

A Flexible Stochastic Method for Solving the MAP Problem in Topic Models

Xuan Bui^{1,2}, Tu Vu¹, Khoat Than¹, and Ryutaro Ichise³

¹ Hanoi University of Science and Technology, Hanoi, Vietnam

² University of Information and Communication Technology, Thai Nguyen, Vietnam

³ National Institute of Informatics, Tokyo, Japan

{vutu201130, thanhxuan1581}@gmail.com

khoattq@soict.hust.edu.vn, ichise@nii.ac.jp

Abstract. The estimation of the posterior distribution is the core problem in topic models, unfortunately it is intractable. There are approximation methods and sampling methods proposed to solve it. However, most of them do not have any clear theoretical guarantee of neither quality nor rate of convergence. Online Maximum a Posteriori Estimation (OPE) is another approach with concise guarantee on quality and convergence rate, in which we cast the estimation of the posterior distribution into convex/non-convex optimization problem. In this paper, we propose a more general and flexible version of OPE, namely Generalized Online Maximum a Posteriori Estimation (G-OPE), which not only enhances the flexibility of OPE in different real-world datasets but also preserves key advantage theoretical characteristics of OPE when comparing to the state-of-the-art methods. We employ G-OPE as inference method for a topic proportion of a document within large text corpora. The experimental and theoretical results show that our new approach performs better than OPE and other state-of-the-art methods.

Key words: Topic models, posterior inference, Online MAP estimation, large-scale learning, non-convex optimization

1 Introduction

Topic models are widely used in text processing and Latent Dirichlet Allocation (LDA)[3] is the core of a large family of probabilistic models. LDA provides an efficient tool to analyze hidden themes in data and helps us recover hidden structures/evolution in big text collections. The key problem in topic models is to compute the posterior distribution of a document given other parameters. The posterior inference problem in topic models is to infer the topic proportion of documents and topics which are distributions over vocabulary. Large datasets or streaming environment contain huge number of documents, hence the problem of estimating topic proportion for an individual document is especially important. The quality of learning method for LDA is determined by the quality of the inference method being employed. Unfortunately, the posterior distribution of a

document is intractable [3]. There are two main approaches to tackle it. One is approximating the intractable distribution by tractable distribution, for example Variational Bayes inference (VB) [3]. The other is a sampling method, which draws numerous of samples from target distribution then estimating the interesting quality from these samples. The well-known method is Collapsed Gibbs Sampling (CGS) [7]. There are also famous methods such as Collapsed Variational Bayes (CVB) [12, 13], CVB0 [2], Stochastic Variational Inference (SVI) [9], etc. To our best knowledge, there are not any mathematical guarantees for quality and convergence rate in existing approaches. Therefore, in practice we do not have any ideas about how to stop the methods we are using but trying, observing and retrying again to reach the best solution.

Another way to solve the posterior distribution is to view it as an optimization problem. To infer about topic proportion of a document is to solve the maximum a posterior of topic proportion given words in this document and all topics of corpus [14]. This optimization problem is usually non-convex and NP-hard in practice [11]. There is very few theoretical contributions in non-convex optimization literature, especially in topic models. Online Maximum a Posteriori Estimation (OPE) [14] which is an online version of Frank-Wolfe algorithm [8] is a stochastic algorithms to solve such kind of non-convex problem. OPE is theoretically guaranteed to converge to a local stationary point at a rate of $\mathcal{O}(1/T)$ where T was the number of iterations [14]. Although OPE is easy to implement, fast convergence and is mathematically guaranteed, it remains some problems. The weakness of OPE is that it is not well adaptive with different datasets because of the uniform distribution in its operation. We will exploit this crucial point to propose a new and more general algorithm based on OPE. When changing its operations, we have to retain the advantage of the original algorithms, that is theoretical guarantees.

Our main contribution is following:

- We propose new algorithm Generalized Online Maximum a Posteriori Estimation (G-OPE) for solving posterior inference problem in topic models based on OPE. G-OPE is more general and flexible than OPE, adapts better in different datasets and preserves the key advantages of original method.
- We employed G-OPE into the existing algorithm Online-OPE [14] to learn LDA in online settings and streaming environment.
- We conduct experiments to demonstrate that Online-GOPE outperforms novel methods to learn LDA.

The rest of this paper is organized as follows. In Section 2, we introduce an overview of posterior inference with LDA and main ideas of existing methods. In Section 3, our new algorithm G-OPE is proposed in details. In Section 4, we conduct experiments with two large datasets with state-of-the-art methods in two different measures. Finally Section 5 is our conclusion.

