

# Dynamic Topic Detection Model by Fusing Sentiment Polarity

Xi Ding<sup>1</sup>, Lanshan Zhang<sup>2</sup>, Ye Tian<sup>1</sup>, Xiangyang Gong<sup>1</sup> and Wendong Wang<sup>1</sup>

<sup>1</sup>State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications  
Beijing, 100876, China

<sup>2</sup>School of Digital Media and Design Arts, Beijing University of Posts and Telecommunications  
Beijing, 100876, China

Email: dingxi515@163.com, zls326@sina.com, {yetian, xygong, wdwang}@bupt.edu.cn

## Abstract

Traditional static topic models mainly focus on the statistical correlation between words, but ignore the sentiment tendency and the temporal properties which may have great effects on topic detection results. This paper proposed an LDA-based dynamic sentiment-topic (DST) model, which could not only detect and track topics but could also analyse the shift of general's sentiment tendency towards certain topic. This model combines the data with the sentiment and dynamic properties of time by maximum likelihood estimation and the sliding window. We use Gibbs sampling method to estimate and update model parameters, and use random EM algorithm for model reasoning. Experiments on real dataset demonstrate that DST model outperforms the existing algorithms.

*Keywords:* dynamic sentiment topic model, sentiment analysis, Gibbs sampling.

## 1 Introduction

With the rapid development of Web 2.0 technology, social media, represented by Facebook, Twitter, Blog and Weibo, has really taken off in the last few years. More and more people choose social media as main information exchange platform to publish and access real-time information. In addition to daily communication functions, the huge amount of user generated content (UGC) often also carries people's attitudes towards certain public events or products/services. Existing researches show that compared with traditional media channels, social media data which conveys public sentiment has more important social and economic value for community management agency, enterprise and the general public. Affected by this, the sentiment analysis and topic detection of unstructured social media data has emerged as a hot research hotspot among the various tasks of social networking analysis area.

Researchers who are concentrating on the area of sentiment analysis mainly focus on studying the sentiment polarity (positive, neutral, negative) classification methods of the social media data. Although a host of research achievements have been disclosed [Pang, B. and Lee, L.][Aue, A. and Gamon, M][Turney, P, D], they are mostly based on supervised learning methods, and there still existed two limitations. First, a lot of manually annotated

samples are needed for parameters adjustment process. Second, sentiment classifier trained from certain topic area often does not apply to other topic areas since the sentiment distribution and topic content are closely related. Besides, an unavoidable problem is that the sentiment word dictionary that was trained from traditional corpus could not be applied to social data that are flexibly expressed with much more emotion icons and disjunctive questions used, considering the huge difference of language style between social media and traditional media. Topic detection derives from TDT (Topic Detection and Tracking) technology, which mainly focuses on detecting and organizing unknown topics from traditional formal expressed text stream. Topic detection task consists of two branches: historical topic detection and online topic detection. For the former, its objective is to dig out the hidden topics from a given corpus with unsupervised clustering means. Each cluster corresponds to a certain topic. While online topic detection determines whether newly arrived text stream belongs to an existing topic or a new topic according to historic information. Compared with newspaper, periodicals and academic report, social media are most non-standard expressed unstructured short texts in real-time formalism. These features make the task of topic detection which takes social media as study object more challenging.

Sentiment analysis and topic detection are taken as two independent research tasks in current social networking analysis fields. However, "sentiment" and "topic" are two highly associated concepts. On one hand, the generation and spread process of sentiment must rely on a certain body, i.e., a specified topic. On the other hand, the change of sentiment would react on its carrier, i.e., certain specified topic, and consequently affect the evolution of topic. Take the event "the loss of communication of flight M370" as example. At the beginning stage, anxiety and trepidation are circulating widely among the people who express concern about this event. As time goes on, the mood gradually evolves into sadness and discontent. Under the influence of this mood, M370 event evolves into a new stage of "relatives' doubt about Malaysian Airline's emergency treatment" from "the lost of aircraft" and "the search and rescue". Thus it could be seen, the probability distribution of topic is affected by sentiment, and these two concepts are highly relevant.

By fusing the polarity of sentiment, a dynamic sentiment-topic (DST) model is proposed based on Latent Dirichlet Allocation (LDA) in this paper. LDA model assumes that, a document is composed of various topics with different probabilistic combination, within which each topic itself is also a probability distribution of a series of words. For a given corpus, LDA detects all hidden

