

# Learning Document Embeddings Along With Their Uncertainties

Santosh Kesiraju , Oldřich Plchot, Lukáš Burget, and Suryakanth V. Gangashetty

**Abstract**—Majority of the text modeling techniques yield only point-estimates of document embeddings and lack in capturing the uncertainty of the estimates. These uncertainties give a notion of how well the embeddings represent a document. We present Bayesian subspace multinomial model (Bayesian SMM), a generative log-linear model that learns to represent documents in the form of Gaussian distributions, thereby encoding the uncertainty in its covariance. Additionally, in the proposed Bayesian SMM, we address a commonly encountered problem of intractability that appears during variational inference in mixed-logit models. We also present a generative Gaussian linear classifier for topic identification that exploits the uncertainty in document embeddings. Our intrinsic evaluation using perplexity measure shows that the proposed Bayesian SMM fits the unseen test data better as compared to the state-of-the-art neural variational document model on (*Fisher*) speech and (*20Newsgroups*) text corpora. Our topic identification experiments show that the proposed systems are robust to over-fitting on unseen test data. The topic ID results show that the proposed model outperforms state-of-the-art unsupervised topic models and achieve comparable results to the state-of-the-art fully supervised discriminative models.

**Index Terms**—Bayesian methods, embeddings, topic identification.

## I. INTRODUCTION

LEARNING word and document embeddings have proven to be useful in a wide range of information retrieval, speech and natural language processing applications [1]–[5].

These embeddings elicit the latent semantic relations present among the co-occurring words in a sentence or *bag-of-words* from a document. Majority of the techniques for learning these embeddings are based on two complementary ideologies, (i) topic modeling, and (ii) word prediction. The former methods are primarily built on top of the *bag-of-words* model and tend to capture higher-level semantics such as topics. The latter techniques capture lower-level semantics by exploiting the contextual information of words in a sequence [6]–[8].

Manuscript received October 4, 2019; revised February 23, 2020 and May 25, 2020; accepted July 15, 2020. Date of publication July 27, 2020; date of current version August 14, 2020. This work was supported by the Czech Ministry of Education, Youth and Sports from the National Programme of Sustainability (NPU II) project “IT4Innovations excellence in science LQ1602”. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Jianfeng Gao (*Corresponding author: Santosh Kesiraju.*)

Santosh Kesiraju is with the Brno University of Technology and International Institute of Information Technology, Hyderabad 500032, India (e-mail: kesiraju@fit.vutbr.cz).

Oldřich Plchot and Lukáš Burget are with the Brno University of Technology, 61200, Czechia (e-mail: ipchot@fit.vutbr.cz; burget@fit.vutbr.cz).

Suryakanth V. Gangashetty is with the International Institute of Information Technology, Hyderabad 500032, India (e-mail: svg@iiit.ac.in).

Digital Object Identifier 10.1109/TASLP.2020.3012062

On the other hand, there is a growing interest towards developing pre-trained language models [9], [10], that are then fine-tuned for specific tasks such as document classification, question answering, named entity recognition, etc. Although these models achieve state-of-the-art results in several NLP tasks; they require enormous computational resources to train [11].

Latent variable models [12] are a popular choice in unsupervised learning; where the observed data is assumed to be generated through the latent variables according to a stochastic process. The goal is then to estimate the model parameters, and also the latent variables. In probabilistic topic models (PTMs), the latent variables are attributed to topics, and the generative process assumes that every document is a distribution over topics and every topic is modeled as a distribution over words in the vocabulary [13]. Recent works showed that auto-encoders could also be seen as generative models for images and text [14], [15]. Generative models allow us to incorporate prior information about the latent variables, and with the help of variational Bayes (VB) techniques [14], [16], [17], one can infer a posterior distribution over the latent variables instead of just point-estimates. The posterior distribution captures the uncertainty of the latent variable estimates while trying to explain (fit) the observed data and our prior belief. In the context of text modeling, these latent variables are seen as embeddings.

In this paper, we present the Bayesian subspace multinomial model (Bayesian SMM) as a generative model for the *bag-of-words* representation of documents. We show that our model can learn to represent each document in the form of a Gaussian distribution, thereby encoding the uncertainty in its covariance. Further, we propose a generative Gaussian classifier that exploits this uncertainty for topic identification (ID). The proposed VB framework can be extended in a straightforward way for subspace  $n$ -gram model [18], that can model  $n$ -gram distribution of words in sentences.

Earlier, (non-Bayesian) SMM was used for learning document embeddings in an unsupervised fashion. They were then used for training linear classifiers for topic ID from spoken and textual documents [19], [20]. However, one of the limitations was that the learned document embeddings (also termed as document  $i$ -vectors) were only point-estimates and were prone to over-fitting, especially for shorter documents. Our proposed model can overcome this problem by capturing the uncertainty of the embeddings in the form of posterior distributions.

Given the significant prior research in PTMs and related algorithms for learning representations, it is crucial to draw precise relations between the presented model and previous works. We

do this from the following viewpoints: (a) graphical models illustrating the dependency of random and observed variables, (b) assumptions of distributions over random variables and their limitations, and (c) approximations made during the inference and their consequences.

The contributions of this paper are as follows: (a) we present Bayesian subspace multinomial model and analyze its relation to popular models such as latent Dirichlet allocation (LDA) [21], correlated topic model (CTM) [22], paragraph vector (PV-DBOW) [8] and neural variational document model (NVDM) [15], (b) we adapt tricks from [14] for faster and efficient variational inference of the proposed model, (c) we combine optimization techniques from [23], [24] and use them to train the proposed model, (d) we propose a generative Gaussian classifier that exploits uncertainty in the posterior distribution of document embeddings, (e) we provide experimental results on both text and speech data showing that the proposed document representations achieve state-of-the-art perplexity scores, and (f) with our proposed classification systems, we illustrate robustness of the model to over-fitting and at the same time obtain superior classification results when compared systems based on state-of-the-art unsupervised models.

We begin with the description of Bayesian SMM in Section II, followed by VB for the model in Section III. The complete VB training procedure and algorithm is presented in Section III-A. The procedure for inferring the document embedding posterior distributions for (unseen) documents is described in Section III-B. Section IV presents a generative Gaussian classifier that exploits the uncertainty encoded in document embedding posterior distributions. Relationship between Bayesian SMM and existing popular topic models is described in Section V. Experimental details are given in Section VI, followed by results and analysis in Section VII. Finally, we conclude and discuss directions for future research in Section VIII.

## II. BAYESIAN SUBSPACE MULTINOMIAL MODEL

The following steps explain the generative process of our model. For each document, a  $K$ -dimensional latent vector  $\mathbf{w}$  is generated from isotropic Gaussian prior with mean  $\mathbf{0}$  and precision  $\lambda$ :

$$\mathbf{w} \sim \mathcal{N}(\mathbf{w} | \mathbf{0}, (\lambda \mathbf{I})^{-1}) \quad (1)$$

The latent vector  $\mathbf{w}$  is a low dimensional embedding ( $K \ll V$ ) of document-specific distribution of words, where  $V$  is the size of the vocabulary. More precisely, for each document, the  $V$ -dimensional vector of word probabilities  $\boldsymbol{\theta}$  is calculated as:

$$\boldsymbol{\eta} = \mathbf{m} + \mathbf{T} \mathbf{w}, \quad (2)$$

$$\boldsymbol{\theta} = \text{softmax}(\boldsymbol{\eta}), \quad (3)$$

where  $\{\mathbf{m}, \mathbf{T}\}$  are parameters of the model. The vector  $\mathbf{m}$  known as universal background model (or bias) represents log uni-gram probabilities of words.  $\mathbf{T}$  known as total variability (or weight) matrix [25], [26] is a low-rank matrix defining subspace of document-specific distributions.

Finally, for each document, a vector of word counts  $\mathbf{x}$  (*bag-of-words*) is sampled from Multinomial distribution:

$$\mathbf{x} \sim \text{Multi}(\boldsymbol{\theta}; N), \quad (4)$$

where  $N$  is the number of words in the document.

$\boldsymbol{\eta}$  from (3) represents the *natural parameters* of the Multinomial distribution. Further, we can see that our model is linear in the space of natural parameters (2). Note that the parameters of any probability distribution under the *exponential family* can be expressed in terms of its *natural parameters* [16].

The above described generative process fully defines our Bayesian model, which we will now use to address the following problems: given training data  $\mathbf{X}$ , we estimate model parameters  $\{\mathbf{m}, \mathbf{T}\}$  and, for any given document  $\hat{\mathbf{x}}$ , we infer a posterior distribution over the corresponding document embedding  $p(\mathbf{w} | \hat{\mathbf{x}})$ . Parameters of such posterior distribution can be then used as a low dimensional representation of the document. Note that such distribution also encodes the inferred uncertainty about such representation.

Using Bayes' rule, the posterior distribution of document embedding  $\mathbf{w}$  is written as:<sup>1</sup>

$$p(\mathbf{w} | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{x} | \mathbf{w}) p(\mathbf{w}) d\mathbf{w}}. \quad (5)$$

In numerator of (5),  $p(\mathbf{w})$  represents the prior distribution of document embeddings (1), and  $p(\mathbf{x} | \mathbf{w})$  represents the likelihood of observed data. According to our generative process, we assume that every document  $\mathbf{x}$  is a sample from Multinomial distribution (4); hence the log-likelihood is computed as follows:

$$\log p(\mathbf{x} | \mathbf{w}) = \sum_{i=1}^V x_i \log \theta_i, \quad (6)$$

$$= \sum_{i=1}^V x_i \log \left( \frac{\exp\{m_i + \mathbf{t}_i \mathbf{w}\}}{\sum_j \exp\{m_j + \mathbf{t}_j \mathbf{w}\}} \right), \quad (7)$$

$$= \sum_{i=1}^V x_i \left[ (m_i + \mathbf{t}_i \mathbf{w}) - \log \left( \sum_{j=1}^V \exp\{m_j + \mathbf{t}_j \mathbf{w}\} \right) \right], \quad (8)$$

where  $\mathbf{t}_i$  represents a row in matrix  $\mathbf{T}$ . The problem arises while computing the denominator in (5). It involves solving the integral over a product of likelihood term containing the softmax function and Gaussian distribution (prior). There exists no analytical form for this integral. This intractability is a generic problem that arises while performing Bayesian inference for mixed-logit models [22], [27], multi-class logistic regression or any other model where the likelihood  $p(\mathbf{x} | \mathbf{w})$  and prior  $p(\mathbf{w})$  are not conjugate to each other [16]. In such cases, one can resort to variational inference and find an approximation to the posterior distribution  $p(\mathbf{w} | \mathbf{x})$ . This approximation to the true posterior

<sup>1</sup>For clarity, explicit conditioning on  $\mathbf{T}$  and  $\mathbf{m}$  is omitted in the subsequent equations.

