

Solving the Cold-Start Problem in Recommender Systems with Social Tags

ZI-KE ZHANG¹, CHUANG LIU^{2,3}, YI-CHENG ZHANG^{1,4} ^(a) and TAO ZHOU^{4,5}

¹ *Department of Physics, University of Fribourg, Chemin du Musée 3, 1700 Fribourg, Switzerland*

² *School of Business, East China University of Science and Technology, Shanghai 200237, P. R. China*

³ *Engineering Research Center of Process Systems Engineering (Ministry of Education), East China University of Science and Technology, Shanghai 200237, P. R. China*

⁴ *Web Sciences Center, University of Electronic Science and Technology of China, Chengdu 610054, P. R. China*

⁵ *Department of Modern Physics, University of Science and Technology of China, Hefei 230026, P. R. China*

PACS 89.20.Ff – Computer science and technology

PACS 89.75.Hc – Networks and genealogical trees

PACS 89.65.-s – Social and economic systems

Abstract. - In this Letter, based on the user-tag-object tripartite graphs, we propose a recommendation algorithm that makes use of social tags. Besides its low cost of computational time, the experimental results on two real-world data sets, *Del.icio.us* and *MovieLens*, show that it can enhance the algorithmic accuracy and diversity. Especially, it provides more personalized recommendation when the assigned tags belong to diverse topics. The proposed algorithm is particularly effective for small-degree objects, which reminds us of the well-known *cold-start problem* in recommender systems. Further empirical study shows that the proposed algorithm can significantly solve this problem in social tagging systems with heterogeneous object degree distributions.

Introduction. – Many complex systems can be well described by networks where nodes represent individuals, and edges denote the relations among them [1–5]. Recently, the personalized recommendation in complex networks has attracted increasing attention from physicists [6–10]. Personalized recommendation aims at finding objects (e.g. books, webpages, music, etc.) that are most likely to be collected by users. For example, classical information retrieval can be viewed as recommending documents with given words [11], and the process of link prediction can be considered as a recommendation problem in unipartite networks [12–14]. The central problem of personalized recommendation can be divided into two parts: one is the estimation of similarity based on the historical records of user activities [15, 16]; the other is the usage of accessorial information (e.g., object attributes) to efficiently filter out irrelevant objects. For the formal task, since computing and storing the similarities of all user pairs is costly, we usually consider only the top- k most similar users [17]. For the latter task, very accurate descriptions of objects may be helpful in filtering irrelevant

objects, however, it is limited to the attribute vocabulary, and, on the other hand, attributes describe global properties of objects which are less helpful to generate personalized recommendations.

Recently, the advent of Web2.0 and its affiliated applications bring a new form of paradigm, *social tagging systems* (or called *collaborative tagging systems*), which introduces a novel platform for users' participation. A social tagging system allows users to freely assign tags to annotate their collections, requires no specific skills for users to participate in, broadens the semantic relations among users and objects, and thus has attracted much attention from the scientific community. Golder *et al.* studied its usage patterns and classified seven kinds of tag functions [18]. Similar to the tagging functions, the keywords and PACS numbers are analyzed to better characterize the structure of co-authorship and citation networks [19, 20]. Furthermore, many efforts have been done to explain the emergent properties of social tagging systems. Cattuto *et al.* [21] proposed a memory-based Yule-Simon model to describe the aging effects and occurrence frequencies of tags. Zhang and Liu [22] proposed an evolutionary hypergraph model, where users not only assign tags to objects but also

