# EVALUATING SCIENTIFIC PUBLICATIONS BY N–LINEAR RANKING MODEL

Vu Le Anh (Ho Chi Minh City, Vietnam)
Hai Vo Hoang (Ho Chi Minh City, Vietnam)
Hieu Le Trung (Da Nang, Vietnam)
Kien Le Trung (Thua Thien Hue, Vietnam)
Jason J. Jung (Seoul, Korea)

Dedicated to András Benczúr on the occasion of his 70th birthday

Communicated by Le Manh Thanh

(Received June 1, 2014; accepted July 1, 2014)

Abstract. Ranking has been applied in many domains using recommendation systems such as search engine, e-commerce, and so on. We will introduce and study N-linear mutual ranking, which can rank n classes of objects at once. The ranking scores of these classes are dependent to the others. For instance, PageRank by Google is a 2-linear ranking model, which ranks the web-pages and links at once. Particularly, we focus to N-star ranking model and demonstrate it in ranking conference and journal problems. We have conducted the experiments for the proposed models to classical ones. The experiments are based on the DBLP dataset, which contains more than one million papers, authors and thousands of conferences and journals in computer science. The experimental results show that N-star ranking model evaluates everything much more detail based on the context of their relationships.

Key words and phrases: N-star ranking, Markov chain, PageRank, Academic ranking, Conference ranking, Ranking algorithms, Prolific ranking, Recommendation systems, Bibliographical database, DBLP

<sup>1998</sup> CR Categories and Descriptors: H.3.3 Information Search and Retrieval

This work is the extended version of the paper "A General Model for Mutual Ranking Systems" in *Intelligent Information and Database Systems*, Lecture Notes in Computer Science **8397**, 2014, pp. 211-220.

## 1. Introduction

Ranking is an interesting but difficult problem on many information processing systems. With a large amount of information, the systems need to adapt efficient ranking schemes to sort out (or to select) only the information which are highly relevant to the users' contexts. Particularly, in the context of *bibliometrics*, a set of given entities can be quantified to compare several evaluation indicators (e.g., popularity and reputation). For example, impact factors (IF) of international journals can be measured by taking into account how many times the papers in the corresponding journals have been cited.

In this work we focus on a system of ranking classes. Their ranking scores have mutual dependencies, which can be expressed by a system of linear equations. Let us explain the ideas by two examples.

PageRank. PageRank is very well-known ranking for website [19], which was applied in Google search engine. We rewrite the original formula by a system of two generic linear equations describing the mutual dependency of ranking of two classes, Web and Link

(1.1) 
$$Link \leftarrow 100\% \times Web,$$

$$(1.2) Web \leftarrow 85\% \times Link + 15\% \times Random.$$

Equation 1.1 says that the rank score of a link is determined by the rank score of the web, in which the link is included. Equation 1.2 says that the rank score of a web is determined by 85% from the rank score of links, which refers to the web; and 15% from randomness. Thus the rank scores of webs and links are mutually dependent. Moreover, we prove that there exists only one system of rank scores that satisfies the above system of linear equations.

Ranking scientific publication. We propose a model for ranking 4 classes Authors (Author), Publications (Pub), Conference (Conf) and Citations (Cite). Their relationships are described by the following system of four generic linear equations

$$(1.3) \qquad \qquad Author \leftarrow 100\% \times Pub,$$

$$(1.4) Conf \leftarrow 100\% \times Pub,$$

(1.5) 
$$Cite \leftarrow 100\% \times Pub,$$

(1.6) 
$$\begin{array}{rcl} Pub & \longleftarrow & 30\% \times Author + 30\% \times Conf \\ & & + 30\% \times Cite + 10\% \times Random. \end{array}$$

Equation 1.3 says that the rank score of each author is determined by the rank scores of his publications. Equation 1.4 says that the rank score of each conference is determined by the rank scores of its publications. Equation 1.5 says that the rank score of each citation is determined by the rank score of the publication, which is the owner of the citation. Equation 1.6 says that the rank score of each publication is determined by 30% from the rank score of its authors; 30% from the rank score of its conference; 30% from the rank scores of the citations; and 10% from randomness.

Both of above ranking systems are described by systems of linear equations called *N*-linear ranking models. Here are the key questions: Does the system of linear equations have a unique solution? How can we compute the solution? And how do the models work in realistic ranking systems? We solve only a part of problems by studying a special case of N-linear ranking model, *N*-star ranking model. Both of two examples are N-star ranking models and there exist unique solutions. Moreover we can estimate it by a loop of computing the linear function.

The main contributions and outline of this paper are as follows.

*N-linear ranking model.* We describe the background of the N-linear ranking model (Section 2). N-linear ranking model is the system of N ranking scores of N classes. The rank scores depend on others by a linear constraint system (Subsection 2.1). We introduce the affect and reflect relation between classes (Subsection 2.1). We explain these definitions in detail by the case study of PageRank (Subsection 2.2).

*N-star ranking model.* We define the N-star ranking model as a N-linear ranking model in which there exists a core class (Section 3). We prove that there exists unique N-star ranking model which satisfies a given linear constraint system (Proposition 3.2). We show that PageRank is a 2-star ranking model (Proposition 3.1). Finally, we describe the algorithm to compute scores of classes based on the linear constraint system.

Ranking bibliographical database. We study two N-star ranking models for the author, publication and conference ranking problem in different contexts (Section 4). The first model is general N-star ranking model for 4 classes: authors, publications, conferences, citations (Definition4.1). In the second model, we simplify the conditions by the assumption that everything is equal (Definition 4.2).

*Experiments.* We do the experiments for the simple N-linear ranking model of authors, publications and conference ranking (Section 5). We have designed the three different datasets to adapt the limit of computing resources (Subsection 5.1). The datasets are classified into two contexts: with/ without citations. We propose the models and the measurements for comparing different ranking

scores in both contexts (Subsection 5.2). We show the results and have discussions over datasets (Subsection 5.3). Our results are quite different from the naive one's and provide us some interesting things.

*Related works and conclusion.* We discuss the related works of N-linear ranking model (Section 6). We discuss the power and applicability of N-linear ranking model (Section 7).

#### 2. Backgrounds

### 2.1. N-linear mutual ranking system

The couple  $(\mathcal{A}, R)$  is called a *ranking system* if (i)  $\mathcal{A} = \{a_1, \ldots, a_n\}$  is a finite set, and (ii)  $R : \mathcal{A} \to [0; +\infty)$ .  $\mathcal{A}$  is called a *class*,  $a \in \mathcal{A}$  is called an *object* of the class  $\mathcal{A}$ , and R is called a *score* of  $\mathcal{A}$ . R is positive if  $R(a) > 0 \quad \forall a \in \mathcal{A}$ .  $n = |\mathcal{A}|$  is the size of  $\mathcal{A}$ .

