

# Time-aware Topic Recommendation Based on Micro-blogs

Huizhi Liang\*, Yue Xu<sup>□</sup>, Dian Tjondronegoro<sup>□</sup>, Peter Christen\*

\*Research School of Computer Science, The Australian National University

<sup>□</sup>Faculty of Science and Technology, Queensland University of Technology

\*{huizhi.liang, peter.christen}@anu.edu.au, <sup>□</sup>{yue.xu, dian}@qut.edu.au

## ABSTRACT

Topic recommendation can help users deal with the information overload issue in micro-blogging communities. This paper proposes to use the implicit information network formed by the multiple relationships among users, topics and micro-blogs, and the temporal information of micro-blogs to find semantically and temporally relevant topics of each topic, and to profile users' time-drifting topic interests. The Content based, Nearest Neighborhood based and Matrix Factorization models are used to make personalized recommendations. The effectiveness of the proposed approaches is demonstrated in the experiments conducted on a real world dataset that collected from Twitter.com.

## Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval-Information Filtering; H.5.3 [Information Interfaces and Presentation]: Group and Organization Interfaces-Collaborative computing

## General Terms

Algorithms, Experimentation

## Keywords

Temporal dynamics, Topic Recommendation, Micro-blogs, Collaborative Filtering, Personalization, Web 2.0

## 1. INTRODUCTION

Micro-blogs is one kind of popular Web 2.0 information. Rather than a pure social network like Facebook, a micro-blogging platform is regarded as an information network [9], where many people use it for information purpose [9]. With the rapid growth of user numbers, there are a large number of topics emerging every day. They not only include a small number of hot/stream topics, but also a large number of less popular ones. To help users solve the information overload issue, it is important to recommend personally interesting topics to users. Although the latest version of Twitter has embedded the function of recommending topics (e.g., hashtags, popular keywords) to users, the academic research of making personalized topic recommendations based on micro-blogs has attracted less attention so far.

Recently, the social tie or social interaction information of micro-blogs have been used to discover communities [14], make recommendations [5]. However, more recent research findings [9] suggest that social tie information may not be very helpful for users who use micro-blogging environment for information

purpose, since users with similar topic interests may be not explicitly connected, and weak ties often can provide access to novel information [9]. Thus, making better use of the content information of micro-blogs is crucial for topic recommendations. On the other hand, micro-blogs contain implicit information network formed by the multiple relationships among users, micro-blogs, and topics. These relationships can be used to find content relevant topics, which is ignored by other approaches.

Another unique feature of micro-blogs is that the topics are temporally associated with each other. The temporal information of topics can help to find relevant topics, which should be considered. Moreover, the fact that the topics of micro-blogging communities change quickly with time [6] makes it necessary to recommend topics that are not only topically appropriate, but also have been talked or published in the recent past. Thus, how to capture the temporal dynamics especially recency information in micro-blogs and profile users' time-sensitive topic interests is very important. In this paper, we propose to make use of the temporal information and the multiple relationships among users, topics and micro-blogs to make personalized topic recommendations.

## 2. RELATED WORK

The research of recommender systems in micro-blogging communities is mainly focusing on recommending news [10], URLs [5], and users to follow [3]. Although some work like [10] considered the temporal dynamics of micro-blogs, how to incorporate the recency information to find temporally associated topics still needs to be explored. Time-aware latent topical models [8] can be used to find the latent topics in micro-blogs. However, as latent topics are usually broad or abstract, the recommendations of specific topics, such as hashtags and keywords, are more applicable in micro-blogging communities. With a list of recommended specific topics such as hashtags and keywords, users can read those micro-blogs that are relevant to the recommended topics and publish micro-blogs with these topics to participate in discussions or conversations directly.

Temporal dynamics in recommender systems are of great importance [1]. For example, Koren [1] modeled the time factors for each user in a factorization model. Xiang et. al. [12] proposed a graph based approach to hybrid users' short- and long-term preferences. However, these approaches were based on users' explicit or implicit ratings and did not consider the content information of items. Moreover, the patterns of temporal variations of micro-blogs [2] are different with the items such as movies, research papers that were used in these approaches. Efron et. al. [7] proposed temporal models to rank recent tweets. Different with the problem of ranking tweets [7], this paper focuses on the recommendation of topics.

## 3. PROBLEM DEFINITION

We define the key concepts that are used in this paper:

- **Users:**  $U = \{u_1, u_2, \dots, u_{|U|}\}$  contains all users in a micro-blogging community who have published micro-blogs.

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.

CIKM'12, October 29–November 2, 2012, Maui, HI, USA.

Copyright 2012 ACM 978-1-4503-1156-4/12/10...\$15.00.

