the probability based on Bag-of-words(BOWs) assumptions. Bag-of-words(BOWs) is a assumption that all words in the document are considered are independently distributed. To represent BOW, let a document collection  $W = (w_1, w_2, \dots, w_D), d \in \{1 \dots D\}$  documents, where  $w_d$  is a single document  $d$  contain words  $w_d = (w_{d1}, w_{d2}, \dots, w_{dN_d}), n \in \{1 \dots N_d\}$  word position, where  $w_{dn}$  is a single word in document  $d$  at position  $n$ . In equation 2.1, we take unigram model as an example[49],

$$p(W|\phi) = \prod_{d=1}^D p(w_d|\phi) = \prod_{d=1}^D \prod_{n=1}^{N_d} p(w_{dn}|\phi) = \prod_{d=1}^D \prod_{n=1}^{N_d} \phi_{w_{dn}} = \prod_{v=1}^V \phi_v^{N_v} \quad (2.1)$$

For sake of convenience, due to documents contain different size of words. It is not feasible to construct matrix representation for modeling. For the reason, we can exchange the representation for BOW from iterating the word occurrence for each word position to iterating the word occurrence in a vocabulary set. In this way, we can define a matrix for BOW:  $W \in \mathbb{R}^{D \times V}$ , with D rows of document and V rows of vocabulary count. Formally, for every document  $W = (w_{d1}, w_{d2}, \dots, w_{dV})$ ,  $v \in \{1, \dots, V\}$  vocabularies,  $w_{dv}$  is the occurrence of a vocabulary v in document d.

## 2.3 LKJ Correlation Distribution

LKJ distribution [27] is a distribution for modeling correlation matrix. The distribution is described as equation 2.2<sup>1</sup>

$$f(C|\eta) = 2^{\sum_{k=1}^{K-1} (2(\eta-1)+K-k)(K-k)} \times \quad (2.2)$$

$$\prod_{k=1}^{K-1} (B(\eta + (K-k-1)/2, \eta + (K-k-1)/2)^{K-k}) (\det(C)^{\eta-1}) \quad (2.3)$$

$B(\cdot, \cdot)$  is beta distribution, and K is the number of variable in correlation matrix.  $\eta$  is concentration parameter for LKJ distribution. When  $\eta = 1$  it is basically an uniform distribution allocated over the correlation matrix. If  $\eta > 0$  and become larger, the density of the matrix concentrates around center. To apply it into normal distribution as covariance, we could apply transformation equation 2.4 and turn it into covariance matrix[2].  $\text{diag}(\sigma)$  is the diagonal elements of the variance vector  $\sigma$ .

$$\Sigma = \text{diag}(\sigma) \cdot C \cdot \text{diag}(\sigma) \quad (2.4)$$

Directly drawing correlation matrix from LKJ distribution is not efficient in reality case. It is common to draw correlation matrix from factorized Cholesky LKJ distribution instead, where the probability density function is described as equation 2.5

$$\text{LKJChol}(L|\eta) \propto |J| \det(LL^\top)^{(\eta-1)} \quad (2.5)$$

$$= \prod_{k=2}^K L_{kk}^{K-k+2\eta-2} \quad (2.6)$$

The lower triangular matrix L is a Cholesky factorization for the correlation matrix iff  $L_{k,k} > 0$

$$\Sigma = \text{diag}(\sigma) \cdot LL^\top \cdot \text{diag}(\sigma) \quad (2.7)$$

similarly, the transformation from LKJ Cholesky matrix to covariance matrix as equation 2.4.

## 2.4 Topic Models

### 2.4.1 Correlated Topic Model

Correlated Topic Model (CTM)[7] is an extension of LDA[9] that utilize the correlation of latent topics, and relates the similar documents together. Instead

---

<sup>1</sup><https://distribution-explorer.github.io/multivariate-continuous/lkj.html>

of Dirichlet distribution, CTM applies multivariate logistic-normal distribution to model the word distribution.

The model contain  $K$  topics distribution as  $\beta_{1:K}$ ,  $z_{n,d}$  is the topic assigned to the  $n$ -th topic and  $d$ -th document.  $\theta_d$  is the corresponding proportion a topic is distributed to  $d$ -th document.  $\mu$  and  $\Sigma$  are the corresponding mean and  $K \times K$  covariance matrix of the distribution between documents.

---

**Algorithm 1:** Generative Process for CTM

---

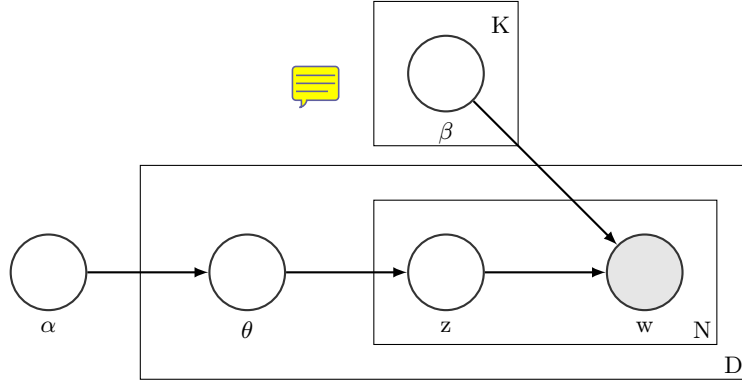
```

1 Initialize  $\mu, \Sigma$ 
2 for document  $d$  in  $D$  do
3   Sample a topic distribution  $\eta_d \sim \mathcal{N}(\mu, \Sigma)$ 
4   for word position  $n$  in  $N_d$  do
5     Sample a topic assignment  $z_{dn} \sim \text{Mult}(f(\eta_d))$ 
6     Sample a word  $w_{dn} \sim \text{Mult}(\beta_{z_d,n})$ 
7   end
8 end

```

---

From Algorithm 8, the parametrization  $\mu, \Sigma$  are initialized. For each document  $d$  in document collection  $D$ , a topic distribution  $\eta_d$  is drawn from normal distribution parametrized  $\mu, \Sigma$ . Then for each word position  $n$  in  $N$  words in document  $d$ , a topic assignment  $z_{dn}$  is drawn from multinomial distribution parametrized  $f(\eta_d)$ , where the transformation  $f(\eta)$  represents the softmax function maps the sample draw from normal distribution to topic proportion  $\theta$ , a topic distribution which points on the simplex that all elements in the vector sums to 1. Finally, a word is sampled from multinomial distribution  $\text{Mult}(\beta_{z_d,n})$ . From figure 2.1 we can see, word and topic-word assignment are in the word and



**Figure 2.1:** Graphical representation for CTM

document plate  $N \times D$ . The document topic proportion  $\eta$  is on the document plate  $D$ . Specifically, the topic word proportion  $\beta$  is on the topic plate  $K$ , which is specified as word distribution selected by topic assignment  $z$ .

**Mathematical Formulation** The joint distribution for CTM is described as follows,

$$p(\eta, z, w | \beta, \mu, \Sigma) = \prod_{d=1}^D p(\eta_d | \mu, \Sigma) \prod_{n=1}^{N_d} p(z_{dn} | \eta_d) p(w_{dn} | z_{dn}, \beta_{1:K})$$

and the ELBO is defined as,

$$\begin{aligned}\mathcal{L} \geq & \sum_{d=1}^D \mathbb{E}_{q_d} [\log p(\eta_d, z_d, w_d | \mu, \Sigma, \beta_{1:K})] - \sum_{i=1}^K \log \text{KL}(q(\eta_i | \lambda_i, \nu_i^2) || p(\eta_i | \mu, \Sigma)) \\ & - \sum_{n=1}^N \log \text{KL}(q(z_n | \phi_n) || p(z_n | \eta_d))\end{aligned}$$

### 2.4.2 Embedded Topic Model

Embedded Topic Model [14] is one of the state-of-art approaches for topic model task. It takes word distribution  $\beta$  as a topic embedding for words. Similar to Word2Vec[30], the word distribution is a softmax function of the inner product of context matrix  $\rho$  and context embedding  $\alpha$ . Specifically, the algorithm equation 2.8

$$\beta \sim \sigma(\rho^\top \alpha) \quad (2.8)$$

The word is drawn from the generative process shown in algorithm 2, for each document, sample a topic distribution  $\theta$  from logistic-normal distribution parameterized with zero mean and identity covariance. Then for each word position  $n$ , the model sample topic assignment  $z_{dn}$  from categorical distribution  $\theta_d$ . Finally, a word is drawn from  $\text{softmax}(\rho^\top \alpha)$  on  $z_{dn}$  the row.

---

**Algorithm 2:** Generative Process for ETM

---

```

1 foreach document  $d \in 1 \dots D$  do
2   Draw document topic distribution  $\theta_d \sim \mathcal{LN}(0, I)$ 
3   foreach word position  $n \in 1 \dots N_d$  do
4     Draw topic assignment  $z_{d,n} \sim \text{Cat}(\theta_d)$ 
5     Draw word  $w_{d,n} \sim \sigma(\rho^\top \alpha)_{z_{dn}}$ 
6   end
7 end
```

---

## 2.5 Representation learning

### 2.5.1 Word Embedding

Word embedding[3] is a kind of representation for words from document collections using a vector formulation. The nature of word embedding is that, the words that having similar meaning have a close distance(in most case euclidean distance), and vice versa. For instance, continuous bag-of-words(CBOW) [30] is a kind of word embeddings converting bag-of-word in to a vector of n-dimension continuous space, which contains the following formulation,

$$w \sim \text{softmax}(\rho^\top \alpha)$$

where  $\rho \in \mathbb{R}^{L \times V}$  is the embedding matrix which a function  $f : \mathbb{R}^V \mapsto \mathbb{R}^L$  maps  $V$  vocabularies into  $L$  dimension of continuous vector space. And  $\alpha$  is the context embedding, which conveniently convert the latent dimension  $L$  to a custom dimension of continuous embedding space  $\mathbb{R}^N$  as  $\tilde{f} : \mathbb{R}^L \mapsto \mathbb{R}^N$ .

### 2.5.2 Transformer

Transformer[42] is a popular neural network architecture in natural language processing. To briefly explain transformer, it is an stacked encoder-decoder architec-



ture. The component that makes transformer stands out of other architectures is the multi-head self-attention mechanism.

In this section we only cover the main components of transformer. The details for transformer can be reviewed in author's blog post<sup>2</sup>.

To define Transformer, equation 2.9 denotes the scaled dot product that transformer use to calculate attention score.  $Q \in \mathbb{R}^{T \times d_k}$ ,  $K \in \mathbb{R}^{T \times d_k}$  and  $V \in \mathbb{R}^{T \times d_v}$  are the query, key, value term vector used to calculate the context vector.  $S$  represents the sequence length and  $d_k$  and  $d_v$  the dimension for the key and value respectively.  $\frac{1}{\sqrt{d_k}}$  is used to scale down attention matrix in which to maintain a proper variance for the attention scores.

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V \quad (2.9)$$

Weight matrices  $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$ ,  $W_i^K \in \mathbb{R}^{d_{model} \times d_k}$  and  $W_i^V \in \mathbb{R}^{d_{model} \times d_v}$  for query, key and value are parameters to be learned. The scaled dot-product attention computes a sequence of vector outputs.  $W^O \in \mathbb{R}^{hd_v \times d_{model}}$ , where  $d_{model}$  is the dimension of key-value pair multiplies number of head.

$$\text{Multihead}(Q, K, V) = \text{concat}(h_1, \dots, h_h) W^O \quad (2.10)$$

$$h_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (2.11)$$

The sub-layer multi-head self-attention is connected to the fully-connected feed-forward network. The outputs of each sub-layer are added to the original output with residual connection applying Layer Normalization as eq. 2.12.

$$x = \text{LayerNorm}(x + \text{Sublayers}(x)) \quad (2.12)$$

## Positional Encoding

One drawback for Multi-Head Attention block is that it does not consider information about word positioning. The positional encoding function maps the sentence sequences into i-dimension of hidden space for each permutation position  $pos$ . And so the model can identify the from the additional positional features space of the input.<sup>3</sup>

$$PE_{(pos,i)} = \begin{cases} \sin \left( \frac{pos}{10000^{i/d_{model}}} \right) & \text{if } i \bmod 2 = 0 \\ \cos \left( \frac{pos}{10000^{(i-1)/d_{model}}} \right) & \text{otherwise} \end{cases}$$

the positional encoding function  $PE_{pos,i}$ , where the  $pos$  represents the permutation and the odd/even dimension are treated on sine/cosine function respectively.

## 2.6 Time Series model

### 2.6.1 Gaussian Process (GP)

Gaussian Process[36] is a versatile algorithm which is commonly used for both supervised and unsupervised learning problems including regression, classification and clustering. It has been demonstrated a strong pursuit on times-series

<sup>2</sup><https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

<sup>3</sup>Description from [https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial\\_notebooks/tutorial6/Transformers\\_and\\_MHAttention.html](https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial_notebooks/tutorial6/Transformers_and_MHAttention.html).

and reinforcement learning problems. A Gaussian Process  $\mathcal{GP}(\cdot, \cdot)$  is denoted by following, where  $m$  is the mean and  $k$  stands for a kernel for covariance function,

$$f \sim \mathcal{GP}(m, k)$$

For a set of input  $x_1, \dots, x_N$ , the joint probability density function  $p(f(x_1), \dots, f(x_N))$  is a normal distribution condition on mean vector  $m(X)$  and covariance matrix  $k(x, x')$  which is a positive semidefinite matrix. In most case the mean vector are specified with zero mean ( $m(X)=0$ ).

$$p(f(X)) = \mathcal{N}(m(X), k(X, X))$$

In particular, in Gaussian Process regression task, a output variable  $y_i$  denoted  $y_i \sim f(x_i) + \epsilon$  where  $\epsilon$  is a Gaussian noise. Therefore, given a set of training input and target,

$$p(y|f) = \mathcal{N}(y; f, \sigma_\epsilon^2 I)$$

Radial Basis Function kernel (namely RBF kernel) is a common function that used as a kernel in GP model. For a kernel  $k(x, x')$ ,

$$k(x, x') = \sigma^2 \exp\left(-\frac{\|x - x'\|^2}{2l^2}\right)$$

where  $l$  is the kernel bandwidth. The ELBO for GP is derived from [18], as

$$\log p(y) \geq \mathbb{E}_{q(f)}[\log p(y|f)] - KL[q(u)||p(u)]$$

where the variational distribution is  $q(f) \approx \int p(f|u)q(u)du$

### 2.6.2 Gaussian Process Latent Variable Model (GPLVM)

Gaussian Process Latent Variable Model [24, 40] is a dimension reduction method turn high dimensional data into low-dimension space. Formally, given a latent variable  $X \in \mathbb{R}^{N \times Q}$ , and high dimensional real valued observations  $Y \in \mathbb{R}^{N \times D}$ , the model induce high dimensional from a low-dimensional mapping  $\mathcal{X} \rightarrow \mathcal{Y}$  such that  $Q \ll D$ . The prior latent variable  $X$  is defined at a standard normal distribution in dimension  $Q$ .