**Definition 2.1.**  $\Omega = \{(\mathcal{A}_i, R_i)\}_{i=1}^N$  is called a N-linear mutual ranking system described by a system  $\{\alpha_{ij}, \beta_i, I_i, W_{ij}\}$  if  $(\mathcal{A}_i, R_i)$  is a ranking system and  $\alpha_{ij}, \beta_i \in [0; +\infty)$ ,  $I_i = (t_u^{(i)})_{n_i}$ ,  $n_i = |\mathcal{A}_i|$ , is a  $n_i$ -dimensional normalized nonnegative real number vector,  $W_{ij} = (\omega_{kl}^{(ij)})_{n_i \times n_j}$  is a nonnegative real number and normalized columns matrix such that for all  $i = 1, \ldots, N$ ,

$$\sum_{j} \alpha_{ij} + \beta_i = 1 \quad and \quad R_i = \sum_{j=1}^{N} \alpha_{ij} W_{ij} R_j + \beta_i I_i$$

 $\{\alpha_{ij}, \beta_i, I_i, W_{ij}\}$  is called a *linear constraint system* of  $\Omega$ .

Note that, generally since  $\sum_{i} \alpha_{ij} + \frac{1}{N} \sum_{j} \beta_{j}$  is different one, a *N*-linear mutual ranking system is not a Markov chain. Let  $a_{iu}, a_{jv}$  be objects in  $\mathcal{A}_i, \mathcal{A}_j$  respectively. Suppose  $\mathcal{C}^*(a_{iu}, a_{jv}) = \alpha_{ij} \omega_{uv}^{(ij)}$ . From the definitions, we have

$$R_i(a_{iu}) = \sum_{j=1}^N \sum_{v=1}^{n_j} \alpha_{ij} \omega_{uv}^{(ij)} R_j(a_{jv}) + \beta_i t_u^{(i)} = \sum_{j=1}^N \sum_{v=1}^{n_j} \mathcal{C}^*(a_{iu}, a_{jv}) R_j(a_{jv}) + \beta_i t_u^{(i)}.$$

 $a_{jv}$  is called affect to  $a_{iu}$  (denoted by  $a_{jv} \to a_{iu}$ ) if  $\mathcal{C}^*(a_{iu}, a_{jv}) > 0$ . Class  $\mathcal{A}_i$ is called total affect and reflect directly to class  $\mathcal{A}_j$  (denoted by  $\mathcal{A}_i \to \mathcal{A}_j$ ) if  $\forall a_{jv} \in \mathcal{A}_j$ :  $\exists a_{iu_1}, a_{iu_2} \in \mathcal{A}_i$ :  $a_{jv} \to a_{iu_2} \land a_{iu_1} \to a_{jv}$ .

**Definition 2.2.** Class  $\mathcal{A}_i$  is called total affect and reflect to class  $\mathcal{A}_j$ , denoted by  $\mathcal{A}_i \rightsquigarrow \mathcal{A}_j$ , if  $\mathcal{A}_i \rightarrow \mathcal{A}_j$  or  $\exists \mathcal{A}_k : \mathcal{A}_i \rightarrow \mathcal{A}_k \land \mathcal{A}_k \rightsquigarrow \mathcal{A}_j$ .

## 2.2. PageRank

We rewrite the PageRank into a 2-linear mutual ranking system as follows:

 $\mathcal{W} = \mathcal{A}_1$  is the class representing for the set of webpages.  $\mathcal{L} = \mathcal{A}_2$  is the class representing for hyperlinks. For each hyperlink  $l \in \mathcal{L}$  from web  $u \in \mathcal{W}$  to web  $v \in \mathcal{W}$ , we denote u = in(l) and v = out(l). For each  $v \in \mathcal{W}$ , we denote:  $IN(v) = \{l \in \mathcal{L} | v = out(l)\}$  and  $N_{out}(v) = |\{l \in \mathcal{L} | v = in(l)\}|$ .

PageRank[19] determined the ranking system of webpages by the following formula:  $\forall v \in \mathcal{W}$ ,

(2.1) 
$$R_w(v) = d \sum_{l \in IN(v), u = in(l)} \frac{R_w(u)}{N_{out}(u)} + \frac{1 - d}{|\mathcal{W}|}.$$

where  $d \in (0, 1)$  is a constant.

Suppose  $W_{21} = (\delta_{kt})_{|\mathcal{L}| \times |\mathcal{W}|}$  is a matrix in which  $\delta_{kt} = \frac{1}{N_{out}(w_t)}$  if  $l_k \in \{l : w_t = in(l)\}$ , otherwise 0.  $W_{21}$  is a nonnegative real number and normalized columns matrix. Suppose  $W_{12} = (\gamma_{tk})_{|\mathcal{W}| \times |\mathcal{L}|}$  is a matrix in which  $\gamma_{tk} = 1$  if  $w_t = out(l_k)$ , otherwise 0. We construct a 2-linear mutual ranking system on two classes  $\mathcal{W}$  and  $\mathcal{L}$  as follows: Let  $\bar{R}_w$  and  $\bar{R}_l$  be scores on the classes  $\mathcal{W}$ ,  $\mathcal{L}$  respectively. They satisfy

(2.2) 
$$\bar{R}_l = W_{21}\bar{R}_w$$
 and  $\bar{R}_w = dW_{12}\bar{R}_l + (1-d)I_{|\mathcal{W}|},$ 

where  $I_{|\mathcal{W}|}$  denotes the  $|\mathcal{W}|$ -dimensional vector in which all its elements are  $1/|\mathcal{W}|$ . It is not difficult to see that (2.2) confirms: for all webs  $v \in \mathcal{W}$ ,

$$\bar{R}_w(v) = d \sum_{l \in IN(v), u = in(l)} \frac{R_w(u)}{N_{out}(u)} + \frac{1 - d}{|\mathcal{W}|}$$

Since the equation (2.1) has the unique solution which is the PageRank score (see in [19]),  $\bar{R}_w$  is the PageRank score  $R_w$ . Vice versa, if  $R_w$  is a solution of (2.2),  $R_w$  should be  $\bar{R}_w$ . Thus, the PageRank score  $R_w$  is totally determined by the equation (2.2), or in other words, PageRank can be presented as the two-linear ranking system described by (2.2).

Note that, since for each link  $l \in \mathcal{L}$ , let u = in(l) and v = out(l), then web u affects to link l,  $(u \to l)$  and link l affects to web v,  $(l \to v)$ . Therefore, the class  $\mathcal{W}$  is total affect and reflect directly to the class  $\mathcal{L}, \mathcal{W} \to \mathcal{L}$ .

#### 3. N-star Ranking Model

**Definition 3.1.** Let  $\Omega = \{(\mathcal{A}_i, R_i)\}_{i=1}^N$  be a N-linear mutual ranking system.  $\Omega$  is called a N-star ranking if

1.  $\exists i: \mathcal{A}_i: (\beta_i > 0) \land (I_i \text{ is positive }) \land (\forall \mathcal{A}_i (j \neq i): \mathcal{A}_i \rightsquigarrow \mathcal{A}_i).$ 

2.  $\forall j \neq i : \alpha_{j1} = 1.$ 

 $\mathcal{A}_i$  is called a core of the system  $\Omega$ .

If  $\Omega = \{(\mathcal{A}_i, R_i)\}_{i=1}^N$  is a N-star ranking system described by a linear constraint system  $\{\alpha_{ij}, \beta_i, I_i, W_{ij}\}, \{\alpha_{ij}, \beta_i, I_i, W_{ij}\}$  is called N-star constraint system of  $\Omega$ .

**Proposition 3.1.** PageRank is a 2-star ranking system.

**Proof.** Because 1 - d > 0 and  $\mathcal{W} \rightsquigarrow \mathcal{L}$ ,  $\mathcal{W}$  is the core of PageRank. The second condition is clear since  $\alpha_{21} = 1$ . Hence, PageRank is a 2-star ranking system.

