



# CNAVER: A Content and Network-based Academic Venue Recommender system<sup>☆</sup>

Tribikram Pradhan<sup>\*</sup>, Sukomal Pal

Department of Computer Science and Engineering, Indian Institute of Technology (BHU), Varanasi, Uttar Pradesh, India



## ARTICLE INFO

### Article history:

Received 24 April 2019

Received in revised form 1 October 2019

Accepted 6 October 2019

Available online 17 October 2019

### Keywords:

Venue recommender system

Social network analysis

Meta-path analysis

Random walk with restart (RWR)

Graph clustering

Rank-based fusion

## ABSTRACT

The phenomenon of rapidly developing academic venues poses a significant challenge for researchers: how to recognize the ones that are not only in accordance with one's scholarly interests but also of high significance? Often, even a high-quality paper is rejected because of a mismatch between the research area of the paper and the scope of the journal. Recommending appropriate scholarly venues to researchers empowers them to recognize and partake in important academic conferences and assists them in getting published in impactful journals. A venue recommendation system becomes helpful in this scenario, particularly when exploring a new field or when further choices are required. We propose CNAVER: A Content and Network-based Academic Venue Recommender system. It provides an integrated framework employing a rank-based fusion of paper-paper peer network (PPPN) model and venue-venue peer network (VVPN) model. It only requires the title and abstract of a paper to provide venue recommendations, thus assisting researchers even at the earliest stage of paper writing. It also addresses cold start issues such as the involvement of an inexperienced researcher and a novel venue along with the problems of data sparsity, diversity, and stability. Experiments on the DBLP dataset exhibit that our proposed approach outperforms several state-of-the-art methods in terms of precision, nDCG, MRR, accuracy,  $F - measure_{macro}$ , average venue quality, diversity, and stability.

© 2019 Elsevier B.V. All rights reserved.

## 1. Introduction

A recommender system recommends different objects based on a user's preferences using various data analysis techniques [1–3]. Academic recommender systems can suggest collaborators [4–6], papers [7–11], citations [12–16], and/or academic venues [15, 17–22]. These systems have been useful to academicians as they objectively provide users with personalized information services [22–25]. Although there have been quite a few works on different academic recommendations, very little body of work exists in the literature on academic venue recommendations, albeit started quite early [19].

The researchers, in general, intend to publish in academic venues that acknowledge high-quality papers and participate in academic conferences or workshops that are relevant to their area of research [26,27]. Among various problems that researchers confront, an important task is to identify appropriate publication venues. The task is nowadays being increasingly difficult

due to the continuous increase in the number of research areas and dynamic change in the scope of journals [19,28]. More collaborations are taking place among disciplines in the research communities, which is leading to reduced compartmentalization at the coarse level but a continuous increase in the number of venues in interdisciplinary areas [19,29]. For example, DBLP<sup>1</sup> dataset, a collection of scientific publication records and their relationship within that collection has 9585 computer science conferences<sup>2</sup> and the number of journals is more than 4152<sup>3</sup> [30].

As the research horizon expands, researchers find it challenging to remain up to date with new findings, even within their disciplines [31]. Moreover, with time, researchers' interests expand, evolve, or adapt in rapidly changing subject areas needing information on appropriate venues in the changed scenario [25]. Increase in interdisciplinary research areas also poses great challenges to research institutes and their libraries as they strive to understand information-seeking behaviors and dynamic information needs of the users [19]. Information specialists need timely and seamless information on researchers' reading priorities to make decisions on venue subscriptions instead of relying only on the venues' impact factor or users' explicit requests.

<sup>☆</sup> No author associated with this paper has disclosed any potential or pertinent conflicts which may be perceived to have impending conflict with this work. For full disclosure statements refer to <https://doi.org/10.1016/j.knosys.2019.105092>.

<sup>\*</sup> Corresponding author.

E-mail addresses: [tpradhan.rs.cse16@itbhu.ac.in](mailto:tpradhan.rs.cse16@itbhu.ac.in) (T. Pradhan), [spal.cse@iitbhu.ac.in](mailto:spal.cse@iitbhu.ac.in) (S. Pal).

<sup>1</sup> <http://dblp.uni-trier.de/db/>

<sup>2</sup> <http://dblp.uni-trier.de/db/conf/>

<sup>3</sup> <http://dblp.uni-trier.de/db/journals/>

On the other hand, researchers also need to know about new venues to remain updated. They usually get updates from colleagues/supervisors, friends, internet, and books but often the information is not sufficiently comprehensive and/or appropriate as their research demands. The researchers, therefore, sometimes end up approaching inappropriate venues resulting in rejections, delays in publication and/or compromise in the quality of the publication. Venue recommendation for either journals or conferences, in particular, has, therefore, become an essential area of research in recent times [32]. Out of many reasons for its increasing importance, some are given below as emerging scenarios [22].

- (a) A researcher from the industry has made a breakthrough in her research area. To collaborate with her peers from academia, she may want to find a suitable academic venue (conference) that she is not very aware of.
- (b) A junior researcher, i.e., a researcher who is at the initial stage of her research and has no or very few publications, intends to extend her research area. But a lack of knowledge about appropriate academic venues becomes a challenge for her to explore newer areas.
- (c) A veteran researcher knows her research area very well, but when she ventures into a new field or works in an interdisciplinary area, she may look for a cross-field venue recommendation.
- (d) A journal may merge with some other related journal with modified scopes and objectives. The researchers may not be aware of such developments.

To recommend a suitable venue of high quality, we need to focus from the perspective of a researcher's needs and development of the particular research area in question on the following issues.

- (i) What are the most relevant venues of publications for a researcher in question?
- (ii) How can a researcher find high-quality venues?
- (iii) What are the most suitable conferences/workshops a researcher should participate in, for a given area?

Most of the existing techniques depend on co-authors' past publications and/or ratings of the venues provided by other researchers to perform such a venue recommendation. A few approaches use a random walk model, topic-based similarity for the same. Based on our literature survey, we attempt to explore the following research questions (RQs).

- RQ1:** To what extent the existing collaborative filtering (CF) based models can recommend suitable venues from the perspective of a researcher's needs?
- RQ2:** Can the current content-based filtering (CBF) approaches recommend suitable venues from the perspective of a researcher's needs?
- RQ3:** Is there any issues with the network-based (NB) recommendation model on venue recommendation?
- RQ4:** To what extent do the existing approaches handle cold-start problems like new researchers, new venues, and other issues such as data sparsity, diversity, and stability, etc.?

Several issues with the existing state-of-the-art techniques have been reported in the literature. For example, content-based and network-based similarity techniques utilize only a single aspect of either content or citations from scientific papers, respectively [9,33]. Very few attempts have been made that considered both content and network, especially in the domain of venue recommendation. To bridge this gap, we would like to propose a hybrid method CNAVER: An improved content-based filtering and improved network-based fusion model to provide a personalized academic venue recommender system.

CNAVER is built on two major components that contribute in parallel, but finally, their contributions are fused together to present a coherent venue recommendation. One is the paper-paper peer network (PPPN) model and the other venue-venue peer network (VVPN) model. While PPPN explores the interaction among papers towards venue recommendation, VVPN actually studies it among publication venues.

