
Rethinking Collapsed Variational Bayes Inference for LDA

Issei Sato

The University of Tokyo

Hiroshi Nakagawa

The University of Tokyo

SATO@R.DL.ITC.U-TOKYO.AC.JP

N3@DL.ITC.U-TOKYO.AC.JP

Abstract

We propose a novel interpretation of the collapsed variational Bayes inference with a zero-order Taylor expansion approximation, called CVB0 inference, for latent Dirichlet allocation (LDA). We clarify the properties of the CVB0 inference by using the α -divergence. We show that the CVB0 inference is composed of two different divergence projections: $\alpha = 1$ and -1 . This interpretation will help shed light on CVB0 works.

1. Introduction

Latent Dirichlet allocation (LDA) (Blei et al., 2003) is a well-known probabilistic latent variable model. It is used to model the co-occurrence of words by using latent variables called topics where a document is represented as a “bag of words”. It has a wide variety of applications in many fields. Originally, the variational Bayes (VB) inference was used for learning LDA. The collapsed variational Bayes (CVB) inference was developed as an alternative deterministic inference for LDA (Teh et al., 2007). The CVB inference is a variational inference improved by marginalizing out parameters as in a collapsed Gibbs sampler (Griffiths & Steyvers, 2004). (Sung et al., 2008) generalized the CVB inference for conjugate-exponential family models, called latent-space variational Bayes (LSVB) inference.

Since the CVB inference requires intractable integrals, Teh et al. (Teh et al., 2007) used a second-order Taylor expansion to perform the integrals. Asuncion et al. (Asuncion et al., 2009; Asuncion, 2010) proposed another approximation that uses only the zero-order information, called the CVB0 inference. The CVB0 inference does not have the drawbacks that

Table 1. Main results: CVB0 is a special case of α -divergence projection. CVB0 is interpreted as follows: The ($\alpha = 1$)-divergence is used to estimate $n_{d,t}$, which is the number of times topic t appears in document d . The ($\alpha = 1$) divergence is used to estimate $n_{t,v}$, which is the number of times word v is generated from topic t . The ($\alpha = -1$) divergence is used to estimate n_t , which is the number of times topic t appears in the all documents. “EP” indicates the expectation propagation proposed for an aspect model in (Minka & Lafferty, 2002). In this table, the approximation by Taylor expansion is not assumed with “CVB”. “Marginalization” indicates marginalizing out parameters of LDA.

Inference	Marginalization	α -divergence
VB	NA	$\alpha \rightarrow 0$
CVB	✓	$\alpha \rightarrow 0$
CVB0	✓	$\alpha = 1$ for $n_{d,t}, n_{t,v}$ $\alpha = -1$ for n_t
EP	NA	$\alpha = 1$

other inferences do: VB contains digamma functions which are computationally expensive, while CVB requires the maintenance of variance counts. In contrast, the stochastic nature of the collapsed Gibbs sampler causes it to converge more slowly than the deterministic algorithms. Asuncion et al.’s empirical results suggest that the CVB0 inference learns models that are as good as or better than those learned by the VB and CVB inferences and the collapsed Gibbs sampler in terms of perplexity. Furthermore, as shown in (Asuncion, 2010), when the asymmetric Dirichlet parameters are estimated over document-topic distribution, the predictive performance of the CVB0 inference clearly outperforms that of the CVB inference.

We have the question of why CVB0 outperformed CVB, even though the approximation of CVB is more accurate than that of CVB0. In this paper, we propose an interpretation of the CVB0 inference for LDA by using the α -divergence. Using the α -divergence helps

Appearing in *Proceedings of the 29th International Conference on Machine Learning*, Edinburgh, Scotland, UK, 2012. Copyright 2012 by the author(s)/owner(s).

clarify the properties of the CVB0 inference. We also experimentally show the performance of the subspecies of the CVB0 inference, which is derived with the α -divergence projection framework. Our analysis of the relationship between existing inference algorithms and α -divergence is summarized in Table 1, the meaning of which is revealed in later sections.

The remainder of this paper is organized as follows. Sections 3 and 4 review LDA and the CVB / CVB0 inference for LDA, respectively. Sections 5 and 6 explain α -divergence and its local projection, respectively. The key sections 7 and 8 describe local α -divergence projection for LDA and its connection to the CVB0 inference. Section 9 introduces other local projections inspired by the CVB0 inference. Section 10 evaluates algorithms in terms of document modeling. Section 11 concludes this paper.

2. Preliminaries

Suppose that we have N documents, V vocabularies, and T topics. $\mathbf{w} = \{\mathbf{w}_d\}_{d=1}^N$ denotes a set of documents and $\mathbf{z} = \{\mathbf{z}_d\}_{d=1}^N$ is a set of assigned topics. $\theta_{d,t}$ denotes the probability of topic t appearing in document d . $\phi_{t,v}$ denotes the probability of word v appearing in topic t .

$n_{d,t}(\mathbf{z})$ denotes the number of observations of topic t in document d . n_d denotes the total number of words in document d . $n_{t,v}(\mathbf{w}, \mathbf{z})$ denotes the number of observations of word v assigned to topic t and $n_{t,\cdot}(\mathbf{z}) = \sum_v n_{t,v}(\mathbf{w}, \mathbf{z})$. For simplicity, we denote them by $n_{d,t}$, $n_{t,v}$ and $n_{t,\cdot}$. The superscription “ $\setminus d, i$ ” denotes the corresponding variables or counts with $w_{d,i}$ and $z_{d,i}$ excluded, e.g., $\mathbf{w}^{\setminus d, i} = \mathbf{w} \setminus \{w_{d,i}\}$, $\mathbf{z}^{\setminus d, i} = \mathbf{z} \setminus \{z_{d,i}\}$, and $n_{t,v}^{\setminus d, i}$ is the number of observations of word v assigned to topic t leaving out $z_{d,i}$.

$\mathbb{E}[x]$ denotes the expectation of x and $\mathbb{V}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$ the variance. $\text{Multi}(\cdot)$ denotes the multinomial distribution. $\text{Dir}(\cdot)$ denotes the Dirichlet distribution.