The ranking scores are determined by the N-star constraint system and the classes.

**Proposition 3.2.** Suppose the classes  $\{\mathcal{A}_i\}_{i=1}^N$  and the N-star constraint system  $\{\alpha_{ij}, \beta_i, I_i, W_{ij}\}$  are given. There exists a unique  $\{R_i\}_{i=1}^N$  in which  $R_i$  is a score on  $\mathcal{A}_i$ , such that  $\Omega = \{(\mathcal{A}_i, R_i)\}_{i=1}^N$  is a N-star ranking described by  $\{\alpha_{ij}, \beta_i, I_i, W_{ij}\}$  and for all i,  $\sum_{a \in \mathcal{A}_i} R_i(a) = 1$ .

**Proof.** Assuming without loss of generality that  $A_1$  is a core of a *N*-star ranking system described by  $\{\alpha_{ij}, \beta_i, I_i, W_{ij}\}$ , the sequence of scores  $R_1, \ldots, R_N$  satisfy the following equations

 $(3.1) R_1 = WR_1 and R_i = W_{i1}R_1 \forall i = 2, \dots, N,$ 

where

(3.2) 
$$W = \alpha_{11}W_{11} + \alpha_{12}W_{12}W_{21} + \dots + \alpha_{1N}W_{1N}W_{N1} + \beta_1\mathbf{I}_1$$

and  $\mathbf{I}_1$  is the  $(n_1 \times n_1)$ -matrix whose columns are  $I_1$ . It is not difficult to infer that because  $W_{1i}$  and  $W_{i1}$  are transition matrices (i.e. nonnegative and normalized columns matrices), the new square matrix  $W_{1i}W_{i1}$  is a stochastic matrix (i.e. a transition and square matrix). The matrices  $W_{11}$  and  $\mathbf{I}_1$  are also stochastic matrices. Since  $\mathbf{I}_1$  has positive entries,  $\beta_1 > 0$  (because  $\mathcal{A}_1$  is the core), and  $\sum_{j} \alpha_{1j} + \beta_1 = 1$ , the matrix W is also a stochastic matrix with positive entries. The Perron-Frobenius theorem (see in [7, 11]) confirms that there exists a unique score  $R_1$  with  $\sum_{a \in \mathcal{A}_1} R_1(a) = 1$  such that

$$R_1 = WR_1$$

From (3.1), the unique existence of  $R_1$  infers the unique existences of  $R_2, \ldots, R_N$ . Moreover, since  $W_{21}, \ldots, W_{N1}$  are normalized columns and  $\sum_{a \in \mathcal{A}_1} R_1(a) = 1$ , we

have  $\sum_{a \in A_i} R_i(a) = 1$  for all i = 2, ..., N. The proposition is proven.

The ranking scores are computed by following algorithm:

**Algorithm :** Finding the sequence scores  $\{R_i\}_{i=1}^N$ 

**Input** :  $\alpha_{ij}$ ,  $\beta_i$ ,  $W_{ij}$ ,  $I_i$ 

**Output :**  $\{R_i\}_{i=1}^N$ 

1. begin

2. Check the N-star ranking model with the core  $\mathcal{A}_1$ 

3. Let 
$$W \leftarrow \alpha_{11}W_{11} + \sum_{i=2}^{N} \alpha_{1i}W_{1i}W_{i1} + \beta_1 \mathbf{I}_1$$

4. Initialize  $R_1^{(0)}$ : uniform distribution, k = 0

5. repeat

6. k = k + 1

7. Update  $R_1^{(k)} \leftarrow W R_1^{(k-1)}$ 

**s. until**  $\|R_1^{(k)} - \overline{R}_1^{(k-1)}\| \le a$  stopping criterion

9. Let  $R_1 = R_1^{(k)}$  and  $R_i = W_{i1}R_1$ , i = 2, ..., N

10. end

## 4. Ranking authors, publications and conferences

In this section, we apply the N-star ranking model for constructing a model to evaluate authors, publications and conferences (journals) in the world of science. Concretely, we consider a four-star ranking model corresponding with four ranking systems:  $(\mathcal{A}, R_a) - \mathcal{A}$  is a set of all scientists who have publications,  $(\mathcal{P}, R_p) - \mathcal{P}$  is a set of all publications,  $(\mathcal{C}, R_c) - \mathcal{C}$  is a set of all sciential conferences and sciential journals, and  $(\mathcal{L}, R_l) - \mathcal{L}$  is a set of all citations between publications.  $R_a$ ,  $R_p$ ,  $R_c$  and  $R_l$  are the scores for each classes  $\mathcal{A}$ ,  $\mathcal{P}$ ,  $\mathcal{C}$  and  $\mathcal{L}$ , respectively.

For each citation  $l \in \mathcal{L}$ , u = in(l) and v = out(l) if l is from  $u \in \mathcal{P}$  to  $v \in \mathcal{P}$ . For each publication  $v \in \mathcal{P}$ , we denote:  $IN(v) = \{l \in \mathcal{L} | v = out(l)\};$  $OUT(v) = \{l \in \mathcal{L} | v = in(l)\}$   $N_{out}(v) = |OUT(v)|$ . If a publication does not cite any publication, we assume that it cites to all publications;  $C(v) \in \mathcal{C}$  is the conference of v;  $A(v) \subseteq \mathcal{A}$  is the set of authors of v. For each author  $a \in \mathcal{A}$ ,  $P(a) = \{v \in \mathcal{P} | a \in A(v)\}$  is a set of publications of a. For each conference  $c \in \mathcal{C}$ ,  $P_c(c) = \{v \in \mathcal{P} | c = C(v)\}$  is a set of publications published in c.

The 4-star ranking system model for ranking authors, publications, conferences and citations is constructed based on some following ideas:

 The score of an author depends only on his publications, and each publication affects to all of its authors: ∀a ∈ A, p ∈ P :

(4.1) 
$$R_a(a) = \sum_{p' \in P(a)} \mathcal{C}^*(a, p') R_p(p') \text{ and } \sum_{a' \in A(p)} \mathcal{C}^*(a', p) = 1$$

If a publication affects equally to its authors (a), (4.1) is rewritten as follows:

 $\forall a \in \mathcal{A}, p \in \mathcal{P}, a' \in A(p):$ 

(4.2) 
$$C^*(a',p) = \frac{1}{|A(p)|} \quad and \quad R_a(a) = \sum_{p' \in P(a)} \frac{R_p(p')}{|A(p')|}.$$

2. The score of a conference depends only on its publications:

(4.3) 
$$\forall c \in \mathcal{C} : R_c(c) = \sum_{p' \in P_c(c)} R_p(p').$$

The score of a citation depends on the citing publication, and each publication affects to all of its citations:
 ∀l ∈ L, p ∈ P, p' = in(l) :

(4.4) 
$$R_l(l) = \mathcal{C}^*(l, p') R_p(p') \text{ and } \sum_{l' \in OUT(p)} \mathcal{C}^*(l', p) = 1.$$

If a publication affects equally to its citations (b), (4.4) is rewritten as follows:

 $\forall l \in \mathcal{L}, p \in \mathcal{P}, p' = in(l), l' \in OUT(p):$ 

(4.5) 
$$C^*(l',p) = \frac{1}{N_{out}(p)} \quad and \quad R_l(l) = \frac{R_p(p')}{|N_{out}(p')|}.$$

