

# Toward More Accurate Bug Triaging with Topic Modeling: Technical Report2

Xin Xia<sup>1,2</sup>, Ying Ding<sup>1</sup>, David Lo<sup>1</sup>, Jafar M. Al-Kofahi<sup>3</sup>,  
Tien N. Nguyen<sup>3</sup>, and Xinyu Wang<sup>2</sup>

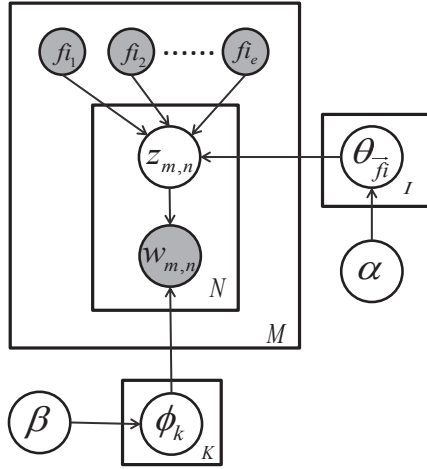
<sup>1</sup>School of Information Systems, Singapore Management University, Singapore

<sup>2</sup>College of Computer Science and Technology, Zhejiang University, China

<sup>3</sup>Electrical and Computer Engineering Department, Iowa State University, USA

xxkidd@zju.edu.cn, {ying.ding.2011, davidlo}@smu.edu.sg,

{jafar, tien}@iastate.edu, and wangxinyu@zju.edu.cn



**Figure 1: The Graphical Model of Multi-feature Topic Model (MTM)**

In this technical report, we describe the detail steps to derive the Gibbs sampling function of MTM. We first describe the generative process of MTM by using more formal description. Next, we present the training phase of MTM, and the formula derivation. Finally, we present the inference of topics for a new bug report.

## 1. GENERATIVE PROCESS

Our proposed *multi-feature topic model (MTM)* is a generative model that considers the features of a bug report. The general graphical model of MTM is shown in Figure 1. The structure is similar to LDA with a few exceptions: first,

we now have additional observed variables which are the features of bug reports such as the products and components of the reports. Also, rather than having a topic distribution per document, we now have a topic distribution per feature combination (denoted as  $\theta$  in the graphical model).

Various features (e.g., product, component, reporter, version, platform, etc.) can be considered in MTM. We denote the features as  $f_{i_1}, f_{i_2}, \dots, f_{i_e}$ . In this work, we only use an instance of MTM with two features: product and component – we set  $f_{i_1}$  to be the product, and  $f_{i_2}$  to be the component. The generative process of MTM is as follows:

1. For a feature combination  $\vec{f}_i = (f_{i_1}, f_{i_2}, \dots, f_{i_e})$  appearing in the corpus, sample the feature-topic distribution  $\theta_{\vec{f}_i}$  from  $\text{Dir}(\alpha)$ ;
2. For each topic  $t = 1, 2, \dots, K$ , sample the topic-word distribution  $\phi_t$  from  $\text{Dir}(\beta)$ ;
3. For each document (i.e., bug report)  $m = 1, 2, \dots, M$ , with feature combination  $\vec{f}_i^m$ :  
For each word  $w_{mn}$  in document  $m$ :  
Sample a topic  $z_{mn}$  from  $\text{Mult}(\theta_{\vec{f}_i^m})$ ;  
Sample a word  $w_{mn}$  from  $\text{Mult}(\phi_{z_{mn}})$ .

## 2. TRAINING PHASE

Table 1 summarizes the notations used in the remaining part of this section. Our goal is to generate the topics (i.e.,  $z_{mn}$ ) for all the words in the  $M$  documents (i.e., bug reports). Gibbs sampling is often used to recover the topics given a textual corpus and a graphical model [1]. Gibbs sampling is a Markov Chain Monte Carlo method, which is widely used in parameter estimation. In Gibbs sampling, we perform many iterations. For each iteration, we perform many steps. At each step, we estimate the value of one latent variable (e.g., the topic of a word in a document (i.e., bug report)) while keeping the values of all other latent variables constant (e.g., the topics of all other words in the documents). The process continues until a certain number of iterations have been performed or the process converges (i.e., there is no more change to the values assigned to the latent variables).