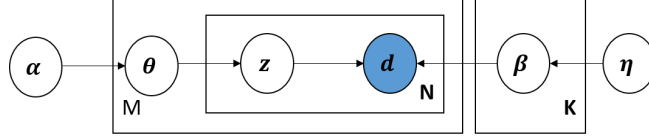


Fig. 1: Latent Dirichlet Allocation

2 Related work

Latent Dirichlet Allocation (LDA) [3] is the basic and famous model in topic modeling. LDA models a collection of text documents as topics. It represents each document as a probability distribution θ_d over topics, and each topic β_k as a probability distribution over words. In Fig.1, K is number of topics, M is number of documents in corpus, N is number of words in each documents. Note that $\theta_d \in \Delta_K$, $\beta_k \in \Delta_V$ ⁴. The generative process for each document d as follows

- draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$
- For the n^{th} word of d :
 - draw topic index $z_{dn} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - draw word $w_{dn} | z_{dn}, \beta \sim \text{Multinomial}(\beta_{z_{dn}})$

The most important problem we need to solve in order to use LDA is to compute the posterior distribution of hidden variable in a given document $p(\theta, z | w, \alpha, \beta)$. However, it is intractable. There are many ways to handle it. Variational Bayesian Inference [3] approximates $p(z_d, \theta_d, d | \beta, \alpha)$ by obtaining a lower bound on the likelihood which is adjustable by variational distributions. CVB and CVB0 deal with $p(z_d, d | \beta, \alpha)$, CGS draws samples from $p(z_d, w | \beta, \alpha)$ to estimate it. Eventually, all methods try to estimate the topic proportion θ_d . In this paper, we infer topic proportion for a document directly by solving the Maximum a Posteriori Estimation (MAP) of θ_d given all words of this documents and model's parameters.

the MAP estimation of topic mixture for a given document d :

$$\theta^* = \arg \max_{\theta \in \bar{\Delta}_K} \Pr(d, \theta | \beta, \alpha) = \arg \max_{\theta \in \bar{\Delta}_K} \Pr(d | \theta, \beta) \Pr(\theta | \alpha) \quad (1)$$

Under the assumption about the generative process, problem (1) is equivalent to the following:

$$\theta^* = \arg \max_{\theta \in \bar{\Delta}_K} \sum_j d_j \log \sum_{k=1}^K \theta_k \beta_{kj} + (\alpha - 1) \sum_{k=1}^K \log \theta_k \quad (2)$$

Within convex/concave optimization, problem (2) is relatively well-studied. In the case of $\alpha \geq 1$, it can easily be shown that the problem (2) is concave,

⁴ $\Delta_N = \{x \in R^N | x_n \geq 0, \sum_n x_n = 1\}$

Algorithm 1 OPE: Online Maximum a Posteriori Estimation**Input:** document \mathbf{d} and model $\{\beta, \alpha\}$ **Output:** θ that maximizes $f(\theta) = \sum_j d_j \log \sum_{k=1}^K \theta_k \beta_{kj} + (\alpha - 1) \sum_{k=1}^K \log \theta_k$ Initialize θ_1 arbitrary in $\bar{\Delta}_K$ **for** $t = 1, 2, \dots, T$ **do**Pick f_t uniformly from $\{\sum_j d_j \log \sum_{k=1}^K \theta_j \beta_{kj}; (\alpha - 1) \sum_{k=1}^K \log \theta_k\}$ $F_t := \frac{2}{t} \sum_{h=1}^t f_h$ $e_t := \arg \max_{\mathbf{x} \in \Delta_K} \langle F'_t(\theta_t), \mathbf{x} \rangle$ $\theta_{t+1} := \theta_t + \frac{e_t - \theta_t}{t}$ **end for**

and therefore it can be solved in polynomial time. Unfortunately, in practice of LDA, the parameter α is often small, says $\alpha < 1$, causing problem (2) to be non-concave. Sontag and Roy showed that problem (2) is NP-hard in the worst case when parameter $\alpha < 1$ [11]. Consider problem (2) as a non-convex optimization problem, the gradient-based method such as Gradient Descent (GD) is ineffective because of the existence of saddle point, hence we need an effective random methods to avoid it.

OPE [14] is an iterative algorithm for problem (2). In the literature of iterative optimization algorithms, in each iteration, they try to build a tractable function that approximates true objective function, then optimize approximating function to reach the next point. The various algorithms have different techniques to build their own approximation. For example, using Jensen's inequality, Expectation-Maximization (EM) [5] or Variational Inference (VI) [3] calculate the Evidence Lower Bound (ELBO) then maximize it. Gradient Descent constructs its quadratic approximation in each step and minimizes the quadratic. OPE solves the problem (2) by constructing an approximate sequence by stochastic way and solve it by Frank-Wolfe update formula [6].