topics by computing the posterior probability of each word that it belongs to a certain topic. LDA model gains tremendous success in topic analysis area. It is widely applied to social networking analysis areas [Lacoste, S, Sha, F and Jordan, M.][Ramage, D., Hall, D., Nallapati, R., and Manning, C.][Wallach, H., Mimno, D., McCallum, A.]. LDA model divides the document generation process into two stages: “document-topic” generative process and “topic-word” generative process. However, in “document-topic-word” three layer architecture, “topic” is affected by the polarity of “sentiment”. The influence is mainly embodied in three aspects. Firstly, “document-topic” probability distribution would change as time passes. The probability distribution of topics within a document at time slot  $t$  is different with that at time slot  $t-1$ , however the probability distribution at time slot  $t$  depends on the probability distribution at time slot  $t-1$ . Secondly, the probability distribution of the sentiment polarity would change as time passes. The probability distribution of the sentiment polarity at time slot  $t$  depends on the probability distribution at time slot  $t-1$ . Thirdly, “topic-word” probability distribution would change as time passes. The probability distribution of word within a topic at time slot  $t$  is different with that at time slot  $t-1$ . However the probability distribution at time slot  $t$  depends on the probability distribution at time slot  $t-1$ . Based on the above considerations, the DST model proposed in this paper takes the influence that sentiment polarity exerts on the probability distribution of topic and the probability distribution of word into account. Also, DST model adjusts the probability distributions dynamically according to the dependencies in temporal dimension. Thus, for the large amount of social media data, this model could not only detect the hidden topics, but can also classify the topics according to their sentiment polarity. Furthermore, the evolution process of topic and its corresponding sentiment would be deduced.

The main contribution of this paper embodies in the following aspects. Firstly, by fusing sentiment polarity into the topic modelling process, DST model could properly reflect the impact that sentiment exerts on topic detection. Secondly, DST model enhances LDA model by introducing the dynamic characteristics, such that it is more appropriate for those social media data with strong real-time characteristics. Thirdly, experiments on a large real dataset demonstrate that DST model is superior to existing algorithms.

The rest part of this paper is organized as follow. Section 2 gives a brief introduction of the related work. In Section 3, we describe the proposed dynamic sentiment-topic model in detail. In the next section, we conduct an empirical analysis of DST model on a large real dataset. Section 5 concludes the whole paper and discusses the future work.

## 2 Related Work

Vast majority of research work has been conducted on topic analysis fields. Traditional topic detection method, such as TDT, aims to segment text stream into different news articles. TDT deals principally with formal text stream. Since the huge gap in language style, expressing way between formal text and social media text, TDT cannot work well. Statistical learning theory is the basis of

most traditional topic detection algorithms. Latent Semantic Indexing (LSI) [Hofmann, T] is a typical statistical learning based topic model, within which Single Value Decomposition (SVD) method is utilized to decompose the document-word matrix, and the subspace with the largest degree of distinction in TF-IDF (Term Frequency-Inverse Document Frequency) feature space then corresponds to the hidden topic. PLSA (Probabilistic Latent Semantic Analysis) which is also called Aspect Model derives from LSI. PLSA enhanced LSI from linear algebra level to probability space. Though PLSA is significant for topic modelling, the disadvantages are also obvious. The number of parameters in PLSA model grows linearly with the size of corpus, such that it may cause over-fitting problem. LDA (Latent Dirichlet Allocation) expands PLSA by introducing Dirichlet as prior distribution. LDA greatly reduces the number of required parameters, meanwhile the over-fitting problem get well solved. However, multinomial distribution based topic model (LDA) still face the over-fitting problem when it comes to short texts since the sparse vocabulary.

In recent years, plenty of improved topic models are proposed based on LDA. DTM (Dynamic Topic Model)[David, M., John, D.] models the natural parameter space with state-space model. It constructs a LDA-like topic model for corpus within each time slot, and these topic models locate in a common markov chain. That means dependency exists between the adjacent topic models, and the topic model within time slot  $t$  evolves from that within time slot  $t-1$ . Unlike DTM, TOT [Xueru, W.,McCallum, A.](Topic over Time) model could depict the long term evolution of a certain topic, and it can also predict the topic distribution within a given time slot. However, this model only takes the relationship between word and time into account, but ignores the influence that sentiment polarity imposes on topic modelling. A newly proposed JST (Joint Sentiment-Topic) model [He, Y. and Lin, C.] is weakly supervised. This model integrates sentiment factor into topic modelling process, and it assumes all words in a document embody sentiment tendency, so each topic also has its own sentiment attributes. Strong sentiment would influent people’s concern about the topic, and consequently change the topic distribution and word distribution. Based on this assumption, JST model could not only detect the topic hidden in text streams, but could also classify the topics according to its corresponding sentiment polarity. However, the evolution of sentiment polarity of a specified topic is not taken into account in JST model. This limitation leads to the inadaptability when deal with social media data which are timeliness and expressed with rich affections.

