

RESEARCH

Open Access



# M-BiRank: co-ranking developers and projects using multiple developer-project interactions in open source software community

Dengcheng Yan<sup>1,2</sup> , Bin Qi<sup>1</sup>, Yiwen Zhang<sup>1\*</sup> and Zhen Shao<sup>1</sup>

\*Correspondence:

zhangyiwen@ahu.edu.cn

<sup>1</sup>School of Computer Science and Technology, Anhui University, Jiulong Road No. 111, Hefei, 230601, China

Full list of author information is available at the end of the article

## Abstract

Social collaborative coding is a popular trend in software development, and such platforms as GitHub provide rich social and technical functionalities for developers to collaborate on open source projects through multiple interactions. Developers often follow popular developers and projects for learning, technical selection, and collaboration. Thus, identifying popular developers and projects is very meaningful. In this paper, we propose a multiplex bipartite network ranking model, M-BiRank, to co-rank developers and projects using multiple developer-project interactions. Firstly, multiple developer-project interactions such as commit, issue, and watch are extracted and a multiplex developer-project bipartite network is constructed. Secondly, a random layer is selected from this multiplex bipartite network and initial ranking scores are calculated for developers and projects using BiRank. Finally, initial ranking scores diffuse to other layers and mutual reinforcement is taken into consideration to iteratively calculate ranking scores of developers and projects in different layers. Experiments on real-world GitHub dataset show that M-BiRank outperforms degree centrality, traditional single layer ranking methods, and multiplex ranking method.

**Keywords:** M-BiRank, Ranking, Multiplex bipartite network, Social collaborative coding, GitHub

## 1 Introduction

Open source software community is now a main driven force of innovations, and plenty of software developers collaborate on millions of open source software projects, among which are many popular software projects that drive the innovations of different fields [1, 2]. For example, deep learning frameworks such as TensorFlow, PyTorch, and MXNet contributed by famous companies simplify the building of deep learning models, which to some extent speed up the innovations in the field of artificial intelligence in both academia and industry [3, 4].

In open source software community, developers from different areas usually take the social collaborative coding paradigm and participate in different portions of a common

project using the social and technical functionalities provided by the community [5, 6]. Taking GitHub as an example, developers from different areas and with different technical backgrounds can collaborate on a project by committing codes and commenting on issues, star/fork a project for improving technical skills or technical selections, and follow other professional developers for keeping pace with new trends. Much like the role of opinion leaders in social networks, influential developers and projects drive the technical trends and the prosperity of open source community. Thus, identifying influential developers and projects will be of great significance. Existing work on influence analysis for open source software community mainly focused on applying traditional unipartite single layer graph ranking methods [7, 8], including PageRank [9] and HITS [10], although many new graph ranking methods for more complex network structures, such as bipartite network [11, 12] and multiplex network [13, 14], have been proposed. On the other hand, existing graph ranking methods have not merged bipartite network and multiplex network as a single network model, which is necessary for our case to model multiple interactions between developers and projects. Potential applications of influence analysis of open source software community would include service recommendations [15–19] and risk assessment [20].

In this paper, we focus on modeling multiple interactions between developers and projects as a multiplex bipartite network and propose a new ranking method based on it in an iterative and mutually enhanced way. The main contributions of this research are many folds:

- We propose a multiplex bipartite network model to represent multiple interactions between developers and projects.
- We propose a new ranking model called M-BiRank on multiplex bipartite network which takes into account the mutual reinforcement between different types of nodes as well as different layers.
- We apply the proposed model to real-world GitHub dataset, showing that our model outperforms baseline ranking models.

The remainder of the paper is organized as follows. Section 2 gives a brief introduction to related works on ranking models from the perspective of network and its applications in software engineering. In Section 3, details about the proposed M-BiRank model are illustrated. Then, the experiment results and discussions are given in Section 4. Finally, we briefly summarize our work and explain future directions in Section 5.

## 2 Related work

Identifying influential nodes in social networks has been a hot topic for decades. Existing works mainly focus on either structural properties or diffusion dynamics. Plenty of structure-based metrics and random walk-based methods have been proposed.

Structure-based metrics usually base itself on some intuition for centrality from either local or global views. Degree is the most common local structure-based centrality metrics. Based on degree, Chen et al. [21] proposed a semi-local centrality metric, considering both the nearest and the next nearest neighbors. Chen et al. [22] further considered the negative impact of local clustering on information diffusion in networks and proposed ClusterRank. In addition to extending degree, several local structure-based centrality metrics are originated from H-index, which is originally used to measure the citation

impact of a scholar or a journal [23]. Zhao et al. [24] first extended H-index concept to networks and defined the h-Degree metrics for weighted networks. Liu et al. [25] combined the H-index of both node itself and its neighbors and proposed a local H-index centrality. Lü et al. [26] revealed the relation among degree, H-index, and coreness and introduced a family of H-indices.

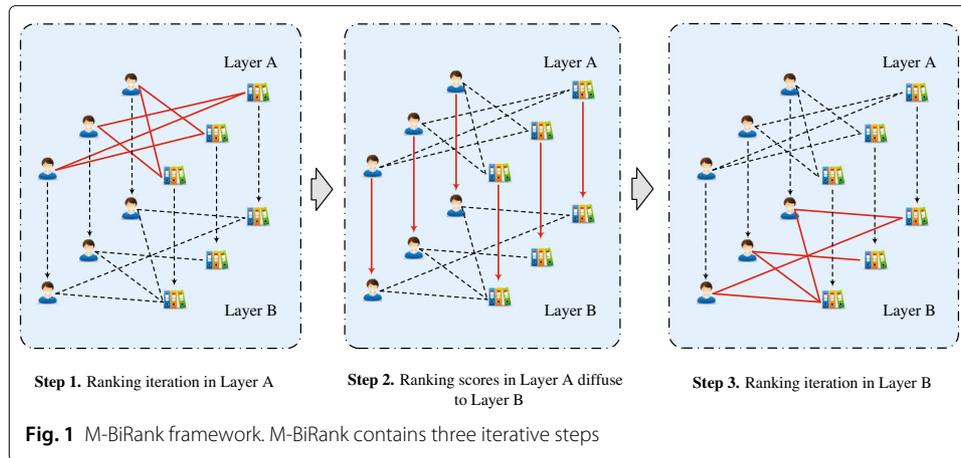
Local structure-based centrality metrics benefit from low computational complexity at the cost of reducing effectiveness. While global structure-based centrality metrics can better identify influential nodes from a global view of the whole network. Earlier researches in sociology introduced several global structure-based centrality metrics, including closeness centrality [27], betweenness centrality [28], and eigenvector centrality [29]. Recently, researches introduced eigenvector centrality to more complex network structures. Wang et al. [30, 31] extended eigenvector centrality to multilayer networks under a framework of tensor decomposition.

Random walk-based methods apply resource diffusion dynamic process in networks and measure node's influence according to the final resource the node obtains at stationary state of the dynamic process. Typical random walk-based ranking methods include PageRank [9] and HITS [10]. To solve the problem of dangling nodes of PageRank, Lü et al. [32] added a ground node to the original network, making the original network connected, and proposed the parameter-free LeaderRank method. Halu et al. [13] extended PageRank to multiplex network and proposed Multiplex PageRank, which included four kinds of intra-layer enhancement mechanism [33]. To address the ranking problem on bipartite networks, He et al. [11] proposed the BiRank method. Instead of modeling pairwise interactions, higher-order network models have recently been proposed and applied to ranking nodes with group interactions. Treating scientific collaboration as a group interaction, Liang et al. [34] modeled it as a hypergraph and proposed HHGBiRank. From another view of higher-order structure of networks, that is motif, Zhao et al. [35] proposed motif-based PageRank.