Details of OPE is in Algorithm 1. At each iteration t , it draws a sample function $f_t(\theta)$ and builds the approximate function $F_t(\theta)$ which is the average of all previous sample function. The most interesting idea behind OPE is that the objective function is the sum of a likelihood and a prior. In each step, it builds an approximate function $F_t(\theta)$ by choosing either likelihood or prior with equal probabilities $\{0.5, 0.5\}$. That means when inferring about the topic proportion of a document, we use either the evidence of the document (likelihood) or knowledge we have known before (prior). This behavior is very natural to human. However, OPE considers likelihood and prior with the same contributions by using uniform distribution. In fact, when humans deal with a new sample, one can rely on more likelihood if we have observed enough evidences, or rely on more prior knowledge if we have been lack of evidences. This simple idea leads us to build a more general and flexible version of OPE by using Bernoulli distribution instead of uniform distribution.

3 Generalized Online Maximum a Posteriori Estimation

In this section, we introduce our new algorithm, namely Generalized Online Maximum a Posteriori Estimation (G-OPE) based on OPE. OPE operates by choosing the likelihood or prior at each step t , then builds the approximation $F_t(\boldsymbol{\theta})$ which is the average of all parts draw from previous steps and current step. In G-OPE, in order to introduce the Bernoulli distribution into the sampling step, we need to modify the likelihood and prior so that the approximation function $F_t(\boldsymbol{\theta}) \rightarrow f(\boldsymbol{\theta})$ as $t \rightarrow \infty$.

Denote

$$g_1(\boldsymbol{\theta}) = \sum_j d_j \log \sum_{k=1}^K \theta_k \beta_{kj} ; g_2(\boldsymbol{\theta}) = (\alpha - 1) \sum_{k=1}^K \log \theta_k$$

then $f(\boldsymbol{\theta}) = g_1(\boldsymbol{\theta}) + g_2(\boldsymbol{\theta})$ with $g_1(\boldsymbol{\theta})$ and $g_2(\boldsymbol{\theta})$ are the likelihood and prior respectively.

Denote

$G_1(\boldsymbol{\theta}) := g_1(\boldsymbol{\theta})/p$; $G_2(\boldsymbol{\theta}) := g_2(\boldsymbol{\theta})/(1-p)$. $G_1(\boldsymbol{\theta})$ and $G_2(\boldsymbol{\theta})$ are the adjusted likelihood and prior respectively.

G-OPE is detailed in Algorithm 2. In Algorithm 2, $f(\boldsymbol{\theta})$ is the true objective function we need to maximize. At step t , $f_t(\boldsymbol{\theta})$ is the sample function we draw from set of adjusted likelihood and prior, $F_t(\boldsymbol{\theta})$ is the approximate function we build. T is number of iterations for whole algorithm. Because G-OPE is stochastic, in theory we consider $T \rightarrow \infty$.

We use Bernoulli distribution with parameter p to replace for uniform distribution in OPE. For each iteration ($t = 1, 2, \dots, T$), we pick $f_t(\boldsymbol{\theta})$ as Bernoulli random variable with probability p from $\{G_1(\boldsymbol{\theta}), G_2(\boldsymbol{\theta})\}$. In statistic theory, as t increases (at least 20) and it is better to choose p not close to 0 or 1.

$$\Pr(f_t(\boldsymbol{\theta}) = G_1(\boldsymbol{\theta})) = p ; \Pr(f_t(\boldsymbol{\theta}) = G_2(\boldsymbol{\theta})) = 1 - p$$

Consider t independent Bernoulli trials with $\{\Pr(f_h = G_1) = p ; \Pr(f_h = G_2) = 1 - p\} \forall h = 1..t$, we build a stochastic approximate sequence

$$F_t := \frac{1}{t} \sum_{h=1}^t f_h, \forall t = 1, 2, \dots, T$$

$F_t(\boldsymbol{\theta})$ is the average of all sample functions drawn until current step. So it is guaranteed to converge to $f(\boldsymbol{\theta})$ as $t \rightarrow \infty$, which will be shown in Theorem 1. The parameter p controls how much likelihood part and prior part contribute to the objective function. We can utilize this point to choose the most suitable p in each circumstance. OPE is a special case of G-OPE when Bernoulli parameter $p = 0.5$. So OPE is not flexible in many datasets. G-OPE adapts well with different datasets, we will show it in the experiment section. In the rest of this section, we will show that G-OPE preserves the key advantage of OPE which is the guarantee of the quality and convergence rate. This character is unknown for the existing methods in posterior estimation in topic models.

