

# Distributed Gibbs Sampling of Latent Dirichlet Allocation: The Gritty Details

Yi Wang  
yiwang@tencent.com

THIS IS AN EARLY DRAFT. YOUR FEEDBACKS ARE HIGHLY APPRECIATED.

## 1 Introduction

I started to write this chapter since November 2007, right after my first MapReduce implementation of the AD-LDA algorithm[?]. I had worked so hard to understand LDA, but cannot find any document that were comprehensive, complete, and contain all necessary details. It is true that I can find some open source implementations of LDA, for example, LDA Gibbs sampler in Java<sup>1</sup> and GibbsLDA++<sup>2</sup>, but code does not reveal math derivations. I read Griffiths' paper [?] and technical report [?], but realized most derivations are skipped. It was lucky to me that in the paper about Topic-over-time[?], an extension of LDA, the authors show part of the derivation. The derivation is helpful to understand LDA but is too short and not self-contained. A primer [?] by Gregor Heinrich contains more details than documents mentioned above. Indeed, most of my initial understanding on LDA comes from [?]. As [?] focuses more on the Bayesian background of latent variable modeling, I wrote this article focusing on more on the derivation of the Gibbs sampling algorithm for LDA. You may notice that the Gibbs sampling updating rule we derived in (38) differs slightly from what was presented in [?] and [?], but matches that in [?].

## 2 LDA and Its Learning Problem

In the literature, the training documents of LDA are often denoted by a long and segmented vector  $W$ :

$$W = \begin{bmatrix} \{w_1, \dots, w_{N_1}\}, & \text{words in the 1st document} \\ \{w_{N_1+1}, \dots, w_{N_1+N_2}\}, & \text{words in the 2nd document} \\ \dots & \dots \\ \{w_{1+\sum_{j=1}^{D-1} N_j}, \dots, w_{\sum_{j=1}^D N_j}\} & \text{words in the } D\text{-th document} \end{bmatrix}, \quad (1)$$

where  $N_d$  denotes the number of words in the  $d$ -th document.<sup>3</sup>

In rest of this article, we consider hidden variables

$$Z = \begin{bmatrix} \{z_1, \dots, z_{N_1}\}, & \text{topics of words in the 1st document} \\ \{z_{N_1+1}, \dots, z_{N_1+N_2}\}, & \text{topics of words in the 2nd document} \\ \dots & \dots \\ \{z_{1+\sum_{j=1}^{D-1} N_j}, \dots, z_{\sum_{j=1}^D N_j}\} & \text{topics of words in the } D\text{-th document} \end{bmatrix}, \quad (3)$$

where each  $z_i \in Z$  corresponds to a word  $w_i \in W$  in (1).

<sup>1</sup><http://arbylon.net/projects/LdaGibbsSampler.java>

<sup>2</sup><http://gibbslda++.sourceforge.net>

<sup>3</sup>Another representation of the training documents are two vectors:  $\mathbf{d}$  and  $W$ :

$$\begin{bmatrix} \mathbf{d} \\ W \end{bmatrix} = \begin{bmatrix} d_1, & \dots, & d_N \\ w_1, & \dots, & w_N \end{bmatrix}, \quad (2)$$

where  $N = \sum_{j=1}^D N_j$ ,  $d_i$  indices in  $\mathcal{D}$ ,  $w_i$  indices in  $\mathcal{W}$ , and the tuple  $[d_i, w_i]^T$  denotes an edge between the two vertices indexed by  $d_i$  and  $w_i$  respectively. Such representation is comprehensive and is used in many implementations of LDA, e.g. [?].

However, many papers (including this article) do not use this representation, because it misleads readers to consider two sets of random variables,  $W$  and  $\mathbf{d}$ . The fact is that  $\mathbf{d}$  is the structure of  $W$  and is highly dependent with  $W$ .  $\mathbf{d}$  is also the structure of the hidden variables  $Z$ , as shown by (3).

### 3 Dirichlet and Multinomial

To derive the Gibbs sampling algorithm with LDA, it is important to be familiar with the conjugacy between Dirichlet distribution and multinomial distribution.