Key contributions of this work are the followings:

- To deal with "cold-start"<sup>4</sup> issues like a new researcher and a new venue, PPPN and VVPN models are fused to provide a personalized venue recommender system. CNAVER works irrespective of researchers' past publication records, instead only focuses on the work at hand. New venues with no citations available are also considered for abstract similarity. New venues are also given an equal chance for inclusion in the final recommendation.
- To address "data sparsity"<sup>5</sup> issues citation network is used to examine the importance of each candidate paper through the cumulative scores of centrality measures such as degree, betweenness, and closeness, etc. Later on, contextual similarity such as LDA on abstract and Doc2Vec on the title are performed to reduce the bibliographic network size and also to increase the relatedness among papers.
- To resolve the issue of "diversity"<sup>6</sup> a fusion model incorporating both PPPN and VVPN model are proposed. Specifically, age-discounted(age-discounted) based Venue2Vec, meta-path features, and biased random walk are incorporated into the NB model to recommend venues from diverse publishers. The proposed system CNAVER can provide a "diversified recommendation". It takes into account both journals and top tier conferences from multiple publishers like Elsevier, Springer, IEEE, ACM, and others.
- To address the issue of "stability"<sup>7</sup> a fusion model CNAVER incorporating both PPPN and VVPN models are proposed. Any network-based approach is known to cause instability in the ranks with time as the introduction of new nodes or edges change the topology and thereby change recommendations [34]. We therefore also took into account content-based approaches at several stages within both PPPN and VVPN pipeline. A separate study on stability proved that it worked and ours was the most stable system, more than the existing ones.
- Comprehensive experiments were conducted using a real-world dataset, i.e., DBLP to evaluate the performance of the proposed system CNAVER. The proposed system outperforms several other state-of-the-art venue recommendation models with substantial improvements in precision@k, nDCG@k, MRR, accuracy,  $F - measure_{macro}$ , average venue-quality (ave-quality), diversity and stability.

This paper is organized as follows. We visit the related literature in Section 2. We provide motivation in Section 3 and then more elaborate problem description in Section 4. Description of different features we adopted in our recommendation framework is provided in Section 5, data preprocessing steps are described in Section 6, contextual similarity calculations are shown in Section 7, description of peer-peer network models are depicted in Section 8 and fusion model is illustrated in Section 9. Experimental details are illustrated in Section 10. We report experimental results and insightful discussions are in Sections 11 and 12 respectively. We finally conclude in Section 13.

<sup>4</sup> Cold start issues mainly indicate the new researchers and new venues in academia.

<sup>5</sup> Sparsity denotes average distance among pairs of related papers.

<sup>6</sup> Diversity means how many different venues are recommended.

<sup>7</sup> Stability denotes resilience to change in ranked recommendations with the introduction of new papers.

## 2. Related work

Adomavicius and Tuzhilin [2] authored a comprehensive review of recommender systems and suggested mainly three types of recommender systems based on their working principles. In addition, we also include network-based recommendation [22]. We attempt to provide here necessary background in the academic recommender systems according to their taxonomy.

### 2.1. Collaborative filtering based recommendation (CF)

Collaborative recommender systems (or collaborative filtering systems) predict the utility of items for a user based on the items previously rated by other users who have similar likings or tastes [2]. In the field of academic recommendations, Yang et al. [35] proposed a model to explore the relationship between publication venues and writing styles using three kinds of stylistic features: lexical, syntactic and structural. In another paper, Yang et al. [36] used a collaborative filtering model incorporating writing style and topic information of papers to recommend venues. Yang et al. [15] proposed another joint multi-relational model (JMRM) of venue recommendation for author-paper pairs.

Hyunh et al. [37] proposed a collaborative knowledge model (CKM) to organize collaborative relationships among researchers. The model quantified the collaborative distance, the similarity of actors before recommendations. Yu et al. [38] proposed a prediction model that used collaborative filtering for a personalized academic recommendation based on the continuity feature of a user's browsing content. Liang et al. [1] proposed a probabilistic approach consolidating user exposure that was modeled as a latent variable, inducing its incentive from data for collaborative filtering. Alhoori et al. [19] recommended scholarly venues taking into account the researcher's reading behavior based on personal references and the temporal factor of when references were added. Trappey et al. [39] presented a new patent recommendation system based on clusters of users having similar patent search behaviors.

### 2.2. Content-based filtering based recommendation (CBF)

In CBF, users are recommended items similar to the ones the user preferred in the past. Kochen et al. [40] proposed a method to recommend journals for authors' manuscripts based on relevance, acceptance rate, and prestige of journals. Medvet et al. [20] considered the title and abstract of papers to recommend scholarly venues considering  $n$ -gram based Canvar-Trenkle, two-steps-LDA, and LDA + clustering to retrieve language profile, a subtopic of papers, and identification of the main topic as a research field.

Errami et al. [41] proposed a model called eTBLAST to recommend journals based on abstract similarity using the  $z$ -score of a set of extracted keywords and weighted formula of "Journal score". Schurmie et al. [42] proposed the Journal/ Author Name Estimator (Jane)<sup>8</sup> on biomedical database MEDLINE to recommend journals based on abstract similarity. They exploited a weighted  $k$ -nearest neighbors and Lucene similarity score to rank articles. Similarly, Wang et al. [9] presented a content-based publication recommender system (PRS) on computer science exploiting soft-max regression and chi-square based feature selection techniques.

Recently, few online services have started providing support for suggesting journals using keywords, title, and abstract matching. These services include Elsevier Journal Finder<sup>9</sup> [43],

Springer Journal Suggester,<sup>10</sup> Edanz Journal Selector<sup>11</sup> and End-Note Manuscript Matcher<sup>12</sup> etc. Elsevier Journal Finder requires only the title and abstract of a paper and uses noun phrases as features and Okapi BM25+ to recommend journals. But, recommendations are restricted to Elsevier publishers only [43].

### 2.3. Network-based recommendation (NB)

On top of the above approaches, the approach based on a network representation of the input data has gained considerable attention in the recent past. Here, a social graph is built among the authors based on co-authorship. An edge exists between two authors if they co-author at least one paper [21,44]. The venue having the highest count among the papers within  $n$ -hops from a given author-node is recommended. Klamma et al. [32] proposed a Social Network Analysis (SNA) based method using collaborative filtering to recognize most similar researchers and rank obscure events by integrating the rating of most similar researchers for the recommendations. Silva et al. [45] proposed a three-dimensional research analytics framework (RAF), incorporating relevance, productivity, and connectivity parameters.

Pham et al. [46] used the number of papers of a researcher in a venue to determine her rating for that venue using the clusters on social networks. Later, Pham et al. [47] presented clustering techniques on a social network of researchers to identify communities to generate venue recommendations. They also applied traditional CF calculations to provide the suggestions. Chen et al. [28] introduced a model AVER to recommend the scholarly venues to a target researcher. This approach utilizes a random walk with restart (RWR) model on the co-publication network incorporating author-author and author-venue relations. Later, Yu et al. [22] extended AVER to personalized academic venue recommendation model PAVE where the topic distribution of researcher's publications and venues were utilized in LDA.

Luong et al. [48] identified suitable publication venues by investigating the co-authorship network, most frequent conferences, normalized scores based on most successive conferences. Luong et al. [21] in another work recommended suitable publication venues by investigating authors' co-authorship networks in a similar field. Xia et al. [18] provided venue recommendations using Pearson correlation and characteristic social information of conference participants to enhance smart conference participation.

### 2.4. Hybrid recommendation (HR)

Hybrid approaches combine collaborative and content-based methods avoiding certain limitations of content-based and collaborative systems. Wang et al. [9] proposed hybrid article recommendations incorporating social tag and friend information. Boukhris et al. [49] suggested a hybrid venue recommendation based on the venues of the co-citers, co-affiliated researchers, the co-authors of the target researcher. It is based on bibliographic data with citation relationships between papers. Minkov et al. [50] introduced a method of recommending future events. Tang et al. [4] introduced a cross-domain topic learning (CTL) model to rank and recommend potential cross-domain collaborators. Xia et al. [18] proposed a socially aware recommendation system for conferences. Similarly, Cohen et al. [5] explored the domain of mining-specific context in a social network to recommend collaborators.