In addition to identifying influential nodes for general purpose social networks, needs have also emerged for open source software community, a special kind of social network. Xuan et al. [8] constructed several social networks based on the communications between developers in Apache and applied degree, PageRank, and HITS for developer ranking. Hu et al. [7] studied the problem of influence identification of developers in GitHub and proposed a Following-Star-Fork-Activity-based approach. Joblin et al. [5] employed several activity counts, centrality metrics, and network structural properties to distinguish core and peripheral developers.

### 3 Method

In this section, we will present a *Multiplex Bipartite Ranking* method, called *M-BiRank*, for co-ranking developers and projects in open source software community. As shown in Fig. 1, the proposed M-BiRank consists of three parts and incorporates two basic assumptions that address the issues in Section 1. We will start by giving the definition of multiplex bipartite network and introducing the notations. Then, thorough explanations and mathematical formulations are given for the two basic assumptions, that is, mutual reinforcement between different types of entities and between different layers of network. Finally, we will introduce the overall algorithm and time complexity analysis of it.



### 3.1 Multiplex bipartite network

**Definition 1** *Multiplex bipartite network.* A multiplex bipartite network is a set of  $M$  bipartite networks  $G^A = (U \cup P, E^A)$  where  $A = 1, \dots, M$ .  $U$  and  $P$  represent two different kinds of nodes and all layers have the same sets of nodes  $U$  and  $P$ , while the set of edges  $E^A$  depends on the layer  $A$ . Specially, there is no edge between same kinds of nodes. A multiplex bipartite network is formally defined as  $G = (G^1, \dots, G^M)$ . Each network  $G^A$  is described by the bipartite weight matrix  $W^A \in \mathbb{R}^{|U| \times |P|}$  with elements  $w_{ij}^A$ , where  $w_{ij}^A > 0$  if there is a link with weight  $w_{ij}^A$  between nodes  $u_i$  and  $p_j$  in layer  $A$ , and  $w_{ij}^A = 0$  otherwise.

In this paper, we model multiple interactions between developers and projects as a developer-project multiplex bipartite network. The notations we will use throughout the article are summarized in Table 1.

### 3.2 Mutual reinforcement between developers and projects

Most ranking methods in networks adopt the intuition as PageRank and HITS that an influential node should be linked by many other influential nodes, which is also applicable in the case of developer-project bipartite network. For example, an elite developer usually participates in popular projects. In open source software community, it is quite a practice to estimate the influence of a developer by how popular the projects she/he participates in are and how much she/he contributes to these popular projects. And a project with influential developers or organizations as major contributors always attracts a large

**Table 1** Notations and explanations

Notation	Explanation
$u_i^A$	Ranking score of developer $i$ in layer $A$
$p_j^A$	Ranking score of project $j$ in layer $A$
$W^A$	Weight matrix of layer $A$ , its element $w_{ij}^A$ represents the edge weight between developer $i$ and project $j$ in layer $A$
$\gamma, \lambda$	Hyperparameters
$u_i^0, p_j^0$	Query vectors for developer $i$ and project $j$
$\theta$	Threshold for SIR simulation
$ U^A ,  P^A $	Numbers of developers and projects in layer $A$
$d_i^A, d_j^A$	Degrees of developer $i$ and project $j$ in layer $A$
$S^A$	Symmetric normalization of weight matrix $W^A$

number of attention. Taking TensorFlow as an example, it got thousands of stars quickly upon its first release in GitHub because it is supported by Google.

This intuition forms our first assumption that a developer (project) should be ranked high if it is connected to high-ranked projects (developers) in a certain layer  $A$ , which can be formulized as follows:

$$p_j^A = \sum_{i=1}^{|U|} w_{ij}^A u_i^A \quad (1)$$

$$u_i^A = \sum_{j=1}^{|P|} w_{ij}^A p_j^A \quad (2)$$

In order to employ a prior belief on nodes' importance and provide better ranking results, we also adopt query vector and symmetric normalization as BiRank. The prior belief on nodes' importance and rankings from network structure are balanced with two parameters  $\gamma$  and  $\lambda$ . The final formulation of the mutual reinforcement between developers and projects is as follows:

$$p_j^A = \gamma \sum_{i=1}^{|U|} \frac{w_{ij}^A}{\sqrt{d_i^A} \sqrt{d_j^A}} u_i^A + (1 - \gamma) p_j^0 \quad (3)$$

$$u_i^A = \lambda \sum_{j=1}^{|P|} \frac{w_{ij}^A}{\sqrt{d_i^A} \sqrt{d_j^A}} p_j^A + (1 - \lambda) u_i^0 \quad (4)$$

### 3.3 Mutual reinforcement between different layers

Besides considering the mutual reinforcement between developers and projects in each single layer, we also take into account the mutual reinforcement between different layers. From our experience as open source software developers, we could firmly assume different interactions between developers and projects reflect different aspects of influence and only a composition of all the aspects could reflect a comprehensive influence of developers and projects. For example, committing code to a project indicates a developer's coding skill and commenting issues of a project may show a developer's design skill or bug-fixing skill. The influence of a developer should be measured by summarization of both coding and design skills.

To implement mutual reinforcement between different layers, we choose to incorporate the ranking scores of developers and projects from the first layer as an enhancement of the query vectors of the second layer. The mathematical formulations are given in Eqs. (5) and (6).

$$p_j^B = \gamma \sum_{i=1}^{|U|} \frac{w_{ij}^B}{\sqrt{d_i^B} \sqrt{d_j^B}} u_i^B + (1 - \gamma) p_j^A \quad (5)$$

$$u_i^B = \lambda \sum_{j=1}^{|P|} \frac{w_{ij}^B}{\sqrt{d_i^B} \sqrt{d_j^B}} p_j^B + (1 - \lambda) u_i^A \quad (6)$$

To be more clearer, we transform Eqs. (5) and (6) to their equivalent matrix form in Eqs. (7) and (8):

$$p^B = \gamma (S^B)^T u^B + (1 - \gamma) p^A \quad (7)$$

$$u^B = \lambda (S^B) p^B + (1 - \lambda) u^A \quad (8)$$

where  $S^B = (D_u^B)^{-\frac{1}{2}} W^B (D_p^B)^{-\frac{1}{2}}$  is the symmetric normalization of weight matrix  $W^B$ .

### 3.4 Overall algorithm

By combing both the mutual reinforcement between developers and projects in each single layer and that between different layers, we finally propose the M-BiRank method to co-rank developers and projects in open source software community and the overall algorithm is shown in Algorithm 1.

