

Stochastic Collapsed Variational Bayesian Inference for Latent Dirichlet Allocation

James Foulds^{1*} Levi Boyles^{1†} Christopher Dubois^{2‡}
Padhraic Smyth^{1§}
Max Welling^{3¶}

¹Department of Computer Science, University of California, Irvine

²Department of Statistics, University of California, Irvine

³Informatics Institute, University of Amsterdam

May 14, 2013

Abstract

In the internet era there has been an explosion in the amount of digital text information available, leading to difficulties of scale for traditional inference algorithms for topic models. Recent advances in stochastic variational inference algorithms for latent Dirichlet allocation (LDA) have made it feasible to learn topic models on large-scale corpora, but these methods do not currently take full advantage of the collapsed representation of the model. We propose a stochastic algorithm for collapsed variational Bayesian inference for LDA, which is simpler and more efficient than the state of the art method. We show connections between collapsed variational Bayesian inference and MAP estimation for LDA, and leverage these connections to prove convergence properties of the proposed algorithm. In experiments on large-scale text corpora, the algorithm was found to converge faster and often to a better solution than the previous method. Human-subject experiments also demonstrated that the method can learn coherent topics in seconds on small corpora, facilitating the use of topic models in interactive document analysis software.

*jfoulds@ics.uci.edu

†lboyles@uci.edu

‡duboisc@ics.uci.edu

§smyth@ics.uci.edu

¶m.welling@uva.nl

1 Introduction

Topic models such as latent Dirichlet allocation (LDA) [7] have become a fixture in modern machine learning. Inference algorithms for topic models provide a low-dimensional representation of text corpora that is typically semantically meaningful, despite being completely unsupervised. Their use has spread beyond machine learning to become a standard analysis tool for researchers in many fields [14, 3, 15]. In the internet era, there is a need for tools to learn topic models at the “web scale”, especially in an industrial setting. For example, companies such as Yahoo! publish a continually updated stream of online articles, and needs to analyse candidate articles for topical diversity and relevance to current trends, which could be facilitated by topic models.

We would therefore like to have the tools to build topic models that scale to such large corpora, taking advantage of the large amounts of available data to create models that are accurate and contain more topics. Traditional inference techniques such as Gibbs sampling and variational inference do not readily scale to corpora containing millions of documents or more. In such cases it is very time-consuming to run even a single iteration of the standard collapsed Gibbs sampling [11] or variational Bayesian inference algorithms [7], let alone run them until convergence. The first few passes through the data for these algorithms are inhibited by randomly initialized values that misinform the updates, so multiple such expensive iterations are required to learn the topics.

A significant recent advance was made by Hoffman et al. [13], who proposed a stochastic variational inference algorithm for LDA topic models. Because the algorithm does not need to see all of the documents before updating the topics, this method can often learn good topics before a single iteration of the traditional batch inference algorithms would be completed. The algorithm processes documents in an online fashion, so it can be applied to corpora of any size, or even to never-ending streams of documents. A more scalable variant of this algorithm was proposed by Mimno et al. [16], which approximates the gradient updates in a sparse way in order to improve performance for larger vocabularies and greater numbers of topics.

A complementary direction that has been useful for improving inference in Latent Dirichlet allocation is to take advantage of its “collapsed” representation, where parameters are marginalized out, leaving only latent variables. It is possible to perform inference in the collapsed space and recover estimates of the parameters afterwards. For existing inference techniques that operate in a batch setting, the algorithms that operate in the collapsed space are more efficient at improving held-out log probability than their uncollapsed counterparts, both per iteration and in wall-clock time per iteration [11, 24, 2]. Reasons for this advantage include the per-token updates which propagate updated information sooner, simpler update equations, fewer parameters to update, no expensive calls to the digamma function, and the avoidance of tightly coupled pairs of parameters which inhibit mixing for Gibbs sampling [10, 2, 24]. For variational infer-

ence, perhaps the most important advantage of the collapsed representation is that the variational bound is strictly better than for the uncollapsed representation, leading to the potential to learn more accurate topic models [24]. The existing online inference algorithms for LDA do not fully take advantage of the collapsed representation – although the sparse online LDA algorithm of Mimno et al. [16] collapses out per-document parameters θ , the topics themselves are not collapsed so there is no improvement in the variational bound.

In this work, we develop a stochastic algorithm for LDA that operates fully in the collapsed space, thus transferring the aforementioned advantages of collapsed inference to the online setting. This facilitates learning topic models both more accurately and more quickly on large datasets. The proposed algorithm is also very simple to implement, requiring only basic arithmetic operations. We show that from another perspective, the algorithm can also be interpreted as a MAP estimation algorithm. This interpretation allows us to prove the convergence of the algorithm. We also explore the benefit of our method on small problems, showing that it is feasible to learn human-interpretable topics in seconds.

2 Background

Probabilistic topic models such as Latent Dirichlet allocation (LDA) [7] use latent variables to encode co-occurrence patterns between words in text corpora, and other bag-of-words data. In the LDA model, there are K topics $\phi_k, k \in \{1, \dots, K\}$, which are discrete distributions over words. For example, a topic on baseball might give high probabilities to words such as “pitcher”, “bat” and “base”. The assumed generative process for the LDA model is

```

Generate each topic  $\phi_k \sim \text{Dirichlet}(\eta), k \in \{1, \dots, K\}$ 
For each document  $j$ 
  Generate a distribution over topics  $\theta_j \sim \text{Dirichlet}(\alpha)$ 
  For each word  $i$  in document  $j$ 
    Sample a topic  $z_{ij} \sim \text{Discrete}(\theta_j)$ 
    Sample the word  $w_{ij} \sim \text{Discrete}(\phi_{z_{ij}})$  .

```

To scale LDA inference to very large datasets, a stochastic variational inference algorithm was proposed by Hoffman et al. [13]. We will discuss its more general form [12], which applies to all graphical models whose parameters can be split into “global” parameters G and “local” parameters L_j pertaining to each data point x_j , with complete conditionals being exponential family distributions. The algorithm examines one data point at a time to learn that data point’s “local” variational parameters, such as θ_j in LDA. It then updates “global” variational parameters, such as topics ϕ_k , via a stochastic natural gradient update. Their general scheme is given in Algorithm 1.

Algorithm 1 Stochastic Variational Inference (Hoffman et al.)

- Input: Data x_1, \dots, x_D , step sizes $\rho_t, t = 1 : M$ (Max iterations)
 - Randomly initialize “global” (e.g. topic) parameters G
 - For $t = 1 : M$
 - Select a random data point (e.g. document) $x_j, j \in \{1, \dots, D\}$
 - Compute “local” (e.g. document-level) variational parameters \mathbf{L}_j
 - $\hat{\mathbf{G}} = D\mathbf{L}_j$
 - $\mathbf{G} := (1 - \rho_t)\mathbf{G} + \rho_t\hat{\mathbf{G}}$
-

For an appropriate local update and sequence of step sizes ρ , this algorithm is guaranteed to converge to the optimal variational solution [12]. In the case of LDA, let λ_k be the parameter vector for a variational Dirichlet distribution on topic ϕ_k . This method computes variational distributions for topic assignments and the distribution over topics for document j using a regular VB update, then for each topic k computes $\hat{\lambda}_k$, an estimate for what λ_k would be if all D documents were identical to document j . It then updates the λ_k ’s via a natural gradient update, which takes the form

$$\lambda_k := (1 - \rho_t)\lambda_k + \rho_t\hat{\lambda}_k . \quad (1)$$

The online EM algorithm of Cappe and Moulines [9] is another general-purpose method for learning latent variable models in an online setting. The algorithm alternates between an ordinary M-step which maximizes the EM lower bound with respect to parameters θ , and a stochastic expectation step, which updates exponential family sufficient statistics \mathbf{s} with an online average

$$\mathbf{s} := (1 - \rho_t)\mathbf{s} + \rho_t\hat{\mathbf{s}}(Y_{n+1}; \theta) , \quad (2)$$

with Y_{n+1} being a new data point, θ being the current parameters, and $\hat{\mathbf{s}}(Y_{n+1}; \theta)$ being an estimate of the sufficient statistics based on these values.

In this article, we show how to perform stochastic variational inference in the collapsed representation of LDA, using an algorithm inspired by the online algorithms of Hoffman et al. and Cappe and Moulines. The new algorithm takes advantage of a fast collapsed inference method called “CVB0” [2] to further improve the efficiency of stochastic LDA inference.