<sup>8</sup> <http://jane.biosemantics.org>

<sup>9</sup> <http://journalfinder.elsevier.com>

<sup>10</sup> <http://journalfinder.com>

<sup>11</sup> <https://www.edanzediting.com/journal-selector>

<sup>12</sup> <http://endnote.com/product-details/manuscript-matcher>

### 3. Motivation

Although the techniques discussed above do provide recommendations reasonably well, they suffer from a lot of issues. Below we discuss the problems that led us to investigate further in connection to our stated research questions (RQs) as presented in Section 1.

#### 3.1. Problems with CF approach (RQ1)

Although CF has been quite popular in the last decade for scholarly venue recommendation, most of them suffer from the following drawbacks.

- (i) CF approaches are less effective when there are not enough ratings present in the researcher-venue matrix. The recommendations may not be useful in the case of a new researcher who lacks publication history.
- (ii) The techniques are not likely to recommend a new venue or a less popular venue as the venue lacks in publication statistics. Therefore, some relevant venues may be missed.
- (iii) Computational cost is high because of an extensive number of articles, venues, and researchers involved are taken into consideration during processing, and thus, scalability is a challenge.
- (iv) Researcher-venue matrix that is at the core of the techniques is exceptionally sparse as most of the researchers publish and cite a few articles and are involved with very few academic venues.

#### 3.2. Problems with CBF approach (RQ2)

CBF or topic-based models use the author's profile, the content of their papers as well as that of the papers published at a specific venue [51]. Most approaches use LDA for topic modeling and rank venues based on the similarity of venues that published similar papers [21]. But, in information retrieval, longer documents get an advantage during the computation of query-document similarity over their shorter counterparts [52]. Other salient issues with CBF approaches are as follows.

- (i) CBF approaches suffer from limited content analysis, which can significantly reduce the quality of recommendation [9, 33]. Most of the time, they require the full text of the paper and thus, are not usable at the early stage of paper-writing [20]. Usually, the abstract is not sufficient to extract the necessary reliable and relevant information.
- (ii) New venues are less likely to be recommended as the models prefer venues with a high number of papers published therein.
- (iii) The models provide a poor recommendation to a new researcher who lacks publication records.
- (iv) The recommendations are heavily biased towards the past area of research of a researcher and therefore not suitable when one changes her area of interest or works in an interdisciplinary field.

#### 3.3. Problems with NB approach (RQ3)

To alleviate the problem of limited content analysis in CBF and the cold-start issue occurs in CF approaches of late, network-based approach (NB) or co-author based approach has been proposed [7,22,28,45,47,48]. In this model, a social graph is built among the authors in light of co-authorship [15,46]. An edge exists between researchers if they co-author no less than one paper [8,21]. A couple of works consider a random walk with restart [22,28]. The venue having the highest count among the papers within  $n$ -hops from the author-node is recommended. A few limitations of the approaches are as follows.

**Table 1**

Type of vertices used in HIN.

No.	Vertices type
1	$P_{main} = \{\text{set of papers that belonging to a particular venue}\}$
2	$P_{ref} = \{\text{set of papers that cited by a } P_{main} \text{ paper}\}$
3	$P_{cite} = \{\text{set of papers that cites a } P_{main} \text{ paper}\}$
4	$A(author) = \{\text{author of any type of paper } (P_{main}, P_{cite}, P_{ref})\}$
5	$T(term) = \{\text{term appearing in titles or abstracts of a } P_{main} \text{ paper}\}$
6	$V(venue) = \{\text{set of any venue where } P_{main} \text{ type papers published}\}$

- (i) Irrespective of actual content, each paper authored by the same set of authors will receive the same recommendation.
- (ii) Recommendations are very poor for a new researcher who does not have any past publication records.
- (iii) It cannot recommend a new venue as the model is based on the publication history of venues.
- (iv) Venues with less popularity among the co-authors of a given author are seldom recommended, although content-wise they may be appropriate.

#### 3.4. Cold-start issues present in the existing approach (RQ4)

Most of the approaches discussed above suffer from various cold start issues for new researchers, new venues, or less popular venues and also other problems like data sparsity, scalability, diversity, and stability, etc. A few limitations of the existing approaches are as follows.

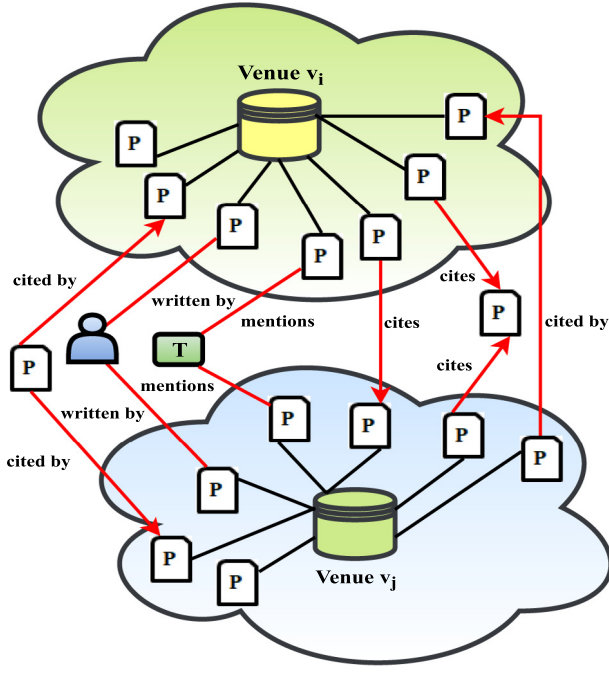
- (i) In the CF-based venue recommender system, data sparsity is a major issue that arises due to the sparseness of the matrix of researcher-venue ratings. It also suffers from the cold-start problem for new researchers.
- (ii) CBF performs poorly due to ambiguity in text comparison and also suffers from cold start issues of new researchers and new venues.
- (iii) Two major issues of CBF approaches are limited content analysis and over-specialization, and due to which, the lack of diversity is a severe problem in this type of approach [3, 9,33].
- (iv) Scalability is also a major challenge in CF and CBF based venue recommender systems as procedures therein are not linear in input-size.
- (v) Most of the time, stability is a severe issue in CF and NB based venue recommender systems.
- (vi) The existing techniques have an undue bias against the new venues and new researchers coming into the system with fewer publication records [21,41,47].

### 4. Problem description and other definitions

**Definition 1.** Heterogeneous Information Network (HIN) [53,54]. It is defined as a directed graph  $G=(\mathcal{N}, \mathcal{L})$  with a node type mapping function  $\delta: \mathcal{N} \rightarrow \mathcal{W}$  and a link type mapping function  $\mu: \mathcal{L} \rightarrow \mathcal{Y}$ . Each node  $n \in \mathcal{N}$  belongs to one particular node type in the node type set  $\mathcal{W}$ :  $\delta(n) \in \mathcal{W}$ , and each link  $l \in \mathcal{L}$  belongs to a particular link type in the link type set  $\mathcal{Y}$ :  $\mu(l) \in \mathcal{Y}$ . Here the types of nodes  $|\mathcal{W}| > 1$  and the type of links  $|\mathcal{Y}| > 1$ .

