

Recommending #-Tags in Twitter

Eva Zangerle, Wolfgang Gassler, Günther Specht

Databases and Information Systems
Institute of Computer Science
University of Innsbruck, Austria
`firstname.lastname@uibk.ac.at`

Abstract. Twitter, currently the most popular microblogging tool available, is used to publish more than 140,000,000 messages a day. Many users use hashtags to categorize their tweets. However, hashtags are not restricted in any way in terms of usage or syntax which leads to a very heterogeneous set of hashtags occurring in the Twitter universe and therefore, decreases the search capabilities. In this paper, we present an approach for the recommendation of highly appropriate hashtags to the user during the creation process. The recommendations aim at encouraging the user to (i) use hastags at all, (ii) use more appropriate hashtags and (iii) avoid the usage of synonymous hashtags. Therefore the vocabulary of hashtags becomes more homogenous regarding both syntax and semantics.

1 Introduction

Social networks have gained significant importance on the web throughout the last years. The most popular microblogging tool, Twitter, has experienced tremendous success lately and has become very important as both a social network and a news media [13]. Twitter enables all registered users to post 140-character messages and follow other users. The users's personal timeline (home-view on the Twitter universe) basically includes all messages – the so-called tweets – of all followed users. The notion of a follower describes a user who follows another user. Vice versa, the notion of a followee describes a user who is followed by another user. Such a connection between users is not reciprocal - user A can follow any other user B without requiring user B to follow user A back. Additionally, all messages are fully accessible to the public. Nowadays, 140,000,000 messages - so-called *tweets* - are posted every day. As reported by Twitter¹, every day more than 400,000 new users join the Twitter network.

The basic motivation of users to join Twitter and participate is manifold [11]. Millions of users use Twitter to keep track of friends and keep friends updated. Users may seek for advice on certain problems or participate in general discussions about certain topics. Some participants follow celebrities or companies in order to stay updated. Many of the active users - those who are not just following other users, but are also actively posting tweets - use Twitter as a medium to

¹ <http://blog.twitter.com/2011/03/numbers.html>

let the world know what they're up to or simply to share some information they consider useful. The probably most important feature of Twitter is the retweet functionality. It enables users to further broadcast tweets they consider worth spreading within the Twitter network. Mostly, the retweeted message remains unchanged. A retweeted message contains "RT: @originaluser" followed by the original message. This retweeting, which was also heavily analysed in [13], can spread an important message all over the world within minutes. Due to the ever increasing amount of Twitter messages and the resulting chaos within the Twittersphere, the microblogging community started to use so-called hashtags as a means for the manual categorization of tweets. The categorization can be used for either searching for certain topics based on the used hashtags or to be able to follow certain conversations about a certain topic on Twitter. The only requirement for hashtags is to start with a hash symbol #. Besides this fact, hashtags do not have to conform to any rules or regulations and can be seen as typical tags as used in common Web 2.0 applications like e.g. Blogs. Hashtags may appear at any arbitrary position within the message and may consist of any arbitrary combination of characters. This makes them easy to use, but at the same time leads to a significant lack of structure and uniformity. During our research, we crawled a data set and analyzed it. In the process we found that that users utilized very different popular hashtags for their tweets about the same topic. For example, the Tour de France (a world-famous bicycle race in France) was very popular. Tweets about this topic contain different hashtags, such as #tdf, #tourdefrance, #cycling or #procycling. Twitter offers its users a search engine which is able to search for keywords, but also for hashtags. Therefore, when searching for discussions about the Tour de France by using the search hashtag #tourdefrance, the user might not be able to retrieve all tweets containing information about the Tour de France. This is due to the fact that other users used the hashtag #tdf, which the user did not specify in the search query. Certainly, tweets containing the hashtag #tdf would also have been a perfect match for the user's query. However, due to the heterogeneous hashtag vocabulary used by the active Twitter community, many synonymous hashtags are used for describing the same semantic information.

In this paper we introduce an approach for the recommendation of hashtags. Our approach computes recommendations based on an analysis of existing tweets by other users and recommends suitable hashtags for the currently entered message to the user. This recommendation mechanism aims at encouraging the user to make use of hashtags and creating a more homogeneous hashtag vocabulary in order to enhance the quality of search result. Additionally, we present general statistics about the use of hashtags within Twitter and an evaluation of our approach.

The remainder of this paper is organized as follows. Section 2 describes the basic concepts of Twitter and hashtags. Section 3 is concerned with the process of hashtag recommendations. Section 4 contains the experiments and evaluations of the presented approach. Subsequently, Section 5 describes important related work and Section 6 concludes the paper.