The Dirichlet distribution is defined as:

$$\text{Dir}(\mathbf{p}; \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{v=1}^{|\boldsymbol{\alpha}|} p_v^{\alpha_v - 1} , \quad (4)$$

where the normalizing constant is the multinomial beta function, which can be expressed in terms of gamma function:

$$B(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^{|\boldsymbol{\alpha}|} \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^{|\boldsymbol{\alpha}|} \alpha_i)} . \quad (5)$$

When Dirichlet distribution is used in LDA to model the prior,  $\alpha_i$ 's are positive integers. In this case, the gamma function degenerates to the factorial function:

$$\Gamma(n) = (n-1)! \quad (6)$$

The multinomial distribution is defined as:

$$\text{Mult}(\mathbf{x}; \mathbf{p}) = \frac{n!}{\prod_{i=1}^K x_i!} \prod_{i=1}^K p_i^{x_i} , \quad (7)$$

where  $x_i$  denotes the number of times that value  $i$  appears in the samples drawn from the discrete distribution  $\mathbf{p}$ ,  $K = |\mathbf{x}| = |\mathbf{p}| = |\boldsymbol{\alpha}|$  and  $n = \sum_{i=1}^K x_i$ . The normalization constant comes from the product of combinatorial numbers.

We say the Dirichlet distribution is the conjugate distribution of the multinomial distribution, because if the prior of  $\mathbf{p}$  is  $\text{Dir}(\mathbf{p}; \boldsymbol{\alpha})$  and  $\mathbf{x}$  is generated by  $\text{Mult}(\mathbf{x}; \mathbf{p})$ , then the posterior distribution of  $\mathbf{p}$ ,  $p(\mathbf{p}|\mathbf{x}, \boldsymbol{\alpha})$ , is a Dirichlet distribution:

$$\begin{aligned} p(\mathbf{p}|\mathbf{x}, \boldsymbol{\alpha}) &= \text{Dir}(\mathbf{p}; \mathbf{x} + \boldsymbol{\alpha}) \\ &= \frac{1}{B(\mathbf{x} + \boldsymbol{\alpha})} \prod_{v=1}^{|\boldsymbol{\alpha}|} p_v^{x_v + \alpha_v - 1} . \end{aligned} \quad (8)$$

Because (8) is a probability distribution function, integrating it over  $\mathbf{p}$  should result in 1:

$$\begin{aligned} 1 &= \int \frac{1}{B(\mathbf{x} + \boldsymbol{\alpha})} \prod_{v=1}^{|\boldsymbol{\alpha}|} p_v^{x_v + \alpha_v - 1} d\mathbf{p} \\ &= \frac{1}{B(\mathbf{x} + \boldsymbol{\alpha})} \int \prod_{v=1}^{|\boldsymbol{\alpha}|} p_v^{x_v + \alpha_v - 1} d\mathbf{p} . \end{aligned} \quad (9)$$

This implies

$$\int \prod_{v=1}^{|\boldsymbol{\alpha}|} p_v^{x_v + \alpha_v - 1} d\mathbf{p} = B(\mathbf{x} + \boldsymbol{\alpha}) . \quad (10)$$

This property will be used in the following derivations.

### 4 Learning LDA by Gibbs Sampling

There have been three strategies to learn LDA: EM with variational inference [?], EM with expectation propagation [?], and Gibbs sampling [?]. In this article, we focus on the Gibbs sampling method, whose performance is comparable with the other two but is tolerant better to local optima.

Gibbs sampling is one of the class of samplings methods known as Markov Chain Monte Carlo. We use it to sample from the posterior distribution,  $p(Z|W)$ , given the training data  $W$  represented in the form of (1). As will be shown by the following text, given the sample of  $Z$  we can infer model parameters  $\Phi$  and  $\Theta$ . This forms a learning algorithm with its general framework shown in Fig. 1.