**Example.** In Fig. 1, we have six types of vertices, such that  $\mathcal{W} = \{P_{main}, P_{ref}, P_{cite}, A, T, V\}$  and six types of edges  $\mathcal{Y}$ . The meaning of each type of vertices is presented in Table 1. In Fig. 1, a  $P_{main}$  paper is either the paper within the set of venue  $v_i$  or venue  $v_j$  and categorized as main paper. Both  $P_{cite}$  and  $P_{ref}$  are considered as non-main papers and could be associated with any venues. Similarly, the meaning of each type of edges is defined in Table 2.





**Fig. 1.** Graphical representation of HIN graph.  $P_{main}$  paper is any  $P$  belongs to either one the disparate venues and latent meta paths between  $P_{main}$  papers may be formed via various vertices types: cited by  $P_{main}$  ( $P_{ref}$ ), cites a  $P_{main}$  ( $P_{cite}$ ), author( $A$ ), term( $T$ ), and venue( $V$ ).

**Definition 2.** Venue-Venue Graph (VVG). Let  $G' = (V', E')$  be the newly generated venue-venue graph (VVG) from HIN based on the similarity score of abstract and title.  $V' = \{v_1, v_2, \dots, v_l\}$ . Each edge  $e = (v_i, v_j) \in E'$  represents a currently similar research scope of  $v_i$  with  $v_j$  based on their past publications. An edge  $e = (v_i, v_j) \in E'$  exists if the similarity score between venues  $v_i$  and  $v_j$  is greater than average similarity score. We weight the edges of the network VVG using content similarity (linear combination of abstract and title) in order to provide a single score as explained in Section 8.2.2.

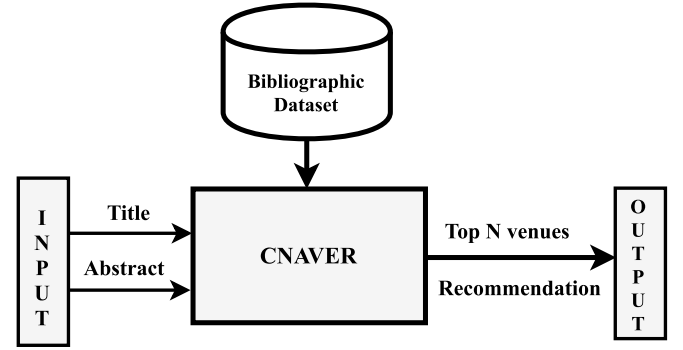
**Definition 3.** Meta-path [55]. A meta-path  $M$  is a path defined on the HIN graph. It joins two or more vertices using one or more edges such that  $M = n_1 \xrightarrow{l_1} n_2 \xrightarrow{l_2} \dots \xrightarrow{l_t} n_{t+1}$ , where the starting and ending vertices are of same vertex type  $P_{main}$ ,  $\delta(n_1) = \delta(n_{t+1})$  and both belong to  $P_{main}$ ,  $P_{main} \in \mathcal{W}$ ,  $\mu(l_1, l_2, \dots, l_t) \in \mathcal{Y}$ .

**Example.** In Fig. 1, There will be a meta path between venue  $v_i$  and venue  $v_j$  via the meta path Venue  $v_i \xrightarrow{\text{publish}} P_{main} \xrightarrow{\text{citedby}} P_{cite} \xrightarrow{\text{cites}} P_{main} \xrightarrow{\text{publishedby}} \text{Venue } v_j$ .

**Definition 4.** Random Walk [56]. A random walk is defined as a node sequence  $S_r = \{v_1, v_2, v_3, \dots, v_l\}$  wherein the  $i$ th node  $v_{i1}$  in the walk is randomly selected from the neighbors of its predecessor  $v_{i2}$ .

**Definition 5.** Citation Network. Let  $G = (V, E)$  be the citation graph, with  $n$  papers.  $V = \{p_1, p_2, \dots, p_n\}$ . In  $G$ , each directed edge  $e = (p_i, p_j) \in E$  represents a citation from  $p_i$  to  $p_j$ .

**Example.** We use the following two phrases to describe the citation network.



**Fig. 2.** The basic block diagram of CNAVER.

**Table 2**

Type of edges used in HIN.

No.	Edges type
1	$n_1 \xrightarrow{\text{written\_by}} n_2 : \delta(n_1) \in \{P_{main}, P_{ref}, P_{cite}\}, \delta(n_2) = A, n_1, n_2 \in N$
2	$n_1 \xrightarrow{\text{contains}} n_2 : \delta(n_1) \in \{P_{main}, P_{ref}, P_{cite}\}, \delta(n_2) = T, n_1, n_2 \in N$
3	$n_1 \xrightarrow{\text{cites}} n_2 : \delta(n_1) \in P_{main}, \delta(n_2) = P_{ref}, n_1, n_2 \in N$
4	$n_1 \xrightarrow{\text{cited\_by}} n_2 : \delta(n_1) \in P_{main}, \delta(n_2) = P_{cite}, n_1, n_2 \in N$
5	$n_1 \xrightarrow{\text{cites}} n_2 : \delta(n_1) \in P_{main}, \delta(n_2) = P_{main}, n_1, n_2 \in N$
6	$n_1 \xrightarrow{\text{cited\_by}} n_2 : \delta(n_1) \in P_{main}, \delta(n_2) = P_{main}, n_1, n_2 \in N$

- (i) References of  $p_i$  represent the set of papers which are referred by the paper  $p_i$ .
- (ii) Citation to  $p_j$  denotes the collection of papers which have used the paper  $p_j$  as a reference.

The rest of the paper, we use the above two phrases to describe the graph around vertex  $p_i$ .

**Definition 6.** Venue Recommendation. Let each paper  $p_i$  published in a particular venue  $v_i$ . So now we have,  $B = \{v_1, v_2, \dots, v_l\}$  be a predefined set of publication venues. Given a input paper (seed paper)  $p_m$ , the venue recommendation task is to recommend a list of suitable publication venues ( $v_1, v_2, \dots, v_N$ ) related to the seed paper  $p_m$ , where the list is ordered from the most relevant to the least relevant.

**Example.** It is essentially a ranking problem. We need to determine the set of papers which are closely related to the seed paper. Venue recommendations are provided if the title and abstract of a seed paper are given to the system as input. The block diagram is depicted in Fig. 2.

## 5. Architecture of CNAVER

We present an overall architecture of the proposed framework alongside its operational strategies. Our goal is to exhibit why a fusion model with a step-insightful layered approach has been chosen as opposed to a flat architecture containing a set of components. As the bibliographic dataset is exceptionally massive in size (2,408,010 papers), if we attempt to recognize the top-most similar papers for each seed paper by looking at contextual similarity against the entire dataset, the overall computational overhead will be high.

### 5.1. Framework of CNAVER

We propose a system comprised of two blocks: Block-I and Block-II as depicted in Fig. 3. To reduce computational overhead

and to make it independent and autonomous of seed papers, particularly Block-I, is developed once for the whole citation network. Later on, we will utilize the seed paper input to interact with Block-II to extract meaningful recommendations from both the PPPN model and VVPN model.

We present a layered architecture where each layer realizes a specialized task. The system consisting of four essential layers, where Layer-1 to Layer-3 belong to Block-I, and Layer-4 to Block-II.