$$p(X) = \prod_{n=1}^N \mathcal{N}(x_n; 0, I_Q)$$

The latent function  $f$  take the kernel  $K_d$  that determine the covariance matrix at  $D$ -dimension

$$p(f|X, \theta) = \prod_{d=1}^D \mathcal{GP}(f_d; 0, K_d)$$

and then it can recover back the observed data through the latent function  $f$  conditioned on the  $\mathcal{GP}$  prior with  $p(Y|f, X)$

$$p(Y|f, X) = \prod_{i=1}^N \prod_{d=1}^D \mathcal{N}(y_{n,d}; f_d(x_n), \sigma_y^2)$$

Here we discuss the variational inference method for the inference process. The variational distribution for GPLVM is given as eq. 2.13 according to GP sparse

approximation in [25]. In the formulation,  $x$  is the latent variable,  $f_d$  is the latent function for the covariance where variational distribution for  $q(X) = \sum_{n=1}^N \mathcal{N}(x_n; m_n, S_n)$  for each  $n^{th}$  row on  $x$ , and  $q(u_d) = \mathcal{N}(u_d|0, K_{MM})$

$$q(\{f_d, u_d\}_{d=1}^D, X) = \prod_{i=1}^N q(x_n) \prod_{d=1}^D p(f_d|u_d, X) q(u_d) \quad (2.13)$$

The latent variable ELBO is given by,

$$\mathcal{L} = \sum_{n,d} \mathbb{E}_{q_\phi(x_n)} \mathbb{E}_{p(f_d|u_d, x_n) q_\lambda(u_d)} [\log \mathcal{N}(y_{n,d}; f_d(x_n), \sigma_y^2)] \quad (2.14)$$

$$- \sum_n \text{KL}(q_\phi(x_n) || p(x_n)) - \sum_d \text{KL}(q_\lambda(u_d) || p(u_d|Z)) \quad (2.15)$$

the  $\phi$  is the local variational parameters,  $\lambda$  is the global variational parameters,  $\theta$  is the kernel hyperparameters and  $\sigma^y$  is the likelihood noise. to avoid heaving computation of  $X$  inside the conditional probability  $p(f_d|X)$ , inducing input has been introduced [25] to replace the  $X$  with inducing variables  $u_d \in \mathbb{R}^M$ , which conditioning on inducing input locations  $z \in \mathbb{R}^{M \times Q}$ ,

### 2.6.3 Dynamic Topic Model

Dynamic Topic Model [8] is a model enables to capture topic information from time-series data set. The model utilize documents from different time and generates the corresponding topic-word representation.

The generative process of the  $d^{th}$  document is the following:

---

**Algorithm 3:** Generative Process for DTM

---

- 1 Sample topics  $\beta^{(t)} \sim \mathcal{N}(\beta^{(t-1)}, \sigma^2 I)$
  - 2 Sample topic proportion mean  $\eta_t \sim \mathcal{N}(\eta_{t-1}, \delta^2 I)$
  - 3 **for** document  $d$  in  $D$  **do**
  - 4     Sample  $\theta_d \sim \mathcal{LN}(\eta_{t_d}, \alpha^2 I)$
  - 5     **for** word position  $n$  in  $N_d$  **do**
  - 6         Sample topic  $z_{d,n} \sim \text{Mult}(\theta_d)$
  - 7         Sample word  $w_{d,n} \sim \text{Cat}(\beta_{z_{d,n}}^t)$
  - 8     **end**
  - 9 **end**
- 

The model assume the topic-word proportion and the are in state-space model that move along the time  $1 \dots T$  with their corresponding variance  $\sigma$  and  $\delta$ .  $\eta_t$  is the latent variable model the prior mean to the topic proportion along the time line.  $\alpha_k^{t-1}$  is the mean of current time  $t$  which take the value previous time step  $t-1$ . The original DTM model approximate the prior by deriving the variational lower bound with Kalman Filter method.  $\gamma^2$  and  $\xi^2$  are the variances for the prior.

## 2.7 Posterior Inference

Since the exact inference of the posterior is intractable in real application, we employed approximation scheme for the posterior inference. The popular approaches are Markov Chain Monte Carlo Method (MCMC) and Variational Inference(VI)[6, 19]. Gibbs sampling is one of the MCMC method and it is fast to compute the approximation and easy to the implementation. Then, Variational EM algorithm to be carried out for maximizing the likelihood over all word in corpus in the document. An alternative way to perform estimation is Monte Carlo method.

### 2.7.1 Variational Inference

Given that posterior approximation is not always practical in real world application. Approximation methods are necessary to be applied. There are two main approaches for the posterior approximation: Markov Chain Monte Carlo (MCMC) and Variational Inference. Variational Inference is a method approximating the posterior in optimization fashion. To give a better intuition, let probability  $p(x)$  depending on a latent variable  $z$  such that  $p(x) = \int p(x|z)p(z)dz$ . We can turn the following posterior inference problem into an optimization problem. Here we derive the bound and perform optimization, by doing some calculation, we can obtain the following lower bound for the likelihood  $\mathcal{L}_i(p, q_i)$

$$\log p(x_i) \geq \mathbb{E}_{z \sim q_i(z)} [\log p(x_i|z) + \log p(z)] + \mathbb{E}_{z \sim q_i(z)} [\log q_i(z)] = \mathcal{L}_i(p, q_i) \quad (2.16)$$

then by deriving the KL-divergence, we can obtain the log probability  $\log p(x_i)$  is the likelihood-term minus the KL-divergence of  $q_i(z)$  and  $p(z|x_i)$ .

$$KL(q_i(z)||p(z|x_i)) = -\mathcal{L}_i(p, q_i) + \log p(x_i) \quad (2.17)$$

$$\mathcal{L}_i(p, q_i) = \log p(x_i) - KL(q_i(z)||p(z|x_i)) \quad (2.18)$$

with this derivation, we turn a posterior approximation problem into an optimization problem by maximizing the evidence lower bound (ELBO).

### 2.7.2 Stochastic Variational Inference

Stochastic Variational Inference (SVI)[19] is a scalable variant of variational inference, which enables mini-batching to split dataset and train for each epoch, then become a standard of optimization for probabilistic models. Two main improvements are made by the SVI: stochastic optimization and noisy gradient.



### 2.7.3 Collapsing Parameters

In the original LDA model, the parameter  $z$  is responsible for sampling topic assignment for each word position in every single document. Collapsing parameters[37] introduced to reduce the latent variable  $z$  in the generative process in order to speed up computation.

$$w_d \sim \prod_{n=1}^{N_d} \text{Cat}(\sigma(\theta_d^\top \beta_{w_{dn}})) \quad (2.19)$$

The trick in equation 2.19 rewrites the original LDA word drawing process, and hence defines a new evidence lower bound for the topic model.

### 2.7.4 Autoencoding Variational Bayes (AEVB)

Generally when we optimize a variational parameter, it is necessary to derive an ELBO and then derive the optimization step for gradient descent. While amortized inference latent variable  $z$  is parameterized by two inference networks  $\mu_\phi(x), \sigma_\phi(x)$ .

$$z = \mathcal{N}(\mu_\phi(x_i), \sigma_\phi(x_i)) \quad (2.20)$$

After deriving the ELBO, we obtain the following likelihood term.

$$\mathcal{L} = \mathbb{E}_{z \sim \mathcal{N}(\mu_\phi(x_i), \sigma_\phi(x_i))} [\log p_\theta(x_i|z)] - KL(q_\phi(z|x_i)||p(z)) \quad (2.21)$$

### 2.7.5 Reparameterization trick

The drawback of amortized inference is that, sampling from normal distribution parameterizing  $\mu_{\phi(x), \sigma_{\phi}(x)}$  could lead to high variance outcome and hamper the inference performance. For the reason, taking reparameterization trick[22] to transform as equation 2.22,

$$z = \mu_{\phi}(x_i) + \epsilon \sigma_{\phi}(x_i), \epsilon \sim \mathcal{N}(0, 1) \quad (2.22)$$

where  $\epsilon$  is a sample from normal distribution  $\mathcal{N}(0, 1)$ . and so the modified ELBO becomes equation 2.23,

$$\mathcal{L} = \mathbb{E}_{\epsilon \sim N(0,1)} [\log p_{\theta}(x_i | \mu_{\phi}(x_i) + \epsilon \sigma_{\phi}(x_i))] - KL(q_{\phi}(z|x_i) || p(z)) \quad (2.23)$$

## 2.8 Evaluation metrics

### 2.8.1 Perplexity

The proposed model will be evaluated with perplexity metric. The metric will examine how well the model can tackle with unseen data. It is equivalent algebraically to the inverse of the geometric mean per-word likelihood. Lower perplexity scores mean better.

$$\text{Perplexity}(D_{test}) = \exp - \frac{\sum_{d=1}^M \sum_{m=1}^{N_d} \log p(w_{dm})}{\sum_{d=1}^M N_d} \quad (2.24)$$

### 2.8.2 Topic Coherence

Topic Coherence[31] measures the quality of the topic

$$TC = \frac{1}{K} \sum_{k=1}^K \frac{1}{45} \sum_{i=1}^{10} \sum_{j=i+1}^{10} f(w_i^{(k)}, w_j^{(k)}) \quad (2.25)$$

where  $\{w_1^{(k)}, \dots, w_{10}^{(k)}\}$  denotes top-10 most likely words in topic k. And function  $f(\cdot, \cdot)$  is the normalized point-wise mutual information.

$$f(w_i, w_j) = \frac{\log \frac{P(w_i, w_j)}{P(w_i)P(w_j)}}{-\log P(w_i, w_j)} \quad (2.26)$$

### 2.8.3 Topic Diversity

In order to compare how the words each topic are differentiate the others. We applied the Topic Diversity metric [14]. Topic Diversity (TD) to be the percentage of unique words in the top 25 words of all topics. Diversity close to 0 indicates redundant topics; diversity close to 1 indicates more varied topics. We define the overall metric for the quality of a model's topics as the product of its topic diversity and topic coherence.

$$TD = \frac{|A \cap B|}{|A \cup B|} \quad (2.27)$$

where  $A$  and  $B$  are top-k words from two topics.

## Chapter 3

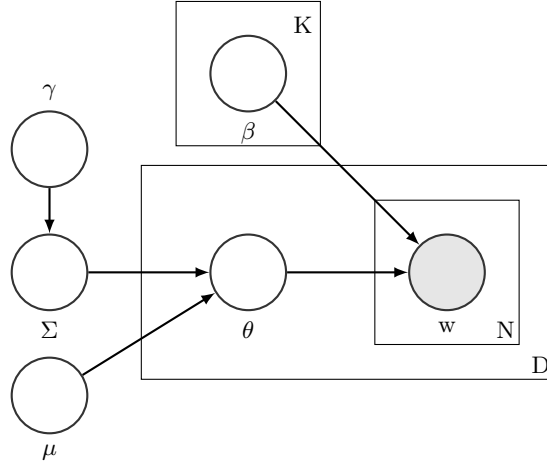
# Transformer embedding with Correlated Topic Model

In this chapter, we give a detailed explanation and procedure of Transformer embedding with Correlated Topic Model(TECTM) to be implemented.

### 3.1 Introduction

Latent dirichlet allocation(LDA)[9] is one of the popular model in topic modeling. However, the model take bag-of-word assumption that all words are independently distributed. Also, the original model take optimization step on entire set of document, where it limits the size of the data set can be trained within a limited memory size and hence scalability is a concern for LDA. Correlated Topic Model(CTM)[7] take consideration between topics by implementing the correlation information over the document-topic proportion. However, the model does not make assumption of prior information on the covariance matrix. Embedded Topic Model(ETM)[14] explore the possibilities the word2vec embedding to be working together with topic model to improve the quality of topic-words generation. However, since word2vec is a simple embedding architecture, the model could be better to be working with embedding that captures positional information.

In this chapter, we propose Transformer embedding with Correlated Topic Model(TECTM), a topic model that takes prior assumption on covariance matrix over the document-topic distribution, and integrate Transformer with the topic model to maintain a better quality of word representations in latent space. We propose LKJ correlation prior into our correlated topic model, where the correlation prior take place to capture correlation between topics by modeling the proportion over document-topic. To take advantage of positional information, we access the possibility the use of Transformer model to improve topic model performance. Transformer takes input sequence from documents and perform scaled dot-product to compute the each token in a sentence relates each other and the importance over a hidden context. In the following of the chapter, we first go through the related works that have been proposed by other authors. Then, we define the proposed model TECTM and derive its inference process. After that, we explain the implementation detail and the algorithmic setting for the experiment. Finally, the results are compared with other state-of-the-art models and discuss the performance our model out perform the other instances.



**Figure 3.1:** Graphical model for **TECTM**

### 3.2 Related works

Correlated Topic Model (CTM)[7] is the original work that proposed to alleviate the problem LDA, which did not utilize the topic information between correlated topics. The proposed model replaced Dirichlet distribution with a logistic-normal prior with covariance matrix to represent the relationship between topics.

There have been a several of works focus on word embedding and topic model. Major of them combined statistical model and embedding approach to model topic distribution. In other words, these method representing a word by mapping every single word into continuous space instead of using a probability distribution as was in typical LDA model. Embedded Topic Model[14] uses word2vec embedding to capture the word representation in latent continuous space. The posterior of the model was approximated by amortized inference. Xun [47] employed words embedding into Correlated Topic Model, the new correlated topic model as Correlated Gaussian Topic Model (CGTM). In their paper they make use of word embedding space and model the correlation between topics by calculation of similarity between words in the embedding space. Similarly, He[16] proposed Correlated Topic Modeling with Topic Embedding (CTMTE), which transformed the topic distribution previously obtained into lower dimension topic embedding space. The correlation between topics were directly computed through the similarity calculation in the vector space. The paper stated it reduces the running time as a scalable framework into large applications.

### 3.3 Model description

The TECTM utilizes the Transformer as embedding to the topic-word representations. To compare with the original topic model, the word-topic distribution  $\beta$  is the First, the topic embedding embeds the vocabulary into L-dimensional space, which is by Transformer embedding. Second, the context embedding maps the embedding into K-dimensional space. In the generative process, the TECTM uses the topic embedding to form a per-topic vector to represent the meaning over the vocabulary. The generative process of the  $d^{th}$  document is the following:

---

**Algorithm 4:** Generative Process for TECTM

---

```
1 Initialize hyperparameters  $\gamma, \mu$ 
2 Sample Cholesky factor  $L \sim \text{LKJChol}(\gamma)$ 
3 for document  $d$  in  $D$  do
4   | Sample topic distribution  $\theta_d \sim \mathcal{LN}(\mu, LL^\top)$ 
5   | for word position  $n$  in  $N_d$  do
6   |   | Sample word  $w_{d,n} \sim \text{Cat}(\theta_d^\top \sigma(\rho^\top \alpha))$ 
7   | end
8 end
```

---

From algorithm 4, starting from step 1, the topic proportion  $\theta_d$  is drawn from the logistic-normal distribution  $\mathcal{LN}(\cdot)$  with zero mean and identical covariance. From Step 2-a, for each word position  $n$  in document  $d$ , a topic assignment to word  $w_{dn}$  is drawn from categorical distribution  $\text{Cat}(\theta_d)$  parameterized by topic proportion  $\theta_d$ . Step 2-b, the model draw a word from embedding of the vocabulary  $\rho$  and the assigned topic embedding  $\alpha_{z_{dn}}$  to draw the observed word from the assigned topic, as given by  $z_{dn}$ . The embedding is applied softmax function to make them topic distribution. The TECTM likelihood uses a matrix of word embedding  $\rho$ , a representation of the vocabulary in a lower dimensional space. In practice, it can either rely on previously fitted embeddings as part of the fitting procedure, it simultaneously finds topics and an embedding space.