---

#### Algorithm 1 M-BiRank Algorithm

---

**Input:**

Weight matrix  $W^A$  and  $W^B$ , query vectors  $u^0, p^0$  and hyper-parameters  $\gamma, \lambda$ ;

**Output:**

Ranking vectors  $u, p$ ;

- 1: Symmetrically normalize  $W^A$  :  $S^A = (D_u^A)^{-\frac{1}{2}} W^A (D_p^A)^{-\frac{1}{2}}$ ;
  - 2: Randomly initialize  $p^A$  and  $u^A$ ;
  - 3: **while** Stopping criteria is not met **do**
  - 4:    $p^A \leftarrow \gamma (S^A)^T u^A + (1 - \gamma) p^0$ ;
  - 5:    $u^A \leftarrow \lambda S^A p^A + (1 - \lambda) u^0$ ;
  - 6: **end while**
  - 7: Symmetrically normalize  $W^B$  :  $S^B = (D_u^B)^{-\frac{1}{2}} W^B (D_p^B)^{-\frac{1}{2}}$ ;
  - 8: Initialize  $p^B$  and  $u^B$  using  $p^A$  and  $u^A$ , respectively;
  - 9: **while** Stopping criteria is not met **do**
  - 10:    $p^B \leftarrow \gamma (S^B)^T u^B + (1 - \gamma) p^A$ ;
  - 11:    $u^B \leftarrow \lambda S^B p^B + (1 - \lambda) u^A$ ;
  - 12: **end while**
  - 13:  $p \leftarrow p^B, u \leftarrow u^B$ ;
  - 14: **return**  $u$  and  $p$ .
- 

### 3.5 Time complexity analysis

The overall time complexity of M-BiRank is a summarization of each layer's time complexity. For each layer, according to Eqs. (7) and (8), the time cost mainly depends on the multiplication of  $(S^B)^T u^B$  and  $(S^B) p^B$ , showing a time complexity of  $O(|U| \cdot |P|)$ . However, most real-world networks are usually very sparse and only non-zero elements (which correspond to existing edges) should be stored and computed regarding matrix multiplication of  $(S^B)^T u^B$  and  $(S^B) p^B$ . Thus, the time complexity of layer  $A$  is  $O(c_A |E_A|)$  and the overall time complexity of M-BiRank is  $O(\sum_{A=1}^M c_A |E_A|)$ , where  $c_A$  denotes the number of iterations in layer  $A$  and  $|E_A|$  denotes the number of edges in layer  $A$ .

## 4 Experiment

In this section, the performance of M-BiRank model is evaluated against the GHTorrent dataset [36].

### 4.1 Datasets

GHTorrent [36] monitors the GitHub public event timeline and retrieves information of developers, projects, and interaction details between them from these events [37]. We choose the GHTorrent dataset as of November 1, 2018, and extract the relationships between developers and projects, both of which mainly belong to PHP community. The steps of data preprocess include the following: (1) choose issues and commits which belong to PHP projects, (2) keep developers and projects which exist in both issues and commits, and (3) construct multiplex developer-project bipartite network, using a function of issue/commit number as edge weight, that is,  $w_{ij}^{\text{issue}} = \log(\#\text{issues}) + 0.3$  or  $w_{ij}^{\text{commit}} = \log(\#\text{commits}) + 0.3$ . The detailed information about the dataset used in our experiment is shown in Table 2.

### 4.2 Evaluation metrics

In order to evaluate and compare the performance of M-BiRank and baseline methods, both correlation analysis and SIR model are adopted.

Correlation analysis mainly focuses on comparing predictions against the ground truth, and Pearson's correlation coefficient (PCC) [38] is chosen. In our experiment, the number of watch of projects and the number of followers of developers are set as the ground truth for the rankings of projects and developers, respectively. PCC reflects the correlation degree of two variables through the linear correlation between vectors, which is defined as follows:

$$\text{PCC}(x, y) = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}} \quad (9)$$

where  $n$  represents the number of elements;  $x_i$  and  $y_i$  represent the  $i$ th element of sample  $x$  and  $y$ , respectively; and the value range of PCC is  $[-1, 1]$ .

However, the ground truth in correlation analysis is some kinds of degree in networks, which is a rough metric in evaluating the influence of developers or projects. To rank more precisely, dynamic models are needed for simulating the influence diffusion process [39]. SIR model [40] is a classical epidemic model and is often used to evaluate the ability of information spreading of a node in social networks. Generally, an influential user with a higher ranking score will spread his/her opinions to more developers. The transmission process of SIR model is shown in Fig. 2, where S (Susceptible), I (Infected), and R (Removed) denote the susceptible, infected, and recovered nodes. At the initial step of the

**Table 2** The statistics of the dataset

Data	Count
Developers	147,105
Projects	126,415
Developer-project commit relations	205,478
Developer-project issue relations	441,558
Developer-project watch relations	68,566
Developer-developer follow relations	10,701



transmission process, several infected nodes are set, and then, the transmission is iteratively repeated until no new nodes are infected [41]. At each step, infected nodes infect its susceptible neighbors with the probability  $\alpha$ , and infected nodes' recovery to removed status with the probability  $\beta$ . So SIR model is suitable for evaluating the ability of information spreading of a node. By applying nodes with highest ranking scores of different ranking methods as the initial infected nodes and comparing the final number of affected nodes (both infected and removed nodes), the effectiveness of different ranking methods can be compared.

### 4.3 Baseline methods

We compare M-BiRank with several baseline methods:

*Degree* [42]. The degrees of developers and projects in different layers of multiplex developer-project bipartite network are calculated and *averaged*.

*PageRank* [9]. PageRank ranks nodes by iteratively propagating scores on the network and is usually suitable for single layer monopartite network. In this experiment, we apply it to multiplex developer-project bipartite network with two different setups. *PageRank-Avg* ignores types of nodes and applies PageRank algorithm directly to different layers of the multiplex developer-project bipartite network. The final ranking score of a node is the average of different layers. *PageRank-Add* merges different layers of multiplex developer-project bipartite network into a single layer of developer-project bipartite network and uses the average edge weights of different layers as edge weights of this single layer bipartite network. Then, we apply PageRank algorithm to this single layer of developer-project bipartite network ignoring types of nodes. Finally, both *PageRank-Avg* and *PageRank-Add* rank developers and projects separately according their final ranking scores. The hyperparameter is set to 0.85.

*BiRank* [11]. BiRank is a propagation-based ranking method on bipartite networks and adopts a normalization strategy in the iterative process. *BiRank-Avg* applies BiRank algorithm to different layers of multiplex developer-project bipartite network separately and averages the ranking scores in different layers as the final ranking scores. *BiRank-Add* firstly merges different layers of multiplex developer-project bipartite network into a single layer of developer-project bipartite network with the average of the edge weights in different layers as edge weights. Both of the hyperparameters are set to 0.85.

*Multiplex PageRank* [13]. Multiplex PageRank considers the impact of the centrality of a node in one layer on that in another layer and introduces nodes' centrality of the preceding layer to current layer in four ways. In this experiment, we choose the Additive Multiplex PageRank and have two different setups, that is, *MPR-Commit* uses the commit layer as the first layer and *MPR-Issue* uses the issue layer as the first layer. The hyperparameter is set to 0.85.

*M-BiRank*. M-BiRank is the method we proposed for ranking nodes in multiplex bipartite network. As the setup in *Multiplex PageRank*, *M-BiRank-Commit* uses the commit

layer as the first layer and *M-BiRank-Issue* uses the issue layer as the first layer. The hyperparameters  $\gamma$  and  $\lambda$  are set to 0.85. Each element of the query vector  $u^0$  ( $p^0$ ) for the corresponding node (developer/project) is set to the sum of its all edges' weights over the total sum of all edges' weights of the whole developer-project bipartite network of the first layer.

#### 4.4 Results

We compare the experimental results of M-BiRank with baseline methods by both correlation analysis and SIR modeling. The hyperparameters for M-BiRank  $\gamma$  and  $\lambda$  are both set to 0.85.

##### 4.4.1 Correlation analysis

In correlation analysis, the follower number of developers and the watch number of projects are set as the ground truth for ranking developers and projects, respectively. Pearson's correlation coefficient (PCC) is calculated between the ranking results from M-BiRank and baseline methods and ground truth rankings. The results are shown in Table 3.