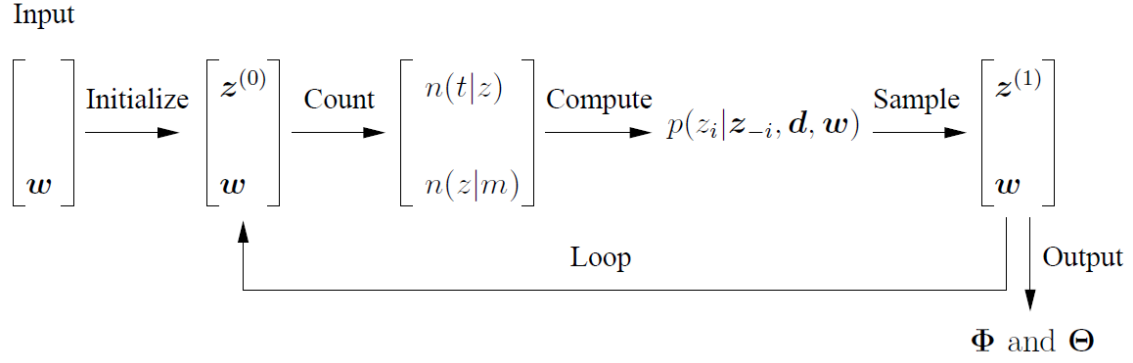


Figure 1: The procedure of learning LDA by Gibbs sampling.

|                             |   |
|-----------------------------|---|
| $N$                         | the number of words in the corpus   |
| $1 \leq i, j \leq N$        | index of words in the corpus  |
| $W = \{w_i\}$               | the corpus, and $w_i$ denotes a word  |
| $Z = \{z_i\}$               | latent topics assigned to words in $W$  |
| $W^{-i} = W \setminus w_i$  | the corpus excluding $w_i$  |
| $Z^{-i} = Z \setminus z_i$  | latent topics excluding $z_i$   |
| $K$                         | the number of topics specified as a parameter   |
| $V$                         | the number of unique words in the vocabulary  |
| $\alpha$                    | the parameters of topic Dirichlet prior   |
| $\beta$                     | the parameters of word Dirichlet prior  |
| $\Omega_{d,k}$              | count of words in $d$ assigned topic $k$ ;<br>$\Omega_d$ denotes the $d$ -th row of matrix $\Omega$ . |
| $\Psi_{k,v}$                | count of word $v$ in corpus assigned $k$ ;<br>$\Psi_k$ denotes the $k$ -th row of matrix $\Psi$ .     |
| $\Omega_{d,k}^{-i}$         | like $\Omega_{d,k}$ but excludes $w_i$ and $z_i$  |
| $\Psi_{k,v}^{-i}$           | like $\Psi_{k,v}$ but excludes $w_i$ and $z_i$  |
| $\Theta = \{\theta_{d,k}\}$ | $\theta_{d,k} = P(z = k d)$ , $\theta_d = P(z d)$ .   |
| $\Phi = \{\phi_{k,v}\}$     | $\phi_{k,v} = P(w = v z = k)$ , $\phi_k = P(w z = k)$ .   |

Table 1: Symbols used in the derivation of LDA Gibbs sampling rule.

In order to sample from  $p(Z|W)$  using the Gibbs sampling method, we need the full conditional posterior distribution  $p(z_i|Z^{-i}, W)$ , where  $Z^{-i}$  denotes all  $z_j$ 's with  $j \neq i$ . In particular, Gibbs sampling does not require knowing the exact form of  $p(z_i|Z^{-i}, W)$ ; instead, it is enough to have a function  $f(\cdot)$ , where

$$f(z_i|Z^{-i}, W) \propto p(z_i|Z^{-i}, W) \quad (11)$$

The following mathematical derivation is for such a function  $f(\cdot)$ .

**The Joint Distribution of LDA** We start from deriving the joint distribution <sup>4</sup>,

$$p(Z, W|\alpha, \beta) = p(W|Z, \beta)p(Z|\alpha) \quad (12)$$

which is the basis of the derivation of the Gibbs updating rule and the parameter estimation rule. As  $p(W|Z, \beta)$  and  $p(Z|\alpha)$  depend on  $\Phi$  and  $\Theta$  respectively, we derive them separately.