### 3.3.1 LKJ Correlation prior

The LKJ Correlation prior  $\text{LKJChol}$  a A Cholesky factor of lower triangle matrix  $L$ , from a decomposed LKJ correlation distribution. The product of lower triangular matrix  $LL^\top$  reconstruct the correlation matrix

$$\Sigma = LL^\top$$

A covariance matrix can be reconstructed in following fashion. In our implementation, we don't scale the covariance matrix with variance  $\sigma$ , which is simply a correlation matrix.

### 3.3.2 Transformer Embeddings

Following the ETM architecture, we modify topic-word distribution as an embedding and put transformer embedding to work into it.

$$\beta \sim \text{softmax}(\rho^\top \alpha) \tag{3.1}$$

Equation 3.1, the topic-word distribution is composed of the dot product of transformer embedding  $\rho \in \mathbb{R}^{L \times V}$  representing the word vector in  $L$ -dimension of continuous space and topic matrix  $\alpha \in \mathbb{R}^{L \times K}$  mapping the  $L$  dimension vector into  $K$ -dimension of topic proportions.

### 3.3.3 Marginal Likelihood

To compute the parameters of the model, we first compute the log-marginal likelihood. In equation 3.2, the marginal likelihood is parameterized by transformer embedding  $\rho$  and topic embedding  $\alpha$ , which is for constructing topic-word proportion  $\beta$  as referenced in equation 3.1,

$$\mathcal{L}(\rho, \alpha) = \sum_{d=1}^D \log p(w_d | \rho, \alpha) \tag{3.2}$$



the marginal probability for  $p(w|\alpha)$  on  $d$ -th document,

$$p(w_d|\rho, \alpha) = \int p(\theta_d|\mu, \Sigma) \prod_{n=1}^{N_d} p(w_{dn}|\theta_d, \rho, \alpha) d\theta_d \quad (3.3)$$

the conditional distribution  $p(w_{dn}|\theta_d, \alpha)$  marginalize out the topic assignment  $z_{dn}$  by collapsing parameters transformation  $w \sim \theta^\top \beta$ ,

$$p(w_{dn}|\theta_d, \rho, \alpha) = \text{Cat} \left( \sum_{k=1}^K \sigma(\theta_{dk} \beta_k) \right) \quad (3.4)$$

as always, computing integral for posterior is intractable, approximate inference is necessary to estimate the true parameter from the integral.

### 3.3.4 Joint Distribution

We give a description of the joint distribution,  $W, Z, \theta$  and  $\Sigma$  are variables and  $\beta, \mu$  and  $\gamma$  are latent variables.  $W$  is the word likelihood from the document collections,  $Z$  represents the topic-word assignment,  $\theta$  models the topic distribution for each document, and  $\Sigma$  is the covariance matrix which depends on the document-topic distribution  $\theta$ .

$$\begin{aligned} p(W, Z, \theta, \Sigma|\beta, \mu, \gamma) &= p(W|Z)p(Z|\theta)p(\theta|\mu, \Sigma) \\ &= p(\Sigma|\gamma) \prod_{d=1}^D p(\theta_d|\mu, \Sigma) \prod_{n=1}^V p(z_{d,n}|\theta_d) p(w_{d,n}|z_{d,n}, \beta) \end{aligned}$$

by taking log on the joint probability, we obtain a objective function for optimization

$$\begin{aligned} \log p(W, Z, \theta, \Sigma|\beta, \mu, \gamma) &= \sum_{d=1}^D \left[ \log p(\theta_d|\mu, \Sigma) + \sum_{n=1}^V [\log p(z_{d,n}|\theta_d) + \log p(w_{d,n}|z_{d,n}, \beta)] \right] \\ &\quad + \log p(\Sigma) \end{aligned}$$

## 3.4 Inference and Estimation

To perform posterior inference, we apply variational inference to transform the log-likelihood function into a lower bounded optimization problem.

### 3.4.1 Variational distribution

The variational distribution for  $q(\theta_d)$  is a multivariate normal distribution in  $\mathbb{R}^D$  that parameterized by mean vector  $\mu(w_d)$ , a inference network take input from bag-of-words  $w_d$  and then outputs means from its weight parameters;  $\Sigma$  is a conditional variational parameter from  $q(\Sigma)$ .

$$q(\theta_d|\Sigma) = \mathcal{N}(\mu(w_d), \Sigma)$$

### 3.4.2 Evidence Lower Bound(ELBO)

To perform Variational inference, it is essential to derive the Evidence Lower Bound (ELBO) first as the objective function for the optimization. By equation

3.5

$$\begin{aligned}
\mathcal{L} &\geq \mathbb{E}_q[\log p(W, Z, \theta, \Sigma)] - \mathbb{E}_q[\log q(Z, \theta, \Sigma)] \\
&= \sum_{d=1}^D \sum_{n=1}^V \mathbb{E}_q[\log p(w_{d,n}|z_{d,n}, \beta)] + \sum_{d=1}^D \sum_{n=1}^V \mathbb{E}_q[\log p(z_{d,n}|\theta_d)] \\
&\quad + \sum_{d=1}^D \mathbb{E}[\log p(\theta_d|\mu, \Sigma)] + \mathbb{E}_q[\log p(\Sigma)] - \sum_{d=1}^D \sum_{n=1}^V \mathbb{E}_q[\log q(z_{d,n}|\alpha_{d,n})] \\
&\quad - \sum_{d=1}^D \mathbb{E}_q[\log q(\theta_d|\lambda_d, \nu_d)] - \mathbb{E}_q[\log q(\Sigma|\gamma)]
\end{aligned} \tag{3.5}$$

**Collapsing Parameters** we can speed up the computation by marginalize out the  $z$ . The probability for word  $w$  could be simplified as  $w_{dn} \sim \text{Cat}(\sigma(\theta_d^\top \beta))$ .

$$\begin{aligned}
\mathcal{L} &\geq \sum_{d=1}^D \int \int q(\theta_d) \log \frac{p(W_d|\theta_d, \beta) p(\theta_d|\mu, \Sigma) p(\Sigma|\gamma)}{q(\theta_d) q(\Sigma)} d\theta_d d\Sigma \\
&= \sum_{d=1}^D (\mathbb{E}_{q(\theta_d)} [\log p(W_d|\theta_d, \beta)] - KL(q(\theta_d)||p(\theta_d|\mu, \Sigma))) - KL(q(\Sigma|\gamma)||p(\Sigma))
\end{aligned} \tag{3.6}$$

Here we define amortized inference [22], a optimization technique which perform inference by defined neural networks.  $\mu_\theta(w)$  is an inference network that takes inputs from the bag-of-words vector  $w_d$ . Then a output is generated by the normal distribution parameterized by mean vector generated by  $\mu_\theta(w)$  and covariance  $\Sigma$ .

$$\begin{aligned}
&= \sum_{d=1}^D \left( \mathbb{E}_{\theta_d \sim \mathcal{N}(\mu_{\theta_d}(w_d), q(\Sigma))} [\log p(w_d|\theta_d, \beta)] - KL(q(\theta_d)||p(\theta_d|\mu, \Sigma)) \right) \\
&\quad - KL(q(\Sigma)||p(\Sigma))
\end{aligned} \tag{3.7}$$

To apply reparameterization trick, we take transformation from normal distribution to  $\theta = \mu + \Sigma^{1/2}\epsilon$  where  $\epsilon \sim N(0, I)$  drawn as  $1 \times K$  vector, as equation 3.8

$$\begin{aligned}
&= \sum_{d=1}^D \left( \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} \left[ \log p(w_d|\sigma(\mu_{\theta_d}(w_d) + \Sigma^{1/2}\epsilon), \beta) \right] - KL(q(\theta_d)||p(\theta_d|\mu, \Sigma)) \right) \\
&\quad - KL(q(\Sigma)||p(\Sigma|\gamma))
\end{aligned} \tag{3.8}$$

The KL-divergence for the logistic-normal distribution is given as equation 3.9 closed-form expression, the KL-divergence between  $q(\theta_d)$  and  $p(\theta_d)$  becomes

$$\text{KL}(q(\theta_d)||p(\theta_d|\mu, \Sigma)) = \frac{1}{2} \left( \text{tr}(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1}(\mu_1 - \mu_0) + \log \frac{|\Sigma_1|}{|\Sigma_0|} - K \right) \tag{3.9}$$

so the ELBO then becomes 3.10

$$\begin{aligned}
\tilde{\mathcal{L}} &= \sum_{d=1}^D \frac{1}{S} \left[ \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ w_d^\top \log \text{Cat}(\sigma((\mu_0(w_d) + \Sigma_0^{1/2}\epsilon)^\top \beta)) \right] \right. \\
&\quad \left. - \frac{1}{2} \left( \text{tr}(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1}(\mu_1 - \mu_0) + \log \frac{|\Sigma_1|}{|\Sigma_0|} - K \right) \right] \\
&\quad - KL(q(\Sigma)||p(\Sigma|\gamma))
\end{aligned} \tag{3.10}$$

The expectation log likelihood term in 3.5 can be efficiently approximated by the Monte Carlo sampling method,

$$\begin{aligned}\mathcal{L} \approx & \sum_{d=1}^D \left[ \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ w_d^\top \log \text{Cat}(\sigma((\mu_0^{(s)}(w_d) + \Sigma_0^{1/2} \epsilon^{(s)}))^\top \beta) \right] \right. \\ & - \frac{1}{2} \left( \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1} (\mu_1 - \mu_0) + \log \frac{|\Sigma_1|}{|\Sigma_0|} - K \right) \Big] \\ & - KL(q(\Sigma) || p(\Sigma | \gamma))\end{aligned}\quad (3.11)$$

where the expectation of the reconstruction loss is taken from a set of sample  $S$  to compute the unbiased estimate of ELBO. We also apply the minibatch to make able the model perform by sub-sampling the document collection. By equation 3.12

$$\begin{aligned}\tilde{\mathcal{L}} \approx & \frac{D}{|\mathcal{B}|} \sum_{d \in \mathcal{D}_{\mathcal{B}}} \left[ \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ w_d^\top \log \text{Cat}(\sigma((\mu_0^{(s)}(w_d) + \Sigma_0^{1/2} \epsilon^{(s)}))^\top \beta) \right] \right. \\ & - \frac{1}{2} \left( \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1} (\mu_1 - \mu_0) + \log \frac{|\Sigma_1|}{|\Sigma_0|} - K \right) \Big] \\ & - KL(q(\Sigma) || p(\Sigma | \gamma))\end{aligned}\quad (3.12)$$

The transformer loss is calculated over the sum of tokens of a selected sentence sequence SEQ from each document, the loss of each token is evaluated in Cross entropy loss.

$$L_{CrossEntropy} = \sum_{d=1}^D \sum_{n=1}^{|\text{SEQ}|} \text{CrossEntropy}(w_{dn}) \quad (3.13)$$

where

$$\text{CrossEntropy}(w) = - \sum_{i=1}^V p(w^{(i)}) \log \hat{p}(w^{(i)}) \quad (3.14)$$

it is a metric for comparing the probability between the word probability  $p(w)$  and the probability from prediction  $\hat{p}(w)$ , which is converted to a probability from one hot vector by mapping it to simplex by softmax function that sum to 1.

### 3.4.3 Optimization step

In algorithm 21, first initialize the model and variational parameters. Then, for each epochs, we obtain the transformer embedding  $\rho$  from transformer. After that, the topic embedding  $\beta$  is computed by taking softmax of dot-product of  $\rho$  and  $\alpha$ . Then a minibatch  $\mathcal{B}$  is selected from the document for optimization. The number of minibatch is the document collection divides minibatch size where  $\#\text{minibatch} = \frac{D}{|\mathcal{B}|}$ . For each minibatch, the model takes a document and sample lower Cholesky matrix from LKJ Cholesky distribution(description see section 2.3). A topic assignment for document  $d$   $\theta_d$  is sampled from logistic-normal distribution  $\mathcal{LN}(\mu, \mathbf{LL}^\top)$ , where  $\mu$  is sampled from half-Cauchy distribution and covariance is a transformation from equation 2.4. For each word position  $n$ , a word is sampled from the softmax of dot-product of transformer embedding  $\rho$  and NN weight  $\alpha$ . After the sampling process for the document collection, we estimate the ELBO loss  $L_{ELBO}$  for the topic model, and the cross entropy loss

$L_{CrossEntropy}$ . Remind that the topic model and transformer take input differently. The topic model part takes bag-of-words input, a document-vocabulary matrix  $D \times V$  counting the occurrence of vocabulary  $v$  in document  $d$ . While transformer take sequence of document as input. To calculate the loss of the model, we sum up the ELBO loss  $L_{ELBO}$  and cross entropy loss for transformer  $L_{CrossEntropy}$ . Then a stochastic gradient is computed by backpropagation. a gradient step to . The process iterates until the maximum iteration is reached.

---

**Algorithm 5:** Training on TECTM

---

```

1 Initialize model and variational parameters
2 for epoch  $i = 1, 2, \dots N$  do
3   Obtain trnasformer embedding  $\rho$ 
4   Compute  $\beta = \text{softmax}(\rho^\top \alpha)$ 
5   Choose a minibatch  $\mathcal{B}$  of documents
6   foreach document  $d$  in  $\mathcal{B}$  do
7     Compute  $\mu_d = \mu_\phi(w_d)$ 
8     Sample  $L \sim \text{LKJChol}(\gamma)$ 
9     Sample  $\theta_d \sim \mathcal{LN}(\mu_d, \Sigma)$  where  $\Sigma = LL^\top$ 
10    foreach word position  $n$  in docuemnt  $N_d$  do
11      | Sample word  $w_{dn} \sim \text{Cat}(\sigma(\theta_d \beta))$ 
12    end
13  end
14  Estimate ELBO loss  $L_{ELBO}$  from Eq. 3.12
15  Compute Transformer loss  $L_{CrossEntropy}$  from Eq. 3.13
16  Compute the total loss  $L = L_{ELBO} + L_{CrossEntropy}$ 
17  Compute the stochastic gradient via backpropagation
18  Take a stochastic gradient step
19  Update model parameters
20  Update variational parameters
21 end

```

---

### 3.5 Results & Evaluations

In this chapter, we perform evaluation on our model and the other algorithms.

#### 3.5.1 Experiment Testing

The experiment will be conducted with a number of existing proposed topic models as mentioned related work section above. We conduct the experiment with those baseline algorithms and evaluate them in terms of accuracy and running time. Some of the source code of competitive were provided by their authors in Github<sup>1</sup>. The outcome result will be extensively studied and conclude the insight behind the algorithms and methodologies. Detail to be stated in section 3.5.3. Our probabilistic part of model implemented using Pyro[4], while the model optimization and transformer implementations are based on PyTorch[34].