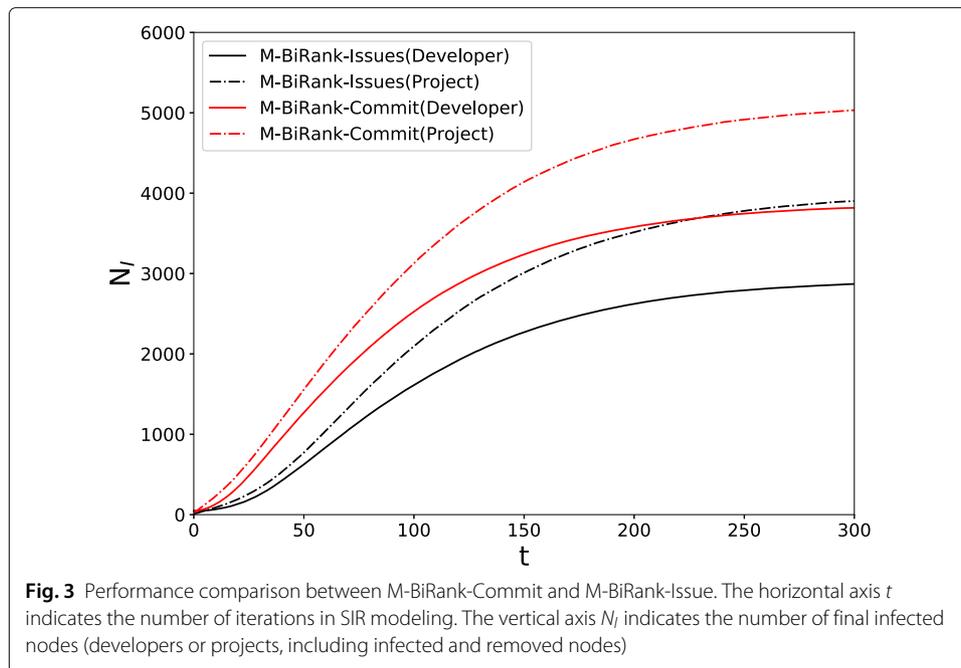
From the results of correlation analysis, we have the following observations:

(1) M-BiRank model we proposed outperforms all the baseline methods for both developer ranking and project ranking. This indicates that it is necessary to model multiple interactions between developers and projects as a multiplex bipartite network, which not only considers mutual enhancement between developers and projects but also takes into account mutual enhancement between different interactions. This highly agrees with real-world practice. For example, a project with elite developers participating in is usually a popular project and a developer participating in popular projects is often an elite developer. Developers have different ways to take part in certain projects such as committing code or solving issues, and different ways are tightly coupled.

(2) Comparing the different settings of M-BiRank itself, M-BiRank-Commit performs better than M-BiRank-Issue in most cases, which means it is better to take the commit layer of the multiplex developer-project bipartite network as the initial layer for M-BiRank model. This also agrees with real-world practice. Issue is a helper function in social collaborative coding which provides a discussion board for software developers about bugs and designs. While commit is a main function during software development for developers, thus, the commit layer is more important. So M-BiRank-Commit performs better in

**Table 3** Correlation analysis

Method	Top 20		Top 50		Top 100	
	Project	Developer	Project	Developer	Project	Developer
Degree	0.441	0.090	0.383	0.173	0.393	0.219
BiRank-Avg	0.402	-0.014	0.442	0.091	0.462	0.016
BiRank-Add	0.398	-0.022	0.422	-0.062	0.487	0.030
PageRank-Avg	0.430	-0.116	0.450	-0.027	0.498	-0.035
PageRank-Add	0.357	-0.111	0.419	-0.067	0.475	-0.004
MPR-Issue	0.350	0.113	0.437	0.160	0.458	0.202
MPR-Commit	0.464	-0.066	0.459	-0.013	0.516	0.017
<b>M-BiRank-Issue</b>	0.374	0.057	0.423	0.068	0.447	<b>0.273</b>
<b>M-BiRank-Commit</b>	<b>0.468</b>	<b>0.152</b>	<b>0.459</b>	<b>0.261</b>	<b>0.541</b>	0.143



identifying more influential developers and projects. In Section 4.4.2, we only compare M-BiRank-Commit with benchmark methods.

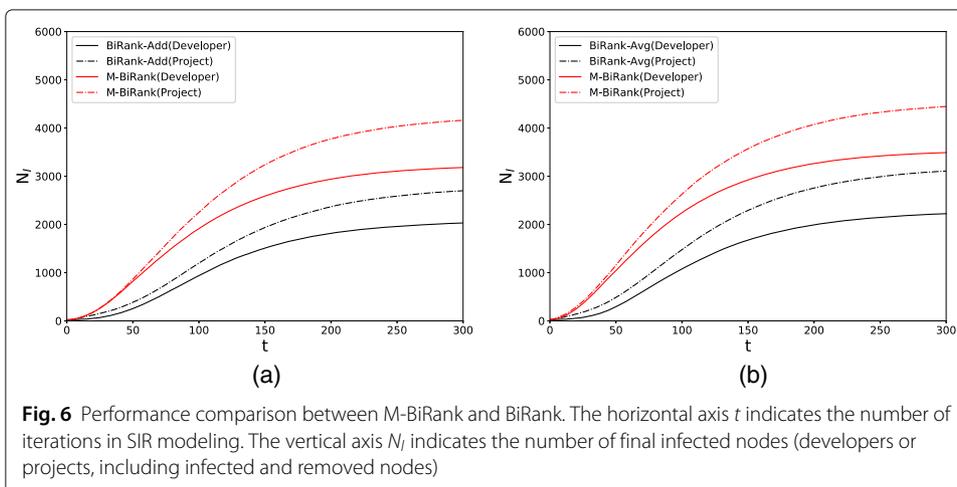
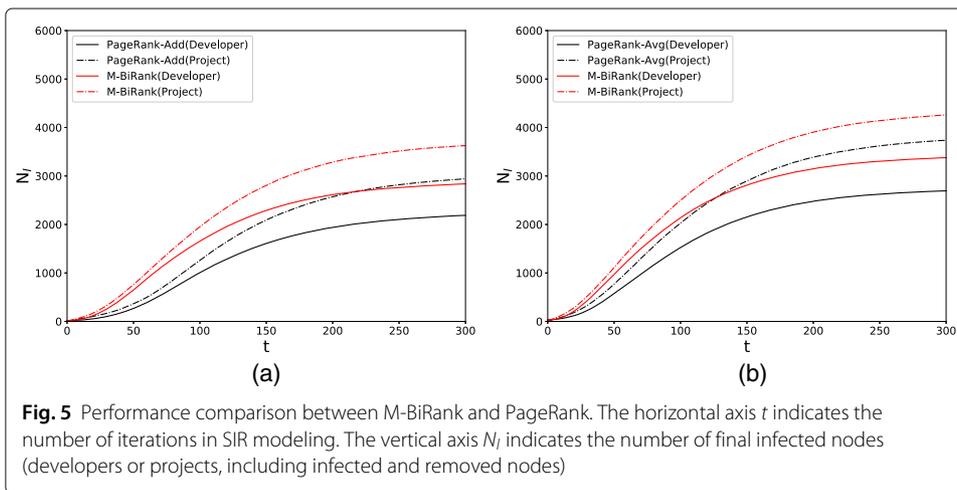
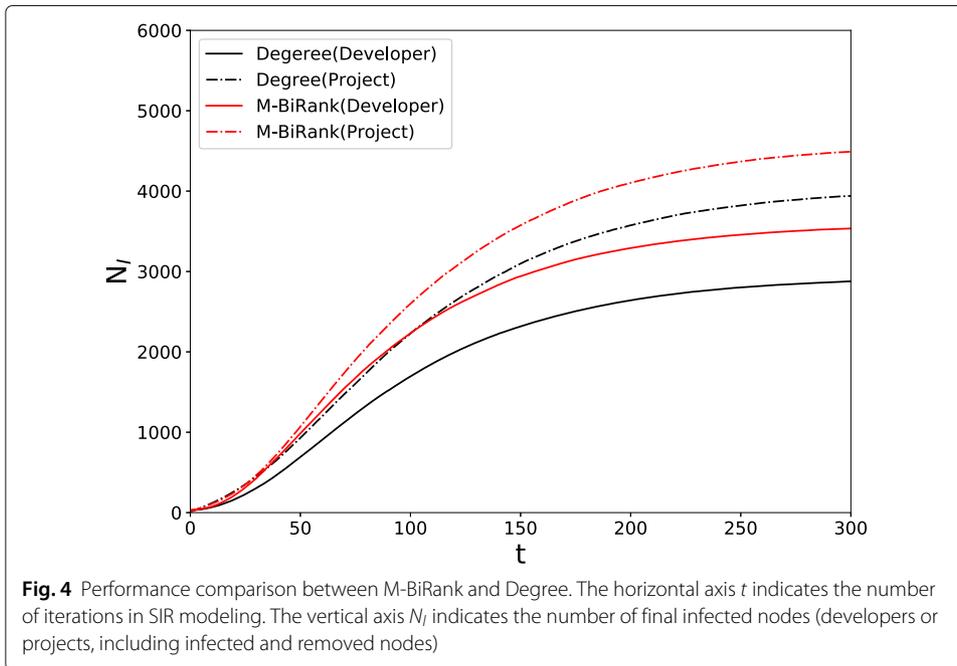
#### 4.4.2 SIR simulation