Four primary layers are portrayed as given underneath:

- (i) Data preprocessing and centrality calculation (**Layer-1**): This layer aims to structure the dataset into a formal model for processing. Mainly it is used for faster extraction of relevant papers and the importance of each candidate papers for further use (**Block I**).
- (ii) Contextual similarity calculation (**Layer-2**): This layer can also be called the feature extraction layer and is mainly introduced to extract required contextual features needed to compute Paper2Vec in PPPN model and Venue2Vec in VVPN model. It is also used to filter only potentially useful papers from Set-II, based on content similarity (**Block I**).
- (iii) Peer-peer network model (**Layer-3**): This layer uses a peer-peer network to process the data and to make a recommendation. The objective of this layer is to reduce computational overhead and to make it independent of seed papers (**Block I**).  
This layer comprises of two distinct models, namely:
  - (a) PPPN model: The main objective is to capture the strength of individual papers and their citation relationship with other papers in a citation network to obtain relevant venues to the seed paper.
  - (b) VVPN model: The main idea behind this model is to capture the similarity (indirect relationship among venues via meta-path analysis) among venues in a heterogeneous bibliographic network to obtain relevant venues to the seed paper.
- (iv) Fusion model (**Layer-4**): To provide a diversified personalized recommendation, the PPPN, and VVPN models are utilized to make predictions individually and later on a fusion model firstly is applied to integrate the strengths of both the models and to reduce their weaknesses (**Block II**).

## 6. Data preprocessing and centrality calculation (Layer-1)

Initially, there were 3,079,007 rows and 7 columns in the combined dataset. After removing duplicate papers, papers with missing fields in the database, etc., we are left with 2,236,968 papers. We also drop papers having fields filled up with inconsistent entries. We also ignore non-textual content from the abstracts of the papers. The detailed statistics of the DBLP data collection are described in Section 10.1. All such papers are checked for their references section. We separately treat the papers having references or not.

- (i) The set of papers where references are available are called Set-I.
- (ii) The set of papers without references are called Set-II.

We generate a citation network only with the Set-I papers. Among the centrality measures, we use degree, betweenness and closeness measures (defined below) among such papers [57,58].

### 6.1. Degree centrality ( $C_D$ )

In a graph, the degree of a node is the number of edges that are adjacent to that node [59]. Higher the number of neighbors of a given node, the higher its impact is. Degree centrality of a paper  $p$  is defined as

$$C_D(p) = \text{indeg}(p) + \text{outdeg}(p) \quad (1)$$

where  $\text{indeg}(p)$  is the number of research articles or papers citing to paper  $p$  and  $\text{outdeg}(p)$  is the number of papers  $p$  is referring to.

For each paper  $p$ , in-degree ( $p$ ) is computed. The papers whose in-degree is greater than or equal to average in-degree of the network are shortlisted for further computation. Later on again, the average score of degree  $\{\text{indeg}(p) + \text{outdeg}(p)\}$  is taken into consideration for removing papers. We adopt such two-stage filtering in order to ensure that: (i) first, highly cited papers are not missed and (ii) no new papers which cite a lot of papers are missed either.

### 6.2. Betweenness centrality ( $C_B$ )

$C_B$  of a node quantifies how frequently the node shows up on different possible shortest paths between any two given nodes. Here  $C_B$  of a paper  $q$  is defined as

$$C_B(q) = \sum_{\substack{p, k, q \in V \\ p \neq k \neq q}} \frac{\sigma_{pk}(q)}{\sigma_{pk}} \quad (2)$$

where  $\sigma_{pk}$  denote the number of shortest paths from  $p$  to  $k$  and  $\sigma_{pk}(q)$  denote the number of shortest paths from  $p$  to  $k$  via  $q$ .

Nodes with high betweenness act as potential deal makers [60].

### 6.3. Closeness centrality ( $C_C$ )

The metric attempts to capture how centrally a node is located vis-a-vis other nodes and is measured as the inverse of total pair-wise distances from the node to all other nodes. Closeness centrality of a node  $p$  is defined as

$$C_C(p) = \frac{1}{\sum_{\substack{q \neq p \\ p \in V}} d_G(p, q)} \quad (3)$$

where  $d_G(p, q)$  denotes the distance between vertices  $p$  and  $q$ , i.e. the minimum length of any path connecting  $p$  and  $q$  in  $G$ .

We have presented a summary of the importance of various centrality measures used in this paper in Table 3. We use the above three measures to shortlist a set of candidate papers for Layer-2 (Contextual similarity calculation). The average score of each measure is used as a threshold to filter important papers from each category. Initially, we remove papers with less than average  $\text{indeg}$  as they are not cited by many and hence less influential. After filtering papers with low in-degree, papers with degree score greater than or equal to average degree scores are finally shortlisted for further computation. The sets were determined individually and merged as a set-based union to consider just the unique papers.

For example, if a very high-quality paper has low in-degree because of its recent publication, the paper may not be considered in degree centrality calculation, but it gets due consideration in Betweenness, and Closeness centrality calculation and, therefore, may qualify based on these measures. This way if a paper lacks in one or more factors in the citation profile, it can qualify through other centrality measures implying a fair chance to all potential papers.