#### 3.5.2 Dataset

To evaluate the performance of the model, we selected *20Newsgroups* and *Reuters-21578* data sets in our evaluation stage. 20Newsgroups consist of 18,846 news

---

<sup>1</sup>For instance, Correlated Topic Model(CTM), <https://github.com/blei-lab/ctm-c>

group documents <sup>2</sup> and the Reuters-21578 includes 10,788 documents in total. Both of the dataset will be preprocessed to remove stop-words and stemming before the evaluation. Both data set were separated into training/testing set for training and evaluation process. *20Newsgroups* data set contains around 20,000 newsgroups documents, which divided into 20 different groups. In the preprocessing stage, we remove the document with only one word. We filter the stop-words, remove the word with special characters. The frequency of words are limited to between 2%-50%. After the preprocessing, the data set was split into 11314, 7532 documents with 5651 vocabularies. *Reuter-21578* data set is a collection of documents from Reuters newswire in 1987. After performing the preprocessing, the processed data set consist of 7769, 3019 documents for train/test documents with 1622 vocabularies. On the data preprocessing stage, we perform tokenization, stopword removal, lemmatization on the documents.

**Transformer learning task** For transformer embedding training, we shape the dataset into a sequence set  $S$  of equal distance of pre-defined sentence length  $SEQLEN$ . Token  $\langle PAD \rangle$  are padded to the sequence when the sentence in  $i$ -index  $S[i] \langle SEQLEN \rangle$ . The out-of-vocab tokens are replaced with  $\langle OOV \rangle$  token. The sequence dataset are to train the transformer embedding as input in the training process. To organize the training task, we define a set of token and target to Transformer learn it explicitly. We configure 15% of tokens within a sentence are to be masked for the training task. First we replace a  $\langle MASK \rangle$  token on the selected word position of the token input. Other tokens will be filled a  $\langle IGNORE \rangle$  token on the corresponding target position to exclude those token position to be calculate in the prediction score.

### 3.5.3 Models

We compare the model performance with a numbers of rivals. We take Latent Dirichlet Allocation (LDA)[9] as the baseline model. Other models include Transformer[42]<sup>3</sup>, ProdLDA[37] and Embedded Topic Model(ETM)[14]. LDA fits the model with with mean-field variational inference <sup>4</sup> ProdLDA, a topic model learning using amortized inference to approximate the variational distribution on document-topic proportion. ETM model, a topic model built on top of ProdLDA model, uses dot-product of word embedding and topic embedding to represent the topic-word distribution  $\beta$ .

### 3.5.4 Algorithmic Settings

To perform posterior inference, we employed Stochastic Variational Inference (SVI) [19] for the optimization problem. We set the minibatch size to 512 documents. For LDA, we applied the model provided from sklearn package (version 0.24.0) <sup>5</sup>. For ETM, we run the experiment with the parameter suggested [14]. For ProdLDA, we perform optimization with inference network architecture as described in the paper [37]. To train the models, every model run on 300 epochs to give the best performance. To perform optimization, we use Adam for the gradient ascent algorithm, and we set the learning rate to 2e-3. we use  $l_2$ -regularization to the 1e-6, We applied the settings from [37] to perform amortized inference. The

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

<sup>3</sup>Not a topic model, but we think it is worth to make comparison still.

<sup>4</sup>Adopted from scikit-learn library

<sup>5</sup>Sklearn website <https://scikit-learn.org/stable/index.html>

inference architecture included 300 dimension of hidden layers. The dimension for embedding  $\rho$  are set to 256. For the Transformer model settings, we define the sequence length to 20, number of head to 8, and 4 layer stacks of transformer encoder and 256 hidden dimension.

### 3.5.5 Quantitative Result

Perplexity has been known for evaluating the fitness of NLP models, however it may not be the best metrics to examine the wellness of a topic model[33] Secondly, in order to compute the perplexity for topic model using AEVB, we follow the previous work, using the variational lower bound to compute the perplexity. For the same reason, the perplexity may not perfectly reflects the true quality of a topic model. Hence, we take the topic coherence [31] as the main measure to evaluate the models. Also, sometimes even a topic model could obtain a high coherence score, the topics could be repeated many times that yields poor result. We also implemented topic diversity [14] as one of the measure of topic model, the measure that evaluate how well each topic contains distinct word from other topics. In this section, we evaluate the model with the following metric : Perplexity, Topic Coherence (TC), Topic Diversity (TD).

#Topic Metrics	k=20			k=50		
	PPL	TC	TD	PPL	TC	TD
LDA	478.8	<b>0.248</b>	0.702	507.1	0.193	0.554
ETM	279.0	0.213	0.500	336.8	0.200	0.211
ProdLDA	796.6	0.188	<b>0.792</b>	521.8	0.180	<b>0.609</b>
TECTM	<b>259.6</b>	0.237	0.586	<b>263.0</b>	<b>0.226</b>	0.367

**Table 3.1:** Result for Reuters-21578 dataset

#Topic Metrics	k=20			k=50		
	PPL	TC	TD	PPL	TC	TD
LDA	2442.7	0.168	<b>0.774</b>	2538.2	0.154	<b>0.713</b>
ETM	<b>1640.5</b>	0.191	0.592	<b>1715.5</b>	0.159	0.337
ProdLDA	6018.5	0.072	0.734	9057.9	0.014	0.703
TECTM	1954.4	<b>0.200</b>	0.630	1944.6	<b>0.194</b>	0.509

**Table 3.2:** Result for 20Newsgroups dataset

**Reuters-21578** On the data set Reuter-21578, our model perform the best in perplexity, with 259.6 and 263.0 when  $k=20$  and  $k=50$  respectively. Our model also performs the best on topic coherence when  $k = 50$ , where it is 0.226, with a competitive topic coherence score when  $k = 20$  as well. To focus on TD score, our model did not beat the best model, but also maintain a good enough TD score to generate variety of topics, with 0.586 and 0.367 when  $k = 20$  and  $k = 50$  respectively. As result apparently our model perform well on several metrics.

**20Newsgroups** On the result from table 3.2, displays that our model has outperform the other model by TC score. On perplexity score, ETM obtain the best score when  $k = 20$  and  $k = 50$ . On the other hand, LDA has the best TD score on both  $k = 20$  and  $k = 50$ . It can When  $k = 20$ , our model has 1954.4 in perplexity score, 0.200 TC score and 0.630 in TD score. When  $k = 50$ , our model

has 1944.6 in perplexity score, 0.194 TC score and 0.509 in TD score. ProdLDA, perform the worse in both perplexity and topic coherence score.

It can be seen that, when the number of topic increases, the topic model performance decrease proportionally. For instance, ProdLDA has a quite well metric score on small data set such as Reuters-21578. However, it performance drag down drastically on a bigger data set such as 20Newsgroups. Also we observe that when our model is capable to maintain a relatively good topic coherence score when the topic increase from the result above. To compare with LDA, it could obtain a high enough TC and TD scores in our experiment settings, however our model does better in TC scores in both topic numbers we conducted investigations.

### 3.5.6 Training

From figure 3.2, 3.3, display the training process of the training loss and log probability by 200 epochs. We observe that the negative log probability is decreasing over epochs. For the metrics, the topic coherence and the topic diversity improves slowly along the training proceeds, the perplexity also decreases along.

### 3.5.7 Qualitative Result

The proposed model will be evaluated with a number of specifically selected topic and examined with their performance separately. The result will be exhaustively compared with other existing models.

From table 3.3, we have selected some topic words each model generated from 20Newsgroups when  $k = 20$ . The topics represent space, operating system, religion, encryption and guns respectively. Our model has shown capability on capturing key words from each topic, such as on topic "space": *nasa*, *space*, *jpl*, *moon*, *earth*, *station*, *flight* are the outputs.

### 3.5.8 Visualization

To clearly demonstrate the representation for the word embedding space inside the parameters that obtained, we ran t-SNE algorithm to map the topic-word representation into 2-dimension continuous space. We selected and zoomed-in a specific topic with its neighboring topics with a red box and displayed over the figures. In figure 3.4 and 3.6, demonstrate the t-SNE visualization of the topic-word distribution for 20Newsgroups for  $k=20$  and  $k=50$ .

We also compare the model with the ETM as displayed on figure 3.5. In figure 3.4, the topic "space"(which consist of words like space, earth, orbit, etc) laying on the 2-dimensional space. The most neighboring topic includes "hardware" and "operating system", which appear the words such as "cable", "cd", "power" for "hardware" topic, and "images", "graphics" and "files" for "operating system" topic. Generally speaking, ETM do not able to distinct the words of topics into visible clusters, as shown on the figure, the upper part of word dots with different topic labels mix together. And inside the topic cluster, there are also some word from other topic mixed in to it. This may implies ETM does not generate good enough topic-words representation to distinguish the difference between topics.

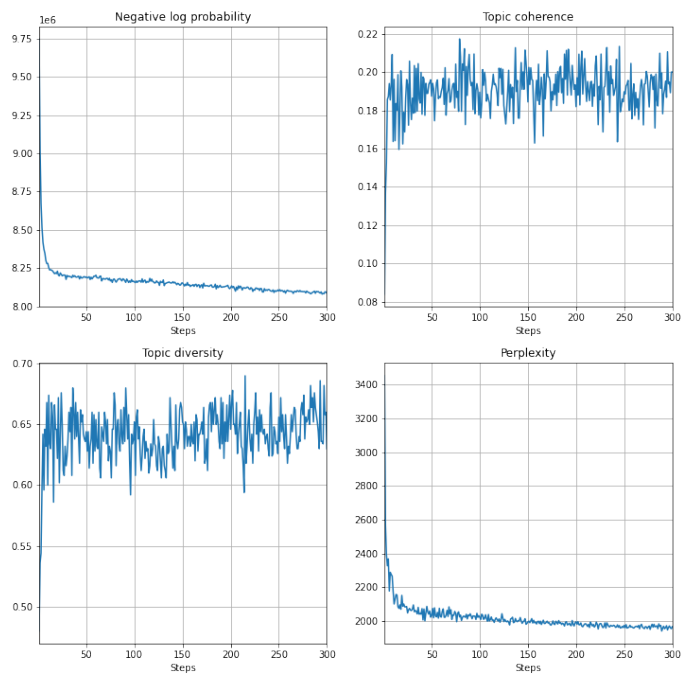


Figure 3.2: 20Newsgroups ( $k=20$ )

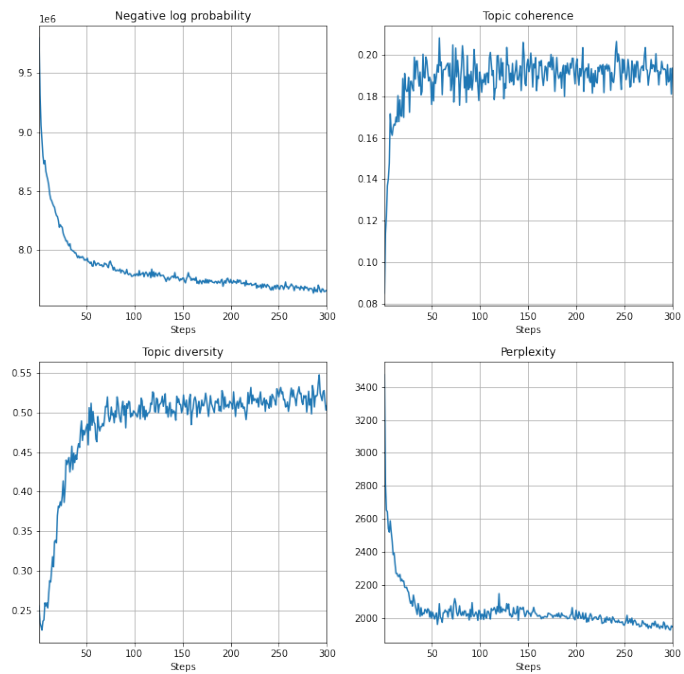
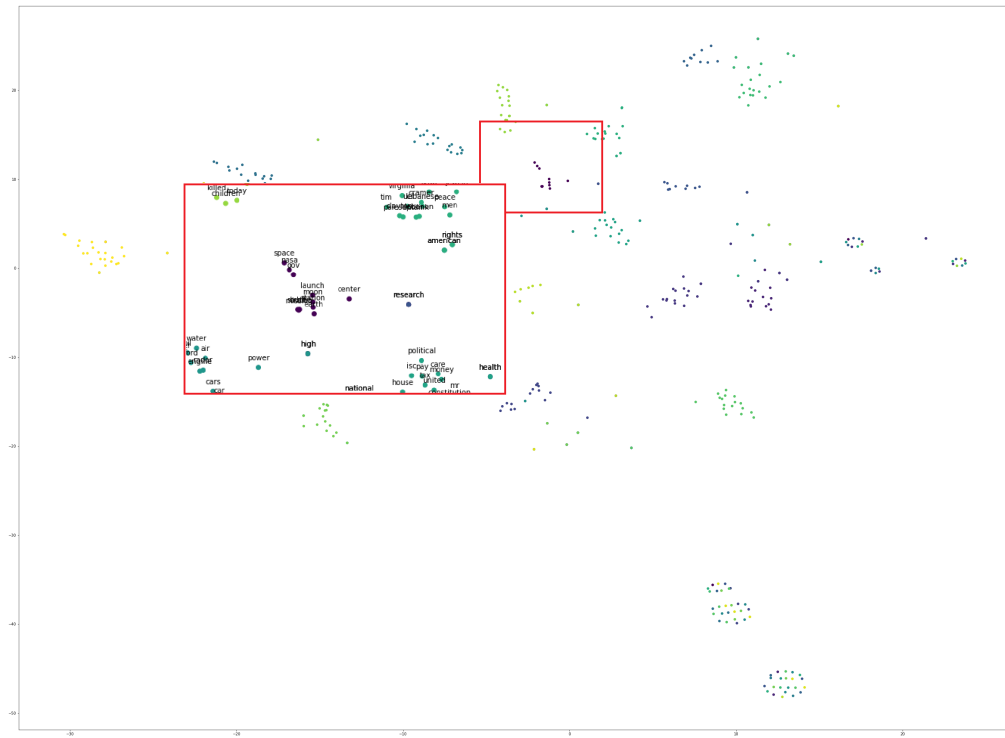
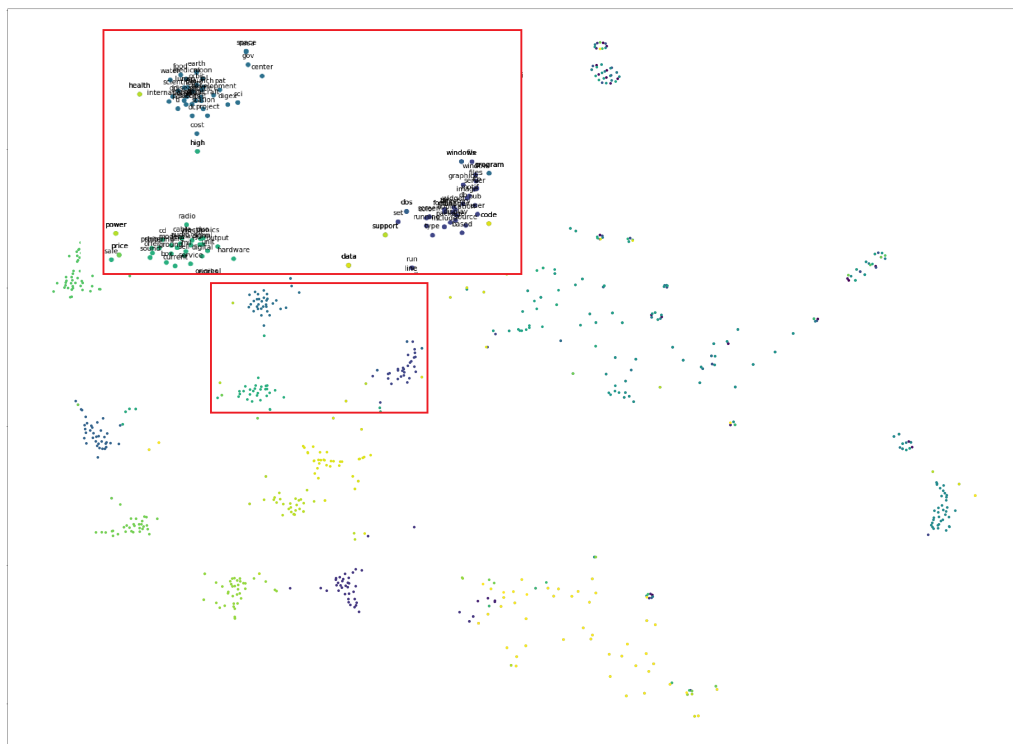