3. Overview of LDA

The following generative process is assumed with LDA. First, document-topic distribution $\boldsymbol{\theta}_d$ and topic-word distribution $\boldsymbol{\phi}_k$ are generated by

$$\boldsymbol{\theta}_d \sim \text{Dir}(\boldsymbol{\gamma}), \quad \boldsymbol{\phi}_t \sim \text{Dir}(\boldsymbol{\beta}), \quad (1)$$

where $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_T)$ is a T -dimensional vector and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_V)$ is a V -dimensional vector.

For each document d , generate the i -th topic $z_{d,i}$ and

word $w_{d,i}$:

$$z_{d,i} \sim \text{Multi}(\boldsymbol{\theta}_d), \quad w_{d,i} \sim \text{Multi}(\boldsymbol{\phi}_{z_{d,i}}). \quad (2)$$

Wallach et al. (Wallach et al., 2009) explored the effects of choosing $\boldsymbol{\gamma}$ and $\boldsymbol{\beta}$ in LDA. They found in Markov chain Monte Carlo (MCMC) simulations that using asymmetric $\boldsymbol{\gamma}$ and symmetric $\boldsymbol{\beta}$ results in better predictive performance for held-out documents. Therefore, we use asymmetric $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_T)$ and symmetric $\boldsymbol{\beta} = (\beta, \dots, \beta)$.

The assignment probability of topic t to the i -th word in document d given $\mathbf{w}^{\setminus d, i}$, $\mathbf{z}^{\setminus d, i}$, $\boldsymbol{\gamma}$ and β is

$$\begin{aligned} & p(z_{d,i} = t | w_{d,i} = v, \mathbf{w}^{\setminus d, i}, \mathbf{z}^{\setminus d, i}, \boldsymbol{\gamma}, \beta) \\ & \propto p(w_{d,i} = v, \mathbf{w}^{\setminus d, i}, \mathbf{z}^{\setminus d, i}, z_{d,i} = t | \boldsymbol{\gamma}, \beta), \\ & \propto p(w_{d,i} = v | z_{d,i} = t, \mathbf{w}^{\setminus d, i}, \mathbf{z}^{\setminus d, i} | \beta) p(z_{d,i} = t | \mathbf{z}^{\setminus d, i}, \boldsymbol{\gamma}), \\ & \propto \frac{n_{t,v}^{\setminus d, i} + \beta}{n_{t,\cdot}^{\setminus d, i} + V\beta} (n_{d,t}^{\setminus d, i} + \gamma_t). \end{aligned} \quad (3)$$

This is used for the collapsed Gibbs sampler.

4. CVB/CVB0 inference for LDA

(Teh et al., 2007) proposed the CVB inference to LDA inspired by the collapsed Gibbs sampler and showed that the CVB-LDA outperformed the VB-LDA in terms of perplexity. They only introduced a variational posterior $q(\mathbf{z})$ by marginalizing out $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$. The free energy of the CVB-LDA is given by

$$\mathcal{F}_{CVB}[q(\mathbf{z})] = \sum_{d=1}^M \sum_{\mathbf{z}_d} q(\mathbf{z}_d) \log \frac{p(\mathbf{w}_d, \mathbf{z}_d | \boldsymbol{\gamma}, \boldsymbol{\beta})}{q(\mathbf{z}_d)}. \quad (4)$$

Thus, the updates for $q(\mathbf{z})$ are obtained by taking derivatives of $\mathcal{F}_{CVB}[q(\mathbf{z})]$ with respect to $\{q(z_{d,i})\}$ and equating to zero:

$$\begin{aligned} & q(z_{d,i} = t) \\ & \propto \exp \mathbb{E}[\log p(w_{d,i} = v, \mathbf{w}^{\setminus d, i}, \mathbf{z}^{\setminus d, i}, z_{d,i} = t | \boldsymbol{\gamma}, \boldsymbol{\beta})]_{q(\mathbf{z}^{\setminus d, i})}, \\ & \propto \exp \left\{ \mathbb{E} \left[\log \frac{n_{t,v}^{\setminus d, i} + \beta}{n_{t,\cdot}^{\setminus d, i} + V\beta} (n_{d,t}^{\setminus d, i} + \gamma_t) \right] \right\}, \\ & \propto \frac{\exp \mathbb{E}[\log(n_{t,v}^{\setminus d, i} + \beta)]}{\exp \mathbb{E}[\log(n_{t,\cdot}^{\setminus d, i} + V\beta)]} \exp \mathbb{E}[\log(n_{d,t}^{\setminus d, i} + \gamma_t)]. \end{aligned} \quad (5)$$

This update equation for $q(\mathbf{z})$ requires approximations to compute intractable expectation. By using the central limit theorem, the expectation should be closely approximated using Gaussian distributions

with means and variances, e.g.,

$$\mathbb{E}[n_{d,t}] = \sum_{i=1}^{n_d} q(z_{d,i} = t), \quad (6)$$

$$\mathbb{V}[n_{d,t}] = \sum_{i=1}^{n_d} q(z_{d,i} = t)(1 - q(z_{d,i} = t)). \quad (7)$$

Moreover, using the second order Taylor expansion, we can approximately calculate

$$\begin{aligned} q(z_{d,i} = t) &\propto \frac{\beta + \mathbb{E}[n_{t,w_{d,i}}^{d,i}]}{V\beta + \mathbb{E}[n_{t,\cdot}^{d,i}]} (\gamma_t + \mathbb{E}[n_{d,t}^{d,i}]) \\ &\exp\left(-\frac{\mathbb{V}[n_{t,w_{d,i}}^{d,i}]}{2(\beta + \mathbb{E}[n_{t,w_{d,i}}^{d,i}])^2} + \frac{\mathbb{V}[n_{t,\cdot}^{d,i}]}{2(V\beta + \mathbb{E}[n_{t,\cdot}^{d,i}])^2}\right) \\ &\exp\left(-\frac{\mathbb{V}[n_{d,t}^{d,i}]}{2(\gamma_t + \mathbb{E}[n_{d,t}^{d,i}])^2}\right), \end{aligned} \quad (8)$$

where the superscription “ d, i ” denotes subtracting $q(z_{d,i} = t)$ and $q(z_{d,i} = t)(1 - q(z_{d,i} = t))$.