2 Hashtags

Hashtagging is a simple and convenient way for users to categorize their own tweets. Such a hashtag within a tweet can simply be specified by adding a hash - '#' - followed by the tag itself. One tweet may also contain multiple hashtags, like in the following example tweet: "Don't forget! Only 7 days till the #SASWeb submission deadline #umap2011 <http://bit.ly/dKgS82>. #recsys #um #adaptivity #web3.0 #ontologies" which was posted by the SASWeb workshop (@sasWeb2011).

The most popular hashtags are either related to long-term popular topics or to current events or topics, e.g. the hashtag #tdf was extensively used during the crawling period as the Tour de France was taking place during this time. Typical long-term topics are e.g. #Apple or #Obama which are featured in thousands of messages a day [13].

2.1 Data Set and Hashtag Analysis

In order to be able to analyse the hashtagging behaviour of Twitter users and to build up a database which forms the basis for all recommendation computations, we had to crawl tweets. Overall, we collected about 16,000,000 tweets from July 2010 until February 2011 via the Twitter Application Programming Interface.

In order to retrieve a diverse and highly representative data set to base our evaluations and analysis on, we decided to use Twitter's API². The basis for our search queries was an English dictionary containing more than 32,000 words. We iterated over the words contained in the dictionary and used them as search keywords for the Twitter Search API. All search results were stored whereby only tweets containing hashtags were used for further analysis. Another approach was to retrieve the public timeline, which basically consists of the ten latest tweets. The timeline is displayed on the Twitter website and is also available via the API. However, these tweets are only updated once a minute. Therefore, only 600 tweets could be retrieved per hour and considering the fact that only 20% of all tweets contain hashtags, this approach was not feasible for crawling a sufficiently large dataset.

After having crawled the data, we had to perform multiple preprocessing steps. This included removing all non-english messages (based on Twitter's language classification mentioned in the metadata of every tweet) and all messages not containing hashtags at all. Furthermore, all messages were transformed to lower-case. Table 1 contains an overview about the crawled data set and its characteristics. Out of the crawled tweets, more than 3 million tweets contained at least one hashtag, which marks 20% of all crawled tweets. The hashtags filtered from all tweets were further analysed in regards to their usage and popularity. Figure 1 displays the long tail distribution of hashtags and their usage. The fact that stands out about this distribution is that 86% of all hashtags within the data set were used within less than five tweets. On the other

² <http://search.twitter.com/search>

hand, the most popular hashtags within the data set (`#jobs`, `#nowplaying`, `#zodiacfacts`, `#news` and `#fb`) were used in 8% of all messages containing hashtags. Another interesting fact is the distribution of the number of hashtags used per tweet which can be seen in Figure 2. We expected the number of hashtags per message to be decreasing steadily. This is mostly the case for messages contains less than 15 hashtags. However, the sudden amplitude at 17 hashtags per message is somewhat surprising. We therefore examined these messages and discovered that these were spam tweets which only contained hashtags and a URL, like e.g. "RT @Bhupesh_tweet: #Quad #loop-http://bit.ly/ciHX2U #retweet #India #Jobs #World #news #canada #ad #win #USA #tdf #oea #hacking #icantstop #sdcc #game". Such tweets typically also feature a high retweet-rate by using a spam network consisting of many Twitter users created for spam purposes.

Characteristic	Value	Percentage
Crawled messages total	16,034,195	100%
Messages containg at least one hashtag	3,209,281	20%
Messages containing no hashtags	12,824,914	80%
Retweets	2,556,617	16%
Direct messages	3,073,948	19%
Hashtags usages total	5,097,545	–
Hashtags distinct	510,170	–
Average number of hashtags per message	1.5884	–
Maximum number of hashtags per message	23	–
Hashtags occurring < 5 times in total	437,266	–
Hashtags occurring < 3 times in total	328,348	–
Hashtags occuring only once	384,187	–

Table 1. Overview about the Crawled Tweets

3 Hashtag Recommendations

The aim of the approach presented in this paper is to find a set of hashtags suitable for any tweet the user enters. These hashtags are then recommended to the user during the creation process of the new tweet. Recommendations are basically be computed by performing the following steps:

1. finding the most similar messages in the crawled data set for the tweet just entered by the user

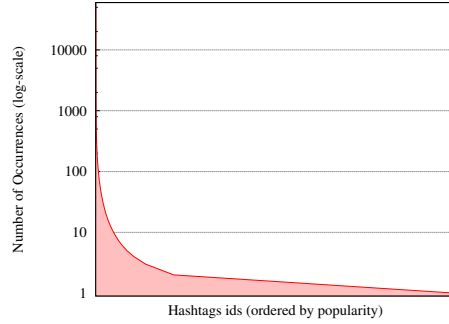


Fig. 1. Long tail of hashtag popularity



Fig. 2. Number of hashtags per message

2. retrieving the set of hashtags used within these most similar messages
3. ranking the computed set of hashtag recommendation candidates

These steps for the computation of hashtag recommendations are discussed in the following sections.