Algorithm 2 G-OPE: Generalized Online maximum a Posteriori Estimation**Input:** document \mathbf{d} and model $\{\beta, \alpha\}$, parameter $p \in (0, 1)$ **Output:** θ that maximizes $f(\theta) = \sum_j d_j \log \sum_{k=1}^K \theta_k \beta_{kj} + (\alpha - 1) \sum_{k=1}^K \log \theta_k$ Initialize θ_1 arbitrary in Δ_K **for** $t = 1, 2, \dots, T$ **do**Pick f_t as Bernoulli distribution from $\{G_1(\theta), G_2(\theta)\}$ where $\{\Pr(f_t(\theta) = G_1(\theta)) = p ; \Pr(f_t(\theta) = G_2(\theta)) = 1 - p\}$ $F_t(\theta) := \frac{1}{t} \sum_{h=1}^t f_h(\theta)$ $e_t := \arg \max_{\mathbf{x} \in \Delta_K} \langle F_t'(\theta_t), \mathbf{x} \rangle$ $\theta_{t+1} := \theta_t + \frac{e_t - \theta_t}{t}$ **end for**

Theorem 1 (Convergence of G-OPE algorithm). *Consider the objective function $f(\theta)$ in Eq.2, given fixed $\mathbf{d}, \beta, \alpha, p$. For G-OPE, with probability 1, the followings hold:*

1. For any $\theta \in \Delta_K$, $F_t(\theta)$ converges to $f(\theta)$ as $t \rightarrow +\infty$.
2. θ_t converges to a local maximal/stationary point of $f(\theta)$ at a rate of $\mathcal{O}(1/t)$.

Proof:

The proof of the rate of convergence is similar to results in [14]. We claim that G-OPE also converges to a local maximum/stationary point. Before the proof, we remind some notations: $B(n, p)$ is the Bernoulli distribution, $N(\mu, \sigma^2)$ is normal distribution. $E(X)$ and $D(X)$ are expectation and variance of random variable X respectively.

The objective function $f(\theta)$ is a non-convex. The criterion used for the convergence analysis is importance in non-convex optimization. For unconstrained problems, the gradient norm $\|\nabla f(\theta)\|$ is typically used to measure convergence, because $\|\nabla f(\theta)\| \rightarrow 0$ captures convergence to a stationary point. However, this criterion can not be used for constrained problems. Instead, we use the "Frank-Wolfe gap" criterion [10].

Denoted

$$g_1(\theta) = \sum_j d_j \log \sum_{k=1}^K \theta_k \beta_{kj} ; \quad g_2(\theta) = (\alpha - 1) \sum_{k=1}^K \log \theta_k$$

and

$$G_1(\theta) := g_1(\theta)/p ; \quad G_2(\theta) := g_2(\theta)/(1 - p)$$

$$f(\theta) = g_1(\theta) + g_2(\theta) = p \cdot G_1(\theta) + (1 - p) G_2(\theta)$$

Pick $f_t(\theta)$ as Bernoulli distribution from $\{G_1(\theta), G_2(\theta)\}$ where

$$\Pr(f_t(\theta) = G_1(\theta)) = p ; \quad \Pr(f_t(\theta) = G_2(\theta)) = 1 - p$$

Let a_t and $t - a_t$ be the number of times that we have already picked $G_1(\theta)$ and $G_2(\theta)$ respectively after t iterations. We have $a_t \sim B(t, p)$ and $E(a_t) =$

$t.p$; $D(a_t) = t.p.(1-p)$. Then $S_t = a_t - t.p \rightarrow N(0, t.p(1-p))$ when $t \rightarrow \infty$. So $S_t/t \rightarrow 0$ as $t \rightarrow \infty$ with probability 1. We have

$$\begin{aligned} F_t &= \frac{1}{t}(a_t G_1 + (t - a_t) G_2) \\ F_t - f &= \frac{S_t}{t}(G_1 - G_2) \\ F'_t - f' &= \frac{S_t}{t}(G'_1 - G'_2) \end{aligned} \quad (3)$$

From Eq.3, we conclude that the $F_t \rightarrow f$ as $t \rightarrow +\infty$ with probability 1. Due to the proof in [14], $\theta_t \rightarrow \theta^*$ with the rate of $\mathcal{O}(1/t)$.