^(a)E-mail:yi-cheng.zhang@unifr.ch

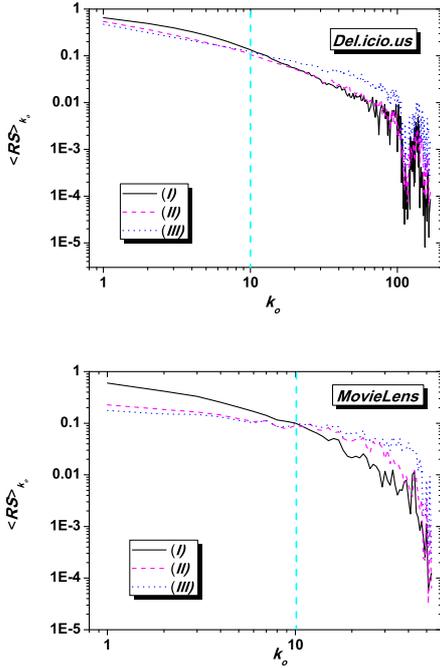


Fig. 1: (Color online) Object-degree-dependent ranking score for the three algorithms in *Del.icio.us* and *MovieLens*. Each data point is obtained by averaging over 50 realizations, each of which corresponds to an independent division of training set and testing set.

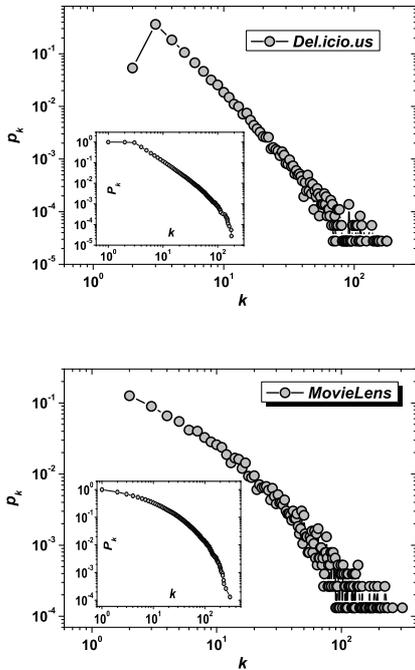


Fig. 2: Object degree distributions of the two data sets. The insets show the accumulative distributions

retrieve objects via tags.

Besides, social tagging systems have already found wide applications in *Recommender Systems*. By considering the tag frequency as weight, Szomszor *et al.* [23] proposed an improved movie recommendation algorithm. Schenkel *et al.* [24] proposed an incremental threshold algorithm taking into account both the social ties among users and semantic relations of different tags, which performs remarkably better than the algorithm without tag expansion. Zhang *et al.* [25] and Shang *et al.* [26] proposed an object-based and user-based hybrid tag algorithm, respectively, harnessing diffusion-based methods to obtain better recommendations. Shang and Zhang [27] considered the tag usage frequency as edge weight in a user-object bipartite network and improved the accuracy of recommendation.

In this Letter, we propose a diffusion-based recommendation algorithm which considers social tags as a bridge connecting users and objects. That is to say, users can efficiently find the target objects via tags. In particular, we consider the usage frequencies of tags as users' personal preference, while the semantic relations between tags and objects as global information. Experimental results show that the present algorithm can significantly improve the recommendation accuracy. Further empirical study shows that the proposed algorithm is especially effective for the objects collected by few users, which reminds us of the well-known *cold-start* problem [28, 29]. Since there is little information available for new objects, social tags can effectively build up relations between existing objects and the new ones. Therefore, the incorporating of tags can remarkably help users find the new (or less popular) yet interesting objects, and thus enhance the overall accuracy. In addition, we employ entropy-based and Hamming-distance-based methods to measure the *inner*- and *inter*-diversity of tag usage patterns, respectively. Experimental results show that there are different tag usage patterns in the two datasets: users assign more diverse tags in *Del.icio.us* than *MovieLens*, and it might shed lights on understanding why the proposed algorithm can enhance the recommendation diversity in *Del.icio.us* largely than *MovieLens*.