4. The score of a publication depends on its citations, its authors and its conference and randomly finding by some reader. Each conference affects to all of its publications. Each author affects to all of its publications. Hence:

$$\forall p \in \mathcal{P}, c \in \mathcal{C}, a \in \mathcal{A}, c' = C(p) :$$

(4.6) 
$$R_{p}(p) = \alpha_{1} \sum_{l' \in IN(p)} R_{l}(l') + \alpha_{2} \sum_{a' \in A(p)} \frac{\mathcal{C}^{*}(p, a')}{\alpha_{2}} R_{a}(a') + \alpha_{3} \frac{\mathcal{C}^{*}(p, c')}{\alpha_{3}} R_{c}(c') + \beta_{p} I_{p},$$

(4.7) 
$$\sum_{p'\in P(a)} \mathcal{C}^*(p',a) = \alpha_2 \quad and \quad \sum_{p'\in P_c(c)} \mathcal{C}^*(p',c) = \alpha_3,$$

where  $\alpha_1, \alpha_2, \alpha_3, \beta_p > 0$  and  $\alpha_1 + \alpha_2 + \alpha_3 + \beta_p = 1$ ,  $I_p$  is a  $|\mathcal{P}|$ -dimensional normalized uniform random vector.

If a conference affects equally to its publications (c) and an author affects equally to its publications (d), the equation (4.6) and (4.7) are rewritten as follows:

 $\forall p \in \mathcal{P}, c \in \mathcal{C}, a \in \mathcal{A}, c' = C(p), p' \in P(a), p'' \in P_c(c) :$   $(4.8) \qquad \mathcal{C}^*(p', a) = \frac{\alpha_2}{|P(a)|} \quad and \quad \mathcal{C}^*(p'', c) = \frac{\alpha_3}{|P_c(c)|},$ 

(4.9) 
$$R_p(p) = \alpha_1 \sum_{l' \in IN(p)} R_l(l') + \alpha_2 \sum_{a' \in A(p)} \frac{R_a(a')}{|P(a')|} + \alpha_3 \frac{R_c(c')}{|P_c(c')|} + \beta_p I_p$$

**Definition 4.1.** The model which is described by equations (4.1), (4.3), (4.4), (4.6) and (4.7) is called the general 4-star ranking model for the ranking authors, publications and conferences problem.

**Definition 4.2.** The model which is described by equations (4.2), (4.3), (4.5), (4.8) and (4.9) is called the simple 4-star ranking model for the ranking authors, publications and conferences problem.

Both of the general and simple 4-star ranking models are N-star ranking systems in which the publication class is the central.

## 5. Experiments

## 5.1. Experiment environments

Computing environments. Our computing resource is a desktop with Intel core duo E7500 (2.8GHz) CPU and 2GB RAM. Our programming language is C#. Because of the limit of computing resource, the number of processed entities is limited by 3 millions.

Datasets. All of datasets are built from the DBLP data sets<sup>\*</sup>, on fields of computer science. We have built program to parse DBLP dataset in XML format to extract the authors, title, and publication venue information from the guides [14, 15]. We mention our readers that DBLP has no information about the citations between publications. The citations are collected from Academic Microsoft<sup>†</sup>. Since a DBLP publication can be both in a conference and another journal and there is a limit of computing resource, we keep only the publications related to conferences and ignore the journals.

We have chosen the 3 following datasets:

- $D_1$  contains all publications for all conferences collected in DBLP. Because of the limit of computing resource, we do not process the citations (which are about over 2 millions). Hence,  $D_1$  is for ranking 3 classes publications, authors and conferences.
- $D_2$  contains all not small publications in conferences which have at least 300 papers. There are about 70% small conferences (the number of publications is smaller than 300) and the total number if their papers are less than 30% (see Table 1). Hence we choose  $D_2$  to study the distorting effect of small-published conferences. By doing experiments both in  $D_1$  and  $D_2$ , we study more exactly about in the case we have no citation information.
- $D_c$  dataset contains publications in database conferences only.  $D_c$  has internal citations, they are cited from a publication of the dataset to another one inside the dataset, too. Hence,  $D_c$  is for ranking 4 classes publications, authors, conferences and citations.

The statistical figures of three datasets are shown in Table 1.

<sup>\*</sup>http://dblp.uni-trier.de accessed on May 2013

<sup>&</sup>lt;sup>†</sup>http://http://academic.research.microsoft.com accessed on March 2014

Datasets	nPubs.	nAuthors.	nConfs.	nCites.
$D_1$	1253997	845295	3351	-
$D_2$	1045888	746504	949	-
$D_c$	77361	87354	644	156429

Table 1. Experimental datasets.

### 5.2. Measurements with/without citations

#### 5.2.1. Models without citations

In the case there is no citation information  $(D_1 \text{ and } D_2 \text{ datasets})$ , we propose two models: Simple DBLP 3-star Ranking (SD3R) model and Simple DBLP 3-star Ranking (SD3R) model for the experiments.

*NPC Model.* Because there is no information of citations, we consider all publications are equal. Thus the naive model just counts the number of publications to evaluate the authors and conferences. Hence,  $\forall a \in \mathcal{A}, c \in \mathcal{C}$ :

(5.1) 
$$R_a(a) ::= |P(a)|,$$

(5.2) 
$$R_c(c) ::= |P_c(c)|.$$

SD3R Model. The model is determined from Definition 4.2. Because there is no information of citations, we omitted equations (4.5), (4.8). For the simplicity, we propose

$$\alpha_1 = 0, \qquad \alpha_2 = \alpha_3, \qquad \beta_p = 1 - 2\alpha_2.$$

The equation (4.9) is rewritten as follows :  $\forall p \in \mathcal{P}, \ c' = C(p),$ 

(5.3) 
$$R_p(p) = \alpha_2 \sum_{a' \in A(p)} \frac{R_a(a')}{|P(a')|} + \alpha_2 \frac{R_c(c')}{|P_c(c')|} + (1 - 2\alpha_2)I_p.$$

# 5.2.2. Models with citations

In the case there is citation information  $(D_c \text{ datasets})$ , we propose the model Simple Citation 4-star Ranking (SC4R) model and compare it to other ranking systems, which are (i) Naive model based on citation counting (NCC) for ranking author and conference; (ii) H-index for ranking authors.

*SC4R Model.* The model is determined from Definition 4.2. For the simplicity, we propose  $\alpha_1 = \alpha_2 = \alpha_3$  and  $\beta_p = 1 - 3\alpha_1$ . The equation (4.9) is

rewritten as follows:  $\forall p \in \mathcal{P}, c' = C(p),$ 

(5.4) 
$$R_p(p) = \alpha_1 \sum_{l' \in IN(p)} R_l(l') + \alpha_1 \sum_{a' \in A(p)} \frac{R_a(a')}{|P(a')|} + \alpha_1 \frac{R_c(c')}{|P_c(c')|} + (1 - 3\alpha_1)I_p.$$

*NCC Model.* We rank the authors and conferences based on only the citations. Hence,  $\forall a \in \mathcal{A}, c \in \mathcal{C}$ :

(5.5) 
$$R_a(a) ::= \sum_{v \in P(a)} |IN(v)|.$$

(5.6) 
$$R_c(c) ::= \sum_{v \in P_c(c)} |IN(v)|.$$

*H-index.* H-Index is introduced in [1, 6]. The index is based on the distribution of citations received by a given researcher's publications. An author a has index h if h of his/her |P(a)| papers have at least h citations each, and the other |P(a)| - h papers have no more than h citations each. We rewrite it. Suppose

$$I(a,k) = |\{v \in P(a) | IN(v) \ge k\}| \qquad (k \in \mathcal{N}, a \in \mathcal{A}).$$

The ranking of an author is defined as follows

(5.7) 
$$R_a(a) ::= Max(\{k \in \mathcal{N} | I(a,k) \ge k\}).$$

## 5.2.3. Measure the difference of ranking scores

In this subsection we study how do evaluate the difference between two ranking scores. Given a set of n objects  $\mathcal{A} = \{\omega_1, \omega_2, \ldots, \omega_n\}$  and two ranking scores  $R_1$ ,  $R_2$  on it. We measure two kinds of differences measurements: concordance measurement and different value measurement.