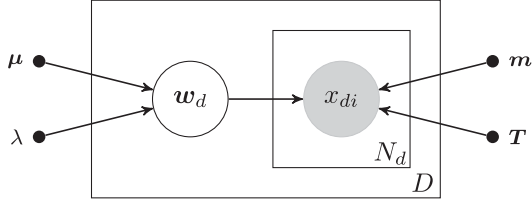


Fig. 1. Graphical model for Bayesian SMM. The arrows indicate the dependency. The rectangular plate with a symbol on the lower right corner denotes the number of repetitions of variables inside the plate. The shaded  $x_{di}$  represents the observed variable (word count),  $w_d$  represents the document-specific latent variable with hyperparameters  $\mu$  and  $\lambda$ . The variables  $m$  and  $T$  represent model parameters.

is referred as variational distribution  $q(w)$  and is obtained by minimizing the Kullback-Leibler (KL) divergence  $D_{\text{KL}}(q || p)$  from the approximate to the true posterior. But, computing the  $D_{\text{KL}}(q || p)$  also requires the functional form of true posterior  $p(w | x)$ , which is intractable. Hence, we take an alternative approach to minimize the KL divergence. We express the log marginal (evidence) of the data as:

$$\log p(x) = \mathbb{E}_q[\log p(x, w)] + H[q] + D_{\text{KL}}(q || p), \quad (9)$$

$$= \mathcal{L}(q) + D_{\text{KL}}(q || p). \quad (10)$$

Here  $H[q]$  represents the entropy of  $q(w)$ . Given the data  $x$ ,  $\log p(x)$  is a constant with respect to  $w$ , and  $D_{\text{KL}}(q || p)$  can be minimized by maximizing  $\mathcal{L}(q)$ , which is known as *Evidence Lower Bound* (ELBO) for a document. This is the standard formulation of variational Bayes [16], where the problem of finding an approximate posterior is transformed into the optimization of the functional  $\mathcal{L}(q)$ .

### III. VARIATIONAL BAYES

In this section, using the VB framework, we derive and explain the procedure for estimating model parameters  $\{m, T\}$  and inferring the variational distribution,  $q(w)$ . Before proceeding, we note that our model assumes that all documents and the corresponding document embeddings (latent variables) are independent. This independence can be seen from the graphical model in Fig. 1. Hence, we derive the inference only for one document embedding  $w$ , given an observed vector of word counts  $x$ .

We chose the variational distribution  $q(w)$  to be Gaussian, with mean  $\nu$  and precision  $\Gamma$ , i.e.,  $q(w) = \mathcal{N}(w | \nu, \Gamma^{-1})$ . The functional  $\mathcal{L}(q)$  now becomes:

$$\mathcal{L}(q) = \mathbb{E}_q[\log p(x, w)] + H[q], \quad (11)$$

$$= \mathbb{E}_q[\log p(x | w)] + \mathbb{E}_q[\log p(w)] + H[q], \quad (12)$$

$$= \underbrace{\mathbb{E}_q[\log p(x | w)]}_{\mathbf{A}} - \underbrace{D_{\text{KL}}(q || p)}_{\mathbf{B}} \quad (13)$$

The term  $\mathbf{B}$  in (13) is the KL divergence from the variational distribution  $q(w)$  to the document-independent prior (1), which can be computed analytically [28] as:

$$D_{\text{KL}}(q || p) = \frac{1}{2} [\lambda \text{tr}(\Gamma^{-1}) + \log |\Gamma| - K \log \lambda + \lambda \nu^T \nu - K], \quad (14)$$

where  $K$  denotes the document embedding dimensionality. The term  $\mathbf{A}$  from (13) is the expectation over log-likelihood of a document (8):

$$\mathbb{E}_q[\log p(x | w)] = \sum_{i=1}^V x_i \left[ (m_i + t_i \nu) - \underbrace{\mathbb{E}_q \left[ \log \left( \sum_{j=1}^V \exp\{m_j + t_j w\} \right) \right]}_{\mathcal{F}} \right]. \quad (15)$$

(15) involves solving the expectation over the log-sum-exp operation (denoted by  $\mathcal{F}$ ), which is intractable. It appears when dealing with variational inference in mixed-logit models [22], [27]. We can approximate  $\mathcal{F}$  with empirical expectation using samples from  $q(w)$ , but  $\mathcal{F}$  is a function of  $q(w)$ , whose parameters we are seeking by optimizing  $\mathcal{L}(q)$ . The corresponding gradients of  $\mathcal{L}(q)$  with respect to  $q(w)$  will exhibit high variance if we directly take samples from  $q(w)$  for the empirical expectation [29]. To overcome this, we will re-parametrize the random variable  $w$  by introducing a differentiable function  $g$  over another random variable  $\epsilon$  [14]. If  $p(\epsilon) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ , then:

$$w = g(\epsilon) = \nu + \mathbf{L} \epsilon, \quad (16)$$

where  $\mathbf{L}$  is the Cholesky factor of  $\Gamma^{-1}$ . Using this re-parametrization of  $w$ , we obtain the following approximation:

$$\mathcal{F} \approx \frac{1}{R} \sum_{r=1}^R \log \left( \sum_{j=1}^V \exp\{m_j + t_j g(\tilde{\epsilon}_r)\} \right), \quad (17)$$

where  $R$  denotes the total number of samples  $\tilde{\epsilon}_r$  from  $p(\epsilon)$ .

Combining (14), (15) and (17), we get the approximation to  $\mathcal{L}(q)$ . We will introduce the document suffix  $d$ , to make the notation explicit:

$$\begin{aligned} \mathcal{L}(q_d) \approx & -D_{\text{KL}}(q_d || p) \\ & + \sum_{i=1}^V x_{di} \left[ (m_i + t_i \nu_d) \right. \\ & \left. - \frac{1}{R} \sum_{r=1}^R \log \left( \sum_{j=1}^V \exp\{m_j + t_j g(\tilde{\epsilon}_{dr})\} \right) \right]. \end{aligned} \quad (18)$$

For the entire training data  $\mathbf{X}$ , the complete ELBO will be simply the summation over all the documents, i.e.,  $\sum_d \mathcal{L}(q_d)$ . Note that the KL divergence term in (18) is always non-negative and is independent of the observed data, hence acts as a regularization term for the embeddings.

#### A. Training

The variational Bayes (VB) training procedure for Bayesian SMM is stochastic because of the sampling involved in the re-parametrization trick (16). Like the standard VB approach [16], we optimize ELBO alternately with respect to  $q(w)$  and

$\{\mathbf{m}, \mathbf{T}\}$ . Since we do not have closed-form update equations, we perform gradient-based updates. Additionally, we regularize rows in matrix  $\mathbf{T}$  while optimizing. Thus, the final objective function becomes:

$$\mathcal{L} = \sum_{d=1}^D \mathcal{L}(q_d) - \omega \sum_{i=1}^V \|\mathbf{t}_i\|_1, \quad (19)$$

where we have added the term for  $\ell_1$  regularization of rows in matrix  $\mathbf{T}$ , with corresponding weight  $\omega$ . The same regularization was previously used for non-Bayesian SMM in [20]. This can also be seen as obtaining a maximum a posteriori estimate of  $\mathbf{T}$  with Laplace priors.

1) *Parameter Initialization*: The vector  $\mathbf{m}$  is initialized to log uni-gram probabilities estimated from training data. The values in matrix  $\mathbf{T}$  are randomly initialized from  $\mathcal{N}(0, 0.001)$ . The prior over latent variables  $p(\mathbf{w})$  is set to isotropic Gaussian distribution with mean  $\mathbf{0}$  and  $\lambda = \{1, 10\}$ . The variational distribution  $q(\mathbf{w})$  is initialized to  $\mathcal{N}(\mathbf{0}, (0.1)\mathbf{I})$ . Later in Section VII-A, we will show that initializing the posterior to a sharper Gaussian distribution helps to speed up the convergence.

2) *Optimization*: The gradient-based updates are done by ADAM optimization scheme [23]; in addition to the following tricks:

We simplified the variational distribution  $q(\mathbf{w})$  by making its precision matrix  $\mathbf{\Gamma}$  diagonal.<sup>2</sup> Further, while updating it, we used the log standard deviation parametrization, which ensures that the variance is always positive:

$$\mathbf{\Gamma}^{-1} = \text{diag}(\exp\{2\boldsymbol{\varsigma}\}). \quad (20)$$

The gradients of the objective (18) w.r.t. the mean  $\boldsymbol{\nu}$  is given as follows:

$$\nabla_{\boldsymbol{\nu}} \mathcal{L} = \left[ \sum_{i=1}^V \mathbf{t}_i^T (x_i - \frac{1}{R} \sum_{r=1}^R \theta_{ir} \sum_{k=1}^V x_k) \right] - \lambda \boldsymbol{\nu}, \quad (21)$$

where

$$\theta_{ir} = \frac{\exp\{m_i + \mathbf{t}_j g(\epsilon_r)\}}{\sum_j \exp\{m_j + \mathbf{t}_j g(\epsilon_r)\}}. \quad (22)$$

The gradient w.r.t log standard deviation  $\boldsymbol{\varsigma}$  is given as:

$$\begin{aligned} \nabla_{\boldsymbol{\varsigma}} \mathcal{L} &= \mathbf{1} - \lambda \exp\{2\boldsymbol{\varsigma}\} \\ &\quad - \sum_{k=1}^V x_k \frac{1}{R} \sum_{r=1}^R \sum_{i=1}^V \theta_{ir} \mathbf{t}_i^T \odot \exp\{\boldsymbol{\varsigma}\} \odot \epsilon_r, \end{aligned} \quad (23)$$

where  $\mathbf{1}$  represents a column vector of ones,  $\odot$  denotes element-wise product, and  $\exp$  is element-wise exponential operation.

The  $\ell_1$  regularization term makes the objective function (19) discontinuous (non-differentiable) at points where it crosses the orthant. Hence, we used sub-gradients and employed orthant-wise learning [24]. The gradient of the objective (19) w.r.t. a row  $\mathbf{t}_i$  in matrix  $\mathbf{T}$  is computed as follows:

<sup>2</sup>This is not a limitation of the model, but only a simplification.

---

### Algorithm 1: Stochastic VB Training.

---

```

1 initialize the model and the variational parameters
2 repeat
3   for  $d = 1 \dots D$  do
4     sample  $\tilde{\epsilon}_{dr} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   $r = 1 \dots R$ 
5     compute  $\mathcal{L}(q_d)$  using (18)
6     compute gradient  $\nabla_{\boldsymbol{\nu}_d} \mathcal{L}$  using (21)
7     compute gradient  $\nabla_{\boldsymbol{\varsigma}_d} \mathcal{L}$  using (23)
8     update  $\boldsymbol{\nu}_d$  and  $\boldsymbol{\varsigma}_d$  using ADAM
9   end
10  compute  $\mathcal{L}$  using (19)
11  compute sub-gradients  $\tilde{\nabla}_{\mathbf{t}_i} \mathcal{L}$  using (24) and (25)
12  update rows in  $\mathbf{T}$  using (26)
13 until convergence or max_iterations

```

---

$$\begin{aligned} \nabla_{\mathbf{t}_i} \mathcal{L} &= -\omega \text{sign}(\mathbf{t}_i) + \sum_{d=1}^D \left[ x_{di} \boldsymbol{\nu}_d^T \right. \\ &\quad \left. - \left[ \left( \sum_{k=1}^V x_{ki} \right) \frac{1}{R} \sum_{r=1}^R \theta_{dir} (\boldsymbol{\nu}_d^T + \epsilon_{dr}^T \odot \exp\{\boldsymbol{\varsigma}^T\}) \right] \right]. \end{aligned} \quad (24)$$

Here,  $\text{sign}$  and  $\exp$  operate element-wise. The sub-gradient  $\tilde{\nabla}_{\mathbf{t}_i} \mathcal{L}$  is defined as:

$$\tilde{\nabla}_{\mathbf{t}_i} \mathcal{L} \triangleq \begin{cases} \nabla_{\mathbf{t}_i} \mathcal{L} + \omega, & t_{ik} = 0, \nabla_{\mathbf{t}_i} \mathcal{L} < -\omega \\ \nabla_{\mathbf{t}_i} \mathcal{L} - \omega, & t_{ik} = 0, \nabla_{\mathbf{t}_i} \mathcal{L} > \omega \\ 0, & t_{ik} = 0, |\nabla_{\mathbf{t}_i} \mathcal{L}| \leq \omega \\ \nabla_{\mathbf{t}_i} \mathcal{L}, & |t_{ik}| > 0 \end{cases}. \quad (25)$$

Finally, the rows in matrix  $\mathbf{T}$  are updated according to:

$$\mathbf{t}_i \leftarrow \mathcal{P}_O(\mathbf{t}_i + \mathbf{d}_i) \quad (26)$$

where  $\mathbf{d}_i$  is the step in ascent direction:

$$\mathbf{d}_i = \eta \text{diag}(\sqrt{\mathbf{s}_i} + \epsilon)^{-1} \mathbf{f}_i. \quad (27)$$

Here,  $\eta$  is the learning rate,  $\mathbf{f}_i$  and  $\mathbf{s}_i$  represent bias-corrected first and second moments (as required by ADAM) of sub-gradient  $\tilde{\nabla}_{\mathbf{t}_i} \mathcal{L}$  respectively.  $\mathcal{P}_O$  represents orthant projection, which ensures that the update step does not cross the point of non-differentiability. It is defined as:

$$\mathcal{P}_O(\mathbf{t}_i + \mathbf{d}_i) \triangleq \begin{cases} 0 & \text{if } t_{ik}(t_{ik} + d_{ik}) < 0, \\ t_{ik} + d_{ik} & \text{otherwise.} \end{cases} \quad (28)$$

The orthant projection introduces explicit zeros in the estimated  $\mathbf{T}$  matrix and results in a sparse solution. Unlike in [20], we do not require to apply the sign projection, because both the gradient  $\tilde{\nabla}_{\mathbf{t}_i} \mathcal{L}$  and step  $\mathbf{d}$  have same sign (point to the same orthant). The stochastic VB training is outlined in Algorithm 1.