Consider

$$\begin{aligned} \langle F'_t(\theta_t), \frac{e_t - \theta_t}{t} \rangle &= \langle F'_t(\theta_t) - f'(\theta_t), \frac{e_t - \theta_t}{t} \rangle + \langle f'(\theta_t), \frac{e_t - \theta_t}{t} \rangle \\ &= \frac{S_t}{t^2} \langle G'_1(\theta_t) - G'_2(\theta_t), e_t - \theta_t \rangle + \langle f'(\theta_t), \frac{e_t - \theta_t}{t} \rangle \end{aligned}$$

Note that $g_1(\theta), g_2(\theta)$ are Lipschitz continuous on $\bar{\Delta}_K$, so is $-f(\theta)$. Hence there exists a constant L such that $\langle f'(z), y - z \rangle \leq f(y) - f(z) + L\|y - z\|^2 \forall y, z \in \bar{\Delta}_K$.

$$\begin{aligned} \langle f'(\theta_t), \frac{e_t - \theta_t}{t} \rangle &= \langle f'(\theta_t), \theta_{t+1} - \theta_t \rangle \leq f(\theta_{t+1}) - f(\theta_t) + L\|\theta_{t+1} - \theta_t\|^2 \\ &= f(\theta_{t+1}) - f(\theta_t) + \frac{L}{t^2} \|e_t - \theta_t\|^2 \end{aligned}$$

Since e_t and θ_t belong to $\bar{\Delta}_K$, $|\langle G'_1(\theta_t) - G'_2(\theta_t), e_t - \theta_t \rangle|$ and $\|e_t - \theta_t\|^2$ are bounded above for any t .

Therefore, there exists a constant $c_1 > 0$ such that

$$\langle F'_t(\theta_t), \frac{e_t - \theta_t}{t} \rangle \leq c_1 \frac{|S_t|}{t^2} + f(\theta_{t+1}) - f(\theta_t) + \frac{c_1 L}{t^2} \quad (4)$$

Summing both sides of Eq.4 for all t we have

$$\sum_{h=1}^t \frac{1}{h} \langle F'_h(\theta_h), e_h - \theta_h \rangle \leq \sum_{h=1}^t c_1 \frac{|S_h|}{h^2} + f(\theta_{t+1}) - f(\theta_1) + \sum_{h=1}^t \frac{c_1 L}{h^2} \quad (5)$$

As $t \rightarrow +\infty$, $f(\theta_t) \rightarrow f(\theta^*)$ due to the continuity of $f(\theta)$. As a result, Eq.5 implies

$$\sum_{h=1}^{+\infty} \frac{1}{h} \langle F'_h(\theta_h), e_h - \theta_h \rangle \leq \sum_{h=1}^{+\infty} c_1 \frac{|S_h|}{h^2} + f(\theta^*) - f(\theta_1) + \sum_{h=1}^{+\infty} \frac{c_1 L}{h^2} \quad (6)$$

Applying proof in [14] and based on the law of large numbers, we have $\sum_{h=1}^{\infty} c_1 \frac{|S_h|}{h^2}$ converges in probability. Moreover, the term $\sum_{h=1}^{+\infty} \frac{1}{h^2}$ is bounded above.

So $\sum_{h=1}^{+\infty} \frac{1}{h} \langle F'_h(\boldsymbol{\theta}_h), \mathbf{e}_h - \boldsymbol{\theta}_h \rangle$ is bounded above.

Because $\mathbf{e}_t = \arg \max_{\mathbf{x} \in \Delta_K} \langle F'_t(\boldsymbol{\theta}_t), \mathbf{x} \rangle$, so $\langle F'_t(\boldsymbol{\theta}_t), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle \geq 0$.

If exists $t_0 > 0, c_3 > 0$ such as $\langle F'_t(\boldsymbol{\theta}_t), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle \geq c_3 \forall t > t_0$ then $\sum_{t=1}^{\infty} \frac{1}{t} \langle F'_t(\boldsymbol{\theta}_t), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle > \sum_{t=1}^{\infty} \frac{c_3}{t}$. And because $\sum_{t=1}^{\infty} \frac{1}{t}$ is not bounded above, so $\sum_{h=1}^{+\infty} \frac{1}{h} \langle F'_h(\boldsymbol{\theta}_h), \mathbf{e}_h - \boldsymbol{\theta}_h \rangle \rightarrow \infty$, which contradicts with the clause we claimed.