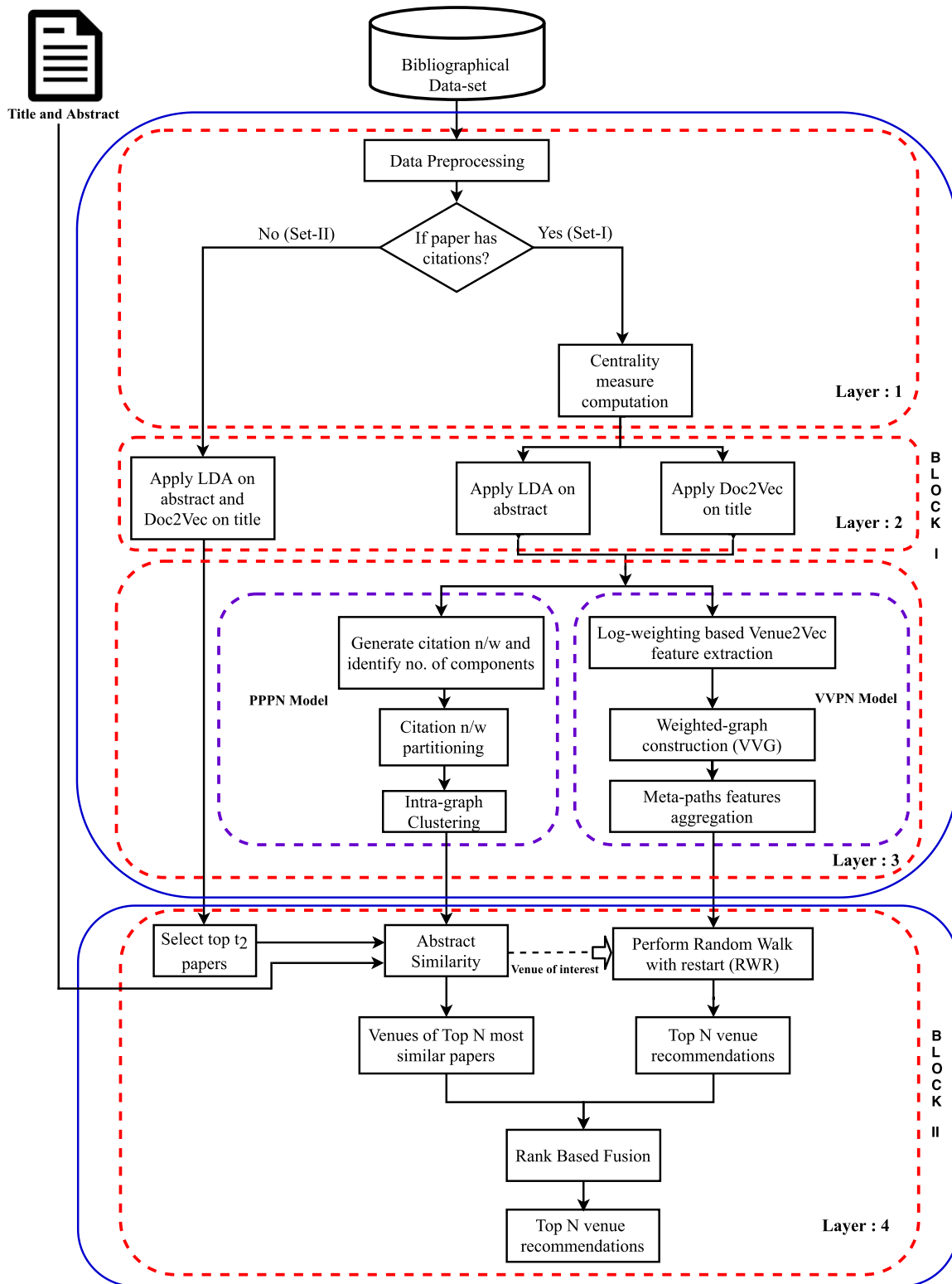


Fig. 3. Architecture of CNAVER.

**Table 3**  
Interpretation of various centrality measures.

Measures	Interpretation in citation networks
Degree	How many papers can this article reach directly?
Betweenness	How likely is this papers to be the most direct route between two papers in the citation network?
Closeness	How fast can this paper reach everyone in the citation network?

## 7. Contextual similarity calculation (Layer-2)

In this module, we mainly extract content-based features to prune the set of papers shortlisted in Layer-1 further. Sometimes it is quite challenging to observe the similarity among papers by looking at only the word level similarities. Even there are cases where the semantic meaning of words is unable to capture the similarity among papers where the use of words in context also needs to be seen. Hence, we need some mechanism specifically to capture the semantic meaning, to discover the hidden patterns and to extract the latent topics other than just identifying words matching. In this paper, we applied LDA on abstract and Doc2Vec on the title to address these above issues.

An abstract typically provides a summary containing the main idea of a paper. We use the LDA model on the abstract to generate the feature description [61]. LDA is used to identify topics automatically and to derive hidden patterns exhibited by a text corpus. We have chosen LDA over other methods due to its simplicity, easiness in implementation, fast computation, ability to discover coherent topics and also to handle diverse topics in a text corpus. We set the number of the topic as parameter  $k$  while mining a paper's topic distribution to perform LDA. It is used as it generates the probability distribution of words and documents based on the co-occurrence of words and documents, which focus on describing their connotative topics.

We also tried LDA on the title, but due to insufficient terms present in titles, it did not perform well to discover hidden patterns. Hence, Doc2Vec is used to extract the feature description from the title of a paper as Doc2Vec captures contextual information of words occurring in titles [62]. It is mainly used to generate sentence/document embeddings [63]. It is chosen over other methods due to its potential to overcome the weaknesses such as the ordering of words, the semantics of the words, data sparsity, and high dimensionality in bag-of-words models and other approaches. We have used Doc2Vec on the title but not on abstract because it was found a little bit expensive to represent each document by a dense vector that is trained to predict surrounding words in contexts sampled from the document.

## 8. Peer-peer network model (Layer 3)

Features so extracted from abstract and title are fused in the next layer to compute:

- (i) *Paper2Vec in PPPN model*: In PPPN model, we would like to explore identifying suitable venues through paper-paper peer network by exploiting the concept of Paper2Vec approach without the age-discounted scheme.
- (ii) *Venue2Vec in VVPN model*: In VVPN model, we would like to see the quality of the recommendation by incorporating venue-venue peer network through the concept of Venue2Vec approach.

### 8.1. The architecture of PPPN model

Due to information overload, it is not practical to check full content similarity to recognize related papers with the seed paper. To address this issue, we are attempting to discover inherent community structures in a bibliographic citation network to understand the network more deeply and reveal interesting relations among the papers.

The process of PPPN model mainly involves four steps: (i) Paper2Vec feature extraction, (ii) Citation network partitioning, (iii) Topic-oriented intra-graph clustering, (iv) Abstract similarity using Okapi BM25+ algorithm.

#### 8.1.1. Paper2vec feature extraction

The results from LDA and Doc2Vec can be considered as two sets of vectors. For each paper  $p_i$ , we get a vector  $A_i$  for abstract similarity and vector  $T_i$  for title similarity. The length of the vector is taken as size  $k$ . We are computing the vectors for both abstract and title only once and later on; we will utilize those vectors to calculate the similarity with seed papers. To avoid repetitive computation, a fixed-length vector is considered in this paper.

$$A_i^p = [a_{1i}, a_{2i}, \dots, a_{ki}] \quad (4)$$

$$T_i^p = [t_{1i}, t_{2i}, \dots, t_{ki}] \quad (5)$$

Using  $A_i^p$  and  $T_i^p$  for a paper  $p_i$ , we compute cosine similarity with their counterpart from the seed paper ( $p_j$ ).

$$Sim\_abstract(p_i, p_j) = \frac{A_i^p \cdot A_j^p}{|A_i^p| |A_j^p|} = \frac{\sum_{b=1}^k (a_{bi} * a_{bj})}{\sqrt{\sum_{b=1}^k a_{bi}^2} * \sqrt{\sum_{b=1}^k a_{bj}^2}} \quad (6)$$

$$Sim\_title(p_i, p_j) = \frac{T_i^p \cdot T_j^p}{|T_i^p| |T_j^p|} = \frac{\sum_{b=1}^k (t_{bi} * t_{bj})}{\sqrt{\sum_{b=1}^k t_{bi}^2} * \sqrt{\sum_{b=1}^k t_{bj}^2}} \quad (7)$$

The overall similarity between a shortlisted paper ( $p_i$ ) and the seed paper ( $p_j$ ) is calculated as a weighted sum of the two similarities.

$$Sim(p_i, p_j) = c * Sim\_abstract(p_i, p_j) + (1 - c) * Sim\_title(p_i, p_j) \quad (8)$$

where  $c \in [0, 1]$  is a tuning parameter.  $Sim(p_i, p_j)$  is used to find similarity with the seed paper (See Algorithm 1). Top  $R$  papers according to the above similarity are also chosen with the topmost paper being paper of interest ( $I$ ) for a given seed paper as discussed in Section 8.1.4.

Generally, researchers cite conceptually related and relevant papers to their work. But all cited papers are not conceptually related to the citing paper, and their corresponding venues may not be similar to the venue of seed paper.

To capture both the strength of connection as well as semantics such as the related topics shared by papers, we apply a hybrid approach of link analysis and topic-oriented intra-graph clustering in a bibliographic citation network.

To reduce the time complexity, we perform intra-graph clustering in two stages:

- (i) To find sub-graphs for the entire citation network found after centrality measure based on modularity<sup>13</sup> maximization.
- (ii) Within a sub-graph apply intra-graph clustering based on both link and content information.

There are other reasons for this two-stage intra-graph clustering. We attempt to cluster the entire citation network found after centrality measures using the Jarvis-Patrick algorithm. But due to unexpected behavior of citation relationship and non-globular nature of papers, the final clusters are found to have a less intra-cluster similarity. Due to irregular dimensionality or sparseness relationship among papers, the clusters found are either very large or clusters with less number of papers or sometimes results with singleton clusters.