According to the definition of LDA, we have

$$p(W|Z, \beta) = \int p(W|Z, \Phi)p(\Phi|\beta)d\Phi \quad (13)$$

where  $p(\Phi|\beta)$  has Dirichlet distribution:

$$p(\Phi|\beta) = \prod_{k=1}^K p(\phi_k|\beta) = \prod_{k=1}^K \frac{1}{B(\beta)} \prod_{v=1}^V \phi_{k,v}^{\beta_v-1} \quad (14)$$

<sup>4</sup>Notations used in this article are identical with those used in [?].

and  $p(W|Z, \Phi)$  has multinomial distribution:

$$p(W|Z, \Phi) = \prod_{i=1}^N \phi_{z_i, w_i} = \prod_{k=1}^K \prod_{v=1}^V \phi_{k,v}^{\Psi_{k,v}}, \quad (15)$$

where  $\Psi$  is a  $K \times V$  count matrix and  $\Psi_{k,v}$  is the number of times that topic  $k$  is assigned to word  $v$ . With  $W$  and  $Z$  defined by (1) and (3) respectively, we can represent  $\Psi_{k,v}$  mathematically by

$$\Psi_{k,v} = \sum_{i=1}^N \mathbb{I}\{w_i = v \wedge z_i = k\}, \quad (16)$$

where  $N$  is the corpus size in number of words. In later sections, we use  $\Psi_k$  to denote the  $k$ -th row of the matrix  $\Psi$ .

Given (15) and (14), (13) becomes

$$p(W|Z, \beta) = \int \prod_{k=1}^K \frac{1}{B(\beta)} \prod_{v=1}^V \phi_{k,v}^{\Psi_{k,v} + \beta_v - 1} d\phi_k. \quad (17)$$

Using the property of integration of a product, which we learned in our colleague time,

$$\int \prod_{k=1}^K f_k(\phi_k) d\phi_1 \dots d\phi_K = \prod_{k=1}^K \int f_k(\phi_k) d\phi_k, \quad (18)$$

we have

$$\begin{aligned} p(W|Z, \beta) &= \prod_{k=1}^K \left( \int \frac{1}{B(\beta)} \prod_{v=1}^V \phi_{k,v}^{\Psi_{k,v} + \beta_v - 1} d\phi_k \right) \\ &= \prod_{k=1}^K \left( \frac{1}{B(\beta)} \int \prod_{v=1}^V \phi_{k,v}^{\Psi_{k,v} + \beta_v - 1} d\phi_k \right). \end{aligned} \quad (19)$$

Using property (10), we can compute the integration over  $\phi_k$  in a close form, so

$$p(W|Z, \beta) = \prod_{k=1}^K \frac{B(\Psi_k + \beta)}{B(\beta)}, \quad (20)$$

where  $\Psi_k$  denotes the  $k$ -th row of matrix  $\Psi$ .

Now we derive  $p(Z|\alpha)$  analogous to  $p(W|Z, \beta)$ . Similar with (14), we have

$$p(\Theta|\alpha) = \prod_{d=1}^D p(\theta_m|\alpha) = \prod_{d=1}^D \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_{d,k}^{\alpha_k - 1}; \quad (21)$$

similar with (15), we have

$$p(Z|\Theta) = \prod_{i=1}^N \theta_{d_i, z_i} = \prod_{d=1}^D \prod_{k=1}^K \theta_{d,k}^{\Omega_{d,k}}; \quad (22)$$

and similar with (20) we have

$$\begin{aligned} p(Z|\alpha) &= \int p(Z|\Theta) p(\Theta|\alpha) d\Theta \\ &= \prod_{d=1}^D \left( \int \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_{d,k}^{\Omega_{d,k} + \alpha_k - 1} d\theta_d \right) \\ &= \prod_{d=1}^D \frac{B(\Omega_d + \alpha)}{B(\alpha)}, \end{aligned} \quad (23)$$

where  $\Omega$  is a count matrix and  $\Omega_{d,k}$  is the number of times that topic  $k$  is assigned to words in document  $d$ ;  $\Omega_d$  denotes the  $d$ -th row of  $\Omega$ . With input data defined by (2),  $\Omega_{d,k}$  can be represented mathematically

$$\Omega_{d,k} = \sum_{i=1}^N \mathbb{I}\{d_i = d \wedge z_i = k\}, \quad (24)$$

where  $N$  is the corpus size by words. In later sections, we use  $\Omega_d$  to denote the  $d$ -th row of the matrix  $\Omega$ .