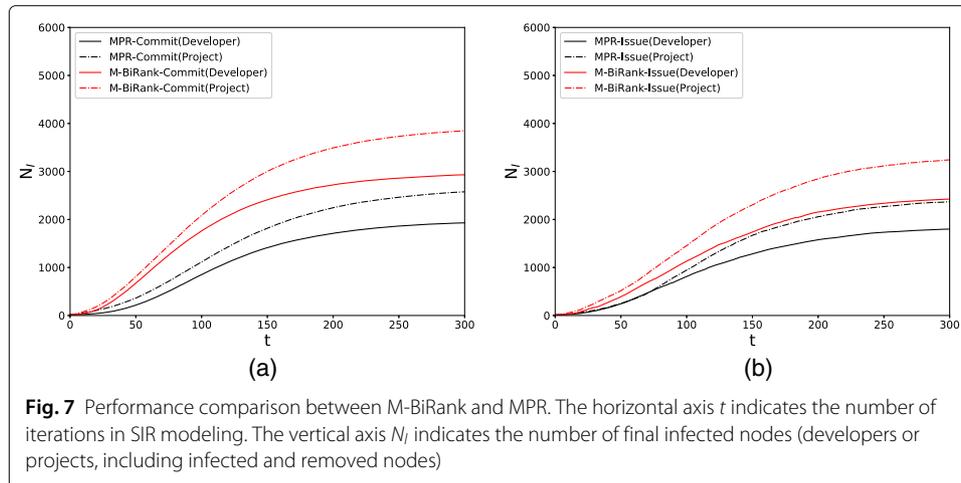
In this section, to evaluate the information spreading ability of top 100 developers ranked by different methods, SIR model is adopted on commit layer of the developer-project multiplex bipartite network. M-BiRank is compared against each baseline method separately. For each comparison, the initial infected nodes (developers) for SIR model are the top 100 developers ranked by each method excluding those ranked top 100 by both methods. During the SIR process, an infected node infects each of its neighbors with probability  $\alpha = 0.005$  simultaneously and recovers to removed state with probability  $\beta = 0.006$ . For each SIR simulation, we run 300 iterations at most and repeat 10 times to average the value of each step. The results are shown in Figs. 3, 4, 5, 6, and 7, and several significant observations are found:

(1) The result of comparison between different settings of M-BiRank itself in Fig. 3 indicates M-BiRank-Commit performs better, which is in perfect accordance with the result found in correlation analysis in Section 4.4.1. Thus, only M-BiRank-Commit is compared against baseline methods in the rest part of this section.

(2) M-BiRank outperforms all the baseline methods in identifying influential developers, which means nodes' types and mutual reinforcement among different interactions play important roles and multiplex bipartite network can model multiple interactions between two different types of nodes more precisely. Specially, the performance difference between M-BiRank and BiRank is larger than that between M-BiRank and Multiplex PageRank (MPR), from which we can conclude that considering mutual reinforcement among different interactions is of more importance than distinguishing nodes' types.

(3) The number of final infected projects is more than that of developers in both M-BiRank and all the baseline methods. According to researches on epidemics on networks,





information spreads faster and broader in networks with shorter average path length. From Table 2, we can see the average degree of projects is larger than that of developers.

#### 4.5 Case study

In addition to correlation analysis and SIR simulation, we further do a detailed case study to show the effectiveness of our model in identifying influential developers and projects. The top 20 developers and projects ranked by our model M-BiRank are listed in Tables 4 and 5, respectively, followed by their ranks in baseline methods.

Table 4 indicates baseline methods, and M-BiRank ranks the first six developers similarly, while some influential PHP developers ranked in top 20 by M-BiRank are not identified or ranked with lower scores by baseline methods. For example, Fabien Potencier (GitHub ID: fabpot) and Taylor Otwell (GitHub ID: taylorotwell), the most active contributors of the two most popular PHP frameworks, Symfony and Laravel, are not identified as influential developers by some of the baseline methods. From GitHub as of June 1, 2020, Symfony and Laravel have 23.3k and 59.4k stars, respectively, and Fabien Potencier and Taylor Otwell have 10.4k and 18.6k followers, respectively. Taylor Otwell has more followers than Fabien Potencier, and Laravel is more popular than Symfony, but Fabien Potencier is ranked higher than Taylor Otwell because Laravel is based on some popular components of Symfony. Thus, we can conclude that Fabien Potencier is more influential than Taylor Otwell.

As for projects, from Table 5, we can see both M-BiRank and baseline methods rank popular PHP frameworks with higher scores. But some important PHP components identified by M-BiRank are not identified as influential projects or ranked with lower scores by baseline methods. For example, *illuminate/database*, a popular ORM library, is ranked with a high score by M-BiRank but is ranked with a lower score by BiRank-Add and PageRank-Add, and is never identified as influential projects by BiRank-Avg, PageRank-Avg, and MPR-Commit. As we know, in modern web development, ORM is quite critical because it is responsible for accessing database.