Concordance measurement. The ranking scores  $R_1$  and  $R_2$  are called "concordant" when large values of  $R_1$  go with large values of  $R_2$  (see in [18]). More precisely, given  $R_1$  and  $R_2$ , two objects ( $\omega_i, \omega_j$ ) are concordant if

$$[R_1(\omega_i) - R_1(\omega_j)][R_2(\omega_i) - R_2(\omega_j)] \ge 0,$$

and discordant if

$$[R_1(\omega_i) - R_1(\omega_j)][R_2(\omega_i) - R_2(\omega_j)] < 0.$$

From this idea, we state that  $R_1$  and  $R_2$  are *similar* if the probability of  $(\omega_i, \omega_j)$  be concordant is high and the probability of  $(\omega_i, \omega_j)$  be discordant is low, and *different* if vice versa. We also propose the *Kendall measure* which is introduced in [12] as a tool to evaluate these quantities. The Kendall measure of  $R_1$  and  $R_2$  is defined as follows (5.8)

$$\mathcal{K}(R_1, R_2) = \frac{1}{\frac{1}{2}n(n-1)} \Big( \sharp \big\{ (\omega_i, \omega_j): \text{ concordant} \big\} - \sharp \big\{ (\omega_i, \omega_j): \text{ discordant} \big\} \Big),$$

where  $\sharp(A)$  denotes the number of elements in the set A. It is clear that  $-1 \leq \mathcal{K}(R_1, R_2) \leq 1$ , and it receives value -1 if  $R_1$  and  $R_2$  are totally different and 1 if  $R_1$  and  $R_2$  are totally similar. We also measure the Spearman's rank correlation coefficient [3], denoted by  $\rho_{\mathcal{A}}(R_1, R_2)$ , between two ranking scores  $R_1, R_2$  of  $\mathcal{A}$ .

Different value measurement. The necessary condition of the measurement is that  $R_1$  and  $R_2$  should be normalized. We measure the different value between  $R_1$  and  $R_2$  as follows:

 $\forall \omega \in \mathcal{A}:$ 

(5.9) 
$$\Delta^{R_1,R_2}(\omega) = |R_1(\omega) - R_2(\omega)|,$$

(5.10) 
$$\%\Delta^{R_1,R_2}(\omega) = \frac{\Delta^{R_1,R_2}(\omega)}{R_1(\omega)}.$$

We also measure  $Avg_{\Delta}(R_1, R_2)$ , average value of  $\Delta^{R_1, R_2}$  over  $\mathcal{A}$ ; determine  $TopN_{\Delta}^{inc}(R_1, R_2)$  and  $TopN_{\Delta}^{des}(R_1, R_2)$  which are the top N increasing and decreasing  $\Delta^{R_1, R_2}$  values of  $\mathcal{A}$ .

### 5.3. Experimental results and discussions

#### 5.3.1. SD3R vs. NPC

**Remark 5.1.** The rank score of conferences in SD3R and NPC models are almost the same, but  $\Delta^{NPC,SD3R}$  reflects how hot the conferences are.

Figure 1 shows that all most of top 20 conference ranking values are the same for both methods and in both datasets  $D_1, D_2$ . There are a slightly different values of top 20 conferences between  $D_1$  and  $D_2$ , since  $D_2$  contains only conferences having more than 300 publications. Figure 3 shows that the average of  $\Delta_c^{SD3R,NPC}$  of conferences is small, less than 0.1 for both  $D_1$  and  $D_2$ . It also emphasizes that SD3R and NPC methods give very similar results when ranking conferences. The Spearman correlation coefficient of ranking



Figure 1. Top 20 conference ranking by SD3R vs. NPC,  $\alpha_2 = 0.45$ 

conferences  $\rho_c(SD3R, NPC)$  near by 1.0 indicates that the conference ranking scores of both methods are perfect monotone.

The top 5 *increasing* conferences (see Figure 2-a) are young, annual events with hot topics, such as remote sensing (IGARSS), computer human interaction (CHI), medical image computing (MICCAI), solid-state circuits (ISSCC), intelligent robot (IROS). The top 5 *decreasing* conferences (see Figure 2-b) are held for over a long time or biennial events, in local community (GI, MFCS) or long exploited topics such as artificial intelligence (IJCAI, AAAI), image processing (IFIP).

**Remark 5.2.** SD3R ranks authors differently and reflects the contribution

Figure 2. Top 5 most different ranking scores of conferences by SD3R vs. NPC over  $D_1$  and  $D_2$ ,  $\alpha_2 = 0.45$ 

	a) most increasing value										
CONF	FULL NAME	YEAR	NPC	SD3R(D1)	$\Delta^{\text{SD3R[D1],NPC}}$	SD3R(D <sub>2</sub> )	$\Delta^{\text{SD3R[D2],NPC}}$				
IGARSS	IEEE International Geoscience and Remote Sensing Symposium	2005	9689	10640.91	951.91	11636.60	1947.60				
СНІ	Computer Human Interaction	1990	8737	9571.49	834.49	10074.77	1337.77				
MICCAI	Medical Image Computing and Computer - Assisted Intervention	1998	3778	4552.62	774.62	5120.85	1342.85				
ISSCC	International Solid-State Circuits Conference	2009	1185	1794.49	609.49	2233.59	1048.59				
IROS	International Conference on Intelligent RObots and Systems	1992	7904	8410.25	506.25	8747.37	843.37				
	b) most decreasing valu	e									
CONF	FULL NAME	YEAR	NPC	SD3R(D <sub>1</sub> )	$\Delta^{\text{SD3R(D1),NPC}}$	SD3R(D <sub>2</sub> )	$\Delta^{\text{SD3R(D2),NPC}}$				
IJCAI	International Joint Conference on Artificial Intelligence (biennial)	1969	5635	5158.60	-476.40	4861.19	-773.81				
IFIP	International Federation for Information Processing (biennial)	1959	2796	2383.08	-412.92	2100.05	-695.95				
GI	GI-Jahrestagung (language spoken is Germany)	1972	4349	4021.90	-327.10	3794.84	-554.16				
AAAI	Conference on Artificial Intelligence (USA)	1980	5990	5689.97	-300.03	5473.98	-516.02				
MFCS	Mathematical Foundations of Computer Science (Czechoslovakia area)	1972	2115	1833.59	-281.41	1621.16	-493.84				

DATASET	Confe	rences	Authors			
	ρ	Avg∆	ρ	Avg∆		
D1	0.997	0.083	0.693	0.343		
D <sub>2</sub>	D <sub>2</sub> 0.993 0		0.690	0.341		

Figure 3. The different ranking scores values by SD3R vs. NPC,  $\alpha_2 = 0.45$ 

# of the author better than NPC. $\Delta^{NPC,SD3R}$ helps us detecting the key authors.

The  $Avg_{\Delta}(SD3R, NPC)$  of authors in  $D_1$  and  $D_2$  are around 34% (see Figure 3). It indicates the rank scores of authors in SD3R and NPC models are significantly different. Figure 4 shows that SD3R is quite different from NPC in the case of the top 20 most different ranking scores of authors.