- **Micro-blogs (e.g., tweets):**  $S = \{s_1, s_2, \dots, s_{|S|}\}$  contains all micro-blog messages generated by users in  $U$ . A micro-blog may contain hashtags, keywords, URL links and others [5].
- **Topics:**  $C = \{c_1, c_2, \dots, c_{|C|}\}$  contains all topics of micro-blogs generated by users in  $U$ . **Hashtags** are given by users to label the topics of their micro-blogs or to participate in group discussions/conversations, denoted as  $H = \{h_1, h_2, \dots, h_{|H|}\}$ . A hashtag  $h_i$  is a keyword (i.e., term)  $k_i$  preceded by a '#' symbol in a micro-blog,  $h_i = \#k_i$ . To differentiate with those keywords that are not used as hashtags, a keyword that has been used as a hashtag by at least one user is defined as a tag term. The **tag term** set  $P$  contains all tag terms obtained from hashtags  $H$ ,  $P = \{k_i | \#k_i \in H\}$ . The **keyword** set that contains all the keywords extracted from micro-blogs  $S$  is defined as  $\mathcal{K}$ , and  $\mathcal{K} = \{k_1, k_2, \dots, k_{|\mathcal{K}|}\}$ . The topics of micro-blogs can be represented by hashtags  $H$ , tag terms  $P$  or keywords  $\mathcal{K}$ .
- **Topic assigning:** is to assign a topic to an item. Similar to social tagging [13][4], hashtagging is a kind of explicit topic assigning behavior, as a user places a hash symbol before a term in a micro-blog to label one topic of this micro-blog. Besides explicit topic assigning behavior, the topics contained in a micro-blog can be regarded as a kind of implicit topic assigning behavior. To be more general, the topic assigning behavior is defined as  $e: U \times C \times S \rightarrow \{0,1\}$ . If a micro-blog  $s_k$  contributed by user  $u_i$  contains or belong to topic  $c_j$ , then  $e(u_i, c_j, s_k) = 1$ , otherwise,  $e(u_i, c_j, s_k) = 0$ .

Let  $u_i \in U$  be a target user,  $C_{u_i}$  be the topic set that  $u_i$  already has,  $t_{c_j}$  be the latest time stamp of  $c_j$ ,  $T_{u_i}$  be the correspondent time stamp set of  $C_{u_i}$ ,  $l(T_{u_i})$  be the latest time stamp of  $T_{u_i}$ ,  $\check{C}_{u_i}$  be the candidate topic set that is unknown to  $u_i$  and is more recent than  $l(T_{u_i})$ . Let  $c_k \in \check{C}_{u_i}$  be a candidate topic,  $\mathcal{A}(u_i, c_k)$  be the prediction score of how much  $u_i$  would be interested in  $c_k$ . The problem of topic recommendation is defined as generating a set of ordered topics  $c_l, \dots, c_m \in \check{C}_{u_i}$  to  $u_i$ , where  $\mathcal{A}(u_i, c_l) \geq \dots \geq \mathcal{A}(u_i, c_m)$  and  $t_{c_l} > l(T_{u_i}), \dots, t_{c_m} > l(T_{u_i})$ .

## 4. THE PROPOSED APPROACHES

### 4.1 The Relationship Modeling

In a micro-blogging community, User-Microblog is the basic relationship. The introducing of explicit (i.e., hashtagging behaviour) or implicit topic assigning behaviours form multiple relationships among users, micro-blogs and topics.

- **User-Microblog relationship:** This includes User-Microblog Mapping and Microblog-User Mapping. Microblog-User Mapping is a one-to-one mapping.
  - 1) *User-Microblog Mapping*  $S_{u_i}: U \rightarrow 2^S, S_{u_i} = \{s_k | \exists c_j \in C, \forall s_k \in S, e(u_i, c_j, s_k) = 1\}$ . It maps a user to his/her generated micro-blogs.
- **User-Topic relationship:** This records each user's topics and the user group of each topic. It includes User-Topic Mapping and Topic-User Mapping.
  - 2) *User-Topic Mapping*  $C_{u_i}: U \rightarrow 2^C, C_{u_i} = \{c_j | \exists s_k \in S, \forall c_j \in C, e(u_i, c_j, s_k) = 1\}$ . It maps a user to a set of topics that are used by the user.
  - 3) *Topic-User Mapping*  $U_{c_j}: C \rightarrow 2^U, U_{c_j} = \{u_i | \exists s_k \in S, \forall u_i \in U, e(u_i, c_j, s_k) = 1\}$ . It maps a topic to a set of users who have used this topic.
- **Topic-Microblog relationship:** This records each micro-blog's topics and the aggregated micro-blogs of each topic.