#### 4.6 Experimental settings discussion

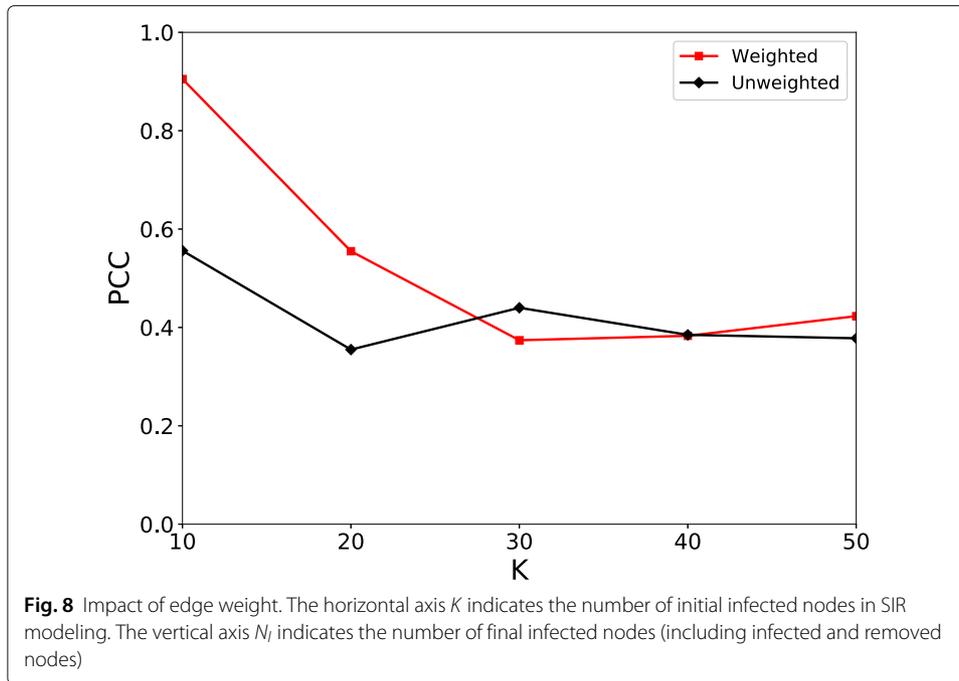
In the experiment, several key settings will affect the performance of M-BiRank and a brief discussion about these settings is shown as follows.

**Table 4** Top 20 developers

Developer	Ranking									
	M-BiRank	Degree	BiRank-Avg	BiRank-Add	PageRank-Avg	PageRank-Add	MPR-Issue	MPR-Commit		
scrutinizer-auto-fixer	1	1	10	2	5	2	1	67		
GrahamCampbell	2	3	2	3	2	3	2	7		
danielbachhuber	3	7	3	6	4	6	9	3		
weierophinney	4	2	13	4	12	4	4	42		
stof	5	12	5	7	6	7	6	4		
freeknuurze	6	5	19	5	21	5	5	83		
Nyholm	7	8	9	11	11	12	10	15		
kartik-v	8	6	35	20	43	20	7	95		
Ocramius	9	7	12	12	13	10	11	19		
cryobone	10	16	8	9	10	9	12	11		
fabpot	11	20	81	34	55	23	8	-		
clue	12	18	7	13	8	13	17	8		
sagikazarmark	13	23	6	10	7	11	15	5		
lsmith77	14	37	16	17	16	15	16	14		
dbu	15	9	24	29	20	18	20	17		
cordoval	16	19	4	8	3	8	19	2		
pippinsplugins	17	14	22	26	24	29	26	25		
taylorotwell	18	43	-	30	66	25	13	-		
dantleech	19	27	17	23	17	21	21	16		
splitbrain	20	13	21	14	26	14	29	47		

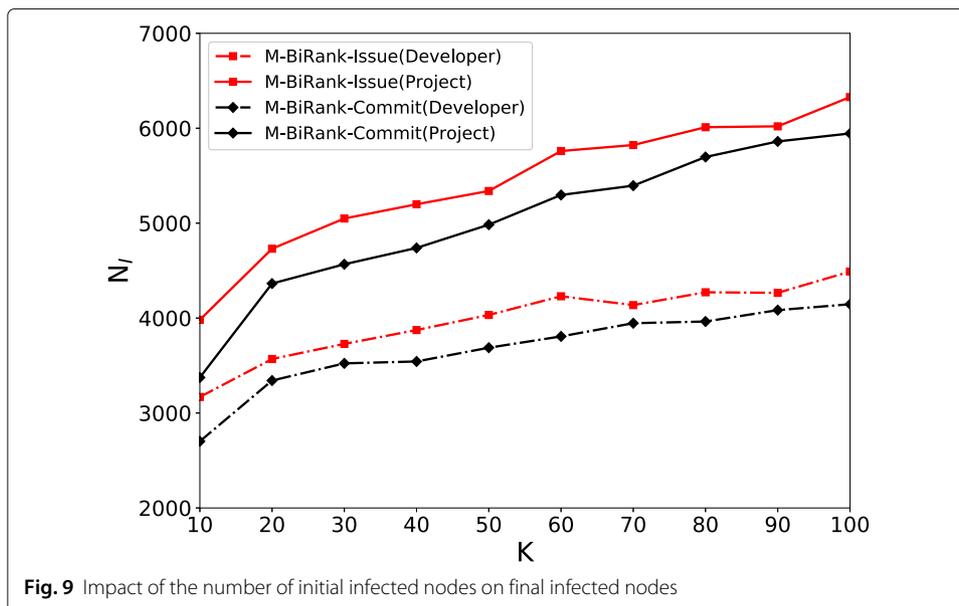
**Table 5** Top 20 projects

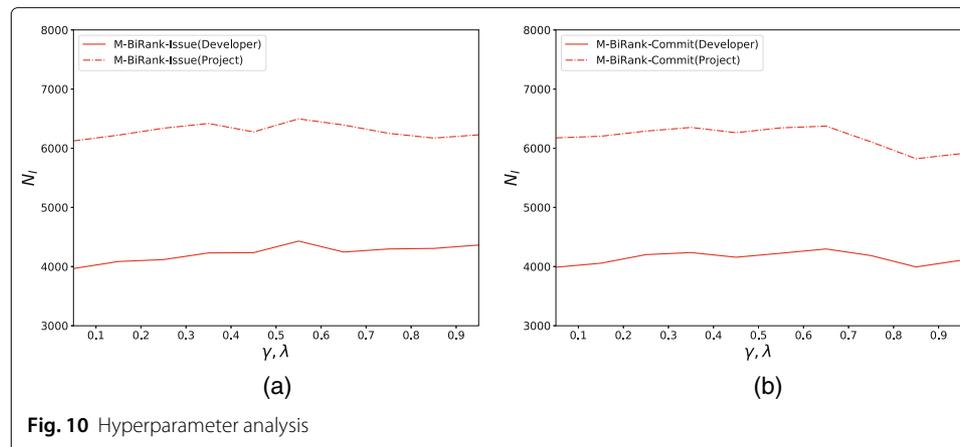
Project	Ranking									
	M-BiRank	Degree	BiRank-Avg	BiRank-Add	PageRank-Avg	PageRank-Add	MPR-Issue	MPR-Commit		
symfony/symfony	1	1	2	2	2	2	1	2		
laravel/framework	2	3	1	1	1	1	2	1		
magento/magento2	2	7	4	4	4	4	3	4		
yiiisoft/yii2	4	2	3	3	3	3	4	3		
owncloud/core	5	12	5	5	5	5	5	6		
joomla/joomla-cms	6	5	6	7	7	7	6	7		
cakephp/cakephp	7	8	8	9	9	9	7	9		
Automattic/jetpack	8	60	13	16	18	19	11	19		
symfony/symfony-docs	9	4	14	21	19	20	9	20		
EllisLab/CodeIgniter	10	6	10	11	11	11	10	12		
pyrocms/pyrocms	11	12	12	12	13	12	14	14		
composer/composer	12	22	7	6	6	6	8	5		
woocommerce/woocommerce	13	18	9	8	8	8	12	8		
phalcon/cphalcon	14	17	11	10	10	10	13	10		
illuminate/database	15	16	-	53	-	55	15	-		
sonata-project/SonataAdminBundle	16	32	18	20	15	18	17	16		
Sylus/Sylus	17	23	16	22	16	22	16	15		
yiiisoft/yii	18	19	20	18	21	17	18	21		
PrestaShop/PrestaShop	19	14	31	32	41	34	24	54		
laravel/laravel	20	35	19	19	22	21	22	18		



First, we will study the impact of edge weight in the ranking process. Both unweighted and weighted developer-project multiplex bipartite networks are constructed, and for weighted case, the interaction times are summed as edge weight. Then, correlation analysis on top  $k$  developers and projects is applied and the results are shown in Fig. 8, from which it can be concluded that weighted developer-project multiplex bipartite network performs better than unweighted case and edge weight plays an important role in identifying influential developers and projects.