We realize that most authors in the top *decreasing* ranking values have a large number of publications (see more details in Figure 5, about five authors in Figure 4-b, Figure 4-d). These authors have a large number of co-authors publications and belonging to many different conferences.

From Figure 4-a and Figure 4-c, we observe that most of top *increasing* authors do not have a big number of publications. Most of their papers have only one author and are published in some specialized conferences. Additionally, in Figure 6, we find out they are really key-persons in their research topics in real life, such as:

- *Emeritus Professor Lotfi A. Zadeh* at the University of California, Berkeley invented the theory of fuzzy sets.
- *Professor Ryotaro Kamimura* at Tokai University specializes on Machine Learning and Pattern Recognition.
- *Ellen M. Voorhees* at NIST, is very famous from international workshops series: the Text REtrieval Conference (TREC), TREC Video (TRECVid), and the Text Analysis Conference (TAC).
- Senior Specialist Toshihiko Yamakami, at ACCESS, Japan Advanced Institute of Science and Technology, is professional on Mobile Social Application.
- *Professor Keqin Li* is from State University of New York at New Paltz, notable for parallel and distributed computing.

Also from Figure 4-a and Figure 4-c, there are some special cases in increasing ranking list, in which authors have high NPC ranking scores. We suppose that it is because almost all of their papers are published in high ranking score conferences. Some example authors are as following:



Figure 4. Top 20 most different ranking scores of authors by SD3R vs. NPC,  $\alpha_2=0.45$ 

Figure 5. Top 5 most decreasing ranking scores of authors having biggest NPC

NAME	NPC (D <sub>1</sub> )	SD3R(D <sub>1</sub> )	$\Delta_{a}^{SD3R,NPC}(D_{1})$	nConfs.(D <sub>1</sub> )	NPC (D <sub>2</sub> )	SD3R(D <sub>2</sub> )	$\Delta_a^{SD3R,NPC}(D_2)$	nConfs.(D <sub>2</sub> )
Wen Gao	528	398.23	-129.77	95	488	373.12	-114.88	81
Wei Wang	510	441.53	-68.47	236	448	394.00	-54.00	186
Wei Liu	467	393.73	-73.27	219	426	362.59	-63.41	183
Mahmut T. Kandemir	427	362.44	-64.56	78	405	350.31	-54.69	66
Mario Piattini Velthius	417	346.01	-70.99	100	315	259.88	-55.12	65

Figure 6. Top 5 most increasing ranking scores of authors having smallest NPC

NAME	NPC (D <sub>1</sub> )	SD3R(D <sub>1</sub> )	$\Delta_{a}^{SD3R,NPC}(D_{1})$	nConfs.(D <sub>1</sub> )	NPC (D <sub>2</sub> )	SD3R(D <sub>2</sub> )	$\Delta_a^{SD3R,NPC}(D_2)$	nConfs.(D <sub>2</sub> )
Lotfi A. Zadeh	53	130.50	77.50	29	45	103.998	58.99811962	21
Ryotaro Kamimura	62	151.46	89.46	14	55	134.444	79.44385132	11
Ellen M. Voorhees	62	133.97	71.97	11	61	139.214	78.21359319	10
Toshihiko Yamakami	70	195.17	125.17	29	55	159.163	104.1627818	21
Keqin Li	74	151.74	77.74	14	70	141.899	71.89933631	12

a) Differen	t ranking v	alues of co	onferences	b) Different ranking values of authors				
R <sub>1</sub> ,R <sub>2</sub>	ρ	κ	Avg∆	R <sub>1</sub> ,R <sub>2</sub>	ρ	κ	Avg∠	
SC4R,NCC	0.910	0.833	1.087	SD4R,NCC	0.610	0.889	0.987	
SC4R,NPC	0.877	0.825	0.787	SD4R,NPC	0.715	0.889	1.054	
NCC,NPC	0.688	0.676	4.901	SD4R,H-index	0.602	0.855	1.631	

Figure 7. Different ranking measurement of SC4R with  $\alpha_1 = 0.3$  vs. NPC, NCC, H-index on  $D_c$ 

- Edwin Hancock is a very well known scientist on computer vision.
- Norman C. Beaulieu, a Canadian engineer and professor in the ECE Department of the University of Alberta is very famous in broadband digital and communications systems.
- *Professor Irith Pomeranz*, affiliated at School of Electrical and Computer Engineering, Purdue University is noble for Computer Engineering VLSI and Circuit Design.

# 5.3.2. SC4R, H-index, NCC, NPC

**Remark 5.3.** In the case when citations are considered, SC4R is more finegrained than NPC and NCC on the conference ranking problem. Moreover, the citations is the main factor making the rank score meaningful.

Because of the naive feature of counting, both NPC and NCC methods treat the citations or publications separately. There may be too many authors having the same NPC or NCC value. With SC4R, all citations and publications are ranked together, thus the SC4R rank scores seem to be more meaningful. The different ranking measurement figures are shown in Figure 7. The Kendall' tau coefficients point out that ranking by SC4R is more concordant with NCC than NPC method, and each pair of them is not quite the same concordant. It confirms again the fact that citation is the main factor for ranking scientific publications.

**Remark 5.4.** SC4R seems to be in the middle of the popularity and the prestige ranking of conferences.

Figure 8 shows the trend of SC4R lines in top ranking value from three methods. In each graph, the SC4R line always fluctuates between NPC (popularity) line and NCC (prestige) line, nearly coincides with the average of them.

Figure 9 shows top 10 rank positions by SC4R methods compared to complied lists of MS RANK [16] and CORE RANK  $^{\ddagger}$ . Let us see three interesting cases:

<sup>&</sup>lt;sup>‡</sup>http://www.core.edu.au/coreportal. Accessed on May -25, 2014.

Figure 8. Top highest ranking value of conferences by NPC, NCC and SC4R on  $D_c$  with  $\alpha_1 = 0.3$ 



- VLDB vs. HICSS. The number of articles published in HICSS is fivefold of VLDB's but HICSS's cited papers are just a fourth of VLDB. In our SC4R method, VLDB stands before HICSS. Their SC4R ranking values are 8403.14 and 7037.91 respectively, totally not far different.
- *PODS and ISWC*. We can easily point out two prestige database conferences, PODS and ISWC, which are not in top 10 in popularity, being 18 and 16 respectively.
- DEXA and ICEIS. We can see two low-ranked conferences by MS RANK and CORE experts assessments are listed in top 10 of SC4R conference rank, DEXA and ICEIS. This can be explained by the experimental dataset  $D_c$  contains many papers as well as citations of these two conferences. DEXA has 3982 papers and that number of ICEIS is 2993.