Seldom researches consider dynamic characteristics of topics and the influence of sentiment polarity simultaneously in topic modelling process. The dJST (dynamic joint sentiment-topic) model [He, Y., Lin, C., Gao, W., and Wong, K.-F] is most close to our work. It derives from LDA model, but combines both dynamic temporal characteristic of topic and author’s sentiment tendency in their model. dJST model estimates the hyper parameters of topic distribution and word distribution. The drawback is that dJST overlooks the fact that the sentiment polarity probability distribution in a document is not

constant, but always changing according to the evolution process of topic.

### 3 DST Model

#### 3.1 Overview

In this section we would detail the DST Model. It is a generative model, and belongs to a form of Bayesian probability model which is a predictive model that learns by the joint probability density distribution of the data, and then finds the conditional probability distribution. The model uses reasonable mathematical methods, such as maximum likelihood estimation and sliding window that makes the sentimental with the dynamic characteristics and time properties reasonably integrate into the existing traditional static text analysis model—LDA model, so we can get the public’s sentimental tendencies of certain topic which is occurred in a period of time and the changes of the sentimental tendencies of certain topic in the trend. Because of the actual situation, the text file generated by social network often goes with the author’s sentiment and sentimental tendencies vary among people which go with dynamic changes over time. The framework of our model has five layers where the document layer contains many types of sentimental tendencies, and topic layer is associated with sentiment layer, and the word layer is associated with sentiment and topic layers. Sentimental tendencies layer is the connecting bridge and plays an important role in the five layers. A graphical model of DST model is represented in Figure 1.

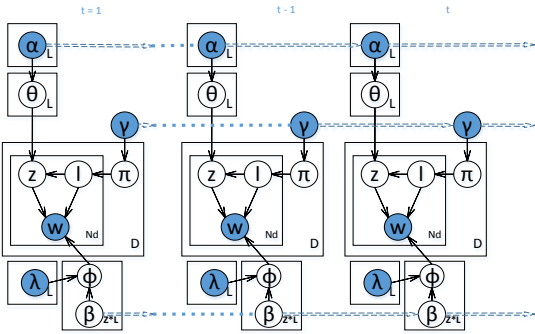


Figure 1: DST MODEL

Model assumes that at each epoch  $t$ , each of text documents set contains sentimental tendencies, and all the sentimental tendencies is divided into  $L$  Classes which express as  $L = \{l_1, l_2, \dots, l_L\}$ . The object processed by the model is text document collections, which are received with their order of publication timestamps preserved and are represented as  $D_t = \{d_{t_1}, d_{t_2}, \dots, d_{t_D}\}$  and the text document collections’ size is  $D_t$ . Appearing in the text document collections at epoch  $t$ , all words constitute a vector of word token whose size is  $V_t$ , and the word token can be expressed as  $V_t = \{W_{t_1}, W_{t_2}, \dots, W_{t_{V_t}}\}$ . Our notations are summarized in Table 1.

Model defines that for each epoch  $t$ , a document  $d_t$  contains  $L$  classes of sentimental tendencies, and the sentimental tendency obeys the sentiment distribution  $\pi_d^t$ . Each sentimental tendencies  $l$  can contain  $Z$  topics, and

Table 1: Notations used in the paper

Symbol	Description
$D_t$	number of documents in epoch $t$
$L$	number of sentiment labels
$Z$	number of topics
$V_t$	number of unique words in the current epoch $t$
$S$	number of time slices
$\gamma^t$	matrix of $D_t \times L$ dimension, row $d$ represents the priors of the mixing proportion of sentiment in document $d_t$
$\alpha_l^t$	matrix of $L \times Z$ dimension, row $l$ represents the priors of the mixing proportion of topics in sentiment label $l$
$\beta_{l,z}^t$	matrix of $L \times Z \times V_t$ dimension, priors for the word distribution conditioned on sentiment labels and topics
$\lambda_l^t$	matrix of $L \times V_t$ dimension which encodes the word prior sentiment polarity information
$\pi_d^t$	parameter notation for the sentiment label mixture proportion for document $d_t$ . $\pi^t = \{\pi_d^t\}_{d=1}^{D_t}$ ( $D_t \times L$ matrix)
$\theta_l^t$	multinomial distribution over topics for the $l$ th sentiment label for document $d_t$ . $\theta^t = \{\{\theta_{d,l}^t\}_{l=1}^L\}_{d=1}^{D_t}$ ( $D_t \times L \times Z$ matrix)
$\phi_{l,z}^t$	multinomial distribution over words for the $l$ th sentiment label and $z$ th topic at epoch $t$ . $\phi^t = \{\{\phi_{l,z}^t\}_{z=1}^Z\}_{l=1}^L$ ( $L \times Z \times V_t$ matrix)
$\delta_{l,z,s}^t$	multinomial word distribution of sentiment label $l$ and topic $z$ with time slice $s$ at epoch $t$ . $\delta_{l,z,s}^t = \{\delta_{l,z,s,w}^t\}_{w=1}^{V_t}$
$\sigma_{l,z,s}^t$	weight vector, $\sigma_{l,z,s}^t = \{\sigma_{l,z,s,1}^t, \dots, \sigma_{l,z,s,S}^t\}$ , each of which determines the contribution of time slice $s$ in computing the priors of $\phi_{l,z}^t$