Figure 3.3: 20Newsgroups ( $k=50$ )





**Figure 3.4:** t-SNE Visualization for our model #Topic:20



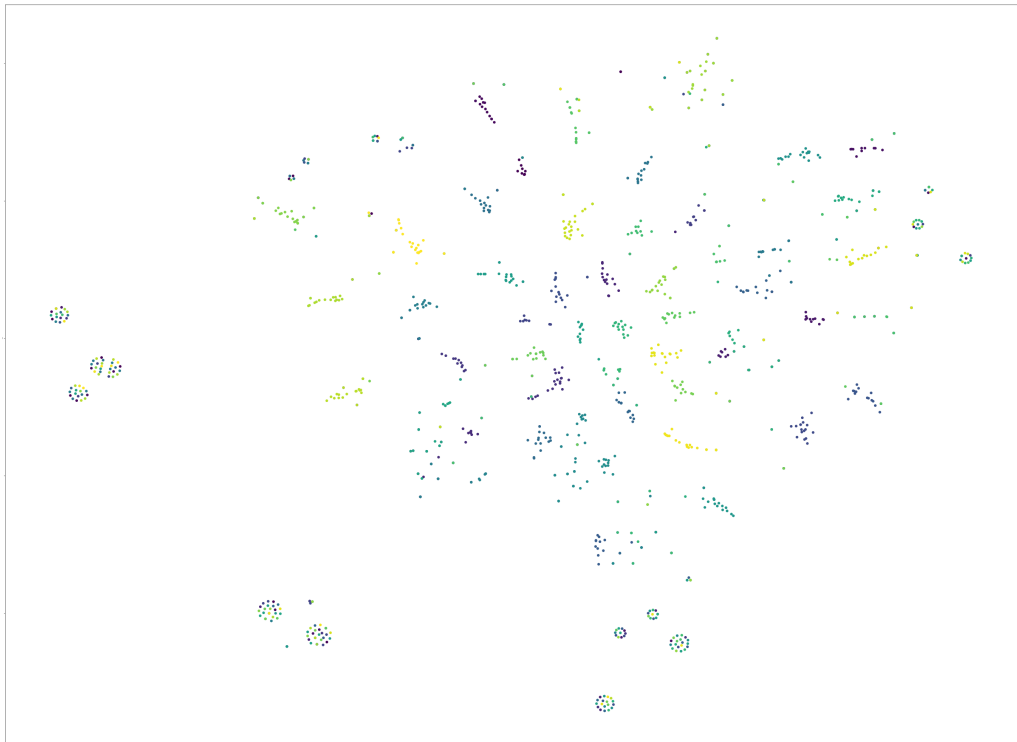
**Figure 3.5:** t-SNE visualization for ETM #Topic:20

Our Model
<b>nasa</b> gov <b>space</b> <b>jpl</b> <b>moon</b> <b>earth</b> <b>station</b> <b>flight</b> research digex <b>windows</b> <b>window</b> problem <b>dos</b> running <b>file</b> <b>mouse</b> mit de ms <b>god</b> <b>jesus</b> <b>christian</b> people faith bible time church good things <b>key</b> chip <b>encryption</b> clipper <b>security</b> <b>privacy</b> government <b>keys</b> public escrow <b>gun</b> people control government <b>guns</b> <b>weapons</b> american make clinton state
ProdLDA
<b>nasa</b> <b>space</b> gov people <b>station</b> <b>time</b> <b>orbit</b> dc program shuttle <b>scsi</b> <b>drive</b> <b>controller</b> max <b>drives</b> <b>ide</b> senior <b>tape</b> time people <b>god</b> <b>jesus</b> <b>atheists</b> <b>christian</b> bible religion <b>atheism</b> <b>christians</b> word truth <b>key</b> <b>crypto</b> session nt chips chip <b>serial</b> dos <b>keys</b> <b>encrypted</b> <b>gun</b> people god religion writes life morality ohio argument question
LDA
<b>space</b> , <b>nasa</b> , gov, access, <b>launch</b> , <b>earth</b> , digex, <b>moon</b> , <b>orbit</b> <b>file</b> , <b>window</b> , <b>program</b> , <b>ftp</b> , <b>files</b> , <b>server</b> , <b>image</b> , <b>graphics</b> , <b>windows</b> <b>god</b> , people, <b>jesus</b> , <b>christian</b> , bible, writes, life, <b>christians</b> , time <b>key</b> , <b>encryption</b> , chip, clipper, <b>keys</b> , <b>security</b> , <b>government</b> , <b>privacy</b> <b>gun</b> , <b>guns</b> , <b>law</b> , <b>police</b> , people, <b>weapons</b> , <b>crime</b> , <b>fbi</b> , control
ETM
<b>space</b> , <b>nasa</b> , gov, mr, president, health, research, year, center <b>windows</b> , <b>file</b> , <b>window</b> , <b>program</b> , <b>files</b> , <b>server</b> , <b>version</b> , dos, <b>image</b> <b>god</b> , people, <b>jesus</b> , <b>christian</b> , <b>israel</b> , bible, <b>jews</b> , <b>christians</b> , <b>israeli</b> <b>key</b> , <b>encryption</b> , chip, clipper, <b>keys</b> , <b>privacy</b> , <b>security</b> , technology, government <b>gun</b> , people, government, <b>law</b> , state, <b>guns</b> , article, <b>weapons</b> , control

Table 3.3: Top-9 words for each topic from 5 topics selected

### 3.5.9 Discussion

We have compared our model with a number of state-of-the-art models on multiple data sets. As we mentioned in previous section, perplexity may not be a good measure for comparing the topic models. Also, normally the perplexity for topic model are calculated using the posterior distribution, while the topic model using Autoencoding Variational Bayes(AEVB) are calculated using the variational lower bound[37]. The results between these models may not able to draw conclusion easily. We have carefully compared the models on topic-model specific metrics like topic coherence and topic diversity. The result exhibits our model has a outstanding performance on obtaining high quality topic words over various of predefined topic numbers, especially in topic coherence score. In two of the metric we compared, topic coherence and topic diversity, our model shows its strength in generating high quality topics and diversified words. The comparison on the quantitative result also demonstrate that our model is robust to generate high quality topics upon the topic number defined increases. In the t-SNE visualization, our model also made a clearer representation on 2-dimensional space which implies a better intra-cluster divergence the model distinguish the difference between topics.



**Figure 3.6:** Visualization #Topics:50

## Chapter 4

# Times Series Topic Retrieval with TECTM

In most of the real-life cases, the context (or formally topic information), that to be mentioned in the media such as news and documents are changing over time. Also, the word meaning and slang at the time may not valid on other span of time. In previous chapter 3, we have implemented a topic model that could capture high-quality topic words. However, the model cannot not distinct the difference of topic representation and the topic-word relation along the time of model. In this chapter, we expand the embedded topic model to deal with times-series task, namely Dynamic Transformer embedded topic model. The model utilizes the time information.

### 4.1 Introduction

Most of the existing topic models are designed for handle unified document set, and such no time specified information are assumed forehand. Dynamic Topic Model(DTM)[8] is a extended version from LDA, which a state-space model introduce latent variable to control the model parameters and proportion over time stamp  $1 \dots T$ . However, the model does not handle correlation information for topics obtained. Dynamic Correlated Topic Model(DCTM)[41] have considered the correlation information between topics though out the time. The author proposed . However, the model estimate the mean of the prior distribution of document-topic proportion with bag-of-words per-document, which may miss out to explore the time-topic relation from exploiting the input of bag-of-word per-time period. Our model demonstrate a better strength in obtaining topic coherence scores. Dynamic Embedded Topic Model(DETM)[13] constructed a times-series Embedded Topic Model by constructing the state-space model to guide the latent variables that control the mean for topic-proportion prior and the topic-word proportion. Likewise, the model does not consider correlation information. Besides, our model refines the word embedding quality on the latent space and hence improves overall topic-word predictive performance.

For this reason, we propose Dynamic Transformer Embedding Correlated Topic Model(DTECTM), a model that we build on top of the model from chapter 3. That is, a model that make use of the transformer embedding to compute the topic words from the document set. Additionally, we manage to capture the time-series information. By using Gaussian Process Latent Variable Model(GPLVM), we infer the topic distribution in different time from the latent variables it obtains. Our model takes the information from GPLVM as residual input and bag-of-words to infer the document-topic proportion. The results are compared with data sets against other state-of-the-art models.

## 4.2 Related works

Moreover, time-series model is one of the most practical for real application in topic modeling, which it is essential to extract keywords along the time line. In chapter 4 we will explore and implement . As an introductory, here we discuss the related works in advance. Blei [8] extended a time-series topic model on top of LDA, namely dynamic topic model(DTM). The model assumes the prior for topic-word proportion is an Markov state-space model along time  $t$ . The posterior is approximated with variational bayes and Kalman filter inference. Henning[17] proposed Kernel topic model, which the model is equipped the gaussian process with kernel covariance as the hyperparameter of Dirichlet prior for document-topic proportion. Tomasi[41] implemented a time-series correlated topic model, using Gaussian process to model for modeling the hyperparameter of topic-word proportion and the mean for document-topic proportion, along with using Wishart process for parametrizing the covariance matrix. Dieng[14] improve a model which on top of the ProDLDA topic model, implemented Word2Vec semantics to further improve the performance on topic coherence and predictive distribution. And respectively, the same author[13] extended the dynamic embedded topic model from previous embedded topic model. The model perform inference on posterior by deriving variational lower bound and amortized inference. The model take assume of topic-word proportion and document topic proportion as a Markov state-space model.

## 4.3 Model description

The DTECTM utilizes the Transformer as embedding to the topic-word representations. To compare with the original topic model, the word-topic distribution  $\beta$  is the dot-product of the Transformer embedding  $\rho$  and the topic embedding  $\alpha$ . First, the topic embedding  $\alpha$  embeds the vocabulary into L-dimensional space, which is by Transformer embedding  $\rho$ . Second, the context embedding maps the embedding into K-dimensional space. In the generative process, the LKJTM uses the topic embedding to form a per-topic vector to represent the meaning over the vocabulary.

To express the way word embedding to be applied in our model,  $\rho \in \mathbb{R}^{V \times L}$  is the Transformer embedding in the model.  $\rho_v$  is a vector represents the embedding of vocabulary on n-th index. To explain the way our model capture time-related information from document set, we here discuss variables change over time. The topic embedding  $\{\alpha_k^{(t)}\}_{t=1}^T \in \mathbb{R}^L$  is a vector distributed at specific k topic. Topic proportion  $\theta_d$  is same as typical topic model, which a simply a vector represent proportion for each topic on document d. The latent variable  $\eta$  decide the topic proportion holds on each timestamp ranged between  $1, \dots, T$ .

### 4.3.1 Generative Model

The generative process is as following:

---

**Algorithm 6:** Generative Process for DTECTM

---

```

1 Initialize hyperparameters
2 Obtain transformer embedding  $\rho$ 
3 for time  $t$  in  $T$  do
4   | Draw topic embedding  $\alpha^{(t)} \sim \mathcal{N}(\alpha^{(t-1)}, \xi^2 I)$ 
5   | Draw topic proportion mean  $\eta_t \sim \mathcal{N}(0, I)$ 
6   | Sample correlation  $L_t \sim \text{LKJChol}(\gamma_t)$ 
7 end
8 for document  $d$  in  $D$  do
9   | Sample topic proportion  $\theta_{t_d} \sim \mathcal{LN}(\eta_{t_d}, \Sigma_{t_d})$ 
10  for word position  $n$  in  $N_d$  do
11    | Sample word  $w_{d,n} \sim \text{Cat}(\sigma(\theta_{t_d}^\top (\rho^\top \alpha^{(t_d)})))$ 
12  end
13 end

```

---

From algorithm 6, first the model draws a topic embedding  $\alpha^{1:T}$  from normal distribution at time  $1, \dots, T$ . At time step 0, the topic embedding initialized at  $\mathcal{N}(0, I)$ . Then a topic mean  $\eta_t \in \mathbb{R}^K$  over timestamps is generated from the Gaussian Process Latent Variable Model (GPLVM), which performs inference a dimensions of topic K from number of vocabularies V dimension. Specifically, taking bag-of-word by time  $w_t$ , which is collected by categorizing the document by time and group them into word count matrix by timestamp. And then a normalization is performed to make sure the words in different timestamp are in same proportion. For each document, draw a topic proportion  $\theta_d$  from logistic-normal distribution  $\mathcal{LN}(\cdot, \cdot)$  condition on topic mean  $\eta_{t_d}$  at the timestamp t of document d, and the variance  $\xi^2 I$ . After that, for each word position n in  $N_d$ , a word is drawn from the dot-product of word embedding  $\rho$  and topic embedding  $\alpha_d^{(t_d)}$  at timestamp  $t_d$ .