### B. Inferring Embeddings for New Documents

After obtaining the model parameters from VB training, we can infer (extract) the posterior distribution of document embedding  $q(\mathbf{w})$  for any given document  $\mathbf{x}$ . This is done by iteratively updating the parameters of  $q(\mathbf{w})$  that maximize  $\mathcal{L}(q)$  from (18). These updates are performed by following the same ADAM optimization scheme as in training.

Note that the embeddings are extracted by maximizing the ELBO, which does not involve any supervision (topic labels). These embeddings which are in the form of posterior distributions will be used as input features for training topic ID classifiers. Alternatively, one can use only the mean of the posterior distributions as point estimates of document embeddings.

#### IV. GAUSSIAN CLASSIFIER WITH UNCERTAINTY

In this section, we will present a generative Gaussian classifier that exploits the uncertainty in posterior distributions of document embedding. Moreover, it also exploits the same uncertainty while computing the posterior probability of class labels. The proposed classifier is called Gaussian linear classifier with uncertainty (GLCU) and is inspired by [30], [31]. It can be seen as an extension to the simple Gaussian linear classifier (GLC) [16].

Let  $\ell = 1 \dots L$  denote class labels,  $d = 1 \dots D$  represent document indices, and  $\mathbf{h}_d$  represent the class label of document  $d$  in one-hot encoding.

GLC assumes that every class is Gaussian distributed with a specific mean  $\boldsymbol{\mu}_\ell$ , and a shared precision matrix  $\mathbf{D}$ . Let  $\mathbf{M}$  denote a matrix of class means, with  $\boldsymbol{\mu}_\ell$  representing a column. GLC is described by the following model:

$$\mathbf{w}_d = \boldsymbol{\mu}_d + \boldsymbol{\varepsilon}_d, \quad (29)$$

where  $\boldsymbol{\mu}_d = \mathbf{M}\mathbf{h}_d$ ,  $p(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon} | \mathbf{0}, \mathbf{D}^{-1})$  and  $\mathbf{w}_d$  represent embedding for document  $d$ . GLC can be trained by estimating the parameters  $\Theta = \{\mathbf{M}, \mathbf{D}\}$  that maximize the class conditional likelihood of all training examples. For a single training example, the likelihood is computed as:

$$p(\mathbf{w}_d | \mathbf{h}_d, \Theta) = \mathcal{N}(\mathbf{w}_d | \boldsymbol{\mu}_d, \mathbf{D}^{-1}). \quad (30)$$

In our case, however, the training examples come in the form of posterior distributions,  $q(\mathbf{w}_d) = \mathcal{N}(\mathbf{w}_d | \boldsymbol{\nu}_d, \boldsymbol{\Gamma}_d^{-1})$  as extracted using our Bayesian SMM. In such case, the proper ML training procedure should maximize the expected class-conditional likelihood, with the expectation over  $\mathbf{w}_d$  calculated for each training example with respect to its posterior distribution  $q(\mathbf{w}_d)$ , i.e.,  $\mathbb{E}_q[\mathcal{N}(\mathbf{w}_d | \boldsymbol{\mu}_d, \mathbf{D}^{-1})]$ .

However, it is more convenient to introduce an equivalent model, where the observations are the means  $\boldsymbol{\nu}_d$  of the posteriors  $q(\mathbf{w}_d)$  and the uncertainty encoded in  $\boldsymbol{\Gamma}_d^{-1}$  is introduced into the model through the latent variable  $\mathbf{y}_d$  as:

$$\boldsymbol{\nu}_d = \boldsymbol{\mu}_d + \mathbf{y}_d + \boldsymbol{\varepsilon}_d, \quad (31)$$

where  $p(\mathbf{y}_d) = \mathcal{N}(\mathbf{y}_d | \mathbf{0}, \boldsymbol{\Gamma}_d^{-1})$ . The resulting model is called GLCU. Since the random variables  $\mathbf{y}_d$  and  $\boldsymbol{\varepsilon}_d$  are Gaussian-distributed, the resulting class conditional likelihood is obtained using the convolution of two Gaussians [16]:

$$p(\boldsymbol{\nu}_d | \mathbf{h}_d, \Theta) = \mathcal{N}(\boldsymbol{\nu}_d | \boldsymbol{\mu}_d, \boldsymbol{\Gamma}_d^{-1} + \mathbf{D}^{-1}). \quad (32)$$

The model parameters for both GLC and GLCU have the same interpretation, i.e., each class is Gaussian distributed with specific mean and a common precision matrix. The difference lies in the evaluation of the likelihood function (30) vs (32).

GLCU can be trained by estimating its parameters  $\Theta$ , that maximize the class conditional likelihood of training data (32). This can be done efficiently by using the following EM algorithm.

#### A. EM Algorithm

In the E-step, we calculate the posterior distribution of latent variables:

$$p(\mathbf{y}_d | \boldsymbol{\nu}_d, \mathbf{h}_d, \Theta) \propto p(\boldsymbol{\nu}_d | \mathbf{y}_d, \mathbf{h}_d, \Theta) p(\mathbf{y}_d) \propto \mathcal{N}(\mathbf{y}_d | \mathbf{u}_d, \mathbf{V}_d^{-1}), \quad (33)$$

where

$$\mathbf{V}_d = \mathbf{D} + \boldsymbol{\Gamma}_d, \quad (34)$$

$$\mathbf{u}_d = [\mathbf{I} + \mathbf{D}^{-1}\boldsymbol{\Gamma}_d]^{-1}(\boldsymbol{\nu}_d - \boldsymbol{\mu}_d). \quad (35)$$

In the M-step, we maximize the auxiliary function  $\mathcal{Q}$  with respect to the model parameters  $\Theta$ . This auxiliary function is the expectation of log joint-probability with respect to  $p(\mathbf{y}_d | \boldsymbol{\nu}_d)$ :

$$\begin{aligned} \mathcal{Q} &= \mathbb{E}_p \left[ \sum_{d=1}^D \log p(\boldsymbol{\nu}_d, \mathbf{y}_d | \Theta) \right] \\ &= \frac{-D}{2} \log |\mathbf{D}| - \frac{1}{2} \left[ \sum_{d=1}^D (\text{tr}(\mathbf{D}\mathbf{V}_d^{-1}) + (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d))^T \right. \\ &\quad \left. \times \mathbf{D} (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d))) \right] + \text{const.} \end{aligned} \quad (36)$$

Maximizing the auxiliary function  $\mathcal{Q}$  w.r.t.  $\Theta$ , we have:

$$\boldsymbol{\mu}_\ell := \frac{1}{|\mathcal{I}_\ell|} \sum_{d \in \mathcal{I}_\ell} (\boldsymbol{\nu}_d - \mathbf{u}_d) \quad \forall \ell = 1 \dots L \quad (38)$$

$$\mathbf{D}^{-1} := \frac{1}{D} \left[ \sum_{d=1}^D (\mathbf{a}_d \mathbf{a}_d^T) + \mathbf{V}_d^{-1} \right], \quad (39)$$

where  $\mathbf{a}_d = \mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d)$ , and  $\mathcal{I}_\ell$  is the set of documents from the class  $\ell$ . To train the GLCU model, we alternate between E-step and M-step until convergence.

#### B. Classification

Given the posterior distribution of a test document embedding  $q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \boldsymbol{\nu}, \boldsymbol{\Gamma}^{-1})$ , we compute the class conditional likelihood according to (32), and the posterior probability of a class  $\mathcal{C}_k$  is obtained by applying the Bayes' rule:

$$p(\mathcal{C}_k | \boldsymbol{\nu}, \boldsymbol{\Gamma}, \Theta) = \frac{p(\boldsymbol{\nu} | \boldsymbol{\mu}_k, \mathbf{D}, \boldsymbol{\Gamma}) p(\mathcal{C}_k)}{\sum_\ell p(\boldsymbol{\nu} | \boldsymbol{\mu}_\ell, \mathbf{D}, \boldsymbol{\Gamma}) p(\mathcal{C}_\ell)}. \quad (40)$$

#### V. RELATED MODELS

In this section, we review and relate some of the popular PTMs and neural network-based document models. We begin with a brief review of LDA [21], a probabilistic generative model for the *bag-of-words* representation of documents.

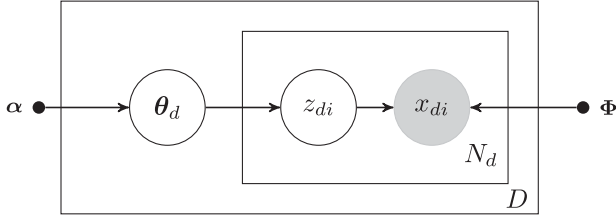


Fig. 2. Graphical model for LDA, where  $x_{di}$  is the observed variable, and  $\theta_d$ ,  $z_{di}$  are the latent variables.

### A. Latent Dirichlet Allocation

Let  $\Phi_{K \times V}$  represent  $K$  topics, where every (row) topic  $\phi_k$  is a discrete distribution over a fixed vocabulary of size  $V$ . LDA assumes that every document  $d$  is generated by a two-step process: first, a document-specific vector (embedding) representing a discrete distribution over  $K$  topics is sampled, i.e.,  $\theta_d \sim \text{Dir}(\alpha)$ . Then, for each word in document  $d$ , a topic indicator variable  $z_i$  is sampled:  $z_i \sim \text{Multi}(\theta_d; 1)$  and the word  $x_i$  is in turn sampled from the topic-specific distribution:  $x_i \sim \text{Multi}(\phi_{z_i}; 1)$ .

The topic ( $\phi$ ) and document ( $\theta$ ) vectors live in  $(V - 1)$  and  $(K - 1)$  simplexes, respectively. For every word  $x_i$  in document  $d$ , there is a discrete latent variable  $z_i$  that tells which topic was responsible for generating the word. This topic-word dependency can be seen from the graphical model in Fig. 2.