Therefore, $\langle F'_t(\boldsymbol{\theta}_t), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle \rightarrow 0$ as $t \rightarrow \infty$.

$$\begin{aligned} \langle F'_t(\boldsymbol{\theta}_t), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle &= \langle f'(\boldsymbol{\theta}_t) + \frac{S_t}{t} (G'_1(\boldsymbol{\theta}_t) - G'_2(\boldsymbol{\theta}_t)), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle \\ &= \langle f'(\boldsymbol{\theta}_t), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle + \langle \frac{S_t}{t} (G'_1(\boldsymbol{\theta}_t) - G'_2(\boldsymbol{\theta}_t)), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle \end{aligned}$$

Because $\frac{S_t}{t} \rightarrow 0$ then $\langle f'(\boldsymbol{\theta}_t), \mathbf{e}_t - \boldsymbol{\theta}_t \rangle \rightarrow 0$. Apply Frank-Wolfe gap criterion, $\boldsymbol{\theta}^*$ is stationary/local maximum of f , which completes the proof ■.

Besides, in the non-convex optimization field, the idea of how to build the approximate function in G-OPE can be utilized in the case of objective function f which is the sum of two parts $f = g + h$. In each step, choose g or h in Bernoulli distribution with parameter p , and adjust p to adapt with different circumstance. Randomness can help algorithms jump out of local minimum/maximum. Therefore, to design new stochastic algorithms, we begin with a deterministic version, add a sequence of approximation in the G-OPE style, working with each approximation at each iteration by deterministic update formula. This is an open idea for our future works.

4 Experiments

In this section, we will investigate the performance of G-OPE in real world datasets. G-OPE can play as the core inference step when learning LDA, we will investigate the performance of G-OPE through the performance of Online-OPE [14] when changing its core inference method. So we derived Online-GOPE. We conducted two experiments. The first one is the effect of parameter p in G-OPE when learning LDA and the second is in comparison Online-GOPE with the current state-of-the-art methods.

4.1 Datasets and Settings

The datasets for our investigation are New York Times and Pubmed⁵ These are very large datasets. The number of documents is large and the size of vocabulary is large also. Details of datasets are presented in Table 1.

To evaluate the performance of learning methods in LDA, we used *Log Predictive Probability (LPP)* and *Normalized Pointwise Mutual Information (NPMI)*

⁵ The datasets were taken from <http://archive.ics.uci.edu/ml/>

Table 1: Two data sets for our experiments

Data sets	No.Documents	No.Terms	No.Train	No.Test
New York Times	300000	141444	290000	10000
Pubmed	330000	100000	320000	10000

measures. These measures is commonly used in topic models. *Predictive Probability* [9] measures the predictiveness and generalization of a model to new data, while NPMI [1, 4] evaluates semantics quality of an individual topic in these models.

Some common parameters is set follow: the number of topics is $K = 100$, the hyper-parameters in LDA model is $\alpha = \frac{1}{K} = 0.01, \eta = \frac{1}{K} = 0.01$. For each inference method, the number of iterations is $T = 50$. We compare the online learning algorithms together and the mini-batch size is $S = |C_t| = 5000$. For the other state-of-the-art methods, the forgetting rate $\kappa = 0.9$, we fixed $\tau = 1$. These chosen parameters is best for online learning LDA in many previous works.

As algorithms we compares are stochastic, so to avoid randomness, we run each method five times, and report the average results.