We encountered a couple of issue,<sup>14</sup> If we try to cluster the entire citation network without applying intermediate network partitioning. To get dense clusters, clusters with varying shapes,

<sup>13</sup> Modularity of a partition is a scalar incentive between  $-1$  and  $1$  that estimates the density of connections inside sub-graphs when contrasted with joins between sub-graphs.

<sup>14</sup> When Jarvis-Patrick algorithm was employed on 32,069 papers shortlisted after centrality measures; few clusters were with the average number of papers more than 700, some with less than 3 papers or even a single paper.



sizes, and densities (either not exactly larger in size nor singleton clusters), and to handle high dimensionality we, therefore, apply network partitioning before applying graph clustering.

---

**Algorithm 1:** Modified Jarvis-Patrick clustering
 

---

**Input:** Observed citation sub-graphs  $S$  with paper-paper connectivity  
**Output:** The algorithm partition input papers into non-hierarchical clusters  
**Initialization:** Let  
 $P = \{p_1, p_2, \dots, p_n\}$  be the set of candidate papers present in sub-graphs  $S$   
 $T$  = User-defined threshold for similarity  
 $F$  = minimum required number of neighbors in common  
**for**  $i \leftarrow 1$  **to**  $|P|$  **do**  
  **for**  $j \leftarrow 1$  **to**  $|P|$  **do**  
    **if**  $(p_i \neq p_j)$  **then**  
      **for**  $k \leftarrow 1$  **to** 11 **do**  
         $c \leftarrow (k-1)*0.1$  /\* param values  $c = \{0, 0.1, \dots, 1\}$  \*/  
         $sim_k(p_i, p_j) \leftarrow \frac{c*S_1 + (1-c)*S_2}{dist(p_i, p_j)}$   
        where,  
         $S_1 \leftarrow abstract\_similarity(p_i, p_j)$  using Eq. (6)  
         $S_2 \leftarrow title\_similarity(p_i, p_j)$  using Eq. (7)  
         $dist(p_i, p_j) \leftarrow$  the minimum hop length between  $p_i$  and  $p_j$   
      **end**  
    **end**  
  **end**  
  **end**  
  **for**  $i \leftarrow 1$  **to**  $|P|$  **do**  
    resultant-set( $p_i$ ) = set of neighbors of  $p_i$   
    =  $\{p_j : sim_k(p_i, p_j) \geq T \text{ for any } k\}$   
  **end**  
  **for**  $i \leftarrow 1$  **to**  $|P|$  **do**  
    **for**  $j \leftarrow 1$  **to**  $|P|$  **do**  
      **if**  $|resultant\_set(p_i) \cap resultant\_set(p_j)| \geq F$  **then**  
        cluster( $p_i$  and  $p_j$ )  
      **end**  
    **end**  
  **end**  
**end**  
**return** identified clusters along with their non-overlapping papers

---

### 8.1.2. Citation network partitioning

We use the Louvain algorithm for graph partitioning [64]. The quality of the partitions is ensured by high modularity scores [65, 66]. This method is chosen over other community detection approaches due to its simplicity, lesser computational time, and better quality of communities (Modularity).

A weighted citation graph  $G = (V, E)$ , where  $i, j \in V$ , an edge  $l(i, j) \in E$  has weight  $w_{i,j}$ . The objective of this step is to partition a citation network  $G$  into a set  $S$  of mutually exclusive and exhaustive sub-graphs  $S_i = (V_i, E_i)$ .

$$\bigcup V_i = V; \quad \forall S_i \in S \quad (9)$$

$$V_i \cap V_j = \emptyset; \quad \forall S_i, S_j \in S \quad (10)$$

The step provides us 293 number of partitions which are almost uniform containing about an equal number of papers.

### 8.1.3. Topic-oriented intra-graph clustering

We consider each partition for further clustering based on link and contextual similarity. A weighted sub-graph  $S_i = (V_i, E_i)$  is divided here into  $n_i$  clusters using Jarvis-Patrick algorithm [67]. The reason behind the selection of Jarvis-Patrick to cluster each sub-graphs found after Louvain algorithm are: It will find tight clusters embedded in loose one. It is mainly good for detecting chain-like or non-globular clusters. The clustering steps are very fast, and the overhead requirement is very low. The capability to find clusters of different shapes, sizes, and densities in high dimensional data.

The objective of this step is to make from each partition coherent clusters of papers that are closely related to each other.

Let  $C_i$  be a set of  $n_i$  number of such clusters for partition  $S_i$ .

$$\bigcup_{j \in \{1, 2, \dots, n_i\}} c_{ij} = C_i \quad (11)$$

$$c_{ij} \cap c_{ik} = \emptyset \quad j \neq k \quad (12)$$

Although Jarvis-Patrick works well in graph clustering it suffers from a problem.<sup>15</sup> It utilizes two parameters: the minimum number of common neighbors ( $F$ ) and the size of the neighbor list ( $T$ ) between a pair of nodes. But these parameters are predefined before applying Jarvis-Patrick and are not generally modified dynamically. Due to these hand-coded or fixed size of the neighbor list ( $T$ ) in citation networks, we are not guaranteed to get clusters with consistent quality. The reason is the non-globular or irregular dimensionality among papers in a citation network.

To address the above issue and to catch a gathering of more similar objects in one cluster, we alter the original Jarvis-Patrick algorithm. A variable-length nearest neighbor list, a proximity threshold is utilized to decide a variable number of neighbors for each paper. All neighbors that pass the similarity threshold are considered as neighbors to this paper. By this alteration, outliers are prevented from joining a cluster while preventing the arbitrary splitting of large clusters is emerging from the limitations imposed by the fixed-length threshold. The detailed steps are given in Algorithm 1. This step provides us 387 number of clusters.

---

**Algorithm 2:** Sub-clusters merging algorithm
 

---

**Input:** Identified sub-clusters along with non-overlapping set of papers  
**Output:** Merging clusters to collect relevant candidate set of papers  
**Initialization:** Let  
 $C = \bigcup_i \{c_{i1}, c_{i2}, \dots, c_{in_i}\}$  be the set of sub-clusters for all the partitions taken together (found after applying Jarvis-Patrick algorithm)  
 $R = \{r_1, r_2, \dots, r_r\}$  be the set of topmost  $r$  similar papers by using Eq. (8)  
 $candidate\_set = \emptyset$   
**for**  $i \leftarrow 1$  **to**  $|R|$  **do**  
  **for**  $j \leftarrow 1$  **to**  $|C|$  **do**  
    **if**  $(r_i \in c_{ij})$  **then**  
       $candidate\_set = candidate\_set \cup c_{ij}$  /\*All papers in  $c_{ij}$  \*/  
    **end**  
     $j \leftarrow j + 1$   
  **end**  
   $i \leftarrow i + 1$   
**end**  
collect the set of identified sub-clusters and merge them  
**return** final  $candidate\_set$

---

### 8.1.4. Abstract similarity using Okapi BM25+ algorithm

Keeping in mind the overall goal to retrieve only conceptually related papers with the seed paper, merging of clusters need to be done before applying abstract similarity. The complete steps are quoted in Algorithm 2. To perform such merging, we need to take after the accompanying rules as given below:

- (i) Select top  $R$  papers considering the cumulative score of abstract and title similarity with seed paper as discussed and examined in Section 7 and Section 8.1.1.
- (ii) Select the topmost similar paper as paper the of interest ( $I$ ) and extract its associated venue as the venue of interest ( $Z$ ).
- (iii) Take individually selected papers ( $R$ ) and identify their corresponding clusters found by the Jarvis-Patrick algorithm.