2.1 CVB0

In the collapsed representation of LDA, we marginalize out topics Θ and distributions over topics Φ , and perform inference only on the topic assignments \mathbf{Z} . The collapsed variational Bayesian inference (CVB) approach of Teh et al. [24] maintains variational

discrete distributions γ_{ij} over the K topic assignment probabilities for each word i in each document j . Teh et al. showed that although the updates for a coordinate ascent algorithm optimizing the evidence lower bound with respect to γ are intractable, an algorithm using approximate updates works well in practice, outperforming the classical VB algorithm in terms of prediction performance. Asuncion et al. [2] showed that a simpler version of this method, called CVB0, is much faster while still maintaining the accuracy of CVB. The CVB0 algorithm iteratively updates each γ_{ij} via

$$\gamma_{ijk} \propto \frac{N_{w_{ij}k}^{\Phi} + \eta_{w_{ij}}}{N_k^{Z^{-ij}} + \sum_w \eta_w} (N_{jk}^{\Theta^{-ij}} + \alpha) , \quad (3)$$

for each topic k , with w_{ij} corresponding to the word token for the j th document’s i th word. The \mathbf{N}^Z , \mathbf{N}^{Θ} and \mathbf{N}^{Φ} variables, henceforth referred to as the “CVB0 statistics”, are variational expected counts corresponding to their indices, and the $-ij$ superscript indicates the exclusion of the current value of γ_{ij} . Specifically, \mathbf{N}^Z is the vector of expected number of words assigned to each topic, \mathbf{N}_j^{Θ} is the equivalent vector for document j only, and each entry w, k of matrix \mathbf{N}^{Φ} is the expected number of times word w is assigned to topic k across the corpus,

$$N_k^Z \triangleq \sum_{ij} \gamma_{ijk} \quad N_{jk}^{\Theta} \triangleq \sum_i \gamma_{ijk} \quad N_{wk}^{\Phi} \triangleq \sum_{ij:w_{ij}=w} \gamma_{ijk} . \quad (4)$$

Note that $\mathbf{N}_j^{\Theta} + \alpha$ is an unnormalized variational estimate of the posterior mean of document j ’s distribution over topics θ_j , and column k of $\mathbf{N}^{\Phi} + \beta$ is an unnormalized variational estimate of the posterior mean of topic ϕ_k . The update for CVB0 closely resembles the collapsed Gibbs update for LDA, but is deterministic.

CVB0 is currently the fastest technique for LDA inference for single-core batch inference in terms of convergence rate [2]. It is also as simple to implement as collapsed Gibbs sampling, and has a very similar update procedure. Sato and Nakagawa [22] showed that the terms in the CVB0 update can be understood as optimizing the α -divergence, with different values of α for each term. The α -divergence is a generalization of the KL-divergence that variational Bayes minimizes, and optimizing it is known as power EP [17]. A disadvantage of CVB0 is that the memory requirements are large as it needs to store a variational distribution γ for every token in the corpus. This can be improved slightly by “clumping” every occurrence of a specific word in each document together and storing a single γ for them.

3 Stochastic CVB0

We would like to exploit the efficiency and simplicity of CVB0, and the improved variational bound of the collapsed representation in a stochastic algorithm. Such an

algorithm should not need to maintain the γ variables, thus circumventing the memory requirements of CVB0, and should be able to provide an estimate for the topics when only a subset of the data have been visited. Recall that the CVB0 statistics \mathbf{N}^Z , \mathbf{N}^Θ and \mathbf{N}^Φ are all that are needed to both perform a CVB0 update and to recover estimates of the topics. So, we want to be able to estimate the CVB0 statistics based on the set of tokens we have observed.

Suppose we have seen a token w_{ij} , and its associated γ_{ij} . The information this gives us about the statistics depends on how the token was drawn. If the token was drawn uniformly at random from all of the tokens in the corpus, the expected value of \mathbf{N}^Z with respect to the sampling distribution is $C\gamma_{ij}$, where C is the number of words in the corpus. For the same sampling procedure, the expectation of the word-topic expected counts matrix \mathbf{N}^Φ is $C\mathbf{Y}^{(ij)}$, where $\mathbf{Y}^{(ij)}$ is a $W \times K$ matrix with the w_{ij} th row being γ_{ij} and with zeros in the other entries. Now if the token was drawn uniformly from the tokens in document j , the expected value of \mathbf{N}_j^Θ is $C_j\gamma_{ij}$.¹

Since we may not maintain the γ 's, we cannot perform these sampling procedures directly. However, with a current guess at the CVB0 statistics we can *update* a token's variational distribution, and observe its new value. We can then use this γ_{ij} to improve our estimate of the CVB0 statistics. This suggests an iterative procedure, alternating between a "maximization" step, approximately optimizing the evidence lower bound with respect to a particular γ_{ij} via CVB0, and an "expectation" step, where we update the expected count statistics to take into account the new γ_{ij} . As the algorithm continues, the γ_{ij} 's we observe will change, so we cannot simply average them. Instead, we can follow Cappe and Moulines [9] and perform an online average of these statistics via Equation 2.

In the proposed algorithm, we process the corpus one token at a time, examining the tokens from each document in turn. For each token, we first compute a new γ_{ij} . We do not store the γ 's, but compute (updated versions of) them as needed via CVB0. This means we must make a small additional approximation in that we cannot subtract current values of γ_{ij} in Equation 3. With large corpora and large documents this difference is negligible. The update becomes

$$\gamma_{ijk} \propto \frac{N_{w_{ijk}}^\Phi + \eta_{w_{ij}}}{N_k^Z + \sum_w \eta_w} (N_{jk}^\Theta + \alpha) . \quad (5)$$

We then use this to re-estimate our CVB0 statistics. We use one sequence of step-sizes ρ^Φ for \mathbf{N}^Φ and \mathbf{N}^Z , and another sequence ρ^Θ for \mathbf{N}^Θ . While we are processing randomly ordered tokens i of document j , we are effectively drawing random tokens from it, so the expectation of \mathbf{N}_j^Θ is $C_j\gamma_{ij}$. We update \mathbf{N}_j^Θ with an online average of

¹Other sampling schemes are possible, which would lead to different algorithms. For example, one could sample from the set of tokens with word index w to estimate \mathbf{N}_w^Φ . Our choice leads to an algorithm that is practical in the online setting.

the current value and its expected value,

$$\mathbf{N}_j^\ominus := (1 - \rho_t^\ominus)\mathbf{N}_j^\ominus + \rho_t^\ominus C_j \gamma_{ij} . \quad (6)$$

Although we process a document at a time, we eventually process all of the words in the corpus. So for the purposes of updating \mathbf{N}^Φ and \mathbf{N}^Z , in the long-run we are effectively drawing tokens from the entire corpus. The expected \mathbf{N}^Φ after observing one γ_{ij} is $C\mathbf{Y}^{(ij)}$, and the expected \mathbf{N}^Z is $C\gamma_{ij}$. In practice, it is too expensive to update the entire \mathbf{N}^Φ after every token, suggesting the use of minibatch updates. The expected \mathbf{N}^Φ after observing a minibatch M is the average of the per-token estimates, and similarly for \mathbf{N}^Z , leading to the updates:

$$\mathbf{N}^\Phi := (1 - \rho_t^\Phi)\mathbf{N}^\Phi + \rho_t^\Phi \hat{\mathbf{N}}^\Phi \quad (7)$$

$$\mathbf{N}^Z := (1 - \rho_t^Z)\mathbf{N}^Z + \rho_t^Z \hat{\mathbf{N}}^Z \quad (8)$$

where $\hat{\mathbf{N}}^\Phi = \frac{C}{|M|} \sum_{ij \in M} \mathbf{Y}^{(ij)}$ and $\hat{\mathbf{N}}^Z = \frac{C}{|M|} \sum_{ij \in M} \gamma_{ij}$. Depending on the lengths of the documents and the number of topics, it is often also beneficial to perform a small number of extra passes to learn the document statistics before updating the topic statistics. We found that one such burn-in pass was sufficient in all of the datasets we tried in our experiments. Pseudo-code for the algorithm, which we refer to as ‘‘Stochastic CVB0’’ (SCVB0) is given in Algorithm 2.