4) *Microblog-Topic Mapping*  $C_{s_k}: S \rightarrow 2^C, C_{s_k} = \{c_j | \exists u_i \in U, \forall c_j \in C, e(u_i, c_j, s_k) = 1\}$ . It maps a  $s_k$  to a set of topics.

5) *Topic-Microblog Mapping*  $S_{c_j}: C \rightarrow 2^S, S_{c_j} = \{s_k | \exists u_i \in U, \forall s_k \in S, e(u_i, c_j, s_k) = 1\}$ . It maps a topic to a set of micro-blogs that contains or belong to this topic.

- **User-Topic-Microblog relationship:** This records each user's personal topic assigning relationships.

6) *(User×Topic)-Microblog Mapping*  $S_{u_i, c_j}: U \times C \rightarrow 2^S, S_{u_i, c_j} = \{s_k | \forall s_k \in S, e(u_i, c_j, s_k) = 1\}$ . It maps a user-topic pair to a set of micro-blogs.

The multiple relationships of micro-blogs can be used to find each user's topic interests and the related topics of each topic, which will be discussed in the following sub sections.

### 4.2 Topic Representation

The process of determining the time-aware related topics of each topic and representing each topic with a set of content relevant and temporally associated topics is called topic representation.

**[Definition 1] (Topic Representation):** represents the time-aware content relevant topics of a given topic  $c_k \in C$  with respect to all users in  $U$ . Let  $w_{k,y}^{c(t)}$  denote the weight of how much topic  $c_k$  is relevant to topic  $c_y \in C$ . The relationship between a topic and a set of topics in time period  $t$  can be defined as the mapping  $\mathcal{R}^{C(t)}: C \rightarrow 2^{C \times [0,1]}$ , such that  $\mathcal{R}^{C(t)}(c_k) = \{(c_y, w_{k,y}^{c(t)}) | c_y \in C\}$ .  $\mathcal{R}^{C(t)}(c_k)$  is called the topic representation of topic  $c_k$ .

For a given topic  $c_k \in C$ , based on the Topic-User Mapping, we can get the user set of this topic denoted as  $U_{c_k}$ . For each user  $u_i \in U_{c_k}$ , a set of micro-blogs containing topic  $c_k$  (i.e.,  $S_{u_i, c_k}$ ) can be obtained based on the personal topic assigning relationship (User×Topic)-Microblog Mapping. In the viewpoint of this user, the topics of these micro-blogs are closely related. Let  $r_{u_i, c_k}^t(c_y)$  denote the time-aware relevance weight of a given topic  $c_k$  and another topic  $c_y \in C$  in terms of user  $u_i$ . It can be estimated based on the average relevance weight of  $c_y$  to every  $s_j \in S_{u_i, c_k}$ . Let  $f_{y,j}$  denote the relevance weight of  $c_y$  to  $s_j$ .  $f_{y,j} = \frac{n_{y,j}}{\sum_{c_z \in C} n_{z,j}}$ , where  $n_{y,j}$  is the number of occurrence of  $c_y$  in  $s_j$ . The exponential decay that shows strong performance in recency ranking of micro-blogs [7] is adopted to measure the recency decay. It is formed by the function  $\lambda e^{-\lambda \Delta t}$ , where  $\lambda$  is the decay rate and  $\Delta t$  is the time in hours (or in days) that has elapsed from the current time. Let  $T_{u_i, c_k}$  be the time stamp set of  $S_{u_i, c_k}$ ,  $l(T_{u_i, c_k})$  be the latest time stamp of  $T_{u_i, c_k}$ ,  $t^*$  be the given latest current time stamp, the recency weight of  $u_i$  for  $c_k$  can be calculated as  $rec(T_{u_i, c_k}) = \lambda e^{-\lambda \cdot |t^* - l(T_{u_i, c_k})|}$ . Let  $|S_{u_i, c_k}|$  be the number of micro-blogs in  $S_{u_i, c_k}$ ,  $r_{u_i, c_k}^t(c_y)$  is calculated as:

$$r_{u_i, c_k}^t(c_y) = \sum_{s_j \in S_{u_i, c_k}} \frac{f_{y,j}}{|S_{u_i, c_k}|} \cdot rec(T_{u_i, c_k}) \quad (1)$$

The overall relevance of two topics  $c_k$  and  $c_y$  can be measured through calculating the sum of the relevance weight of  $c_k$  and  $c_y$  for all the users of  $U_{c_k}$ . However, the importance of one user  $u_i \in U_{c_k}$  for the topic representation of  $c_k$  may be different. Assuming each micro-blog is equally important, the more micro-blogs of topic  $c_k$  are contributed by user  $u_i$ , the more important  $u_i$  is for the topic representation of  $c_k$ . Let  $\mathcal{P}_{c_k}(u_i)$  denote the importance weight of  $u_i$  to the topic representation of  $c_k$ ,