Data. – The empirical data used in this paper include:

(i) *Del.icio.us* – one of the most popular social bookmarking web sites, which allows users not only to store and organize personal bookmarks (URLs), but also to look into other users' collections and find what they might be interested in by simply keeping track of the baskets with social tags; (ii) *MovieLens* – a movie rating system, where each user votes movies in five discrete ratings 1-5. A tagging function is added in from January 2006. In both data sets, we remove the isolated nodes and guarantee that each user has collected at least one object, each object has been collected by at least two users, assigned by at least two tags, and each tag is used by at least twice by every adjacent user. Table 1 summarizes the basic statistics of the purified data sets.

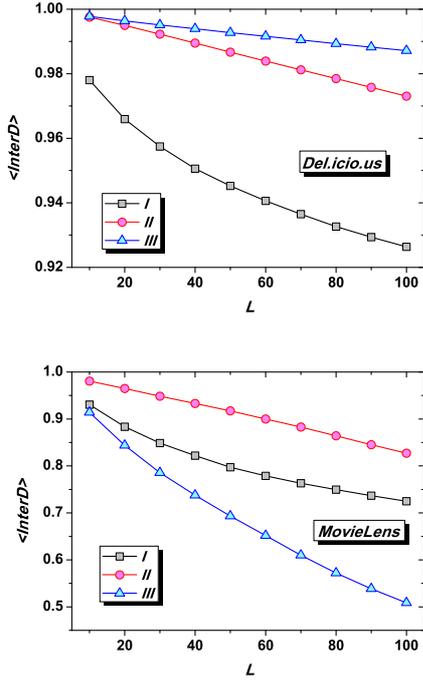


Fig. 3: (Color online) $\langle InterD \rangle$ as a function of the length of recommendation list for the three algorithms in *Del.icio.us* and *MovieLens*.

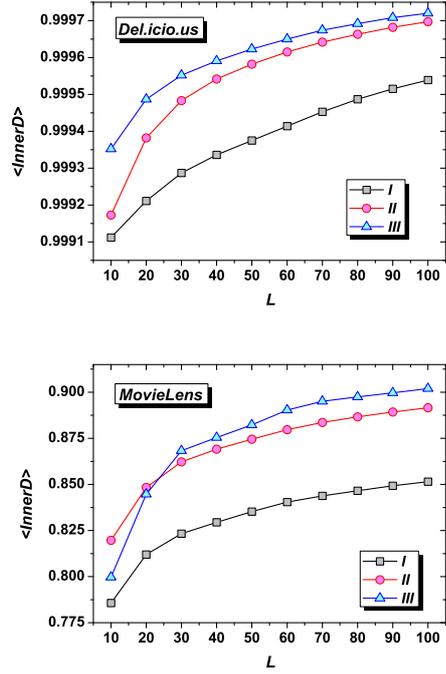


Fig. 4: (Color online) $\langle InnerD \rangle$ as a function of the length of recommendation list for the three algorithms in *Del.icio.us* and *MovieLens*.

Every data set is consisted of many entries, and each follows the form $\mathbb{F} = \{\text{user}, \text{object}, \text{tag}_1, \text{tag}_2, \dots, \text{tag}_t\}$, where t is the number of tags assigned to this object by this user. Then each data set is randomly divided into two parts: the training set, is treated as known information, while the testing set is used for testing. In this Letter, the training set always contains 90% of entries and the remaining 10% of entries constitute the testing set.

Algorithms. – A recommender system considered in this Letter consists of three sets, respectively of users $U = \{U_1, U_2, \dots, U_n\}$, objects $O = \{O_1, O_2, \dots, O_m\}$, and tags $T = \{T_1, T_2, \dots, T_r\}$. The tripartite graph representation can be described by three matrices, A , A' and A'' for user-object, object-tag and user-tag relations. If U_i has collected O_j , we set $a_{ij} = 1$, otherwise $a_{ij} = 0$. Analogously, we set $a'_{jk} = 1$ if O_j has been assigned by the tag T_k , and $a'_{jk} = 0$ otherwise. Furthermore, the users' preferences on tags can be represented by a weighted matrix A'' , where a''_{ik} is the number of times that U_i has adopted T_k .