Given (20) and (23), the joint distribution (12) becomes:

$$\begin{aligned} p(Z, W | \alpha, \beta) &= p(W | Z, \beta) p(Z | \alpha) \\ &= \prod_{k=1}^K \frac{B(\Psi_k + \beta)}{B(\beta)} \cdot \prod_{d=1}^D \frac{B(\Omega_d + \alpha)}{B(\alpha)} \end{aligned} \quad (25)$$

**The Gibbs Updating Rule** With (25), we can derive the Gibbs updating rule for LDA:

$$p(z_i = k | Z^{-i}, W, \alpha, \beta) = \frac{p(z_i = k, Z^{-i}, W | \alpha, \beta)}{p(Z^{-i}, W | \alpha, \beta)} . \quad (26)$$

Because  $z_i$  depends only on  $w_i$ , (26) becomes

$$p(z_i | Z^{-i}, W, \alpha, \beta) \propto \frac{p(Z, W | \alpha, \beta)}{p(Z^{-i}, W^{-i} | \alpha, \beta)} \quad \text{where } z_i = k . \quad (27)$$

The numerator in (27) is (25); whereas the denominator has a similar form:

$$\begin{aligned} p(Z^{-i}, W^{-i} | \alpha, \beta) &= p(W^{-i} | Z^{-i}, \beta) p(Z^{-i} | \alpha) \\ &= \prod_{k=1}^K \frac{B(\Psi_k^{-i} + \beta)}{B(\beta)} \cdot \prod_{d=1}^D \frac{B(\Omega_d^{-i} + \alpha)}{B(\alpha)} , \end{aligned} \quad (28)$$

where  $\Psi_k^{-i}$  denotes the  $k$ -th row of the count  $K \times V$  count matrix  $\Psi^{-i}$ , and  $\Psi_{k,v}^{-i}$  is the number of times that topic  $k$  is assigned to word  $w$ , but with the  $i$ -th word and its topic assignment excluded. Similarly,  $\Omega_d^{-i}$  is the  $d$ -th row of count matrix  $\Omega^{-i}$ , and  $\Omega_{d,k}^{-i}$  is the number of words in document  $d$  that are assigned topic  $k$ , but with the  $i$ -th word and its topic assignment excluded. In more details,

$$\begin{aligned} \Psi_{k,v}^{-i} &= \sum_{\substack{1 \leq j \leq N \\ j \neq i}} \mathbb{I}\{w_j = v \wedge z_j = k\} , \\ \Omega_{d,k}^{-i} &= \sum_{\substack{1 \leq j \leq N \\ j \neq i}} \mathbb{I}\{d_j = d \wedge z_j = k\} , \end{aligned} \quad (29)$$

where  $d_j$  denotes the document to which the  $j$ -th word in the corpus belongs. Some properties can be derived from the definitions directly:

$$\begin{aligned} \Psi_{k,v} &= \begin{cases} \Psi_{k,v}^{-i} + 1 & \text{if } v = w_i \text{ and } k = z_i; \\ \Psi_{k,v}^{-i} & \text{all other cases.} \end{cases} \\ \Omega_{d,k} &= \begin{cases} \Omega_{d,k}^{-i} + 1 & \text{if } d = d_i \text{ and } k = z_i; \\ \Omega_{d,k}^{-i} & \text{all other cases.} \end{cases} \end{aligned} \quad (30)$$

From (30) we have

$$\begin{aligned} \sum_{v=1}^V \Psi_{z_i,v} &= 1 + \sum_{v=1}^V \Psi_{z_i,v}^{-i} \\ \sum_{k=1}^K \Omega_{d_i,k} &= 1 + \sum_{k=1}^K \Omega_{d_i,k}^{-i} \end{aligned} \quad (31)$$

Using (30) and (31), we can simplify (28):

$$p(z_i = k | Z^{-i}, W, \alpha, \beta) = \frac{B(\Psi_k + \beta)}{B(\Psi_k^{-i} + \beta)} \cdot \frac{B(\Omega_d + \alpha)}{B(\Omega_d^{-i} + \alpha)} . \quad (32)$$

where  $d$  denotes the document that contains  $w_i$ .