(Asuncion et al., 2009) showed the usefulness of an approximation using only zero-order information, called the CVB0 inference. The update using only zero-order information is given by

$$q(z_{d,i} = t) \propto \frac{\beta + \mathbb{E}[n_{t,w_{d,i}}^{d,i}]}{V\beta + \sum_v \mathbb{E}[n_{t,v}^{d,i}]} (\gamma_t + \mathbb{E}[n_{d,t}^{d,i}]). \quad (9)$$

We derive this CVB0 inference by using α -divergence, which enables us to reveal the relationship among other inference algorithms.

5. α -Divergence

This section reviews α -divergence. A readable introduction is provided in (Minka, 2005).

Let our task be to approximate a complex probabilistic distribution $p(\mathbf{x})$ where $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$. We approximate $p(\mathbf{x})$ as $q(\mathbf{x})$, which is a simple probabilistic distribution, such as fully factorized distribution, i.e., $q(\mathbf{x}) = \prod_{i=1}^n q(x_i)$. A basic approach to obtaining $q(\mathbf{x})$ is to minimize information divergence such as the Kullback-Leibler divergence:

$$KL[p||q] = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} + \int (q(\mathbf{x}) - p(\mathbf{x})) d\mathbf{x}, \quad (10)$$

where $p(\mathbf{x})$ and $q(\mathbf{x})$ do not need to be normalized. By using the KL-divergence, the estimation of $q(\mathbf{x})$ is defined by the KL-projection of $p(\mathbf{x})$ onto a family of $q(\mathbf{x})$ as follows:

$$q^*(\mathbf{x}) = \operatorname{argmin}_{q(\mathbf{x})} KL[p(\mathbf{x})||q(\mathbf{x})]. \quad (11)$$

α -divergence is a generalization of the KL divergence (Amari, 1985; Trottni & Spezzaferrri, 2002; Zhu & Rohwer, 1995), indexed by $\alpha \in (-\infty, \infty)$. The α parameter can be used in different ways by different authors. In this paper, we define α -divergence by the convention used in (Minka, 2005):

$$D_\alpha[p||q] = \frac{\int \alpha p(\mathbf{x}) + (1 - \alpha)q(\mathbf{x}) - p(\mathbf{x})^\alpha q(\mathbf{x})^{1-\alpha} d\mathbf{x}}{\alpha(1 - \alpha)}, \quad (12)$$

where $p(\mathbf{x})$ and $q(\mathbf{x})$ do not need to be normalized. If $p = q$, α -divergence is zero. Some special cases are

$$D_{-1}[p||q] = \frac{1}{2} \int \frac{(q(x) - p(x))^2}{p(x)} dx \quad (13)$$

$$\lim_{\alpha \rightarrow 0} D_\alpha[p||q] = KL[q(\mathbf{x})||p(\mathbf{x})] \quad (14)$$

$$D_{0.5}[p||q] = 2 \int (\sqrt{q(x)} - \sqrt{p(x)})^2 dx \quad (15)$$

$$\lim_{\alpha \rightarrow 1} D_\alpha[p||q] = KL[p(\mathbf{x})||q(\mathbf{x})] \quad (16)$$

$$D_2[p||q] = \frac{1}{2} \int \frac{(p(x) - q(x))^2}{q(x)} dx. \quad (17)$$

The case $\alpha = 0.5$ is known as the Hellinger distance, and $\alpha = 2$ is the χ^2 distance. Since $\alpha = -1$ swaps the position of p and q of the χ^2 distance, we call the case $\alpha = -1$ “the inverse χ^2 distance”, which is the key divergence in this paper.

6. Local α -divergence projection

In this section, we introduce a local divergence projection-based inference.

Suppose that the approximate distribution $q(\mathbf{x})$ is fully factorized. We derive the update $q(x_i)$ minimizing α -divergence as follows. Taking derivatives of α -divergence (12) with respect to $q(x_i)$ and equating them to zero, we obtain the following fixed point iteration equations:

$$q(x_i) \propto \mathbb{E} \left[\left(\frac{p(\mathbf{x})}{q(\mathbf{x}^{\setminus i})} \right)^\alpha \right]_{q(\mathbf{x}^{\setminus i})}^{\frac{1}{\alpha}} \quad (18)$$

In many cases, this update is intractable and thus we introduce an approximation for Eq. (18).

Since Eq. (18) is

$$q(x_i) \propto \mathbb{E} \left[\left(p(x_i|\mathbf{x}^{\setminus i}) \frac{p(\mathbf{x}^{\setminus i})}{q(\mathbf{x}^{\setminus i})} \right)^\alpha \right]_{q(\mathbf{x}^{\setminus i})}^{\frac{1}{\alpha}}, \quad (19)$$

we replace $p(\mathbf{x}^{\setminus i})$ with $q(\mathbf{x}^{\setminus i})$, obtaining

$$q(x_i) \propto \mathbb{E} \left[p(x_i|\mathbf{x}^{\setminus i})^\alpha \right]_{q(\mathbf{x}^{\setminus i})}^{\frac{1}{\alpha}}. \quad (20)$$

In the case $\alpha = 1$, the update (20) is similar to belief propagation, and the factorized neighbors algorithm (Rosen-Zvi et al., 2005).

The update (20) means that it locally minimize α -divergence, i.e., for each i ,

$$q^*(x_i) = \operatorname{argmin}_{q(x_i)} D_\alpha[p(x_i|\mathbf{x}^{\setminus i})q(\mathbf{x}^{\setminus i})||q(\mathbf{x})]. \quad (21)$$

In the case $\alpha = 1$, i.e., KL divergence, this local projection-based inference is equal to the EP algorithm. We describe the connection of this α -divergence projection with the CVB0 inference in the next section.