Subsequently, we introduce the proposed algorithm, as well as two baseline ones: (I) user-object diffusion [8]; (II) user-object-tag diffusion [25]; (III) user-tag-object diffusion. Given a target user U_i , the above three algorithms will generate final score of each object, f_j , that are pushed into recommendation resource for him/her, are described as following:

(I) Supposing that a kind of resource is initially located

on objects. Each object averagely distributes its resource to all neighboring users, and then each user redistributes the received resource to all his/her collected objects. The final resource vector for the target user U_i , \vec{f} , after the two-step diffusion is:

$$f_j = \sum_{l=1}^n \sum_{s=1}^m \frac{a_{lj} a_{ls} a_{is}}{k(U_l) k(O_s)}, \quad j = 1, 2, \dots, m, \quad (1)$$

where $k(U_l) = \sum_{j=1}^m a_{lj}$ is the number of collected objects for user U_l , and $k(O_s) = \sum_{i=1}^n a_{is}$ is the number of neighboring users for object O_s .

(II) The initial resources are set as same as I, but each object equally distributes its resource to all neighboring tags, and then each tag redistributes the received resource to all its neighboring objects. Thus, the final resource vector, \vec{f}' , is:

$$f'_j = \sum_{l=1}^r \sum_{s=1}^m \frac{a'_{jl} a'_{ls} a_{is}}{k'(T_l) k'(O_s)} \quad j = 1, 2, \dots, m, \quad (2)$$

where $k'(T_l) = \sum_{j=1}^m a'_{jl}$ is the number of neighboring objects for tag T_l , $k'(O_s) = \sum_{l=1}^r a'_{sl}$ is the number of neighboring tags for object O_s .

(III) Different from I and II, here, the initial resources are located on tags according to their frequencies used by the target user U_i . Then each tag distributes the initial

Table 1: Basic statistics of the two data sets. n , m , r are the total numbers of users, objects and tags, respectively. $\langle k \rangle$, $\langle k' \rangle$ and $\langle k'' \rangle$ denote the average number of objects collected by a user, tags assigned by an object and tags adopted by a user respectively. *Del.* and *Mov.* represent the data sets *Del.icio.us* and *MovieLens*, respectively.

Data	n	m	r	$\langle k \rangle$	$\langle k' \rangle$	$\langle k'' \rangle$
<i>Del.</i>	4902	36224	10584	43.85	38.82	286.86
<i>Mov.</i>	648	1590	1382	15.04	19.89	22.89

resource directly to all its neighboring objects. Thus, the final resource vector, f_j'' , reads:

$$f_j'' = \sum_{l=1}^r \frac{a'_{jl} a''_{il}}{k'(T_l)} \quad (3)$$

After we obtain the final score of objects, all the objects that U_i has not collected are ranked in a descending order, and the top L objects will be recommended to U_i .

Comparing with algorithms I and II, the advantages of algorithm III are threefold. Firstly, since social tags highly reflect users' personal preferences, algorithm III is promisingly expected to generate more personalized recommendation. Secondly, the one-step diffusion can clearly save computational time especially for large-scale data. Thirdly, algorithm III reveals the essential role of tags: building a bridge between users and objects, helping users retrieve and organize collections without the limit of hierarchical structure and vocabulary of words.

Metrics. — To give solid and comprehensive evaluation of the proposed algorithm, we employ three different metrics that characterizing the accuracy and diversity of recommendations.