During inference, the generative process is inverted to obtain posterior distribution over latent variables,  $p(\theta, z | x, \alpha, \Phi)$ , given the observed data  $x$  and prior belief  $\alpha$ . Since the true posterior is intractable, Blei *et al.* [21] resorted to the variational inference, which finds an approximation to the true posterior as a variational distribution  $q(\theta, z)$ . Further, mean-field approximation was made to make the inference tractable, i.e.,  $q(\theta, z) = q(\theta) \prod_i q(z_i)$ .

In the original model proposed by Blei *et al.* [21], the parameters  $\Phi$  were obtained using the maximum likelihood approach. The choice of Dirichlet distribution for  $q(\theta)$  simplifies the inference process because of the Dirichlet-Multinomial conjugacy. However, the assumption of Dirichlet distribution causes limitations to the model, i.e.,  $q(\theta)$  cannot capture correlations between topics in each document. This was the motivation for Blei and Lafferty [22] to model documents with Gaussian distributions, and the resulting model is called correlated topic model (CTM).

### B. Correlated Topic Model

The generative process for a document in CTM [22] is the same as in LDA, except for document vectors are now drawn from Gaussian:

$$\eta \sim \mathcal{N}(\eta | \mu, (\lambda \mathbf{I})^{-1}), \quad (41)$$

$$\theta = \text{softmax}(\eta). \quad (42)$$

In this formulation, the document embeddings  $\eta$  are no longer in the  $(K - 1)$  simplex; instead they are dependent through the logistic normal. This is the same as in our proposed Bayesian SMM (1). The advantage is that the document vectors can model the correlations in topics. The topic distributions over vocabulary  $\Phi$ , however, still remained discrete. In Bayesian SMM, the

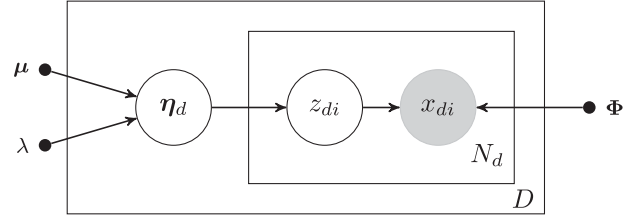


Fig. 3. Graphical model for CTM, where  $x_{di}$  is the observed variable, and  $\eta_d$ ,  $z_{di}$  are the latent variables.

topic-word distributions ( $T$ ) are not discrete; hence it can model the correlations between words and (latent) topics [22]. The variational inference in CTM is similar to that of LDA, including the mean-field approximation, because of the discrete latent variable  $z$  (Fig. 3). An additional problem is dealing with the non-conjugacy. More specifically, it is the intractability while solving the expectation over log-sum-exp function (see  $\mathcal{F}$  from (15)). Blei and Lafferty [22] used Jensen's inequality to form an upper bound on  $\mathcal{F}$ , and this in-turn acted as a lower bound on ELBO. In our proposed Bayesian SMM, we also encountered the same problem, and we approximated  $\mathcal{F}$  using the re-parametrization trick (Section III). There exist similar approximation techniques based on quasi-Monte Carlo sampling [27].

Unlike in LDA or CTM, Bayesian SMM does not require to make a mean-field approximation, because the topic-word mixture is not discrete, thus eliminating the need for discrete latent variable  $z$ .

### C. Subspace Multinomial Model

SMM is a log-linear model; originally proposed for modeling discrete prosodic features for the task of speaker verification [25]. Later, it was used for phonotactic language recognition [32] and eventually for topic identification and document clustering [19], [20]. A similar model was proposed by Maas *et al.* [33] for unsupervised learning of word representations. One of the significant differences among these works is the type of regularization used for matrix  $T$ .

Another difference is in obtaining embeddings  $w_d$  for a given test document. Maas *et al.* [33] obtained them by projecting the vector of word counts  $x_d$  onto the matrix  $T$ , i.e.,  $w_d = T x_d$ , whereas the authors from [19], [20] extracted the embeddings by maximizing regularized log-likelihood function. However, the embeddings extracted using SMM are prone to over-fitting, especially when the observed documents are short. Our Bayesian SMM overcomes this problem by capturing the uncertainty of document embeddings in the posterior distribution. Our experimental analysis in Section VII-C illustrates the robustness of Bayesian SMM.

### D. Paragraph Vector

Paragraph vector bag-of-words (PV-DBOW) [8] is also a log-linear model, which is trained stochastically to maximize the likelihood of a set of words from a given document. SMM can be seen as a special case of PV-DBOW since it maximizes the likelihood of all the words in a document.

TABLE I  
DATA SPLITS FROM FISHER PHASE 1 CORPUS, WHERE EACH DOCUMENT  
REPRESENTS ONE SIDE OF THE CONVERSATION

| Set               | # docs. | Duration (hrs.) |
|-------------------|---------|-----------------|
| ASR training      | 6208    | 553             |
| Topic ID training | 2748    | 244             |
| Topic ID test     | 2744    | 226             |

### E. Neural Network-Based Models

Neural variational document model (NVDM) is an adaptation of variational auto-encoders for document modeling [15]. The encoder models the posterior distribution of latent variables given the input, i.e.,  $p_{\theta}(z|x)$ , and the decoder models distribution of input data given the latent variable, i.e.,  $p_{\theta}(x|z)$ . In NVDM, the authors used *bag-of-words* as input, while their encoder and decoders are two-layer feed-forward neural networks. The decoder part of NVDM is similar to Bayesian SMM, as both the models maximize expected log-likelihood of the data, assuming Multinomial distribution. In simple terms, Bayesian SMM is a decoder with a single feed-forward layer. For a given test document, in NVDM, the approximate posterior distribution of latent variables is obtained directly by forward propagating through the encoder; whereas in Bayesian SMM, it is obtained by iteratively optimizing ELBO. The experiments in Section VII show that the posterior distributions obtained from Bayesian SMM represent the data better as compared to the ones obtained directly from the encoder of NVDM.

### F. Sparsity in Topic Models

Sparsity is often one of the desired properties in topic models [34], [35]. Sparse coding inspired topic model was proposed in [36], where the authors have obtained sparse representations for both documents and words.  $\ell_1$  regularization over rows in  $T$  matrix of SMM ( $\ell_1$  SMM) was observed to yield better results when compared to LDA, STC and  $\ell_2$  regularized SMM ( $\ell_2$  SMM) [20]. The relation between SMM and sparse additive generative model (SAGE) [34] was explained in [19]. In [37], the authors proposed an algorithm to obtain sparse document embeddings (called sparse composite document vector (SCDV)) from pre-trained word embeddings. In our proposed Bayesian SMM, we introduce sparsity into the model parameters  $T$  by applying  $\ell_1$  regularization and using orthant-wise learning.

## VI. EXPERIMENTS

### A. Datasets

We have conducted experiments on two benchmark datasets [19], [38], [39] from speech and NLP communities. The first one is *Fisher* speech corpus,<sup>3</sup> which is a collection of 5850 conversational telephone speech recordings with a closed set of 40 topics. Each conversation is approximately 10 minutes long with two sides of the call and is supposedly about one topic. We considered each side of the call (recording) as an independent document, which resulted in a total of 11700 documents. Table I

presents the details of data splits; they are the same as used in earlier research [19], [40], [41]. Our pre-processing involved removing punctuation and special characters, but we did not remove any stop words. Using Kaldi open-source toolkit [42], we trained a sequence-discriminative DNN-HMM automatic speech recognizer (ASR) system [43] to obtain automatic transcriptions. The ASR system resulted in 18% word-error-rate on a held-out test set. We report experimental results on both manual and automatic transcriptions. The vocabulary size while using manual transcriptions was 24854, for automatic, it was 18292, and the average document length is 830, and 856 words respectively.

The text corpus used is *20Newsgroups*,<sup>4</sup> which is a benchmark dataset for evaluating topic models [15], [44]–[46]. It contains 11314 training and 7532 test documents over 20 topics. Our pre-processing involved removing punctuation and words that do not occur in at least two documents, which resulted in a vocabulary of 56433 words. The average document length is 290 words.

### B. Hyperparameters of Bayesian SMM

In our topic ID experiments, we observed that the embedding dimension ( $K$ ) and regularization weight ( $\omega$ ) for rows in matrix  $T$  are the two important hyperparameters. The embedding dimension was chosen from  $K = \{100, \dots, 800\}$ , and regularization weight from  $\omega = \{0.0001, \dots, 10.0\}$ .

### C. Proposed Topic ID Systems

Our Bayesian SMM is an unsupervised model trained iteratively by optimizing the ELBO; it does not necessarily correlate with the performance of topic ID. This is valid for SMM, NVDM or any other generative model trained without supervision. A typical way to overcome this problem is to have an early stopping mechanism (ESM), which requires to evaluate the topic ID accuracy on a held-out (or cross-validation) set at regular intervals during the training. It can then be used to stop the training earlier if needed.

Using the above-described scheme, we trained three different classifiers: (i) Gaussian linear classifier (GLC), (ii) multi-class logistic regression (LR), and, (iii) Gaussian linear classifier with uncertainty (GLCU). Note that GLC and LR cannot exploit the uncertainty in the document embeddings; and are trained using only the mean parameter  $\nu$  of the posterior distributions; whereas GLCU is trained using the full posterior distribution  $q(w)$ , i.e., along with the uncertainties of document embeddings as described in Section IV. GLC and GLCU do not have any hyperparameters to tune, while the  $\ell_2$  regularization weight for the parameters of LR was tuned using cross-validation experiments. Our code is available online.<sup>5</sup>

### D. Baseline Topic ID Systems

1) *NVDM*: Since NVDM and our proposed Bayesian SMM share similarities, we chose to extract the embeddings from

<sup>3</sup>[Online]. Available: <https://catalog.ldc.upenn.edu/LDC2004S13>

<sup>4</sup>[Online]. Available: <http://qwone.com/~jason/20Newsgroups/>

<sup>5</sup>[Online]. Available: <https://github.com/BUTSpeechFIT/BaySMM>

NVDM and use them for training linear classifiers. Given a trained NVDM model, embeddings for any test document can be extracted just by forward propagating through the encoder. Although this is computationally cheaper, one needs to decide when to stop the training, as a fully converged NVDM may not yield optimal embeddings for discriminative tasks such as topic ID. Hence, we used the same early stopping mechanism, as described in the earlier section. We used the same three classifier pipelines (LR, GLC, GLCU) as we used for Bayesian SMM. Our architecture and training scheme are similar to ones proposed in [15], i.e., two feed-forward layers with either 500 or 1000 hidden units and {sigmoid, ReLU, tanh} activation functions. The latent dimension was chosen from  $K = \{100, \dots, 800\}$ . The hyperparameters were tuned based on cross-validation experiments.

2) *SMM*: Our second baseline system is non-Bayesian SMM with  $\ell_1$  regularization over the rows in  $T$  matrix, i.e.,  $\ell_1$  SMM. It was trained with hyperparameters such as embedding dimension  $K = \{100, \dots, 800\}$ , and regularization weight  $\omega = \{0.0001, \dots, 10.0\}$ . The embeddings obtained from SMM were then used to train GLC and LR classifiers. Note that we cannot use GLCU here, because SMM yields only point-estimates of embeddings. We used the same early stopping mechanism to train the classifiers. The experimental analysis in Section VII-C shows that Bayesian SMM is more robust to over-fitting when compared to SMM and NVDM, and does not require an early stopping mechanism.