$\mathcal{P}_{c_k}(u_i) = \frac{|S_{u_i, c_k}|}{|S_{c_k}|}$ . Moreover, similar to the *idf* weighting approach, the popularity of  $c_y$  in all topic representations should be considered. By considering the importance weight of  $u_i$  and the popularity of  $c_y$ ,  $w_{k,y}^{c(t)}$  can be calculated as:

$$w_{k,y}^{c(t)} = \sum_{u_i \in U} \mathcal{P}_{c_k}(u_i) \cdot \tau_{u_i, c_k}^t(c_y) \cdot itf(c_y) \quad (2)$$

Where  $itf(c_y)$  is the inverse topic frequency of  $c_y$ ,  $itf(c_y) = 1 / \log(e + |N_{c_y}|)$ , where  $e$  is a constant approximately equal to 2.72,  $|N_{c_y}|$  is the number of topics that have been described by  $c_y$ , and  $0 < itf(c_y) \leq 1$ . The mapping  $\mathcal{R}^{c(t)}(c_k)$  can be viewed as vector  $\mathcal{R}^{c(t)}(c_k) = \langle w_{k,1}^{c(t)}, \dots, w_{k,|C|}^{c(t)} \rangle$  for topics  $\langle c_1, \dots, c_{|C|} \rangle$ .

### 4.3 User Profiling

User profiles are used to describe users' interests and preferences information. The process of finding time-aware topic preferences of each user is called user representation. It is defined as below:

**[Definition 2] (User Representation):** represents the time-aware topic preferences of each user. Let  $w_{i,y}^{u(t)}$  denote the weight of how much the user  $u_i$  is interested in topic  $c_y \in C$ . The relationship between a user and a set of topics in time period  $t$  can be defined as the mapping  $\mathcal{R}^{u(t)}: U \rightarrow 2^{C \times [0,1]}$ , such that  $\mathcal{R}^{u(t)}(u_i) = \{(c_y, w_{i,y}^{u(t)}) \mid c_y \in C\}$ .  $\mathcal{R}^{u(t)}$  is called the user representation of  $u_i$ .

To calculate  $w_{i,y}^{u(t)}$ , we first calculate how much the user is interested in  $c_x \in C_{u_i}$ . Since the number of micro-blogs that contain  $c_x$  indicates how strong this user is interested in  $c_x$ , we use the ratio between the number of micro-blogs that contain  $c_x$  and generated by  $u_i$ , and the total number of micro-blogs generated by user  $u_i$ , to measure the preference weight of  $u_i$  to  $c_x$ , denoted as  $\mathcal{P}_{u_i}(c_x)$ .  $\mathcal{P}_{u_i}(c_x) = \frac{|S_{u_i, c_x}|}{|S_{u_i}|}$ . The higher the value of  $\mathcal{P}_{u_i}(c_x)$ , the more the user is interested in  $c_x$ . Based on Equation 1, we can get the time-aware relevance weight  $r_{u_i, c_x}^t(c_y)$  of  $c_x$  and  $c_y$  in terms of  $u_i$ . As discussed in Section 4.2,  $r_{u_i, c_x}^t(c_y)$  considers the recency weight of each topic  $c_x$  for the user representation of  $u_i$ . The older the time stamp of the micro-blogs generated by user  $u_i$  with the topic  $c_x$  are, the less important  $c_x$  is for the user representation of  $u_i$ . Thus, we can measure each user  $u_i$ 's preferences to the topic  $c_y \in C$  through calculating the product of  $\mathcal{P}_{u_i}(c_x)$  and  $r_{u_i, c_x}^t(c_y)$ . Considering the inverse topic frequency of each topic, the weight  $w_{k,y}^{c(t)}$  can be calculated as:

$$w_{i,y}^{u(t)} = \sum_{c_x \in C} \mathcal{P}_{u_i}(c_x) \cdot r_{u_i, c_x}^t(c_y) \cdot iuf(c_y) \quad (3)$$

Where  $iuf(c_y)$  is the inverse user frequency of  $c_y$ ,  $iuf(c_y) = 1 / \log(e + |U_{c_y}|)$ ,  $|U_{c_y}|$  is the number of users that have been described by  $c_y$ ,  $0 < iuf(c_y) \leq 1$ .  $\mathcal{R}^{u(t)}(u_i)$  can be viewed as vector  $\mathcal{R}^{u(t)}(u_i) = \langle w_{i,1}^{u(t)}, \dots, w_{i,|C|}^{u(t)} \rangle$  for topics  $\langle c_1, \dots, c_{|C|} \rangle$ .

## 4.4 Personalized Recommendation

In this section, based on the user and topic representations, three kinds of recommendation approaches are proposed.