Further simplification can be achieved by substituting

$$B(\mathbf{x}) = \frac{\prod_{k=1}^{\dim \mathbf{x}} \Gamma(x_k)}{\Gamma(\sum_{k=1}^{\dim \mathbf{x}} x_k)} , \quad (33)$$

and we have

$$p(z_i = k | Z^{\neg i}, w_i = v, W^{\neg i}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\prod_{v=1}^V \Gamma(\Psi_{k,v} + \beta_v)}{\Gamma(\sum_{v=1}^V \Psi_{k,v} + \beta_v)} \cdot \frac{\prod_{k=1}^K \Gamma(\Omega_{d,k} + \alpha_k)}{\Gamma(\sum_{k=1}^K \Omega_{d,k} + \alpha_k)} . \quad (34)$$

Using (30), we can remove most factors in the products of (34) and get

$$p(z_i = k | Z^{\neg i}, w_i = v, W^{\neg i}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\frac{\Gamma(\Psi_{k,v} + \beta_v)}{\Gamma(\sum_{v=1}^V \Psi_{k,v} + \beta_v)}}{\frac{\Gamma(\Psi_{k,v}^{\neg i} + \beta_v)}{\Gamma(\sum_{v=1}^V \Psi_{k,v}^{\neg i} + \beta_v)}} \cdot \frac{\frac{\Gamma(\Omega_{d,k} + \alpha_k)}{\Gamma(\sum_{k=1}^K \Omega_{d,k} + \alpha_k)}}{\frac{\Gamma(\Omega_{d,k}^{\neg i} + \alpha_k)}{\Gamma(\sum_{k=1}^K \Omega_{d,k}^{\neg i} + \alpha_k)}} . \quad (35)$$

Here it is important to know that

$$\Gamma(y) = (y-1)! \quad \text{if } y \text{ is a positive integer} , \quad (36)$$

so we expand (35) and using (30) to remove most factors in the factorials. This leads us to the Gibbs updating rule for LDA:

$$p(z_i = k | Z^{\neg i}, w = v, W^{\neg i}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\Psi_{k,v} + \beta_{w_i} - 1}{\left[ \sum_{v=1}^V \Psi_{k,v} + \beta_t \right] - 1} \cdot \frac{\Omega_{d,k} + \alpha_k - 1}{\left[ \sum_{k=1}^K \Omega_{d,k} + \alpha_z \right] - 1} . \quad (37)$$

Note that the denominator of the second factor at the right hand side of (38) does not depend on  $z_i$ , the parameter of the function  $p(z_i | Z^{\neg i}, W, \boldsymbol{\alpha}, \boldsymbol{\beta})$ . Because the Gibbs sampling method requires only a function  $f(z_i) \propto p(z_i)$ , c.f. (11), we can use the following updating rule in practice:

$$p(z_i = k | Z^{\neg i}, w = v, W^{\neg i}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\Psi_{k,v} + \beta_{w_i} - 1}{\left[ \sum_{v=1}^V \Psi_{k,v} + \beta_t \right] - 1} \cdot [\Omega_{d,k} + \alpha_k - 1] . \quad (38)$$

**Parameter Estimation** By the definition of  $\phi_{k,v}$  and  $\theta_{d,k}$ , we have

$$\begin{aligned} \phi_{k,v} &= \frac{\Psi_{k,v} + \beta_t}{\left( \sum_{v'=1}^V \Psi_{k,v'} + \beta_{v'} \right)} , \\ \theta_{m,k} &= \frac{\Omega_{d,k} + \alpha_k}{\left( \sum_{k=1}^K \Omega_{d,k} + \alpha_z \right)} . \end{aligned} \quad (39)$$

**The Algorithm** With (38), we can realize the learning procedure shown in Fig. 1 using Algorithm. 1.

In this algorithm, we use a 2D array  $\text{NWZ}[\mathbf{w}, \mathbf{z}]$  to maintain  $\Psi$  and  $\text{NZM}[\mathbf{z}, \mathbf{m}]$  for  $\Omega$ . Also, to accelerate the computation, we use a 1D array  $\text{NZ}[\mathbf{z}]$  to maintain  $n(z) = \sum_{v=1}^{|\mathcal{W}|} n(t; z)$ . Note that we do not need an array  $\text{NM}[\mathbf{d}]$  for  $n(z) = \sum_{k=1}^{|\mathcal{Z}|} n(z, d)$ . Although it appears in (37), but it is not necessary in (38).