An optional optimization to the above algorithm is to only perform one update for each distinct token in each document, and scale the update by the number of copies in the document. This process, often called ‘‘clumping’’, is standard practice for fast implementations of all LDA inference algorithms, though it is only exact for uncollapsed algorithms, where the z_{ij} ’s are D-separated by θ_j . Suppose we have observed w_{tj} , which occurs m_{tj} times in document j . Plugging Equation 6 into itself m_{tj} times and noticing that all but one of the resulting terms form a geometric series, we can see that performing m_{tj} updates for \mathbf{N}_j^\ominus while holding γ_{ij} fixed is equivalent to

$$\mathbf{N}_j^\ominus := (1 - \rho_t^\ominus)^{m_{tj}} \mathbf{N}_j^\ominus + C_j \gamma_{ij} (1 - (1 - \rho_t^\ominus)^{m_{tj}}) . \quad (9)$$

4 An Alternative Perspective: MAP Estimation

In the SCVB0 algorithm, because the γ ’s are not maintained we must approximate Equation 3 with Equation 5, neglecting the subtraction of the previous value of γ_{ij} from the CVB0 statistics when updating γ_{ij} . It can be shown that this approximation results in an algorithm which is equivalent to an EM algorithm for MAP estimation, due to Asuncion et al. [2], which operates on an unnormalized parameterization of LDA. Therefore, the approximate collapsed variational updates of SCVB0 can also be

Algorithm 2 Stochastic CVB0

- Randomly initialize \mathbf{N}^Φ , \mathbf{N}^Θ ; $\mathbf{N}^Z := \sum_w \mathbf{N}_w^\Phi$
 - $\hat{\mathbf{N}}^\Phi := \mathbf{0}$; $\hat{\mathbf{N}}_\bullet := \mathbf{0}$
 - For each document j
 - For zero or more “burn-in” passes
 - * For each token i
 - Update γ_{ij} (Equation 5)
 - Update \mathbf{N}_j^Θ (Equation 6)
 - For each token i
 - * Update γ_{ij} (Equation 5)
 - * Update \mathbf{N}_j^Θ (Equation 6)
 - * $\hat{\mathbf{N}}_{w_t} := \hat{\mathbf{N}}_{w_t} + C\gamma_{ij}$
 - * $\hat{\mathbf{N}}^Z := \hat{\mathbf{N}}^Z + C\gamma_{ij}$
 - If minibatch finished
 - * Update \mathbf{N}^Φ (Equation 7)
 - * Update \mathbf{N}^Z (Equation 8)
 - * $\hat{\mathbf{N}}^\Phi := \mathbf{0}$; $\hat{\mathbf{N}}^Z := \mathbf{0}$
-

understood as MAP estimation updates. Using this interpretation of the algorithm, we now give an alternative derivation of SCVB0 as a version of Cappe and Moulines’ online EM algorithm [9] as applied to MAP estimation for LDA, thus providing an alternative perspective on the algorithm.

In particular, iterating the following update optimizes an EM lower bound on the posterior probability of the parameters:

$$\bar{\gamma}_{ijk} \propto \frac{\bar{N}_{w_{ijk}}^\Phi + \eta - 1}{\bar{N}_k^Z + W(\eta - 1)} (\bar{N}_{jk}^\Theta + \alpha - 1), \quad (10)$$

where $\bar{\gamma}_{ijk} \triangleq Pr(z_{ij} = k | \bar{N}^\Phi, \bar{N}^Z, \bar{N}^\Theta, w_{ij})$ are EM “responsibilities”, and the other variables, which we will refer to as *EM statistics*, are aggregate statistics computed from sums of these responsibilities,

$$\bar{N}_k^Z = \sum_{ij} \bar{\gamma}_{ijk} \quad \bar{N}_{jk}^\Theta = \sum_i \bar{\gamma}_{ijk} \quad \bar{N}_{wk}^\Phi = \sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk}. \quad (11)$$

Upon completion of the algorithm, MAP estimates of the parameters can be recovered by

$$\hat{\phi}_{wk} = \frac{\bar{N}_{wk}^{\Phi} + \eta - 1}{\bar{N}_k^Z + W(\eta - 1)} \quad \hat{\theta}_{jk} = \frac{\bar{N}_{jk}^{\Theta} + \alpha - 1}{C_j + K\alpha - K}, \quad (12)$$

where C_j is the length of document j . A sketch of the derivation for this algorithm, which we will refer to as *unnormalized MAP LDA* (MAP_LDA_U), is given in Appendix A. Note that if we identify the EM statistics and responsibilities with CVB0 statistics and variational distributions, the SCVB0 update in Equation 5 is identical to Equation 10 but with the hyper-parameters adjusted by one.

We now adapt online EM to this setting. In its general form, the online EM algorithm performs maximum likelihood estimation by alternating between updating an online estimate of the expected sufficient statistics for the complete-data log-likelihood and optimizing parameter estimates via a regular EM M-step. We consider this algorithm as applied to an unnormalized parameterization of LDA, where the parameters of interest are estimates \hat{N}^{Φ} , \hat{N}^{Θ} , \hat{N}^Z of the EM statistics, which are related to Θ and Φ via Equation 12. We also adapt the algorithm to perform MAP estimation, and to operate with stochasticity at the word-level instead of the document-level. The resulting algorithm is procedurally identical to SCVB0.

First, let us derive the expected sufficient statistics. Written in exponential family form, the complete data likelihood for a word w_{ij} and its topic assignment z_{ij} is

$$\begin{aligned} & \exp \left(\sum_{wk} [w_{ij} = w][z_{ij} = k] \log \left(\frac{\hat{N}_{wk}^{\Phi} + \eta - 1}{\hat{N}_k^Z + W(\eta - 1)} \right) \right. \\ & \quad \left. + \sum_k [z_{ij} = k] \log \left(\frac{\hat{N}_{jk}^{\Theta} + \alpha - 1}{N_j + K(\alpha - 1)} \right) \right) \\ & \propto \exp \left(\sum_{wk} [w_{ij} = w][z_{ij} = k] \log(\hat{N}_{wk}^{\Phi} + \eta - 1) \right. \\ & \quad \left. + \sum_k [z_{ij} = k] \log(\hat{N}_{jk}^{\Theta} + \alpha - 1) \right. \\ & \quad \left. - \sum_k [z_{ij} = k] \log(\hat{N}_k^Z + W(\eta - 1)) \right), \quad (13) \end{aligned}$$

where $[a = b]$ is a Kronecker delta function, equal to one if $a = b$ and zero otherwise, and \hat{N} variables denote current estimates, not necessarily synchronized with $\bar{\gamma}$. We can see that the sufficient statistics are the delta functions (and products of delta functions),

$$\begin{aligned} S^{(w)}(w_{ij}, z_{ij}) = & ([w_{ij} = 1][z_{ij} = 1], \dots, [w_{ij} = W][z_{ij} = K], \\ & [z_{ij} = 1], \dots, [z_{ij} = K], [z_{ij} = 1], \dots, [z_{ij} = K])^{\top}, \quad (14) \end{aligned}$$

and the expected sufficient statistics, given current parameter estimates, are appropriate entries of $\bar{\gamma}$,

$$\bar{s}^{(w)}(w_{ij}, z_{ij}) = ([w_{ij} = 1]\bar{\gamma}_{ij1}, \dots, [w_{ij} = W]\bar{\gamma}_{ijK}, \bar{\gamma}_{ij1}, \dots, \bar{\gamma}_{ijK})^\top . \quad (15)$$

Cappe and Moulines normalize the likelihood, and the sufficient statistics, by the number of data points n , so that n need not be specified in advance. However, since we are performing MAP estimation, unlike the MLE algorithm described by Cappe and Moulines, we need to estimate the unnormalized expected sufficient statistics for the entire corpus in order to maintain the correct scale relative to the prior. This can be achieved by scaling the per-word expected sufficient statistics by appropriate constants to match the size of the corpus (or document, for per-document statistics)

$$\bar{s}'^{(w)}(w_{ij}, z_{ij}) = (C[w_{ij} = 1]\bar{\gamma}_{ij1}, \dots, C[w_{ij} = W]\bar{\gamma}_{ijK}, C_j\bar{\gamma}_{ij1}, \dots, C_j\bar{\gamma}_{ijK}, C\bar{\gamma}_{ij1}, \dots, C\bar{\gamma}_{ijK})^\top . \quad (16)$$

Notice that the average of these corpus-wide expected sufficient statistics, computed across all tokens in the corpus, is equal to the EM statistics, i.e. the parameters to be optimized in the M-step. Collecting them into appropriate matrices, we can write the the expected sufficient statistics as

$$\hat{\bar{s}} = (\bar{\mathbf{N}}^\Theta, \bar{\mathbf{N}}^\Phi, \bar{\mathbf{N}}^Z) . \quad (17)$$

In fact, optimizing the EM objective function with respect to the parameters, we find that the M-step assigns the parameter estimate EM statistics to be consistent with the EM statistics computed in the E-step (CF Appendix A),

$$\hat{\bar{\mathbf{N}}}^\Theta := \bar{\mathbf{N}}^\Theta \quad \hat{\bar{\mathbf{N}}}^\Phi := \bar{\mathbf{N}}^\Phi \quad \hat{\bar{\mathbf{N}}}^Z := \bar{\mathbf{N}}^Z . \quad (18)$$