### 4.3.2 Joint Distribution

To describe the joint distribution for the model, equation 4.1

$$\begin{aligned}
p(W, \theta, \Sigma, \eta, \alpha | \rho, \gamma) = & \prod_{d=1}^D \left[ p(\theta_d | \eta_{t_d}, \Sigma_{t_d}) \prod_{n=1}^{N_d} p(w_{d,n} | \theta_d, \rho, \alpha^{(t_d)}) \right] \setminus \\
& \prod_{t=1}^T \left[ p(\eta_t) p(\Sigma_t | \gamma_t) \prod_{k=1}^K p(\alpha_k^{(t)} | \alpha_k^{(t-1)}) \right] \quad (4.1)
\end{aligned}$$

For the bag-of-word input, we have V vocabularies over D documents. A number of K topics are defined to be introduced in the model. Each document is labeled a timestamp t over a time span between  $1, \dots, T$ .  $\theta_d \in \mathbb{R}^K$  is the topic proportion for document d.  $\eta_t \in \mathbb{R}^K$  is the topic proportion corresponds to time t of document d.  $w_d \in \mathbb{R}^V$  is the bag-of-word of distribution at document d.  $w_t \in \mathbb{R}^V$  is the bag-of-word of distribution at time t, which we arranged the documents in different timestamp and group them into bag-of-word representation in dimension  $\mathbb{R}^{V \times T}$ . In particular,  $w_t$  is normalized to attain the same ratio of tokens on each single time stamp.  $\alpha_k^{(t)} \in \mathbb{R}^V$  is the topic embedding for topic k at time t. The topic embedding demonstrates the vocabulary representations on the

specific time stamp for the document.  $\rho \in \mathbb{R}^{W \times L}$  is the transformer embedding that maps the words into L dimension of continuous latent space.

$$p(w_{d,n}|\theta_d, \rho, \alpha^{(t_d)}) = \text{Cat} \left( \sigma \left( \sum_{k=1}^K \theta_{dk} (\rho^\top \alpha^{(t_d)})_{k, w_{dn}} \right) \right) \quad (4.2)$$

The  $p(\Sigma_{t_d}|\gamma_{t_d})$  draws the covariance prior to the document-topic proportion  $\theta_d$ . The LKJ correlation prior generates a positive definite matrix of  $\mathbb{R}^{K \times K}$

$$p(\Sigma_t|\gamma_t) = \text{LKJCorr}(\gamma_t) \quad (4.3)$$

where  $\Sigma_{t_d}$  can be decomposed in form of the product of two lower triangular matrices. And such the correlation prior distribution can be factorized as a LKJCorrChol parameterized by the concentration variable  $\gamma_t$ , where it produces lower triangular matrices  $L_t$  that can be converted into covariance matrix by transformation mentioned in 2.4.

$$\text{LKJCorr}(\gamma_t) = \Sigma_t = L_t L_t^\top \quad (4.4)$$

$$L_t \sim \text{LKJCorrChol}(\gamma_t) \quad (4.5)$$

The document-topic proportion  $\theta_d$  is a logistic-normal distribution that conditioned on the mean prior  $\eta_t$  and covariance matrix  $\Sigma_t$ . The distribution ensure the proportion vector  $\theta_d \in \mathbb{R}^K$  in every document d are constrained into a simplex that values are summed to 1.

$$p(\theta_d|\eta_{t_d}, w_d) = \mathcal{LN}(\eta_{t_d}, \Sigma_{t_d}) \quad (4.6)$$

The time-to-topic proportion is a Gaussian Process Latent Variable Model(GPLVM), which takes an input of latent variable from lower dimension space to recover the observed data. We take the latent variable to induce the relation between time-topic representation. The  $\eta_{1:T} \in \mathbb{R}^K$  is a latent variables learnt from GPLVM that contributes the time-topic proportion as a mean for variable  $\theta$ . We can use the latent variable  $\eta$  to recover the observed data  $w_{1:T}$ . The latent function is then computed using the kernel  $\mathcal{K}_\theta$  parameterized by the hyperparameter  $\theta$  and the information from  $\eta$ .  $\mathcal{K}_\theta \in \mathbb{R}^{V \times V}$  is the kernel function that determine the shape of covariance matrix, which controls the shape of outcome distribution. Then the observed bag-of-words  $w$  is recovered back by normal distribution taking mean of  $f_n(\eta_t)$  and variance  $\sigma^2$ .

$$p(\eta_t) = \mathcal{N}(0, I) \quad (4.7)$$

$$p(f_n|\eta_t, \theta) = \prod_{n=1}^V \mathcal{N}(0, \mathcal{K}_\theta) \quad (4.8)$$

$$p(w_{t,n}|f_n, \eta_t) = \prod_{t=1}^T \prod_{n=1}^V \mathcal{N}(f_n(\eta_t), \sigma^2) \quad (4.9)$$

The topic-word proportion is contributed as  $p(\alpha^{(t)}|\alpha^{(t-1)}) = \mathcal{N}(\alpha^{(t-1)}, \gamma^2 I)$ , which is a Markov chain conditioned on variance parameterization  $\xi^2$ . The topic embedding  $p(\alpha_k^{(t)}|\alpha_k^{(t-1)})$  is a random walk go from time step  $1 \dots T$ , each time step  $t$  take the prior mean from the sampled value from previous time step  $t-1$  and Gaussian noise  $\xi^2$  as variance.

$$p(\alpha_k^{(t)}|\alpha_k^{(t-1)}) = \mathcal{N}(\alpha_k^{(t-1)}, \xi^2 I) \quad (4.10)$$

## 4.4 Inference and Estimation

Since the posterior is intractable to compute, we apply variational inference to approximate the parameters for the log-marginal likelihood.

### 4.4.1 Variational Distribution

To begin with, we first setup variational distribution to approximate the parameter of the model.

$$q(\theta, \eta, \Sigma, \alpha) = \prod_{d=1}^D q(\theta_d | w_d, \eta_{t_d}, \Sigma_{t_d}) \prod_{t=1}^T q(\eta_t | w_t) q(\Sigma_t | \gamma_t) \prod_{k=1}^K q(\alpha_k^{(t)}) \quad (4.11)$$

The variational distribution for topic proportion  $q(\theta_d | \eta_{t_d}, w_d)$  is logistic-normal distribution. We applied amortized inference to approximate the model, which the mean and covariance matrix is generated by two inference network  $\mu_\phi$  and  $\sigma_\phi$  taking bag-of-word input  $w_d$  at each document  $d$  and residual input  $\eta_{t_d}$  end up  $\mathbb{R}^{D \times (V+K)}$  dimension in the input space. The output for time-topic proportion is applied to the residual connection on the amortized inference of document-topic proportion  $\theta$ . The inference network for variational parameter  $\phi$  take the residual input from stacked input: both  $\mathbb{R}^V$  bag-of-words vectors  $w_d$ , and the  $\mathbb{R}^K$  time-topic proportion  $\eta_{t_d}$ . Then we transform the using reparameterization trick. To perform amortized inference for  $\theta$ , we apply Layer Normalization[1] to normalize the input vectors. The LayerNorm performs normalization over features, which enables a better training. Also, it is more stable than batch normalization for training while the batch size is pretty small. And thus it is able to maintain a lower variance than batch normalization does throughout the training loop.

$$q(\theta_d | \eta_{t_d}, w_d, \Sigma_{t_d}) = \mathcal{N}(\mu_\phi(x), \Sigma_{t_d}) \quad (4.12)$$

$$x = \text{LayerNorm}([w_d, \eta_{t_d}]) \quad (4.13)$$

The variational distribution for  $q(\eta_t | w_t)$  is basically the normal distribution parametrized by two inference network, which takes the input from bag-of-word  $w_d$  and the variable from [25],  $L$  is the dimension for latent input space,  $V$  is the token size.

$$q(\eta_t | w_t) = \int \prod_{i=1}^V q(w_i) \prod_{d=1}^L p(f_d | u_d, W) q(u_d) du_d \quad (4.14)$$

The variational distribution for  $q(\alpha_k^{(t)})$  is the embedding is parameterized mean  $\mu_\varphi$  and  $\sigma_\varphi$  conditioned on local parameter  $\varphi$ . Notice that  $\varphi$  is not a variational parameter as usual amortized inference.

$$q(\alpha^{(t)}) = \mathcal{N}(\mu_\varphi^{(t)}, \sigma_\varphi^{(t)}) \quad (4.15)$$

both variational distribution  $q(\theta_d | \eta_{t_d}, w_d)$  and  $q(a^{(t)})$  applied reparameterization trick [22] with transformation  $\mathcal{N}(\mu, \sigma) \approx \mu + \epsilon \sigma^{1/2}, \epsilon \sim \mathcal{N}(0, 1)$ , to avoid high variance on the variational variables.

### 4.4.2 Evidence lower bound (ELBO)

We take the log marginal likelihood from eq. 4.1. Noted the given equation could derive the ELBO by simply applying the Jensen inequality. For sake of



simplicity, we decompose the marginal likelihood two parts, the document-topic proportion part and the time-topic proportion part. Then summing up the terms as eq. 4.16. In such way, we are able to derive a ELBO without calculating the KL-divergence for  $q(\eta_{t_d})$  and  $p(\eta_{t_d}|w_t)$ . And so we could obtain the ELBO for  $\log p(\eta_{t_d}|w_t)$  conveniently from [40].

$$\begin{aligned}
\mathcal{L} &\geq \mathbb{E}_q[\log p(W, \theta, \eta, \Sigma, \alpha|\rho, \gamma)] - \mathbb{E}_q[\log q(\theta, \eta, \Sigma, \alpha)] \\
&= \sum_{d=1}^D \sum_{n=1}^V \mathbb{E}_q[\log p(w_{d,n}|\theta_d, \rho, \alpha^{(t_d)})] + \sum_{d=1}^D \mathbb{E}[\log p(\theta_d|\eta_{t_d}, \Sigma_{t_d})] \\
&\quad + \sum_{t=1}^T \mathbb{E}_q[\log p(\eta_t)] + \sum_{t=1}^T \mathbb{E}_q[\log p(\Sigma_t)] + \sum_{t=1}^T \sum_{k=1}^K \mathbb{E}_q[\log p(\alpha_k^{(t)}|\alpha_k^{(t-1)})] \\
&\quad - \sum_{d=1}^D \mathbb{E}_q[\log q(\theta_d|\mu_\phi(x_d), \Sigma_t)] - \sum_{t=1}^T \mathbb{E}_q[\log q(\eta_t|w_t)] - \sum_{t=1}^T \mathbb{E}_q[\log q(\Sigma_t|\gamma_t)] \\
&\quad - \sum_{t=1}^T \sum_{k=1}^K \mathbb{E}_q[\log q(\alpha_k^{(t)}|\alpha_k^{(t-1)})] \tag{4.16}
\end{aligned}$$

and by rearranging the terms, the ELBO can be represented as follows, where the first term is the reconstruction loss, and the remaining is the KL-divergence between the prior and variational distribution for its variational parameters.

$$\begin{aligned}
&= \sum_{d=1}^D \sum_{n=1}^V \mathbb{E}_q[\log p(w_{d,n}|\theta_d, \rho, \alpha^{(t_d)})] - \sum_{d=1}^D \text{KL}(q(\theta_d|w_d, \eta_{t_d}, \Sigma_{t_d})||p(\theta_d|\eta_{t_d}, \Sigma_{t_d})) \\
&\quad - \sum_{t=1}^T \text{KL}(q(\eta_t|w_t)||p(\eta_t)) - \text{KL}(q(\Sigma_t|\gamma_t)||p(\Sigma_t)) \\
&\quad - \sum_{k=1}^K \text{KL}(q(\alpha_k^{(t)}|\alpha_k^{(t-1)})||p(\alpha_k^{(t)}|\alpha_k^{(t-1)})) \tag{4.17}
\end{aligned}$$

For sake of simplicity, we separate the ELBO into three parts. We first derive the ELBO for the document-topic model part, denoted as  $\mathcal{L}_1$ , we can obtain  $p(w_d, \theta_d, \alpha|\eta_{t_d}) = p(w_d|\theta_d, \alpha^{(t_d)})p(\theta_d|\eta_{t_d})p(\alpha)$  by factorization.

$$\begin{aligned}
\mathcal{L}_1 &= \sum_{d=1}^D \sum_{n=1}^V \mathbb{E}_q[\log p(w_{d,n}|\theta_d, \rho, \alpha^{(t_d)})] - \sum_{d=1}^D \text{KL}(q(\theta_d|w_d, \eta_{t_d}, \Sigma_{t_d})||p(\theta_d|w_d, \eta_{t_d}, \Sigma_{t_d})) \\
&\quad - \sum_{t=1}^T \text{KL}(q(\Sigma_t^\phi|\gamma_t)||p(\Sigma_t)) \tag{4.18}
\end{aligned}$$

To speed up computation, we apply mini-batching as previous chapter

$$\tilde{\mathcal{L}}_1 \approx \frac{|\mathcal{D}|}{|\mathcal{B}|} \sum_{d \in \mathcal{D}_B} \left[ \sum_{n=1}^V \mathbb{E}_q[\log p(w_{d,n}|\theta_d, \rho, \alpha^{(t_d)})] - \text{KL}(q(\theta_d|w_d, \eta_{t_d}, \Sigma_{t_d})||p(\theta_d|w_d, \eta_{t_d}, \Sigma_{t_d})) \right] \tag{4.19}$$

$$- \sum_{t=1}^T \text{KL}(q(\Sigma_t^\phi|\gamma_t)||p(\Sigma_t)) \tag{4.20}$$

The first term is the expected likelihood term for reconstructing the word  $w_{dn}$  from the model, where  $p(w_{d,n}|\theta_d, \rho, \alpha^{(t_d)})$  is the log likelihood probability parameterized by the variational parameters  $\theta_d, \alpha^{t_d}$  and transformer embedding  $\rho$ ,

which  $\sigma(\cdot)$  is the softmax function. The topic-word proportion is a dot product of transformer embedding  $\rho$  and  $\alpha^{(t_d)}$

$$\mathbb{E}_{q(\theta_d)q(\alpha)}[\log p(w_{d,n}|\theta_d, \rho, \alpha^{(t_d)})] = \mathbb{E}_{\theta_d \sim \mathcal{N}(\mu_\phi(x), \Sigma_{t_d}^\phi)}[w_{d,n} \text{Cat}(\sigma(\theta_d^\top (\rho^\top \alpha^{(t_d)})))] \quad (4.21)$$

and then apply the reparameterization trick to maintain a low-variance gradient estimate to the likelihood term, the transformation  $\theta_d = \mu_\phi(x) + (\Sigma_{t_d}^\phi)^{1/2}\epsilon$ , where  $\epsilon \sim \mathcal{N}(0, I)$  is gaussian noise variance.

$$= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ w_{d,n} \text{Cat}(\sigma((\mu_\phi(x_d) + (\Sigma_{t_d}^\phi)^{1/2}\epsilon)^\top (\rho^\top \alpha^{(t_d)}))) \right] \quad (4.22)$$

The second term is the KL-divergence between  $p(\theta_d)$  and  $q(\theta_d)$ , the distributions are logistic-normal distributions so can be represented in closed-form. So we simply derive the KL-divergence by substitution, the variable  $x$  is the parameter stacked with bag-of-word  $w_d$  and residual input  $\eta_{t_d}$  from 4.12

$$\begin{aligned} & \text{KL}(q(\theta_d|\mu_\phi(x), \Sigma_{t_d}^\phi) || p(\theta_d|\eta_{t_d}, \Sigma_{t_d})) \\ &= \frac{1}{2} \left( \text{tr}(\Sigma_{t_d}^{-1} \Sigma_{t_d}^\phi) + (\eta_{t_d} - \mu_\phi(x_d))^\top \Sigma_{t_d}^{-1} (\eta_{t_d} - \mu_\phi(x_d)) + \log \frac{|\Sigma_{t_d}|}{|\Sigma_{t_d}^\phi|} - K \right) \end{aligned} \quad (4.23)$$