The input of this algorithm is not in the form of (1); instead, it is in a more convenient form: with each document represented by a set of words it includes, the input is given as a set of documents. Denote the  $n$ -th word in the  $m$ -th document by  $w_{m,n}$ , the corresponding topic is denoted by  $z_{m,n}$  and corresponding document is  $m$ . With these correspondence known, the above derivations can be easily implemented in the algorithm.

```

zero all count variables NWZ, NZM, NZ ;
foreach document  $m \in [1, D]$  do
  foreach word  $n \in [1, N_m]$  in document  $m$  do
    sample topic index  $z_{m,n} \sim \text{Mult}(1/K)$  for word  $w_{m,n}$ ;
    increment document-topic count:  $\text{NZM}[z_{m,n}, m]++$  ;
    increment topic-term count:  $\text{NWZ}[w_{m,n}, z_{m,n}]++$  ;
    increment topic-term sum:  $\text{NZ}[z_{m,n}]++$  ;
  end
end
while not finished do
  foreach document  $m \in [1, D]$  do
    foreach word  $n \in [1, N_m]$  in document  $m$  do
       $\text{NWZ}[w_{m,n}, z_{m,n}]--$ ,  $\text{NZ}[z_{m,n}]--$ ,  $\text{NZM}[z_{m,n}, m]--$  ;
      sample topic index  $\tilde{z}_{m,n}$  according to (40) ;
       $\text{NWZ}[w_{m,n}, \tilde{z}_{m,n}]++$ ,  $\text{NZ}[\tilde{z}_{m,n}]++$ ,  $\text{NZM}[\tilde{z}_{m,n}, m]++$  ;
    end
  end
  if converged and  $L$  sampling iterations since last read out then
    read out parameter set  $\Theta$  and  $\Phi$  according to (39) ;
  end
end
end

```

**Algorithm 1:** The Gibbs sampling algorithm that learns LDA.



Figure 2: The ground-truth  $\Phi$  used to synthesize our testing data.

The sampling operation in this algorithm is not identical with the full conditional posterior distribution (38); instead, it does not contain those “−1” terms:

$$p(z_i) = \frac{\Psi_{z_i, w_i} + \beta_{w_i}}{\left[ \sum_{v=1}^V \Psi_{k,v} + \beta_t \right]} \cdot [\Omega_{d_i, z_i} + \alpha_{z_i}] . \quad (40)$$

This makes it elegant to maintain the consistency of sampling new topics and updating the counts (NWZ, NZ and NZM) using three lines of code: decreasing the counts, sampling and increasing the counts.<sup>5</sup>

## 5 Experiments

The synthetic image data mentioned in [?] is useful to test the correctness of any implementation of LDA learning algorithm. To synthesize this data set, we fix  $\Phi = \{\phi_k\}_{k=1}^{K=10}$  as visualized in Fig. 2 (also, Fig. 1 in [?]), set  $\alpha = [1]$ , the number of words/pixels in each document/image  $d = 100$ . Because every image has  $5 \times 5$  pixels, the vocabulary size is 25.<sup>6</sup>

Fig. 3 shows the result of running Gibbs sampling to learn an LDA from this testing data set, where the left pane shows the convergence of the likelihood, and the right pane shows the updating process of  $\Phi$  along the iterations of the algorithm. The 20 rows of images in the right pane are visualizations of  $\Phi$  estimated at the iterations of 1, 2, 3, 5, 7, , 10, 15, 20, 25, 50, 75, 100, 125, 150, 175, , 200, 250, 300, 400, 499. From this figure we can see that since iteration 100, the estimated  $\Phi$  is almost identical to Fig. 1(a) of [?].

<sup>5</sup>This trick illustrates the physical meaning of those “−1” terms in (38). I learn this trick from Heinrich’s code at <http://arbylon.net/projects/LdaGibbsSampler.java>.

<sup>6</sup>All these settings are identical with those described in [?].