1. *Ranking Score (RS)* [8].— In the present case, for each entry in the testing set (i.e. a user-object pair), *RS* is defined as the rank of the object, divided by the number of all uncollected objects for the corresponding user. Apparently, the less the *RS*, the higher accuracy the algorithm is. The average ranking score $\langle RS \rangle$ is given by averaging over all entries in the testing set.
2. *Inter Diversity (InterD)* [8, 31].— *InterD* measures the differences of different users' recommendation lists, thus can be understood as the inter-user diversity. Denote O_R^i the set of recommended objects for user U_i , then

$$InterD = \frac{2}{n(n-1)} \sum_{i \neq j} \left(1 - \frac{|O_R^i \cap O_R^j|}{L} \right), \quad (4)$$

where $L = |O_R^i|$ is the length of recommendation list. In average, greater or less *InterD* mean respectively greater or less personalization of users' recommendation lists.

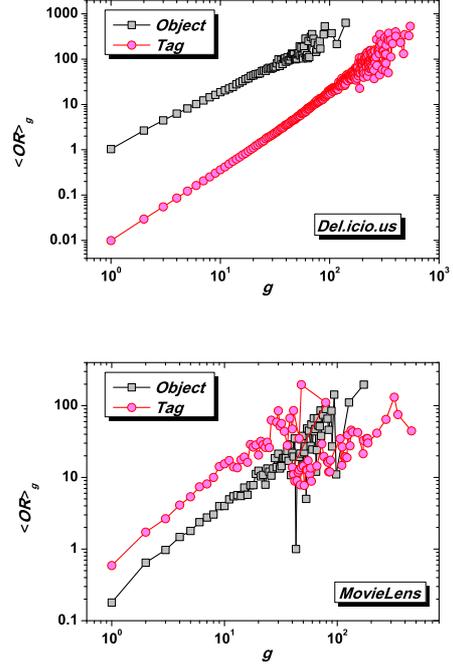


Fig. 5: (Color online) $\langle OR \rangle$ as a function of g for the two data sets. The black squares represent $\langle OR \rangle$ for objects and the red circles are $\langle OR \rangle$ of tags, respectively.

3. *Inner Diversity (InnerD)* [31].— *InnerD* measures the differences of objects within a user's recommendation list, thus can be considered as the inner-user diversity. It reads,

$$InnerD = 1 - \frac{2}{nL(L-1)} \sum_{i=1}^n \sum_{j \neq l, j, l \in O_R^i} S_{jl}, \quad (5)$$

where $S_{jl} = \frac{|\Gamma_{O_j} \cap \Gamma_{O_l}|}{\sqrt{|\Gamma_{O_j}| \times |\Gamma_{O_l}|}}$ is the cosine similarity between objects O_j and O_l , where Γ_{O_j} denotes the set of users having collected object O_j . In average, greater or less *InnerD* suggests respectively greater or less topic diversification of users' recommendation lists.

Results. — To make clear the role of social tags, a microscopic picture of algorithmic accuracy is very helpful. Especially, since social tags are used to describe the objects, we would like to see the dependence of accuracy on object degree, namely the number of users collecting it. Given an object degree k_o , the degree-dependent average ranking score, denoted by $\langle RS \rangle_{k_o}$, is defined as the mean positions averaged over all the entries in the testing set with object degree equal to k_o .

In Table 2 and Table 3, we give the overall $\langle RS \rangle$ of the three algorithms for the observed data sets. It indicates that the $\langle RS \rangle$ is significantly enhanced by the present algorithm. Fig. 1 reports the correlation between accuracy

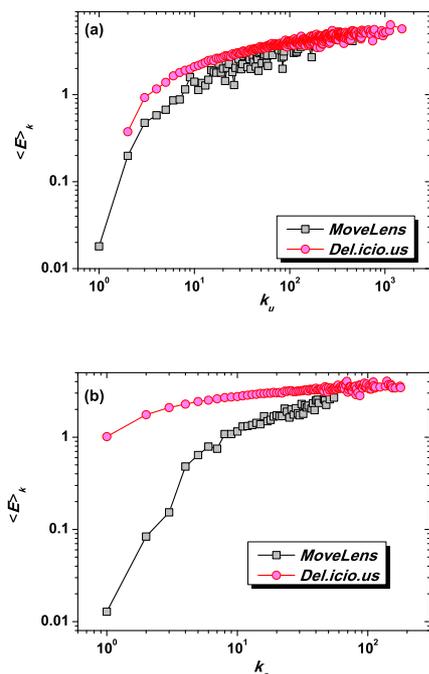


Fig. 6: (Color online) (a) $\langle E \rangle$ as a function of user degree; (b) $\langle E \rangle$ as a function of object degree, respectively.