3.1 Similarity of Tweets

In order to be able to determine similar tweets, a similarity measure for the comparison of two at most 140 character long messages has to be introduced. This metric is used to rank the results gathered from searching similar tweets within the crawled data set. These similar messages are subsequently considered to contain valuable hashtag recommendation candidates. A straightforward solution is to use the term frequency - inverse document frequency measure for the comparison of tweets. In order to be able to use tf/idf for the computation of the similarity of tweets, the formula stated in Equation (1) is used.

$$tf_idf_{t,d} = tf_{t,d} * idf_t \quad (1)$$

$$tf_{t,d} = n_{t,d} \quad (2)$$

$$idf_t = \log \frac{|D|}{|\{d : t \in d\}|} \quad (3)$$

In the case of searching a set of tweets, the set D of documents which have to be searched is the set of tweets in the system. The term frequency basically is the number of occurrences of a term t within a given document d (tweet). The inverse document frequency (idf) constitutes the importance of a term t within the whole set of documents which are searched. This is computed by taking the number of all documents ($|D|$) within the index and dividing it by the number of documents which contain the searched term ($|\{d : t \in d\}|$). The computation of the tf/idf measure for a given search query (in our case the tweet inserted by the user), is subsequently accomplished by computing the sum of all tf/idf of all terms t occurring within the search query d : $\sum_{t \in d} tf_idf(t)$. Furthermore, the final score is increased if more of the terms of the query are matched. The final set of similar tweets (those obtaining the highest tf/idf-based score ratings) is restricted to a set of tweets having a score above a certain threshold corresponding to the total number of results and the specified limit of total results.

3.2 Ranking

After having obtained the set of the most similar messages to the tweet the user just entered, the hashtags are extracted from these tweets. These hashtags are referred to as hashtag recommendation candidates throughout the remainder of this paper. The ranking of these hashtag recommendation candidates is crucial for the success of recommendations. This is due to the fact that both the cognition of the user and the space available for displaying the recommendations is limited. In most cases a set of 5-10 recommendations is most appropriate which also correspond to the capacity of short-term memory (Miller, 1956). Therefore the top- k recommendations are shown to the user, where k denotes the size of the set of recommended hashtags presented to the user. This restricted set is based on the set of all hashtags which were extracted from the most similar messages to the newly created tweet. To present the most suitable top- k hashtags to the user, the recommendation candidates have to be ranked. For our approach, we evaluated three ranking methods, which can be summarized as follows:

- *OverallPopularityRank*: This ranking approach is based on the popularity of the hashtag recommendation candidates. It basically considers the number of occurrences of the respective hashtag within our data set. The more popular a hashtag is overall, the higher the resulting rank of the hashtag.
- *RecommendationPopularityRank*: This ranking method basically counts the occurrences of each hashtag within the set of recommendation candidates. The higher the number of occurrences, the more (similar) messages contain this hashtag. Therefore, it is likely that the hashtag is suitable for the tweet the user just entered.

- *SimilarityRank*: This ranking method is based on the similarity value between the tweet entered by the user and the tweet which provides a hashtag recommendation candidate. The more similar the messages are, the more likely it is that the hashtags contained in this similar message are suitable for the tweet entered by the user. In the case that multiple tweets contain the hashtag which has to be ranked, the similarity of the most similar tweet is used. As a metric for the similarity of tweets, we used tf/idf as described in 3.1.

4 Evaluation

A recommendation engine prototype implementing this approach has been developed based on Apache’s Lucene³ fulltext index. We used the fulltext index to store the crawled tweets which enabled us to find the most similar messages by using Lucene’s Search Index.