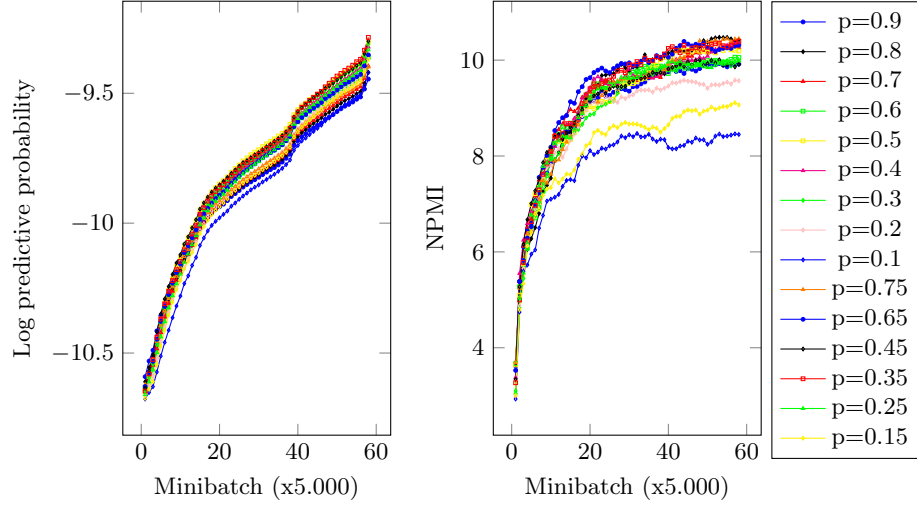
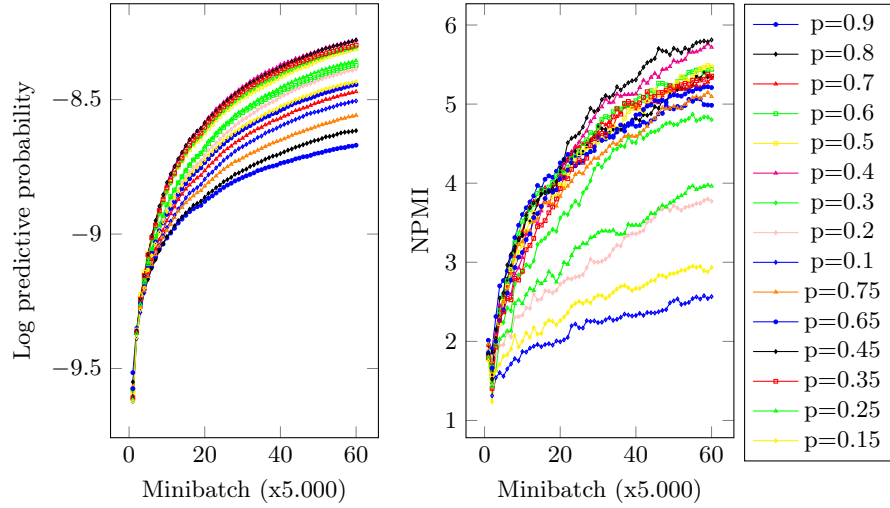
The script of experiments is that: for the first experiment, we run Online-GOPE with different values of parameter p then choose the best one. In the second experiment, we compare Online-GOPE obtained with the best parameter p to some methods in learning LDA such as VB, CVB, CGS, OPE.

4.2 The effect of parameter p

In this experiment, we investigate how important the value of parameter p is. Because $p \in (0, 1)$, and p is good if it is not close to 0 and 1. So we choose p respectively in $\{0.1, 0.15, \dots, 0.9\}$, then run Online-GOPE in two datasets. We report the performance of Online-GOPE in Fig.2 and Fig.3. We can easily observe that p affects very much in the performance in terms of both measures. In Fig.2, Online-GOPE reaches the best performance on New York Times for LPP measure at $p = 0.35$ and for NPMI measure at $p = 0.75$. In Fig.3, Online-GOPE reaches the best performance on Pubmed for LPP measure at $p = 0.4$, for NPMI measure at $p = 0.45$.

This results support our idea about the contributions of likelihood part and prior part of topic proportion inference for a document. The different dataset has the suitable value of p . If we want to get the best performance on the generalization or on semantics quality of topics, we have different p to choose. Therefore G-OPE is very flexible in the real world dataset.

The good values of p depend on how much likelihood part and prior part possess in total. The likelihood depends on the length of the documents. In our datasets, the average length of a document in New York Times is 329 while the average length of a document in Pubmed is 65. That explains why we have different best values of p for each dataset.

Fig. 2: Online-GOPE with different values of p on New York TimesFig. 3: Online-GOPE with different values of p on Pubmed

4.3 Comparison of G-OPE with novel algorithms

In this experiment, we compare Online-GOPE with the best value of p in previous experiment to the original Online-OPE and other methods: Online-VB, Online-CVB, Online-CGS. All of these algorithms try to learn the topics over the words β or variational parameters λ . The difference among these algorithms is the inner inference procedures.