Then, experimental settings for SIR simulation are discussed. It can be seen from Fig. 9 that the more initial infected nodes are set, the more final infected nodes. It is also obvious





**Fig. 10** Hyperparameter analysis

that the same number of top  $k$  projects being set as initial infected nodes will result in more final infected nodes.

Finally, the hyperparameters  $\gamma$  and  $\lambda$  of our proposed M-BiRank model are analyzed. For simplicity, we consider the condition that  $\gamma$  and  $\lambda$  are equal. It can be concluded from Fig. 10 that both prior belief of developers' (projects') importance and rankings from network structure play roles in the final rankings of developers (projects) and their contributions to final rankings are approximately equal.

## 5 Conclusions

In this work, we study the problem of identifying influential developers and projects in open source software community. We model multiple interactions between developers and projects as a multiplex bipartite network and propose an iterative refinement ranking method M-BiRank by incorporating the mutual reinforcement between developers and projects as well as between multiple developer-project interactions. The proposed M-BiRank is evaluated against four baseline methods on real-world GitHub dataset. Extensive experimental analysis and case study show M-BiRank significantly outperforms baseline methods in both correlation analysis and SIR simulation.

The general idea behind the proposed M-BiRank is modeling multiple kinds of entities and interactions in open source software community into a single network and incorporating mutual reinforcement between different kinds of entities as well as between different types of interactions when ranking. As we know, there are other entities such as blogs and organizations in addition to developers and projects in open source software community and plenty of interactions between them such as user-user following and project-project dependency. In future work, more entities and interactions could be introduced and modeled as a heterogeneous information network and mutual reinforcement in ranking would be generalized using meta-path.

### Abbreviations

M-BiRank: Multiplex bipartite ranking; PCC: Pearson's correlation coefficient; MPR: Multiplex PageRank

### Acknowledgements

None.

### Authors' contributions

Dengcheng Yan and Bin Qi designed the model. Bin Qi conducted the experiment and prepared the figures. Dengcheng Yan wrote the entire article. All authors read and approved the final manuscript.

**Funding**

This work is supported by the National Natural Science Foundation of China (Grant No. 61872002), the University Natural Science Research Project of Anhui Province (Grant No. KJ2019A0037), the University Collaborative Innovation Program of Anhui Province (Grant No. GXXT-2019-013), and the Doctoral Scientific Research Foundation of Anhui University (Grant No. Y040418194).

**Availability of data and materials**

The datasets used and/or analyzed during the current study are available from the corresponding author on reasonable request.

**Competing interests**

The authors declare that they have no competing interests.

**Author details**

<sup>1</sup>School of Computer Science and Technology, Anhui University, Jiulong Road No. 111, Hefei, 230601, China. <sup>2</sup>Institutes of Physical Science and Information Technology, Anhui University, Jiulong Road No. 111, Hefei, 230601, China.

Received: 10 June 2020 Accepted: 29 September 2020

Published online: 27 October 2020

**References**

1. Y. Chen, N. Zhang, Y. Zhang, X. Chen, W. Wu, X. S. Shen, Toffee: task offloading and frequency scaling for energy efficiency of mobile devices in mobile edge computing. *IEEE Trans. on Cloud Comput.* (2019). <https://doi.org/10.1109/TCC.2019.2923692>
2. Y. Zhang, J. Pan, L. Qi, Q. He, Privacy-preserving quality prediction for edge-based iot services. *Future Gener. Comput. Syst.* **114**, 336–348 (2021). <https://doi.org/10.1016/j.future.2020.08.014>
3. Y. Chen, N. Zhang, Y. Zhang, X. Chen, Dynamic computation offloading in edge computing for internet of things. *IEEE Internet Things J.* **6**(3), 4242–4251 (2019). <https://doi.org/10.1109/JIOT.2018.2875715>
4. Y. Chen, Y. Zhang, Y. Wu, L. Qi, X. Chen, X. Shen, Joint task scheduling and energy management for heterogeneous mobile edge computing with hybrid energy supply. *IEEE Internet Things J.* **7**(9), 8419–8429 (2020). <https://doi.org/10.1109/JIOT.2020.2992522>
5. M. Joblin, S. Apel, C. Hunsen, W. Mauerer, in *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*, Classifying developers into core and peripheral: an empirical study on count and network metrics (IEEE, Buenos Aires, 2017), pp. 164–174. <https://doi.org/10.1109/ICSE.2017.23>
6. D.-C. Yan, Z.-W. Wei, X.-P. Han, B.-H. Wang, Empirical analysis on the human dynamics of blogging behavior on github. *Physica A.* **465**, 775–781 (2017). <https://doi.org/10.1016/j.physa.2016.08.054>
7. Y. Hu, S. Wang, Y. Ren, K.-K. R. Choo, User influence analysis for github developer social networks. *Expert Syst. Appl.* **108**, 108–118 (2018). <https://doi.org/10.1016/j.eswa.2018.05.002>
8. Q. Xuan, C. Fu, L. Yu, Ranking developer candidates by social links. *Adv. Complex Syst.* **17**, 1550005 (2014). <https://doi.org/10.1142/S0219525915500058>
9. L. Page, S. Brin, R. Motwani, T. Winograd, The pagerank citation ranking: bringing order to the web. Technical report (1999). Stanford InfoLab
10. J. M. Kleinberg, Authoritative sources in a hyperlinked environment. *J. ACM.* **46**(5), 604–632 (1999). <https://doi.org/10.1145/324133.324140>
11. X. He, M. Gao, M.-Y. Kan, D. Wang, Birank: towards ranking on bipartite graphs. *IEEE Trans. Knowl. Data Eng.* **29**(1), 57–71 (2016). <https://doi.org/10.1109/TKDE.2016.2611584>
12. H. Liu, H. Kou, C. Yan, L. Qi, Keywords-driven and popularity-aware paper recommendation based on undirected paper citation graph. *Complexity* (2020). <https://doi.org/10.1155/2020/2085638>
13. A. Halu, R. J. Mondragón, P. Panzarasa, G. Bianconi, Multiplex pagerank. *PLoS One.* **8**(10), 78293 (2013). <https://doi.org/10.1371/journal.pone.0078293>
14. J. Iacovacci, C. Rahmede, A. Arenas, G. Bianconi, Functional multiplex pagerank. *EPL.* **116**(2), 28004 (2016). <https://doi.org/10.1209/0295-5075/116/28004>
15. Y. Zhang, Y. Zhou, F. Wang, Z. Sun, Q. He, Service recommendation based on quotient space granularity analysis and covering algorithm on spark. *Knowledge-Based Syst.* **147**, 25–35 (2018). <https://doi.org/10.1016/j.knosys.2018.02.014>
16. Y. Zhang, G. Cui, S. Deng, F. Chen, Y. Wang, Q. He, Efficient query of quality correlation for service composition. *IEEE Trans. Serv. Comput.* (2018). <https://doi.org/10.1109/TSC.2018.2830773>
17. Y. Zhang, K. Wang, Q. He, F. Chen, S. Deng, Z. Zheng, Y. Yang, Covering-based web service quality prediction via neighborhood-aware matrix factorization. *IEEE Trans. Serv. Comput.* (2019). <https://doi.org/10.1109/TSC.2019.2891517>
18. Y. Zhang, C. Yin, Q. Wu, Q. He, H. Zhu, Location-aware deep collaborative filtering for service recommendation. *IEEE Trans. Syst., Man, Cybern. Syst.* (2019). <https://doi.org/10.1109/TSMC.2019.2931723>
19. L. Qi, X. Wang, X. Xu, W. Dou, S. Li, Privacy-aware cross-platform service recommendation based on enhanced locality-sensitive hashing. *IEEE Transactions on Network Science and Engineering* (2020). <https://doi.org/10.1109/TNSE.2020.2969489>
20. C. Zhou, A. Li, A. Hou, Z. Zhiwang, Z. Zhang, P. Dai, F. Wang, Modeling methodology for early warning of chronic heart failure based on real medical big data. *Expert Syst. Appl.* **151**, 113361 (2020). <https://doi.org/10.1016/j.eswa.2020.113361>
21. D. Chen, L. Lü, M.-S. Shang, Y.-C. Zhang, T. Zhou, Identifying influential nodes in complex networks. *Physica A.* **391**(4), 1777–1787 (2012). <https://doi.org/10.1016/j.physa.2011.09.017>
22. D.-B. Chen, H. Gao, L. Lü, T. Zhou, Identifying influential nodes in large-scale directed networks: the role of clustering. *PLoS One.* **8**(10) (2013). <https://doi.org/10.1371/journal.pone.0077455>

