



# Inferring user interests in microblogging social networks: a survey

Guangyuan Piao<sup>1</sup>  · John G. Breslin<sup>1</sup>

Received: 20 December 2017 / Accepted in revised form: 11 August 2018 / Published online: 18 August 2018  
© Springer Nature B.V. 2018

## Abstract

With the growing popularity of microblogging services such as Twitter in recent years, an increasing number of users are using these services in their daily lives. The huge volume of information generated by users raises new opportunities in various applications and areas. Inferring user interests plays a significant role in providing personalized recommendations on microblogging services, and also on third-party applications providing social logins via these services, especially in cold-start situations. In this survey, we review user modeling strategies with respect to inferring user interests from previous studies. To this end, we focus on four dimensions of inferring user interest profiles: (1) *data collection*, (2) *representation* of user interest profiles, (3) *construction and enhancement* of user interest profiles, and (4) the *evaluation* of the constructed profiles. Through this survey, we aim to provide an overview of state-of-the-art user modeling strategies for inferring user interest profiles on microblogging social networks with respect to the four dimensions. For each dimension, we review and summarize previous studies based on specified criteria. Finally, we discuss some challenges and opportunities for future work in this research domain.

**Keywords** User modeling · User interests · User profiles · Social web · Microblogging · Twitter · Social networks · Information filtering · Recommender systems · Personalization · Survey

---

✉ Guangyuan Piao  
guangyuan.piao@insight-centre.org  
<https://parklize.github.io>  
John G. Breslin  
<http://www.johnbreslin.com>

<sup>1</sup> Insight Centre for Data Analytics, Data Science Institute, National University of Ireland Galway, Galway, Ireland

## 1 Introduction

Microblogging<sup>1</sup> social networks such as Twitter<sup>2</sup> and Facebook<sup>3</sup> are being widely used in our daily lives. Twitter and Facebook have 328 million and 2 billion monthly active users,<sup>4,5</sup> which shows the popularity of these services. The abundant information generated by users in OSNs creates new opportunities for inferring user interest profiles, which can be used for providing personalized recommendations to those users either on those OSNs or on third-party services allowing social login functionality<sup>6</sup> from the same OSNs. Social login is a technology which allows visitors to a website to log in using their OSN accounts rather than having to register a new one.<sup>7</sup> A recent survey showed that over 94% of 18–34 year olds have used social login via Twitter, Facebook, etc.<sup>8</sup> With the continued widespread development of the social login functionality, inferring user interest profiles from their OSN activities plays a central role in many applications for providing personalized recommendations with the permission of those users, especially for cold-start users who have joined those services recently.

In the literature, there have been many studies that focused on inferring user interest profiles with different purposes such as providing personalized recommendations with respect to news (Abel et al. 2011b; Gao et al. 2011), research articles (Große-Bölting et al. 2015; Nishioka and Scherp 2016), and Points Of Interest (POI) (Abel et al. 2012). Despite the popularity of inferring user interests in OSNs, there is a lack of an extensive review on user modeling strategies for inferring user interest profiles in OSNs. To our knowledge, only one related short survey (Abdel-Hafez and Xu 2013) has been formally published. Abdel-Hafez and Xu (2013) provided a general overview of user modeling in social media websites which includes all types of OSNs without focusing on a specific type. As a result, the details of user modeling techniques for microblogging websites were not presented in Abdel-Hafez and Xu (2013). For example, including OSNs such as Delicious<sup>9</sup> and Flickr<sup>10</sup> which are based on *folksonomies* (folks taxonomies) together with microblogging OSNs for a single survey presents some difficulties due to the volume of literature on *folksonomy*-based user modeling (e.g., Abel 2011; Carmagnola et al. 2008; Hung et al. 2008; Mezghani et al. 2012; Szomszor et al. 2008, to name a few). In addition, the survey conducted by Abdel-Hafez and Xu (2013) does not cover studies from recent years. In this survey, we focus in particular on user modeling strategies in microblogging OSNs in terms of several user modeling dimensions, and analyze over 50 studies including more recent ones (see “Appendix A” for details of the surveyed studies).

---

<sup>1</sup> <https://en.wikipedia.org/wiki/Microblogging>.

<sup>2</sup> <https://twitter.com/>.

<sup>3</sup> <https://www.facebook.com/>.

<sup>4</sup> <https://www.omnicoreagency.com/twitter-statistics/>.

<sup>5</sup> <https://www.omnicoreagency.com/facebook-statistics/>.

<sup>6</sup> [https://en.wikipedia.org/wiki/Social\\_login](https://en.wikipedia.org/wiki/Social_login).

<sup>7</sup> <https://hbr.org/2011/10/social-login-offers-new-roi-fr>.

<sup>8</sup> <http://www.gigya.com/blog/why-millennials-demand-social-login/>.

<sup>9</sup> <https://del.icio.us/>.

<sup>10</sup> <https://www.flickr.com/>.

STEP 2 OF 4

Continue

What are you interested in?

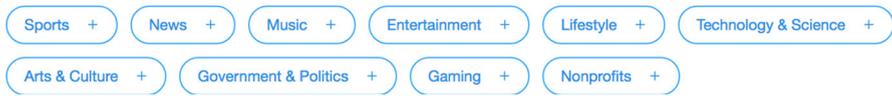


Fig. 1 Explicit information from users during signing up on Twitter

There has been a varied set of terms used to denote inferring user interests in the literature, such as “user (interest) modeling/profiling/detection”, “inferring/modeling/predicting user interests”. User modeling/profiling, as a broad term, may refer to different meanings without a specific definition. A general definition of *user profiling* given by Zhou et al. (2012) is “the process of acquiring, extracting and representing the features of users”. Similarly, in Brusilovsky et al. (2007), the *user model* is defined in the context of adaptive systems as “a representation of information about an individual user that is essential for an adaptive system to provide the adaptation effect”. Based on a specific definition of what the *features* and *information* are in these definitions by Zhou et al. (2012) and Brusilovsky et al. (2007), the corresponding user models/profiles and the process of obtaining them might be different.

Rich (1979) along with Cohen and Perrault (1979) and Perrault et al. (1978), where the terms *user model* and *user modeling* can be traced back to, also pointed out the need for classifying your user model as it might refer to several different things without a proper definition. Three major dimensions were used in Rich (1979) for classifying user models:

- Are they models of a canonical user or are they models of individual users?
- Are they constructed explicitly by the user themselves or are they abstracted by the system on the basis of the user’s behavior?
- Do they contain short-term or long-term information?

Explicit information denotes the information which requires direct input by users such as surveys or forms, which will impose an additional burden on the users. Figure 1 shows an example of collecting *explicit* information about user interests during sign up on Twitter for the first time.

1.1 Definition of user modeling in this survey

In the context of research on inferring user interests on OSNs, most studies have focused on exploiting *implicit* information such as the posts of users in order to infer user interest profiles. Based on the classification criteria from Rich (1979), user models discussed in this survey are about individual users constructed implicitly based on their activities. For the third criterion used in Rich (1979), there is no clear cut option as both short- and long-term information have been used in different user modeling strategies in the literature. In addition, user models can refer to various types of information relevant for each user in the domain of OSNs. For example, they might

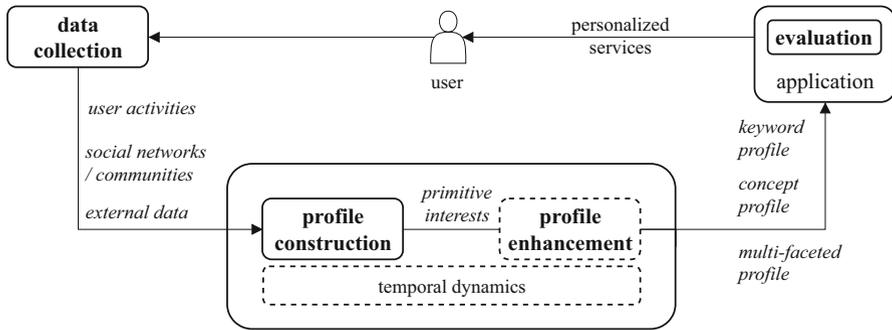


Fig. 2 Overview of user profile-based personalization process

contain basic information such as age, gender, country, etc., or keywords that represent their interests. In this paper, we focus particularly on user models with respect to user interests. Although several terms such as “user model” and “user profile” have been used interchangeably in the literature, here we formally define these terms as follows:

**Definition 1 (User Model)** A user model is a (data) structure that is used to capture certain *characteristics* about an individual user, and a *user profile* is the actual representation in a given user model. The process of obtaining the user profile is called *user modeling*.

Given this definition of a user model and the classification criteria from Rich (1979), user model in this survey aims to capture user *interests* with respect to an *individual* user *implicitly* based on *long-term* or *short-term* knowledge via a user modeling strategy, to derive the interest profile of that user.

Figure 2 presents an overview of the modified user profile-based personalization process from Abdel-Hafez and Xu (2013) and Gauch et al. (2007). We modified the process from Abdel-Hafez and Xu (2013) in order to reflect different aspects of user modeling strategies proposed in previous studies in the context of OSNs in detail. For example, we focus on data collection from *user activities*, *social networks/communities* or *external data* of an OSN instead of *explicit* or *implicit* feedback as most previous studies have focused on exploiting *implicit* information for inferring user interests. The modified user profile-based personalization process consists of three main phases. The first step is collecting data which will be used for inferring user interests. Subsequently, user interest profiles are constructed based on the data collected. We use *primitive interests* (Kapanipathi et al. 2014) to denote the interests directly extracted from the collected data. Those primitive interests can either be used as the final output of a profile constructor or can be further enhanced, e.g., based on background knowledge from Knowledge Bases (KBs) such as Wikipedia.<sup>11</sup> The output of the profile constructor is user interest profiles represented based on a predefined representation of interest profiles, e.g., word-based user interest profiles. Finally, the constructed user profiles are evaluated, and can be used in specific applications such as recommender systems for personalized recommendations.

<sup>11</sup> [www.wikipedia.org](http://www.wikipedia.org).

**Table 1** Online social networks used for previous studies

OSNs (# of studies)	Examples
Twitter (47)	Chen et al. (2010), Lu et al. (2012), Kapanipathi et al. (2011, 2014), Piao and Breslin (2016b, c, d, 2017a, b), Besel et al. (2016a, b), Abel et al. (2011a, b, c, 2012, 2013a), Siehndel and Kawase (2012), Michelson and Macskassy (2010), Bhattacharya et al. (2014), Orlandi et al. (2012), Hannon et al. (2012), Jiang and Sha (2015), Budak et al. (2014), Faralli et al. (2015b, 2017), Weng et al. (2010), Zarrinkalam and Kahani (2015), Zarrinkalam et al. (2016), Narducci et al. (2013), Xu et al. (2011), Garcia Esparza et al. (2013), Nishioka and Scherp (2016), Nishioka et al. (2015), Gao et al. (2011), Vu and Perez (2013), Phelan et al. (2009), Peñas et al. (2013), Sang et al. (2015), Karatay and Karagoz (2015), Kanta et al. (2012), O'Banion et al. (2012), Nechaev et al. (2017), Lim and Datta (2013), Große-Bölting et al. (2015), Trikha et al. (2018), Spasojevic et al. (2014), Jipmo et al. (2017)
Facebook (7)	Kang and Lee (2016), Orlandi et al. (2012), Kapanipathi et al. (2011), Narducci et al. (2013), Bhargava et al. (2015), Ahn et al. (2012), Spasojevic et al. (2014)
LinkedIn (2)	Kapanipathi et al. (2011), Spasojevic et al. (2014)
Google+ <sup>a</sup> (1)	Spasojevic et al. (2014)

<sup>a</sup><https://plus.google.com/>

In this paper, we mainly discuss four dimensions of the user modeling process: (1) *data collection*, (2) *representation* of user interest profiles, (3) *profile construction and enhancement*, and (4) the *evaluation* of the constructed user interest profiles. In summary, the contribution of this paper is threefold.

- First, we provide a detailed review of user modeling approaches on microblogging services in terms of the three phases in Fig. 2 with the following focuses:
  1. *What information is used for inferring user interest profiles?*
  2. *How are the user interest profiles represented?*
  3. *How are the user interest profiles constructed?*
  4. *How are the constructed user profiles evaluated?*
- Second, we summarize the approaches with respect to these focuses based on specified criteria to be specified later on.
- Finally, we discuss the challenges and opportunities based on the strengths and weaknesses of different approaches.

Table 1 provides a summary of OSNs used for the works discussed in this survey. As we can see from the table, Twitter has been widely used due to its popularity and the higher degree of openness. Other OSNs such as Facebook or LinkedIn<sup>12</sup> need to gain the permissions of users to access their data. Therefore, users have to be recruited for conducting an experiment, which results in less studies using these OSNs. In contrast to other studies, the study from Klout,<sup>13</sup> Inc. (Spasojevic et al. 2014), which is a social

<sup>12</sup> <https://www.linkedin.com/>.

<sup>13</sup> <https://klout.com/>.

**Table 2** Purposes of user modeling in OSNs from previous studies

Purpose	Examples
Predicting user interests	Kapanipathi et al. (2014), Kang and Lee (2016), Michelson and Macskassy (2010), Budak et al. (2014), Bhattacharya et al. (2014), Besel et al. (2016a, b), Orlandi et al. (2012), Narducci et al. (2013), Bhargava et al. (2015), Garcia Esparza et al. (2013), Vu and Perez (2013), Ahn et al. (2012), Abel et al. (2011c) Zarrinkalam et al. (2016), Ahn et al. (2012), Spasojevic et al. (2014), Jipmo et al. (2017), Faralli et al. (2017), Jiang and Sha (2015), Xu et al. (2011), Peñas et al. (2013), Lim and Datta (2013)
News recommendations	Abel et al. (2011b), Gao et al. (2011), Zarrinkalam and Kahani (2015), Sang et al. (2015), Kanta et al. (2012), O'Banion et al. (2012)
URL recommendations	Chen et al. (2010), Abel et al. (2011a), Piao and Breslin (2016a, b, c, d, 2017a, b)
Publication recommendations	Nishioka and Scherp (2016), Große-Bölting et al. (2015)
Tweet recommendations	Lu et al. (2012), Sang et al. (2015), Karatay and Karagoz (2015), Trikha et al. (2018)
Researcher recommendations	Nishioka et al. (2015)
POI recommendations	Abel et al. (2012)
User recommendations and classifications	Faralli et al. (2015b)
Concealing user interests	Nechaev et al. (2017)

media platform that aggregates and analyzes data from multiple OSNs, leveraged all the OSNs listed in Table 1. As different design choices can be made for user modeling with different purposes, Table 2 provides an overview of the purpose of user modeling in each study. As we can see from the table, the majority of the previous studies have been conducted with the purpose of predicting user interests followed by recommending different types of content such as news, URLs, publications, and tweets.

Table 3 is a conceptual framework for discussing user modeling strategies proposed in the related work and to act as a “guide” to the rest of this survey. The rest of this paper is organized as follows. In Sect. 2, we discuss what kind of information has been collected for inferring user interests. Section 3 introduces various representations of user interest profiles proposed in the literature. In Sect. 4, we review how user profiles have been constructed based on different dimensions such as considering the temporal dynamics of user interests. In Sect. 5, we discuss how those constructed user profiles have been evaluated in the literature. Finally, we conclude the paper with some discussions of opportunities and challenges with respect to user modeling on microblogging OSNs in Sect. 6.

**Table 3** Conceptual framework for discussing the related work in this survey**Data collection**

1. Using user activities
2. Using the social networks/communities of a user
3. Using external data

**Representation of user interest profiles**

1. Keyword profiles
2. Concept profiles
3. Multi-faceted profiles

**Construction and enhancement of user interest profiles**

1. Profile construction with weighting schemes
  - Heuristic approaches
  - Probabilistic approaches
2. Profile enhancement
  - Leveraging hierarchical knowledge
  - Leveraging graph-based knowledge
  - Leveraging collective knowledge
3. Temporal dynamics
  - Constraint-based approaches
  - Interest decay functions

**Evaluation**

## 2 Data collection

### 2.1 Overview

This section of the survey discusses the first stage of user modeling, which is the data collection. In the context of OSNs, there are various information sources for collecting data in order to infer user interest profiles such as user information including the tweets or profiles with respect to a user and information from that user's social network. The information used for user modeling is important as it might directly affect later stages such as the representation and construction of user interest profiles, and the quality of final profiles. The discussion is carried out over the criteria of whether the information is collected from a *user's activities* or the *social networks/communities* of that user from the target microblogging platform (where the target users come from) or *external data*. Given Twitter is the largest microblogging social networking platform and is the most used OSNs in the literature as depicted in Table 1, here we mainly focus on inferring user interest profiles on Twitter.

#### 2.1.1 Using user activities

A straightforward way of inferring user interests for a target user is leveraging information from the user's activities in OSNs. Take Twitter as an example, a user can have

different activities such as posting, re-tweeting, liking or replying to a tweet. Users can also describe themselves in their profiles or follow other people on Twitter which might reveal their interests. Therefore, we can leverage these user activities to infer user interests. This could be analyzing data from the posts, profiles or following activities of users. For instance, we can assume that a user is interested in `Microsoft` if the user mentions `Microsoft` frequently in the tweets or is following the Twitter account `@Microsoft`. However, inferring user interests from their activities such as posting tweets or re-tweeting requires users to be active, which is not always the case. For example, Gong et al. (2015) reported that a significant portion of Twitter users are *passive ones* who keep following other users in order to consume information on Twitter but who do not generate any content.

### 2.1.2 Using the social networks/communities of a user

Leveraging information from the social networks/communities of a user can be useful to infer user interest profiles, especially for *passive users* who have little activity but who keep following other users to receive information. In this case, the generated content such as the posts and the profiles of users in a user's social network can be used for inferring that user's interests. For example, if many followees of a user post tweets with respect to `Microsoft` frequently or belong to a common community related to `Microsoft`, we can assume that the user is interested in `Microsoft` as well.