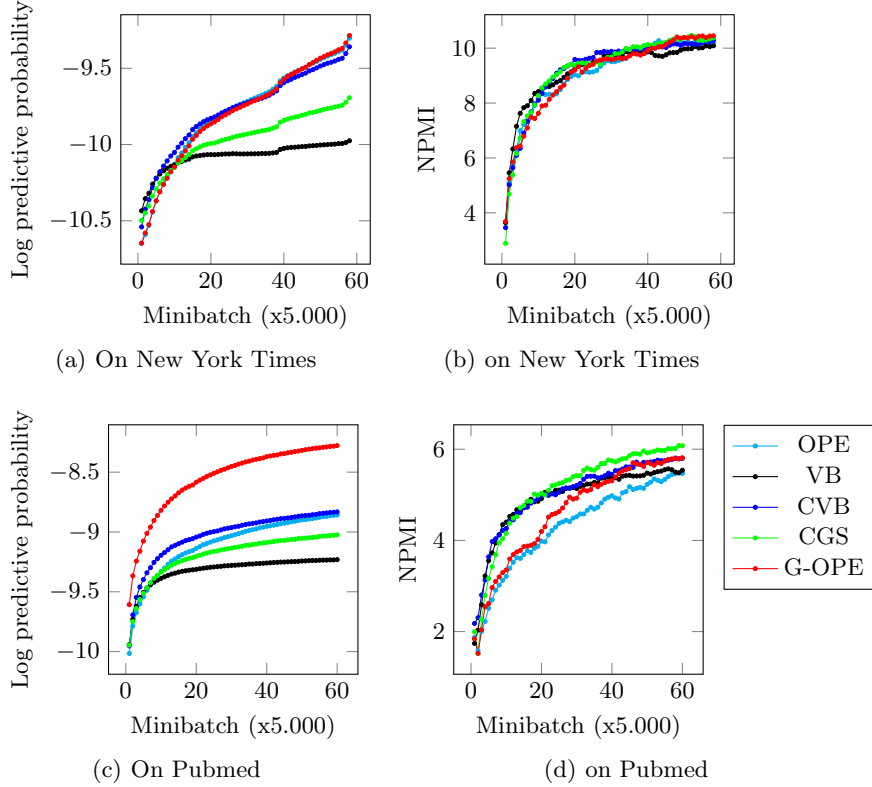


Fig. 4: Online-GOPE compares with Online-OPE, Online-VB, Online-CVB and Online-CGS. Higher is better.

The results is shown in Fig.4. With suitable parameter p , we obtained G-OPE which was better than OPE, VB, CVB, and CGS on LLP measure. For NPMI measure, all algorithms perform the same, but G-OPE is one of the tops. This results show that Online-GOPE performs better than not only original OPE, but also the current novel methods. G-OPE works well because of the right choose of controlled parameter p .

5 Conclusion

We have discussed how posterior inference for individual texts in topic models can be done efficiently with our method. In theory, G-OPE remains the guarantee on quality and convergence rate of original OPE algorithm which is the most important character among existing state-of-the-art inference methods. In practice, the parameter p of Bernoulli distribution in our method is a flexible way to deal with different datasets. Besides, the spiritual idea in building approxima-

tion functions from G-OPE can be easily extended to a wide class of maximum a posteriori estimation or non-convex problems. By exploiting G-OPE carefully, we have derived an efficient method Online-GOPE for learning LDA from data streams or large corpora. As a result, it is the good candidate to help us to work with text streams and big data.

Acknowledgement

This research is funded by Thai Nguyen University of Information and Communication Technology (ICTU) under grant number T2018-07-01.

References

1. N. Aletras and M. Stevenson. Evaluating topic coherence using distributional semantics. In *Proceedings of the 10th International Conference on Computational Semantics*, pages 13–22, 2013.
2. A. Asuncion, M. Welling, P. Smyth, and Y. W. Teh. On smoothing and inference for topic models. In *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence*, pages 27–34. AUAI Press, 2009.
3. D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
4. G. Bouma. Normalized (pointwise) mutual information in collocation extraction. *German Society for Computational Linguistics Language Technology*, pages 31–40, 2009.
5. A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B*, pages 1–38, 1977.
6. M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics*, 3(1-2):95–110, 1956.
7. T. L. Griffiths and M. Steyvers. Finding scientific topics. In *Proceedings of the National academy of Sciences*, volume 101, pages 5228–5235. National Acad Sciences, 2004.
8. E. Hazan and S. Kale. Projection-free online learning. In *Proceedings of Annual International Conference on Machine Learning*, 2012.
9. M. D. Hoffman, D. M. Blei, C. Wang, and J. W. Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
10. S. J. Reddi, S. Sra, B. Póczos, and A. J. Smola. Stochastic frank-wolfe methods for nonconvex optimization. In *Proceedings of 54th Annual Allerton Conference on Communication, Control, and Computing*, pages 1244–1251. IEEE, 2016.
11. D. Sontag and D. Roy. Complexity of inference in latent dirichlet allocation. In *Proceedings of Advances in Neural Information Processing System*, 2011.
12. Y. W. Teh, K. Kurihara, and M. Welling. Collapsed variational inference for hdp. In *Proceedings of Advances in Neural Information Processing Systems*, pages 1481–1488, 2007.
13. Y. W. Teh, D. Newman, and M. Welling. A collapsed variational bayesian inference algorithm for latent dirichlet allocation. In *Proceedings of Advances in Neural Information Processing Systems*, pages 1353–1360, 2006.
14. K. Than and T. Doan. Guaranteed algorithms for inference in topic models. *arXiv preprint arXiv:1512.03308*, 2015.