### 4.4.1 Content based Model

The content based approach is popularly used to recommend items that have similar contents to each target user's topic

interests. The content similarity between  $u_i$  and candidate topic  $c_k \in \tilde{C}_{u_i}$  can be calculated by the similarity of vector  $\mathcal{R}^{u(t)}(u_i)$  and  $\mathcal{R}^{c(t)}(c_k)$ . This paper uses the Cosine similarity to measure the content matching value of  $u_i$  and  $c_k$ . Similarly, the recency of  $c_k$  should be considered,  $rec(T_{c_k}) = \lambda e^{-\lambda \cdot |t^* - l(T_{c_k})|}$ , where  $T_{c_k}$  is the time stamp set of  $c_k$ , and  $l(T_{c_k})$  is the latest time stamp of  $T_{c_k}$ . The content matching between  $u_i$  and  $c_k$  is defined as:

$$sim_{u,c}^t(u_i, c_k) = cosine(\mathcal{R}^{u(t)}(u_i), \mathcal{R}^{c(t)}(c_k)) \cdot rec(T_{c_k}) \quad (4)$$

The prediction score that measures how much  $u_i$  will be interested in  $c_k$  can be calculated based on their content matching value.

$$\mathcal{A}_c(u_i, c_k) = sim_{u,c}^t(u_i, c_k) \quad (5)$$

### 4.4.2 User based K-Nearest-Neighborhood Model

Typically, the similarity of two users  $u_i$  and  $u_j$  can be measured by the similarity of their user profiles (i.e., user representations). In a micro-blogging community, the discussion topics change with time quickly. For a given user  $u_i$ , the active time period of each peer user of  $u_i$  should not be ignored. Let  $rec(T_{u_j})$  denote the recency weight of each peer user  $u_j$ ,  $rec(T_{u_j}) = \lambda e^{-\lambda \cdot |t^* - l(T_{u_j})|}$ , where  $T_{u_j}$  is the time stamp set of the micro-blogs of  $u_j$ , and  $l(T_{u_j})$  is the latest time stamp of  $T_{u_j}$ . The similarity of  $u_i$  and  $u_j$  is:

$$sim_u^t(u_i, u_j) = cosine(\mathcal{R}^{u(t)}(u_i), \mathcal{R}^{u(t)}(u_j)) \cdot rec(T_{u_j}) \quad (6)$$

Different from the traditional neighbourhood based models, as the active time period of  $u_i$  and  $u_j$  may be different, their similarity values are not necessary symmetric (i.e.,  $sim_u^t(u_i, u_j) \neq sim_u^t(u_j, u_i)$ ). We linearly combine the neighbourhood based and the content based approach. The prediction score of  $u_i$  for  $c_k$  is:

$$\mathcal{A}_u(u_i, c_k) = \alpha_1 \cdot \sum_{u_j \in \tilde{N}(u_i) \cap U_{c_k}} \omega \cdot sim_u^t(u_i, u_j) + \alpha_2 \cdot \mathcal{A}_c(u_i, c_k),$$

where  $\tilde{N}(u_i)$  is the neighbourhood of  $u_i$ ,  $\omega = \frac{1}{|\tilde{N}(u_i) \cap U_{c_k}|}$  is used

to smooth the value of  $\sum_{u_j \in \tilde{N}(u_i) \cap U_{c_k}} sim_u^t(u_i, u_j)$  to facilitate linear combination.  $0 \leq \alpha_1 \leq 1$ ,  $0 \leq \alpha_2 \leq 1$  and  $\alpha_1 + \alpha_2 = 1$ .

### 4.4.3 Matrix Factorization Model

The Matrix Factorization Model is typically used to predict the rating score of a user to a given item based on users' explicit rating data [1]. It also can be applied on binary user behavior data after generating negative samples from missing values randomly [11]. Although there is no explicit ratings to topics in a micro-blogging environment, users' topic preferences derived from their micro-blogs can be viewed as users' implicit ratings to topics. In this paper, users' topic preferences that calculated based on the user profiling approach discussed in Section 4.3, are used as positive samples (i.e.,  $w_{i,y}^{u(t)} > 0$ ). We also generated negative samples (i.e.,  $w_{i,y}^{u(t)} = 0$ ) for each user. Let  $|C_{u_i}|$  be the number of topics that  $u_i$  has, we randomly choose  $|C_{u_i}|$  topics that  $u_i$  has not shown interests in as this user's negative samples in the training set. As the task is to recommend Top  $N$  new topics to users, we extend the test set with  $M$  number of randomly selected negative samples. The Top  $N$  topics with highest prediction scores will be recommended to  $u_i$ . Let  $d_{ik}$  denote the prediction score of  $u_i$ 's preferences to  $c_k$ , similar to [1][11], it can be calculated as:

$$d_{ik} = \mu + b_{u_i} + b_{c_k} + p_{u_i}^T \cdot q_{c_k} \quad (7)$$

Where  $\mu$  is the average preference value for all topics. The parameters  $b_{u_i}$  and  $b_{c_k}$  indicate the deviations of  $u_i$  and  $c_k$ , respectively.  $p_{u_i}$  is the  $g$ -dimensional latent factor vector of  $u_i$ ,  $q_{c_k}$  is the  $g$ -dimensional latent factor vector of  $c_k$ . Let  $F^+$  denote all the positive samples and  $F^-$  be all negative samples sampled from missing values. A simple gradient descent technique was applied to minimize the following cost function:

$$\sum_{(u_i, c_k) \in F^+ \cup F^-} (w_{i,k}^{u(t)} - d_{ik}) + \varphi(\|p_{u_i}\|^2 + \|q_{c_k}\|^2 + b_{u_i}^2 + b_{c_k}^2)$$

Similar to neighborhood based model, the final prediction score is linearly combined with the content based approach.

$$\mathcal{A}_l(u_i, c_k) = \beta_1 \cdot \gamma \cdot d_{ik} + \beta_2 \cdot \mathcal{A}_c(u_i, c_k) \quad (8)$$

Where  $\gamma$  is a parameter to control the influence of  $d_{ik}$  to facilitate linear combination.  $0 \leq \beta_1 \leq 1$ ,  $0 \leq \beta_2 \leq 1$  and  $\beta_1 + \beta_2 = 1$ .

## 5. EXPERIMENTAL DESIGN

### 5.1 Recommendation Task

Specifically, the topic recommendation task can include the recommendation of hashtags, tag terms and keywords. In this paper, we focus on hashtag recommendation task. As each user and each hashtag can be represented with a set of related hashtags, tag terms and keywords, the proposed models are:

- **HM**: hashtag model. Each user and each hashtag is represented by hashtags respectively:  $\mathcal{R}^{C(t)}: H \rightarrow 2^{H \times [0,1]}$ ,  $\mathcal{R}^{u(t)}: U \rightarrow 2^{H \times [0,1]}$ . The proposed three kinds of recommendation approaches based on hashtag model are: (a) *HM-Content*: content based approach. (b) *HM-User*: neighbourhood based approach. (c) *HM-MF*: Matrix Factorization approach.
- **TM**: tag term model. Tag terms are used to represent each user and each hashtag:  $\mathcal{R}^{C(t)}: H \rightarrow 2^{P \times [0,1]}$ ,  $\mathcal{R}^{u(t)}: U \rightarrow 2^{P \times [0,1]}$ . The proposed three kinds of recommendation approaches based on this model are *TM-Content*, *TM-User*, and *TM-MF*.
- **KM**: keywords model. Keywords are used to represent each user and each hashtag:  $\mathcal{R}^{C(t)}: H \rightarrow 2^{K \times [0,1]}$ ,  $U \rightarrow 2^{K \times [0,1]}$ . The proposed three kinds of recommendation approaches based on this model are *KM-Content*, *KM-User*, and *KM-MF*.

### 5.2 Data Preparation

The experiments were conducted on a real world data that crawled from Twitter.com. We randomly selected 6,000 users who have used hashtags in their tweets and collected each user's tweets from April 19, 2011 to April 25, 2011. In the crawled raw dataset, nearly 16% of tweets contain hashtags. To avoid too sparse dataset, we only selected those users who have used at least 5 hashtags and their English tweets. The dataset D was split into training and test set. The statistical features of dataset D are shown in Table 1.

Table 1. Statistics of Dataset D

	Training Set	Test Set
D	4,673 Users 191,720 Tweets 38,621 Hashtags 38,701 Tag terms 141,849 Keywords 19/Apr/2011~24/Apr/2011	1,274 Users 11,808 Tweets 4,301 Hashtags 8,095 Tag terms 23,102 Keywords 25/Apr/2011

### 5.3 Experimental Setup

The users that appeared in both training and test set were selected as the test user set. Each test user's topics that appeared in the test

set but did not occur in the training set of this test user was used as this user's test topics. For a test user, a list of ordered topics that he/she has not used in his/her training set will be generated. If a topic in the recommendation list was in the test user's test topic set, then this recommended topic was counted as a hit. We adopt *Precision* and *Recall*, and the *HitRatio* and *HitTopics* to evaluate the accuracy. For a given test user, if the recommended topics got at least one hit topic for this user, then this user is counted as a hit user. *HitRatio* denotes the total number of hit users over all test users, while *HitTopics* denotes the total number of hit topics of all test users. The parameters of the proposed approaches are set after intensive experiments. The exponential rate  $\lambda = 0.01$ ,  $\Delta t$  was the elapsed time in hours. For neighbourhood based approach,  $K=100$ ,  $\alpha_1=0.1$ ,  $\alpha_2=0.9$ . For Matrix Factorization approach, the parameter settings are:  $M=100$ ,  $g=60$ ,  $\varphi=0.004$ , maximum iteration step is 40,  $\gamma=0.005$ ,  $\beta_1=0.4$ ,  $\beta_2=0.6$ .