At each step of Gibbs sampling, we need to calculate the probability of assigning topic  $k$  to each word  $w_{mn}$  based

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

on the topic assignments of all other words and all other variables of the graphical model:

$$p(z_{mn} = k | \mathbf{Z}_{-mn}, \mathbf{W}, \mathbf{FI}, \alpha, \beta)$$

Here,  $\mathbf{Z}$  denotes the topic assignment of all words made so far,  $\mathbf{Z}_{-mn}$  denotes the topic assignment of all words excluding the  $n^{th}$  word in the  $m^{th}$  document,  $\mathbf{W}$  denotes all the words in the bug report collections, and  $\mathbf{FI}$  denotes the feature information collection of all reports (i.e., feature combinations of each bug report in the collection). We refer to this probability as the *Gibbs sampling function*. The derivation of the Gibbs sampling functions used in the training phase (i.e., Equation (12)) is as follows:

The Gibbs sampling function for the training phase can be calculated by Bayes theorem as follows:

$$\begin{aligned} & p(z_{mn} = k | \mathbf{Z}_{-mn}, \mathbf{W}, \mathbf{FI}, \alpha, \beta) \\ = & \frac{p(z_{mn}=k, \mathbf{Z}_{-mn}, \mathbf{W} | \mathbf{FI}, \alpha, \beta)}{p(\mathbf{Z}_{-mn}, \mathbf{W} | \mathbf{FI}, \alpha, \beta)} \\ = & \frac{p(\mathbf{Z}, \mathbf{W} | \mathbf{FI}, \alpha, \beta)}{p(\mathbf{Z}_{-mn}, \mathbf{W}_{-mn} | \mathbf{FI}, \alpha, \beta) p(w_{mn} | \mathbf{Z}_{-mn}, \mathbf{W}_{-mn}, \mathbf{FI}, \alpha, \beta)} \\ \propto & \frac{p(\mathbf{Z}, \mathbf{W} | \mathbf{FI}, \alpha, \beta)}{p(\mathbf{Z}_{-mn}, \mathbf{W}_{-mn} | \mathbf{FI}, \alpha, \beta)} \end{aligned} \quad (1)$$

Since  $p(w_{mn} | \mathbf{Z}_{-mn}, \mathbf{W}_{-mn}, \mathbf{FI}, \alpha, \beta)$  is not related to the value of  $z_{mn}$ , we simplify it. Based on our graphical model, the characteristic joint distribution of multi-feature topic model is:

$$p(\mathbf{Z}, \mathbf{W} | \alpha, \beta, \mathbf{FI}) = p(\mathbf{W} | \mathbf{Z}, \beta) p(\mathbf{Z} | \alpha, \mathbf{FI}) \quad (2)$$

The first term at the right part of the above equation is the same as the first term of the characteristic joint distribution of LDA (c.f. [1]), which is:

$$p(\mathbf{W} | \mathbf{Z}, \beta) = \prod_{k=1}^K \frac{\text{Del}(C_{k\cdot}^{KV} + \beta)}{\text{Del}(\beta)} \quad (3)$$

where  $C_{k\cdot}^{KV}$  is the  $k^{th}$  row of matrix  $C^{KV}$ , which is a vector representing the number of times each word is assigned to topic  $k$ .  $\text{Del}$  is the Dirichlet Delta function, c.f., [1], which can be calculated as:

$$\text{Del}(\beta) = \frac{\prod_{w=1}^V \Gamma(\beta)}{\Gamma(\beta V)} \quad (4)$$