For the second part, we derive the ELBO for GPLVM, which the detailed derivation for the ELBO has been discussed in the original paper [40] as equation 4.24.  $w = \{w_t\}_{t=1}^T \in \mathbb{R}^{T \times V}$  is the observed data, which bag-of-words with respect to the word count by that timestamps over the documents.  $\eta = \{\eta_t\}_{t=1}^T \in \mathbb{R}^{T \times K}$ , the latent variable distributes the topic proportion over time. Known that the dimension reduction is performed, as  $K \ll V$ , where the defined number of topic is supposed to be much smaller than the size of vocabularies.  $f_d$  is the latent function which takes  $M$  inducing point  $u_d$ , where  $M$  is the number of inducing points defined for the training process. And where  $u_d$  is conditioned at input locations  $Z \in \mathbb{R}^{M \times K}$ .  $\phi$  is the local variational parameters.  $\lambda$  is the global variational parameters.  $\sigma_w$  is the gaussian noise. Specifically, we only extract the terms  $\mathcal{L}_2$ , the kl-divergence for  $\eta$  and  $u$ , that to be contributed in the ELBO of the whole model.

$$\begin{aligned} \log p(w_t|\eta_{t_d}) &\geq \sum_{d=1}^V \sum_{t=1}^T \mathbb{E}_{q_\phi(\eta_t)} \mathbb{E}_{p(f_d|u_d, \eta_t)q_\lambda(u_d)} [\log \mathcal{N}(w_{d,t}; f_d(\eta_t), \sigma_w^2)] \\ &\quad - \underbrace{\sum_{t=1}^T \text{KL}(q_\phi(\eta_t) || p(\eta_t)) - \sum_{d=1}^V \text{KL}(q_\lambda(u_d) || p(u_d|Z))}_{\mathcal{L}_2} \end{aligned} \quad (4.24)$$

The KL-divergence for  $\alpha^{(t)}$  in time  $t$  is a closed-form in normal distribution. And so the equation can be derived as 4.25. The variational distribution for  $q(\alpha_k^{(k)})$  is parametrized by two inference network  $\mu_\varphi$  and  $\sigma_\varphi$ , where  $\varphi$  is a local variational parameter. And the prior  $p(\alpha_k^{(t)}|\alpha_k^{(t-1)})$  at time  $t$  takes the mean from previous step  $\alpha_k^{(t-1)}$  with variance  $\xi^2$ . In initial step, the mean for  $\alpha^{(1)}$  located at 0.

$$\mathcal{L}_3 = \text{KL}(q(\alpha_k^{(t)}) || p(\alpha_k^{(t)}|\alpha_k^{(t-1)})) = \frac{1}{2} \left( \log \frac{\xi^2}{\sigma_\varphi^2} + \frac{\sigma_\varphi^2 + (\mu_\varphi - \alpha_k^{(t-1)})^2}{\xi^2} - 1 \right) \quad (4.25)$$

By assembling the ELBO terms from above, and substituting the KL terms  $\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3$  from eq. 4.25, 4.24, 4.23 into 4.17, the ELBO equation becomes follows

$$\begin{aligned}
\log p(w|\theta, \alpha) &\geq \frac{|\mathcal{D}|}{|\mathcal{B}|} \sum_{d \in \mathcal{D}_{\mathcal{B}}} \sum_{n=1}^V \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ w_{d,n} \sigma((\mu_{\phi}(x_d) + (\Sigma_{t_d}^{\phi})^{1/2} \epsilon)^{\top} (\rho^{\top} \alpha^{(t_d)}))_{d,n} \right] \\
&\quad - \frac{1}{2} \frac{|\mathcal{D}|}{|\mathcal{B}|} \sum_{d \in \mathcal{D}_{\mathcal{B}}} \left( \text{tr}(\Sigma_{t_d}^{-1} \Sigma_{t_d}^{\phi}) + (\eta_{t_d} - \mu_{\phi}(x_d))^{\top} \Sigma_{t_d}^{-1} (\eta_{t_d} - \mu_{\phi}(x_d)) + \log \frac{|\Sigma_{t_d}|}{|\Sigma_{t_d}^{\phi}|} - K \right) \\
&\quad - \frac{1}{2} \sum_{t=1}^T \sum_{k=1}^K \left( \log \frac{\xi^2}{\sigma_{\varphi}^2} + \frac{\sigma_{\varphi}^2 + (\mu_{\varphi} - \alpha_k^{(t-1)})^2}{\xi^2} - 1 \right) \\
&\quad - \sum_{t=1}^T \text{KL}(q(\Sigma_t^{\phi}) || p(\Sigma_t | \gamma_t)) \\
&\quad - \sum_{t=1}^T \text{KL}(q_{\phi}(\eta_t) || p(\eta_t)) - \sum_{d=1}^V \text{KL}(q_{\lambda}(u_d) || p(u_d | Z)) \\
&= \mathcal{L}
\end{aligned} \tag{4.26}$$

By the above derivation of ELBO, we can compute the unbiased gradient with Monte Carlo sampling.

---

**Algorithm 7:** Training on DTECTM

---

```

1 Initialize weights, hyperparameters
2 for epoch  $1, \dots, N$  do
3   for time  $t$  in  $1 \dots T$  do
4     Sample topic embedding  $\alpha^{(t)}$  from eq. 4.15
5     Sample time-topic proportion  $\eta_t$  from eq. 4.14
6     Sample  $L_t \sim \text{LKJChol}(\gamma_t)$ 
7   end
8   Choose a minibatch  $\mathcal{B}$  of documents
9   for document  $d$  in minibatch do
10    Compute the topic proportion  $\theta_d$  from eq. 4.12
11    for word  $n$  in document  $d$  do
12      Sample the word  $w_{d,n}$  from eq.4.22
13    end
14  end
15  Estimate ELBO loss  $\mathcal{L}_{\text{ELBO}}$  from Eq. 4.26
16  Compute Transformer loss  $\mathcal{L}_{\text{CrossEntropy}}$ 
17  Compute the unbiased gradient estimate
18  Compute the stochastic gradient via backpropagation
19  Take a stochastic gradient step with Adam
20  Update the model and variational parameters
21 end

```

---

**Algorithm** The procedure for the model training is described in algorithm 7. To begin with, the parameters are initialized. For each epochs  $1, \dots, N$ , the topic embedding  $\alpha^{(t)}$  in every single time stamp  $t$  are computed. To perform stochastic variational inference, we divide data set into smaller data batch  $\mathcal{B}$ . For each document  $d$  in  $\mathcal{B}$ , we sample time information  $\eta_{t_d}$  from 4.14 to the specific time stamp  $t_d$  the document  $d$  belongs to. Then we compute the topic proportion  $\theta_d$ . For each word position  $n$  in document, a word is then to be drawn as  $w_{d,n}$

The ELBO loss is being computed by the sum of document-topic proportion part from equation 4.18 and time-topic proportion part from equation 4.24. Following that, the transformer loss is computed by cross entropy error, To optimizer the model, we compute unbiased gradient estimate from the model The procedure continue repeating until the maximum iteration is reached.

## 4.5 Experiment and results

In this chapter, Dynamic Correlated Topic Model(DCTM)[41], and Dynamic Embedded Topic Model(DETM)[13] to compare with our model.<sup>1</sup>

### 4.5.1 Experiment settings

**Datasets** We select the UN DEBATES as one of the testing corpus for the experiment. It is a collection of transcript from the official of UN member countries expressing about the government’s perspective over the world issues at the time. After preprocessing, It contains 46 years time span of data, with 7507 documents and 6831 tokens in total. We selected 6005 for training, 1402 documents for testing and 100 documents for validation. Second dataset we selected is NEURIPS conference paper dataset. The dataset contains conference papers ranged from 1987-2019. After preprocessing, it contains 9677 papers in total and 9182 tokens. Within the dataset, we pick 7345 documents for training, 1737 for testing and 100 for validation.

**Data pre-processing** To prepare the datasets for training, we pre-process the documents and turn them into useful corpus. We remove the special characters and stop words, and perform tokenization to split document sentences into a list of tokens. In order to train the model, we have to leverage the datasets to feed-in different model. Specifically, the data are shaped into different forms. First, the bag-of-word  $w_{1:D} \in \mathbb{R}^{D \times V}$  is a matrix consist of the word count for vocabularies exist in every document. The document frequency for tokens are set to 100 documents minimum and 50% at max. We also created a time-vocabularies word count matrix for the training of  $\eta$ . The dataset  $w_{1:T} \in \mathbb{R}^{T \times V}$  holds the word count for the vocabulary set over the time span  $t = 1, \dots, T$ .

**Transformer learning task** See 3.5.2.

**Models** We maintain the default settings for DCTM [41], with 0.1 kernel bandwidth and 1 for kernel amplitude.<sup>2</sup> For D-ETM model, we follow the default settings instructed in [13]. We set the variance of on the prior to 0.005.<sup>3</sup>

**Algorithm configurations** Following the parameter settings from [8], the variance of prior are set to  $\xi^2 = 0.005$  on  $\alpha \sim \mathcal{N}(\cdot, \cdot)$ . We set the dimension for the transformer embedding is 256, and so for the hidden dimension for both transformer encoder and decoder. Each transformer encoder and transformer decoder consist of two layers and 2 heads for scaled dot-product. The sequence length the transformer are set to 20. For the gaussian process latent variable model, we selected zero mean with squared exponential function (RBF

<sup>1</sup>DTM was tested in our experiment, nonetheless, did not yield result after 10 hours of running time.

<sup>2</sup><https://github.com/spotify-research/dctm/>

<sup>3</sup><https://github.com/adjidieng/DETM/>

kernel) as the covariance prior, which allows providing adaptive change on topic-proportion against the possible rapid topic changes from documents. The number of inducing point are set to 50. We follow the setting from [41] and set the length scale of kernel as 0.1.

#### 4.5.2 Results

**Training** In the training stage, we perform black box variational inference to estimate the unbiased gradient estimator with Monte Carlo sampling for intractable variational lower bound.

**Quantitative results** The result of models are to compare in terms of perplexity, Topic Coherence (TC), Topic Diversity (TD). To calculate the topic coherence and topic diversity, we average down the scores over time span  $1 \dots T$ . We put the models into UN debates and NIPS dataset for comparison.

In addition to both of the results above, we also put the model we construct from chapter 3 to compare with the time-series version models in this chapter. Since the model does not consider time series information, it only predicts a single set of topic words and thus for the TC and TD score.

	Perplexity	TC	TD
DCTM	4205.1	-0.017	<b>0.953</b>
D-ETM	<b>2842.1</b>	0.076	0.427
Our model(ch3)	2746.3	0.152	0.643
<b>Our model</b>	4103.8	<b>0.092</b>	0.429

**Table 4.1:** Result on UN debates dataset,  $k=30$

	Perplexity	TC	TD
DCTM	4898.4	-0.092	<b>0.966</b>
D-ETM	<b>2282.3</b>	0.120	0.517
Our model(ch3)	1960.8	0.201	0.733
<b>Our model</b>	5461.0	<b>0.136</b>	0.502

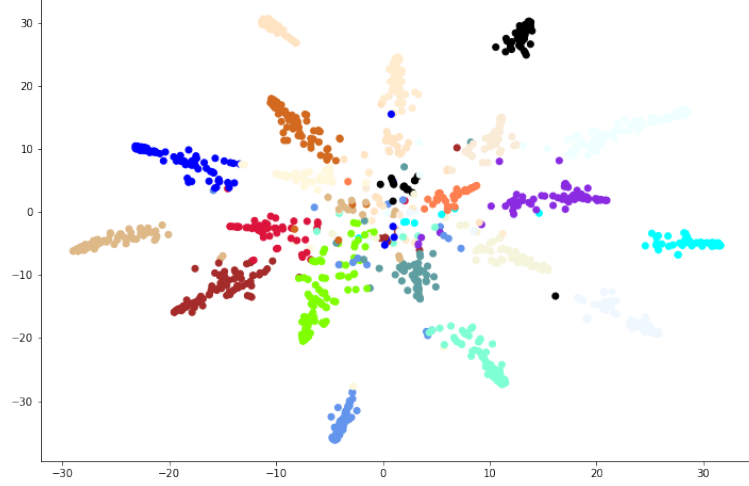
**Table 4.2:** Result on NeurIPS dataset (1987-2019),  $k=30$

From table 4.1, we have compared our model with the latest model. We only highlight the times-series model with the highest score. On the UN debate dataset, we conduct the experiment with number of topic  $k = 30$ . In terms of perplexity, the D-ETM model has performed a better score. On the other hand, our model obtain a better performance in terms of topic coherence and topic diversity, which is 0.092 and 0.429 respectively. In particular, UN debates dataset tends to contain more diversified topic mentions in different years and word usages by leaders from different countries. The DCTM holds the best topic diversity value, but the poor result in perplexity and topic coherence score make it meaningless. On NIPS dataset, table 4.2 contains the scores that we obtained from each model we conducted experiment. From the result, when  $k = 30$ , our model resulted 0.136 in topic coherence score and 0.502 in topic diversity. To compare with other instances, our model perform the best in topic coherence score. D-ETM has the lowest perplexity score, which perform the best in predictive task. On the other hand, our model has apparently beats the other models in both TC and TD scores.

## Qualitative results

**UN debates** We put the document-topic proportion obtained from the model and run t-SNE algorithm to transform the document-topic proportion into 2-dimensional continuous space. Shown in figure 4.1, different colors of dots represent a specific topic. As can be seen, apparently the documents can be classified in to interpretable cluster according to their topic.

From the result we obtained, we select top-6 topics have the highest proportion

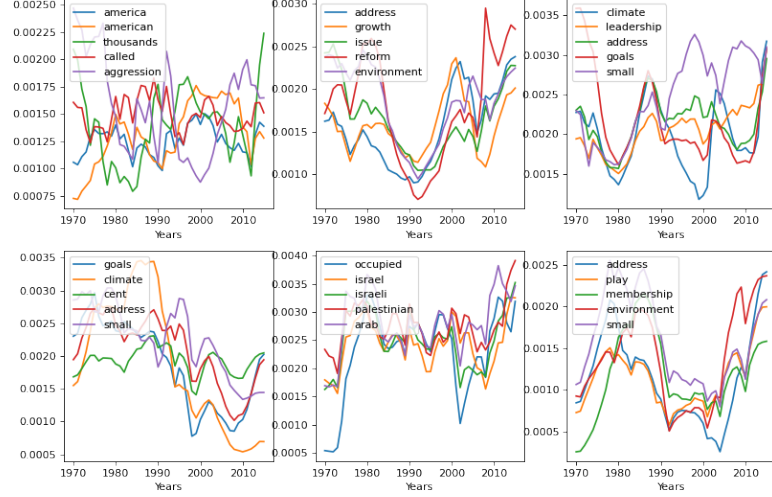


**Figure 4.1:** t-SNE visualization for documents labeled by topic(UN debates)

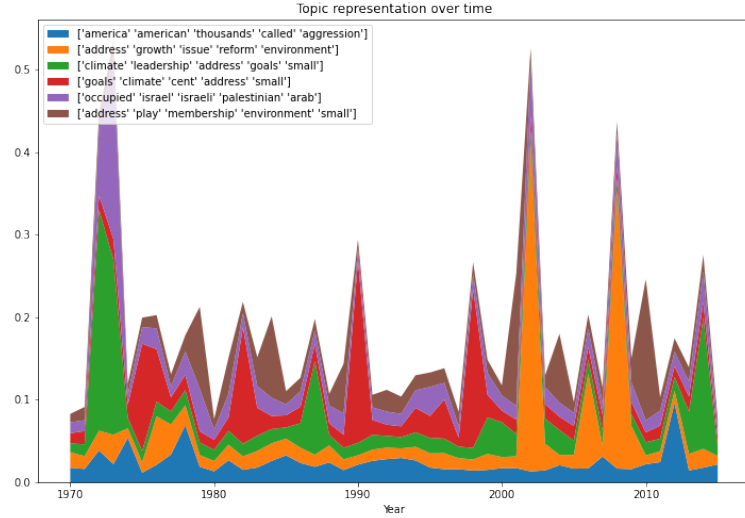
to the document set and display them on figure 4.2. And for those topics, we extract top-5 words having the highest proportion. Then we track those words how they change their proportion to the corresponding topic. As demonstrated, the topic words in each topic selected doesn't diverse each other very much in proportion, resulted a coherent trend that the word are more likely adhere to the topic.