7. CVB0 as α -divergence projection

In this section, we derive the CVB0 inference by using the local α -divergence projection. First, we describe how the case $\alpha = 1$, i.e. EP, cannot be applied for the collapsed LDA. Second, we derive a divergence projection applicable to the collapsed LDA and explain the relationship between this projection and the CVB0 inference.

We apply Eq. (21) with $\alpha = 1$ (EP) to the collapsed LDA. For each $z_{d,i}$, we perform

$$q^*(z_{d,i}) = \operatorname{argmin}_{q(z_{d,i})} KL[p(z_{d,i}|\mathbf{w}, \mathbf{z}^{\setminus d,i})q(\mathbf{z}^{\setminus d,i})||q(\mathbf{z})]. \quad (22)$$

The update for $q(z_{d,i})$ is

$$\begin{aligned} q(z_{d,i} = t) &\propto \mathbb{E} \left[p(z_{d,i} = t | w_{d,i} = v, \mathbf{w}^{\setminus d,i}, \mathbf{z}^{\setminus d,i}) \right]_{q(\mathbf{z}^{\setminus d,i})}, \\ &\propto \mathbb{E} \left[(n_{d,t}^{\setminus d,i} + \gamma_t) \frac{n_{t,v}^{\setminus d,i} + \beta}{n_{t,\cdot}^{\setminus d,i} + V\beta} \right]_{q(\mathbf{z}^{\setminus d,i})}. \end{aligned} \quad (23)$$

The problem is that we cannot analytically execute this expectation. (Asuncion, 2010) derived Eq.(23) in a different way where he changed the CVB free energy by moving the logarithm out of the expectations, and pointed out the relationship between Eq.(23) and the CVB0 inference, which inspired this work. However, the intractable expectation in Eq.(23) was not executed. This intractability makes interpreting the CVB0 inference difficult.

Here, we derive another approach by using the α -divergence projection. The key idea is to construct $q(z_{d,i})$ by using the novel three parameters .

We define $q(z_{d,i})$ as follows:

$$q(z_{d,i} = t) \propto a(z_{d,i})b(z_{d,i})c(z_{d,i}) \quad (24)$$

$$a(z_{d,i} = t) = \tilde{n}_{d,t}^{\setminus d,i} + \gamma_t, \quad (25)$$

$$b(z_{d,i} = t) = \tilde{n}_{t,v}^{\setminus d,i} + \beta, \quad (26)$$

$$c(z_{d,i} = t) = \frac{1}{\tilde{n}_{t,\cdot}^{\setminus d,i} + V\beta}, \quad (27)$$

where we do not assume that $\tilde{n}_{d,i}^{\setminus d,i}$, $\tilde{n}_{t,v}^{\setminus d,i}$ and $\tilde{n}_{t,\cdot}^{\setminus d,i}$ are expected counts, i.e., these are parameters of $q(z_{d,i})$.

We also define

$$q^{\setminus a}(z_{d,i}) = b(z_{d,i})c(z_{d,i}), \quad (28)$$

$$q^{\setminus b}(z_{d,i}) = a(z_{d,i})c(z_{d,i}) \quad (29)$$

$$q^{\setminus c}(z_{d,i}) = a(z_{d,i})b(z_{d,i}). \quad (30)$$

Since our definition of α -divergence does not require normalization of the probabilistic distribution, we can introduce the following local projection:

$$\begin{aligned} a^*(z_{d,i} = t) &= \\ &\operatorname{argmin}_{a(z_{d,i})} D_\alpha[(n_{d,t}^{\setminus d,i} + \gamma_t)q^{\setminus a,d,i}(\mathbf{z})||a(z_{d,i})q^{\setminus a,d,i}(\mathbf{z})], \end{aligned} \quad (31)$$

where $q^{\setminus a,d,i}(\mathbf{z}) = q^{\setminus a}(z_{d,i})q(\mathbf{z}^{\setminus d,i})$. Solving the above optimization (see Appendix A), we obtain

$$a^*(z_{d,i} = t) = \mathbb{E} \left[(n_{d,t}^{\setminus d,i} + \gamma_t)^\alpha \right]_{q(\mathbf{z}^{\setminus d,i})}^{\frac{1}{\alpha}}. \quad (32)$$

As in $a(z_{d,i})$, we obtain $b^*(z_{d,i})$ and $c^*(z_{d,i})$ by locally minimizing the α -divergence:

$$\begin{aligned} b^*(z_{d,i} = t) &= \\ &\operatorname{argmin}_{b(z_{d,i})} D_\alpha[(n_{t,v}^{\setminus d,i} + \beta)q^{\setminus b,d,i}(\mathbf{z})||b(z_{d,i})q^{\setminus b,d,i}(\mathbf{z})], \end{aligned} \quad (33)$$

$$\begin{aligned} c^*(z_{d,i} = t) &= \\ &\operatorname{argmin}_{c(z_{d,i})} D_\alpha\left[\frac{1}{(n_{t,\cdot}^{\setminus d,i} + V\beta)}q^{\setminus c,d,i}(\mathbf{z})||c(z_{d,i})q^{\setminus c,d,i}(\mathbf{z})\right]. \end{aligned} \quad (34)$$

Thus, we have

$$b^*(z_{d,i} = t) = \mathbb{E} \left[(n_{t,v}^{\setminus d,i} + \beta)^\alpha \right]_{q(\mathbf{z}^{\setminus d,i})}^{\frac{1}{\alpha}}, \quad (35)$$

$$c^*(z_{d,i} = t) = \mathbb{E} \left[\left(\frac{1}{n_{t,\cdot}^{\setminus d,i} + V\beta} \right)^\alpha \right]_{q(\mathbf{z}^{\setminus d,i})}. \quad (36)$$

When we use α -divergence projection with $\alpha = 1$ for

estimating $a(z_{d,i})$ and $b(z_{d,i})$, we have

$$a^{(\alpha=1)}(z_{d,i} = t) = \mathbb{E} \left[n_{d,t}^{d,i} + \gamma_t \right]_{q(\mathbf{z} \setminus d,i)} = \mathbb{E}[n_{d,t}^{d,i}] + \gamma_t, \quad (37)$$

$$b^{(\alpha=1)}(z_{d,i} = t) = \mathbb{E} \left[n_{t,v}^{d,i} + \beta \right]_{q(\mathbf{z} \setminus d,i)} = \mathbb{E}[n_{t,v}^{d,i}] + \beta. \quad (38)$$