23. H. Liu, H. Kou, C. Yan, L. Qi, Link prediction in paper citation network to construct paper correlation graph. *EURASIP J. Wireless Com. Network.* **2019**(1) (2019). <https://doi.org/10.1186/s13638-019-1561-7>
24. S. X. Zhao, R. Rousseau, Y. Y. Fred, h-degree as a basic measure in weighted networks. *J. Informetr.* **5**(4), 668–677 (2011). <https://doi.org/10.1016/j.joi.2011.06.005>
25. Q. Liu, Y.-X. Zhu, Y. Jia, L. Deng, B. Zhou, J.-X. Zhu, P. Zou, Leveraging local h-index to identify and rank influential spreaders in networks. *Physica A.* **512**, 379–391 (2018). <https://doi.org/10.1016/j.physa.2018.08.053>
26. L. Lü, T. Zhou, Q.-M. Zhang, H. E. Stanley, The h-index of a network node and its relation to degree and coreness. *Nat. Commun.* **7**, 10168 (2016). <https://doi.org/10.1038/ncomms10168>
27. G. Sabidussi, The centrality index of a graph. *Psychometrika.* **31**(4), 581–603 (1966). <https://doi.org/10.1007/BF02289527>
28. L. C. Freeman, A set of measures of centrality based on betweenness. *Sociometry.* **40**(1), 35–41 (1977). <https://doi.org/10.2307/3033543>
29. K. Stephenson, M. Zelen, Rethinking centrality: methods and examples. *Soc. Networks.* **11**(1), 1–37 (1989). [https://doi.org/10.1016/0378-8733\(89\)90016-6](https://doi.org/10.1016/0378-8733(89)90016-6)
30. D. Wang, H. Wang, X. Zou, Identifying key nodes in multilayer networks based on tensor decomposition. *Chaos.* **27**(6), 063108 (2017). <https://doi.org/10.1063/1.4985185>
31. D. Wang, X. Zou, A new centrality measure of nodes in multilayer networks under the framework of tensor computation. *Appl. Math. Model.* **54**, 46–63 (2018). <https://doi.org/10.1016/j.apm.2017.07.012>
32. L. Lü, Y.-C. Zhang, C. H. Yeung, T. Zhou, Leaders in social networks, the delicious case. *PLoS One.* **6**(6) (2011). <https://doi.org/10.1371/journal.pone.0021202>
33. W. Zhong, X. Yin, X. Zhang, S. Li, W. Dou, R. Wang, L. Qi, Multi-dimensional quality-driven service recommendation with privacy-preservation in mobile edge environment. *Comput. Commun.* **157**, 116–123 (2020). <https://doi.org/10.1016/j.comcom.2020.04.018>
34. R. Liang, X. Jiang, in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, Scientific ranking over heterogeneous academic hypernetwork (AAAI Press, Phoenix, Arizona, 2016), pp. 20–26
35. H. Zhao, X. Xu, Y. Song, D. L. Lee, Z. Chen, H. Gao, Ranking users in social networks with motif-based pagerank. *IEEE Trans. Knowl. Data Eng.* (2019). <https://doi.org/10.1109/TKDE.2019.2953264>
36. G. Gousios, in *Proceedings of the 10th Working Conference on Mining Software Repositories*, The ghtorrent dataset and tool suite (IEEE, San Francisco, 2013), pp. 233–236. <https://dl.acm.org/doi/10.5555/2487085.2487132>
37. X. Xu, S. Fu, L. Qi, X. Zhang, Q. Liu, Q. He, S. Li, An IoT-oriented data placement method with privacy preservation in cloud environment. *J. Netw. Comput. Appl.* **124**, 148–157 (2018). <https://doi.org/10.1016/j.jnca.2018.09.006>
38. J. Hauke, T. Kossowski, Comparison of values of Pearson's and Spearman's correlation coefficients on the same sets of data. *Quaestiones geographicae.* **30**(2), 87–93 (2011). <https://doi.org/10.2478/v10117-011-0021-1>
39. J. Li, T. Cai, K. Deng, X. Wang, T. Sellis, F. Xia, Community-diversified influence maximization in social networks. *Inf. Syst.* **92**, 101522 (2020). <https://doi.org/10.1016/j.is.2020.101522>
40. R. Yang, B.-H. Wang, J. Ren, W.-J. Bai, Z.-W. Shi, W.-X. Wang, T. Zhou, Epidemic spreading on heterogeneous networks with identical infectivity. *Phys. Lett. A.* **364**(3–4), 189–193 (2007). <https://doi.org/10.1016/j.physleta.2006.12.021>
41. X. Xu, R. Mo, F. Dai, W. Lin, S. Wan, W. Dou, Dynamic resource provisioning with fault tolerance for data-intensive meteorological workflows in cloud. *IEEE Trans. Ind. Inform.* **16**(9), 6172–6181 (2019). <https://doi.org/10.1109/TII.2019.2959258>
42. P. Bonacich, Factoring and weighting approaches to status scores and clique identification. *J. Math. Sociol.* **2**(1), 113–120 (1972). <https://doi.org/10.1080/0022250X.1972.9989806>
43. R. Pastor-Satorras, C. Castellano, P. Van Mieghem, A. Vespignani, Epidemic processes in complex networks. *Rev. Mod. Phys.* **87**(3), 925 (2015). <https://doi.org/10.1103/RevModPhys.87.925>

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen® journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

---

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)

---