We therefore do not need to store parameter estimates $\hat{\bar{\mathbf{N}}}$ separately from expected sufficient statistics $\bar{\mathbf{N}}$, as M-step updated parameter estimates are always equal to the expected sufficient statistics from the E-step. Inserting Equation 16 into Equation 2 and using separate step size schedules for document statistics and topic statistics, the online E-step after processing token w_{ij} is given by

$$\bar{\mathbf{N}}_j^\Theta := (1 - \rho_t^\Theta)\bar{\mathbf{N}}_j^\Theta + \rho_t^\Theta C_j \gamma_{ij} \quad (19)$$

$$\bar{\mathbf{N}}_w^\Phi := (1 - \rho_t^\Phi)\bar{\mathbf{N}}^\Phi + \rho_t^\Phi C \gamma_{ij}[w_{ij} = w] , \forall w \quad (20)$$

$$\bar{\mathbf{N}}^Z := (1 - \rho_t^Z)\bar{\mathbf{N}}^Z + \rho_t^Z C \gamma_{ij} , \quad (21)$$

with $\bar{\gamma}_{ij}$ computed via Equation 10. The online EM algorithm we have just derived is procedurally identical to SCVB0 with minibatches of size one, identifying EM responsibilities and statistics with SCVB0 responsibilities and statistics, and with the hyper-parameters adjusted by one. Under this interpretation, an alternative name for SCVB0 might be *stochastic unnormalized MAP LDA* (S_MAP_LDA_U).

5 Convergence Analysis

The MAP estimation interpretation of SCVB0 is the interpretation that is most amenable to convergence analysis, since MAP_LDA_U exactly optimizes a well-defined objective function, while CVB0 has approximate updates. In this section, the notation will follow the MAP interpretation of the algorithm. We have the following theorem:

Theorem 5.1 *If $0 < \rho_t^\Phi \leq 1 \forall t$, $0 < \rho_t^\Theta \leq 1 \forall t$, $\sum_{t=1}^\infty \rho_t^\Phi = \infty$, $\lim_{t \rightarrow \infty} \rho_t^\Phi = 0$, $\sum_{t=1}^\infty \rho_t^\Theta = \infty$, and $\lim_{t \rightarrow \infty} \rho_t^\Theta = 0$, then in the limit as the number of iterations t approaches infinity SCVB0 converges to a stationary point of the MAP objective function.*

Proof Consider the MAP_LDA_U algorithm, with an update schedule alternating between a full E-step, i.e. updating every $\bar{\gamma}_{ij}$, and a full M-step, i.e. synchronizing the (parameter estimate) EM statistics with the $\bar{\gamma}$'s. The $\bar{\gamma}$'s do not depend on each other given the EM statistics, so we do not need to maintain them between iterations. We can thus view this version of MAP_LDA_U as operating on just the EM statistics. For each EM statistic $c \in \{\hat{\mathbf{N}}_1^\Theta, \dots, \hat{\mathbf{N}}_D^\Theta, \hat{\mathbf{N}}^\Phi, \hat{\mathbf{N}}^Z\}$, let $f_c(X, \hat{s}) : S_c \rightarrow S_c$ be a mapping from a current value to the updated value after such an iteration, i.e. performing an E-step to estimate the $\bar{\gamma}$'s, then using these to update the parameter estimates in the M-step, where X is the full corpus and S_c is the space of possible assignments for EM statistic c .

Let $\hat{s} = (\hat{\mathbf{N}}_1^\Theta, \dots, \hat{\mathbf{N}}_D^\Theta, \hat{\mathbf{N}}^\Phi, \hat{\mathbf{N}}^Z)$ be an assignment of the EM statistics, with \hat{s}_c referring to EM statistic c , and let $\hat{s}^{(t)}$ be the EM statistics at word iteration t of the SCVB0 algorithm. Furthermore, let $\bar{s}_c(w^{(t+1)}, \hat{s})$ be the estimate of $f_c(X, \hat{s})$ based on the word $w^{(t+1)}$ examined at step $t + 1$, as per the right hand side of the SCVB0 update equations. Note that $E[\bar{s}_c(w^{(t+1)}, \hat{s})] = f_c(X, \hat{s})$, where the expectation is with respect to the sampling of $w^{(t+1)}$. Finally, let $\xi^{(t+1)} = \bar{s}_c(w^{(t+1)}, \hat{s}^{(t)}) - f_c(X, \hat{s}^{(t)})$ be the stochastic error made at step $t + 1$, and observe that $E[\xi^{(t+1)}] = 0$. We can rewrite the SCVB0 updates for each EM statistic c as

$$\begin{aligned}
\hat{s}_c^{(t+1)} &= (1 - \rho_{t+1}^c) \hat{s}_c^{(t)} + \rho_{t+1}^c \bar{s}_c(w^{(t+1)}, \hat{s}) \\
&= \hat{s}_c^{(t)} + \rho_{t+1}^c (-\hat{s}_c^{(t)} + \bar{s}_c(w^{(t+1)}, \hat{s})) \\
&= \hat{s}_c^{(t)} + \rho_{t+1}^c (f_c(X, \hat{s}^{(t)}) - \hat{s}_c^{(t)} + \bar{s}_c(w^{(t+1)}, \hat{s}) - f_c(X, \hat{s}^{(t)})) \\
&= \hat{s}_c^{(t)} + \rho_{t+1}^c (f_c(X, \hat{s}^{(t)}) - \hat{s}_c^{(t)} + \xi^{(t+1)}) .
\end{aligned} \tag{22}$$

In this form, we can see that iterating each of the SCVB0 updates corresponds to a Robbins-Monro stochastic approximation (SA) algorithm [21] for finding the zeros of $f_c(X, \hat{s}^{(t)}) - \hat{s}_c^{(t)}$, i.e. the fixed points of MAP_LDA_U for \hat{s}_c . Since MAP_LDA_U is an EM algorithm, its fixed points are the stationary points of the posterior probability of the parameters, as recovered via Equation 12.

Theorem 2.3 of Andreiu et al. [1] states that under mild conditions, the existence of a Lyapunov function, along with a boundedness condition, implies that such a Robins-Monro algorithm will converge with step size schedules such as those above. In the context of an SA algorithm, a Lyapunov function can be understood as an “objective function” which, in the absence of stochastic noise, the SA would improve monotonically if small enough steps were taken in the direction of the updates. In Appendix B, we show that the negative of the Lagrangian of the EM lower bound is a Lyapunov function of the overall SCVB0 algorithm and the set of fixed points of the EM algorithm. The boundedness condition, namely that the state variables stay within a compact subset of the state space, follows by observing that $0 < \|\bar{s}_c(x^{(t+1)}, \hat{s})\|_1 \leq C$ for every EM statistic c , so if the initial state also satisfies this, by convexity \hat{s} will always have its L1 norm similarly bounded. Having demonstrated that the assumptions required by Theorem 2.3 of Andreiu et al. hold, the convergence result follows.

6 Experiments

In this section we describe an experimental analysis of the proposed SCVB0 algorithm with comparison to the stochastic variational Bayes algorithm of Hoffman et al., hereafter referred to as SVB. As well as performing an analysis on several large-scale problems, we also investigate the effectiveness of the stochastic LDA inference algorithms at learning topics in near real-time on small corpora.

6.1 Large-Scale Experiments

We studied the performance of the algorithms on three large corpora. The corpora are:

- *PubMed Central*: A corpus of full-text scientific articles from the open-access PubMed Central database of scientific literature in the biomedical and life sciences. After processing to remove rare words and stopwords, the corpus contained approximated 320M tokens across 165,000 articles, with a vocabulary size of around 38,500 words.
- *New York Times*: A corpus containing 1.8 million articles from the New York Times, published between 1987 and 2007. After processing, the corpus had a dictionary of about 50,000 words and contained 475M distinct tokens.

- *Wikipedia*: This collection contains 4.6 million articles from the online encyclopedia Wikipedia. We used the dictionary of 7,700 words extracted by Hoffman et al. for their experiments on an earlier extracted Wikipedia corpus. There were 811M tokens in the corpus.