4.1 Test Setup

The evaluation was done on a CentOS release 5.1 machine with 8 GB of RAM. The evaluation of the hashtag recommendation approaches was conducted by performing a leave-one-out test. This test was based on the data set described in Section 2.1. Based on the crawled data set, we built a fulltext index comprising all 3.2 mil. cleaned messages without hashtags of this data set. From this index, we randomly chose 10,000 messages with less than six hashtags for each test run. For each of these messages, the contained hashtags were removed from the message and the resulting string was used as the input tweet for the recommendation engine. Naturally, the currently used tweet was removed from the Lucene Index and was not considered for the computation of recommendation candidates. Additionally, no retweets were used as test input tweets as search for similar messages would return an identical retweeted message which would obviously distort the evaluation results.

Based on the hashtag recommendations computed by the recommendation engine, we evaluated the three ranking methods described in 3.2. This was done by computing the precision and recall values of the top- k recommendations with $k = 1, k = 2, \dots, k = 10$ as described in the next section.

4.2 Precision and Recall

For the evaluation of the quality of the computed recommendations, we chose to use the precision and recall values of the recommendations. These metrics are defined as follows:

$$precision(\mathcal{H}_{rec}) = \frac{|\mathcal{H}_{rec} \cap \mathcal{H}_{orig}|}{|\mathcal{H}_{rec}|} \quad (4)$$

³ <http://lucene.apache.org/>

$$recall(\mathcal{H}_{rec}) = \frac{|\mathcal{H}_{rec} \cap \mathcal{H}_{orig}|}{|\mathcal{H}_{orig}|} \quad (5)$$

where $\mathcal{H}_{original}$ is the set of original hashtags which were removed from the original tweet and $\mathcal{H}_{recommended}$ is the set of top- k recommendations. We performed ten test runs for each ranking method with $k = 1, k = 2, \dots, k = 10$. Each test run computed the respective average recall and precision value of 10,000 test tweets. Thus, the evaluation is based on the computation of 100,000 top- k recommendation sets for each ranking method.

4.3 Results

The experiments conducted showed that the approach is feasible of recommending suitable hashtags. The recall values for the top- k recommended hashtags can be seen in Figure 3. In this figure, the recall values for k (the number of recommended hashtags) being between 1 and 10 has been evaluated for the three considered ranking methods. This Figure shows that ranking based on the overall popularity of the hashtag (OverallPopularityRank) and also based on the popularity of the hashtag within the hashtag recommendation candidates (RecommendationPopularityRank) do not perform well. In contrast, SimilarityRank (ranking based on the similarity of the original tweet and the tweet containing the recommendation candidate) is able to perform significantly better. This ranking method leads to promising recall values which are well above the 40% mark for $k > 2$.

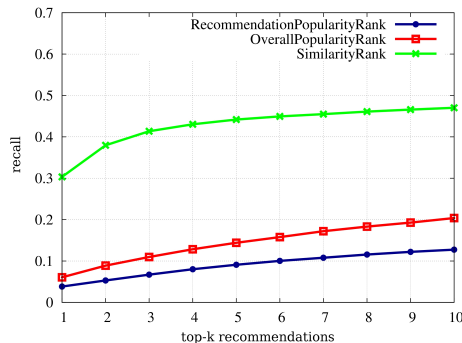


Fig. 3. Recall depending on number of Recommendations

The precision values for the computed recommendation sets decrease with an increasing k . This is due the fact, that we only use test tweets with at most 5 hashtags per message. Therefore even a set of 10 recommendations featuring a recall value of 100% only results in a precision of 50% as five of the ten

recommended hashtags are not applicable as the original message only features five hashtags.

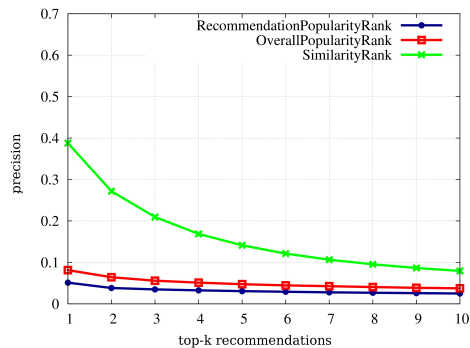


Fig. 4. Precision depending on number of Recommendations

Overall, the evaluations showed that our approach is suitable for the recommendation of hashtags. Another fact which can be derived from the evaluations is that our approach shows the best performance when restricting the set of recommended hashtags to $k = 5$, as the recall value does not improve much with additional recommendations and the precision value is still reasonable.

5 Related Work

The recommendation of Twitter hashtags can benefit from various other fields of research. These areas are (i) tagging of online resources, (ii) traditional recommender systems, (iii) social network analysis and (iv) Twitter analysis. However, to the best of our knowledge, there is no other approach aiming at recommending hashtags to Twitter users.