### 2.1.3 Using external data

The ideal length of a post on any OSN ranges between 60 and 140 characters for better user engagement.<sup>14</sup> Analyzing microblogging services such as Twitter is challenging due to their nature of generating short, noisy texts. Understanding those short messages plays a key role in user modeling in microblogging services. To this end, previous studies have investigated leveraging external data such as the content of embedded links/URLs in a tweet, in order to enrich the short text for a better understanding of it. Haewoon et al. (2010) showed that most of the topics on Twitter are about news which could also be found in mainstream news sites. In this regard, some researchers have proposed linking microblogs to news articles and exploring the content of news articles in order to understand short texts in microblogging services better.

## 2.2 Review

### 2.2.1 Using user activities

The posts generated by users are the most common source of information for inferring user interests. Take Twitter as an example, the tweets or retweets of users provide a great amount of data that might implicitly indicate what kinds of topics a user might be interested in. Therefore, using the post streams of target users for inferring user

---

<sup>14</sup> <https://goo.gl/j97H1R>.

interest profiles has been widely studied in the literature regardless of the different manners for how user interests are represented. For instance, Kapanipathi et al. (2014) extracted Wikipedia entities from the tweet streams of users while Chen et al. (2010) extracted keywords from them. Inferring user interests based on users' posts requires users to be active, i.e., continuously generating content. On the one hand, there is an increasing number of users leveraging OSNs to seek the information they need, e.g., one in three Web users look for medical information, and over half of surveyed users consume news in OSNs<sup>15</sup> (Sheth and Kapanipathi 2016). On the other hand, there is also a rise of passive users in OSNs. For example, two out of five Facebook users only browse information without active participation within the platform<sup>16</sup> (Besel et al. 2016a), and Gong et al. (2015) reported that a significant portion of Twitter users are *passive ones* who consume information on Twitter without generating any content. Therefore, it is also important to infer user interest profiles for those *passive users* in OSNs.

Some studies pointed out that exploring posts for inferring user interests is computationally ineffective and unstable due to the changing interests of users (Besel et al. 2016a, b; Faralli et al. 2015b, 2017; Nechaev et al. 2017). Instead of analyzing posts to infer user interests, these studies proposed using the *followeeship* information of users, which can infer more stable user interest profiles as the relationships of common users tend to be stable (Myers and Leskovec 2014). In this line of work, *topical followees* that can be mapped to Wikipedia entities often need to be identified, e.g., identifying the followee account @messi10stats on Twitter as `wiki17:Lionel_Messi`. One of the problems with these approaches based on topical followees is that only a small portion of users' followees are topical ones. The authors from Faralli et al. (2015b) and Piao and Breslin (2017a) both showed that, on average, only 12.7% and 10% of followees of users in their datasets can be linked to Wikipedia entities. Therefore, a lot of information from followees that do not have corresponding Wikipedia entities is missed. For example, based on the topical-followees approach we cannot infer any interests for a user who is following @Alice who has a biography as “*User Modeling and Recommender Systems researcher*”.

**Pros and cons** Analyzing user activities for inferring user interests collects data from users themselves which can reflect their interests better compared to inferring from their social networks which will be discussed later. However, it requires users actively generate content in order to infer their interests from their generated content such as tweets, retweets, and likes on Twitter. Although leveraging the *topical-followees* approach can be used for inferring user interests for passive users, the usage of followees' information is limited.

## 2.2.2 Using the social networks/communities of a user

To cope with some problems such as inferring user interest profiles for passive users, information from social networks such as tweets from followees or followers or posts

<sup>15</sup> <http://bit.ly/pewsnnews>.

<sup>16</sup> <http://www.corporate-eye.com/main/facebook-growing-problem-passive-users/>.

<sup>17</sup> The prefix `wiki` denotes <https://en.wikipedia.org/wiki/>.

from Facebook friends can be utilized for inferring user interests for *passive users* as well as *active ones*. All aforementioned activities used for inferring a user's interests can be analyzed with respect to a user's social network as well for inferring that user's interests. For instance, Chen et al. (2010) and Budak et al. (2014) explored the tweets of target users and their followees to infer user interests. Although using posts generated by users is of great potential for mining user interests, it also faces some challenges due to the short and noisy nature of microblogs. Compared to the aforementioned topical-followees approach, information from the social networks of users such as their followees can provide much more information. Returning to the example of inferring user interests for a user who is following @Alice in the previous subsection, we can infer this user is interested in *User Modeling and Recommender Systems* based on the biography of @Alice—"User Modeling and Recommender Systems researcher". In Piao and Breslin (2017a), the authors proposed leveraging *biographies* of followees to extract entities instead of mapping followees to Wikipedia entities, and showed the improvement of inferred user interest profiles in the context of URL recommendations.

*List membership*, which is a kind of "tagging" feature on Twitter, has been explored as well. A list membership is a topical list or community which can be generated by any user on Twitter, and the creator of the list can freely add other users to the topical list. For instance, a user @Bob might create a topical list named "Java" and add his followees who have been frequently tweeting about news on this topic. Therefore, if a user @Alice is following users who have been added into many topical lists related to the topic Java, it might suggest that @Alice is interested in this topic as well. Kim et al. (2010) studied the usage of Twitter lists and confirmed that lists can serve as good groupings of Twitter users with respect to their characteristics based on a user study. Based on the study, the authors also suggested that the Twitter list can be a valuable information source in many application domains including recommendations. In this regard, several studies have exploited list memberships of followees to infer user interest profiles (Bhattacharya et al. 2014; Hannon et al. 2012; Piao and Breslin 2017b).

User interests might be following global trends in some trends-aware applications such as news recommendations. To investigate it, Gao et al. (2011) proposed interweaving global trends and personal user interests for user modeling. In addition to leveraging the tweets of a target user for inferring user interests, the authors constructed a trend profile based on all tweets in the dataset in a certain time period. Afterwards, the final user interest profile was built by combining the two profiles. The results showed that combined user interest profiles can improve the performance of news recommendations while the first profile based on personal tweets plays a more significant role in the combination.

**Pros and cons** On the one hand, a lot of data can be collected from the social networks of users, which is useful in the case of when inferring user interest profiles for passive users who do not generate much content but who keep following other users. On the other hand, it is difficult to distinguish the activities of a user's followees that are relevant to the interests of that user. For example, the followees of a user can tweet a

wide range of topics that they are interested in, and the user is not always interested in all those topics.

### 2.2.3 Using external data

One of the challenges of inferring user interests from OSNs is that the generated content is often short and noisy (Bontcheva and Rout 2014). To better understand the short texts of microblogging services such as tweets, external information beyond the target platform has been explored on top of the information sources discussed in the previous subsections. For instance, Abel et al. (2011b, c, 2013a) proposed linking tweets to news articles and extract the *primitive interests* of users based on their tweets as well as the content of related news articles. Several strategies were proposed in Abel et al. (2011c), which were later on developed as a Twitter-based User Modeling Service (TUMS, Tao et al. 2012). However, it requires maintaining up-to-date news streams from mainstream news providers such as CNN<sup>18</sup> in order to link tweets to relevant news articles. Instead, Abel et al. (2011a) and Piao and Breslin (2016c) leveraged the content of the embedded URLs in tweets. Hannon et al. (2012) used a third-party service Listorious,<sup>19</sup> which is a service providing annotated tags of list memberships on Twitter, for inferring user interest profiles. Given a target user  $u$ , the authors construct  $u$ 's interest profile based on the tags of list memberships with respect to the user.

With the popularity of different OSNs, users nowadays tend to have multiple OSN accounts across various platforms (Liu et al. 2013). In this context, some of the previous studies have investigated exploiting user interest profiles from other OSNs for cross-system user modeling. For instance, Orlandi et al. (2012) and Kapanipathi et al. (2011) presented user modeling applications that can aggregate different user interest profiles from various OSNs. However, the evaluation of aggregated user interest profiles has not been provided. Abel et al. (2012) investigated cross-system user modeling with respect to POI, and showed that the aggregation of Twitter and Flickr user data yields the best performance in terms of POI recommendations compared to modeling users separately based on a single platform. The result is in line with another study by them which aggregated user interest profiles on social tagging systems such as Delicious,<sup>20</sup> StumbleUpon,<sup>21</sup> and Flickr (Abel et al. 2013b).

The work from Klout (Spasojevic et al. 2014), which allows their users to add multiple OSN identities on their services, showed many insights on aggregating user information from multiple information sources in different OSNs for inferring user interests. The authors pointed out that using user-generated content (UGC) alone leads to a high precision but low recall for topic recommendations, and therefore, other information sources such as the ones from followees are needed. They also observed that the overlap of a user's interests from different OSNs is very small, which shows that a user may not reveal all his/her interests on any single OSN alone due to the different

---

<sup>18</sup> <http://edition.cnn.com/>.

<sup>19</sup> <http://listorious.com>, not available at the time of writing.

<sup>20</sup> <https://www.delicious.com>.

<sup>21</sup> <https://www.stumbleupon.com>.

**Table 4** Information used for collecting data for inferring user interest profiles

User Activities	Social Networks/Communities	External Data	Examples
✓			Lu et al. (2012), Kapanipathi et al. (2014), Kang and Lee (2016), Piao and Breslin (2016b, d), Orlandi et al. (2012), Weng et al. (2010), Michelson and Macskassy (2010), Siehndel and Kawase (2012), Zarrinkalam and Kahani (2015), Zarrinkalam et al. (2016), Jiang and Sha (2015), Narducci et al. (2013), Xu et al. (2011), Nishioka and Scherp (2016), Nishioka et al. (2015), Peñas et al. (2013), Sang et al. (2015), O'Banion et al. (2012), Große-Bölting et al. (2015), Jipmo et al. (2017), Trikha et al. (2018), Bhargava et al. (2015), Ahn et al. (2012), Besel et al. (2016a, b), Faralli et al. (2015b, 2017), Lim and Datta (2013), Nechaev et al. (2017), Vu and Perez (2013)
	✓		Phelan et al. (2009), Piao and Breslin (2017a, b), Bhattacharya et al. (2014)
		✓	Spasojevic et al. (2014)
✓	✓		Chen et al. (2010), Budak et al. (2014), Karatay and Karagoz (2015), Gao et al. (2011), Kanta et al. (2012)
✓		✓	Piao and Breslin (2016c), Garcia Esparza et al. (2013), Abel et al. (2011a, b, c, 2012, 2013a), Orlandi et al. (2012), Kapanipathi et al. (2011), Hannon et al. (2012)

characteristics of OSNs. Therefore, aggregating users' information in different OSNs leads to a better understanding of their interests (Spasojevic et al. 2014).

**Pros and cons** Leveraging external data such as the content of embedded URLs in a tweet can provide a better understanding of short microblogs, and exploring information from other OSNs of users can reveal their interests better compared to exploring a single OSN. Nevertheless, analyzing external data requires an additional effort and it is not always available. In addition, external data can also have irrelevant content with respect to user interests and might introduce some noise.

### 2.3 Summary and discussion

In this section, we reviewed different information sources that have been used for collecting data in order to infer user interest profiles. Table 4 summarizes information sources used for inferring user interest profiles in the literature. As we can see from Table 4, user activities have been used widely for inferring user interest profiles in microblogging social networks in previous studies.

Although there have been many information sources used for inferring user interests, the comparison of different data sources for inferring user interest profiles has been less explored. Some approaches have utilized different aspects of information of followees such as *topical followees*, *biographies*, or *list memberships* (e.g., Besel et al. 2016a, b; Bhattacharya et al. 2014; Hannon et al. 2012; Piao and Breslin 2017a). However, it has not been clearly shown in these studies if these approaches perform better than exploiting users' posts. The usefulness of user interest profiles built from various information sources might be different depending on different applications. For instance, Chen et al. (2010) showed that user interest profiles based on the user's own streams perform better than profiles based on followee streams in the context of URL recommendations on Twitter. However, those profiles based on followee streams might be more useful for recommending followees.

In addition, combining different information sources have shown its efficiency in a few studies (e.g., Abel et al. 2012; Piao and Breslin 2017b). However, how to combine different information sources for inferring user interests, and whether there is a synergistic effect on application performance by the combination might require more study. For instance, user interests extracted from different data sources can be either aggregated into a single user interest profile (e.g., Abel et al. 2012; Orlandi et al. 2012) or remain as separate profiles (Piao and Breslin 2017b) to measure the preference score of a candidate item for recommendations. Also, combining different data sources has mainly been studied for aggregating user interests from multiple OSNs. Instead, combining different data sources inside the target platform might be useful for inferring user interests as well, e.g., combining extracted user interests from different information sources of followees and users.

### 3 Representation of user interest profiles

#### 3.1 Overview

In this section, we provide an overview of how user interest profiles have been represented in the different approaches. Here we first provide an overview of user representations for personalized information access that was introduced in Gauch et al. (2007), and *multi-faceted profiles* which have been proposed in several studies in the literature. We then carry out the review based on three different types of representations in the context of inferring user interest profiles in OSNs in the literature, which include (1) *keyword profiles*, (2) *concept profiles*, and (3) *multi-faceted profiles*.

In Gauch et al. (2007), the authors defined three types of user representations for personalized information access:

- keyword profiles;
- concept profiles;
- semantic network profiles.

**Keyword profiles** In this representation of user interest profiles, each *keyword* or a *group of keywords* can be used for representing a topic of interest. This approach was predominant in every adaptive information retrieval and filtering system and is

still popular in these areas (Brusilovsky et al. 2007). When using each keyword for representing user interests, the importance of each word with respect to users can be measured using a defined weighting scheme such as TF-IDF (Term Frequency · Inverse Document Frequency) from information retrieval (Salton and McGill 1986). In the case of using groups of keywords for representing user interests, the user interest profiles can be represented as a probability distribution over some topics, and each topic is represented as a probability distribution over a number of words. The topics can be distilled using topic modeling approaches such as Latent Dirichlet Allocation (LDA) (Blei et al. 2003), which is an unsupervised machine learning method to learn topics from a large set of documents.

**Concept profiles** Concept-based user profiles are represented as conceptual nodes (concepts) and their relationships, and the concepts usually come from a pre-existing knowledge base (Gauch et al. 2007). They can be useful for dealing with the problems that keyword profiles have. For example, WordNet (Miller 1995) groups related words together in concepts called *synsets*, which has been proved useful for dealing with *polysemy* in other domains. For example, Stefani (1998) used WordNet synsets for representing user interests in order to provide personalized website access instead of using keywords as they are often not enough for describing someone's interests. Another type of concept is *entities with URIs* (Uniform Resource Identifiers). For instance, this involves using `dbr22:Apple_Inc.` to denote the company Apple, which is disambiguated based on the context of the word *apple* in a text such as tweet and linked to knowledge bases such as Wikipedia or DBpedia (Auer et al. 2007). DBpedia is the semantic representation of Wikipedia and it has become one of the most important and interlinked datasets on the Web of Data, which indicates a new generation of technologies responsible for the evolution of the current Web from a Web of interlinked documents to a Web of interlinked data (Heath and Bizer 2011). To facilitate reading, we use DBpedia concepts to denote concepts from Wikipedia or DBpedia.

**Semantic network profiles** This type of profile aims to address the polysemy problem of keyword-based profiles by using a weighted semantic network in which each node represents a specific word or a set of related words. This type of profile is similar to concept profiles in the sense of the representation of conceptual nodes and the relationships between them, despite the fact that the concepts in semantic network profiles are learned (modeled) as part of user profiles by collecting positive/negative feedback from users (Gauch et al. 2007). As most previous works have focused on implicitly constructing user interest profiles in microblogging services, this type of profile has not been used in the domain of user modeling in microblogging services.

**Multi-faceted profiles** Based on these representation strategies, user interest profiles can include different aspects of user interests such as interests inferred from their tweets, profiles or list memberships. These different aspects of user interests can be combined to construct a single user interest profile or maintained separately as several

---

<sup>22</sup> The prefix `dbr` denotes <http://dbpedia.org/resource/>.

user interest profiles for a target user. Although it is common to use a single representation with respect to a user interest profile, the *polyrepresentation theory* (Ingwersen 1994) based on a cognitive approach indicates that the overlaps between a variety of aspects or contexts with respect to a user within the information retrieval process can decrease the uncertainty and improve the performance of information retrieval. Based on this theory, White et al. (2009) studied polyrepresentation of user interests in the context of a search engine. The authors combined five different aspects/contexts of a user for inferring user interests, and showed that polyrepresentation is viable for user interest modeling.

## 3.2 Review

### 3.2.1 Keyword profiles

Similar to other adaptive information retrieval and filtering systems, representing user interests using *keywords* or *groups of keywords* is popular in OSNs as well. For instance, Chen et al. (2010) and Bhattacharya et al. (2014) represented user interest profiles by using vectors of weighted keywords from the tweets and the descriptions of list memberships of users, respectively. Despite the huge volume of information from UGC, extracting keywords from microblogs for inferring user interest profiles is challenging due to the nature of short and noisy messages (Liao et al. 2012).

As an alternative approach, another special type of keyword such as *tags* and *hashtags*<sup>23</sup> has been used for inferring user interest profiles. In contrast to the words mined from the short texts of microblogs, keywords from tags/hashtags might be more informative and categorical in nature. Abel et al. (2011a, b) investigated hashtag-based user interest profiles by extracting hashtags from the tweets of users, and Hannon et al. (2012) leveraged keywords from the tags of users' list memberships for representing their interest profiles.

Topics distilled from topic modeling approaches such as LDA are also popular for representing user interest profiles. A topic has associated words with their probabilities with respect to the topic. For example, an information technology-related topic can have some top associated words such as “google, twitter, apple, web”. Weng et al. (2010) used LDA to distill 50 topics and represented each user as a probability distribution over these topics. In Abel et al. (2011b, c, 2013a), the authors also used topics for representing user interests where those topics were extracted by ready-to-use NLP (Natural Language Processing) APIs such as OpenCalais.<sup>24</sup>

**Pros and cons** Keyword profiles are the simplest to build, and do not rely on external knowledge from a knowledge base. One of the drawbacks of the keyword-based user profiles is *polysemy*, i.e., a word may have multiple meanings which cannot be distinguished by using keyword-based representation. In addition, these keyword-based approaches lack semantic information and cannot capture relationships among these words, and the assumption of topic modeling approaches that a document has

<sup>23</sup> <https://en.wikipedia.org/wiki/Hashtag>.

<sup>24</sup> <http://www.opencalais.com/>.

rich information is not the case for microblogs (Zarrinkalam 2015). Spasojevic et al. (2014) further pointed out that topic modeling approaches cannot provide a scalable solution for inferring topics for millions of users which include a great number of passive users.