Similarly,  $p(\mathbf{Z} | \alpha, \mathbf{FI})$  is derived as follows,

$$\begin{aligned} p(\mathbf{Z} | \alpha, \mathbf{FI}) &= \int p(\mathbf{Z}, \theta | \alpha, \mathbf{FI}) d\theta \\ &= \int p(\mathbf{Z} | \theta, \mathbf{FI}) p(\theta | \alpha, \mathbf{FI}) d\theta \end{aligned} \quad (5)$$

and

$$\begin{aligned} p(\theta | \alpha, \mathbf{FI}) &= \prod_{\vec{f}_i} p(\theta_{\vec{f}_i} | \alpha) \\ &= \prod_{\vec{f}_i} \frac{1}{\text{Del}(\alpha)} \times \prod_{k=1}^K \theta_{\vec{f}_i, k}^{\alpha-1} \end{aligned} \quad (6)$$

$$\begin{aligned} p(\mathbf{Z} | \theta, \mathbf{FI}) &= \prod_{m=1}^M \prod_{n=1}^{L_m} \prod_{k=1}^K \theta_{\vec{f}_i^m, k}^{I(z_{mn}=k)} \\ &= \prod_{m=1}^M \prod_{k=1}^K \theta_{\vec{f}_i^m, k}^{C(m, k)} \end{aligned} \quad (7)$$

where  $L_m$  is the length of  $m^{th}$  document and  $I(\cdot)$  is an indicator function which returns 1 when the condition is true,  $C(m, k)$  is document-topic count, i.e., the number of times topic  $k$  is assigned to the words in the  $m^{th}$  document. Combining the equations (5), (6) and (7), we derive:

$$\begin{aligned} p(\mathbf{Z} | \alpha, \mathbf{FI}) &= \int \prod_{\vec{f}_i} \frac{1}{\text{Del}(\alpha)} \prod_{k=1}^K \theta_{\vec{f}_i, k}^{C_{\vec{f}_i, k}^{IK} + \alpha - 1} d\theta \\ &= \prod_{\vec{f}_i} \frac{\text{Del}(C_{\vec{f}_i, \cdot}^{IK} + \alpha)}{\text{Del}(\alpha)} \end{aligned} \quad (8)$$

In the above equation,  $C_{\vec{f}_i, \cdot}^{IK}$  is the feature-topic count matrix (see Table 1), and  $C_{\vec{f}_i, k}^{IK}$  is the number of times words inside document with feature information  $\vec{f}_i = (f_{i1}, \dots, f_{ie})$  is assigned to topic  $k$ , and  $C_{\vec{f}_i, \cdot}^{IK}$  is a vector of number of times each topic appears in documents with feature information  $\vec{f}_i$ . From equations (3) and (8), we derive:

$$\begin{aligned} p(\mathbf{Z}, \mathbf{W} | \alpha, \beta, \mathbf{FI}) &= \prod_{k=1}^K \frac{\text{Del}(C_{k\cdot}^{KV} + \beta)}{\text{Del}(\beta)} \times \\ &\quad \prod_{\vec{f}_i} \frac{\text{Del}(C_{\vec{f}_i, \cdot}^{IK} + \alpha)}{\text{Del}(\alpha)} \end{aligned} \quad (9)$$

Similarly, we can get the denominator of Equation (1):

$$\begin{aligned} & p(\mathbf{Z}_{-mn}, \mathbf{W}_{-mn} | \alpha, \beta, \mathbf{FI}) = \\ & \prod_{k=1}^K \frac{\text{Del}(C_{k\cdot, -mn}^{KV} + \beta)}{\text{Del}(\beta)} \times \prod_{\vec{f}_i} \frac{\text{Del}(C_{\vec{f}_i, \cdot, -mn}^{IK} + \alpha)}{\text{Del}(\alpha)} \end{aligned} \quad (10)$$