The recommendation of tags of online resources like images, bookmarks or bibliographic entries is directly related to our approach. Such approaches can be based on the co-occurrence of tags, like e.g. in [14, 20]). The notion of co-occurrence of tags describes the fact that two tags are used to tag the same photo. Therefore, only partly tagged photos can be subject to tag recommendations. Based on these relatively simple approaches, the paper by Rae *et al.* [17] proposes a method for Flickr tag recommendations which takes different contexts into account. Rae distinguishes four different contexts for the computation of recommendations: (i) the user’s previously used tags, (ii) the tags of the user’s contacts, (iii) the tags of the users which are members of the same groups as the user and (iv) the collectively most used tags by the whole community. A similar approach has also been facilitated by Garg and Weber in [6]. Another example for recommendations of tags is based on the BibSonomy platform which

basically allows its users to tag bibliographic entries [14]. This approach extracts tags which might be suitable for the entry from the title of the entry, the tags previously used for the entry and tags previously used by the current user. Based on these resources, the authors propose different approaches for merging these sets of tags. The resulting set is subsequently recommended to the user. Jäschke *et al.* [10] propose a collaborative filtering approach for the computation of tag recommendations. This computation is based on a graph consisting of the users, their tags and the tagged resources. After having constructed this graph, a PageRank-like ranking algorithm (called FolkRank) is applied. Furthermore, [2,15] are mainly concerned with the motivation of users to tag resources. John Hannon *et al.* [7] developed the Twittomender system which facilitates an approach for the recommendation of followees. This is done by creating profiles of users and applying a collaborative filtering approach to these profiles. The Twittomender system also provides search functionality (based on arbitrary keywords) which returns profile information about the found users like e.g. the latest popular keywords used by the specific user or his latest tweet.

Another approach directly connected to Twitter and recommendations is described by Phelan *et al.* [16]. In this approach, Twitter is used for the recommendation of news articles. In particular, Twitter is used to rank the news stories originating from various RSS feeds based on the user's tweets, the user's friends tweets or the public most recent tweets. Also, Jilian Chen *et al.* [5] focused on recommendations based on tweets. In this case, interesting URLs are recommended to the user. Romero *et al.* [19] analyzed how hashtags spread within the Twitter Universe. The hashtags were analyzed with regards to how a hashtag might be used by a user who is exposed to this hashtag by his followers and followees. The authors categorized the top-500 hashtags used within their data set and found that the adoption of hashtags is dependent on the category of the hashtags. E.g. multiple exposure to a hashtag for political or sports topics lead to the adoption of the hashtag with a higher probability than in any other hashtag category.

Kwak *et al.* [13] did a thorough analysis of the Twitter universe focusing on information diffusion within the network. Further analysis of Twitter messages are also contained in [3,11,12,21]. There have been numerous papers throughout the last years addressing the social aspects of Twitter and social online networks in general. Huberman *et al.* [9] found that the Twitter network basically consists of two networks: one dense network consisting of all followers and followees and one sparse network consisting of the actual friends of users. Huberman defines a friend of a user as another Twitter user with whom the user exchanged at least two directed messages. [4] contains an analysis of the retweet messages and [8] is concerned with how Twitter might be suitable for collaboration by exchanging direct messages.

As for the recommender system facilitated in our approach, many publications are focused around collaborative filtering. The papers by Resnick [18] and Adomavicius [1] provide a very good overview about the field of collaborative filtering.

6 Conclusion

In this paper, we presented an approach for the recommendation of hashtags within the Twitter microblogging application. The presented algorithm is based on the analysis of similar tweets and the hashtags contained in these tweets. Our evolutions were based on a self-crawled data set consisting of 12 million tweets. The preliminary evaluations showed promising results as the recall values of the recommendations are about 45-50%. Future work will include integrating the social graph of Twitter users for the recommendation. Furthermore, the ranking of hashtag recommendation candidates is also subject to further research and improvements. The enhancement of the recommendations of synonymous hashtags based on a semantic analysis for the exclusion of synonymous hashtags and their recommendation is also part of future work.