### 3.2.2 Concept profiles

To address some problems of keyword-based approaches, researchers have proposed leveraging *concepts* from KBs such as DBpedia for representing user interests. One of the advantages of leveraging KBs is that we can exploit the background knowledge of these concepts to infer user interests which might not be captured if using keyword-based approaches. For instance, a big fan of the Apple company would be interested in any brand-new products from Apple even the names of these products have never been mentioned in the user's primitive interests (Lu et al. 2012). Concepts from various types of KBs have been leveraged for different purposes of user modeling, such as the ones from simple concept taxonomies with respect to news (Kang and Lee 2016), domain-specific KBs such as STW,<sup>25</sup> ACM CCS, and Medical Subject Headings<sup>26</sup> (MeSH) (Große-Bölting et al. 2015; Nishioka and Scherp 2016; Nishioka et al. 2015), and cross-domain KBs such as DBpedia (Abel et al. 2011a, b, c; Faralli et al. 2015b; Lu et al. 2012; Piao and Breslin 2016b, c, d, 2017a, b). In the following, we discuss some details of the representation strategy using DBpedia concepts which have been the most widely used for representing user interest profiles.

**Entity-based profiles** This approach extracts entities from information sources such as a user's tweets, and uses these entities to represent user interest profiles. Take the following real-word tweet as an example (Michelson and Macskassy 2010):

“#Arsenal winger Walcott: Becks is my England inspiration: <http://tinyurl.com/37zyjsc>”,

there are four entities such as `dbr:Arsenal_F.C.`, and `dbr:Theo_Walcott` within the tweet, which can be used for constructing entity-based user interest profiles. However, this approach is difficult to infer more specific interests which might need to be represented by combining multiple related entities or interests that cannot be found in a knowledge base. To address this issue, some studies have proposed representing each topic of interest as a *conjunction of multiple entities*, which are correlated on Twitter in a certain timespan (Zarrinkalam and Kahani 2015; Zarrinkalam et al. 2016). These sets of entities for representing a topic of interest can be learned via unsupervised approaches in a similar manner to learning topics with topic modeling approaches for keyword-based profiles.

**Category-based profiles** An alternative approach is using DBpedia *categories*, which represents more general user interests compared to using DBpedia *entities*.

<sup>25</sup> <http://zbw.eu/stw>.

<sup>26</sup> <https://www.nlm.nih.gov/mesh/>.

Returning to the example in the previous paragraph, the categories of the mentioned entities in that tweet such as `dbr:Category:English_Football_League` can be used for representing the topic of interests instead of those entities. One can also choose the level or depth of categories in a KB for representing user interest profiles or use all categories related to primitive interests. The top-level DBpedia categories can refer to general ones such as `dbr:Category:Sports` and `dbr:Category:Health` compared to the categories in a lower level such as `dbr:Category:English_Football_League`. For example, Michelson and Macskassy (2010) and Nechaev et al. (2017) used top-level categories to represent user interest profiles while other studies (Faralli et al. 2017; Kapanipathi et al. 2014; Flati et al. 2014, etc.) used hierarchical categories to represent user interest profiles. Figure 3 shows an example of category-based representation of user interests based on extracted entities from followers' account names, which is called *Twixonomy* (Faralli et al. 2017).

**Hybrid representations** Each aforementioned representation has its strengths and weaknesses. In terms of entity- or category-based representations, extracting entities with URIs is a fundamental step for constructing either *entity-* or *category-based* user interest profiles. However, the task of extracting entities is non-trivial (Kapanipathi et al. 2014) due to the noisy, informal language of microblogs (Ritter et al. 2011). In addition, knowledge bases might be out-of-date for emerging concepts on microblogging services, and therefore cannot capture these concepts during the entity extraction process. To overcome the drawbacks of using a single interest format, *hybrid representations* based on various interest formats have been explored as well. Instead of using only entities or categories for representing user interests, hybrid approaches combine different interest formats for constructing user profiles (Faralli et al. 2015b; Nishioka and Scherp 2016; O'Banion et al. 2012; Piao and Breslin 2016b, c, 2017a, b).

For example, O'Banion et al. (2012) used categories as well as entities to represent user interest profiles. Piao and Breslin (2016c, d) proposed a hybrid approach using both DBpedia entities and WordNet synsets for representing user interests in order to capture user interests that might be missed due to the problem with entity recognition in microblogs.

**Pros and cons** On the one hand, concept-based approaches present the semantics between concepts and can leverage background knowledge about concepts for propagating user interest profiles. On the other hand, these approaches rely on pre-existing or pre-constructed KBs which might be not always available in or lack of coverage with respect to some domains.

### 3.2.3 Multi-faceted profiles

Multi-faceted profiles model multiple aspects for a target user based on different information sources or using different representation strategies in order to derive a comprehensive view of that user. The assumption here is that different aspects of users may complement each other and improve the inferred user interest profiles.

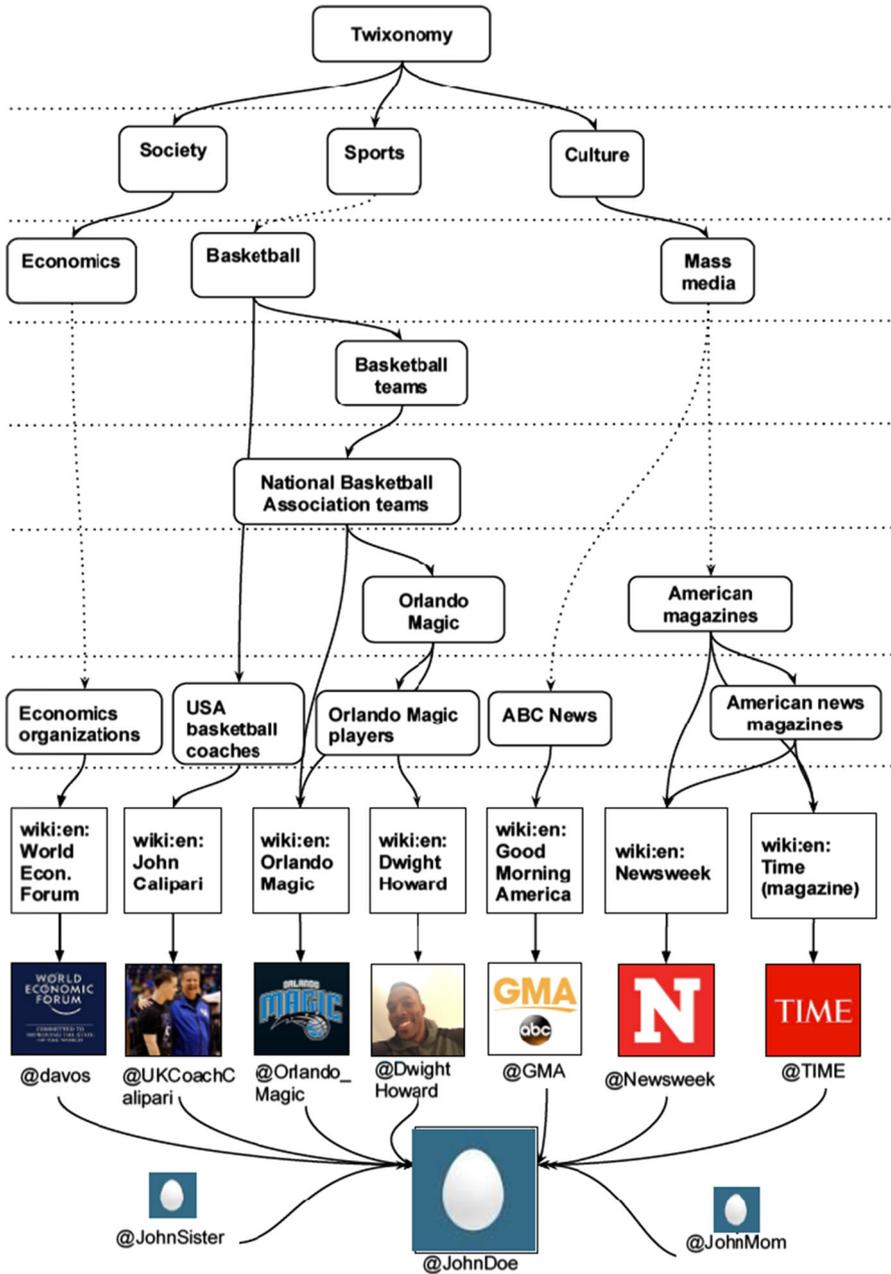
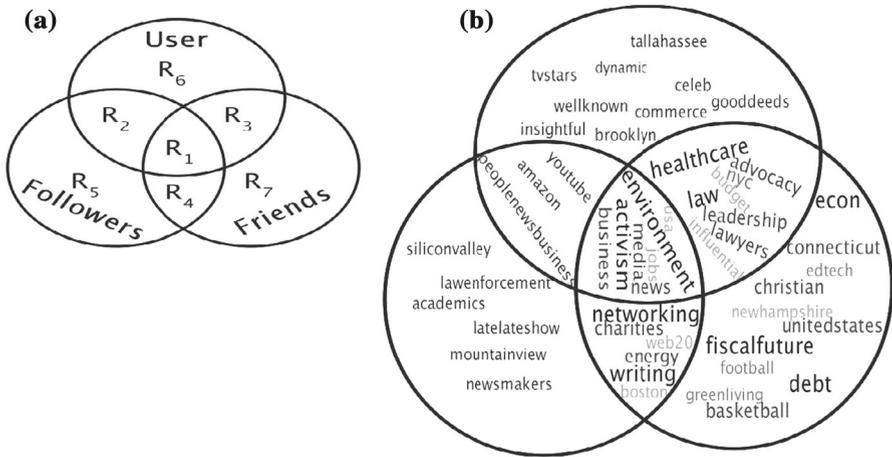


Fig. 3 An example of Twixonomy for a single user (Faralli et al. 2017)



**Fig. 4** **a** Intensional and extensional profile regions. **b** Barack Obama’s profile showing the tags associated with Obama and his followees (friends in the figure) and followers (Hannon et al. 2012)

Hannon et al. (2012) proposed a multi-faceted user profile which includes user interests from target users, their followees, and followers. Figure 4 shows an example from Hannon et al. (2012) for representing user interests, where user interests are represented based on the tags of list memberships of users, followees, or followers provided by a third-party service. The figure shows that user interests inferred from different aspects can complement each other and lead to a better understanding of a target user. However, they did not evaluate the effectiveness of multi-faceted profiles in the context of personalized recommendations and left it as a future work.

The authors in Lu et al. (2012) and Chen et al. (2010) both constructed two keyword-based user interest profiles for each user. In Chen et al. (2010), two keyword-based user interest profiles were built based on the tweets of users and those of their followees for recommending URLs on Twitter. The results in Chen et al. (2010) showed that using user interest profiles based on the tweets of users performs better than using those based on the tweets of their followees. Lu et al. (2012) proposed using DBpedia entities and the affinity of other users to construct two user interest profiles for recommending tweets on Twitter. For a given user, the first user profile was represented as a vector of DBpedia entities, which were extracted from the user’s tweets. Both of these studies did not investigate the synergistic effect of combining these two aspects compared to considering a single aspect of users. More recently, Piao and Breslin (2017b) showed that leveraging concept-based profiles from the biographies and list memberships of followees can complement each other and improve the URL recommendation performance on Twitter.

**Pros and cons** Multi-faceted profiles provide a comprehensive view of a user with respect to his/her interests and can improve recommendation performance. On the other hand, multiple information sources have to be explored for constructing multi-faceted profiles.

### 3.3 Summary and discussion

In this section, we reviewed various ways of representing user interests such as using *keywords*, various types of *concepts*, and some multi-faceted approaches. Table 5 shows a summary of different representations of user interests adopted by previous studies.

Those different representations of user interests might work differently depending on the application where these user profiles are used. For example, we usually have to construct item profiles in the same way as constructing user interest profiles in order to measure the similarity between them for providing recommendations. The entity-based representation strategies for user interests might be appropriate for recommending items with long content, e.g., news or URL recommendations as the content of them is usually long. In contrast, these representation strategies might not work well for recommending items with short descriptions such as tweets due to the difficulty of extracting entities from them. For example, the low recall of entities on Twitter has been reported in both Kapanipathi et al. (2014) and Piao and Breslin (2016c) using several state-of-the-art NLP APIs. In a recent study (Manrique and Mariño 2017), the authors also showed that 30% of the titles of a research article cannot extract any entity at all. Some hybrid approaches such as combining word- and concept-based representations might be useful in this case.

In addition, different facets should be considered carefully for constructing multi-faceted profiles in the context of item recommendations. Each facet of multi-faceted profiles can have different importance for the recommended items, and leveraging completely unrelated facets might introduce noise to the constructed profiles. For example, Piao and Breslin (2017b) showed that different weights are required for different facets in order to achieve the best performance in URL recommendations on Twitter. Abel et al. (2013b) showed that it is helpful to have sufficient overlap between different facets of multi-faceted profiles for tag recommendations in a cold start.

It is also worth noting that the structure of user interest profiles can be different even with the same user interest format. Take a category-based user interest profile as an example, it can be a *vector*, *taxonomy* or *graph* by retaining the hierarchical or general relationships among categories. Also, the final profile extracted from the same structure can be different. For instance, both user interest profiles proposed in Faralli et al. (2017) (see Fig. 3) and Kapanipathi et al. (2014) were represented as a *taxonomy* at first, but were used differently for the final representation of user interests. In Faralli et al. (2017), entities or categories in different levels were used separately as an interest vector for representing a user, e.g., using categories that were two hops away from the user's primitive interests as the final interest profile. However, using a specific abstraction level of the category taxonomy for all users does not consider that different users might have different depths or expertise levels in terms of a topic of interests. In contrast, Kapanipathi et al. (2014) sorted all categories in the taxonomy of a user based on their weights for representing the user's interest profile. The different usages of the category taxonomy indicate some opportunities and challenges. On the one hand, the taxonomy structure of user interests is flexible enough to extract different abstraction levels of user interests or an overview of them. On the other hand, it has

**Table 5** Representation of user interest profiles

Keyword profiles	Concept profiles	Multi-faceted	Examples
✓		✓	Chen et al. (2010), Hannon et al. (2012)
	✓	✓	Lu et al. (2012), Piao and Breslin (2017b), Spasojevic et al. (2014)
✓	✓		Abel et al. (2011a, b, c, 2013a), Kanta et al. (2012)
✓			Weng et al. (2010), Xu et al. (2011), Sang et al. (2015), Bhattacharya et al. (2014), Vu and Perez (2013), Phelan et al. (2009)
	✓		Kapanipathi et al. (2014), Besel et al. (2016a, b), Faralli et al. (2017), Siehndel and Kawase (2012), Michelson and Macskassy (2010), Piao and Breslin (2016b, c, d, 2017a), Karatay and Karagoz (2015), Kang and Lee (2016), Abel et al. (2012), Narducci et al. (2013), Orlandi et al. (2012), Kapanipathi et al. (2011), Jipmo et al. (2017), Zarrinkalam and Kahani (2015), Zarrinkalam et al. (2016), Bhargava et al. (2015), Garcia Esparza et al. (2013), Nishioka and Scherp (2016), Nishioka et al. (2015), Jiang and Sha (2015), Große-Bölting et al. (2015), Gao et al. (2011), Nechaev et al. (2017), Budak et al. (2014), Peñas et al. (2013), Trikha et al. (2018), O'Banion et al. (2012), Ahn et al. (2012), Lim and Datta (2013)

not been investigated which type of user interest profile obtained from the taxonomy structure is better.

## 4 Construction and enhancement of user interest profiles

### 4.1 Overview

So far we have focused our discussion on collecting data from various sources for inferring user interests, and different representations for interest profiles. In this section, we provide details on how user interest profiles of a certain representation can be constructed based on the collected data. The overview of the construction and enhancement of user interest profiles is carried out based on three criteria:

- profile construction with weighting schemes;
- profile enhancement;
- temporal dynamics of user interests.

Based on a defined representation of user interest profiles, a profile constructor aims to determine the weights of user interest formats such as words or concepts in user

profiles with a certain *weighting scheme*. The weights of interest formats denote the importance of these interests with respect to a user. In Sect. 4.2.1, we review different weighting schemes based on various information sources such as users' posts or their followers, etc.

Primitive interest profiles, e.g., entity-based user profiles, can be further enhanced by using background knowledge from knowledge bases. For instance, this can be achieved by inferring category-based user interest profiles on top of the extracted entities from the data collected. Section 4.2.2 describes the approaches leveraging knowledge bases for enhancing primitive interest profiles.

User interests can change over time in OSNs. For instance, a user interest profile built during the last 2 weeks might be totally different from one built from 2 years ago. In Sect. 4.2.3, we look at whether or not the temporal dynamics of user interests have been considered when constructing user interest profiles, and if yes, how they have been incorporated during the construction process.

## 4.2 Review

### 4.2.1 Profile construction with weighting schemes

The output of a profile constructor is a primitive user interest profile represented by weighted interests based on a predefined representation. A *weighting scheme* is a function or process to determine the weights of user interests.

**Heuristic approaches** A common and simple weighting scheme is using the frequency of an interest  $i$  (e.g., a keyword or an entity) to denote the importance of  $i$  with respect to a user  $u$ , which can be formulated as below when the data source is  $u$ 's posts:

$$TF_u(w_i) = \text{frequency of } i \text{ in } u\text{'s posts.} \quad (1)$$

Despite its simplicity, this approach has been widely used in the literature, particularly in entity-based user interest representations (Abel et al. 2011c; Kapanipathi et al. 2014; Tao et al. 2012). Interests represented as concepts such as entities extracted from tweets might come with their confidence scores, and these scores can be incorporated into a weighting scheme. For instance, Jiang and Sha (2015) used TF with the confidence scores of extracted entities from tweets as their weighting scheme.