It is impressive that most of the top 10 of SC4R conference ranking list are the most noble in database field following the assessment of CORE experts and Microsoft Academic Search system.

**Remark 5.5.** SC4R seems to reflect the contribution of the author better than others from combining prestige and popularity criteria.

CONF	NAME	SC4R Rank	NPC Rank	NCC Rank	MS RANK	CORE RANK
VLDB	Very Large Data Bases	1	6	1	1	Α*
HICSS	Hawaii International Conference on System Sciences	2	1	5	5	Α
SIGMOD	International Conference on Management of Data	3	8	2	2	Α*
ICDE	International Conference on Data Engineering	4	3	3	3	Α*
PODS	Symposium on Principles of Database Systems	5	18	4	4	A*
DEXA	Database and Expert Systems Applications	6	2	11	18	В
CIKM	International Conference on Information and Knowledge Managem	7	5	7	7	Α
ER	International Conference on Conceptual Modeling	8	9	6	9	Α
ICEIS	International Conference on Enterprise Information Systems	9	4	32	51	С
ISWC	International Semantic Web Conference	10	16	9	6	Α

Figure 9. Top 10 SC4R rank position of conferences on  $D_c$  with  $\alpha_1 = 0.3$ 

Figure 10-a shows the top 20 highest NPC ranking author values. H-index line is the lowest line and is separated with other lines. It can be explained that many authors have published a big number of papers which get a few citations. We remind that H-index is proposed for the combination of popularity and prestige value of an author. Let us see other three graphs of Figure 10, which show the ranking values by NCC (focus on prestige), H-index (combining popularity and prestige) and SC4R. We found that H-index nearly coincides with NPC line. We suppose that H-index extremely falls into quantity pole in these case. Meanwhile, the NCC line, standing for the prestige of author via high citation, is always the highest line in Fig. 10-b, Fig. 10-c, Fig. 10d. SC4R does not fall extremely into popularity (NPC, H-index) or prestige (NCC) poles. Thus SC4R is the good method to reflecting the contribution of author in which combines prestige and popularity criteria.

Let us look details for some special authors in Fig. 11.

- Professor Jennifer Widom and Alon Halevy. They fall in the case of high citation, high H-index, high publications but lower SC4R. The citation number of Professor Jennifer Widom is nearly double than that of Michael Stonebraker, a pioneer of data base research and technology. She also has H-index higher than Michael Stonebraker but her number of publications is lower than his one. We can see her SD4R value is a bit lower than Michael Stonebraker's SD4R value. We assume the result is because many publications citing her papers do not get high score in SC4R ranking, namely not good quality.
- Jim Gray, Antonin Guttman and Umeshwar Dayal. They fall in the opposite situation in which all authors have low NPC, NCC and H-index values, but high SD3R score. This situation can be explained by their publications are cited by almost all prestige people and high quality publications and are published in noble conferences.

We review the top 20 highest SC4R value people and find out that they all are very famous in database field with many valuable scientific works.



Figure 10. Top 20 ranking value of author by NPC, NCC, H-index\* and SC4R on  $D_c$  with  $\alpha_1 = 0.3$ 

Figure 11. Top 20 SC4R ranking value on  $D_c$  with  $\alpha_1 = 0.3$ 

Author	NPC	NCC	H-index	SC4R	NPC Rank	NCC rank	H-index Rank	SC4R Rank
Michael Stonebraker	96	1653	24	603	27	11	7	1
Rakesh Agrawal	79	2237	22	564	50	4	10	2
Jennifer Widom	86	3150	30	522	40	1	2	3
Hector Garcia-Molina	100	1792	23	501	23	9	8	4
David Dewitt	87	2446	26	459	37	3	4	5
Jeffrey D. Ullman	59	1496	18	419	124	17	25	6
Surajit Chaudhuri	103	1771	25	416	21	10	5	7
Alon Halevy	91	2742	32	415	33	2	1	8
Jim Gray	30	695	11	414	503	62	111	9
Antonin Guttman	5	986	4	398	6568	32	1455	10
Dan Suciu	96	2147	25	394	26	5	6	11
H. V. Jagadish	118	2072	26	380	10	6	3	12
Serge Abiteboul	112	1440	20	380	13	19	12	13
Christos Faloutsos	90	1610	19	345	34	12	17	14
Umeshwar Dayal	107	832	17	332	19	41	28	15
Raghu Ramakrishnan	105	1487	20	332	20	18	13	16
Jeffrey Naughton	84	1834	23	320	44	8	9	17
Gerhard Weikum	143	995	18	303	5	31	23	18
Divesh Srivastava	128	1969	21	300	7	7	11	19
Jiawei Han	143	1061	16	299	4	30	32	20

143

### 6. Related work

*Our prior work.* Our work is the extended version of the prior work [9]. In this version, we have extended the experiments and represent everything more in detail. Concretely, we do the experiments for the case the citations are considered. We also compare our ranking with the H-index, which is the most famous ranking scores for authors recently. Moreover, we improve the results by giving more discussions.

Web-pages ranking. The ranking problem occurs and develops quickly with the era of the Internet and big data. One of the most famous ranking problem is ranking web-pages. A brief overview of this problem can be found at Dilip Kumar Sharma et al. [4]. The recent survey on web-pages ranking algorithms was conducted by Mercy Paul Selvan et al. [13]. They categorized these ranking algorithms into three groups: i) Link analysis algorithms, ii) Personalized web search ranking algorithms and iii) Page Segmentation algorithms. The first group is highlighted by two notable classical algorithms, naming PageRank [19] utilized by Google search engine and Hyperlink-Induced Topic Search (HITS) [8] developed by Jon Kleinberg. Since the hyperlink structure among the webpages is easily represented as a web graph, the PageRank of each web-page can be measured (see Sect. 2.2 for more details). HITS was a precursor to PageRank. The idea behind HITS algorithm classifies the webs into two classes: (i) hubs, served as large directories point to (ii) authoritative pages. A good hub represented a page that pointed to many other pages, and a good authority represented a page that was linked by many different hubs. The model can be rewritten into 2-linear ranking model. The advantages and disadvantages of link analysis algorithms are also discussed in works of Dilip Kumar Sharma et al. [4] and Mercy Paul Selvan et al. [13].

Link-based object ranking. The basic problems of link-based object ranking can be found at the survey on link mining of Getoor and Christopher P. Diehl [10]. Zaiqing Nie et al. proposed Poprank model [23] to rank the popularity of objects based on their web popularity and the object relationship graph. It extends from PageRank and uses the Popularity Propagation Factor to express the relationship between classes of objects. The model is based on the Markov chain model which can be applied in the N-linear mutual ranking systems. Yizhou Sun et al. proposed NetClus algorithm [21] that utilizes links between objects of multi-typed to rank cluster of multi-typed heterogeneous networks. This work is the successor of their previous work, RankClus, which can rank and cluster one-typed objects mutualy [22]. Their model has the same idea as ours when limiting only within a star network schema and giving rank distribution for each type of objects. But our approach is different from theirs. NetClus's objects in the center class (target type) only relates to those belonging to other classes (attribute types). So it can not be applied to compute on the citation network. Other complex ranking systems have been already explored using a different formalism for ranking or classification in heterogeneous networks. For example, the quantium ranking [5] is based on quantium navigation. Their formula is comes from the quantium theory and quite different to ours.