References

1. G. Adomavicius and A. Tuzhilin. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE transactions on knowledge and data engineering*, 17(6):734–749, 2005.
2. M. Ames and M. Naaman. Why we tag: motivations for annotation in mobile and online media. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, CHI '07, pages 971–980, New York, NY, USA, 2007. ACM.
3. S. Asur, B. Huberman, G. Szabo, and C. Wang. Trends in Social Media: Persistence and Decay. *Arxiv preprint arXiv:1102.1402*, 2011.
4. D. Boyd, S. Golder, and G. Lotan. Tweet, tweet, retweet: Conversational aspects of retweeting on twitter. In *hicss*, pages 1–10. IEEE Computer Society, 1899.
5. J. Chen, R. Nairn, L. Nelson, M. Bernstein, and E. Chi. Short and tweet: experiments on recommending content from information streams. In *Proceedings of the 28th international conference on Human factors in computing systems*, pages 1185–1194. ACM, 2010.
6. N. Garg and I. Weber. Personalized, interactive tag recommendation for flickr. In *Proceedings of the 2008 ACM conference on Recommender systems*, RecSys '08, pages 67–74, New York, NY, USA, 2008. ACM.
7. J. Hannon, M. Bennett, and B. Smyth. Recommending twitter users to follow using content and collaborative filtering approaches. In *RecSys '10: Proceedings of the fourth ACM conference on Recommender systems*, pages 199–206, New York, NY, USA, 2010. ACM.
8. C. Honeycutt and S. C. Herring. Beyond microblogging: Conversation and collaboration via twitter. In *HICSS*, pages 1–10. IEEE Computer Society, 2009.
9. B. Huberman, D. Romero, and F. Wu. Social networks that matter: Twitter under the microscope. *First Monday*, 14(1):8, 2009.
10. R. Jaeschke, L. Marinho, A. Hotho, L. Schmidt-Thieme, and G. Stumme. Tag Recommendations in Folksonomies. In J. Kok, J. Koronacki, R. Lopez de Mantaras, S. Matwin, D. Mladenic, and A. Skowron, editors, *Knowledge Discovery in Databases: PKDD 2007*, volume 4702 of *Lecture Notes in Computer Science*, pages 506–514. Springer Berlin / Heidelberg, 2007.
11. A. Java, X. Song, T. Finin, and B. Tseng. 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*, pages 56–65. ACM, 2007.
12. B. Krishnamurthy, P. Gill, and M. Arlitt. A few chirps about twitter. In *Proceedings of the first workshop on Online social networks*, pages 19–24. ACM, 2008.
 13. H. Kwak, C. Lee, H. Park, and S. Moon. What is Twitter, a social network or a news media? In *Proceedings of the 19th international conference on World wide web*, pages 591–600. ACM, 2010.
 14. M. Lipczak and E. Milius. Learning in efficient tag recommendation. In *Proceedings of the fourth ACM conference on Recommender systems*, RecSys '10, pages 167–174, New York, NY, USA, 2010. ACM.
 15. C. Marlow, M. Naaman, D. Boyd, and M. Davis. HT06, tagging paper, taxonomy, Flickr, academic article, to read. In *Proceedings of the seventeenth conference on Hypertext and hypermedia*, page 40. ACM, 2006.
 16. O. Phelan, K. McCarthy, and B. Smyth. Using twitter to recommend real-time topical news. In *Proceedings of the third ACM conference on Recommender systems*, pages 385–388. ACM, 2009.
 17. A. Rae, B. Sigurbjörnsson, and R. van Zwol. Improving tag recommendation using social networks. In *Adaptivity, Personalization and Fusion of Heterogeneous Information*, RIAO '10, pages 92–99, Paris, France, France, 2010. Le Centre de Hautes Etudes Internationales d'Informatique Documentaire.
 18. P. Resnick and H. Varian. Recommender systems. *Communications of the ACM*, 40(3):58, 1997.
 19. D. M. Romero, B. Meeder, and J. M. Kleinberg. Differences in the mechanics of information diffusion across topics: idioms, political hashtags, and complex contagion on twitter. In S. Srinivasan, K. Ramamritham, A. Kumar, M. P. Ravindra, E. Bertino, and R. Kumar, editors, *WWW*, pages 695–704. ACM, 2011.
 20. B. Sigurbjörnsson and R. Van Zwol. Flickr tag recommendation based on collective knowledge. In *Proceeding of the 17th international conference on World Wide Web*, pages 327–336. ACM, 2008.
 21. S. Ye and S. Wu. Measuring Message Propagation and Social Influence on Twitter. com. In *Social Informatics: Second International Conference, Socinfo 2010, Laxenburg, Austria, October 27-29, 2010, Proceedings*, page 216. Springer-Verlag New York Inc, 2010.