One problem with TF is that common words or entities which appear frequently in many users' interest profiles and may not be important as user interests. TF-IDF is another common weighting scheme to cope with this problem. The IDF score of  $i$  with respect to a user  $u$  based on  $u$ 's tweets can be measured as below (Chen et al. 2010):

$$IDF_u(i) = \log \left[ \frac{\# \text{ all users}}{\# \text{ users using } i \text{ at least once}} \right]. \quad (2)$$

Instead of using users for measuring the IDF score of an interest, IDF has been applied in other ways as well. For example, Nishioka and Scherp (2016) applied IDF with randomly retrieved tweets from the streaming API of Twitter, and Gao et al. (2011) applied IDF to value the specificity of an interest within a given period of time. It is worth noting that the IDF weighting can also be applied after the *profile enhancement* process (e.g., Nishioka and Scherp 2016; Piao and Breslin 2016c).

More sophisticated approaches can be applied for weighting user interests. In Vu and Perez (2013), the authors compared different weighting schemes such as TF-IDF, TextRank (Mihalcea and Tarau 2004), and TI-TextRank which was proposed by the authors by combining TF-IDF and TextRank. Based on a user study, the authors showed that TI-TextRank performs best for ranking keywords from the tweets of users.

In the context of OSNs, specific approaches have to be devised for constructing user interest profiles by exploiting their social networks such as followees on Twitter (Chen et al. 2010; Lu et al. 2012). To this end, several methods have been proposed. For example, Chen et al. (2010) first retrieved a set of *high-interest words* for followees as follows in order to build a user profile based on followees' tweets: First, keyword-based user interest profiles were created using the TF-IDF weighting scheme based on the tweets of followees, which are called *self-profiles*. Next, for each *self-profile* for followees of  $u$ , they picked all words that have been mentioned at least once, and selected the top 20% of words based on their occurrences. In addition, the words that are not in other followees' profiles were removed. Subsequently, the weight of each word in the set of *high-interest words* was measured as below:

$$FTF_u(i) = \# u' \text{ s followees who have } i \quad (3)$$

*as one of their high – interest words.*

Similar approaches of  $FTF_u(i)$  were adopted in Piao and Breslin (2017b) and Bhattacharya et al. (2014) but by exploring the list memberships of followees instead of their tweets for extracting user interests.

An alternative approach for aggregating the weights of interests in the followees' profiles is normalizing each followee's profiles and then aggregating those normalized weights for building user interest profiles (Piao and Breslin 2017b; Spasojevic et al. 2014). In Piao and Breslin (2017b), the authors showed that this simple alternative approach performs better compared to  $FTF_u(i)$  for weighting entities extracted from the list memberships of followees when using inferred user interest profiles for URL recommendations on Twitter. These approaches assume that each followee is equally important when aggregating their interest profiles for building the user interest profile of a target user. However, some followees' profiles can be more important compared to others with respect to the target user. In Karatay and Karagoz (2015), the authors incorporated the relative ranking scores of social networks into their weighting scheme to weight the entities of users.

**Probabilistic approaches** The aforementioned approaches focus on interests such as entities appearing in users' posts, however, not all the entities related to a post explicitly appear in that post. In this regard, some approaches extracted interests such as entities by measuring the similarity between a post and an entity. For instance, Lu et al. (2012)

and Narducci et al. (2013) used the Explicit Semantic Analysis (ESA) (Gabrilovich and Markovitch 2007) algorithm, which is designed to compute the similarity between texts, for obtaining the weights of entities for each tweet of a user. Those weights of entities were then aggregated for constructing entity-based primitive interests of users. Ahn et al. (2012) quantified the degree of an interest, i.e., a Facebook entity, based on two factors: (1) the familiarity with each social neighbor, and (2) the similarity between the topic distributions of a social content and an interest. *Social content* is the combined text of a post and its comments between users, and the topic distributions of it is obtained using LDA.

The weights of user interests have also been learned in unsupervised ways in the literature. For instance, Weng et al. (2010) treated tweet histories of each user as a big document, and used LDA to learn topic distributions for each user. Trikha et al. (2018) and Zarrinkalam et al. (2017) also used LDA to infer topic distributions for each user in time intervals where a topic is a set of DBpedia entities. Similarly, user interest profiles were represented as topic vectors where each topic is a set of temporally correlated entities on Twitter in Zarrinkalam and Kahani (2015). To this end, an entity graph based on their temporal correlation as defined by the authors was constructed, and the topics in a time interval were extracted using some existing community detection algorithms such as the *Louvain* method (Rotta and Noack 2011). The Louvain method is a simple and efficient algorithm for community detection, and relies upon a heuristic for optimizing modularity which quantifies the density of the links inside of the communities as compared to the links between communities. Subsequently, each topic  $z$  was transformed into a set of weighted entities using the *degree centrality* of an entity in the topic (community). Finally, they obtained the weight of a topic based on the weight of an entity  $c$  with respect to the topic and the frequency of  $c$  in  $u$ 's tweets.

Budak et al. (2014) proposed a probabilistic generative model to infer user interest profiles which are represented as an interest probability distribution over ODP (Open Directory Project<sup>27</sup>) categories. In their proposed approach, the authors considered three aspects such as (1) the posts of a target user, (2) the activeness of the user, and (3) the influence of friends. They assumed that time is divided into fixed time steps, and transformed the problem into inferring the probability of a user being interested in each of the interests, given a social network that evolves over time including posts and social network information. Sang et al. (2015) also proposed a probabilistic framework for inferring user interest profiles. Differing from Budak et al. (2014), Sang et al. (2015) assumed users have long- and short-term interest (topic) distributions. Long-term interests denote stable preferences of users while short-term interests denote user preferences over short-term topics of events in OSNs. However, they did not consider users' social networks.

In contrast to the aforementioned approaches, which assume all tweets posted by users are related to their interests, Xu et al. (2011) proposed a modified author-topic model (Rosen-Zvi et al. 2004) for distinguishing interest-related and unrelated tweets when learning the topic distributions of users.

---

<sup>27</sup> <https://en.wikipedia.org/wiki/DMOZ>.

#### 4.2.2 Profile enhancement

One of the advantages of constructing primitive interest profiles using concepts such as entities is that they can be further enhanced by external knowledge to deliver the final interest profiles. The approaches used in the literature for enhancing primitive user interests have mainly leveraged *hierarchical*, *graph-based*, or *collective* knowledge.

**Leveraging hierarchical knowledge** One line of approach for enhancing entity-based primitive interest profiles is apply an adapted *spreading activation* (Collins and Loftus 1975) function on a hierarchical knowledge base. For example, Kapanipathi et al. (2014) proposed representing user interest profiles as Wikipedia categories based on a hierarchical knowledge base, which is a refined Wikipedia category system built by the authors. The user interest profiles were then constructed using the hierarchical knowledge base with the following two steps. First, Wikipedia entities in users' tweets were extracted as their primitive interests. Second, these entities were used as activated nodes for applying an adapted spreading activation function on the hierarchical knowledge base in order to infer weighted categories for representing user interest profiles.

The spreading activation function proposed by Kapanipathi et al. (2014) can be applied to any case where a set of entities and a hierarchical knowledge base are available. Therefore, many studies that followed have adopted this function but with different approaches for extracting entities or with different hierarchical knowledge bases (Besel et al. 2016a, b; Große-Bölting et al. 2015; Nishioka and Scherp 2016; Piao and Breslin 2017a). For instance, Nishioka and Scherp (2016) extracted entities and applied the spreading activation function on STW, which is a hierarchical knowledge base from the economics domain. Große-Bölting et al. (2015) investigated several spreading activation functions including the one proposed in Kapanipathi et al. (2014) with the ACM CCS concept taxonomy in the computer science domain. The results showed that using a basic spreading activation function provides the best user interest profiles compared to using other ones in the context of research article recommendations.

Besel et al. (2016a, b) extracted entities by mapping followees' Twitter accounts to Wikipedia entities, and used WiBi (Flati et al. 2014) as their hierarchical knowledge base for applying the spreading activation function proposed in Kapanipathi et al. (2014). Similarly, Faralli et al. (2015b) also mapped followees' Twitter accounts to Wikipedia entities, and used them as users' primitive interests for propagation with WiBi. However, a simpler propagation strategy was adopted in Faralli et al. (2015b). In Faralli et al. (2017), the authors extended their previous work (Faralli et al. 2015a) and proposed a methodology to build *Twixonomy*, which is a Wikipedia category taxonomy. *Twixonomy* is built by using a graph pruning approach based on a variant of Edmonds optimal branching (Edmonds 1968). The authors showed that the proposed approach can generate a more accurate taxonomy compared to the one proposed in Kapanipathi et al. (2014). As we mentioned in Sect. 2.2.1, one issue with these approaches mapping followees' accounts to Wikipedia entities is that only a limited percentage of followees' accounts can be mapped to corresponding entities. For example, Faralli et al. (2015b) and Piao and Breslin (2017a) reported that only 12.7% and 10%

of followers' accounts can be mapped to Wikipedia entities. In this regard, Piao and Breslin (2017a) considered the use of followers' *biographies* for extracting entities, and applied two different propagation strategies; one is the spreading activation function from Kapanipathi et al. (2014), and the other is an interest propagation strategy exploring the DBpedia knowledge graph which will be discussed later on (Piao and Breslin 2016b).

Instead of using refined hierarchical knowledge from Wikipedia, some studies have explored other types of hierarchical knowledge bases as well. Kang and Lee (2016) proposed mapping news categories to tweets for constructing user interest profiles. The authors leveraged news categories from two popular news portals in South Korea (Naver News<sup>28</sup> and Nate News<sup>29</sup>) to build their category taxonomy. This taxonomy consists of 8 main categories and 58 sub-categories, and each category consists of all news articles in the two news corpuses. To assign categories to a tweet, each tweet and news category are represented as a term vector where the weights of terms are calculated using TF-IDF first. As there might be a semantic gap between terms in social media and news portals, the authors leveraged Wikipedia to transform the term vectors of tweets and news categories into a same vector space. The top two news categories to each tweet based on the cosine similarity between their vectors, and these news categories of a user's tweets are then aggregated to construct the final user interest profiles.

Jiang and Sha (2015) leveraged external knowledge sources such as DBpedia, Freebase (Bollacker et al. 2008), and Yago (Suchanek et al. 2007) for constructing a topic hierarchy tree, which is a hierarchical knowledge base consists of over 1000 topics distributed in 5 levels. However, the details for obtaining the topic hierarchy tree were not discussed in their study. The topic hierarchy tree used in Klout service is also bootstrapped using Freebase and Wikipedia, which consists of 3 levels with 15, around 700, and around 9000 concepts in each level, respectively (Spasojevic et al. 2014). In Bhargava et al. (2015), the authors manually built a category taxonomy based on Facebook Page categories and the Yelp<sup>30</sup> category list. The category taxonomy in Bhargava et al. (2015) consists of three levels with 8, 58, and 137 categories in each level, respectively. The authors used features such as entities, hashtags, and document categories which can be extracted from Facebook *likes* and UGC as users' primitive interests, and then measured the confidence of each concept in the category taxonomy based on these features using the Semantic Textual Similarity system (Han et al. 2013).

**Leveraging graph-based knowledge** Instead of leveraging hierarchical knowledge, many studies have leveraged graph-based knowledge for enhancing user profiles. For example, Michelson and Macskassy (2010) exploited Wikipedia categories directly for propagating a user's primitive interests. The authors summed the scores of a category which appeared in multiple depths in the category graph. Differing from exploring the categories of a specified depth (Michelson and Macskassy 2010), Siehndel and Kawase (2012) represented user interest profiles using 23 top-level categories

---

<sup>28</sup> <http://news.naver.com/>.

<sup>29</sup> <http://news.nate.com//>.

<sup>30</sup> <https://www.yelp.com/>.

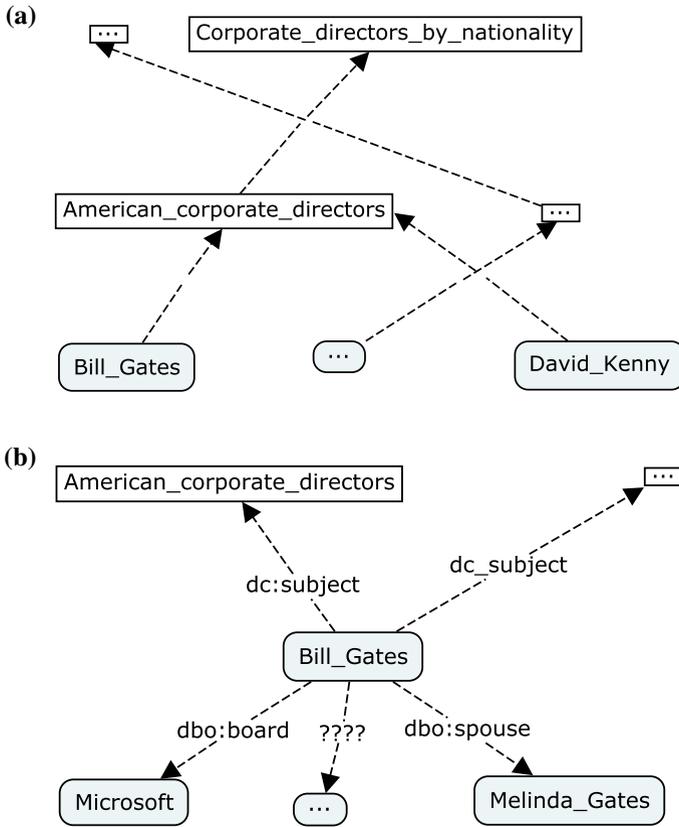


Fig. 5 Examples of WiBi taxonomy and DBpedia graph (Piao and Breslin 2017a). a WiBi taxonomy, b DBpedia graph

of the root node `Category:Main_Topic_Classifications` in Wikipedia. The Wikipedia entities in users’ tweets were extracted as their *primitive interests*, and these entities were then propagated up to the 23 top-level categories with a discounting strategy for the propagation.

With the advent of large, cross-domain Knowledge Graphs (KGs) such as DBpedia, different approaches leveraging background knowledge from KGs have been investigated. A knowledge graph is a knowledge base which consists of an ontology and instances of the classes in the ontology (Färber et al. 2015). The difference between a hierarchical category taxonomy such as WiBi and a knowledge graph such as DBpedia is displayed in Figure 5 (Piao and Breslin 2017a). As we can see from the figure, for an entity, DBpedia goes beyond just categories to provide related entities via the entity’s properties/edges. Depending on the propagation strategies for those entities in a user’s primitive interests, different aspects, e.g., *related entities*, *categories* or *classes* of the entities can be leveraged for the propagation. For example, Peñas et al. (2013) enriched categories in users’ primitive interests using similar categories defined by the `categorySameAs` relationship in DBpedia. Abel et al. (2012) proposed using back-

ground knowledge from DBpedia for propagating user interest profiles with respect to POI. The authors considered entities that were two hops away from a user's primitive interests and that were related to places. However, this approach did not consider any discounting strategy for the weights of propagated user interests. In Orlandi et al. (2012), the authors leveraged DBpedia categories one hop away from of the entities in a user's primitive interests using a discounting strategy for propagating user interests.

Although Orlandi et al. (2012) leveraged DBpedia as the knowledge base instead of Wikipedia, they still exploited categories only, which makes no difference between using DBpedia and Wikipedia. To investigate other aspects of DBpedia such as related entities and classes of primitive interests, Piao and Breslin (2016b) studied three approaches such as *category-*, *class-*, and *property-based* propagation strategies. This study found that exploiting categories and related entities via different properties of primitive interests provides the best performance compared to using corresponding categories only in the context of URL recommendations on Twitter.

An alternative graph for propagating entity-based user interest profiles is the Wikipedia entity graph. Compared to the DBpedia graph, where the edges between two entities are predefined properties in an ontology, the edges in the Wikipedia entity graph denote the mentions of the other entities in a Wikipedia entity (article). Lu et al. (2012) exploited a Wikipedia entity graph to enhance the entity-based primitive interests. Different from exploiting Wikipedia categories, the intuition behind this approach is that if a user is interested in iPhone, the user might be interested in other products from Apple, instead of being interested in other mobile phones in the same category such as Smartphones. To this end, the authors used the ESA algorithm to extract entities from the tweets of users as their primitive interests, and then expanded these entities using a random walk on the Wikipedia entity graph.

In Jipmo et al. (2017), the authors assumed there are a set of interests  $i \in I$ , e.g., Sports, Politics, etc., which the user modeling system needs to measure the corresponding weights for each interest. After building a bag of entities based on the ones extracted from a user's tweets, the relevance score of an interest  $i$  is measured as below, which can be seen as a spreading activation approach with some constraints:

$$S_i^u = \sum_{a \in BOE_u} \frac{1}{\min\{dist(a, c), c \in BOC_i\}} \quad (4)$$

where  $BOE_u$  denotes the bag of entities extracted from  $u$ 's tweets, and  $BOC_i$  denotes a set of categories containing the name of  $i$  in their titles. For example, for an interest sports,  $BOC_i$  consists of categories such as Category:Sports by year, Category:Sports in France, etc.  $dist(a, c)$  refers to the length of the shortest directed path from  $a$  to  $c$  in the Wikipedia graph.

**Leveraging collective knowledge** More recently, some studies proposed leveraging collective knowledge powered by the great amount of interest profiles of all users in a dataset, and enhancing a user profile with other related interests identified as frequent patterns in all profiles using frequent pattern mining (FPM). FPM was designed to find frequent patterns (itemsets or a set of items that appear together in a transaction dataset

frequently). In the context of user modeling, previous studies have treated each user interest as an item, each interest profile as a transaction, and all user interest profiles as the transaction dataset (Faralli et al. 2015b; Trikha et al. 2018). Trikha et al. (2018) leverages frequent pattern mining techniques to identify topic sets. Here, a topic set consists of the topics frequently appear together in user profiles. Afterwards, the other topics in the topic sets that contain the topics in a user's profile are added into that user's profile as well.

Take an example from Trikha et al. (2018), a topic set identified via FPM might consist of two topics  $z_1$  and  $z_2$ , where  $z_1 = \{\text{Mixtape, Hip\_hop\_music, Rapping, Kanye\_West, Jay-Z, Remix}\}$  and  $z_2 = \{\text{Lady\_Gaga, Song, Album, Concert, Canadia\_Hot\_100}\}$ .  $z_1$  refers to the topic about hip hop music produced by two American rappers Jay-Z and Kanye\\_West while  $z_2$  represents the topic about Lady\\_Gaga's concert in Canada. As these two topics frequently appear together in user interest profiles, the users who are interested in  $z_1$  might be also interested in  $z_2$  even  $z_2$  is not in their primitive interests. In contrast to Trikha et al. (2018), Faralli et al. (2015b) did not directly enhance user interest profiles with other interests that occur together frequently, but used FPM for user classification and recommendation. It is worth noting that both Faralli et al. (2015b) and Trikha et al. (2018) used the FP-Growth algorithm (Han and Pei 2000) for frequent pattern mining in their studies.