When we use α -divergence projection with $\alpha = -1$ for estimating $c(z_{d,i})$, we have

$$\begin{aligned} c^{(\alpha=-1)}(z_{d,i} = t) &= \mathbb{E} \left[\left(\frac{1}{n_{t,\cdot}^{d,i} + V\beta} \right)^{-1} \right]_{q(\mathbf{z} \setminus d,i)}, \\ &= \mathbb{E} \left[n_{t,\cdot}^{d,i} + V\beta \right]_{q(\mathbf{z} \setminus d,i)}^{-1}, \\ &= \frac{1}{\mathbb{E}[n_{t,\cdot}^{d,i}] + V\beta}. \end{aligned} \quad (39)$$

Therefore, we have the following update for $q(z_{d,i})$

$$\begin{aligned} q(z_{d,i} = t) &\propto a^{(\alpha=1)}(z_{d,i})b^{(\alpha=1)}(z_{d,i})c^{(\alpha=-1)}(z_{d,i}), \\ &= (\mathbb{E}[n_{d,t}^{d,i}] + \gamma) \frac{\mathbb{E}[n_{t,v}^{d,i}] + \beta}{\mathbb{E}[n_{t,\cdot}^{d,i}] + V\beta}. \end{aligned} \quad (40)$$

Although the updates are performed in order, i.e., update a^* given b and c , b^* given a^* and c , and c^* given a^* and b^* , this update is equal to the CVB0 update in Eq.(9).

8. Discussion

In this section, we explain why the CVB0 inference outperforms the CVB inference. To sum up this discussion, in the CVB0 inference, the ‘‘zero-forcing effect’’ works only with the $n_{t,\cdot}$ estimation, while in the CVB inference it works with the $q(z)$ estimation.

The previous section showed that the CVB0 inference is composed of the three projections with a mixture of $\alpha = 1$ and $\alpha = -1$:

$$D_1 = KL[(n_{d,t}^{d,i}(\mathbf{z}) + \gamma_t)q^{a_{d,i}}(\mathbf{z})||q(\mathbf{z})], \quad (41)$$

$$D_1 = KL[(n_{t,v}^{d,i}(\mathbf{z}) + \beta)q^{b_{d,i}}(\mathbf{z})||q(\mathbf{z})], \quad (42)$$

$$D_{-1} \left[\frac{1}{(n_{t,\cdot}^{d,i}(\mathbf{z}) + V\beta)} q^{c_{d,i}}(\mathbf{z}) || q(\mathbf{z}) \right]. \quad (43)$$

This projection-based update with a different divergence measure reveals the properties of the CVB0 inference. Ideally, we use the ($\alpha = 1$)-divergence projection, i.e., $D_1[p||q] = KL[p||q]$, but the integrals

$\mathbb{E}[\frac{1}{n_{t,\cdot}^{d,i} + V\beta}]$ are not easy to evaluate. Instead, we use the inverse χ^2 divergence $D_{-1}[p||q]$ for estimating $c(z_{d,i})$.

$D_{-1}[p||q] = \frac{1}{2} \int \frac{(q(x)-p(x))^2}{p(x)} dx$ is known as a zero-forcing divergence (Minka, 2005) which emphasizes q to be small when p being small, i.e., $p(x) = 0$ forces $q(x) = 0$, which means that it avoids ‘‘false positive’’. In our case (43), the zero-forcing effect on the $n_{t,\cdot}$ estimation means that the emphasis in the estimation is on high-frequency topics or low-frequency topics tend to be estimated as zero in an entire corpus. We think that affecting $n_{t,\cdot}^{d,i}$ matters much less than affecting $n_{d,t}^{d,i}$ and $n_{t,v}^{d,i}$ throughout a whole corpus in LDA. We explain the zero-forcing effect of CVB0 in more detail in the next section.

Returning to Eq.(20), i.e., $q(x_i) \propto \mathbb{E} [p(x_i|\mathbf{x}^{\setminus i})^\alpha]_{q(\mathbf{x}^{\setminus i})}^{\frac{1}{\alpha}}$, we describes the relationship between the CVB inference and α -divergence projection. First, we introduce the following theorem:

Theorem 1 (Liapunov’s inequality) *If x is a non-negative random variable, and we have two real numbers $\alpha_2 > \alpha_1$, then*

$$\mathbb{E}[x^{\alpha_2}]^{\frac{1}{\alpha_2}} \geq \mathbb{E}[x^{\alpha_1}]^{\frac{1}{\alpha_1}}. \quad (44)$$

and

$$\lim_{\alpha \rightarrow 0} \mathbb{E}[x^\alpha]^{\frac{1}{\alpha}} = \exp \mathbb{E}[\log(x)]. \quad (45)$$

Using Eq.(20) and Theorem 1, we obtain

$$q(x_i) \propto \lim_{\alpha \rightarrow 0} \mathbb{E} \left[p(x_i|\mathbf{x}^{\setminus i})^\alpha \right]_{q(\mathbf{x}^{\setminus i})}^{\frac{1}{\alpha}} = \exp(\mathbb{E}[\log p(x_i|\mathbf{x}^{\setminus i})]) \quad (46)$$

This is the variational inference minimizing $KL[q||p]$.