## 6. RESULTS AND DISCUSSIONS

### 6.1 Results of the Proposed Approaches

Figure 1 shows the results of different topic and user representation models for the proposed content based approaches. The Top 5 *Precision* and *Recall* results of *HM-Content*, *TM-Content*, *KM-Content* are shown in this graph.

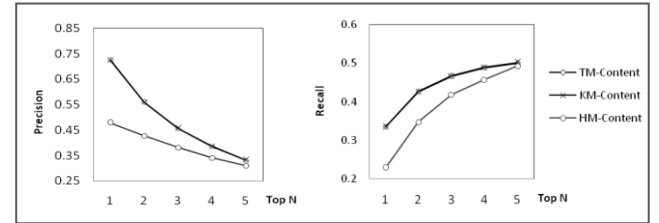


Figure 1. Results of Different Representation Models

Figure 1 shows that the content based approaches adopted tag term model (i.e., *TM-Content*) performed better than the approach based on hashtag model (i.e., *HM-Content*). It can be explained that quite a number of tweets not only contain hashtags but also contain tag terms that have been used by other users as hashtags. Although some users did not explicitly put hash symbols before these terms, they have similar topic interests with those users who have explicitly hashtagged these terms. Thus, more related topics can be obtained, which will help to find potential interested hashtags for each user. Moreover, *TM-Content* has similar performances with the approach based on keywords model (i.e., *KM-Content*). As the number of keywords usually is much larger than the number of tag terms, tag term model is computationally more efficient than the keywords model. Tag terms can be viewed as user selected document features of micro-blogs. With high accuracy and relatively low computation complexity, overall, the tag term model performed better than the other two models.

The comparison of the three proposed recommendation approaches based on tag term model are shown in Figure 2.

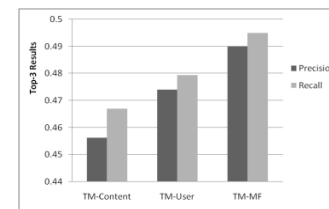


Figure 2. Results of Different Recommendation Approaches based on Tag Term Model

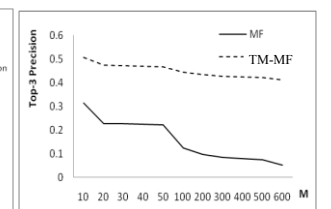


Figure 3. Results of Matrix Factorization Models with Different M Value models

Figure 2 shows that *TM-MF* performed better than the other two approaches. Figure 3 shows that with the increase of *M* value, both the Top-3 precision results of the proposed Matrix Factorization approach *TM-MF*, and *MF*, the matrix factor approach that does not combine content matching results, decreased. This is because that adding more negative samples in the test set usually will increase the error rate. Thus, this approach may unfairly take the advantage of the large proportion of positive samples in the test set, when only a very small number of negative samples were added to the test set. The content based approaches have less parameters and are computationally more efficient.

## 6.2 Comparison with Baseline Models

In this set of experiments, we compared the accuracy values of *TM-Content* with those related state-of-the-art temporal and non-temporal baseline methods. For fair comparison, the temporal baseline approaches adopted the same exponential decay function.

- *CF-User*: This is the standard user based collaborative filtering (CF) approach [13]. The similarity of two users was calculated based on the overlap of their hashtags.
- *tf-idf*: each user and topic are weighted by *tf-idf* approach [5].
- *MF-tf-idf*: It is based on the *tf-idf* weighted user topic profiles. It is inspired by the work [11]. No recency weighting, and the topic and user representation approaches are adopted.
- *MostPopular*: recommend the most popular hashtags to users.
- *MostRecent*: recommend the most recent hashtags to users.
- *MostRecentPopular*: recommend the most recent and popular hashtags to users.
- *ContentRecency*: This approach is based on *tf-idf* approach and inspired by the work of ranking recent information [7].

The Top-10 *HitRatio* and *HitTopics* results of these baseline models are shown in Figure 4. The Top 5 precision and recall results of *TM-Content*, and two better performed baseline models *MostRecentPopular* and *ContentRecency* are shown in Figure 5.