these topics all obey the topic distribution  $\theta_l^t$ , and each of these topics is composed of  $V_t$  words which obey the word distributed  $\phi_{l,z}^t$ . In order to study the dynamic variable data, we assume that the distribution of current sentiment at the current epoch are influenced by sentiment in the past (so the Dirichlet hyper parameter  $\gamma^t$  corresponding to the distribution  $\pi_d^t$  obeys the gamma distribution, ie  $\gamma^t \sim \text{gamma}(\xi\gamma^{t-1}, \xi)$ ); and at current epoch  $t$ , the topic distribution corresponding to the document  $d$  and the sentimental tendencies  $l$  is affected by the topic in the past (ie the parameters  $\alpha_l^t$  with the corresponding Dirichlet distribution  $\theta_l^t$  obeys gamma distribution, ie  $\alpha_l^t \sim \text{gamma}(\mu\alpha_l^{t-1}, \mu)$ ); same time the current sentiment-topic specific word distributions  $\phi_{l,z}^t$  at epoch  $t$  are generated according to the word distributions at previous epochs (the word distributions at previous epochs is reflected in two parameters introduced  $\delta_{l,z}^t$ 、 $\sigma_{l,z,s}^t$ . These

two parameters take different values at different time scales. So our stipulate

$$\beta_{l,z}^t = \Sigma \delta_{l,z}^t * \sigma_{l,z,s}^t \quad (1)$$

Document generation process is mainly determined by the topic and emotional tendencies mark. Formal definition of document process in the model shown in Figure 1 is as follows:

- 
- (1) For each sentiment label  $l=1,\dots,L$
- For each topics  $z=1,\dots,Z$
- At the scale S, Compute  $\beta_{l,z}^t = \Sigma \delta_{l,z}^t * \sigma_{l,z,s}^t$  ;
- Select the matrix of  $L \times Z$  word distribution  $\varphi_{l,z}^t \sim Dir(\beta_{l,z}^t)$  ;
- (2) For each document  $d_i=1,\dots,D_i$
- According to  $\gamma^t \sim \text{gamma}(\xi\gamma^{t-1}, \xi)$  , compute  $\gamma^t$  ,Choose a sentiment distribution  $\pi_d^t \sim Dir(\overline{\gamma}_d^t)$  ;
- For each sentiment label  $l$  under document  $d$ , according to  $\alpha_l^t \sim \text{gamma}(\mu\alpha_l^{t-1}, \mu)$  , compute  $\alpha_l^t$  , choose a topic distribution  $\theta_l^t \sim Dir(\overline{\alpha}_{l,z}^t)$  ;
- For each word  $W=1,\dots,V$  in document  $d_i$
- Choose a sentiment label  $l_w \sim Mult(\pi_d^t)$  ;
- Choose a topic  $z_w \sim Mult(\theta_{l_w}^t)$  ;
- Choose a word  $W_w \sim Mult(\varphi_{l_w,z_w}^t)$  ;
- (3) Repeat the above process can get the whole document
- 

**Figure 2:** Document generation process

### 3.2 Inference Learning

We present a stochastic EM algorithm to sequentially update model parameters at each epoch using the newly obtained document set and the derived evolutionary parameters. At each iteration step, we use the collapsed Gibbs sampling to estimate potential distribution of sentiment label  $l$  and topics  $z$ , and use the maximum likelihood estimation method to estimate the hyper parameters.

The total probability of the model for the document set  $V_t$  at epoch  $t$  given the evolutionary parameters  $\delta_{l,z}^t$  ,  $\sigma_{l,z,s}^t$  and the previous model parameter is:

$$P(w^t, l^t, z^t | \overline{\gamma}^t, \overline{\alpha}_l^t, \overline{\delta}_{l,z}^t, \overline{\sigma}_{l,z,s}^t) = P(l^t | \overline{\gamma}^t) P(z^t | l^t, \overline{\alpha}_{l,z}^t) P(w^t | l^t, z^t, \overline{\delta}_{l,z}^t, \overline{\sigma}_{l,z,s}^t) \quad (2)$$

For the first term on the right-hand side of Equation (2), by integrating out  $\pi$  , we obtain

$$P(\overline{l}^t | \overline{\gamma}^t) = \prod_d \frac{\Gamma(\sum_l \gamma_l^t)}{\prod_l \Gamma(\gamma_l^t)} \cdot \frac{\prod_l \Gamma(N_{d,l}^t + \gamma_l^t)}{\Gamma(N_{d,l}^t + \sum_l \gamma_l^t)} \quad (3)$$

That is equal to  $p(\overline{l}^t | \overline{\gamma}^t) = \prod_{d=1}^{D_t} \frac{\Delta(n_{d,l}^t + \gamma_d^t)}{\Delta(\gamma_d^t)}$ . Where  $D_t$  is the total number of documents in epoch  $t$ ,  $N_{d,l}^t$  is the number of times sentiment label  $l$  being assigned to some word tokens in document  $d$  at epoch  $t$ ,  $N_d^t = \sum_l N_{d,l}^t$ . It is a constant which can be derived from statistical data input and  $\Gamma$  is the gamma function.