In LDA, we have

$$\begin{aligned} q(z_{d,i} = t) &\propto \lim_{\alpha \rightarrow 0} \mathbb{E} \left[p(z_{i,d}|w_{d,i} = v, \mathbf{w}^{\setminus d,i}, \mathbf{z}^{\setminus d,i})^\alpha \right]_{q(\mathbf{z} \setminus d,i)}^{\frac{1}{\alpha}} \\ &= \exp(\mathbb{E}[\log p(z_{d,i}|w_{d,i} = v, \mathbf{w}^{\setminus d,i}, \mathbf{z}^{\setminus d,i})]) \\ &\propto \exp \mathbb{E} \left[\log \frac{n_{d,t}^{d,i} + \gamma_t}{n_{d,\cdot}^{d,i} + \sum_t \gamma_t} \frac{n_{t,v}^{d,i} + \beta}{n_{t,\cdot}^{d,i} + V\beta} \right]_{q(\mathbf{z} \setminus d,i)}, \\ &\propto \exp(\mathbb{E}[\log(n_{d,t}^{d,i} + \gamma_t)]) \frac{\exp(\mathbb{E}[\log(n_{t,v}^{d,i} + \beta)])}{\exp(\mathbb{E}[\log(n_{t,\cdot}^{d,i} + V\beta)])} \end{aligned} \quad (47)$$

The update Eq.(47) is the same update as the CVB inference in Eq.(5). ($\alpha \rightarrow 0$)-divergence is also known to induce the zero-forcing effect.

9. Subspecies inspired by CVB0

In this section, we consider other projection-based algorithms that help clarify the property of the zero-forcing effect in CVB0.

9.1. CVB with $(\alpha = 1)$ -divergence

From our view point, the CVB0 inference is composed of two different-type divergence projections: $\alpha = 1, -1$. We consider using only $\alpha = 1$ for the projections. To do this, we have to calculate the expectation given by

$$c^{(\alpha=1)}(z_{d,i} = t) = \mathbb{E} \left[\frac{1}{n_{t,\cdot}^{\setminus d,i} + V\beta} \right]_{q(\mathbf{z}^{\setminus d,i})}. \quad (48)$$

Since we cannot derive the analytical solution for this expectation, we propose two approximation methods. The first is a stochastic approximation called sample averaging given by

$$\tilde{c}^{(\alpha=1)}(z_{d,i} = t) = \frac{1}{S} \sum_{s=1}^S \frac{1}{n_{t,\cdot}^{\setminus d,i}(\mathbf{z}^{(s)}) + V\beta}, \quad (49)$$

where S denotes the number of samples and $\mathbf{z}^{(s)}$ is the s -th samples generated from $q(\mathbf{z})$. This method is accurate but not practical when S takes a large value. We use this approximation to investigate the accuracy of the next approximation.

The second is a deterministic approximation that uses the same approximation of CVB with the second-order Taylor expansion and Gaussian approximation given by

$$\hat{c}^{(\alpha=1)}(z_{d,i} = t) = \frac{1}{\mathbb{E}[n_{t,\cdot}^{\setminus d,i}] + V\beta} + \frac{\mathbb{V}[n_{t,\cdot}^{\setminus d,i}]}{(\mathbb{E}[n_{t,\cdot}^{\setminus d,i}] + V\beta)^3}. \quad (50)$$

As shown in the experiments (Sec.10), we find that the second term of Eq.(50) is vanishingly small. $\hat{c}^{(\alpha=1)}$ in Eq.(50) is calculated as $\frac{1}{\mathbb{E}[n_{t,\cdot}^{\setminus d,i}] + V\beta} \left(1 + \frac{\mathbb{V}[n_{t,\cdot}^{\setminus d,i}]}{(\mathbb{E}[n_{t,\cdot}^{\setminus d,i}] + V\beta)^2} \right)$. We find $\frac{\mathbb{V}[n_{t,\cdot}^{\setminus d,i}]}{(\mathbb{E}[n_{t,\cdot}^{\setminus d,i}] + V\beta)^2} = O(1/n)$ in many cases where n denotes the number of all words (tokens). For example, the variance takes the largest value when $q(z_{d,i} = t) = 1/2$ for all d and i . In this case, $\mathbb{E}[n_t] = n/2$ and $\mathbb{V}[n_t] = n(1 - 1/2)/2 = n/4$. Therefore, we consider $c^{(\alpha=-1)}$ is similar to $c^{(\alpha=1)}$, which means that CVB0 is rarely affected by the zero-forcing effect.

9.2. Type-base CVB0 Inference

We derive a type-based inference as an application of our framework. In a type-based inference, we only estimate the probabilistic distribution for each type in a document not each token; this is beneficial for computation cost and memory usage.

We exclude all counts of word v from document d , denoted by superscription " $\setminus d, v$ ". The appearance probability of word v given $\mathbf{w}_d^{\setminus d,v}$ and $\mathbf{z}_d^{\setminus d,v}$ is

$$p(w_{d,*} = v | \mathbf{w}^{\setminus d,v}, \mathbf{z}^{\setminus d,v}) = \sum_{t=1}^T \frac{n_{d,t}^{\setminus d,v} + \gamma_t}{n_{d,\cdot}^{\setminus d,v} + \sum_t \gamma_t} \frac{n_{t,v}^{\setminus d,v} + \beta}{n_{t,\cdot}^{\setminus d,v} + V\beta} \quad (51)$$

Moreover, we have

$$p(z_{d,v} = t | \mathbf{w}^{\setminus d,v}) \propto \mathbb{E} \left[\frac{n_{d,t}^{\setminus d,v} + \gamma_t}{n_{d,\cdot}^{\setminus d,v} + \sum_t \gamma_t} \frac{n_{t,v}^{\setminus d,v} + \beta}{n_{t,\cdot}^{\setminus d,v} + V\beta} \right]_{p(\mathbf{z}^{\setminus d,v} | \mathbf{w}^{\setminus d,v})} \quad (52)$$

Here, we consider obtaining an approximation distribution $q(z_{d,v})$. Instead of $z_{d,i}$, we define $q(\mathbf{z})$ factorized by using $q(z_{d,v})$, i.e., $q(z_{d,i}) = \sum_{v=1}^V q(z_{d,v}) \delta(w_{d,i} = v)$ and $q(\mathbf{z}) = \prod_{v=1}^V q(z_{d,v})^{n_{d,v}}$.

The update of $q(z_{d,v})$ is obtained by

$$q(z_{d,v} = t) \propto \mathbb{E} \left[(n_{d,t}^{\setminus d,v} + \gamma_t) \frac{n_{t,v}^{\setminus d,v} + \beta}{n_{t,\cdot}^{\setminus d,v} + V\beta} \right]_{q(\mathbf{z}^{\setminus d,v})}, \quad (53)$$

which is derived by minimizing the α -divergence as in $q(z_{d,i})$.