3) *ULMFiT*: The third baseline system is the universal language model fine-tuned for classification (ULMFiT) [9]. The pre-trained<sup>6</sup> model consists of 3 BiLSTM layers. Fine-tuning the model involves two steps: (a) fine-tuning LM on the target dataset and (b) training classifier (MLP layer) on the target dataset. We trained several models with various drop-out rates. More specifically, the LM was fine-tuned for 15 epochs,<sup>7</sup> with drop-out rates from  $\{0.2, \dots, 0.6\}$ . The classifier was fine-tuned for 50 epochs with drop-out rates from  $\{0.2, \dots, 0.6\}$ . A held-out development set was used to tune the hyperparameters (drop-out rates, and fine-tuning epochs).

4) *TF-IDF*: The fourth baseline system is a standard term frequency-inverse document frequency (TF-IDF) based document representation, followed by multi-class logistic regression (LR). Although TF-IDF is not a topic model, the classification performance of TF-IDF based systems is often close to state-of-the-art systems [19]. The hyperparameter ( $\ell_2$  regularization weight) of LR was selected based on 5-fold cross-validation experiments on training set.

## VII. RESULTS AND DISCUSSION

### A. Convergence Rate of Bayesian SMM

We observed that the posterior distributions extracted using Bayesian SMM are always much sharper than standard Normal distribution. Hence we initialized the variational distribution

<sup>6</sup>[Online]. Available: <https://github.com/fastai/fastai>

<sup>7</sup>Fine-tuning LM for higher number of epochs degraded the classification performance.

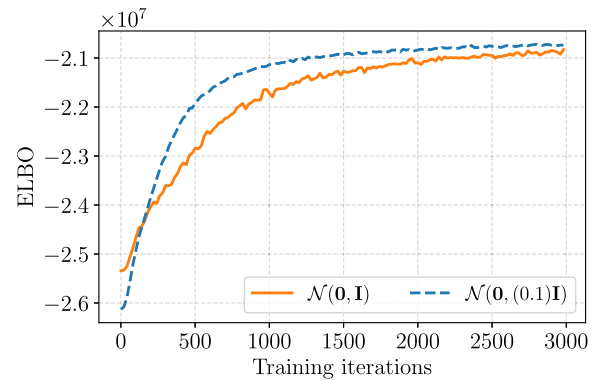


Fig. 4. Convergence of Bayesian SMM for various initializations of variational distribution. The model was trained on *20NewsGroups* corpus with  $K = 100$ , and  $\omega = 1$ .

to  $\mathcal{N}(\mathbf{0}, (0.1)\mathbf{I})$  to speed up the convergence. Fig. 4 shows objective (ELBO) plotted for two different initializations of variational distribution. Here, the model was trained on *20NewsGroups* corpus, with the embedding dimension  $K = 100$ , regularization weight  $\omega = 1.0$  and prior set to standard Normal. We can observe that the model initialized to  $\mathcal{N}(\mathbf{0}, (0.1)\mathbf{I})$  converges faster as compared to the one initialized to standard Normal. In all the further experiments, we initialized<sup>8</sup> both the prior and variational distributions to  $\mathcal{N}(\mathbf{0}, (0.1)\mathbf{I})$ .

### B. Perplexity

Perplexity is inversely proportional to the log-likelihood of the data. When computed on the test data, it gives a notion of how well the model explains (fits) the test (unseen) data. Perplexity computed on test data is a standard way of evaluating language models [47], [48]. Since topic models built on bag-of-words are equivalent to unigram language models, perplexity is seen as an intrinsic measure for topic models [15], [49]. It is computed as an average of every test document according to:

$$\text{PPL}_{\text{DOC}} = \exp \left\{ \frac{-1}{D} \sum_{d=1}^D \frac{\log p(\mathbf{x}_d)}{N_d} \right\}, \quad (46)$$

or for an entire test corpus according to:

$$\text{PPL}_{\text{CORPUS}} = \exp \left\{ - \frac{\sum_{d=1}^D \log p(\mathbf{x}_d)}{\sum_{d=1}^D N_d} \right\}, \quad (47)$$

where  $N_d$  is the number of word tokens in document  $d$ .

In our case,  $\log p(\mathbf{x})$  from (10) cannot be evaluated, because the KL divergence from variational distribution  $q$  to the true posterior  $p$  cannot be computed; as the true posterior is intractable (5). We can only compute  $\mathcal{L}(q)$ , which is a lower bound on  $\log p(\mathbf{x})$ . Thus, the resulting perplexity values act as upper bounds. This is true for NVDM [15] or any other model in the VB framework, where the true posterior is intractable [16]. We estimated  $\mathcal{L}(q)$  from (18) using 32 samples, i.e.,  $R = 32$ , in order to compute perplexity. We used the same number of samples for

<sup>8</sup>One can introduce hyper-priors and learn the parameters of prior distribution.



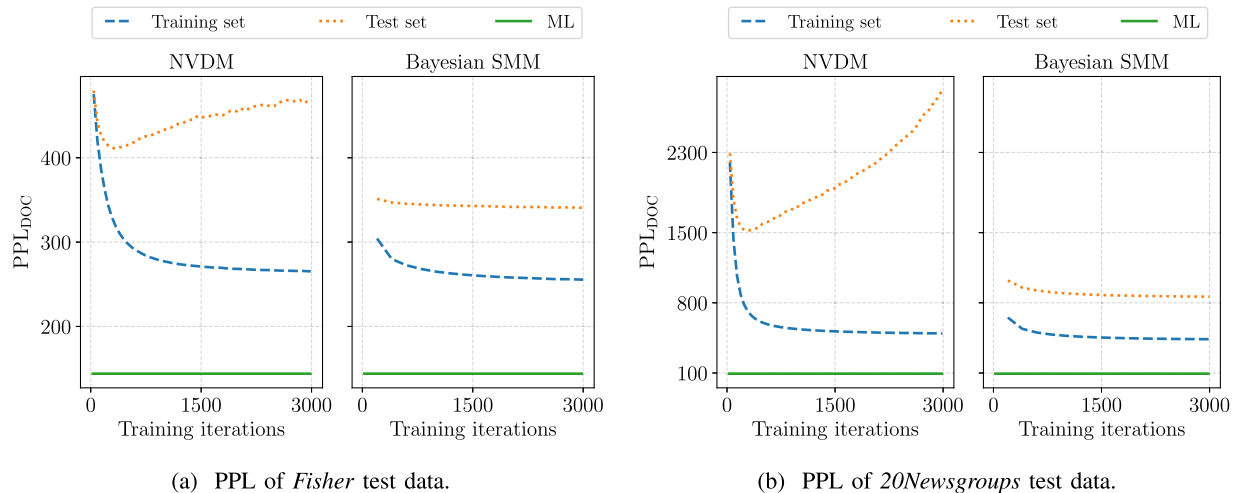


Fig. 5. Comparison of training and test data perplexities obtained using Bayesian SMM and NVDM for both *Fisher* and *20Newsgroups* datasets. The horizontal solid green line shows the test data perplexity computed using the maximum likelihood (ML) probabilities estimated on the test data. The latent (embedding) dimension was set to 200 for both the models.

TABLE II  
COMPARISON OF PERPLEXITY (PPL) RESULTS ON *20NEWSGROUPS*. THE VALUES IN THE BRACKETS INDICATE RESULTS WITH A LIMITED VOCABULARY OF 2000 WORDS

| Model        | $K$ | PPL <sub>CORPUS</sub> | PPL <sub>DOC</sub> |
|--------------|-----|-----------------------|--------------------|
| NVDM         | 50  | 1287 (769)            | 1421 (820)         |
| NVDM         | 200 | 1387 (852)            | 1519 (870)         |
| Bayesian SMM | 50  | <b>1043 (629)</b>     | <b>1064 (639)</b>  |
| Bayesian SMM | 200 | <b>882 (519)</b>      | <b>851 (515)</b>   |
| ML estimate  | -   | 153 (90)              | 93 (42)            |

the baseline NVDM. We observed that perplexity values are consistent when estimated with  $R \geq 16$ . In [15], the authors used 20 samples.

We present the comparison of *20Newsgroups* test data perplexities obtained using Bayesian SMM and NVDM in Table II. It shows the perplexities of *20Newsgroups* corpus under full and a limited vocabulary of 2000 words [15]. We also show the perplexity computed using the maximum likelihood probabilities estimated on the test data. It acts as the lower bound on the test perplexities. NVDM was shown [15] to achieve superior perplexity scores when compared to LDA, docNADE [50], Deep Auto Regressive Neural Network models [51]. To the best of our knowledge, our model achieves state-of-the-art perplexity scores on *20Newsgroups* corpus under limited and full vocabulary conditions.

In further investigation, we trained both Bayesian SMM and NVDM until convergence. At regular checkpoints during the training, we froze the model, extracted the embeddings for both training and test data, and computed the perplexities; shown in Figs. 5(a) and 5(b). We can observe that both the Bayesian SMM and NVDM fit the training data equally well (low perplexities). However, in the case of NVDM, the perplexity of test data increases after a certain number of iterations; suggesting that NVDM fails to generalize and over-fits on the training data.

In the case of Bayesian SMM, the perplexity of the test data decreases and remains stable, illustrating the robustness of our model.

### C. Early Stopping Mechanism for Topic ID Systems

The embeddings extracted from a model trained purely in an unsupervised fashion does not necessarily yield optimum results when used in a supervised scenario. As discussed earlier in Section VI-C and VI-D, an early stopping mechanism (ESM) during the training of an unsupervised model (e.g., NVDM, SMM, and Bayesian SMM) is required to get optimal performance from the subsequent topic ID system. The following experiment illustrates the idea of ESM:

We trained SMM, Bayesian SMM and NVDM on *Fisher* data until convergence. At regular checkpoints during the training, we froze the model, extracted the embeddings for both training and test data. We chose GLC for SMM, GLCU for NVDM, and Bayesian SMM as topic ID classifiers. We then evaluated the topic ID accuracy on the cross-validation<sup>9</sup> and test sets. Fig. 6 shows the topic ID accuracy on cross-validation and test sets obtained at regular checkpoints for all the three models. The circular dot (●) represents the best cross-validation score and the corresponding test score that is obtained by employing ESM. In the case of (non-Bayesian) SMM, the test accuracy drops significantly after a certain number of iterations; suggesting the strong need of ESM. The cross-validation accuracies of NVDM and Bayesian SMM are similar and remain consistent over the iterations. However, the test accuracy of NVDM is relatively lower and decreases over the iterations. On the other hand, the test accuracy of Bayesian SMM increases and stays consistent. It shows the robustness of our proposed model, which besides does not require any ESM. In all the further topic ID experiments, we

<sup>9</sup>5-fold cross-validation on training set.

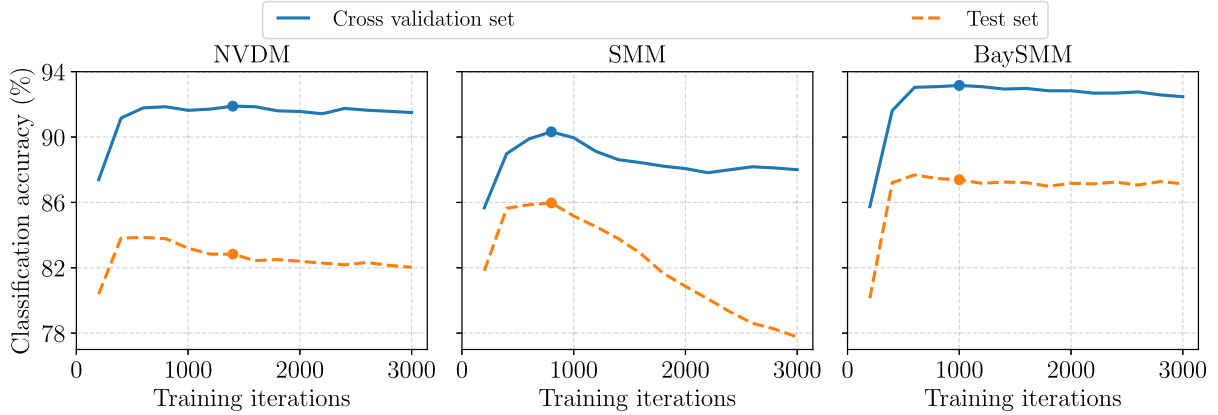


Fig. 6. Performance of topic ID systems on *Fisher* data at various checkpoints during model training. The circular dot (●) represents the best cross-validation score and the corresponding test score obtained using the early stopping mechanism (ESM). The embedding dimension was set to 100 for all the models.