For the second term, by integrating out  $\theta$  , we obtain:

$$p(\overline{z}^t | \overline{\alpha}_k^t, \overline{l}^t) = \prod_d \prod_k \frac{\Gamma(\sum_z \alpha_{k,z}^t) \prod_z \Gamma(N_{d,k,z}^t + \alpha_{k,z}^t)}{\prod_z \Gamma(\alpha_{k,z}^t) \Gamma(N_{d,k,z}^t + \sum_z \alpha_{k,z}^t)} \quad (4)$$

That is equal to  $p(\overline{z}^t | \overline{\alpha}_k^t, \overline{l}^t) = \prod_{d=1}^{D_t} \prod_{k=1}^L \frac{\Delta(N_{d,k,z}^t + \alpha_{k,z}^t)}{\Delta(\alpha_{k,z}^t)}$ .

Where  $N_{d,k,z}^t$  is the number of times a word from document  $d$  being associated with topic  $z$  and sentiment label  $l$  at epoch  $t$ . It is a constant which can be derived from statistical data input and  $N_{d,k}^t = \sum_z N_{d,k,z}^t$ .

For the last term, by integrating out  $\varphi$  , we obtain:

$$p(\overline{w}^t | \overline{l}^t, \overline{z}_k^t) = \prod_k \prod_z \frac{\Gamma(\sum_s \delta_{k,z,s}^t) \prod_w \Gamma(N_{k,z,w}^t + \sum_s \delta_{k,z,s}^t \sigma_{k,z,s,w}^t)}{\prod_w \Gamma(\sum_s \delta_{k,z,s}^t \sigma_{k,z,s,w}^t) \Gamma(N_{k,z}^t + \sum_s \delta_{k,z,s}^t)} \quad (5)$$

That is equal to  $p(\overline{w}^t | \overline{l}^t, \overline{z}_k^t) = \prod_{k=1}^L \prod_{z=1}^Z \frac{\Delta(n_{k,z}^t + \delta_{k,z,s}^t \sigma_{k,z,s,w}^t)}{\Delta(\delta_{k,z,s}^t \sigma_{k,z,s,w}^t)}$ . Where  $N_{k,z,w}^t$  is

the number of times word  $w$  appeared in topic  $z$  and with sentiment label  $l$  at epoch  $t$ . It is a constant which can be derived from statistical data input and  $N_{k,z}^t = \sum_w N_{k,z,w}^t$ .

Gibbs sampling is the use of Markov Chain Monte Carlo method. According to the prior distribution of model parameters and statistical data which is given on the pre-S slots, it sequentially samples each variable of interest, sentiment label  $l_t$  and topic  $z_t$  through the case that it must meet the detailed balance condition. To meet this condition, let the index  $x=(d,n,t)$  and the subscript  $\setminus x$  denote a quantity that excludes counts in word position  $n$  of document  $d$  in epoch  $t$ , the conditional posterior for  $z_x$  and  $l_x$  by marginalizing out the random variables  $\varphi$  ,  $\theta$  , and  $\pi$  is:

$$p(z_x = j, l_x = k | \overline{w}^t, \overline{z}_{\setminus x}^t, \overline{l}_{\setminus x}^t, \overline{\alpha}_l^t, \overline{\gamma}^t, \overline{\delta}^t, \overline{\sigma}^t) \propto \frac{N_{k,j,w_j \setminus x}^t + \sum_s \delta_{k,j,s,w_j}^t \sigma_{k,j,w_j}^t \cdot N_{d,k,j \setminus x}^t + \alpha_{k,k}^t \cdot N_{d,k \setminus x}^t + \gamma_k^t}{N_{k,j \setminus x}^t + \sum_s \delta_{k,j,s}^t \cdot N_{d,k \setminus x}^t + \sum_j \alpha_{k,j}^t \cdot N_{d \setminus x}^t + \sum_k \gamma_k^t} \quad (6)$$

There are two types of parameters to be estimation which are evolutionary parameters and hyper parameters.