### 4.2.3 Temporal dynamics of user interests

User interests in OSNs can change over time, and many studies have been conducted in order to investigate the temporal dynamics of user interests in OSNs. For example, Jiang and Sha (2015) showed that the similarity of current user interest profiles with the profiles at the beginning of the observation period of their dataset is the lowest while the similarity of current profiles with the ones built in the last month is the highest. Similarly, Abel et al. (2011b) showed that a user interest profile built in an earlier week differs more from the current profile compared to one built recently. In order to incorporate the temporal dynamics of user interests into user modeling strategies, there are mainly two types of approaches: (1) *constraint-based* approaches, and (2) *interest decay functions*.

**Constraint-based approaches** Constraint-based approaches extract user interest profiles based on specified constraints, e.g., using a *temporal constraint* to build user interest profiles based on their tweets posted in the last 2 weeks or using an *item constraint* to construct user profiles based on the last 100 tweets of the users. For example, Abel et al. (2011b) investigated several temporal constraints such as *long- and short-term*, and *weekend* in their user modeling strategies on Twitter for a news recommender system. *Long-term* profiles extract user interests from entire historical tweets of users while *short-term* profiles extract user interests from tweets posted within the last 2 weeks. They showed that long-term entity-based profiles outperform short-term ones in the context of news recommendations. User interests can be different within different time frames such as during the week or on the weekends. The experimental results in Abel et al. (2011b) also showed that entity-based interest pro-

files based on their tweets posted on weekends can outperform long-term profiles for recommending news on weekends.

Some interests of users such as professional interests are stable while other interests such as the ones related to a certain event can be temporary. A user modeling strategy can apply temporal dynamics selectively to different information sources based on their characteristics. This type of strategy has been adopted in practical user modeling systems such as the one in Klout (Spasojevic et al. 2014), in which a 90 day window is used for capturing the temporal dynamics of user interests for some temporal information sources, and an all-time window is used for more permanent sources such as professional interests.

Nishioka and Scherp (2016) compared both constraint-based approaches and interest decay functions for constructing user interest profiles on Twitter in the context of publication recommendations. Differing from the results in the domain of news (Abel et al. 2011b), results from Nishioka and Scherp (2016) showed that a constraint-based approach constructing user interest profiles within a certain period performs better than using an interest decay function in the context of publication recommendations.

**Interest decay functions** Constraint-based approaches include interests which meet predefined constraints, and exclude other interests completely. Instead of constructing user interest profiles in a certain period (e.g., short-term), or based on temporal patterns (e.g., weekends), interest decay functions aim at including all the interests of a user but decaying old ones. The intuition behind those interest decay functions is that a higher weight should be given to recent interests than old ones.

A popular type of interest decay function applies exponential decay to user interests. For example, the interest decay function from Orlandi et al. (2012) is defined as follows:

$$x(t) = x_0 \cdot e^{-t/\beta} \quad (5)$$

Here,  $x(t)$  is the decayed weight at time  $t$ , and  $x_0$  denotes the initial weight (at time  $t = 0$ ). This interest decay function also has an initial time window (7 days), and the interests in the time window are not discounted. The authors in Orlandi et al. (2012) set  $\beta = 360$  days and  $\beta = 120$  days for their experiment, and showed that using  $\beta = 360$  days performs better than using  $\beta = 120$  days in terms of an evaluation based on a user study. We use `decay(Orlandi)` to denote this approach in this study. A similar decay function was used in Bhargava et al. (2015) and Nishioka and Scherp (2016), where a weight for the last update was used instead of initial weight (Bhargava et al. 2015). In O'Banion et al. (2012), the authors also used an exponential decay function:  $x(t) = x_0 \cdot 0.9^d$  where  $d$  is the difference in days between the current date and the date that a concept was mentioned.

Abel et al. (2011a) also proposed a time-sensitive interest decay function, which is denoted by `decay(Abel)` in this survey. The weight of an entity  $e$  with respect to a user  $u$  at a specific time is measured as below.

$$w(e, time, T_{tweets,u,e}) = \sum_{t \in T_{tweets,u,e}} \left( 1 - \frac{|time - time(t)|}{max_{time} - min_{time}} \right)^d \tag{6}$$

where  $T_{tweets,u,e}$  denotes the set of tweets mentioning  $e$  that have been posted by  $u$ .  $time(t)$  denotes the timestamp of a given tweet  $t$ , and  $max_{time}$  and  $min_{time}$  denote the highest (youngest) and lowest (oldest) timestamp of a tweet in  $T_{tweets,u,e}$ . In addition, the parameter  $d$  determines the influence of the temporal distance ( $d = 4$  in Abel et al. 2011a). In contrast to the aforementioned exponential decay functions, this approach incorporates the age of an entity  $e$  at the recommendation time, and the time span of  $e$  with respect to  $u$ .

In order to compare different interest decay functions in the context of user modeling in OSNs, Piao and Breslin (2016b) investigated three interest decay functions for constructing user interest profiles on Twitter including `decay(Abel)` and `decay(Orlandi)`. The other one is a modified interest decay function from Ahmed et al. (2011), which was used in advertisement recommendations on web portals (i.e., Yahoo!<sup>31</sup>). The modified interest decay function used in Piao and Breslin (2016b) is defined as follows:

$$w_{ik}^t = \mu_{2week} w_{ik}^{t,week} + \mu_{2month} w_{ik}^{t,month} + \mu_{all} w_{ik}^{t,all} \tag{7}$$

where  $\mu_{2week} = \mu$ ,  $\mu_{2month} = \mu^2$  and  $\mu_{all} = \mu^3$  where  $\mu = e^{-1}$ . This decay function combines three levels of abstractions where the decay of user interests in each abstraction is  $\mu$  times the previous abstraction. We use `decay(Ahmed)` to denote this approach in this survey. Piao and Breslin (2016b) conducted a comparative study of user interest profiles constructed based on the three aforementioned interest decay functions and the profiles based on *short-* and *long-term* periods. Those interest profiles were then evaluated in the context of URL recommendations. The results showed that using `decay(Ahmed)` and `decay(Orlandi)` have competitive performance in terms of URL recommendations, and perform better than using `decay(Abel)` as well as *short-* and *long-term* profiles which were constructed without any interest decay. In addition, the experimental results indicate that although the performance increases by giving a higher weight to recent user interests, it starts decreasing once the weight of recent interests is too high. That is, although applying the decay function to recent user interests increases the performance, we still need the old history in order to provide the best performance in the context of URL recommendations.

Instead of considering the temporal dynamics of user interests with respect to individual users, global trends in an OSN can be incorporated into a user modeling strategy. In Gao et al. (2011), the authors combined user interests from tweets of a target user (user profiles) and of all users (trend profiles) for constructing user interest profiles. The TF weighting scheme is used for constructing user profiles. For trend profiles, they applied a time-sensitive TF-IDF (t-TF-IDF) weighting scheme to concepts:

$$w_{t-TF-IDF}(I_j, c) = w_{TF-IDF}(I_j, c) \cdot (1 - \hat{\sigma}(c)) \tag{8}$$

<sup>31</sup> <https://yahoo.com/>.

where  $w_{TF-IDF}(I_j, c)$  denotes the TF-IDF score of a concept  $c$  in a given time interval  $I_j$ , and  $\hat{\sigma}(c)$  denotes the normalized standard deviation of timestamps of tweets that refer to  $c$ . Kanta et al. (2012) further incorporated location-aware trends into the trend-aware user modeling approach in Gao et al. (2011) to improve the performance of inferred user interest profiles in the context of news recommendations.

### 4.3 Summary and discussion

This section reviewed a number of approaches for constructing and enhancing user interest profiles. Table 6 summarizes the approaches discussed in this section in terms of the three dimensions: (1) weighting schemes for constructing primitive interests, (2) approaches for incorporating the temporal dynamics of user interests, and (3) profile enhancement methods.

As we can see from the table, many studies have incorporated the temporal dynamics of user interests in their user modeling strategies. Among interest decay functions, exponential decay functions such as `decay(Orlandi)` have been adopted widely. When incorporating the temporal dynamics of user interests, it is important to choose constraint-based approaches or interest decay functions based on the purpose of user modeling. For instance, when using inferred user interest profiles for recommending items such as news or URLs in OSNs, interest decay functions perform better than constraint-based approaches such as short- and long-term profiles (Piao and Breslin 2016b). However, the results from Nishioka and Scherp (2016) indicate that a constraint-based approach based on a certain period for profiling outperforms the one applying exponential decay for building user profiles in the context of a publication recommender system. One possible explanation is that user interests change differently with respect to different domains. For example, user interests should be adapted to their recent interests for news or URL recommendations, however, user interests with respect to research may not.

Jiang and Sha (2015) also pointed out that users have two types of interests; (1) *stable interests* (which they call primary interests in Jiang and Sha 2015), and (2) secondary interests. The stable interests of a user are original preferences inherent to that user, such as programmers who like efficient algorithms or lawyers who like debate, etc. (Jiang and Sha 2015). In contrast, secondary interests are temporary ones which closely follow hot topics or events in a specific timespan. This is in line with the user modeling strategy used in Klout (Spasojevic et al. 2014), which applies a short-term window for capturing user interests that are temporary and uses a long-term window for more stable user interests.

Different types of knowledge from various knowledge bases have been leveraged for enhancing the primitive interests of users. The diversity of KBs and the different structures of hierarchical KBs indicate the complexity of representing knowledge in KBs as well. Table 7 summarizes the differences between hierarchical KBs used in the literature. For instance, the constructed Wikipedia category taxonomy in Kapanipathi et al. (2014) consists of 15 levels with 802,194 categories while the topic hierarchy tree built by Jiang and Sha (2015) consists of 5 levels with over 1000 topics. The topic hierarchy tree used in Klout has 3 levels which consists of 15 main categories, around

**Table 6** Approaches for constructing and enhancing user interest profiles

Weighting scheme Heuristics	Temporal dynamics		Profile enhancement		Examples
	Prob. inference	Constraint-based	Int. dec. functions	Hie. Gra. Col.	
✓					Hannon et al. (2012), Kapanipathi et al. (2011), Lim and Datta (2013), Narducci et al. (2013), Abel et al. (2011c), Bhattacharya et al. (2014), Zarrinkalam and Kahani (2015), Zarrinkalam et al. (2016), Chen et al. (2010), Garcia Esparza et al. (2013), Ahn et al. (2012), Yu and Perez (2013), Phelan et al. (2009), Karatay and Karagoz (2015), Weng et al. (2010), Xu et al. (2011) Abel et al. (2011b, 2013a)
✓		✓			Gao et al. (2011), Kanta et al. (2012), Abel et al. (2011a), O'Banion et al. (2012)
✓			✓		Sang et al. (2015)
✓		✓		✓	Große-Börling et al. (2015)
✓				✓	Lu et al. (2012), Abel et al. (2012), Peñas et al. (2013), Piao and Breslin (2016d, 2017b)
✓				✓	Kapanipathi et al. (2014), Jipmo et al. (2017), Michelson and Maeskassy (2010), Kang and Lee (2016), Besel et al. (2016a, b), Faralli et al. (2017), Nechaev et al. (2017), Jiang and Sha (2015), Stehndel and Kawase (2012)

Table 6 continued

Weighting scheme Heuristics	Temporal dynamics		Profile enhancement		Examples		
	Prob. inference	Constraint-based	Int. dec. functions	Hie.		Gra.	Col.
✓				✓	✓		Piao and Breslin (2017a)
✓				✓	✓	✓	Faralli et al. (2015b)
✓	✓	✓		✓		✓	Trikha et al. (2018)
✓	✓	✓	✓	✓			Bhargava et al. (2015)
✓				✓			Nishioka et al. (2015), Budak et al. (2014)
✓		✓			✓		Spasojevic et al. (2014)
✓			✓		✓		Orlandi et al. (2012), Piao and Breslin (2016b, c)
✓		✓	✓	✓			Nishioka and Scherp (2016)

Prob. probability, Int. dec. interest decay, Hie. hierarchical knowledge, Gra. graph-based knowledge, Col. collective knowledge

**Table 7** The structures of hierarchical knowledge bases for representing topics in different studies

Study	# Levels	# Topics	Details
Kapanipathi et al. (2014)	15	802,194	N/A
Jiang and Sha (2015)	5	~ 1000	N/A
Spasojevic et al. (2014)	3	~ 10,000	15 → ~ 700 → ~ 9000
Kang and Lee (2016)	2	66	8 → 58
Bhargava et al. (2015)	3	203	8 → 58 → 137

The final column shows the number of concepts in each level of the corresponding hierarchical knowledge base

700 sub-categories, and around 9000 entities (Spasojevic et al. 2014). A concept taxonomy built manually by referring to external websites such as news portals or Facebook Page categories has less complexity compared to a taxonomy based on KBs such as Wikipedia. For example, the category taxonomy built based on news portals (Kang and Lee 2016) has 8 main categories and 58 sub-categories. The one built based on Facebook and Yelp categories (Bhargava et al. 2015) also has 8 and 58 categories for the top-2 levels with an additional 137 categories in its third level. We can observe that the hierarchical knowledge bases used in practice or built based on taxonomies used in practice tend to have a small number of levels (2–5). Applying a spreading activation function, even the same one, to those different taxonomies might have different results. There is a lack of comparison of different hierarchical knowledge bases and their effect in the context of inferring user interest profiles.

Furthermore, although some studies investigated the comparison between using different KBs such as Wikipedia categories and the DBpedia graph, there was no comparative study on exploiting the Wikipedia entity graph (Lu et al. 2012), categories in other KBs such as ODP, and the DBpedia graph. In addition, despite the fact that different KBs might be useful in different domains (Nguyen et al. 2015), enhancing user interests based on other KBs such as Wikidata (Vrandečić and Krötzsch 2014), or BabelNet (Navigli and Ponzetto 2012) has not been fully explored.

## 5 Evaluation approaches

### 5.1 Overview

In this section, we describe evaluation approaches used for evaluating different user interest profiles that are generated by different user modeling strategies in the literature. User modeling is one of the main building blocks in many adaptive systems such as recommender systems. Many previous studies on the evaluation of adaptive systems suggested that it is important to evaluate different blocks separately in order to identify the problems in the adaptive systems (Brusilovsky et al. 2001; Paramythis et al. 2010). Gena and Weibelzahl (2007) provided a list of methods for evaluating adaptive systems, where some of them can be used for evaluating the quality of user modeling component as well. These evaluation methods include (1) *questionnaires*, (2) *interviews*, and (3) *logging use*.

**Questionnaires** Questionnaires consist of pre-defined questions, which can be in different styles such as scalar or multi-choice, and ranked (Gena and Weibelzahl 2007). In our context, this approach can be used for collecting users' explicit feedback about their interest profiles for evaluation. To this end, this approach requires recruiting users for the experiment of building user interest profiles with their OSN accounts. At the end of the experiment, these users can provide feedback on user interest profiles constructed by different user modeling strategies.

**Interviews** The second approach is used to collect users' opinions and experiences, preferences and behavior motivations (Gena and Weibelzahl 2007) with respect to adaptive systems. Interviews can be used after building users' interest profiles to gather their opinion such as satisfaction and accuracy about the inferred user interest profiles. Compared to questionnaires, interviews are more flexible but more difficult to be administered. Therefore, this method has not been exploited for evaluating user modeling strategies in the literature.

**Extrinsic evaluation (logging use)** This approach uses the actions of users in the context of adaptive systems for evaluation, e.g., whether a user liked a recommend item in a recommender system. This can be considered an extrinsic way of evaluating user interest profiles in terms of the performance of applications where these profiles are applied. For example, one common approach is using constructed user interest profiles as an input to a recommender system, and adopting some well-established evaluation metrics of recommender systems for measuring the quality of user interest profiles indirectly. Manual analysis is sometimes used together with other evaluation approaches. In this case, the authors present some examples of user interest profiles built for several users (e.g., some representative users on Twitter such as *Barack Obama*), and discuss the quality of profiles with respect to these users.

## 5.2 Review

### 5.2.1 Evaluation based on Questionnaires

A common approach for evaluating constructed user interest profiles is based on a user study with questionnaires. For example, Narducci et al. (2013) evaluated user interest profiles built for 51 users from Facebook and Twitter based on their feedback on two aspects: *transparency* and *serendipity* using a 6-point discrete rating scale. The first aspect aims to evaluate to what extent the keywords in the profile reflect personal interests, and the second one aims to measure to what extent the profile contains unexpected interesting topics. Similarly, Kapanipathi et al. (2014) recruited 37 users and built category-based user interest profiles based on their tweets on Twitter. Afterwards, the 37 users provided explicit feedback, e.g., Yes/Maybe/No with respect to the categories in those profiles. Similar approaches have been used in Bhattacharya et al. (2014), Besel et al. (2016a, b), Budak et al. (2014), and Orlandi et al. (2012). However, instead of recruiting volunteers for an experiment, the authors in Budak et al. (2014) first inferred user interest profiles for 500 randomly chosen users on Twitter,

and emailed them using the email addresses in their profiles to get feedback about their inferred interests. Instead of using the feedback from target users for inferred user interest profiles, Kang and Lee (2016) and Michelson and Macskassy (2010) labeled user interests themselves or used recruited annotators.

Explicit feedback can be obtained in a system which has user interest profiles that can be modified by users. For example, Garcia Esparza et al. (2013) implemented a stream filtering system where users are represented based on 18 defined categories such as `Music` and `Sports`. For evaluation, the authors asked each participant to give explicit feedback on their profiles by deleting or adding categories that they felt were incorrect or missing.