We explored predictive performance versus wall-clock time between SCVB0 and SVB. To compare the algorithms fairly, we implemented both of them in the fast high-level language Julia [5]. Our implementation of SVB closely follows the python implementation provided by Hoffman, and has several optimizations not mentioned in the original paper including handling the latent topic assignments z implicitly, “clumping” of like tokens and sparse updates of the topic matrix. Our algorithm was implemented as it is written in Algorithm 2, using the clumping optimization but with no additional algorithmic optimizations. Specifically, neither implementation used the complicated optimizations taking advantage of sparsity that are exploited by the Vowpal Wabbit implementation of SVB² and in the variant of SVB proposed by Mimno [16], but instead represent a “best-effort” attempt to implement each algorithm efficiently yet following the spirit of the original pseudo-code.

In all experiments, each algorithm was trained using minibatches of size 100. We used a step-size schedule of $\frac{s}{(\tau+t)^\kappa}$ for document iteration t , with $s = 10$, $\tau = 1000$ and $\kappa = 0.9$. For SCVB0, the document parameters were updated using the same schedule with $s = 1$, $\tau = 10$ and $\kappa = 0.9$. We used LDA hyper-parameters $\alpha = 0.1$ and $\eta = 0.01$ for SCVB0. For SVB, we tried both these same hyperparameter values as well as shifting by 0.5 as recommended by [2] to compensate for the implicit bias in how uncollapsed VB treats hyper-parameters. We used a single pass to learn document parameters for SCVB0, and tried both a single pass and five passes for SVB.

For each experiment we held out 10,000 documents and trained on the remaining documents. We split each test document in half, estimated document parameters on one half and computed the log-probability of the remaining half of the document. Figures 1 through 3 show held-out log-likelihood versus wall-clock time for each algorithm.

For the PubMed Central data, we found that all algorithms perform similarly after about an hour, but prior to that SCVB0 is better, indicating that SCVB0 makes better use of its time. All algorithms perform similarly per-iteration (see Figure 4), but SCVB0 is able to benefit by processing more documents in the same amount of time. The per-iteration plots for the other datasets were similar.

Our experiments show that SCVB0 shows a more substantial benefit when employed on larger datasets. In both the New York Times and Wikipedia experiments SCVB0 converged to a better solution than SVB for any of its parameter settings. Furthermore, SCVB0 outperforms SVB throughout the run.

²https://github.com/JohnLangford/vowpal_wabbit/wiki

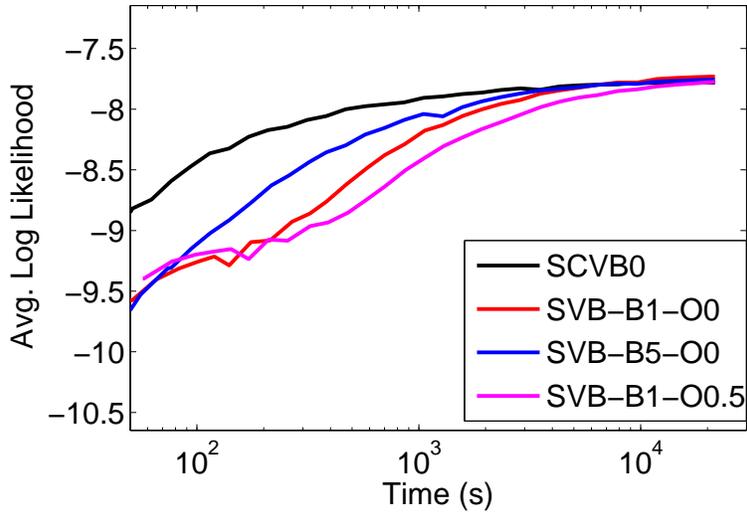


Figure 1: Log-likelihood vs Time for the PubMed Central experiments. SVB-B x -O y corresponds to running SVB with x burn-in passes and with hyper-parameters offset from $\alpha = 0.1$ and $\eta = 0.01$ by y .

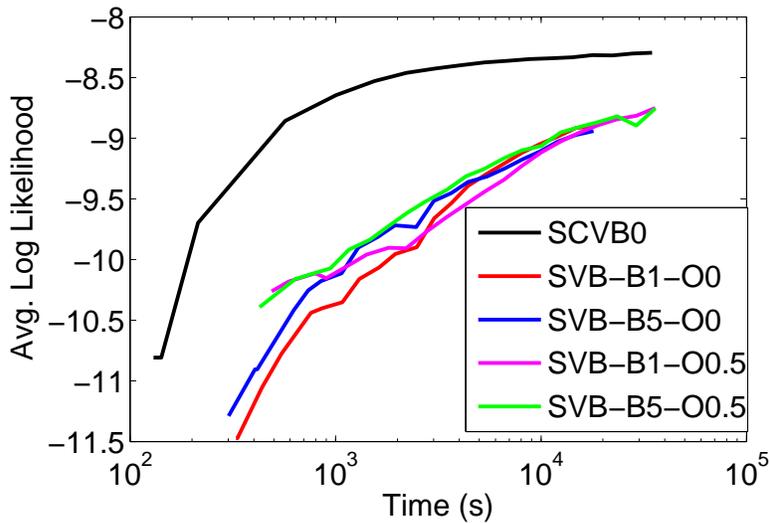


Figure 2: Log-likelihood vs Time for the New York Times experiments. SVB-B x -O y corresponds to running SVB with x burn-in passes and with hyper-parameters offset from $\alpha = 0.1$ and $\eta = 0.01$ by y .

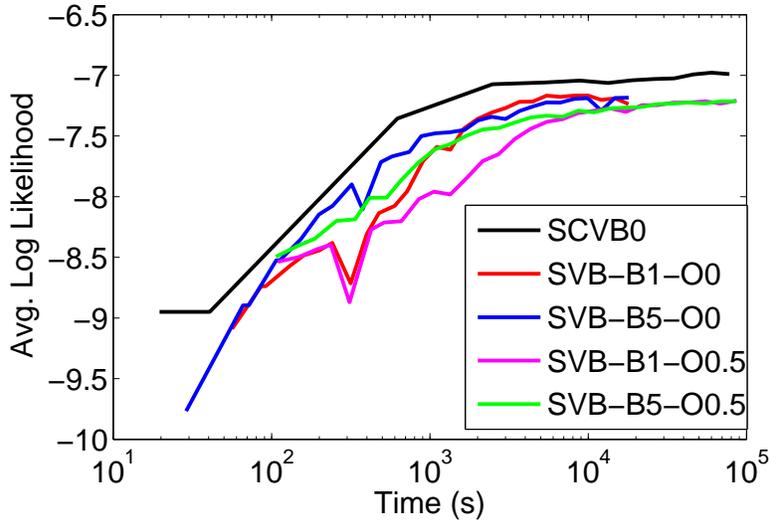


Figure 3: Log-likelihood vs Time for the Wikipedia experiments. SVB-B x -O y corresponds to running SVB with x burn-in passes and with hyper-parameters offset from $\alpha = 0.1$ and $\eta = 0.01$ by y .

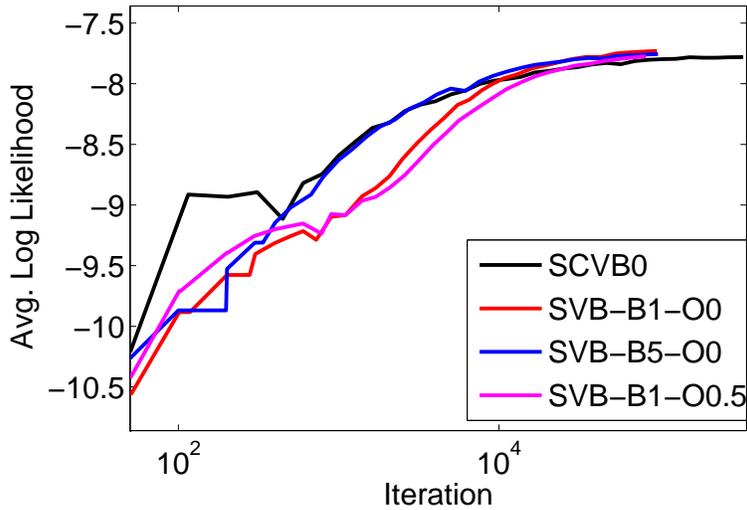


Figure 4: Log-likelihood vs Iteration for the PubMed Central experiments. SVB-B x -O y corresponds to running SVB with x burn-in passes and with hyper-parameters offset from $\alpha = 0.1$ and $\eta = 0.01$ by y .

6.2 Small-Scale Experiments

Stochastic algorithms for LDA have previously only been used on large corpora, however they have the potential to be useful for finding topics very quickly on small corpora as well. The ability to learn interpretable topics in a matter of seconds is very beneficial for exploratory data analysis (EDA) applications, where a human is in the loop. Near real-time topic modeling opens the way for the use of topic models in interactive software tools for document analysis.