#### A. The evolutionary parameters estimation

We use the fixed-point iteration method to estimate the weight vector  $\delta_{l,z}^t$  directly from data by maximizing the joint distribution in Equation (2). The update formula is:

$$(\delta_{k,j,s}^t)^{new} \leftarrow \frac{\delta_{k,j,s}^t \sum_w \sigma_{d,k,j,w}^t A}{B} \quad (7)$$

Where,

$$A = \phi(N_{d,k,j,w}^t + \sum_s \delta_{k,j,s}^t \sigma_{d,k,j,w}^t) - \phi(\sum_s \delta_{k,j,s}^t \sigma_{d,k,j,w}^t) \quad (8)$$

$$B = \phi(N_{d,k,j}^t + \sum_s \delta_{k,j,s}^t) - \phi(\sum_s \delta_{k,j,s}^t) \quad (9)$$

And  $\phi(\cdot)$  is the digamma function defined by

$$\phi(x) = \frac{\partial \log \Gamma(x)}{\partial x}.$$

The evolutionary parameter  $\sigma_{d,k,j,w}^t$  accounts for the historical word distributions at different time slices. So we estimate  $\{\sigma_{d,k,j,w}^t\}_{w=1}^V$  the word distribution in topic  $j$  and sentiment label  $k$  at time slice  $s$ , which can be calculated as follows:

$$\sigma_{d,k,j,w}^t = \frac{C_{k,j,s,w}^t}{\sum_w C_{k,j,s,w}^t} \quad (10)$$

Where  $C_{k,j,s,w}^t$  is the expected number of times word  $w$  is assigned to sentiment label  $k$  and topic  $j$  at time slice  $s$ . At different time scales, its value is different. Thus  $C_{k,j,s,w}^t$  can be obtained directly from the count  $\hat{N}_{k,j,w}^t$ , that is  $C_{k,j,s,w}^t = \hat{N}_{k,j,w}^{t-s}$ , the expected number of times word  $w$  is associated with sentiment label  $l$  and topic  $z$  at epoch  $t'$ , which can be calculated by:

$$\hat{N}_{k,j,w}^t = N_{k,j,w}^t \cdot \frac{N_{k,j,w}^t + \sum_s \delta_{k,j,s}^t \sigma_{d,k,j,w}^t}{N_{k,j}^t + \sum_s \delta_{k,j,s}^t} \quad (11)$$

Where  $N_{k,j,w}^t$  is the observed count for the number of times word  $w$  is associated with sentiment label  $l$  and topic  $z$  at epoch  $t'$ .

## B. Hyper parameters ( $\alpha^t, \gamma^t$ ) estimation

We estimate  $\alpha^t$  and  $\gamma^t$  from data using maximum-likelihood as part of the online stochastic EM algorithm.

A common practice for the implementations of topic models is to use symmetric Dirichlet hyper parameters. However, it has been found that an asymmetric Dirichlet prior over the per-document topic proportions has substantial advantages over a symmetric prior. So when first entering a new epoch, we initialize the asymmetric  $\alpha_1^t = 0.01, \gamma^t = 0.01$ . Afterwards according that  $\gamma^t$  obeys

gamma distribution  $\gamma^t \sim \text{gamma}(\xi \gamma^{t-1}, \xi)$  which lead to the probability function of  $\gamma^t$  is

$$p(\gamma^t | \gamma^{t-1}, \xi) = \prod_t \frac{\xi \xi \gamma^{t-1} (\gamma^t)^{\xi \gamma^{t-1} - 1} \exp(-\xi \gamma^t)}{\Gamma(\xi \gamma^{t-1})} \quad \text{and } \alpha_1^t \text{ obeys}$$

gamma distribution  $\alpha_1^t \sim \text{gamma}(\mu \alpha_1^{t-1}, \mu)$  which lead to the probability function of  $\gamma^t$  is

$$p(\alpha_k^t | \alpha_k^{t-1}, \mu) = \prod_z \frac{\mu \mu \alpha_{k,z}^{t-1} (\alpha_{k,z}^t)^{\mu \alpha_{k,z}^{t-1} - 1} \exp(-\mu \alpha_{k,z}^t)}{\Gamma(\mu \alpha_{k,z}^{t-1})}$$

For every 40 Gibbs sampling iterations,  $\alpha_1^t$  and  $\gamma^t$  are learned directly from data using maximum-likelihood estimation.

$$(\alpha_{k,z}^t)^{new} \leftarrow \frac{\mu \alpha_{k,z}^{t-1} - 1 + \alpha_{k,z}^t \sum_d (\phi(N_{d,k,z}^t + \alpha_{k,z}^t) - \phi(\alpha_{k,z}^t))}{\mu + \sum_d (\phi(N_{d,k}^t + \sum_z \alpha_{k,z}^t) - \phi(\sum_z \alpha_{k,z}^t))} \quad (12)$$