TABLE III  
COMPARISON OF RESULTS ON *FISHER* TEST SETS, FROM EARLIER PUBLISHED WORKS, OUR BASELINES AND PROPOSED SYSTEMS. ★ INDICATES A PURE DISCRIMINATIVE MODEL

| Systems      | Model        | Classifier | Manual transcriptions |             | Automatic transcriptions |             |
|--------------|--------------|------------|-----------------------|-------------|--------------------------|-------------|
|              |              |            | Accuracy (%)          | CE          | Accuracy (%)             | CE          |
| Prior works  | BoW [40]     | NB         | 87.61                 | -           | -                        | -           |
|              | TF-IDF [19]  | LR         | 86.41                 | -           | -                        | -           |
| Our Baseline | TF-IDF       | LR         | 86.59                 | 0.93        | 86.77                    | <b>0.94</b> |
|              | ULMFiT ★     | MLP        | 86.41                 | <b>0.50</b> | 86.08                    | <b>0.50</b> |
|              | $\ell_1$ SMM | LR         | 86.81                 | 0.91        | 87.02                    | 1.09        |
|              | $\ell_1$ SMM | GLC        | 85.17                 | 1.64        | 85.53                    | 1.54        |
|              | NVDM         | LR         | 81.16                 | 0.94        | 83.67                    | 1.15        |
|              | NVDM         | GLC        | 84.47                 | 1.25        | 84.15                    | 1.22        |
| Proposed     | NVDM         | GLCU       | 83.96                 | 0.93        | 83.01                    | 0.97        |
|              | Bayesian SMM | LR         | <b>89.91</b>          | <b>0.89</b> | <b>88.23</b>             | 0.95        |
|              | Bayesian SMM | GLC        | <b>89.47</b>          | 1.05        | <b>87.23</b>             | 1.46        |
|              | Bayesian SMM | GLCU       | <b>89.54</b>          | <b>0.68</b> | <b>87.54</b>             | <b>0.77</b> |

report classification results for Bayesian SMM without ESM; while the results for SMM and NVDM are with ESM.

#### D. Topic ID Results

This section presents the topic ID results in terms of classification accuracy (in %) and cross-entropy (CE) on the test sets. Cross-entropy gives a notion of how confident the classifier is about its prediction. A well-calibrated classifier tends to have lower a cross-entropy.

Table III presents the classification results on *Fisher* speech corpora with manual and automatic transcriptions, where the first two rows are the results from earlier published works. Hazen [40], used discriminative vocabulary selection followed by a naive Bayes (NB) classifier. Having a limited (small) vocabulary is the major drawback of this approach. Although we have used the same training and test splits, May *et al.* [19] had a slightly larger vocabulary than ours, and their best system is similar to our baseline TF-IDF based system. The remaining rows in Table III show our baselines and proposed systems. We can see that our proposed systems achieve consistently better accuracies; notably, GLCU, which exploits the uncertainty in

document embeddings has much lower cross-entropy than its counterpart GLC. To the best of our knowledge, the proposed systems achieve the best classification results on *Fisher* corpora with the current set-up, i.e., treating each side of the conversation as an independent document. It can be observed that ULMFiT has the lowest cross-entropy among all the systems.

Table IV presents classification results on *20Newsgroups* dataset. The first three rows give the results as reported in earlier works. Pappagari *et al.* [39], proposed a CNN-based discriminative model trained to jointly optimize categorical cross-entropy loss for classification task along with binary cross-entropy for verification task. Sparse composite document vector (SCDV) [37] exploits pre-trained word embeddings to obtain sparse document embeddings, whereas neural tensor skip-gram model (NTSG) [52] extends the idea of a skip-gram model for obtaining document embeddings. The authors in (SCDV) [37] have shown superior classification results as compared to paragraph vector (PV-DM, PV-DBOW), LDA, NTSG, and other systems. The next rows in Table IV present our baselines and proposed systems. We see that the topic ID systems based on Bayesian SMM and logistic regression is better than all the other models, except for the discriminative CNN model. We can

TABLE IV  
COMPARISON OF RESULTS ON *20Newsgroups* FROM EARLIER PUBLISHED WORKS, OUR BASELINES AND PROPOSED SYSTEMS. \* INDICATES A PURE DISCRIMINATIVE MODEL

| Systems       | Model        | Classifier | Accuracy (%) | CE          |
|---------------|--------------|------------|--------------|-------------|
| Prior works   | CNN [39] *   | -          | <b>86.12</b> | -           |
|               | SCDV [37]    | SVM        | <b>84.60</b> | -           |
|               | NTSG-1 [52]  | SVM        | 82.60        | -           |
| Our Baselines | TF-IDF       | LR         | 84.47        | <b>0.73</b> |
|               | ULMFiT *     | MLP        | 83.06        | 0.89        |
|               | $\ell_1$ SMM | LR         | 82.01        | <b>0.75</b> |
|               | $\ell_1$ SMM | GLC        | 82.02        | 1.33        |
|               | NVDM         | LR         | 79.57        | 0.86        |
|               | NVDM         | GLC        | 77.60        | 1.65        |
|               | NVDM         | GLCU       | 76.86        | 0.88        |
| Proposed      | Bayesian SMM | LR         | <b>84.65</b> | <b>0.53</b> |
|               | Bayesian SMM | GLC        | 83.22        | 1.28        |
|               | Bayesian SMM | GLCU       | 82.81        | 0.79        |

also see that all the topic ID systems based on Bayesian SMM are consistently better than variational autoencoder inspired NVDM, and (non-Bayesian) SMM.

In general, discriminative classifiers tend to perform better than generative classifiers, since discriminative classifiers model  $p(y|x)$  directly, whereas generative classifiers model  $p(x, y)$  to determine  $p(y|x)$ . However, in the presence of lower training data, generative classifiers tend to be better [53]. For instance, the cross-entropy results on *Fisher* test set (Table III) show that GLCU is better than LR. The same is not observed on *20Newsgroups* test set, where the training data is  $5\times$  more when compared to *Fisher*. Moreover, Gaussian based generative classifiers are much faster to train and can be easily adapted to newer classes when compared to discriminative classifiers such as logistic regression.

The advantages of the proposed Bayesian SMM are summarized as follows: (a) the document embeddings are Gaussian distributed which enables to train simple generative classifiers like GLC, or GLCU; that can be extended to newer classes easily, (b) although the Bayesian SMM is trained in an unsupervised fashion, it does not require any early stopping mechanism to yield optimal topic ID results; document embeddings extracted from a fully converged model can be directly used for classification tasks without any fine-tuning.

### E. Uncertainty in Document Embeddings

The uncertainty captured in the posterior distribution of document embeddings correlates strongly with the size of the document. The trace of the covariance matrix of the inferred posterior distributions gives us the notion of uncertainty. Fig. 7 shows an example of the uncertainty captured in the embeddings. Here, the Bayesian SMM was trained on *20Newsgroups* with an embedding dimension of 100.

## VIII. CONCLUSION

We have presented a generative model for learning document representations (embeddings) and their uncertainties. Our

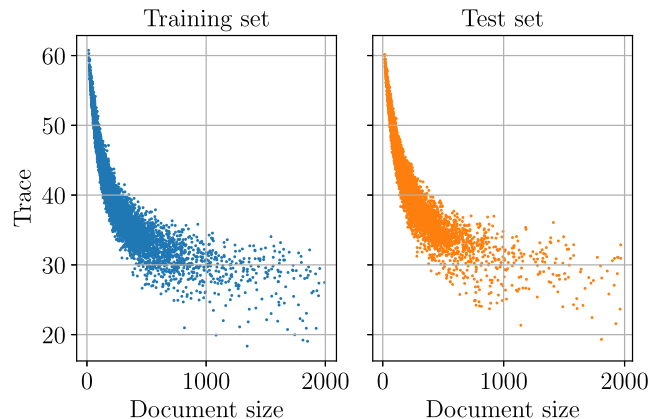


Fig. 7. Uncertainty (trace of covariance of posterior distribution) captured in the document embeddings of *20Newsgroups* dataset.

proposed model achieved state-of-the-art perplexity results on the standard *20Newsgroups* and *Fisher* datasets. Next, we have shown that the proposed model is robust to over-fitting and unlike in SMM and NVDM, it does not require any early stopping mechanism for topic ID. We proposed an extension to simple Gaussian linear classifier that exploits the uncertainty in document embeddings and achieves better cross-entropy scores on the test data as compared to the simple GLC. Using simple linear classifiers on the obtained document embeddings, we achieved superior classification results on Fisher speech *20Newsgroups* text corpora. We also addressed a commonly encountered problem of intractability while performing variational inference in mixed-logit models by using the re-parametrization trick. This idea can be translated straightforwardly to the subspace  $n$ -gram model for learning sentence embeddings, and also for learning word embeddings along with their uncertainties. The proposed Bayesian SMM can be extended to have topic-specific priors for document embeddings, which enables to encode topic label uncertainty explicitly in the document embeddings. There exist other scoring mechanisms that exploit the uncertainty in embeddings [54], which we plan to explore in our future works.

## APPENDIX A GRADIENTS OF LOWER BOUND

The variational distribution is Gaussian with the following parameterization:

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \boldsymbol{\nu}, \text{diag}(\exp\{2\boldsymbol{\zeta}\})). \quad (45)$$

The lower bound for a single document is:

$$\begin{aligned} \mathcal{L}_d \approx & -\frac{1}{2} \left[ \lambda \text{tr}(\text{diag}(\exp\{2\boldsymbol{\zeta}\})) - \log |\text{diag}(\exp\{2\boldsymbol{\zeta}\})| \right. \\ & \left. - K \log \lambda + \lambda \boldsymbol{\nu}^\top \boldsymbol{\nu} - K \right] \\ & + \sum_{i=1}^V x_i \left[ (m_i + \mathbf{t}_i \boldsymbol{\nu}) - \frac{1}{R} \sum_{r=1}^R \log \left( \sum_{j=1}^V \exp\{m_j + \mathbf{t}_j g(\boldsymbol{\epsilon}_r)\} \right) \right], \end{aligned} \quad (46)$$

where

$$g(\boldsymbol{\epsilon}) = \boldsymbol{\nu} + \text{diag}(\exp\{\boldsymbol{\varsigma}\})\tilde{\boldsymbol{\epsilon}}. \quad (47)$$

It is convenient to have the following derivatives:

$$\frac{\partial g(\boldsymbol{\epsilon})}{\partial \boldsymbol{\nu}} = \mathbf{I}. \quad (48)$$

$$\begin{aligned} \frac{\partial(\mathbf{t}_i g(\boldsymbol{\epsilon}))}{\partial \boldsymbol{\varsigma}} &= \text{diag}(\mathbf{t}_i^\top) \text{diag}(\exp\{\boldsymbol{\varsigma}\}) \text{diag}(\tilde{\boldsymbol{\epsilon}}) \\ &= \mathbf{t}_i^\top \odot \exp\{\boldsymbol{\varsigma}\} \odot \tilde{\boldsymbol{\epsilon}}. \end{aligned} \quad (49)$$