We investigated the performance of the stochastic algorithms in the small-scale scenario using a corpus of 1740 scientific articles from the machine learning conference NIPS, between 1987 and 1999. We ran the two stochastic inference algorithms for 5 seconds each, using the parameter settings from the previous experiments but with 20 topics. Each algorithm was performed ten times. In the five seconds of training, SCVB0 was typically able to examine 3300 documents, while SVB was typically able to examine around 600 documents.

With the EDA application in mind, we performed a human-subject experiment in the vein of the experiments proposed by Chang and Blei [8]. The sets of topics returned by each run were randomly assigned across seven human subjects. The participants were all machine learning researchers with technical expertise in the subjects of interest to the NIPS community. The subjects did not know which algorithms generated which runs. The top ten words of the topics in each run were shown to the subjects, who were given the following instructions:

Here are 20 collections of related words. Some words may not seem to “belong” with the other words. Count the total number of words in each collection that don’t “belong”.

This task finds the number of “errors” that a topic model inference algorithm makes, relative to human judgement. It was found that the SCVB0 algorithm had 0.76 errors per topic on average, with a standard deviation of 1.1, while SVB had 1.6 errors per topic on average, with standard deviation 1.2. A one-sided two sample t-test rejected the hypothesis that the means of the errors per topic were equal, with significance level $\alpha = 0.05$. Example topics are shown in Table 1.

We also performed a similar experiment on Amazon Turk involving 52 people using the New York Times corpus. We ran the two stochastic inference algorithms for 60 seconds each using the same parameter settings as above but with 50 topics. Each user was presented with 20 random topics from each algorithm. Example topics are shown in Table 2. Again, the subjects did not know which algorithms generated each set of topics. We included two easy questions with obvious answers and removed results from users who did not answer them correctly. Comparing the number of “errors” for SCVB0 to SVB for each user, we find that SCVB0 had significantly fewer errors for the sampled population at the $\alpha = .05$ level using a paired t-test, with p-value $< .001$.

SCVB0			SVB		
receptor	data	learning	model	results	visual
protein	classification	function	set	learning	data
secondary	vector	network	data	distribution	activity
proteins	class	neural	training	information	saliency
transducer	classifier	networks	learning	map	noise
binding	set	time	error	activity	similarity
concentration	algorithm	order	parameters	time	model
odor	feature	error	markov	figure	neural
morphology	space	dynamics	estimate	networks	representations
junction	vectors	point	speech	state	functions

Table 1: Randomly selected example topics after five seconds running time on the NIPS corpus.

7 Discussion / Related Work

Connections can be drawn between SCVB0 and other methods in the literature. The SCVB0 scheme is reminiscent of the online EM algorithm of Cappe and Moulines [9], which also alternates between per data-point parameter updates and online estimates of the expected values of sufficient statistics. Online EM optimizes the EM lower bound on the log-likelihood in the M-step and computes online averages of exponential family sufficient statistics, while SCVB0 (approximately) updates the mean-field evidence lower bound in the M-step and computes online averages of sufficient statistics required for a CVB0 update in the E-step.

The SCVB0 algorithm also has a very similar structure to SVB, alternating between passes through a document (the optional “burn-in” passes) to learn document parameters, and updating variables associated with topics. However, SCVB0 is stochastic at the word-level while SVB is stochastic at the document level. In the general framework of Hoffman et al., inference is performed on “local” parameters specific to a data point, which are used to perform a stochastic update on the “global” parameters. For SVB, the document parameters Θ_j are local parameters for document j , and topics are global parameters. For SCVB0, the γ_{ij} ’s are local parameters for a word, and both document parameters N^Θ and topic parameters N^Φ are global parameters. This means that updates to document parameters can be made before processing all of the words in the document.

The incremental algorithm of Banerjee and Basu [4] for MAP inference in LDA is also closely related to the proposed algorithm. They estimate topic probabilities for each word in sequence, and update MAP estimates of Φ and Θ incrementally, using the expected assignments of words to topics in the current document. SCVB0 can

SCVB			SVB		
county	station	league	president	year	mr
district	company	goals	midshipmen	cantatas	company
village	railway	years	open	edward	mep
north	business	club	forrester	computing	husbands
river	services	clubs	archives	main	net
area	market	season	iraq	years	state
east	line	played	left	area	builder
town	industry	cup	back	withdraw	offense
lake	stations	career	times	households	obscure
west	owned	team	saving	brain	advocacy

Table 2: Randomly selected example topics after sixty seconds running time on the NYT corpus.

be understood as the collapsed, stochastic variational version of Banerjee and Basu’s incremental uncollapsed MAP estimation algorithm. Interpreting SCVB0 as a MAP estimation algorithm, SCVB0 is the online EM algorithm for MAP estimation operating on the unnormalized representation of LDA, while Banerjee and Basu’s algorithm is the incremental EM algorithm operating on the usual normalized representation of LDA.

Another stochastic algorithm for LDA, due to Mimno et al. [16], operates in a partially collapsed space, placing it in-between SVB and SCVB0 in terms of representation. Their algorithm collapses out Θ but does not collapse out Φ . Estimates of online natural gradient update directions are computed by performing Gibbs sampling on the topic assignments of the words in each document, and averaging over the samples. The gradient estimate is non-zero only for word-topic pairs which occurred in the samples. When implemented in a sparse way, the updates scale sub-linearly in the number of topics, causing large improvements in high-dimensional regimes, however these updates are less useful in when the number of topics is smaller (around 100 or less). The performance gains of this algorithm depend on a careful implementation that takes advantage of the sparsity. For SCVB0, the minibatch updates are sparse in the rows (words), so some performance enhancements along the lines of those used by Mimno et al. are likely to be possible.

There has been a substantial amount of work on speeding up LDA inference in the literature. Porteous et al. [20] improved the efficiency of the sampling step for the collapsed Gibbs sampler, and [25] explore a number of alternatives for improving the efficiency of LDA. The Vowpal Wabbit system for fast machine learning, due to John Langford and collaborators, has a version of SVB that has been engineered to be extremely efficient. Parallelization is another approach for improving the efficiency of topic models. Newman et al. [19] introduced an approximate parallel algorithm

for LDA where data is distributed across multiple machines, and an exact algorithm for an extension of LDA which takes into account the distributed storage. Smola and Narayanamurthy developed an efficient architecture for parallel LDA inference [23], using a distributed (key, value) storage for synchronizing the state of the sampler between machines.

8 Conclusions

We have introduced SCVB0, an algorithm for performing fast stochastic collapsed variational inference in LDA, and shown that it outperforms stochastic VB on several large document corpora, converging faster and often to a better solution. The algorithm is relatively simple to implement, with intuitive update rules consisting only of basic arithmetic operations. We also found that the algorithm was effective at learning good topics from small corpora in seconds, finding topics that were superior than those of stochastic VB according to human judgement.

There are many directions for future work. The speed of the method could potentially be improved by exploiting sparsity, using techniques such as those employed by Mimno et al. [16]. Furthermore, the collapsed representation facilitates the use of the parallelization techniques explored by Newman et al. in [19]. Finally, SCVB0 could be incorporated into an interactive software tool for exploring the topics of document corpora in real-time.

9 Acknowledgments

We would like to thank Arthur Asuncion for many helpful discussions.

References

- [1] C. Andrieu, É. Moulines, and P. Priouret. Stability of stochastic approximation under verifiable conditions. *SIAM Journal on control and optimization*, 44(1):283–312, 2005.
- [2] A. Asuncion, M. Welling, P. Smyth, and Y. Teh. On smoothing and inference for topic models. In *Uncertainty in Artificial Intelligence*, 2009.
- [3] D. C. Atkins, T. N. Rubin, M. Steyvers, M. A. Doeden, B. R. Baucom, and A. Christensen. Topic models: A novel method for modeling couple and family text data. *Journal of Family Psychology*, 6:816–827, 2012.