In contrast to obtaining explicit feedback on inferred user interest profiles, a user study can be conducted on the performance of a specific application where those inferred user interest profiles play an important role. For example, Chen et al. (2010) conducted a user study with respect to a URL recommender system on Twitter, which is based on the inferred user interest profiles. Therefore, instead of directly giving feedback on the constructed user interest profiles, the users participating in the study were given URL recommendations, and they marked each URL as one of their interests or not. Similarly, Nishioka and Scherp (2016) obtained explicit feedback from users on publication recommendations based on their interest profiles. These user studies can also be considered as extrinsic evaluation, which we will discuss in the next section, as they are not evaluating user interest profiles directly.

**Pros and cons** Evaluation approaches based on the explicit feedback of profiled users with respect to their interest profiles would arguably be the most direct and accurate way for evaluating those profiles. However, this also requires recruiting volunteers and imposes an extra burden for users, and therefore limits the number of participants for evaluation (e.g., 37 users were recruited for evaluation in Kapanipathi et al. 2014).

### 5.2.2 Extrinsic evaluation

To evaluate the quality of inferred user interest profiles without imposing an extra burden on users, offline evaluation in terms of the performance of a specific application has been used. In this case, user interest profiles are used as an input to an application such as a news recommender system where these profiles play an important role. Afterwards, different profiles created by different user modeling strategies are compared in terms of the recommendation performance using each profile. The recommendation performance can be evaluated by well-established evaluation metrics for recommender systems such as *mean reciprocal rank* (MRR) which denotes at which rank the first item relevant to the user occurs on average, *success at rank N* (S@N), which stands for the mean probability that a relevant item occurs within the top-N recommendations, and well-known *precision* and *recall*. For a complete list of evaluation metrics and their details we refer the reader to Bellogn et al. (2017) and Herlocker et al. (2004) respectively.

For instance, Abel et al. (2011b) evaluated three different user modeling strategies in terms of S@N and MRR in the context of news recommendations, and Spasojevic et al. (2014) evaluated their user modeling strategy in terms of precision and recall

in the context of topic recommendations on Klout. Similarly, Sang et al. (2015) also evaluated user interest profiles in terms of news recommendations in addition to tweet recommendations. Piao and Breslin (2016b,c,d, 2017a,b) evaluated different user modeling strategies in the context of URL recommendations on Twitter where the set of ground truth URLs is those shared by users on Twitter in the last 2 weeks. In Faralli et al. (2015b), the authors evaluated user interest profiles in terms of user classifications and recommendations. For the classification task, the user interest profiles were used for classifying each user to the appropriate label, e.g., Starbucks fan. For the recommendation task, the authors evaluated the performance of leveraging different hierarchical levels of interests with respect to interest recommendations using itemset mining.

In contrast to previous studies which have focused on inferring user interest profiles, Nechaev et al. (2017) focused on users' privacy and evaluated different followee-suggestion strategies for concealing user interests which can be inferred from users' activities in OSNs based on state-of-the-art user modeling strategies.

**Pros and cons** Extrinsic evaluation provides an offline setting for evaluating inferred user interest profiles. Therefore, it facilitates the evaluation process of different user modeling strategies as these strategies are evaluated based on a collected dataset (or logs). However, this approach does not directly evaluate the inferred user interest profiles, and lacks the opinions of users with respect to the inferred interest profiles. There are other evaluation approaches used in some studies besides the aforementioned two methods. For example, Abel et al. (2011c) compared the number of distinct entities and topics in user interest profiles for evaluating news-based enrichment of their tweets. In Faralli et al. (2017), the authors run two experiments to evaluate their approach of building interest taxonomies. First, they compared their approach against other approaches proposed for constructing user interest taxonomies using other gold standard taxonomies. Second, they provided samples of generated user interest profiles, and compared inferred Wikipedia categories with respect to several users based on different user modeling strategies. Similarly, Xu et al. (2011) evaluated their topic modeling approach by comparing it against other topic modeling methods in terms of *perplexity*, and then discussed some user interest profiles produced by different approaches. User interest profiles have also been used for specific applications such as followee, tweet, and news recommendations (Chen et al. 2012; Hong et al. 2013; Phelan et al. 2009; Weng et al. 2010), where user modeling strategies were not evaluated or compared to other alternatives.

### 5.3 Summary and discussion

In this section, we reviewed different evaluation approaches that have been used in the literature for evaluating constructed user interest profiles. Table 8 provides a summary of previous studies in terms of evaluation methods.

Evaluating user interest profiles based on a user study is important for understanding different aspects of user interests, e.g., abstraction levels of user interests. For example, Orlandi et al. (2013) studied the specificity of user interests and evaluated it based on

**Table 8** Evaluation approaches for constructed user interest profiles

Questionnaires	Extrinsic evaluation	Examples
✓		Kapanipathi et al. (2014), Kang and Lee (2016), Michelson and Macskassy (2010), Budak et al. (2014), Bhattacharya et al. (2014), Besel et al. (2016a, b), Orlandi et al. (2012), Narducci et al. (2013), Bhargava et al. (2015), Garcia Esparza et al. (2013), Vu and Perez (2013), Ahn et al. (2012), Chen et al. (2010), Nishioka and Scherp (2016)
	✓	Abel et al. (2011a, b, c, 2012), Chen et al. (2010), Zarrinkalam and Kahani (2015), Sang et al. (2015), Kanta et al. (2012), O'Banion et al. (2012), Piao and Breslin (2016b, c, d, 2017a, b), Lu et al. (2012), Sang et al. (2015), Gao et al. (2011), Karatay and Karagoz (2015), Trikha et al. (2018), Nishioka et al. (2015), Große-Bölting et al. (2015), Zarrinkalam et al. (2016), Ahn et al. (2012), Spasojevic et al. (2014), Jipmo et al. (2017), Faralli et al. (2015b), Nechaev et al. (2017)

a user study, which showed that users prefer to give a higher score over non-specific entities. However, the extra effort of recruiting users and gaining feedback from them is time consuming, and limits the scale of users for evaluation. The evaluation in terms of the performance of a specific application has the advantage of its offline setting and using a relatively larger number of users compared to a user study. Both evaluation approaches can be used in an appropriate way for designing and evaluating user modeling strategies. For example, based on a user study on the specificity of user interests (Orlandi et al. 2013), we can design ways to incorporate the feedback from users' preferences regarding non-specific entities into a user modeling strategy, and evaluate the strategy at a large scale in offline settings based on a collected dataset such as the one from Twitter.

One of the challenges of the offline evaluation in terms of the performance of a specific application is the lack of benchmarks that are freely available (Faralli et al. 2015b). Despite the openness of some microblogging services such as Twitter, it is time consuming to collect all data used in different user modeling approaches, e.g., tweets, list memberships, biographies of followees/followers in addition to the information about users. In addition, different datasets with different user sizes might produce different results even using the same user modeling strategies for comparison. It is also important to evaluate different user interest profiles in the context of different applications beyond a specific one. For example, in Manrique and Mariño (2017), the authors showed that user interest profiles based on different user modeling strategies perform differently in the context of recommending articles based only on titles, abstracts, and full texts. Although the study (Manrique and Mariño 2017) is in the context of research article recommendations, it is highly likely that different user interest profiles from microblogging services will have different levels of performance based on the applications in which these profiles are applied.

## 6 Conclusions and future directions

In previous sections, we reviewed the state-of-the-art approaches used in different user modeling stages for inferring user interest profiles, which is beneficial both for researchers who are interested in user modeling in the social networks domain as well as those researchers in some other domains. It is also useful for third-party application providers who aim to utilize user interest profiles via social login functionalities in terms of providing personalized services for their users.

In this final section, we conclude this paper in Sect. 6.1 with respect to the four dimensions of inferring user interest profiles: (1) data collection, (2) representations of user interest profiles, (3) construction and enhancement of user interest profiles, and (4) the evaluation of the constructed profiles. In Sect. 6.2, we first review what progress has been made to date since Abdel-Hafez and Xu (2013), and then outline some opportunities and challenges for inferring user interests on microblogging social networks which we envision can inspire future directions in this research field.

### 6.1 Conclusions

To sum up, user activities such as the tweets posted by users are the most widely used information source for inferring user interests. However, many recent studies have started exploring other information sources such as the social networks of users as an alternative to user activities as the passive usage of OSNs is on the rise. Regarding the representations of user interest profiles, a clear tendency of leveraging concepts such as DBpedia entities or categories can be observed given their advantages of using background knowledge about those concepts from a KB. In addition to leveraging the hierarchical or graph-based knowledge of a KB for enriching user interests, several recent studies also have shown the effectiveness of leveraging collective knowledge for enriching user interest profiles (Faralli et al. 2015b; Trikha et al. 2018). With respect to incorporating the temporal dynamics of user interests, there is no single best method for inferring user interests with different purposes. Instead, one should choose constraint-based or interest decay functions based on the application needs, and the characteristics of items. For evaluating user interest profiles, both questionnaires and extrinsic evaluation strategies have been adopted at comparable levels of popularity.

### 6.2 Future directions

In Abdel-Hafez and Xu (2013), the authors proposed three future directions with respect to user modeling in OSNs, which requires (1) more dynamicity, (2) more enrichment, and (3) more comprehensiveness. On the one hand, we observe that there have been many efforts towards the second direction. These efforts include leveraging the collective knowledge powered by all users (Faralli et al. 2015b; Trikha et al. 2018) for enriching the interest profiles of each user, and the comparison between different KBs for enriching user interests (Piao and Breslin 2017a). On the other hand, the first and third directions proposed by Abdel-Hafez and Xu (2013) have not made much progress. For example, Abdel-Hafez and Xu (2013) proposed incorporating

more dynamicity with respect to user interest profiles with some assumptions such as different topics might decay with different speed, and the interest weights of each user can have different weights in different context. On top of the directions proposed by Abdel-Hafez and Xu (2013) and the recent studies we reviewed in this paper, we further proposed several future directions which are related to:

- mining user interests;
- multi-faceted user interests;
- comprehensive user modeling;
- evaluation of user modeling strategies.

**Mining user interests** To better infer user interests, researchers have proposed various approaches such as enriching short content, filtering noise in UGC, and exploring social networks. Many studies have adopted traditional weighting schemes from information retrieval such as TF or TF-IDF to somehow filter the noise in UGC for mining user interests. However, some studies have shown that incorporating some special characteristics of the services (e.g., temporal dynamics, short content) into the design of a weighting scheme can improve the quality of user interest profiles. For example, TI-TextRank which combines TF-IDF and TextRank performs better than either of them on their own as a weighting scheme for user modeling on Twitter. In this regard, more weighting schemes adapted towards microblogging services should be investigated, e.g., combining different weighting schemes used in the literature. Furthermore, mining interest-related items from data sources such as posts (e.g., Xu et al. 2011) can be useful as microblogging services have multiple usages such as information seeking, sharing and social networking (Java et al. 2007).

In addition, more sophisticated approaches for understanding the semantics of UGC are required. For example, for those approaches that rely on extracted entities for inferring user interest profiles, extracting entities from microblogs is a fundamental step which is challenging by itself. Only a few studies have considered the uncertainty (confidence) of the extracted entities, which we think might impact the overall quality of the primitive interests of users as well as the enhanced ones. Moreover, most approaches have extracted explicitly mentioned entities based on NLP APIs such as Tag.Me,<sup>32</sup> Aylien,<sup>33</sup> OpenCalais, etc. However, there can be many entities implicitly mentioned in tweets. In Perera et al. (2016), the authors showed that over 20% of mentions of movies are implicit references, e.g., a tweet referring the movie *Gravity*—“ISRO sends probe to Mars for less money than it takes Hollywood to make a movie about it”. It shows that advanced methods for extracting entities, such as the one proposed in Perera et al. (2016), have great potential to improve the quality of user modeling. Also, considering the context of a microblog might be useful when extracting entities instead of just considering the single microblog of a user. The context might refer to some previous microblogs posted by the user, or other microblogs with the same hashtag in the microblogging service. For example, Shen et al. (2013) showed that the quality of entity extraction can be improved by incorporating user interests as contextual information. Furthermore, promising results from recent studies (Far-

<sup>32</sup> <https://tagme.d4science.org/tagme/>.

<sup>33</sup> <https://aylien.com/>.

alli et al. 2015b; Trikha et al. 2018) indicate that leveraging collective knowledge via frequent pattern mining approaches is also effective in inferring implicit user interests.

**Multi-faceted user interests** There exists various aspects/views of users based on different dimensions of user modeling such as the data source, representation level, and temporal dynamics of user interests. Although many studies represent an individual user using a single user interest profile, we believe that multi-faceted user interest profiles should be given more attention as some previous studies have also shown their efficiency compared to a single model. It is not necessary to maintain several user interest profiles for a single user, but a single model can also be built with relevant information from different aspects, and a view/aspect made for the user based on the information needs for different applications. GeniUS (Gao et al. 2012) is a good example in this regard, which is a user modeling library that stores concept-based user interest profiles using the RDF<sup>34</sup> format (a W3C recommendation) with widely used ontologies such FOAF (Brickley and Miller 2012), SIOC,<sup>35</sup> and WI.<sup>36</sup> In GeniUS, user interest profiles are represented as DBpedia entities and enriched by background knowledge such as the type (domain) of an entity from DBpedia. Therefore, the constructed profile is flexible enough to retrieve its sub-profiles with respect to specific domains (e.g., *MUSIC*), which is useful for recommending domain-specific items. The idea is that, for example, we only need your music-related interest profile in the context of music recommendations. The results in Gao et al. (2012) indicate that domain-specific profiles clearly outperform the whole user profiles for domain-specific tweet recommendations in terms of six different domains. Although GeniUS only considers different views of users in terms of topical domains, the same idea can be extended to other views. For instance, different user profiles can be extracted dynamically with different approaches for incorporating temporal dynamics, e.g., retrieving short-term profiles for recommending tweets during an event, which might be more useful compared to using long-term profiles. Also, multiple user interest profiles in terms of representation level using different interest formats have been used in other domains such as personal assistants (Guha et al. 2015), which can be useful for user modeling in microblogging services as well. In Guha et al. (2015), several user interest profiles based on different representations such as keywords and Freebase entities were constructed.

**Comprehensive user modeling** In the previous survey on user modeling (Abdel-Hafez and Xu 2013), the authors also suggested that more comprehensive user modeling strategies should be investigated by considering different dimensions of user modeling together. Many of the previous studies have ignored some of the dimensions such as temporal dynamics (e.g., Phelan et al. 2009). Investigating the synergistic effect of different dimensions is important for developing better user modeling strategies, which is crucial for the performance of applications. To this end, several research questions should be answered such as “which combinations of different approaches

---

<sup>34</sup> <https://www.w3.org/RDF/>.

<sup>35</sup> <http://sioc-project.org/>.

<sup>36</sup> <http://smiy.sourceforge.net/wi/spec/weightedinterests.html>.

in each dimension can provide the best user interest profiles” or “does a dimension really matter in the context of the combination for providing the best performance?”. For example, Piao and Breslin (2016c) showed that a rich representation of user interests (using WordNet synsets and DBpedia entities) and enriching short content with the text of embedded URLs are the most important factors followed by temporal dynamics in the context of URL recommendations on Twitter. However, enhancing user interest profiles has little effect when we have a rich representation or enriched content of microblogs. Similar results have been observed in the context of inferring research interests of users based on their publications (Manrique and Mariño 2017). The results in Manrique and Mariño (2017) indicate that enhancing primitive interests can improve the performance when only short texts (e.g., titles) are available but not in the case when longer texts (e.g., full texts of publications) are available. We believe that these studies are good starting points for some future works, e.g., using different user interest profiles for different data sources instead of using a single representation of an individual user for the combination.

In addition, other user modeling dimensions which have been proposed in other domains can be considered in the social media domain as well. For example, a *scrutable* user model proposed in the context of teaching, which aims to let users have the right and possibility to have access to and control their user profiles (Carmagnola et al. 2011; Holden and Kay 1999; Kay 2006), can be a promising dimension to be incorporated into user modeling strategies in OSNs and merits further investigation and evaluation.

**Evaluation of user modeling strategies** As we mentioned in Sect. 5.3, the lack of common benchmarks and datasets hinders comparison with other approaches, which ends up with several studies directly comparing to results reported in previous studies (Faralli et al. 2015b). This does not reflect a correct comparison due to the difference of datasets in terms of platforms as well as user sizes. However, it is also challenging due to the regulations of microblogging services such as Twitter,<sup>37</sup> and the differences in data sources used in each study. Another possible direction is providing all proposed approaches as user modeling libraries that are publicly available, in the same way as GeniUS and TUMS,<sup>38</sup> so that other researchers can easily reimplement the approaches proposed in previous studies for comparison.

It is also important to evaluate inferred user interest profiles in terms of multiple tasks or different settings to understand the strengths and weaknesses of different user interest profiles. For instance, Nishioka et al. (2015) showed that considering the temporal dynamics of user interests has a positive influence on a computer science dataset but not on a medicine dataset. Manrique and Mariño (2017) showed that different user modeling strategies work differently for different types of texts that are available in the context of research article recommendations. In this regard, evaluating the performance of different user modeling strategies based on different datasets or settings can provide a clear understanding of when to use what types of user profiles, which is important for researchers in different domains as well as third-party appli-

<sup>37</sup> Twitter restrict developers from sharing the content of tweets, see <https://developer.twitter.com/en/developer-terms/agreement-and-policy>.

<sup>38</sup> Both GeniUS and TUMS are available at <http://www.wis.ewi.tudelft.nl/tweetum/>.

cation providers with different types of content to be personalized. A recent work by Tommaso et al. (2018) provides a user interests dataset which is useful in this context. It includes half million Twitter users with an average of 90 multi-domain preferences per user on music, books, etc., where those preferences are extracted from multiple platforms based on the messages of those Twitter users who also use Spotify,<sup>39</sup> Goodreads,<sup>40</sup> etc.

Finally, previous studies have adopted accuracy and ranking metrics such as precision, recall, and MRR for the extrinsic evaluation of inferred user interest profiles. However, non-accuracy metrics such as serendipity, novelty, and diversity have received increasing attention in recommender systems (Bellogn et al. 2017; Kamin-skas and Bridge 2016). Therefore, it is worth investigating the effect of different user modeling strategies and their inferred interest profiles in the context of recommender systems in terms of those non-accuracy metrics.