Using the local α -divergence projection with $\alpha = 1$ for $n_{d,t}^{\setminus d,v} + \gamma_t$ and $n_{t,v}^{\setminus d,v} + \beta$, and $\alpha = -1$ for $\frac{1}{n_{t,\cdot}^{\setminus d,v} + V\beta}$, we have

$$q(z_{d,v} = t) \propto (\mathbb{E}[n_{d,t}^{\setminus d,v}] + \gamma_t) \frac{\mathbb{E}[n_{t,v}^{\setminus d,v}] + \beta}{\mathbb{E}[n_{t,\cdot}^{\setminus d,v}] + V\beta}. \quad (54)$$

We call this update the type-based CVB0 (TCVB0) inference.

10. Experiments

We compared CVB0 with its subspecies on document modeling in terms of perplexity to investigate the effect of $\alpha = -1$. All results are averaged values from five experimental runs with random initialization. We set the number of iterations to 100 for each inference.

We use a fixed point equation for updating γ introduced in (Minka, 2000). We set $\beta = 0.01$ because

(Asuncion et al., 2009) showed that CVB0-LDA with $\beta = 0.01$ worked well when compared with other settings ($\beta = 0.01$ was also used in (Griffiths & Steyvers, 2004)).

In this section’s figures, “CVB” indicates the second order approximation of the CVB inference.

“CVB1s” indicates the stochastic approximation in Eq.(49) with $S = 50$. “CVB1d” indicates the deterministic approximation in Eq.(50).

We used four sets of text data with different properties. The first was ‘NIPS corpus (NIPS)’ from which the number of documents was $N = 1,500$ and the vocabulary size was $V = 12,245$. The second was “The Wall Street Journal (WSJ)” from which we randomly chose $N = 5,000$ ($V = 38,272$) documents. The third was “Enron email corpus (Enron)” from which we randomly chose $N = 5,000$ ($V = 14,758$) documents. The fourth was “20 news group corpus (20ng)” from which we randomly chose $N = 5,000$ ($V = 13,176$). Stop words were eliminated.

The comparison metric we used for document modeling was the perplexity used by (Teh et al., 2007; 2008) that indicates the prediction performance for held-out words. We randomly split the words in a document into training words \mathbf{w}_d^{train} (80%) and test words \mathbf{w}_d^{test} (20%).

Figure 1 shows the experimental results. The bar graph indicates the results for test set perplexity in terms of ($T = 40, 80, 120$) in each corpus. CVB0, CVB1s, CVB1d and TCVB0 outperformed CVB in terms of perplexity. Although we compared VB with others, we eliminated the VB results to clarify the differences of inference algorithms because CVB outperformed VB and the VB results change the scale of a bar-graph in some corpora.

The performances of CVB1s and CVB1d were similar to that of CVB0. Since the results of CVB1d were similar to those of CVB1s, the approximation used in CVB1d seemed to be accurate. When we analyzed $\frac{\mathbb{V}[n_{i,\cdot}]}{(\mathbb{E}[n_{i,\cdot}] + V\beta)^2}$ in Sec.9.1, the maximum value in all corpus when $T = 120$ was about $3.17e^{-4}$, which is negligible compared with 1. Therefore, as discussed in Sec.9.1, CVB0 was not affected by the zero-forcing effect. We believe this is the reason CVB0 worked better than CVB. Moreover, the performance of TCVB0 was similar to that of CVB0. Consequently, the TCVB0 inference was practical.

11. Conclusion

In this paper, we reviewed existing inference algorithms of LDA in terms of the α -divergence projection. We showed that the CVB0 inference is composed of ($\alpha = 1, -1$)-divergence projections and that $\alpha = -1$ is similar to $\alpha = 1$ in LDA, which means that CVB0 is not affected by the zero-forcing effect in LDA. Combining the marginalization of parameters and the heterogeneous α -divergence projection is useful because it is easy to apply to other topic models learned by the collapsed Gibbs sampler. Future work is to develop an online-update extension, such as that by (Hoffman et al., 2010; Sato et al., 2010; Wang et al., 2011). From the relationship between EP and assumed density filtering, we can extend the local α -divergence projection into an online algorithm, which leads to the online CVB0 inference. A convergence analysis is also important remaining work.

A. Derivation for Eq.(32)

Taking derivatives of

$$D_\alpha[(n_{d,t}^{d,i} + \gamma_t)q^{a_{d,i}}(\mathbf{z})|a(z_{d,i})q^{a_{d,i}}(\mathbf{z})],$$

with respect to $a(z_{d,i})$ and equating them to zero,

$$0 = \sum_{\mathbf{z} \setminus d,i} q^{a_{d,i}}(\mathbf{z}) - a(z_{d,i})^{-\alpha} \sum_{\mathbf{z} \setminus d,i} (n_{d,t}^{d,i} + \gamma_t)^\alpha q^{a_{d,i}}(\mathbf{z}),$$

and we obtain the following fixed point iteration equations:

$$\begin{aligned} a(z_{d,i}) &= \left[\frac{\sum_{\mathbf{z} \setminus d,i} (n_{d,t}^{d,i} + \gamma_t)^\alpha q^{a_{d,i}}(\mathbf{z})}{\sum_{\mathbf{z} \setminus d,i} q^{a_{d,i}}(\mathbf{z})} \right]^{\frac{1}{\alpha}}, \\ &= \left[\frac{\sum_{\mathbf{z} \setminus d,i} (n_{d,t}^{d,i} + \gamma_t)^\alpha b(z_{d,i})c(z_{d,i})q(\mathbf{z} \setminus d,i)}{\sum_{\mathbf{z} \setminus d,i} b(z_{d,i})c(z_{d,i})q(\mathbf{z} \setminus d,i)} \right]^{\frac{1}{\alpha}}, \\ &= \left[\frac{b(z_{d,i})c(z_{d,i}) \sum_{\mathbf{z} \setminus d,i} (n_{d,t}^{d,i} + \gamma_t)^\alpha q(\mathbf{z} \setminus d,i)}{b(z_{d,i})c(z_{d,i}) \sum_{\mathbf{z} \setminus d,i} q(\mathbf{z} \setminus d,i)} \right]^{\frac{1}{\alpha}}. \end{aligned} \quad (55)$$