and object degree. The ranking score decays with the increasing k_o for all the three algorithms. In addition, the three curves intersect around $k_o=10$, which is a relatively small value considering the heterogeneous object-degree distribution shown in Fig. 2. From Fig. 1, it is seen that the algorithmic accuracy of algorithm III is better than that of algorithms I or II for $k_o \leq 10$, but worse when $k_o > 10$ (see also Table 2 and Table 3), which reminds us of the well-known cold-start problem in recommender systems: how to recommend the unpopular and/or new objects to users? It is very difficult for a user to be aware of these *cold objects* by random surfing since they are not hot items, and for a recommender system to recommend them to right places since there are usually insufficient information about them. In fact, there are 90.04% and 69.35% objects with $k_o \leq 10$ in *Del.icio.us* and *MovieLens*, respectively. Therefore, a successful recommender system has to make reasonable recommendations of cold objects. Comparing with the algorithms I and II, the present one can effectively help users find those cold objects via social tags.

Fig. 3 and Fig. 4 show the experimental results of $\langle InterD \rangle$ and $\langle InnerD \rangle$, respectively. In Fig. 3, $\langle InterD \rangle$ is enhanced only for *Del.icio.us*. The reason for small $\langle InterD \rangle$ of algorithm III in *MovieLens* is that there are only movies in that data set, and thus a comparatively small number of tags are used with huge overlapping. The overlapping ratio, OR , of tags for users to assign to the

Table 2: Algorithmic accuracy for *Del.icio.us*. $\langle RS \rangle_{k_o \leq 10}$ is the average ranking score over objects with degree equal or less than 10, and $\langle RS \rangle_{k_o > 10}$ is the average ranking scores over objects with degree greater than 10. Each value is obtained by averaging over 50 realizations, each of which corresponds to an independent division of training set and testing set.

Algorithms	$\langle RS \rangle$	$\langle RS \rangle_{k_o \leq 10}$	$\langle RS \rangle_{k_o > 10}$
I	0.276	0.369	0.054
II	0.209	0.275	0.049
III	0.196	0.249	0.068

same objects, is defined as:

$$OR_g = \frac{1}{N_g} \sum_{i \neq j, G(i,j)=g} OR(i,j), \quad (6)$$

where N_g is the number of user pairs (i,j) such that $i \neq j$, and $G(i,j) = g$ denotes the number of common objects collected by users i and j . $OR(i,j)$ is defined as the total number of tag agreements on the same objects for user pair (i,j) . Similar definition can also be used to quantify the overlapping ratio of objects collected by users with the same tags. Clearly, larger OR indicates smaller diversity, and vice versa. Fig. 5 shows the correlation between $\langle OR \rangle_g$ and g . One can see that $\langle OR \rangle_g$ of tags is smaller than that of objects in *Del.icio.us*, while it is not the case for *MovieLens*. In a word, social tags can help generate more diverse recommendation only if the tags are themselves used in a diverse way.