---

## ALGORITHM: Gibbs sampling procedure for DST MODEL

---

**Input:** Number of topics  $Z$ , number of sentiment labels  $L$ , number of time slices  $S$ , word prior polarity transformation matrix  $\lambda$ , epoch  $t \in \{1, 2, \dots, \text{maxEpochs}\}$ , a stream of documents  $D_t = \{d_{t_1}, d_{t_2}, \dots, d_{t_D}\}$ ;

**Output:** DST model;

---

Sort documents according to their time stamps;

**FOR**  $t = 1$  to  $\text{maxEpochs}$  **do**

**IF**  $t == 1$  **THEN**

        Set  $\beta_{l,z}^t = \lambda \times 0.01$ ;

**End**

**ELSE**  $\delta_{l,z}^t, \sigma_{l,z,s}^t$

        Set  $\delta_{l,z}^t = \delta_{l,z}^{t-1}$ ;

        Set  $\sigma_{l,z,s}^t = 1/S$ ;

        Set  $\beta_{l,z}^t = \sum \delta_{l,z}^t * \sigma_{l,z,s}^t$ ;

**END**

        Set  $\gamma^t = (0.05 \times \text{Average document length})/L$ ;

        Set  $\alpha^t = (0.05 \times \text{Average document length})/(L \times Z)$ ;

        Initialize  $\pi^t, \theta^t, \phi^t$ , and all count variables;

        Initialize sentiment label and topic

        assignment randomly for all word tokens in  $D_t$ ;

**FOR**  $i = 1$  to  $\text{max Gibbs Sampling Iterations}$  **DO**

$[\pi^t, \theta^t, \phi^t, l_t, z_t] = \text{GibbsSampling}(D_t, \alpha^t,$

$\beta_{l,z}^t, \gamma^t)$ ;

**FOR every 40 Gibbs sampling iterations DO**

            Update  $\alpha^t$  using Equation (12) ;

            Update  $\gamma^t$  using Equation (13) ;

            Update  $\delta_{l,z}^t$  using Equation (7) ;

            Set  $\beta_{l,z}^t = \sum \delta_{l,z}^t * \sigma_{l,z,s}^t$ ;

**END**

**FOR every 200 Gibbs sampling iterations DO**

            Update  $\pi^t, \theta^t, \phi^t$  with the new sampling results;

**END**

**END**

    Update  $\beta_{l,z}^t$ , using Equation (10) ;

**END**

---

**Figure 3: Online Stochastic EM Algorithm**

$$(\gamma^t)^{new} \leftarrow \frac{\xi \gamma^{t-1} - 1 + \gamma^t \sum_d (\phi(N_{d,1}^t + \gamma^t) - \phi(\gamma^t))}{\xi + \sum_d (\phi(N_d^t + \sum_i \gamma_i^t) - \phi(\sum_i \gamma_i^t))} \quad (13)$$

The detail description of the online stochastic EM algorithm for the DST model is given in Algorithm 1.

## 4 Experiment

To evaluate the proposed solution, we conducted experiments using a real dataset collected from Sina Weibo.

### 4.1 Dataset Preparation

We implement a Sina Weibo crawler using the APIs provided by Sina Corp. (<http://open.weibo.com/>) to collect the experiment dataset.

From January 23, 2013 to February 23, 2013, the data crawling process continued for 31 days. Table 1 illustrates the statistics of the dataset. Considered that many public

events happed in the period of December 19, 2011 to December 26, 2011, we select the Weibo data generated in that time period as our test dataset.

**Table 2: Statistics of Collected Dataset**

Statistical Measure	Value
Number of Users	531,935
Number of Posts	25,838,961
Number of Comments	14,498,416

**4.2 DST Model V.S Existing Models**

In order to evaluate the effectiveness of our proposed model in topic detection, we compared the performance of DST with LDA and DST with metric of perplexity.

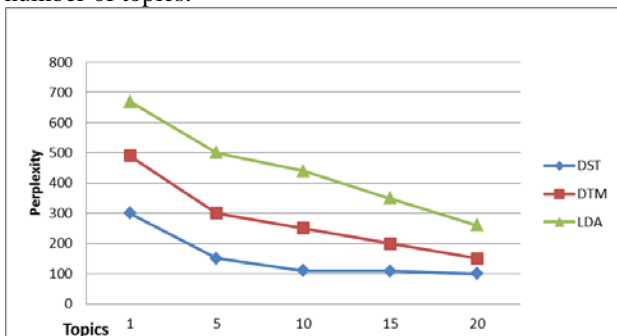
Perplexity is a metric which was mostly used in language modelling. It could measure a language model’s prediction ability when applied on an uncharted dataset. Perplexity is defined as the reciprocal geometric mean of the likelihood of a test corpus given a trained model’s Markov Chain state  $M$ . In this paper, we adopted the definition of perplexity as below [He, Y., Lin, C., Gao, W., and Wong, K.-F]:

$$Perplexity = P(D_t | M) = \exp \left( - \frac{\sum_{d=1}^{D_t} \log p(\tilde{w}_d' | M)}{\sum_{d=1}^{D_t} \tilde{N}_d'} \right) \quad (14)$$

Here, this metric means the per-word predictive perplexity of the unseen test set  $D_t$  at each epoch  $t$  based on the previously trained model  $M$ .