Since $\sum_{\mathbf{z} \setminus d,i} q(\mathbf{z} \setminus d,i) = 1$, we have

$$a(z_{d,i}) = \left[\sum_{\mathbf{z} \setminus d,i} (n_{d,t}^{d,i} + \gamma_t)^\alpha q(\mathbf{z} \setminus d,i) \right]^{\frac{1}{\alpha}}. \quad (56)$$

References

- Amari, S. *Differential-Geometrical Methods in Statistic*. Springer, New York, 1985.
- Asuncion, A. Approximate mean field for dirichlet-based models. In *Topic Models Workshop, ICML*. 2010.

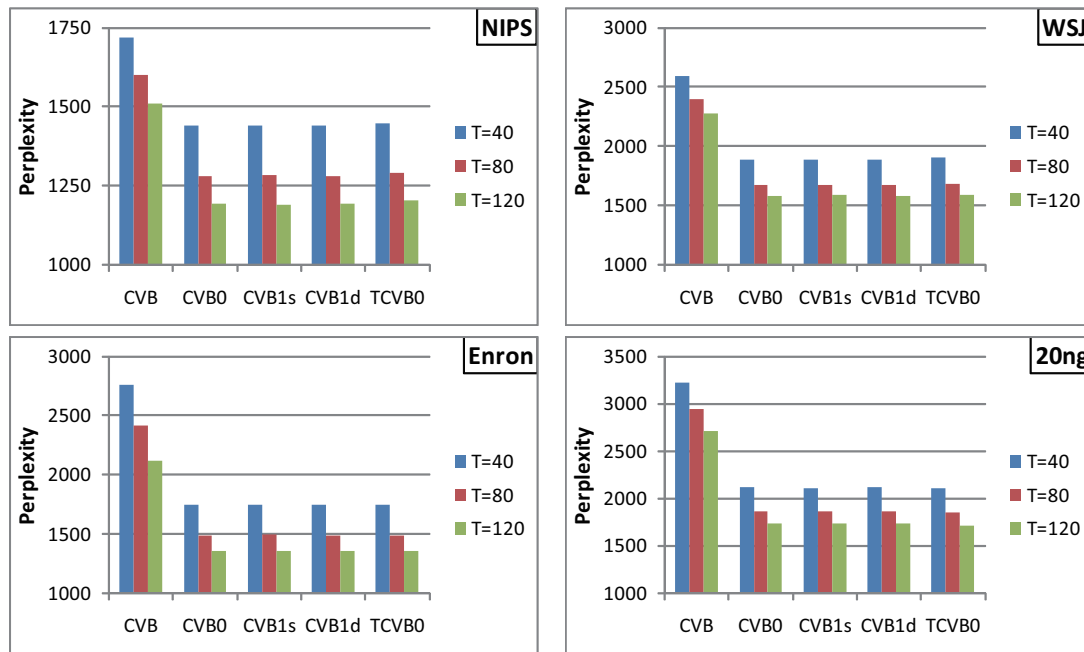


Figure 1. Experiment results for document modeling in four datasets. T denotes the number of topics. Lower perplexity indicates better performance.

Asuncion, A., Welling, M., Smyth, P., and Teh, Y. W. On smoothing and inference for topic models. In *Proceedings of the International Conference on Uncertainty in Artificial Intelligence*, 2009.

Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3: 993–1022, 2003.

Griffiths, T. L. and Steyvers, M. Finding scientific topics. *Proc Natl Acad Sci U S A*, 101 Suppl 1:5228–5235, 2004. ISSN 0027-8424.

Hoffman, Matthew D., Blei, David M., and Bach, Francis R. Online learning for latent dirichlet allocation. In *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010*, pp. 856–864, 2010.

Minka, T. and Lafferty, J. Expectation-Propagation for the generative aspect model. In *Proceedings of the 18th Conference on Uncertainty in Artificial Intelligence (UAI)*. Morgan Kaufmann, 2002.

Minka, Thomas. Divergence measures and message passing. Technical report, Microsoft Research, 2005.

Minka, Thomas P. Estimating a dirichlet distribution. Technical report, Microsoft, 2000.

Rosen-Zvi, Michal, Jordan, Michael I., and Yuille, Alan L. The dir hierarchy of approximate inference. In *Proceedings of the 21st Conference in Uncertainty in Artificial Intelligence*, pp. 493–500, 2005.

Sato, Issei, Kurihara, Kenichi, and Nakagawa, Hiroshi. Deterministic single-pass algorithm for lda. In *Advances in Neural Information Processing Systems 23*, pp. 2074–2082. 2010.

Sung, Jaemo, Ghahramani, Zoubin, and Bang, Sung-Yang. Latent-space variational bayes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30:2236–2242, 2008. ISSN 0162-8828.

Teh, Yee Whye, Newman, David, and Welling, Max. A collapsed variational Bayesian inference algorithm for latent Dirichlet allocation. In *Advances in Neural Information Processing Systems 19*, 2007.

Teh, Yee Whye, Kurihara, Kenichi, and Welling, Max. Collapsed variational inference for hdp. In *Advances in Neural Information Processing Systems 20*, 2008.

Trottini, M. and Spezzaferri, F. A generalized predictive criterion for model selection. In *Canadian Journal of Statisticse*, 2002.

Wallach, Hanna, Mimno, David, and McCallum, Andrew. Rethinking lda: Why priors matter. In *Advances in Neural Information Processing Systems 22*, pp. 1973–1981. 2009.

Wang, Chong, Paisley, John William, and Blei, David M. Online variational inference for the hierarchical dirichlet process. *Journal of Machine Learning Research - Proceedings Track*, 15:752–760, 2011.

Zhu, Huaiyu and Rohwer, Richard. Information geometric measurements of generalisation. Technical report, Aston University, 1995.