- [4] A. Banerjee and S. Basu. Topic models over text streams: A study of batch and online unsupervised learning. In *SIAM Data Mining*, 2007.
- [5] J. Bezanson, S. Karpinski, V. B. Shah, and A. Edelman. Julia: A fast dynamic language for technical computing. *CoRR*, abs/1209.5145, 2012.
- [6] C. M. Bishop et al. *Pattern recognition and machine learning*, volume 1. springer New York, 2006.
- [7] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation. *The Journal of Machine Learning Research*, 3:993–1022, 2003.
- [8] J. Boyd-Graber, J. Chang, S. Gerrish, C. Wang, and D. Blei. Reading tea leaves: How humans interpret topic models. In *Proceedings of the 23rd Annual Conference on Neural Information Processing Systems*, 2009.
- [9] O. Cappé and E. Moulines. On-line expectation–maximization algorithm for latent data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):593–613, 2009.
- [10] B. Carpenter. Integrating out multinomial parameters in latent Dirichlet allocation and naive bayes for collapsed Gibbs sampling. Technical report, LingPipe, 2010.
- [11] T. Griffiths and M. Steyvers. Finding scientific topics. *Proceedings of the National Academy of Sciences of the United States of America*, 101(Suppl 1):5228, 2004.
- [12] M. Hoffman, D. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *arXiv preprint arXiv:1206.7051*, 2012.
- [13] M. D. Hoffman, D. M. Blei, and F. Bach. Online learning for latent Dirichlet allocation. *Advances in Neural Information Processing Systems*, 23:856–864, 2010.
- [14] D. Mimno. Computational historiography: Data mining in a century of classics journals. *Journal on Computing and Cultural Heritage (JOCCH)*, 5(1):3, 2012.
- [15] D. Mimno. Reconstructing pompeian households. *Uncertainty in Artificial Intelligence*, 2012.
- [16] D. Mimno, M. Hoffman, and D. Blei. Sparse stochastic inference for latent Dirichlet allocation. In *Proceedings of the International Conference on Machine Learning*, 2012.
- [17] T. Minka. Power EP. Technical report, Microsoft Research, Cambridge, UK, 2004.

- [18] R. M. Neal and G. E. Hinton. A view of the em algorithm that justifies incremental, sparse, and other variants. In M. Jordan, editor, *Learning in graphical models*, pages 355–368. Springer, 1998.
- [19] D. Newman, A. Asuncion, P. Smyth, and M. Welling. Distributed algorithms for topic models. *The Journal of Machine Learning Research*, 10:1801–1828, 2009.
- [20] I. Porteous, D. Newman, A. Ihler, A. Asuncion, P. Smyth, and M. Welling. Fast collapsed gibbs sampling for latent dirichlet allocation. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 569–577, 2008.
- [21] H. Robbins and S. Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.
- [22] I. Sato and H. Nakagawa. Rethinking collapsed variational Bayes inference for LDA. *Proceedings of the International Conference on Machine Learning*, 2012.
- [23] A. Smola and S. Narayanamurthy. An architecture for parallel topic models. *Proceedings of the VLDB Endowment*, 3(1-2):703–710, 2010.
- [24] Y. Teh, D. Newman, and M. Welling. A collapsed variational bayesian inference algorithm for latent Dirichlet allocation. *Advances in Neural Information Processing Systems*, 19:1353, 2007.
- [25] L. Yao, D. Mimno, and A. McCallum. Efficient methods for topic model inference on streaming document collections. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 937–946. ACM, 2009.

A Derivation of the Unnormalized MAP Algorithm

Here, we give a more complete derivation of the MAP_LDA_U algorithm than is provided in Asuncion et al [2], and show that using a certain ordering of the EM updates results in an algorithm which is very similar to CVB0.

MAP estimation aims to maximize the log posterior probability of the parameters,

$$\begin{aligned}
 \log Pr(\Theta, \Phi|w, \eta, \alpha) &= \sum_j \log Pr(w_j|\Theta_j, \Phi) \\
 &+ \sum_{jk} (\alpha - 1) \log(\theta_{jk}) + \sum_{wk} (\eta - 1) \log(\phi_{wk}) + \text{const.} \quad (23)
 \end{aligned}$$

This objective function cannot readily be optimized directly via, e.g., a gradient update, since the log-likelihood term and its gradient require an intractable sum over z inside the logarithm. Instead, EM may be performed. A standard Jensen’s inequality argument gives the EM objective function as described by Neal and Hinton [18], which, when applied to the MAP estimation problem, is a lower bound $\mathcal{L}(\Theta, \Phi, \bar{\gamma})$ on the posterior probability (CF Bishop [6]),

$$\log Pr(\Theta, \Phi|X) \geq \mathcal{L}(\Theta, \Phi, \bar{\gamma}) \triangleq R(\Theta, \Phi, \bar{\gamma}) - \sum_{ijk} \bar{\gamma}_{ijk} \log \bar{\gamma}_{ijk} , \quad (24)$$

where

$$\begin{aligned} R(\Theta, \Phi; \Theta^{(t)}, \Phi^{(t)}) &= \sum_{wk} \left(\sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1 \right) \log \phi_{wk} \\ &+ \sum_{jk} \left(\sum_i \bar{\gamma}_{ijk} + \alpha - 1 \right) \log \theta_{jk} + \text{const} \end{aligned} \quad (25)$$

is the expected complete data log-likelihood, plus terms arising from the prior, and the $\bar{\gamma}_{ijk}$ ’s are E-step “responsibilities”,

$$\bar{\gamma}_{ijk} \triangleq Pr(z_{ij} = k | \Theta, \Phi, w_{ij}) \propto Pr(w_{ij} | z_{ij} = k, \Theta, \Phi) Pr(z_{ij} = k | \Theta, \Phi) = \phi_{w_{ij}k} \theta_{jk} . \quad (26)$$

Adding Lagrange terms $-\sum_k \lambda_k^\Phi (\sum_w \phi_{wk} - 1)$ and $-\sum_j \lambda_j^\Theta (\sum_k \theta_{jk} - 1)$, taking derivatives and setting to zero, we obtain the following M-step updates:

$$\phi_{wk} : \propto \sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1 \qquad \theta_{jk} : \propto \sum_i \bar{\gamma}_{ijk} + \alpha - 1 . \quad (27)$$

It is possible to reparameterize the EM algorithm for LDA in terms of unnormalized counts of the EM “responsibilities” instead of Θ and Φ [2], which we refer to as the *EM statistics*. Their definitions are given in Equation 11.

Substituting these values into the M-step updates above, then substituting the optimal (M-step updated) parameter assignments into the EM bound and rearranging, we obtain a reparameterization of the EM bound

$$\begin{aligned}
\log Pr(\Theta, \Phi|X) \geq & \sum_{wk} \left(\sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1 \right) \log(\hat{N}_{wk}^{\Phi} + \eta - 1) \\
& + \sum_{jk} \left(\sum_i \bar{\gamma}_{ijk} + \alpha - 1 \right) \log(\hat{N}_{jk}^{\Theta} + \alpha - 1) \\
& - \sum_k \left(\sum_{ij} \bar{\gamma}_{ijk} + W(\eta - 1) \right) \log(\hat{N}_k^Z + W(\eta - 1)) \\
& - \sum_{ijk} \bar{\gamma}_{ijk} \log \bar{\gamma}_{ijk} + \text{const}
\end{aligned} \tag{28}$$

where $\hat{\mathbf{N}}^{\Phi}$, $\hat{\mathbf{N}}^{\Theta}$ and $\hat{\mathbf{N}}^Z$ are current estimates of the EM statistics, not necessarily synchronized with the $\bar{\gamma}$'s. To derive M-step updates for this reparameterized formulation, we first add Lagrangian terms to enforce the constraints that each of the EM statistics sums to the number of words in the corpus C , $-\lambda_{\Phi}(\sum_{wk} \hat{N}_{wk}^{\Phi} - C)$, $-\lambda_{\Theta}(\sum_{jk} \hat{N}_{jk}^{\Theta} - C)$, $\lambda_Z(\sum_k \hat{N}_k^Z - C)$. In the following, we derive the update for \hat{N}_{wk}^{Φ} ; the derivation is similar for the other parameters. We take derivatives with respect to each parameter and set them to zero, and plug the constraint equations back into the resulting equations:

$$\begin{aligned}
\frac{\sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1}{\hat{N}_{wk}^{\Phi} + \eta - 1} - \lambda_{\Phi} &= 0 \tag{29} \\
\sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1 &= \lambda_{\Phi}(\hat{N}_{wk}^{\Phi} + \eta - 1) \\
\hat{N}_{wk}^{\Phi} &= \frac{\sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1}{\lambda_{\Phi}} - (\eta - 1) \\
C = \sum_{wk} \hat{N}_{wk}^{\Phi} &= \frac{1}{\lambda_{\Phi}} \sum_{wk} \left(\sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1 \right) - KW(\eta - 1).
\end{aligned}$$

Solving for the Lagrange multipliers, they turn out to be one:

$$\lambda_{\Phi} = \frac{\sum_{wk} \sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + KW(\eta - 1)}{C + KW(\eta - 1)} = \frac{C + KW(\eta - 1)}{C + KW(\eta - 1)} = 1.$$

Plugging this back into Equation 29 (in the case of \hat{N}_{wk}^{Φ}), we obtain M-step updates which synchronize the EM statistics with their definitions in Equation 11 (i.e. the update in Equation 18). Note that after the M-step, $\sum_w \hat{N}_{wk}^{\Phi} = \hat{N}_k^Z \forall k$, and we did not need to enforce this explicitly in the algorithm.