Accordingly, in figure 4.3 we also track the change of those topic over time. Different color in the cumulative graph represents the topics evolve along the time period.

**NeurIPS dataset** To investigate the word trends change over time, table 4.3 visualize the words by 5 years interval. In particular, we selected a topic corresponding to reinforcement learning and pick top-10 words for each time stamp. By observation, we see that the topic are coherent in several keywords like CONTROL, ACTION and STATE. On the other hand, some other keywords also highlight their importance upon specific time period. For instance, words like MARKOV, POLICY and DISCOUNT are more likely to appear before; word REINFORCEMENT appears since 2012. For sake of interpretability, we also highlight the first appearance in the year for those keywords related to reinforcement learning. On figure4.4, displays the word trend from top-6 topics over the time. We have selected 3 representative tokens from top-10 words for each topic. And observe how the words trends though the years. On figure 4.5, we provide a stack plot for how



**Figure 4.2:** Word trend for top-6 topics (UN debates)

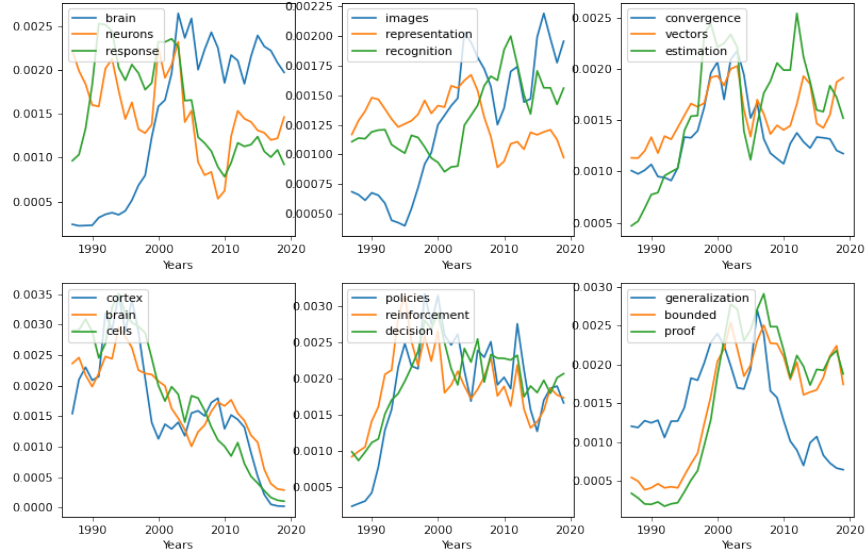


**Figure 4.3:** Topic trend for top-6 topics (UN debate)

those top-6 topics changed over time. Apparently it demonstrates how the topics are inclining or declining. For example, the topic related to "reinforcement" is gaining more popularity over time. Besides, the topic about "cortex, cells" is being less important by the years.

#### 4.5.3 Discussion

In the result, we observe that our model has been out perform the D-ETM model in TC and TD scores. Apparently the topic words in our model generated are both consistent and coherent over time. On the other hand, we notice that our model does not beat the state-of-the-art model in perplexity. Which we speculate that the model trade off the loss from training transformer bring down the predictive



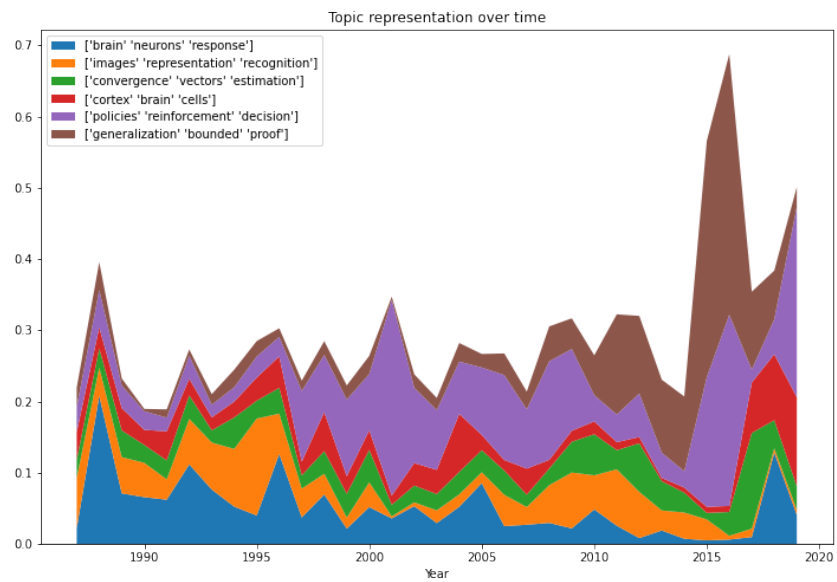
**Figure 4.4:** Word trend for top-6 topics (NeurlPS dataset)

Year	Topic: Reinforcement Learning
-	policy reward action goal control agent states actions reinforcement search
1987	control position world simulated initial search change environment modification action
1992	temporal watkins cambridge dynamic controller sutton states control actions action
1997	states actions markov decision control dynamic discount action discounted policy
2002	transition decision dynamic reward markov policies states actions policy action
2007	artificial intelligence rewards policies programming states action reward actions policy
2012	reinforcement decision markov states mdp action policies reward actions policy
2017	action markov policies reinforcement transition expected control states policy reward

**Table 4.3:** Word trend in topic reinforcement learning (5 years interval)

performance of the model. From analyzing the trend of the word and word group from each topic, we investigate how good the words obtain from each topic evolve over time. We convince that our model can obtain coherent topic from the document set and retrieve diverse keywords during the time span.





**Figure 4.5:** Topic trend for top-6 topics (NeurIPS dataset)

## Chapter 5

### Conclusions

In this thesis, we have proposed a state-of-the-art topic model, Transformer Embedded Correlated Topic Model(TECTM), and with its extension for time-series data, Dynamic Transformer Embedded Correlated Topic Model(DTECTM). The result has shown TECTM a better performance in returning high quality topics compared with other state-of-the-art models. Our model also demonstrates a capability in classifying substantial number of topics.

We also expanded the model to handle time-series data. The model was integrated with Gaussian process latent variable model, which make able the model to capture the time-series information from document set. We also expose our model has a competitive performance compared to other state of the art models.

In the studies though our thesis, we carried a series of experiments with varies of metrics to validate the models. Yet our model is capable to obtain a high quality topics in terms of topic coherence and topic diversity. However, we still notice that DTECTM model does not produce a good enough perplexity score compared with other models. So it is expected to improve the predictivity of the model as a future work.

Their are more improvement to the models. First, Nonparametric Bayesian method did not consider in the research due to the time limitation. Also, since graph model and document shares similarity on power law, we expect to explore the possibility to use graph machine learning knowledge to work with the model.

## References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. In *arXiv:1607.06450 [cs, stat]*.
- [2] John Barnard, Robert McCulloch, and Xiao-Li Meng. Modeling covariance matrices in terms of standard deviations and correlations, with application to shrinkage. pp. 1281–1311. Publisher: JSTOR.
- [3] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic language model. Vol. 3, pp. 1137–1155.
- [4] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal probabilistic programming. Vol. 20, No. 28, pp. 1–6.
- [5] David M. Blei. Probabilistic topic models. Vol. 55, No. 4, pp. 77–84.
- [6] David M. Blei and Michael I. Jordan. Variational inference for dirichlet process mixtures. Vol. 1, No. 1, pp. 121–143. International Society for Bayesian Analysis.
- [7] David M. Blei and John D. Lafferty. A correlated topic model of science. Vol. 1, No. 1, pp. 17–35. Institute of Mathematical Statistics.
- [8] David M. Blei and John D. Lafferty. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning, ICML '06*, pp. 113–120. Association for Computing Machinery.
- [9] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. Vol. 3, pp. 993–1022.
- [10] Jordan Boyd-Graber, Yuening Hu, and David Mimno. Applications of topic models. Vol. 11, No. 2, pp. 143–296.
- [11] Jonathan Chang and David Blei. Relational topic models for document networks. In *Artificial intelligence and statistics*, pp. 81–88. PMLR.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186. Association for Computational Linguistics.
- [13] Adji B. Dieng, Francisco J. R. Ruiz, and David M. Blei. The dynamic embedded topic model. In *arXiv:1907.05545 [cs, stat]*.

- [14] Adji B. Dieng, Francisco J. R. Ruiz, and David M. Blei. Topic modeling in embedding spaces. In *arXiv:1907.04907 [cs, stat]*.
- [15] Nick. Nicholas Charles. Raff Edward Eren, E. Maksim. Solovyev. COVID-19 literature clustering. event-place: University of Maryland Baltimore County (UMBC), Baltimore, MD, USA.
- [16] Junxian He, Zhiting Hu, Taylor Berg-Kirkpatrick, Ying Huang, and Eric P. Xing. Efficient correlated topic modeling with topic embedding. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 225–233. ACM.
- [17] Philipp Hennig, David Stern, Ralf Herbrich, and Thore Graepel. Kernel topic models. In *Artificial Intelligence and Statistics*, pp. 511–519. PMLR. ISSN: 1938-7228.
- [18] James Hensman, Nicolo Fusi, and Neil D Lawrence. Gaussian processes for big data. 2013.
- [19] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. Vol. 14, No. 5.
- [20] Thomas Hofmann. Probabilistic latent semantic analysis. In *arXiv:1301.6705 [cs, stat]*.
- [21] Donghwa Kim, Deokseong Seo, Suhyouon Cho, and Pilsung Kang. Multi-co-training for document classification using various document representations: TF-IDF, LDA, and doc2vec. Vol. 477, pp. 15–29.
- [22] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *arXiv:1312.6114 [cs, stat]*.
- [23] Thomas K Landauer and Susan T Dumais. A solution to plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. Vol. 104, No. 2, p. 211. American Psychological Association.
- [24] Neil D Lawrence. Gaussian process latent variable models for visualisation of high dimensional data. In *Advances in neural information processing systems*, Vol. 2, p. 5. Citeseer.
- [25] Neil D. Lawrence. Learning for larger datasets with the gaussian process latent variable model. In *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, pp. 243–250. PMLR. ISSN: 1938-7228.
- [26] Daniel D. Lee and H. Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. Vol. 401, No. 6755, pp. 788–791. Bandiera\_abtest: a Cg\_type: Nature Research Journals Number: 6755 Primary\_atype: Research Publisher: Nature Publishing Group.
- [27] Daniel Lewandowski, Dorota Kurowicka, and Harry Joe. Generating random correlation matrices based on vines and extended onion method. Vol. 100, No. 9, pp. 1989–2001.

- [28] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pretraining approach. In *arXiv:1907.11692 [cs]*.
- [29] Jon D. Mcauliffe and David M. Blei. Supervised topic models. In J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pp. 121–128. Curran Associates, Inc.
- [30] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119.
- [31] David Mimno, Hanna Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. Optimizing semantic coherence in topic models. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pp. 262–272. Association for Computational Linguistics.
- [32] Christopher E. Moody. Mixing dirichlet topic models and word embeddings to make lda2vec. In *arXiv:1605.02019 [cs]*.
- [33] David Newman, Jey Han Lau, Karl Grieser, and Timothy Baldwin. Automatic evaluation of topic coherence. In *Human language technologies: The 2010 annual conference of the North American chapter of the association for computational linguistics*, pp. 100–108.
- [34] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch.
- [35] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *arXiv:1802.05365 [cs]*.
- [36] Carl Edward Rasmussen and Christopher K. I. Williams. Gaussian processes for machine learning.
- [37] Akash Srivastava and Charles Sutton. Autoencoding variational inference for topic models. In *arXiv:1703.01488 [stat]*.
- [38] Yee Whye Teh. A tutorial on dirichlet processes and hierarchical dirichlet processes.
- [39] Yee Whye Teh, Michael I Jordan, Matthew J Beal, and David M Blei. Hierarchical dirichlet processes. Vol. 101, No. 476, pp. 1566–1581.
- [40] Michalis Titsias and Neil D Lawrence. Bayesian gaussian process latent variable model. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 844–851. JMLR Workshop and Conference Proceedings.
- [41] Federico Tomasi, Praveen Chandar, Gal Levy-Fix, Mounia Lalmas-Roelleke, and Zhenwen Dai. Stochastic variational inference for dynamic correlated topic models. In *Conference on Uncertainty in Artificial Intelligence*, pp. 859–868. PMLR.

- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, {\textbackslash}Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008.
- [43] Hanna M. Wallach. Topic modeling: beyond bag-of-words. In *Proceedings of the 23rd international conference on Machine learning - ICML '06*, pp. 977–984. ACM Press.
- [44] Weiqing Wang, Hongzhi Yin, Ling Chen, Yizhou Sun, Shazia Sadiq, and Xiaofang Zhou. ST-SAGE: A spatial-temporal sparse additive generative model for spatial item recommendation. Vol. 8, No. 3, pp. 48:1–48:25.
- [45] Xuerui Wang, Andrew McCallum, and Xing Wei. Topical n-grams: Phrase and topic discovery, with an application to information retrieval. In *Seventh IEEE International Conference on Data Mining (ICDM 2007)*, pp. 697–702. ISSN: 2374-8486.
- [46] Ke Xu, Yi Cai, Huaqing Min, Xushen Zheng, Haoran Xie, and Tak-Lam Wong. UIS-LDA: a user recommendation based on social connections and interests of users in uni-directional social networks. In *Proceedings of the International Conference on Web Intelligence*, pp. 260–265. ACM.
- [47] Guangxu Xun, Yaliang Li, Wayne Xin Zhao, Jing Gao, and Aidong Zhang. A correlated topic model using word embeddings. In *IJCAI*, pp. 4207–4213.
- [48] Liang Yang, Fan Wu, Junhua Gu, Chuan Wang, Xiaochun Cao, Di Jin, and Yuanfang Guo. Graph attention topic modeling network. In *Proceedings of The Web Conference 2020*, pp. 144–154. ACM.
- [49] 岩田具治. トピックモデル. 講談社.