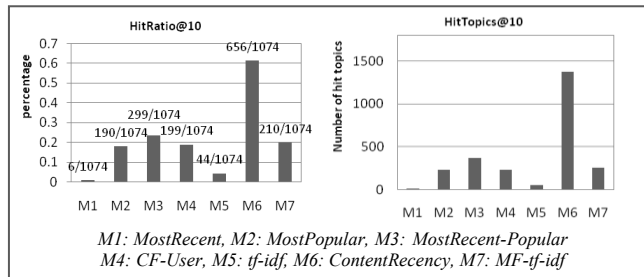


Figure 4. HitRatio and HitTopics Results

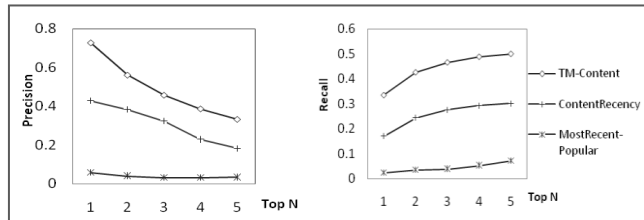


Figure 5. Precision and Recall results of selected models

As shown in Figure 4 and Figure 5, the proposed approach *TM-Content* performed the best. Compared with other approaches, the *MostRecent* approach had the worst performances. This suggests that the accuracy of recommendations may be extremely low if we only recommend the most recent topics to users. The *MostPopular* approach also failed to work well, as only a very small number of

topics were popularly used by all users. The results also suggest that it is very important to make personalized recommendations based on users' topic interests, while it is not enough to just recommend those hot streaming topics to users. The proposed approaches had the best performance. They rely on the multiple relationships among topics, users and tweets, and effectively used the recency information to find the time-aware content relevant topics of each topic and the time-aware topic preferences of each user. Moreover, the time-aware neighborhood formation of users and topics, and time-aware content matching between a user and a topic also contributed to find potentially interested topics for users.

## 7. CONCLUSIONS AND FUTURE WORK

In this paper, we discussed how to make personalized topic recommendations based on micro-blogs. Rather than using explicit social tie information, this paper focuses on making use of the implicit information network formed by the multiple relationships among users, topics and micro-blogs and the temporal information of micro-blogs, to expand topics and profile users. Furthermore, the content, neighborhood and Matrix Factorization based recommendation approaches are presented. The results of hashtag recommendation task show that the proposed approaches are effective. Future work will explore how to incorporate social influence of users to recommend topics.

## 8. ACKNOWLEDGEMENT

This research was carried out as part of the activities of, and funded by, the Smart Services Cooperative Research Centre (CRC) through the Australian Government's CRC Programme (Department of Innovation, Industry, Science and Research). The authors also would like to thank the support of The Australian National University and Queensland University of Technology.

## 9. REFERENCES

- [1] Koren, Y. Collaborative Filtering with Temporal Dynamics. KDD'09, 447-456
- [2] Yang, J., Leskovec, J. Patterns of temporal variation in online media. WSDM '11, 177-186
- [3] Hannon, J., etc., Recommending twitter users to follow using content and collaborative filtering approaches. RecSys '10, 199-206
- [4] Liang, H., Xu, Y., Li, Y., etc., Personalized Recommender System Based on Item Taxonomy and Folksonomy. CIKM'10, 1641-1644.
- [5] Chen, J., Nairn, R., Nelson, L., Bernstein, M., Chi, E., Short and tweet: experiments on recommending content from information streams. CHI '10, 1185-1194
- [6] Lin, J., Snow, R., etc., Smoothing techniques for adaptive online language models: topic tracking in tweet streams. KDD'11, 422-429
- [7] Efron, M., and Golovchinsky, G., Estimation methods for ranking recent information. SIGIR '11, 495-504
- [8] Hong, L., Yin, D., Guo, J., etc., Tracking Trends: Incorporating Term Volume into Temporal Topic Models. KDD'11, 484-492
- [9] Chen, J., Nairn, R., Chi, E., Speak Little and Well: Recommending Conversations in Online Social Streams. CHI'11, 2011, 217-226
- [10] Abel, F., Gao, Q., Houben, G., etc., Analyzing User Modeling on Twitter for Personalized News Recommendations. UMAP'11, 1-12
- [11] Lai, S., Xiang, L., etc., Hybrid Recommendation Models for Binary User Preference Prediction Problem. KDD Cup workshop'11, 2011
- [12] Xiang, L., Yuan, Q., etc., Temporal Recommendation on Graphs via Long- and short-term preference fusion. KDD'10, 723-731
- [13] Liang, H., Xu, Y., etc., Connecting Users and Items with Weighted Tags for Personalized Item Recommendations. HT'10, 51-60.
- [14] Sachan, M., Contractor, D., Faruque, T., Subramaniam, L., Using Content and Interactions for Discovering Communities in Social Networks. WWW'12, 331-340