The E-step finds the expected value of the complete-data log-likelihood, as encoded by the responsibilities $\bar{\gamma}_{ij}$. Plugging in the estimates of Θ and Φ from Equation 12 into Equation 26 gives us the update in Equation 10. Alternatively, adding Lagrange terms $\sum_{ij} \lambda_{ij} (\sum_k \bar{\gamma}_{ijk} - 1)$ to the bound to enforce the constraint that the $\bar{\gamma}$'s sum to one, setting the derivatives to zero then solving for $\bar{\gamma}_{ij}$ also gives us Equation 10.

The standard EM algorithm alternates between complete E and M-steps, i.e. updating all of the $\bar{\gamma}_{ij}$'s, followed by synchronizing the EM statistics with the responsibilities. However, the EM algorithm can be viewed as a coordinate ascent algorithm on the lower bound objective function, and partial E and M-steps also improve this bound [18]. In our case, both updating a single $\bar{\gamma}_{ij}$, and subsequently synchronizing the EM statistics to reflect the new value (partial E and M-steps, respectively) are coordinate ascent updates which improve the EM lower bound in Equation 28. So an algorithm that iteratively performs the update in Equation 10 for each token (a partial E-step), while continuously keeping the EM statistics in synch with the $\bar{\gamma}_{ij}$'s as in Equation 11 (a partial M-step), is equivalent to the above EM algorithm but merely performing the coordinate ascent updates in a different order. This algorithm is very similar to CVB0, but using Equation 10 instead of Equation 3.

B Lyapunov Function

A Lyapunov function is a function which gives a stochastic analogue of the monotonicity property of the EM algorithm, the existence of which is a standard argument for the stability and convergence of a stochastic approximation algorithm. Theorem 2.3 of Andreiu et al. [1] states that under mild conditions, convergence is assured for a Robins-Monro SA algorithm endowed with a Lyapunov function with certain properties. Andreiu et al. consider an SA with state space Θ for finding $h(\theta) = \mathbf{0}$, where Θ is an open subset of \mathbb{R}^n , and $h : \Theta \rightarrow \mathbb{R}^n$. They require the existence of a continuously differentiable function $w : \Theta \rightarrow [0, \infty)$, the Lyapunov function, such that the following conditions hold:

- (i) There exists $M_0 > 0$ such that

$$\mathcal{L} \triangleq \{\theta \in \Theta, \langle \nabla w(\theta), h(\theta) \rangle = 0\} \subset \{\theta \in \Theta, w(\theta) < M_0\}$$

- (ii) There exists $M_1 \in (M_0, \infty]$ such that $\{\theta \in \Theta, w(\theta) < M_1\}$ is a compact set,
- (iii) For any $\theta \in \Theta$ \mathcal{L} , $\langle \nabla w(\theta), h(\theta) \rangle < 0$
- (iv) $w(\mathcal{L})$ has an empty interior.

In our case, recall that in Section 5 we showed that the SCVB0 updates for each of the EM statistics c corresponds to a Robbins-Monro SA for finding the zeros of $f_c(X, \hat{s}^{(t)}) - \hat{s}_c^{(t)}$, i.e. the fixed points of MAP_LDA_U for \hat{s}_c . In the overall algorithm, $\theta = (\hat{N}^\Theta, \hat{N}^\Phi, \hat{N}^Z)$, and $h(\theta)$ is the direction of the M-step update that we would take if we were to first perform a full E-step. Finding $h(\theta) = \mathbf{0}$, as the SA algorithm is designed to do, corresponds to finding the fixed points of the MAP_LDA_U EM algorithm, which are at the stationary points of the posterior distribution of the parameters, i.e. the objective function for MAP estimation.

We will now show that minus the EM lower bound, augmented with Lagrange terms, is a Lyapunov function of the overall algorithm. As we found in Appendix A, if we include Lagrange constraints in the EM bound to ensure that the EM statistics sum to C , set the gradient to zero and solve for the Lagrange multipliers, the Lagrange multipliers turn out to equal one. Substituting this value into the Lagrangian and dropping constant terms, we have our candidate function

$$\begin{aligned}
-w(\hat{N}^\Theta, \hat{N}^\Phi, \hat{N}^Z) \triangleq & \sum_{wk} \left[\left(\sum_{ij:w_{ij}=w} \bar{\gamma}_{ijk} + \eta - 1 \right) \log(\hat{N}_{wk}^\Phi + \eta - 1) - \hat{N}_{wk}^\Phi \right] \\
& + \sum_{jk} \left[\left(\sum_i \bar{\gamma}_{ijk} + \alpha - 1 \right) \log(\hat{N}_{jk}^\Theta + \alpha - 1) - \hat{N}_{jk}^\Theta \right] \\
& - \sum_k \left[\left(\sum_{ij} \bar{\gamma}_{ijk} + W(\eta - 1) \right) \log(\hat{N}_k^Z + W(\eta - 1)) - \hat{N}_k^Z \right] \\
& - \sum_{ijk} \bar{\gamma}_{ijk} \log \bar{\gamma}_{ijk} , \tag{30}
\end{aligned}$$

where $\bar{\gamma}$ are E-step estimates computed from the current EM statistics – note that $w(\theta)$ is not a function of them. We want to show that conditions (i) – through (iv) hold for $w(\theta)$.

Condition (ii) holds because the EM statistics have L1 norm bounded by C . Condition (iv) holds by Sard’s theorem. The key conditions are (i) and (iii), which involve the directional derivative of $w(\theta)$ at θ along $h(\theta)$, $\langle \nabla w(\theta), h(\theta) \rangle$ (where we have appended the EM statistics so that θ is a vector). This is the instantaneous change in $w(\theta)$ in the direction of the EM update. Note that a step with a step-size multiplier of one in the direction $h(\theta)$ is guaranteed by the monotonicity of EM to improve the (Lagrangian of the) lower bound, and thereby lower $w(\theta)$. However, we have to check that an infinitesimal step in that direction also improves this function.

Suppose we are not at a fixed point of EM, i.e. $h(\theta) \neq \mathbf{0}$. Fixing $\bar{\gamma}$ to E-step-updated values based on θ , we know from the derivation of the M-step update that the Lagrangian of the EM lower bound has a unique maximum at the M-step-updated value, located at $\theta + h(\theta)$. Since this maximum is unique and there are no other stationary points, each point in the direction $h(\theta)$ of the maximum has an increasingly

large value of the Lagrangian of the EM bound, holding $\bar{\gamma}$ fixed. These values computed with $\bar{\gamma}$ fixed to its current value are a lower bound on the Lagrangian $-w(\theta)$ at those points, as $w(\theta)$ is computed using E-step updated $\bar{\gamma}$'s which must strictly improve the EM lower bound relative to the current (or any other) $\bar{\gamma}$. So every point on the line segment between θ and $\theta + h(\theta)$ has a strictly higher value of the Lagrangian $-w(\theta)$ than at θ , i.e. $-w(\theta + \lambda h(\theta)) - (-w(\theta)) > 0, \forall \lambda \in (0, 1]$. This implies that $\langle \nabla w(\theta), h(\theta) \rangle = \lim_{\lambda \rightarrow 0} \frac{w(\theta + \lambda h(\theta)) - w(\theta)}{\lambda} < 0$, and (iii) holds.³

At a fixed point of MAP_LDA_U, which (due to the properties of EM) happens IFF the algorithm is at a stationary point of the MAP objective function, $h(\theta) = \mathbf{0}$, and it can be shown by inspection that $\nabla w(\theta) = \mathbf{0}$ only under these conditions also. In this case, the directional derivative $\langle \nabla w(\theta), h(\theta) \rangle = 0$, and $\theta \in \mathcal{L}$. So $\theta \in \mathcal{L}$ IFF θ is at a stationary point of the MAP objective function, and (i) holds. Along with a boundedness condition demonstrated in Section 5, Theorem 2.3 of Andreiu et al. [1] now gives us that with an appropriate sequence of step sizes, in the limit as the number of iterations approaches infinity the distance from \mathcal{L} is zero.

³The directional derivative of w at θ along v is defined to be $\lim_{\lambda \rightarrow 0} \frac{w(\theta + \lambda v) - w(\theta)}{\lambda}$. If w is differentiable at θ , the directional derivative equals $\langle \nabla w(\theta), v \rangle$.