The Gamma function  $\Gamma(a)$  has an important property:

$$\Gamma(a) = a \times \Gamma(a - 1) \quad a > 1 \quad (11)$$

Based on this, we can get the final Gibbs sampling function:

$$\begin{aligned} p(z_{mn} = k | \mathbf{Z}_{-mn}, \mathbf{W}, \mathbf{FI}, \alpha, \beta) &\propto \frac{C_{kv, -mn}^{KV} + \beta}{\sum_{v'} C_{kv', -mn}^{KV} + V\beta} \\ &\quad \times \frac{C_{\vec{f}_i, k, -mn}^{IK} + \alpha}{\sum_{k'} C_{\vec{f}_i, k', -mn}^{IK} + K\alpha} \end{aligned} \quad (12)$$

Here,  $C_{kv, -mn}^{KV}$  is the number of times word  $v$  is assigned to topic  $k$  excluding  $w_{mn}$  (i.e., the  $n^{th}$  word in the  $m^{th}$  document).  $C_{\vec{f}_i, k, -mn}^{IK}$  is the number of times topic  $k$  is assigned to a word inside a document with feature combination  $\vec{f}_i$  excluding  $w_{mn}$ . In other words,

Table 1: Symbols Associated with Multi-feature Topic Model

Notation	Type	Description
$M$	scalar	Numbers of documents (i.e., bug reports) in the document collection.
$K$	scalar	Numbers of topics.
$V$	scalar	Number of unique terms in the documents.
$e$	scalar	Number of features.
$\alpha$	scalar	Dirichlet prior, hyperparameter for the topic distribution for each feature combination.
$\beta$	scalar	Dirichlet prior, hyperparameter for the word distribution for each topic.
$E$	vector	Different feature combinations of bug reports in the document collection.
$I$	scalar	Number of feature combinations.
$\vec{f}_i$	vector	Vector representation of a feature combination, i.e., $\vec{f}_i \in E$ , and $\vec{f}_i = (f_{i_1}, f_{i_2}, \dots, f_{i_e})$
$Z$	vector	Topic assignment of all words.
$\mathbf{FI}$	vector	Feature combinations of bug reports in the document collection.
$\mathbf{W}$	vector	All words in the bug reports in the document collection.
$\phi_k$	vector	Word distribution for topic $k$ .
$\theta_{\vec{f}_i}$	vector	Topic distribution for feature combination $\vec{f}_i$ .
$\vec{f}_i^m$	vector	Feature combination of the $m^{th}$ bug report, i.e., $\vec{f}_i^m = (f_{i_1}^m, \dots, f_{i_e}^m)$ .
$z_{mn}$	scalar	Topic of the $n^{th}$ word in the $m^{th}$ bug report.
$w_{mn}$	scalar	$n^{th}$ word in the $m^{th}$ bug report.
$C^{KV}$	$K \times V$ matrix	Topic-word count matrix, i.e., a matrix containing the number of times various topics are assigned to various words in the document collection. Also, $C_{kv}^{KV}$ is the number of times word $v$ in the vocabulary is assigned to topic $k$ .
$C^{IK}$	$I \times K$ matrix	Feature-topic count matrix, i.e., a matrix containing the number of times various topics are assigned to various feature combinations. A topic is assigned to a feature combination if it is assigned to a word inside a document with that feature combination. Also, $C_{\vec{f}_i, k}^{IK}$ is the number of times topic $k$ is assigned to a word inside a document with feature combination $\vec{f}_i$ .

$$C_{kv, -mn}^{KV} = \begin{cases} C_{kv}^{KV} - 1 & \text{if } w_{mn} = v \text{ and } z_{mn} = k \\ C_{kv}^{KV} & \text{otherwise} \end{cases} \quad (13)$$

$$C_{\vec{f}_i, k, -mn}^{IK} = \begin{cases} C_{\vec{f}_i, k}^{IK} - 1 & \text{if } \vec{f}_i^m = \vec{f}_i \text{ and } z_{mn} = k \\ C_{\vec{f}_i, k}^{IK} & \text{otherwise} \end{cases} \quad (14)$$