In the experiments we studied the influence of the topic number settings on the DST model performance. With the number of time slices fixed at  $S = 4$ , we vary the topic number  $T \in \{1, 5, 10, 15, 20\}$ .

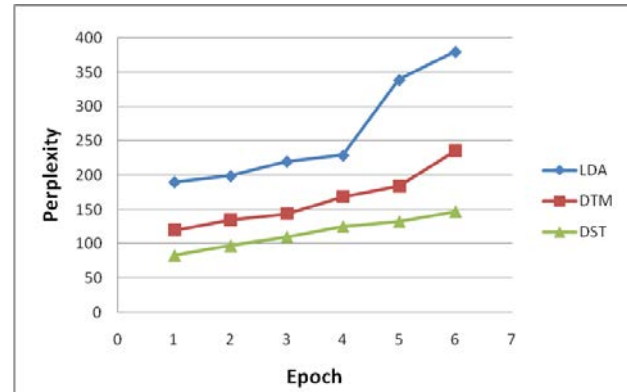
Figure 4 shows the average per-word perplexity over epochs with different number of topics. DST has lower perplexities than all the other models with the increased number of topics.



**Figure 4: Perplexity VS. umber of topics**

The average perplexity for each epoch with the number of time slices set to 4 and the number of topics set to 15 for the DST models is shown in Figure 5. In addition, we also plot the perplexity results of LDA and DTM. LDA only uses the data in the previous epoch for training and hence it does not model dynamics while DTM uses all past data for model learning. We set the number of topics to 15 for both

DTM and DST. For LDA, the number of topics was set to 3 corresponding to positive, negative, and neutral sentiment labels.



**Figure 5: Perplexity vs. Number of epochs**

Figure 5 shows that LDA has the highest perplexity values followed by DTM and DST. The perplexity gap between DTM and the DST models increases with the increasing number of epochs. The variants of DST models have quite similar perplexities.

**5 Conclusion**

In this paper, we proposed a dynamic Sentiment Topic (DST) model, which could models the topics and its corresponding sentiment polarities from huge amount of user generated content data. By taking the sentiment effects into account, DST could properly reflect the impact that sentiment exerts on topic detection. Also, it enhances LDA model by introducing the dynamic characteristics, such that it is more appropriate for those social media data with strong real-time characteristics. Experiments results on a large real dataset show that compared with existing models, DST performances better in terms of perplexity. In future work, we may work on reducing the temporal and spatial overhead of this model.

**Acknowledgements**

This research was partially supported by National Natural Science Foundation of China under Grant No. 61402045, Specialized Research Fund for the Doctoral Program of Higher Education under Grant No. 20130005110011, and Fundamental Research Funds for the Central Universities under Grant No. 2014RC1301.

**References**

Pang, B., Lee, L.(2008): Opinion mining and sentiment analysis. *Found. Trends Inf. Retr.* **2**(2): 1–135.  
 Aue, A., Gamon, M. (2005): Customizing sentiment classifiers to new domains: a case study. *Proc. of RANLP*.  
 Turney, P, D. (2001): Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews. *Proc. of ACL*, 417–424.  
 Lacoste, S, Sha, F and Jordan, M. (2008): DiscLDA: Discriminative learning for dimensionality reduction and classification. *Proc. of NIPS*.  
 Ramage, D., Hall, D., Nallapati, R., and Manning, C. (2009): Labeled LDA: A supervised topic model for credit attribution in multi-labeled corpora. *Proc. of EMNLP:248-256*.

- Wallach, H., Mimno, D., McCallum, A. (2009): Rethinking LDA: Why priors matter. *Proc. of Topic Models: Text and Beyond Workshop in Neural Information Processing Systems Conference*.
- Hofmann, T. (1999): *Probabilistic latent semantic indexing*. *Proc. of SIGIR*:50-57.
- David, M., John, D.,(2006): Dynamic Topic Models. *Proc. of the 23rd International Conference on Machine Learning*.
- Xueru, W.,McCallum, A.(2005): Topics over Time: A Non-Markov Continuous-Time Model of Topical Trends. *Proc. of ACM SIGKDD*.
- He, Y. and Lin, C. (2012): Online sentiment and topic dynamics tracking over the streaming data. *Proc. of the ASE/IEEE International Conference on Social Computing (SocialCom'12)*.
- He, Y., Lin, C., Gao,W., and Wong, K.-F.(2013): Dynamic joint sentiment-topic model. *ACM Trans. Intell. Syst.Technol.* 5, 1, Article 6.