Fig. 4 shows that $\langle InnerD \rangle$ is generally improved by our proposed algorithm, indicating that it can help users broaden their horizons. Except for *MovieLens* with very small L . It is again resulted from the narrow choice of tags in *MovieLens*. Recently, the *Shannon entropy* is widely used to quantify network diversity in social sharing networks [32] and social economics [33]. In the Letter, we also employ it to measure individual usage pattern of tags:

$$E(U_i) = - \sum_t p_{i;t} \ln(p_{i;t}), \quad (7)$$

where $p_{i;t}$ is the probability for tag t used by user U_i . Then the dependence of entropy on user degree, E_k , is given by averaging all the $E(U_i)$ with $k(U_i) = k$. Similar definition can be used to quantify the dependence of entropy for objects. Clearly, Larger E_k means that the users are more willing to use diverse topics of tags, or the objects are more likely to be assigned to more diverse tags, and vice versa. Fig. 6 shows that E of *Del.icio.us* are greater than that of *MovieLens* for both users and objects, indicating that *Del.icio.us* is a more diverse system than *MovieLens*, and further giving a reasonable explanation why algorithm III can obtain better $InnerD$ in *Del.icio.us* than *MovieLens*.

Conclusions and Discussion. – In this Letter, we proposed a recommendation algorithm making use of social tags. This algorithm, considers the frequencies of tags

Table 3: Algorithmic accuracy for *MovieLens*.

Algorithms	$\langle RS \rangle$	$\langle RS \rangle_{k_o \leq 10}$	$\langle RS \rangle_{k_o > 10}$
I	0.207	0.307	0.039
II	0.130	0.168	0.055
III	0.123	0.146	0.070

as user preferences on different topics and tag-object links as semantical relations between them. Experimental results demonstrated that the proposed algorithm outperforms the two baseline algorithms in both accuracy and diversity. The present algorithm outperforms others especially for the objects with small degrees ($k_o \leq 10$), which constitute the majority of objects. Therefore, the incorporating of social tags could be, to some extent, helpful in solving the cold-start problem of recommender systems.

Recently, besides the accuracy, the significance of diversity has attracted more and more attention in information filtering [10]. Experimental results in this Letter demonstrated that a wide-range adoption of social tags can enhance the diversity of recommendation. Therefore, we strongly encourage recommender systems to add tagging functions and users to organize their collections by using tags. However, despite the significant role of tags, the polysemy and synonymy problems [18] might result in coarse and inaccurate performance, the tag clustering technique [34] is hopefully to provide a promising way to generate multi-scale recommendations and eventually obtain the best performance.

* * *

This work is partially supported by the Swiss National Science Foundation (Project 200020-121848). Z.-K.Z. and T.Z. acknowledge the National Natural Science Foundation of China under the grant no. 60973069.