**Acknowledgements** This publication has emanated from research conducted with the financial support of Science Foundation Ireland (SFI) under Grant Number SFI/12/RC/2289 (Insight Centre for Data Analytics). Thanks for the anonymous reviewers and the editor for their constructive feedback to improve this work.

## Appendices

### Appendix A The list of surveyed works

#### A.1 Search strategy

In order to draw up a list of search terms, the basic terms are extracted from primary articles are retrieved. After that, other search terms are obtained iteratively based on the keywords that were used interchangeably within the retrieved articles. Overall, the final list of terms used for searching articles is presented in Table 9. These search terms (ST) are used for constructing sophisticated search strings. For example, the search string can be constructed as ST1 AND ST3 while ST1 is a compound term from Term1 and Term2 (e.g., inferring user interests). Initial searches with these search terms for titles and abstracts from electronic databases can obtain many relevant articles but may not be sufficient (Kitchenham 2004). In this regard, additional article candidates are obtained by checking the reference list from primary studies that are relevant, and searching relevant journals and conference proceedings. Abdel-Hafez and Xu (2013) provided a review of user modeling in social media websites in 2013, which includes some approaches with respect to inferring user interests in the context of microblogging social networks. In addition to those approaches mentioned in Abdel-Hafez and Xu (2013), we also review recent user modeling approaches for inferring user interests.

---

<sup>39</sup> <https://www.spotify.com>.

<sup>40</sup> <https://www.goodreads.com/>.

**Table 9** Search terms used in the search strategy of this survey

	Term1	Term2
ST1	Inferring, modeling, predicting	(User) interests
ST2	User (interest)	Modeling, profiling, detection
ST3	Social, online, twitter, microblogging	

## A.2 Selection criteria

In order to assess and select relevant articles from primary studies, inclusion and exclusion criteria should be defined based on the research questions (Kitchenham 2004). The inclusion criteria are as follows:

1. Published in English from 2004.
2. Studies on microblogging social networks.
3. Focus on user modeling strategies for inferring user interest profiles.

On the other hand, exclusion criteria can be defined as follows:

1. Studies that were not peer-reviewed or published.
2. Studies related to user modeling but not focus on microblogging social networks.
3. Studies related to user modeling, but not focus on inferring user interests.

Finally, inclusion or exclusion decisions are made for the fully obtained articles and those papers that only meet our criteria are selected. As a result, 51 articles are selected in this survey. These articles are distributed from 2010 to 2018, and the majority of them were published in conferences or workshops such as WI, UMAP, CIKM, and ECIR.

## A.3 Surveyed studies

The surveyed 51 works are retrieved from different journals, conferences, and workshops, mainly in the user modeling, recommender systems, and Web related fields as follows:

1. Journals
  - ACM SIGAPP Applied Computing Review: Besel et al. (2016b)
  - Web Semantics: Science, Services and Agents on the World Wide Web: Faralli et al. (2017)
  - Social Network Analysis and Mining: Faralli et al. (2015b)
  - Information Systems: Kang and Lee (2016)
  - Procedia Computer Science: Jiang and Sha (2015)
2. Conference proceedings
  - **WI** (IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology): Ahn et al. (2012), Gao et al. (2011), Peñas et al. (2013), Xu et al. (2011), Zarrinkalam and Kahani (2015)

- **UMAP** (Conference on User Modeling Adaptation and Personalization): Abel et al. (2011b), Hannon et al. (2012), Narducci et al. (2013)
- **CIKM** (ACM International Conference on Information and Knowledge Management): Piao and Breslin (2016d), Sang et al. (2015), Vu and Perez (2013)
- **ECIR** (European Conference on Information Retrieval): Piao and Breslin (2017a), Trikha et al. (2018), Zarrinkalam et al. (2016)
- **ISWC** (International Conference on Semantic Web): Abel et al. (2011c), Siehndel and Kawase (2012)
- **IUI** (International Conference on Intelligent User Interfaces): Bhargava et al. (2015), Garcia Esparza et al. (2013)
- **RecSys** (ACM Conference on Recommender Systems): Bhattacharya et al. (2014), Phelan et al. (2009)
- **SEMANTiCS** (International Conference on Semantic Systems): Orlandi et al. (2012), Piao and Breslin (2016b)
- **HT** (ACM Conference on Hypertext and Social Media): Piao and Breslin (2017b)
- **SIGIR** (International ACM Conference on Research and Development in Information Retrieval): Chen et al. (2010)
- **AAAI** (AAAI Conference on Artificial Intelligence): Lu et al. (2012)
- **KDD** (Knowledge Discovery and Data Mining): Spasojevic et al. (2014)
- **IJCAI** (International Joint Conference on Artificial Intelligence): Abel et al. (2013a)
- **ICWE** (International Conference on Web Engineering): Abel et al. (2012)
- **WebSci** (International Web Science Conference): Abel et al. (2011a)
- **ESWC** (Extended Conference on Semantic Web): Kapanipathi et al. (2014)
- **EKAW** (International Conference on Knowledge Engineering and Knowledge Management): Piao and Breslin (2016c)
- **ICSC** (IEEE International Conference on Semantic Computing): Große-Börling et al. (2015)
- **SAC** (ACM Symposium on Applied Computing): Besel et al. (2016a)
- **WSDM** (ACM International Conference on Web Search and Data Mining): Weng et al. (2010)
- **JCDL** (Joint Conference on Digital Libraries): Nishioka and Scherp (2016)
- **i-KNOW** (International Conference on Knowledge Technologies and Data-driven Business): Nishioka et al. (2015)
- **SPIM** (International Conference on Semantic Personalized Information Management: Retrieval and Recommendation): Kapanipathi et al. (2011)
- **OpenSym** (International Symposium on Open Collaboration): Lim and Datta (2013)
- **ADMA** (Advanced Data Mining and Applications): Jipmo et al. (2017)

### 3. Workshop proceedings

- **AND** (Workshop on Analytics for Noisy Unstructured Text Data): Michelson and Macskassy (2010)
- **Micropost** (Workshop on Making Sense of Microposts): Karatay and Karagoz (2015)

- **SMAP** (Workshop on Semantic and Social Media Adaptation and Personalization): Kanta et al. (2012)
- **RSWeb** (Workshop on Recommender Systems and the Social Web): O’Banion et al. (2012)
- **BlackMirror** (Workshop on Re-coding Black Mirror): Nechaev et al. (2017)

#### 4. Others

- Tech Report: Budak et al. (2014).

## References

- Abdel-Hafez, A., Xu, Y.: A survey of user modelling in social media websites. *Comput. Inf. Sci.* **6**(4), 59–71 (2013)
- Abel, F.: Contextualization, user modeling and personalization in the social web—from social tagging via context to cross-system user modeling and personalization. PhD thesis, Leibniz University of Hanover (2011)
- Abel, F., Gao, Q., Houben, G.J., Tao, K.: Analyzing temporal dynamics in Twitter profiles for personalized recommendations in the social web. In: *Proceedings of the 3rd International Web Science Conference*, Koblenz, Germany, pp. 1–8. ACM (2011a)
- Abel, F., Gao, Q., Houben, G.J., Tao, K.: Analyzing user modeling on Twitter for personalized news recommendations. In: *User Modeling, Adaption and Personalization*, Girona, Spain, pp. 1–12. Springer (2011b)
- Abel, F., Gao, Q., Houben, G.J., Tao, K.: Semantic enrichment of Twitter posts for user profile construction on the social web. In: *The Semantic Web: Research and Applications: 8th Extended Semantic Web Conference, ESWC 2011, Heraklion, Crete, Greece*, pp. 375–389. Springer (2011c)
- Abel, F., Hauff, C., Houben, G.J., Tao, K.: Leveraging user modeling on the social web with linked data. In: *Web Engineering: 12th International Conference, ICWE 2012, Berlin, Germany*, pp. 378–385. Springer (2012)
- Abel, F., Gao, Q., Houben, G.J., Tao, K.: Twitter-based user modeling for news recommendations. In: *Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence, IJCAI ’13, Beijing, China*, pp. 2962–2966. AAAI Press (2013a)
- Abel, F., Herder, E., Houben, G.J., Henze, N., Krause, D.: Cross-system user modeling and personalization on the social web. *User Model. User Adapt. Interact.* **23**(2–3), 169–209 (2013b)
- Ahmed, A., Low, Y., Aly, M., Josifovski, V., Smola, A.J.: Scalable distributed inference of dynamic user interests for behavioral targeting. In: *Proceedings of the 17th International Conference on Knowledge Discovery and Data Mining*, San Diego, CA, USA, pp. 114–122. ACM (2011)
- Ahn, D., Kim, T., Hyun, S.J., Lee, D.: Inferring user interest using familiarity and topic similarity with social neighbors in Facebook. In: *Proceedings of the 2012 IEEE/WIC/ACM International Joint Conferences on Web Intelligence and Intelligent Agent Technology, WI-IAT ’12, Washington, DC, USA*, vol. 01, pp. 196–200. IEEE Computer Society (2012)
- Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: DBpedia: a nucleus for a web of open data. In: *The Semantic Web: 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference*, Busan, Korea, pp. 722–735. Springer (2007)
- Bellogn, A., Said, A.: Rate CTRCT, gain DCGDC, error MAEMA, precision MAPMA, learning MLM, error RRMS. Recommender systems evaluation. (2017). <http://ir.ii.uam.es/~alejandro/2017/esnam.pdf>. Accessed 10 June 2018
- Besel, C., Schlötterer, J., Granitzer, M.: Inferring semantic interest profiles from Twitter followees: does Twitter know better than your friends? In: *Proceedings of the 31st Annual ACM Symposium on Applied Computing, SAC ’16, New York, NY, USA*, pp. 1152–1157. ACM (2016a)
- Besel, C., Schlötterer, J., Granitzer, M.: On the quality of semantic interest profiles for online social network consumers. *ACM SIGAPP Appl. Comput. Rev.* **16**(3), 5–14 (2016b)
- Bhargava, P., Brdiczka, O., Roberts, M.: Unsupervised modeling of users’ interests from their Facebook profiles and activities. In: *Proceedings of the 20th International Conference on Intelligent User Inter-*

- faces, IUI '15, New York, NY, USA, pp. 191–201. ACM (2015). <https://doi.org/10.1145/2678025.2701365>
- Bhattacharya, P., Zafar, M.B., Ganguly, N., Ghosh, S., Gummadi, K.P.: Inferring user interests in the Twitter social network. In: Proceedings of the 8th ACM Conference on Recommender Systems, RecSys'14, New York, NY, USA, pp. 357–360. ACM (2014)
- Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent Dirichlet allocation. *J. Mach. Learn. Res.* **3**(Jan), 993–1022 (2003)
- Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: a collaboratively created graph database for structuring human knowledge. In: Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data, pp. 1247–1250. ACM (2008)
- Bontcheva, K., Rout, D.: Making sense of social media streams through semantics: a survey. *Semant. Web* **5**(5), 373–403 (2014). <https://doi.org/10.3233/SW-130110>
- Brickley, D., Miller, L.: FOAF vocabulary specification 0.98. (2012). <http://xmlns.com/foaf/spec/>. Accessed 10 Dec 2017
- Brusilovsky, P., Karagiannidis, C., Sampson, D.: The benefits of layered evaluation of adaptive applications and services. In: Empirical Evaluation of Adaptive Systems. Proceedings of Workshop at the Eighth International Conference on User Modeling, UM2001, pp. 1–8 (2001)
- Brusilovsky, P., Kobsa, A., Nejdl, W.: The Adaptive Web: Methods and Strategies of Web Personalization, vol. 4321. Springer, Berlin (2007)
- Budak, C., Kannan, A., Agrawal, R., Pedersen, J.: Inferring user interests from microblogs. Technical report, Microsoft (2014)
- Carmagnola, F., Cena, F., Console, L., Cortassa, O., Gena, C., Goy, A., Torre, I., Toso, A., Venero, F.: Tag-based user modeling for social multi-device adaptive guides. *User Model. User Adapt. Interact.* **18**(5), 497–538 (2008). <https://doi.org/10.1007/s11257-008-9052-2>
- Carmagnola, F., Cena, F., Gena, C.: User model interoperability: a survey. *User Model. User Adapt. Interact.* **21**(3), 285–331 (2011). <https://doi.org/10.1007/s11257-011-9097-5>
- Chen, J., Nairn, R., Nelson, L., Bernstein, M., Chi, E.: Short and tweet: experiments on recommending content from information streams. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, Atlanta, GA, USA, pp. 1185–1194. ACM (2010)
- Chen, K., Chen, T., Zheng, G., Jin, O., Yao, E., Yu, Y.: Collaborative personalized tweet recommendation. In: SIGIR '12: Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval, Portland, Oregon, USA, pp. 661–670. ACM (2012)
- Cohen, P.R., Perrault, C.R.: Elements of a plan-based theory of speech acts. *Cogn. Sci.* **3**(3), 177–212 (1979)
- Collins, A.M., Loftus, E.F.: A spreading-activation theory of semantic processing. *Psychol. Rev.* **82**(6), 407 (1975)
- Edmonds, J.: Optimum branchings. In: Dantzig, G.B., Veinott, A.F. (eds.) *Mathematics and the Decision Sciences*, pp. 335–345. American Mathematical Society, Providence (1968)
- Faralli, S., Stilo, G., Velardi, P.: Large scale homophily analysis in Twitter using a Twixonomy. In: Proceedings of the 24th International Conference on Artificial Intelligence, Buenos Aires, Argentina, pp. 2334–2340. AAAI Press (2015a)
- Faralli, S., Stilo, G., Velardi, P.: Recommendation of microblog users based on hierarchical interest profiles. *Soc. Netw. Anal. Min.* **5**(1), 1–23 (2015b)
- Faralli, S., Stilo, G., Velardi, P.: Automatic acquisition of a taxonomy of microblogs users' interests. *Web Semant. Sci. Serv. Agents. World Wide Web* (2017). <https://doi.org/10.1016/j.websem.2017.05.004>
- Färber, M., Ell, B., Menne, C., Rettinger, A.: A comparative survey of DBpedia, Freebase, OpenCyc, Wikidata, and YAGO. *Semant. Web J.* **1**, 1–26 (2015)
- Flati, T., Vannella, D., Pasini, T., Navigli, R.: Two is bigger (and better) than one: the Wikipedia bitaxonomy project. In: 52nd Annual Meeting of the Association for Computational Linguistics, ACL, Baltimore, MD, USA, pp. 945–955. Association for Computational Linguistics (ACL) (2014)
- Gabrilovich, E., Markovitch, S.: Computing semantic relatedness using wikipedia-based explicit semantic analysis. In: Proceedings of the 20th International Joint Conference on Artificial Intelligence, Hyderabad, India, pp. 1606–1611. Morgan Kaufmann (2007)
- Gao, Q., Abel, F., Houben, G.J., Tao, K.: Interweaving trend and user modeling for personalized news recommendation. In: Proceedings of the 2011 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology, WI-IAT '11, Washington, DC, USA, vol. 01, pp. 100–103. IEEE Computer Society (2011)

- Gao, Q., Abel, F., Houben, G.J.: Genius: generic user modeling library for the social semantic web. In: *The Semantic Web*, pp. 160–175. Springer (2012)
- Garcia Esparza, S., O'Mahony, M.P., Smyth, B.: CatStream: categorising tweets for user profiling and stream filtering. In: *Proceedings of the 2013 International Conference on Intelligent User Interfaces, IUI '13*, New York, NY, USA, pp. 25–36. ACM (2013)
- Gauch, S., Speretta, M., Chandramouli, A., Micarelli, A.: User profiles for personalized information access. In: *The Adaptive Web*, pp. 54–89. Springer, Berlin (2007)
- Gena, C., Weibelzahl, S.: Usability engineering for the adaptive web. In: *The Adaptive Web*, pp. 720–762. Springer (2007)
- Gong, W., Lim, E.P., Zhu, F.: Characterizing silent users in social media communities. In: *ICWSM (2015)*
- Große-Böling, G., Nishioka, C., Scherp, A.: Generic process for extracting user profiles from social media using hierarchical knowledge bases. In: *2015 IEEE International Conference on Semantic Computing (ICSC) (2015)*. <https://doi.org/10.1109/ICOSC.2015.7050806>
- Guha, R., Gupta, V., Raghunathan, V., Srikant, R.: User modeling for a personal assistant. In: *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining—WSDM '15*, New York, NY, USA, pp. 275–284. ACM Press (2015)
- Haewoon, K., Changhyun, L., Hosung, P., Sue, M.: What is Twitter, a social network or a news media? In: *Proceedings of the 19th International Conference on World Wide Web*, Raleigh, NC, USA. ACM (2010)
- Han, J., Pei, J.: Mining frequent patterns by pattern-growth: methodology and implications. *ACM SIGKDD Explor. Newslett.* **2**(2), 14–20 (2000)
- Han, L., Kashyap, A.L., Finin, T., Mayfield, J., Weese, J.: UMBC\_EBIQUITY-CORE: semantic textual similarity systems. In: *The Second Joint Conference on Lexical and Computational Semantics*, Atlanta, GA, USA, pp. 44–52. Association for Computational Linguistics (2013)
- Hannon, J., McCarthy, K., O'Mahony, M.P., Smyth, B.: A multi-faceted user model for Twitter. In: *User Modeling, Adaptation, and Personalization: 20th International Conference, UMAP 2012*. Montreal, Canada, pp. 303–309. Springer (2012)
- Heath, T., Bizer, C.: Linked data: evolving the web into a global data space. *Synthesis lectures on the semantic web: theory and technology*, vol. 1, no. 1, pp. 1–136 (2011)
- Herlocker, J.L., Konstan, J.A., Terveen, L.G., Riedl, J.T.: Evaluating collaborative filtering recommender systems. *ACM Trans. Inf. Syst.* **22**(1), 5–53 (2004). <https://doi.org/10.1145/963770.963772>
- Holden, S., Kay, J.: *The Scrutable User Model and Beyond*. Basser Department of Computer Science, University of Sydney, Sydney (1999)
- Hong, L., Doumith, A.S., Davison, B.D.: Co-factorization machines: modeling user interests and predicting individual decisions in Twitter. In: *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining, WSDM '13*, New York, NY, USA, pp. 557–566. ACM (2013)
- Hung, C.C., Huang, Y.C., Hsu, J.Y.j., Wu, D.K.C.: Tag-based user profiling for social media recommendation. In: *Proceedings of the 21st ACM Conference on Hypertext and Hypermedia*, pp. 151–156 (2008)
- Ingwersen, P.: Polyrepresentation of information needs and semantic entities elements of a cognitive theory for information retrieval interaction. In: *Proceedings of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, Dublin, Ireland, pp. 101–110. Springer (1994)
- Java, A., Song, X., Finin, T., Tseng, B.: Why we Twitter: understanding microblogging usage and communities. In: *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 Workshop on Web Mining and Social Network Analysis*, San Jose, CA, USA, pp. 56–65. ACM (2007)
- Jiang, B., Sha, Y.: Modeling temporal dynamics of user interests in online social networks. *Proc. Comput. Sci.* **51**, 503–512 (2015)
- Jipmo, C.N., Quercini, G., Bennacer, N.: FRISK: a multilingual approach to find twitterR InterestS via wiKipedia BT. In: *Advanced Data Mining and Applications: 13th International Conference, ADMA 2017*, Singapore, Nov 2017, Proceedings, pp. 243–256. Springer, Cham (2017). [https://doi.org/10.1007/978-3-319-69179-4\\_17](https://doi.org/10.1007/978-3-319-69179-4_17)
- Kaminskas, M., Bridge, D.: Diversity, serendipity, novelty, and coverage: a survey and empirical analysis of beyond-accuracy objectives in recommender systems. *ACM Trans. Interact. Intell. Syst.* **7**(1), 2:1–2:42 (2016). <https://doi.org/10.1145/2926720>
- Kang, J., Lee, H.: Modeling user interest in social media using news media and Wikipedia. *Inf. Syst.* **65**, 52–64 (2016)