### A. Derivatives of the Parameters of Variational Distribution

Taking derivative of the objective function (46) with respect to mean parameter  $\boldsymbol{\nu}$  and using (48):

$$\begin{aligned} \frac{\partial \mathcal{L}_d}{\partial \boldsymbol{\nu}} &= -\lambda \boldsymbol{\nu} + \sum_{i=1}^V x_i \left[ \mathbf{t}_i^\top - \frac{1}{R} \sum_{r=1}^R \sum_{k=1}^V \mathbf{t}_k^\top \mathbf{I} \right. \\ &\quad \left. \times \underbrace{\frac{\exp\{m_k + \mathbf{t}_k g(\boldsymbol{\epsilon}_r)\}}{\sum_j \exp\{m_j + \mathbf{t}_j g(\boldsymbol{\epsilon}_r)\}}}_{\theta_{kr}} \right] \end{aligned} \quad (50)$$

$$= \left[ \sum_{i=1}^V x_i \mathbf{t}_i^\top - \sum_{i=1}^V \mathbf{t}_i^\top \frac{1}{R} \sum_{r=1}^R \theta_{ir} \sum_{k=1}^V x_k \right] - \lambda \boldsymbol{\nu} \quad (51)$$

$$\boxed{\nabla_{\boldsymbol{\nu}} \mathcal{L} = \left[ \sum_{i=1}^V \mathbf{t}_i^\top \left( x_i - \frac{1}{R} \sum_{r=1}^R \theta_{ir} \sum_{k=1}^V x_k \right) \right] - \lambda \boldsymbol{\nu}.} \quad (52)$$

Taking the derivative of the objective function (46) with respect to  $\boldsymbol{\varsigma}$  and using (49):

$$\begin{aligned} \frac{\partial \mathcal{L}_d}{\partial \boldsymbol{\varsigma}} &= -\frac{1}{2} [2\lambda \exp\{2\boldsymbol{\varsigma}\} - 2\mathbf{I}] \\ &\quad + \sum_{i=1}^V x_i \left[ -\frac{1}{R} \sum_{r=1}^R \sum_{k=1}^V \mathbf{t}_k^\top \tilde{\boldsymbol{\epsilon}}_r^\top \underbrace{\frac{\exp\{m_k + \mathbf{t}_k g(\boldsymbol{\epsilon}_r)\}}{\sum_j \exp\{m_j + \mathbf{t}_j g(\boldsymbol{\epsilon}_r)\}}}_{\theta_{kr}} \right] \\ &= \mathbf{1} - \lambda \exp\{2\boldsymbol{\varsigma}\} \\ &\quad - \left[ \sum_{i=1}^V x_i \frac{1}{R} \sum_{r=1}^R \sum_{k=1}^V \mathbf{t}_k^\top \odot \exp\{\boldsymbol{\varsigma}\} \odot \tilde{\boldsymbol{\epsilon}}_r^\top \theta_{kr} \right] \end{aligned} \quad (53)$$

$$\boxed{\nabla_{\boldsymbol{\varsigma}} \mathcal{L} = \mathbf{1} - \lambda \exp\{2\boldsymbol{\varsigma}\} - \left[ \left( \sum_{i=1}^V x_i \right) \frac{1}{R} \sum_{r=1}^R \sum_{k=1}^V \theta_{kr} \mathbf{t}_k^\top \odot \exp\{\boldsymbol{\varsigma}\} \odot \tilde{\boldsymbol{\epsilon}}_r \right].} \quad (54)$$

### B. Derivatives of the Model Parameters

Taking the derivative of complete objective (19) with respect to a row  $\mathbf{t}_k$  from matrix  $\mathbf{T}$ :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{t}_k} &= \frac{\partial}{\partial \mathbf{t}_k} \sum_{d=1}^D \sum_{i=1}^V x_{di} \left[ (m_i + \mathbf{t}_i \boldsymbol{\nu}_d) \right. \\ &\quad \left. - \frac{1}{R} \sum_{r=1}^R \log \left( \sum_{j=1}^V \exp\{m_j + \mathbf{t}_j g(\boldsymbol{\epsilon}_r)\} \right) \right] \\ &\quad - \omega \sum_{i=1}^V \|\mathbf{t}_i\|_1 \end{aligned} \quad (55)$$

$$\begin{aligned} &= \sum_{d=1}^D \left[ x_{dk} \boldsymbol{\nu}_d^\top - \sum_{i=1}^V x_{di} \frac{1}{R} \sum_{r=1}^R g(\boldsymbol{\epsilon}_{dr})^\top \right. \\ &\quad \left. \times \underbrace{\frac{\exp\{m_i + \mathbf{t}_k g(\boldsymbol{\epsilon}_{dr})\}}{\sum_j \exp\{m_j + \mathbf{t}_j g(\boldsymbol{\epsilon}_{dr})\}}}_{\theta_{dkr}} \right] \\ &\quad - \omega \text{sign}(\mathbf{t}_k) \end{aligned} \quad (56)$$

$$\begin{aligned} &= \sum_{d=1}^D \left[ x_{dk} \boldsymbol{\nu}_d^\top - \sum_{i=1}^V x_{di} \frac{1}{R} \sum_{r=1}^R g(\boldsymbol{\epsilon}_{dr})^\top \theta_{dkr} \right] \\ &\quad - \omega \text{sign}(\mathbf{t}_k) \end{aligned} \quad (57)$$

$$\begin{aligned} \nabla_{\mathbf{t}_k} \mathcal{L} &= \sum_{d=1}^D \left[ x_{dk} \boldsymbol{\nu}_d^\top - \left[ \left( \sum_{i=1}^V x_{di} \right) \frac{1}{R} \sum_{r=1}^R \theta_{dkr} g(\boldsymbol{\epsilon}_{dr})^\top \right] \right] \\ &\quad - \omega \text{sign}(\mathbf{t}_k). \end{aligned} \quad (58)$$

## APPENDIX B

### EM ALGORITHM FOR GLCU

#### E-STEP

Obtaining the posterior distribution of latent variable  $p(\mathbf{y}_d | \boldsymbol{\nu}_d, \Theta)$ . Using the results from [28] (p. 41, (358)):

$$\begin{aligned} \log p(\mathbf{y}_d | \boldsymbol{\nu}_d, \mathbf{h}_d, \Theta) &= \log p(\boldsymbol{\nu}_d | \mathbf{y}_d, \mathbf{h}_d) + \log p(\mathbf{y}_d) - \log p(\boldsymbol{\nu}_d) \\ &= \log \mathcal{N}(\boldsymbol{\nu}_d | \boldsymbol{\mu}_d + \mathbf{y}_d, \mathbf{D}^{-1}) \\ &\quad + \log \mathcal{N}(\mathbf{y}_d | \mathbf{0}, \boldsymbol{\Gamma}_d^{-1}) + \text{const} \\ &= -\frac{1}{2} (\boldsymbol{\nu}_d - (\boldsymbol{\mu}_d + \mathbf{y}_d))^\top \mathbf{D} (\boldsymbol{\nu}_d - (\boldsymbol{\mu}_d + \mathbf{y}_d)) \\ &\quad - \frac{1}{2} \mathbf{y}_d^\top \boldsymbol{\Gamma}_d \mathbf{y}_d + \text{const} \\ &= -\frac{1}{2} (\mathbf{y}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d))^\top \mathbf{D} (\mathbf{y}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d)) \\ &\quad - \frac{1}{2} \mathbf{y}_d^\top \boldsymbol{\Gamma}_d \mathbf{y}_d + \text{const} \\ &= \mathcal{N}(\mathbf{y}_d | \mathbf{u}_d, \mathbf{V}_d^{-1}) \end{aligned}$$

where  $\mathbf{u}_d$  is simplified as:

$$\begin{aligned}\mathbf{u}_d &= (\mathbf{D} + \mathbf{\Gamma}_d)^{-1}(\mathbf{D}(\boldsymbol{\nu}_d - \boldsymbol{\mu}_d) + \mathbf{\Gamma}_d \mathbf{0}) \\ &= [\mathbf{D}^{-1}(\mathbf{D} + \mathbf{\Gamma}_d)]^{-1}(\boldsymbol{\nu}_d - \boldsymbol{\mu}_d)\end{aligned}$$

resulting in:

$$\mathbf{u}_d = [\mathbf{I} + \mathbf{D}^{-1}\mathbf{\Gamma}_d]^{-1}(\boldsymbol{\nu}_d - \boldsymbol{\mu}_d) \quad (59)$$

$$\mathbf{V}_d = \mathbf{D} + \mathbf{\Gamma}_d \quad (60)$$

*M-STEP*

Maximizing the auxiliary function:

$$\Theta^{\text{new}} = \underset{\Theta}{\operatorname{argmax}} \mathcal{Q}(\Theta, \Theta^{\text{old}}) \quad (61)$$

$$q(\mathbf{y}) = p(\mathbf{y} | \mathbf{w}, \Theta^{\text{old}}). \quad (62)$$

Using the results from [28] [p. 43, (378)], the auxiliary function  $\mathcal{Q}(\Theta, \Theta^{\text{old}})$  is computed as:

$$\begin{aligned}\mathcal{Q}(\Theta, \Theta^{\text{old}}) &= \mathbb{E}_q \left[ \sum_{d=1}^D \log p(\boldsymbol{\nu}_d, \mathbf{y}_d) \right] \\ &= \sum_{d=1}^D \mathbb{E}_q [\log p(\boldsymbol{\nu}_d | \mathbf{y}_d)] + \mathbb{E}_q [\log p(\mathbf{y}_d)] \\ &= \sum_{d=1}^D \mathbb{E}_q [\log \mathcal{N}(\boldsymbol{\nu}_d | \boldsymbol{\mu}_d + \mathbf{y}_d, \mathbf{D}^{-1})] + \text{const} \\ &= \frac{D}{2} \log |\mathbf{D}| - \frac{1}{2} \sum_{d=1}^D \left[ \mathbb{E}_q [(\boldsymbol{\nu}_d - (\boldsymbol{\mu}_d + \mathbf{y}_d))^{\top} \mathbf{D} (\boldsymbol{\nu}_d - (\boldsymbol{\mu}_d + \mathbf{y}_d))] \right] + \text{const} \\ &= \frac{D}{2} \log |\mathbf{D}| - \frac{1}{2} \sum_{d=1}^D [\operatorname{tr}(\mathbf{D}\mathbf{V}_d^{-1}) \\ &\quad + (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d))^{\top} \mathbf{D} (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d))] \end{aligned}$$

Maximizing the auxiliary function  $\mathcal{Q}$  with respect to model parameters  $\Theta = \{\mathbf{M}, \mathbf{D}\}$

Taking the derivative with respect to each column  $\boldsymbol{\mu}_\ell$  in  $\mathbf{M}$  and equating it to zero:

$$\begin{aligned}\frac{\partial \mathcal{Q}}{\partial \boldsymbol{\mu}_\ell} &= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_\ell} \sum_{d \in \mathcal{I}_\ell} [(\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_\ell))^{\top} \mathbf{D} (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_\ell))] \\ &= -\frac{1}{2} \sum_{d \in \mathcal{I}_\ell} 2\mathbf{D} (\boldsymbol{\mu}_\ell - (\boldsymbol{\nu}_d - \mathbf{u}_d)) \\ &= -\mathbf{D} \left( \sum_{d \in \mathcal{I}_\ell} \boldsymbol{\mu}_\ell - \sum_{d \in \mathcal{I}_\ell} (\boldsymbol{\nu}_d - \mathbf{u}_d) \right) \quad (63)\end{aligned}$$

$$\boldsymbol{\mu}_\ell = \frac{1}{|\mathcal{I}_\ell|} \sum_{d \in \mathcal{I}_\ell} (\boldsymbol{\nu}_d - \mathbf{u}_d) \quad (64)$$