In the above equation,  $C_{kv}^{KV}$  denotes the number of times word  $v$  is assigned to topic  $k$ .  $\vec{f}_i^m$  denotes the feature combination of the  $m^{th}$  document, i.e.,  $\vec{f}_i^m = (f_{i_1}^m, \dots, f_{i_e}^m)$ .  $C_{\vec{f}_i, k}^{IK}$  denotes the number of times topic  $k$  is assigned to a word inside a document with feature combination  $\vec{f}_i$ .

---

**Algorithm 1** The Training Phase of MTM.

---

**Input:** Document collection  $\mathcal{C}$ , Feature information collection of all reports  $\mathbf{FI}$ , Hyperparameters  $\alpha, \beta$ , Topic number  $K$ , Maximum number of iteration  $maxIters$ .

**Output:** Topic assignment for all of the words in  $\mathcal{C}$ :  $Z$

**Initialization:**

```

1: for  $m \leftarrow 1, \dots, M$  do
2:   for  $n \leftarrow 1, \dots, L_m$  do
3:     Sample topic  $z_{mn}$  from  $\text{Mult}(1/K, \dots, 1/K)$ 
4:   end for
5: end for
6: Initialize  $C^{KV}$  and  $C^{IK}$  matrices
Sampling:
7: repeat
8:   for  $m \leftarrow 1, \dots, M$  do
9:     for  $n \leftarrow 1, \dots, L_m$  do
10:      Randomly assign a topic to  $w_{mn}$  according to Equation (12)
11:     end for
12:   end for
13:   Update  $C^{KV}$  and  $C^{IK}$  matrices
14: until Converge or iterated more than  $maxIters$  steps

```

---

The Gibbs sampling algorithm for the training phase of MTM is shown in Algorithm 1. In the algorithm, we have two matrices  $C^{KV}$  and  $C^{IK}$  which are updated at each Gibbs sampling step.  $C^{KV}$  is a topic-word count matrix, i.e., a matrix containing the number of times various topics are assigned to various words in the document collection  $\mathcal{C}$ , and  $C^{IK}$  is a feature-topic count matrix, i.e., a matrix containing the number of times various topics are assigned to various feature combinations. A topic is assigned to a feature combination if it is assigned to a word inside a document with that feature combination. In the initialization part (Lines 1 to 6), we first randomly assign a topic to each word in our corpus and update the two matrices according to this random assignment. In the sampling part (Lines 7 to 13), we perform the following steps iteratively:

1. For each word  $w_{mn}$ , we assign a topic randomly to it according to the probability computed using Equation (12). (Line 10).
2. After going through all words in our corpus, we update the two matrices (i.e.,  $C^{KV}$  and  $C^{IK}$ ) according to the new topic assignments (Line 13).
3. This process repeats many times until convergence or a predefined number of iterations has been reached (Line 14).

### 3. INFERENCE PHASE

In the inference phase, we infer the topic distribution of a new bug report  $m_{new}$ . We fix the topic assignments of words in the training bug reports, and use Gibbs sampling to iterate through the terms in the new bug report enough number of times to get their topic assignments. At each step of Gibbs sampling, we need to calculate the probability of assigning topic  $k$  to each word in  $m_{new}$ . In the inference

part, we estimate the topic assignments of words in a new document. For a new document  $m_{new}$  with feature combination  $f_i^{new}$ , we can estimate the topic assignment of its  $n^{th}$  word. The joint probability of words and topics for  $m_{new}$  is:

$$p(\mathbf{Z}^{new}, \mathbf{W}^{new}, \mathbf{Z}, \mathbf{W} | \alpha, \beta, \vec{f}_i^{new}, \mathbf{FI}) \\ = p(\mathbf{W}^{new}, \mathbf{W} | \mathbf{Z}^{new}, \mathbf{Z}, \beta) p(\mathbf{Z}^{new}, \mathbf{Z} | \alpha, \vec{f}_i^{new}, \mathbf{FI}) \quad (15)$$