- Kanta, M., Simko, M., Bieliková, M.: Trend-aware user modeling with location-aware trends on Twitter. In Proceedings of 7th International Workshop on Semantic and Social Media Adaptation and Personalization, SMAP (2012)
- Kapanipathi, P., Orlandi, F., Sheth, A., Passant, A.: Personalized filtering of the Twitter stream. In: Proceedings of the Second International Conference on Semantic Personalized Information Management: Retrieval and Recommendation, vol. 781, pp. 6–13. CEUR-WS.org, Bonn, Germany (2011)
- Kapanipathi, P., Jain, P., Venkataramani, C., Sheth, A.: User interests identification on Twitter using a hierarchical knowledge base. In: The Semantic Web: Trends and Challenges, Anissaras, Crete, Greece, pp. 99–113. Springer (2014)
- Karatay, D., Karagoz, P.: User interest modeling in Twitter with named entity recognition. In: Making Sense of Microposts (# Microposts 2015), Florence, Italy, pp. 17–20 (2015)
- Kay, J.: Scrutable adaptation: because we can and must. In: International Conference on Adaptive Hypermedia and Adaptive Web-Based Systems, pp. 11–19. Springer (2006)
- Kim, D., Jo, Y., Moon, I.C., Oh, A.: Analysis of Twitter lists as a potential source for discovering latent characteristics of users. In: ACM CHI Workshop on Microblogging, Atlanta, GA, USA, p. 4. Citeseer (2010)
- Kitchenhams, B.: Procedures for Performing Systematic Reviews, vol. 33, pp. 1–26. Keele University, Keele (2004)
- Liao, Y., Moshtaghi, M., Han, B., Karunasekera, S., Kotagiri, R., Baldwin, T., Harwood, A., Pattison, P.: Mining micro-blogs: opportunities and challenges. In: Abraham, A. (ed.) Computational Social Networks, pp. 129–159. Springer, Berlin (2012)
- Lim, K.H., Datta, A.: Interest classification of twitter users using Wikipedia. In: Proceedings of the 9th International Symposium on Open Collaboration, WikiSym '13, Hong Kong, China, pp. 22:1–22:2. ACM (2013)
- Liu, J., Zhang, F., Song, X., Song, Y.I., Lin, C.Y., Hon, H.W.: What's in a name? An unsupervised approach to link users across communities. In: Proceedings of the Sixth ACM International Conference on Web Search and Data Mining, Rome, Italy, pp. 495–504. ACM (2013)
- Lu, C., Lam, W., Zhang, Y.: Twitter user modeling and tweets recommendation based on Wikipedia concept graph. In: Workshops at the Twenty-Sixth AAAI Conference on Artificial Intelligence, Toronto, ON, Canada (2012)
- Manrique, R., Mariño, O.: How does the size of a document affect linked open data user modeling strategies? In: Proceedings of the International Conference on Web Intelligence, WI '17, New York, NY, USA, pp. 1246–1252. ACM (2017)
- Mezghani, M., Zayani, C.A., Amous, I., Gargouri, F.: A user profile modelling using social annotations: a survey. In: Proceedings of the 21st International Conference on World Wide Web, WWW '12 Companion, New York, NY, USA, pp. 969–976. ACM (2012)
- Michelson, M., Macskassy, S.A.: Discovering users' topics of interest on Twitter: a first look. In: Proceedings of the 4th Workshop on Analytics for Noisy Unstructured Text Data, Toronto, ON, Canada, pp. 73–80. ACM (2010)
- Mihalcea, R., Tarau, P.: TextRank: bringing order into texts. In: Lin, D., Wu, D. (eds.) Proceedings of EMNLP 2004, Barcelona, Spain, pp. 404–411. Association for Computational Linguistics (2004)
- Miller, G.A.: WordNet: a lexical database for English. *Commun. ACM* **38**(11), 39–41 (1995)
- Myers, S.A., Leskovec, J.: The bursty dynamics of the Twitter information network. In: Proceedings of the 23rd International Conference on World Wide Web, Seoul, Korea, pp. 913–924. ACM (2014)
- Narducci, F., Musto, C., Semeraro, G., Lops, P., Gemmis, M.: Leveraging encyclopedic knowledge for transparent and serendipitous user profiles. In: User Modeling, Adaptation, and Personalization: 21st International Conference, pp. 350–352. Springer, Berlin (2013)
- Navigli, R., Ponzetto, S.P.: BabelNet: the automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artif. Intell.* **193**, 217–250 (2012)
- Nechaev, Y., Corcoglioniti, F., Giuliano, C.: Concealing interests of passive users in social media. In: The Re-coding Black Mirror 2017 Workshop Co-located with 16th International Semantic Web Conference (ISWC 2017), Vienna, Austria (2017)
- Nguyen, P.T., Tomeo, P., Di Noia, T., Di Sciascio, E.: Content-based recommendations via DBpedia and Freebase: a case study in the music domain. In: International Semantic Web Conference, pp. 605–621 (2015)

- Nishioka, C., Scherp, A.: Profiling vs. time vs. content: what does matter for top-k publication recommendation based on Twitter profiles? In: Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries, JCDL '16, New York, NY, USA, pp. 171–180. ACM (2016)
- Nishioka, C., Große-Bölting, G., Scherp, A.: Influence of time on user profiling and recommending researchers in social media. In: Proceedings of the 15th International Conference on Knowledge Technologies and Data-Driven Business, i-KNOW '15, New York, NY, USA, pp. 9:1–9:8. ACM (2015)
- O'Banion, S., Birnbaum, L., Hammond, K.: Social media-driven news personalization. In: Proceedings of the 4th ACM RecSys Workshop on Recommender Systems and the Social Web, Dublin, Ireland, pp. 45–52. ACM (2012)
- Orlandi, F., Breslin, J., Passant, A.: Aggregated, interoperable and multi-domain user profiles for the social web. In: Proceedings of the 8th International Conference on Semantic Systems, Graz, Austria, pp. 41–48. ACM (2012)
- Orlandi, F., Kapanipathi, P., Sheth, A., Passant, A.: Characterising concepts of interest leveraging linked data and the social web. In: Proceedings of the 2013 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT), WI-IAT '13, Washington, DC, USA, vol. 01, pp. 519–526. IEEE Computer Society (2013)
- Paramythis, A., Weibelzahl, S., Masthoff, J.: Layered evaluation of interactive adaptive systems: framework and formative methods. *User Model. User Adapt. Interact.* **20**(5), 383–453 (2010)
- Peñas, P., del Hoyo, R., Vea-Murguía, J., González, C., Mayo, S.: Collective knowledge ontology user profiling for Twitter—automatic user profiling. In: Proceedings of the 2013 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT) (2013)
- Perera, S., Mendes, P.N., Alex, A., Sheth, A.P., Thirunarayan, K.: Implicit entity linking in tweets BT—the semantic web. In: Sack, H., Blomqvist, E., D'Aquin, M., Ghidini, C., Ponzetto, S.P., Lange, C. (eds.) *Latest Advances and New Domains: 13th International Conference, ESWC 2016*, pp. 118–132. Springer, Cham (2016)
- Perrault, C.R., Allen, J.F., Cohen, P.R.: Speech acts as a basis for understanding dialogue coherence. In: Proceedings of the 1978 Workshop on Theoretical Issues in Natural Language Processing, pp. 125–132. Association for Computational Linguistics (1978)
- Phelan, O., McCarthy, K., Smyth, B.: Using Twitter to recommend real-time topical news. In: Proceedings of the Third ACM Conference on Recommender Systems, RecSys '09, New York, NY, USA, pp. 385–388. ACM (2009)
- Piao, G., Breslin, J.J.G.: Analyzing aggregated semantics-enabled user modeling on Google+ and Twitter for personalized link recommendations. In: UMAP 2016—Proceedings of the 2016 Conference on User Modeling Adaptation and Personalization, Halifax, NS, Canada, pp. 105–109. ACM (2016a) <https://doi.org/10.1145/2930238.2930278>
- Piao, G., Breslin, J.J.G.: Exploring dynamics and semantics of user interests for user modeling on Twitter for link recommendations. In: Proceedings of the 12th International Conference on Semantic Systems, Leipzig, Germany, 13–14 Sept 2016, pp. 81–88. ACM (2016b). <https://doi.org/10.1145/2993318.2993332>
- Piao, G., Breslin, J.J.G.: Interest representation, enrichment, dynamics, and propagation: a study of the synergetic effect of different user modeling dimensions for personalized recommendations on Twitter. In: LNAI, Bologna, Italy, vol. 10024. Springer (2016c). [https://doi.org/10.1007/978-3-319-49004-5\\_32](https://doi.org/10.1007/978-3-319-49004-5_32)
- Piao, G., Breslin, J.J.G.: User modeling on Twitter with WordNet Synsets and DBpedia concepts for personalized recommendations. In: International Conference on Information and Knowledge Management, Proceedings, Indianapolis, IN, USA, 24–28 Oct 2016, pp. 2057–2060. ACM (2016d). <https://doi.org/10.1145/2983323.2983908>
- Piao, G., Breslin, J.J.G.: Inferring user interests for passive users on Twitter by leveraging followee biographies. In: LNCS, Aberdeen, UK, vol. 10193. Springer (2017a). [https://doi.org/10.1007/978-3-319-56608-5\\_10](https://doi.org/10.1007/978-3-319-56608-5_10)
- Piao, G., Breslin, J.J.G.: Leveraging followee list memberships for inferring user interests for passive users on Twitter. In: HT 2017—Proceedings of the 28th ACM Conference on Hypertext and Social Media, Prague, Czech Republic. ACM Press (2017b). <https://doi.org/10.1145/3078714.3078730>
- Rich, E.: User modeling via stereotypes. *Cogn. Sci.* **3**(4), 329–354 (1979)

- Ritter, A., Clark, S., Etzioni, O.: Named entity recognition in tweets: an experimental study. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, Edinburgh, UK, pp. 1524–1534. Association for Computational Linguistics (2011)
- Rosen-Zvi, M., Griffiths, T., Steyvers, M., Smyth, P.: The author-topic model for authors and documents. In: Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence, USA, UAI '04, Arlington, VA, pp. 487–494. AUA Press (2004)
- Rotta, R., Noack, A.: Multilevel local search algorithms for modularity clustering. *J. Exp. Algorithmics* **16**, 2.3:2.1–2.3:2.27 (2011). <https://doi.org/10.1145/1963190.1970376>
- Salton, G., McGill, M.J.: Introduction to Modern information Retrieval. McGraw-Hill, New York (1986)
- Sang, J., Lu, D., Xu, C.: A probabilistic framework for temporal user modeling on microblogs. In: Proceedings of the 24th ACM International Conference on Information and Knowledge Management, CIKM '15, New York, NY, USA, pp. 961–970. ACM (2015). <https://doi.org/10.1145/2806416.2806470>
- Shen, W., Wang, J., Luo, P., Wang, M.: Linking named entities in tweets with knowledge base via user interest modeling. In: Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '13, New York, NY, USA, pp. 68–76. ACM (2013). <https://doi.org/10.1145/2487575.2487686>
- Sheth, A., Kapanipathi, P.: Semantic filtering for social data. *IEEE Internet Comput.* **20**(4), 74–78 (2016)
- Siehdnel, P., Kawase, R.: TwikiMe!: user profiles that make sense. In: Proceedings of the 2012th International Conference on Semantic Web (Posters and Demonstrations Track), ISWC-PD'12, vol. 914, pp. 61–64. CEUR-WS.org (2012)
- Spasojevic, N., Yan, J., Rao, A., Bhattacharyya, P.: LASTA: large scale topic assignment on multiple social networks. In: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA, pp. 1809–1818. ACM (2014). <https://doi.org/10.1145/2623330.2623350>
- Stefani, A.: Personalizing access to web sites: the SiteIF project. In: Proceedings of the 2nd Workshop on Adaptive Hypertext and Hypermedia HYPERTEXT (1998)
- Suchanek, F.M., Kasneci, G., Weikum, G.: Yago: a core of semantic knowledge. In: Proceedings of the 16th International Conference on World Wide Web, pp. 697–706. ACM (2007)
- Szomszor, M., Alani, H., Cantador, I., O'Hara, K., Shadbolt, N.: Semantic modelling of user interests based on cross-folksonomy analysis. In: The Semantic Web—ISWC 2008, Lecture Notes in Computer Science, SE-40, vol. 5318, pp. 632–648. Springer, Berlin (2008)
- Tao, K., Abel, F., Gao, Q., Houben, G.J.: TUMS: Twitter-based user modeling service. In: García-Castro, R., Fensel, D., Antoniou, G. (eds.), The Semantic Web: ESWC 2011 Workshops, vol. 7117, chap. 22, pp. 269–283. Springer, Berlin (2012)
- Tommaso, G.D., Faralli, S., Stilo, G., Velardi, P.: Wiki-MID: a very large multi-domain interests dataset of Twitter users with mappings to Wikipedia. In: The 17th International Semantic Web Conference. Springer (2018)
- Trikha, A.K., Zarrinkalam, F., Bagheri, E.: Topic-association mining for user interest detection. In: The 40th European Conference on Information Retrieval. Springer (2018)
- Vrandečić, D., Krötzsch, M.: Wikidata: a free collaborative knowledgebase. *Commun. ACM* **57**(10), 78–85 (2014)
- Vu, T., Perez, V.: Interest mining from user tweets. In: Proceedings of the 22nd ACM International Conference on Information and Knowledge Management, CIKM '13, New York, NY, USA, pp. 1869–1872. ACM (2013)
- Weng, J., Lim, E.P., Jiang, J., He, Q.: TwitterRank: finding topic-sensitive influential Twitterers. In: Proceedings of the Third ACM International Conference on Web Search and Data Mining, WSDM '10, New York, NY, USA, pp. 261–270. ACM (2010)
- White, R.W., Bailey, P., Chen, L.: Predicting user interests from contextual information. In: Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '09, New York, NY, USA, pp. 363–370. ACM (2009)
- Xu, Z., Ru, L., Xiang, L., Yang, Q.: Discovering user interest on Twitter with a modified author-topic model. In: Proceedings of the 2011 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology, vol. 01, Washington, DC, USA, pp. 422–429. IEEE Computer Society (2011)
- Zarrinkalam, F.: Semantics-enabled user interest mining. In: Gandon, F., Sabou, M., Sack, H., D'Amato, C., Cudré-Mauroux, P., Zimmermann, A. (eds.) The Semantic Web. Latest Advances and New Domains, SE-54, Lecture Notes in Computer Science, vol. 9088, pp. 817–828. Springer (2015)

- Zarrinkalam, F., Kahani, M.: Semantics-enabled user interest detection from Twitter. In: 2015 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, Singapore, pp. 469–476 (2015)
- Zarrinkalam, F., Fani, H., Bagheri, E., Kahani, M.: Inferring implicit topical interests on Twitter. In: European Conference on Information Retrieval, pp. 479–491, Padua, Italy. Springer (2016)
- Zarrinkalam, F., Fani, H., Bagheri, E., Kahani, M.: Predicting users' future interests on Twitter. In: European Conference on Information Retrieval, pp. 464–476. Springer (2017)
- Zhou, X., Xu, Y., Li, Y., Josang, A., Cox, C.: The state-of-the-art in personalized recommender systems for social networking. *Artif. Intell. Rev.* **37**(2), 119–132 (2012)

**Guangyuan Piao** is a Ph.D. student at the Insight Centre for Data Analytics (formerly DERI) at the National University of Ireland Galway. He received his B.Sc. in Computer Science from Jilin University, China, and received his M.Eng. degree in Information and Industrial Engineering from Yonsei University, South Korea. His main research interests include User Modeling, Recommender Systems, and Knowledge Graph. His current research focuses on semantics-aware user modeling and recommender systems leveraging knowledge graphs and latent semantics.

**John G. Breslin** is a Senior Lecturer in Electrical and Electronic Engineering at the College of Science and Engineering at the National University of Ireland Galway, where he is Director of the TechInnovate/AgInnovate programmes. John has taught electronic engineering, computer science, innovation and entrepreneurship topics during the past two decades. He is also a Co-Principal Investigator at the Insight Centre for Data Analytics, and a Funded Investigator at Confirm Smart Manufacturing and VistaMilk. He has written 190 peer-reviewed academic publications (h-index of 37, 5500 citations, best paper awards from DL4KGS, SEMANTICS, ICEGOV, ESWC, PELS), and co-authored the books “The Social Semantic Web” and “Social Semantic Web Mining”. He co-created the SIOC framework, implemented in hundreds of applications (by Yahoo, Boeing, Vodafone, etc.) on at least 65,000 websites with 35 million data instances.