

# How Shall We Catch People’s Concerns in Micro-blogging?

Heng Gao, Qiudan Li, Hongyun Bao, Shuangyong Song  
 State Key Laboratory of Management and Control for Complex Systems  
 Institute of Automation, Chinese Academy of Sciences , Beijing 100190 , China  
 {heng.gao, qiudan.li, hongyun.bao, shuangyong.song}@ia.ac.cn

## ABSTRACT

In micro-blogging, people talk about their daily life and change minds freely, thus by mining people’s interest in micro-blogging, we will easily perceive the pulse of society. In this paper, we catch what people are caring about in their daily life by discovering meaningful communities based on probabilistic factor model (PFM). The proposed solution identifies people’s interest from their friendship and content information. Therefore, it reveals the behaviors of people in micro-blogging naturally. Experimental results verify the effectiveness of the proposed model and show people’s social life vividly.

**Categories and Subject Descriptors:** H.2.8 [Database Management]: Database Application – *Data mining*, H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval – *Information filtering*

## General Terms

Algorithms, Design, Experimentation.

## Keywords

Micro-blogging, probabilistic factor model, user communities

## 1. INTRODUCTION

Nowadays, we’re seeing micro-blogging growing rapidly, and it is playing important roles in our social life [1, 2]. In China, Sina-weibo, a representative micro-blogging system that is similar to Twitter, is getting popular among people of all ages. People are fond of using it to know what’s happening around them: they talk about their daily life, follow their interested people, forward their interested tweets, and so on. All of the above behaviors reflect people’s interest from different angles. But how do we know exactly what people are caring about with the popular social media in a unified way? In general, people may just care about a certain number of topics and want to be involved in the community where they share similar interest. Thus, it’s necessary to mine such valid meaningful communities to characterize people’s behavior in micro-blogging. Meanwhile, people can utilize it to find their most concerning topics and most interest-related friends. In this paper, we propose a unified solution to help people find such valuable communities.

Typically, in micro-blogging, people with similar interest may post on the similar topics. Besides, if user A is interested in user B, A may follow B, then A will easily capture what tweets B has posted, by commenting on and retweeting B’s tweets, A will get convenient access to communicating with B. Existing studies have mainly focused on finding community structure based on user-content relationship or users’ friendship with each other, respectively [3, 4]. Although these approaches are beneficial, it may bring better performance when integrating all these information together. As probabilistic factor model provides a

natural way to combine multiple resources together, in this paper we propose a community mining method based on probabilistic factor model to discover meaningful communities. In the following of the paper, we will illustrate our method in detail, along with the interesting community mining results.

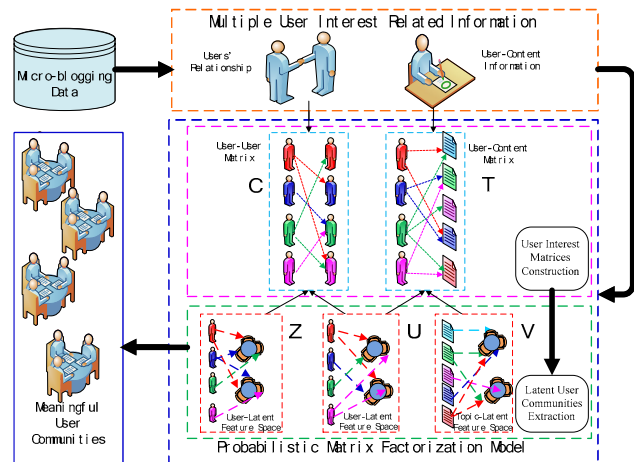


Fig.1 The framework of mining meaningful user communities

## 2. METHOD

The framework of our model is shown in Fig. 1. We define a meaningful user community as a group containing people who form close friendship during their participation of particular topics. Our community discovery model firstly constructs users’ friendship matrix and user-content profile to represent users’ interest from different angles, then, the latent community dimensions are extracted simultaneously from the resource above to obtain users’ interest distribution, via the extracted matrices, we will easily catch people’s concerns with a macro-perspective.

**User Interest Matrices Construction** In Sina-weibo, posting tweets on their interested topics and following people who they are interested in are two main manners for users to represent their concerns, which reflect their interest. Hence, we refer to the user-topic matrix  $T$  and users’ mutual social network matrix  $C$  respectively, as users’ different interest representation. Noticing that  $T \in R^{m \times n}$ ,  $C \in R^{m \times m}$ , where  $m, n$  denote the amount of users and topics, respectively. Every element  $T_{ij}$  in matrix  $T$  represents the frequency user  $i$  has posted on topic  $j$ . Meanwhile, in matrix  $C$  we set  $C_{pq} = 1$  when user  $p$  has followed user  $q$ , otherwise,  $C_{pq} = 0$ . In this way, we formulate the multiple resources of Sina-weibo in order to mine meaningful user communities effectively.

**Latent User Communities Extraction** Based on the user interest matrices  $C$  and  $T$ , we employ the probabilistic factor model proposed in [5], to capture the underlying close relationship between users and the latent communities. As illustrated in Fig.1, we need to find the best low-dimensional matrices  $U, V$  and  $Z$ , such that  $C \approx UZ^T$ ,  $T \approx UV^T$ , where  $U, Z$  are the low-dimensional

Copyright is held by the author/owner(s).  
 WWW 2012 Companion, April 16–20, 2012, Lyon, France.  
 ACM 978-1-4503-1230-1/12/04.

user latent feature space with matrix size  $m \times d$ , and  $V$  denotes the low-dimensional topic latent feature space with matrix size  $n \times d$ , in which  $d$  represents the community dimensionality. By introducing two Gamma parameters  $\alpha$  and  $\beta$ , we constrain the elements in the extracted latent feature space non-negative, which will make our community mining results well-grounded, then we turn the probabilistic factor model into the optimization problem of maximizing the following generalized objective function:

$$L = \alpha C | X \alpha T | Y \alpha U | \alpha, \beta \alpha V | \alpha, \beta \alpha Z | \alpha, \beta \quad (1)$$

where  $X = UZ^T$ ,  $Y = UV^T$ , in this way, we can find the best low-dimensional matrices  $U$ ,  $V$  and  $Z$  with the partial derivative of the objective function (1).