Taking the derivative with respect to shared precision matrix  $\mathbf{D}$  and equating it to zero:

$$\begin{aligned}\frac{\partial \mathcal{Q}}{\partial \mathbf{D}} &= \frac{D}{2} \mathbf{D}^{-1} - \frac{1}{2} \left( \sum_{d=1}^D \mathbf{V}_d^{-1} \right)^{\top} \\ &\quad - \frac{1}{2} \left( \sum_{d=1}^D (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d)) (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d))^{\top} \right)^{\top} \quad (65)\end{aligned}$$

$$\begin{aligned}\mathbf{D}^{-1} &= \frac{1}{D} \left[ \sum_{d=1}^D \mathbf{V}_d^{-1} \right. \\ &\quad \left. + \sum_{d=1}^D (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d)) (\mathbf{u}_d - (\boldsymbol{\nu}_d - \boldsymbol{\mu}_d))^{\top} \right] \quad (66)\end{aligned}$$

## REFERENCES

- [1] X. Wei and W. B. Croft, "LDA-based document models for ad-hoc retrieval," in *Proc. 29th Annu. Int. ACM SIGIR*, Aug. 2006, pp. 178–185.
- [2] T. Mikolov and G. Zweig, "Context dependent recurrent neural network language model," in *Proc. IEEE SLT Workshop*, Dec. 2012, pp. 234–239.
- [3] J. Wintrop and S. Khudanpur, "Limited resource term detection for effective topic identification of speech," in *Proc. IEEE ICASSP*, May 2014, pp. 7118–7122.
- [4] X. Chen *et al.*, "Recurrent neural network language model adaptation for multi-genre broadcast speech recognition," in *Proc. Interspeech. ISCA*, Sep. 2015, pp. 3511–3515.
- [5] K. Beneš, S. Kesiraju, and L. Burget, "i-Vectors in language modeling: An efficient way of domain adaptation for feed-forward models," in *Proc. Interspeech. ISCA*, 2018, pp. 3383–3387.
- [6] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. Adv. NIPS*, Dec. 2013, pp. 3111–3119.
- [7] J. Pennington, R. Socher, and C. D. Manning, "GloVe: Global vectors for word representation," in *Proc. Conf. EMNLP, ACL*, Oct. 2014, pp. 1532–1543.
- [8] Q. V. Le and T. Mikolov, "Distributed representations of sentences and documents," in *Proc. ICML*, Jun. 2014, pp. 1188–1196.
- [9] J. Howard and S. Ruder, "Universal language model fine-tuning for text classification," in *Proc. 56th Annu. Meet. ACL*, Melbourne, Australia: ACL, Jul. 2018, pp. 328–339.
- [10] M. Peters *et al.*, "Deep contextualized word representations," in *Proc. NAACL: HLT. ACL*, Jun. 2018, pp. 2227–2237.
- [11] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. Conf. NAACL: HLT, Volume 1 (Long and Short Papers)*, Minneapolis, Minnesota: Assoc. for Comput. Linguist., Jun. 2019, pp. 4171–4186.
- [12] C. Bishop, "Latent variable models," in *Proc. Learn. Graph. Models*. MIT Press, Jan. 1999, pp. 371–403.
- [13] D. M. Blei, "Probabilistic topic models," *Commun. ACM*, vol. 55, no. 4, pp. 77–84, Apr. 2012.
- [14] D. P. Kingma and M. Welling, "Auto-Encoding variational Bayes," in *Proc. 2nd ICLR*, 2014.
- [15] Y. Miao, L. Yu, and P. Blunsom, "Neural variational inference for text processing," in *Proc. 33rd ICML*, ser. ICML'16. JMLR.org, 2016, pp. 1727–1736.
- [16] C. M. Bishop, *Pattern Recognit. Mach. Learn. (Inform. Sci. Statist.)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [17] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic back propagation and approximate inference in deep generative models," in *Proc. 31st ICML*, ser. Proc. of Mach. Learning Res., vol. 32. Beijing, China: PMLR, 22–24 Jun. 2014, pp. 1278–1286.
- [18] M. Soufifar, L. Burget, O. Pichot, S. Cumani, and J. Cernocký, "Regularized subspace n-gram model for phonotactic vector extraction," in *Proc. Interspeech. ISCA*, Aug. 2013, pp. 74–78.

- [19] C. May, F. Ferraro, A. McCree, J. Wintrobe, D. Garcia-Romero, and B. V. Durme, "Topic identification and discovery on text and speech," in *Proc. Conf. EMNLP*, Sep. 2015, pp. 2377–2387.
- [20] S. Kesiraju, L. Burget, I. Szöke, and J. Černocký, "Learning document representations using subspace multinomial model," in *Proc. Interspeech*. ISCA, Sep. 2016, pp. 700–704.
- [21] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *JMLR*, vol. 3, pp. 993–1022, 2003.
- [22] D. M. Blei and J. D. Lafferty, "Correlated topic models," in *Adv. Neural Inform. Process. Syst. NIPS*, Dec. 2005, pp. 147–154.
- [23] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. 3rd Int. Conf. Learn. Representations*, May 2015.
- [24] G. Andrew and J. Gao, "Scalable Training of L1-Regularized log-linear models," in *Proc. 24th ICML*. New York, USA: ACM, 2007, pp. 33–40.
- [25] M. Kockmann, L. Burget, O. Glembek, L. Ferrer, and J. Černocký, "Prosodic speaker verification using subspace multinomial models with intonation compensation," in *Proc. Interspeech*. ISCA, Sep. 2010, pp. 1061–1064.
- [26] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Trans. Audio, Speech & Lang. Process.*, vol. 19, no. 4, pp. 788–798, 2011.
- [27] N. Depraetere and M. Vandebroek, "A comparison of variational approximations for fast inference in mixed logit models," *Comput. Statist.*, vol. 32, no. 1, pp. 93–125, 2017.
- [28] K. B. Petersen and M. S. Pedersen, "The Matrix Cookbook," Nov. 2012.
- [29] J. Paisley, D. M. Blei, and M. I. Jordan, "Variational Bayesian inference with stochastic search," in *Proc. 29th Int. Conf. Mach. Learn.*, ser. ICML12. Madison, WI, USA: Omnipress, 2012, pp. 1363–1370.
- [30] P. Kenny, T. Stafylakis, P. Ouellet, M. J. Alam, and P. Dumouchel, "PLDA for speaker verification with utterances of arbitrary duration," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, May 2013, pp. 7649–7653.
- [31] S. Cumani, O. Plchot, and R. Fér, "Exploiting i-vector posterior covariances for short-duration language recognition," in *Proc. Interspeech*, no. 9. ISCA, 2015, pp. 1002–1006.
- [32] M. Soufifar *et al.*, "iVector approach to phonotactic language recognition," in *Proc. Interspeech*. ISCA, Aug. 2011, pp. 2913–2916.
- [33] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *Proc. 49th Annu. Meet. ACL: Human Lang. Technol.*, Jun. 2011, pp. 142–150.
- [34] J. Eisenstein, A. Ahmed, and E. P. Xing, "Sparse additive generative models of text," in *Proc. 28th ICML*. USA: Omnipress, 2011, pp. 1041–1048.
- [35] M. V. S. Shashanka, B. Raj, and P. Smaragdis, "Sparse overcomplete latent variable decomposition of counts data," in *Proc. Neural Inf. Process. Syst.*, Dec. 2007, pp. 1313–1320.
- [36] J. Zhu and E. P. Xing, "Sparse topical coding," in *Proc. 27th Conf. UAI*, Jul. 2011, pp. 831–838.
- [37] D. Mekala, V. Gupta, B. Paranjape, and H. Karnick, "Scdv: Sparse composite document vectors using soft clustering over distributional representations," in *Proc. Conf. Empir. Methods Natural Language Process*. Copenhagen, Denmark: ACL, Sep. 2017, pp. 659–669.
- [38] S. Kesiraju *et al.*, "Topic identification of spoken documents using unsupervised acoustic unit discovery," in *IEEE Int. Conf. Acoust., Speech Signal Process.*, New Orleans, LA, USA, 2017, pp. 5745–5749.
- [39] R. Pappagari, J. Villalba, and N. Dehak, "Joint verification-identification in end-to-end multi-scale CNN framework for topic identification," in *Proc. IEEE ICASSP*, Apr. 2018, pp. 6199–6203.
- [40] T. J. Hazen, F. Richardson, and A. Margolis, "Topic identification from audio recordings using word and phone recognition lattices," in *Proc. IEEE Workshop ASRU*, Dec. 2007, pp. 659–664.
- [41] T. J. Hazen, "MCE training techniques for topic identification of spoken audio documents," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 19, no. 8, pp. 2451–2460, Nov. 2011.
- [42] D. Povey *et al.*, "The Kaldi speech recognition toolkit," in *Proc. IEEE Workshop ASRU*. IEEE Signal Processing Society, Dec. 2011.
- [43] K. Veselý, A. Ghoshal, L. Burget, and D. Povey, "Sequence-discriminative training of deep neural networks," in *Proc. Interspeech*. ISCA, Aug. 2013, pp. 2345–2349.
- [44] Y. Miao, E. Grefenstette, and P. Blunsom, "Discovering discrete latent topics with neural variational inference," in *Proc. 34th Int. Conf. Mach. Learn. - Volume 70*, ser. ICML17. JMLR.org, 2017, pp. 2410–2419.
- [45] N. Srivastava, R. Salakhutdinov, and G. E. Hinton, "Modeling documents with deep Boltzmann machines," in *UAI*, Aug. 2013.
- [46] A. Srivastava and C. A. Sutton, "Autoencoding variational inference for topic models," in *Proc. 5th Int. Conf. Learn. Representations, Toulon, France*, Apr. 24–26, 2017, *Conf. Track Proc. OpenReview.net*, 2017.
- [47] Y. Bengio, R. Ducharme, P. Vincent, and C. Janvin, "A neural probabilistic language model," *J. Mach. Learn. Res.*, vol. 3, no. null, pp. 1137–1155, Mar. 2003.
- [48] D. Jurafsky and J. H. Martin, *Speech and Language Process. (2nd Edition)*. USA: Prentice-Hall, Inc., 2009.
- [49] N. Srivastava, R. Salakhutdinov, and G. Hinton, "Modeling documents with a deep Boltzmann machine," in *Proc. Twenty-Ninth Conf. UAI*, ser. UAI'13. Arlington, Virginia, United States: AUAI Press, 2013, pp. 616–624.
- [50] H. Larochelle and S. Lauly, "A neural autoregressive topic model," in *Proc. Adv. NIPS*, Dec. 2012, pp. 2717–2725.
- [51] A. Mnih and K. Gregor, "Neural variational inference and learning in belief networks," in *Proc. 31th ICML*, Jun. 2014, pp. 1791–1799.
- [52] P. Liu, X. Qiu, and X. Huang, "Learning context-sensitive word embeddings with neural tensor skip-gram model," in *Proc. 24th Int. Conf. Artif. Intell.*, ser. IJCAI'15. AAAI Press, 2015, pp. 1284–1290.
- [53] A. Y. Ng and M. I. Jordan, "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes," in *Adv. Neural Inform. Process. Syst. 14*. MIT Press, 2002, pp. 841–848.
- [54] N. Brümmer, A. Silnova, L. Burget, and T. Stafylakis, "Gaussian meta-embeddings for efficient scoring of a heavy-tailed PLDA model," in *Proc. Odyssey the Speaker Lang. Recognit. Workshop*, 2018, pp. 349–356.