Similar to the derivation of the Gibbs sampling function for the training part (i.e., Equation (12)), we can derive Equation (16). This probability is given by the following:

---

**Algorithm 2** The Inference Phase of MTM.

---

**Input:** Words of training data  $\mathbf{W}$ , Topic assignments of words in training data  $\mathbf{Z}$ , Feature information collection of the training data  $\mathbf{FI}$ , Features of the new document  $\vec{f}_i^{new}$ , Words in the new document (i.e., bug report)  $\mathbf{W}^{new}$ , Maximum number of iteration  $maxIters$ .

**Output:** Topic assignment of all the words in the new bug report:  $\mathbf{Z}^{new}$

**Initialization:**

```

1: for  $n \leftarrow 1, 2, \dots, L_{new}$  do
2:   Sample topic  $z_n$  from  $\text{Mult}(1/K, \dots, 1/K)$ .
3: end for
4: Initialize  $C'^{KV}$  and  $C'^K$  matrices
Sampling:
5: repeat
6:   for  $n \leftarrow 1, 2, \dots, L_{new}$  do
7:     Randomly assign a topic to  $n^{th}$  word according to Equation (16)
8:   end for
9:   Update  $C'^{KV}$  and  $C'^K$  matrices
10: until Converge or iterated more than  $maxIters$  steps

```

---

$$p(z_n^{new} = k | \mathbf{Z}_{-n}^{new}, \mathbf{Z}, \mathbf{W}^{new}, \mathbf{W}, \vec{f}_i^{new}, \mathbf{FI}, \alpha, \beta) \\ \propto \frac{C_{kv}^{KV} + C'_{kv, -n}^{KV} + \beta}{\sum_{v'} [C_{kv'}^{KV} + C'_{kv', -n}^{KV}] + V\beta} \times \\ \frac{C_{\vec{f}_i^{new}, k}^{IK} + C'_{k, -n}^{IK} + \alpha}{\sum_{k'} [C_{\vec{f}_i^{new}, k'}^{IK} + C'_{k', -n}^{IK}] + K\alpha} \quad (16)$$

In the above equation,  $\mathbf{Z}_{-n}^{new}$  denotes the current topic assignment of all words excluding the  $n^{th}$  word in the new bug report.  $\mathbf{W}^{new}$  denotes all words in the new bug report.  $C'^{KV}$  and  $C'^K$  represent the topic-word count matrix and feature-topic count matrix for the new bug report  $m_{new}$ .  $C'_{kv, -n}^{KV}$  denotes the number of times word  $v$  is assigned to topic  $k$  in the new bug report excluding the  $n^{th}$  word,  $C'_{k, -n}^K$  is the number of times topic  $k$  is assigned to the new bug report excluding the  $n^{th}$  word.

Algorithm 2 presents the inference phase of MTM by using Gibbs sampling. In the algorithm, we have two matrices  $C'^{KV}$  and  $C'^K$  which are updated at each Gibbs sampling step. In the initialization part (Lines 1 to 4), we first randomly assign a topic to each word in the new bug report and update the two matrices according to this random assignment. In the sampling part (Lines 5 to 10), we perform the following steps iteratively:

1. For the  $n^{th}$  word, we assign a topic randomly to it according to the probability computed using equation (16) (Line 7).
2. After going through all words in the new bug report, we update the two matrices (i.e.,  $C'^{KV}$  and  $C'^K$ ) according to the new topic assignments (Line 9).
3. This process repeats many times until convergence or a predefined number of iterations has been reached (Line 10).

## 4. REFERENCES

- [1] G. Heinrich, "Parameter estimation for text analysis," Web: <http://www.arbylon.net/publications/text-est.pdf>, 2005.