**Meaningful User Communities Discovery** After we get the extraction results, it's easy to mine meaningful user communities: Each element  $u_{ik}$  ( $k = 1, \dots, d$ ) in  $U$  encodes the preference of user  $i$  to latent community  $k$ , and each  $v_{jk}$  in  $V$  can be interpreted as the affinity of topic  $j$  to the latent community  $k$ . The advantages of our model lie in: (i) It unifies the user-following relationship and user-content information simultaneously, which can help us find the meaningful communities effectively. (ii) We can regulate the weight of matrices  $C$  and  $T$  conveniently, thus helping us balance the biased impact of noises on the latent community discovery. (iii) By leading in the Gamma distribution parameters  $\alpha$  and  $\beta$ , we make elements in the user community matrices non-negative, which will make our experimental results more explainable and meaningful.

### 3. EXPERIMENTS

#### 3.1 Dataset and Parameter Settings

To evaluate the performance of our community mining model, we build dataset including user-following relationship and user-content information with time interval of 16 days from October 29<sup>th</sup>, 2011 to November 13<sup>th</sup>, 2011. After removal of users who post less than one tweet per day, we get 1879 users and 1640 topics. The influence parameter  $c$  was empirically set to be 0.5 to evenly weight the matrix  $C$  and matrix  $T$ . We also tune the parameter  $d$  which denotes the latent community dimension size from 4 to 30, and find the best performance at  $d=12$  eventually. As to the Gamma distribution parameters  $\alpha$  and  $\beta$ , we set them as the best performance value of 10, 0.02 respectively.

#### 3.2 Results and Discussions

To demonstrate the validity of the proposed model, non-negative matrix factorization (NMF) method is conducted as the baseline, which only considers the user-content information. Looking into the community mining results shown in Table 1, we delightedly find that our model reveals more interesting phenomena of people's concerns in micro-blogging: First of all, our method obtains more personalized clustering results, as shown in part I, both models mine the user community which concerns the topic of South Korea Street Shot very much, besides, our model successfully takes in people who focus on South Korean clothing that often appears on the South Korea Street Shot. Secondly, our method merges user communities whose topics of concern are very close, as shown in part II, two separate user communities concerning similar leisure topics show up through NMF approach, however, considering users' close friendship, our method

reasonably merges them. Thirdly, our method mines the meaningful community which can't be mined with NMF relying on user-content information, seeing Part III, our method mines people who care the Chinese drama *The Rhino in Love* very much, which is very popular but doesn't show up through NMF approach. Thus, by integrating users' interest of different aspects, we find more meaningful user communities, which will help us catch people's concerns better in micro-blogging.

**Table 1. Community mining results in micro-blogging**

Interesting Phenomena	People's Concerns	
	NMF Based Model	Our Integrating Model
Part I	<i>South Korea Street Shot</i>	<i>South Korea Street Shot &amp; South Korean Clothing</i>
Part II	<i>Hey Gossip   Today's Fun</i>	<i>Hey Gossip &amp; Today's Fun</i>
Part III	Null	<i>The Rhino in Love</i>

Besides, we use the mean value of soft modularity metric  $Q_s$  and users' cosine similarity based on content information to evaluate our community mining results. The higher mean value  $\bar{\mu}$  means the closer friendship and closer topic interest of people falling into the same community. Table 2 shows the performance of our method compared to the NMF based model.

**Table 2. Performance evaluation on community mining results**

Evaluation Metric	NMF Based Model	Our Unified Model
$\bar{\mu}$	0.1211	<b>0.2718</b>

From Table 2, we can see that our unified model outperforms the baseline NMF-based model, which only considers the user-content information. The results reveal the validity of catching people's concerns from their multi-interest distribution. It is due to our model could contain more valuable latent clustering information than the baseline.

### 4. CONCLUSION

In this paper, we propose a PFM-based community discovery method by mining people's interest from their friendship network and content information. Preliminary experiments have proved the validity and effectiveness of our approach. The interesting clustering results will help us understand people's concerns and feel the pulse of our society easily. In the future, we would like to make advertising recommendations based on our existing work in micro-blogging, which will be very promising and interesting. This research is supported by the NNSFC project 61172106 and the BJNSF project 4112062.

### 5. REFERENCES

- [1] Singh, V.K., Gao, M., Jain, R. 2010. Situation Detection and Control Using Spatiotemporal Analysis of Microblogs. In WWW'10, pp. 1181-1182.
- [2] Wu, S., Hofman, J.M., Mason, W.A., Watts, D.J. 2011. Who Says What to Whom on Twitter. In WWW'11, pp. 705-714.
- [3] Lin, Y., Sun, J., Castro, P., Konuru, R., Sundaram, H., Kelliher, A. 2009. MetaFac: Community Discovery via Relational Hypergraph Factorization. In SIGKDD '09, pp. 527-536.
- [4] Psorakis, I., Roberts, S., Ebdon, M. 2011. Overlapping Community Detection using Bayesian Nonnegative Matrix Factorization. Phys. Rev. E 83, 066114 (2011).
- [5] Ma, H., Liu, C., King, I., Lyu, M.R. 2011. Probabilistic Factor Models for Web Site Recommendation. In SIGIR'11, pp.265-274.