Bibliometric ranking. Many methods of assessing the intellectual impact, reputation and influence of scientists, journals, conferences have been proposed over years. In a very recent scientometric study [17], Paul Benjamin Lowry et al. compare expert assessments to bibliometric measures for determining a tiered structure of information systems (IS) journals. They categorize the assessing journal quality methods into three methodology paths: (i) bibliometricbased, (ii) expert-based and (iii) non-validated approach. One of their noticeable conclusion is that bibliometrics can be a complete, less expensive and more efficient substitute for expert assessment. For bibliometric approaching, many detailed definitions of citation metrics can also be found in this paper, such as the most used ISI Impact Factor, h-index and it variants like g-index, Another study of bibliometric graph-based algorithms focus on e-index... ranking researchers was conducted by Xiaorui Jiang et al. [20]. They compare sophisticated citation analysis algorithms like PageRank, SARA, CoRank, FutureRank, P-Rank, BiRank with some simpler methods like citation count and sum of paper ranks, similar to the way we evaluate the experiment results. Further information about bibliometrics and web-based citation analysis can be seen on [2].

# 7. Conclusion and future works

We have introduced and studied N-linear ranking systems. The mutual relationships between ranking objects are described by a system of linear equations. A N-linear mutual ranking system is a N-star ranking system if it has a core class which affects and reflects all other classes in the system. The rank scores of the N-star ranking system are unique and computed by a loop of computing the linear function. We have pointed out that PageRank is a 2-star ranking. It has two classes: the web-pages (a core class) and links.

We have introduced and studied a general and a simple 4-star ranking models for ranking authors, publications, conferences. A general model is a generic one. In a simple model, we consider each publication, author, conference, citation is equal. We have conducted the experiments for the models in two contexts with/without citations. The experimental results are based on the DBLP dataset. We have compared the difference of the values and the concordance between the proposed ranking systems (SD3R and SC4R) and naive, classical ranking systems (NPC, NCC, H-index). We have found that:

- In the case citations are not considered, the rank score of conferences in SD3R and NPC model are almost the same, but  $\Delta^{NPC,SD3R}$  reflects how hot the conferences are. SD3R ranks authors differently and reflects the contribution of the author better than NPC.
- The citations are the main factors for measuring the prestige of publications, conferences and authors. SC4R seems to reflect the rank of the author, the conference better than others from combining prestige and popularity criteria.

As future work, we are planning to i) develop the current system to retrieve big data including publications and citations; ii) study how to combine the proposed ranking system with class keywords, which are tagged in the publications; iii) investigate the time series in N-star ranking and the trend prediction problem, and iv) apply N-star ranking systems in various ranking problems, e.g., business ranking, event ranking, and so on.

## References

- Sidiropoulos, A., D. Katsaros and Y. Manolopoulos, Generalized Hirsch h-index for disclosing latent facts in citation networks, *Scientometrics*, 72 (2) (2006), 253–280.
- [2] Cronin, B., Bibliometrics and beyond: some thoughts on web-based citation analysis, *Journal of Information Science*, 27 (1) (2001), 1–7.
- [3] Spearman, C., "General intelligence," Objectively determined and measured, The American Journal of Psychology, 15 (2) (1904), 201–292.
- [4] Sharma, D.K. and A.K. Sharma, A comparative analysis of web page ranking algorithms, *International Journal on Computer Science and En*gineering, (2010), 2670–2676.
- [5] Snchez-Burillo, E., J. Duch, J. Gmez-Gardenes and D. Zueco, Quantum navigation and ranking in complex networks, Scientific reports 2, 2012.
- [6] Hirsch, J.E., An index to quantify an individual's scientific research output, Proc. of the National Academy of Sciences of the United States of America, 102 (46), (2005), 16569-16572.

- [7] Keener, J., The Perron-Frobenius theorem and the ranking of football teams, SIAM Review, 35 (1) (1993), 80–93.
- [8] Kleinberg, J., Authoritative sources in a hyperlinked environment, Journal of the ACM, 46 (5) (1999), 604-632.
- [9] Vu, L.A., V.H. Hai, L.T. Kien, L.T. Hieu and J.J. Jason, A general model for mutual ranking systems, *Intelligent Information and Database* Systems, Lecture Notes in Computer Science 8397, 2014, 211–220.
- [10] Getoor, L. and C.P. Diehl, Link mining: a survey, ACM SIGKDD Explorations Newsletter, 7 (2) (2005), 3-12.
- [11] Kien, L.T., L.T. Hieu, T.L. Hung and L.A. Vu, MpageRank: The stability of web graph, Vietnam Journal of Mathematics, 37 (2009), 475– 489.
- [12] Kendall, M.G., A new measure of rank correlation, *Biometrika*, (1938), 81–93.
- [13] Selvan, M.P., A.C. Sekar and A.P. Dharshin, Survey on web page ranking algorithms, *International Journal of Computer Applications*, 41 (19) (2012), 1–7.
- [14] Ley, M., DBLP Some lessons learned, *PVLDB*, **2** (2) (2009), 1493–1500.
- [15] Ley, M. and P. Reuther, Maintaining an online bibliographical database: The problem of data quality, EGC 2006, 5–10.
- [16] Microsoft Corporation, Microsoft Academic Search, http://academic.research.microsoft.com/ (June-26-2013)
- [17] Lowry, P.B., G.D. Moody, J. Gaskin, D.F. Galletta, S.L. Humpherys, J.B. Barlow and D.W. Wilson, Evaluating journal Quality and the Association for Information Systems Senior Scholars' Journal basket via bibliometric measures: Do expert journal assessments add value?, *The MIS Quarterly*, **37** (4) (2013), 993–1012.
- [18] Nelsen, R.B., An Introduction to Copulas, 2nd ed., Springer Series in Statistics, Springer, New York, 2006.
- [19] Brin, S. and L. Page, The anatomy of a large-scale hypertextual web search engine, Proc. 7th International World Wide Web Conference, 1998, 107–117.
- [20] Jiang, X., X. Sun and H. Zhuge, Graph-based algorithms for ranking researchers: not all swans are white!, *Scientometrics*, 96 (3) (2013), 743– 759.
- [21] Sun, Y., Y. Yu and J. Han, Ranking-based clustering of heterogeneous information networks with star network schema, *Proc. 15th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, 2009, 797-806.
- [22] Sun, Y., J. Han, P. Zhao, Z. Yin, H. Cheng and T. Wu, Rankclus: integrating clustering with ranking for heterogeneous information network analysis, *Proc. the 12th Int. Conf. on Extending Database Technology: Advances in Database Technology*, ACM, 2009, 565–576.

[23] Nie, Z., Y. Zhang, J.-R. Wen and W.-Y. Ma, Object-level ranking: Bringing order to web objects, Proc. 14th Int. Conf. on World Wide Web, ACM, 2005, 567–574.

# Vu Le Anh

Nguyen Tat Thanh University Ho Chi Minh City, Vietnam e-mail lavu@ntt.edu.vn

# Vo Hoang Hai

Information Technology College Ho Chi Minh City, Vietnam e-mail vohoanghai2@gmail.com

Hieu Le Trung Duy Tan University Da Nang, Vietnam e-mail hieukien820gmail.com

# Kien Le Trung

Hue University of Science Thua Thien Hue, Vietnam e-mail hieukien@hotmail.com

# Jason J. Jung

Chung-Ang University Seoul, Korea e-mail j2jung@gmail.com