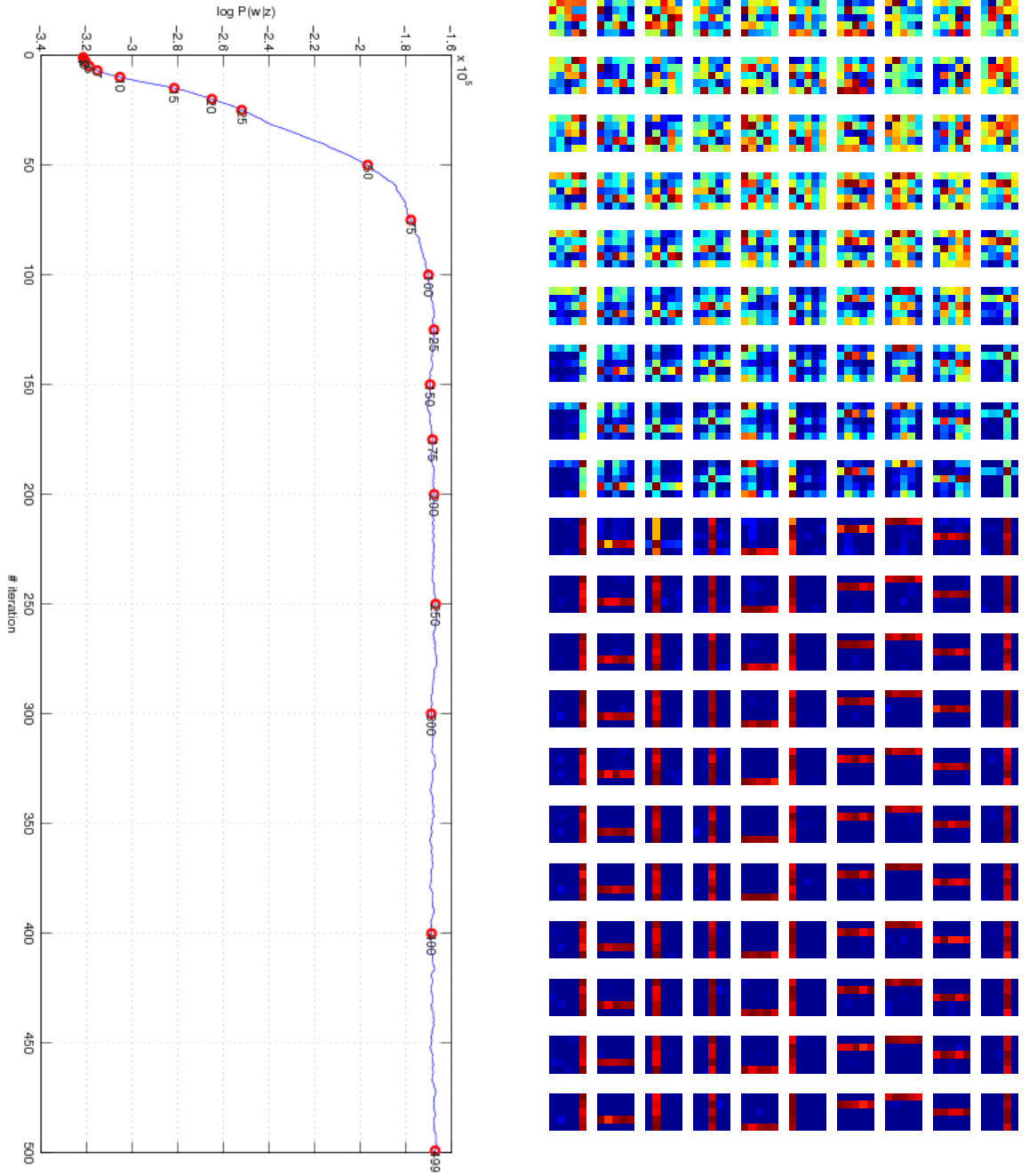


Figure 3: The result of running the Gibbs sampling algorithm to learn an LDA from synthetic data. Each row of 10 images in the right pane visualizes the  $\Phi = \{\phi_k\}_{k=1}^{K=10}$  estimated at those iterations indicated by red circles in the left pane.

## 6 Acknowledgement

Thanks go to Gregor Heinrich and Xuerui Wang for their gentle explanation of some math details.