REFERENCES

- [1] ALBERT R. and BARABÁSI A.-L., *Rev. Mod. Phys.*, **74** (2002) 47.
- [2] DOROGOVTSSEV S. N. and MENDES J. F. F., *Adv. Phys.*, **51** (2002) 1079.
- [3] NEWMAN M. E. J., *SIAM Rev.*, **45** (2003) 167.
- [4] BOCCALETTI S., LATORA V., MORENO Y., CHAVEZ M. and HUANG D.-U., *Phys. Rep.*, **424** (2006) 175.
- [5] COSTA L. DA F., RODRIGUES F. A., TRAVIESOR G. and BOAS P. R. U., *Adv. Phys.*, **56** (2007) 167.
- [6] ZHANG Y.-C., BLATTNER M. and YU Y.-K., *Phys. Rev. Lett.*, **99** (2007) 154301.
- [7] ZHANG Y.-C., MEDO M., REN J., ZHOU T., LI T. and YANG F., *EPL*, **80** (2007) 68003.
- [8] ZHOU T., REN J., MEDO M. and ZHANG Y.-C., *Phys. Rev. E*, **76** (2007) 046115.
- [9] ZHOU T., JIANG L.-L., SU R.-Q. and ZHANG Y.-C., *EPL*, **81** (2007) 58004.
- [10] ZHOU T., KUSCSIK Z., LIU J.-G., MEDO M., WAKELING J. R. and ZHANG Y.-C., *Proc. Natl. Acad. Sci. U.S.A.*, **107** (2010) 4511.
- [11] SALTON G. and MCGILL M. J., *Introduction to Model Information Retrieval* (McGraw-Hill, Auckland, 1983).
- [12] ZHOU T., LÜ L. and ZHANG Y.-C., *Eur. Phys. J. B*, **71** (2009) 623-630.
- [13] LÜ L. and ZHOU T., *EPL*, **89** (2010) 18001.
- [14] LIU W. and LÜ L., *EPL*, **89** (2010) 58007.
- [15] BALABANOVIĆ M. and SHOHAM Y., *Commun. ACM*, **40** (1997) 72.
- [16] SARWAR B., KARYPIS G., KONSTAN J. and RIEDL J., *Proc. the 10th Intl. Conf. WWW* (ACM Press, New York) 2001, pp. 295-305.
- [17] DESHPANDE M. and KARYPIS G., *ACM Trans. Inf. Syst.*, **22** (2004) 143.
- [18] GOLDBER S. A. and HUBERMAN B. A., *J. Info. Sci.*, **32** (2006) 198.
- [19] PALLA G., FARKAS I. J., POLLNER P., DERÉYI I. and VICSEK T., *New J. Phys.*, **10** (2008) 123026.
- [20] ZHANG Z.-K., LÜ L., LIU J.-G. and ZHOU T., *Eur. Phys. J. B*, **66** (2008) 557.
- [21] CATTUTO C., LORETO V. and PIETRONERO L., *Proc. Natl. Acad. Sci. USA*, **104** (2007) 1461.
- [22] ZHANG Z.-K. and LIU C., arXiv: 1003.1931.
- [23] SZOMSZOR M., CATTUTO C., ALANI H., OHARA K., BALDASSARRI A., LORETO V. and SERVEDIO V. D. P., *Proc. the 4th Euro. Semantic Web Conf.* (Innsbruck, Austria) 2007, pp. 71-84.
- [24] SCHENKEL R., CRECELIUS T., KACIMI M., MICHEL S., NEUMANN T., PARREIRA J. X. and WEIKUM G., *Proc. the 31st Annual Intl. ACM SIGIR Conf. Res. Dev. Info. Retr.* (ACM Press, New York) 2008, pp. 523-530.
- [25] ZHANG Z.-K., ZHOU T. and ZHANG Y.-C., *Physica A*, **389** (2010) 179.
- [26] SHANG M.-S., ZHANG Z.-K., ZHOU T. and ZHANG Y.-C., *Physica A*, **389** (2010) 1259.
- [27] SHANG M.-S. and ZHANG Z.-K., *Chin. Phys. Lett.*, **26** (2009) 118903.
- [28] SCHEIN A.I., POPESCU A., UNGAR L.H. and PENNOCK D.M., *Proc. 2001 SIGIR Workshop Recomm. Syst.* (New Orleans, LA) 2001.
- [29] SCHEIN A. I., POPESCU A., UNGAR L. H. and PENNOCK D. M., *Proc. 25th Annual Intl. ACM SIGIR Conf. Research and Development in Information Retrieval* (ACM Press, New York) 2002, pp. 253-260.
- [30] PARK S.T., PENNOCK D., MADANI O., GOOD N. and DE-COSTE D., *Proc. 12th ACM SIGKDD Intl. Conf. Knowledge Discovery and Data Mining* (ACM Press, New York) 2002, pp. 705-711.
- [31] ZHOU T., SU R.-Q., LIU R.-R., JIANG L.-L., WANG B.-H. and ZHANG Y.-C., *New J. Phys.*, **11** (2009) 123008.
- [32] LAMBIOTTEA R. and AUSLOOSB M., *Eur. Phys. J. B*, **50** (2006) 183.
- [33] EAGLE N., MACY M. and CLAXTON R., *Science*, **328** (2010) 1029.
- [34] SHEPITSSEN A., GEMMELL J., MOBASHER B. and BURKE R., *Proc. the 2008 ACM Conf. Recomm. Syst.* (ACM Press, New York) 2008, pp. 259-266.