<sup>15</sup> Between any two papers  $A$  and  $B$ ;  $A$  may have a high number of neighbors while  $B$  having very few due to the fixed size of neighbor lists. Now for the minimum number of common neighbor ( $F$ ) and size of the neighbor list ( $T$ ),  $A$  and  $B$  cannot come to a cluster although they are semantically quite close and related papers in a bibliographic citation network.

**Table 4**  
Research topic distribution of venue  $v_i$ .

Year	Topic <sub>1</sub>	Topic <sub>2</sub>	Topic <sub>3</sub>	Topic <sub>4</sub>	Topic <sub>5</sub>
2008	0.4	0.3	0.2	0	0.1
2009	0	0.3	0.2	0.4	0.1
2010	0.1	0	0.6	0.2	0.1
2011	0.5	0.2	0.2	0	0.1
2012	0.3	0.3	0	0.2	0.2

**Table 5**  
Weighted score of topic distribution of venue  $v_i$ .

Year	Topic <sub>1</sub>	Topic <sub>2</sub>	Topic <sub>3</sub>	Topic <sub>4</sub>	Topic <sub>5</sub>
2008	0.15	0.11	0.07	0	0.03
2009	0	0.12	0.08	0.17	0.04
2010	0.05	0	0.3	0.1	0.05
2011	0.31	0.12	0.12	0	0.06
2012	0.3	0.3	0	0.2	0.2

- (iv) Extract all papers present in those identified clusters (assume  $t_1$ ) and merge them with the selected top papers from set-II (assume  $t_2$ ).

So after getting top R similar paper, merging of clusters is done by using Algorithm 2. In our experiment, we have generally considered 80–120 ( $t_1 + t_2$ ) papers to check the abstract similarity with the seed paper. It has been experimentally observed that there are 65–105 papers ( $t_1$ ) present after the merging of clusters.

To address the deficiency of Okapi BM25 in its term frequency (TF) normalization component, i.e., the TF normalization is not lower bounded properly, in this paper, we adapted Okapi BM25+ (a variant of Okapi BM25) to compute the abstract similarity of  $P_{seed}$ , and  $P_{test}$  papers. It is specifically applied to retrieve only conceptually related papers with seed paper. Okapi BM25+ is based on the probabilistic retrieval framework [68], whose weighting based similarity score can be expressed as follows. Abstract similarity ( $P_{seed}, P_{test}$ ) =

$$\sum_{t \in P_{seed} \cap P_{test}} \ln \left( \frac{P - wf + 0.5}{wf + 0.5} \right) \cdot \left( \frac{(n_1 + 1).cf}{n_1(1 - r + r \frac{wl}{avwl}) + cf} + \delta \right) \cdot \frac{(n_3 + 1).qcf}{n_3 + qcf} \quad (13)$$

where,  $cf$  is the term  $t$ 's frequency in testing paper ( $P_{test}$ ),  $qcf$  is the term's frequency in seed paper ( $P_{seed}$ ),  $P$  is the total number of papers identified ( $t_1 + t_2$ ),  $wf$  is the number of testing papers that hold the term  $t$ ,  $wl$  is the length of abstract (in bytes),  $avwl$  is the average abstract length of papers in each components,  $n_1$  (between 1.0–2.0),  $r$  (usually 0.75),  $n_3$  (between 0–1000), and the value of  $\delta$  is a constant (usually 1.0).

The papers are sorted and ranked in decreasing order of their similarity score with the seed paper. The ranked papers are used to fetch the venues in the same order and suggest user-specified top  $N$  (usually  $N \neq t_1$  or  $t_2$ ) unique venues.

## 8.2. The architecture of VVPN model

We are attempting to discover inherent community structures in a Author-Author Graph (AAG) to understand the network more profoundly and reveal interesting relationships shared among venues. To measure the topic distribution of venues to capture their respective current scope, age-discounted based Venue2Vec is proposed.

The process of VVPN model mainly involves six steps: (i) Venues scope variation with time, (ii) Venue2Vec edge weighting,

(iii) Generation of the venue-venue graph (VVG), (iv) Combining meta-path features, (v) Computing meta-path edge weights as features, and (vi) Recommendation of biased RWR model.

### 8.2.1. Venues scope variation with time

Researchers usually desire to contact those venues which are currently publishing similar research papers. Hence, topic distribution and title embeddings in recent years can describe the current scope of a venue more accurately. Table 4 displays the topic distribution of venue  $v_i$ . To quantify a venue's scope, we initially categorize their publications year-wise to capture the topic distribution of venue using their published papers as depicted in Table 4.

To capture the variation of the scope of venues, we apply LDA based topic modeling on abstract and Doc2Vec on the title of papers published in venues. LDA gives the year wise topic distribution of the venues and Doc2Vec returns a vector for each year based on contextual information from venues published titles. The results from LDA and Doc2Vec can be considered as two sets of vectors.  $L_i^v$  represents the vector of year-wise topic distribution vectors and  $D_i^v$  represents the vector of year-wise title embeddings vectors as depicted in Eq. (14) and in Eq. (15) respectively. The years considered are 2000, 2001, ..., 2012. Each year-wise vector is again a vector of  $k$  different topics as given in Eqs. (19) and (20).

$$L_i^v = [L_{2000i}^v, L_{2001i}^v, \dots, L_{2012i}^v] \quad (14)$$

$$D_i^v = [D_{2000i}^v, D_{2001i}^v, \dots, D_{2012i}^v] \quad (15)$$

Now, we employ a weighted addition of vectors from each set to get one vector for abstract similarity and one vector for title similarity. We use age-discounted scheme (inverse log-weighting scheme) to give more weight to the current year vectors, and the weight reduces in the decreasing order of years. For each venue  $v_i$ , we get a vector  $A_i^v$  for abstract similarity and vector  $T_i^v$  for title similarity as depicted in Eq. (16) and in Eq. (17) respectively.

$$A_i^v = \sum_{y_i \in Y} \frac{L_{y_i}^v}{\log_2(y_o - y_i + 2)}, \text{ and} \quad (16)$$

$$T_i^v = \sum_{y_i \in Y} \frac{D_{y_i}^v}{\log_2(y_o - y_i + 2)} \text{ where} \quad (17)$$

$$Y = \{2000, \dots, 2012\} \text{ and } y_o \text{ is the latest year in } Y. \quad (18)$$

$$L_{y_i}^v = [a_{1i}, a_{2i}, \dots, a_{ki}] \quad (19)$$

$$D_{y_i}^v = [a_{1i}, a_{2i}, \dots, a_{ki}] \quad (20)$$

Using  $A_i^v$  and  $T_i^v$  for a venue  $v_i$ , we compute cosine similarity with their counterpart from the seed paper as discussed in next section Venue2Vec edge weighting.

**Example.** Table 4 shows the initial topic distributions for five topics of venue  $v_i$  and Table 5 shows the topic distribution after age-discounted weighting scheme being applied. Eq. (21) shows the topic distribution vector of venue  $v_i$  in year 2010. The age-discounted vector is given by Eq. (22) (latest year = 2012).

$$A_{2010}^v = [0.1, 0, 0.6, 0.2, 0.1] \quad (21)$$

$$\frac{A_{2010}^v}{\log_2(4)} = [0.05, 0, 0.3, 0.1, 0.05] \quad (22)$$

Furthermore, we adopt a weighted addition of vectors to obtain the final vector, as given in Table 5. The final vector  $A_i$  for venue  $v_i$  after weighted addition will be:

$$A_i^v = [0.81, 0.65, 0.57, 0.47, 0.38] \quad (23)$$

**Table 6**  
Meta-paths used in VVPN model.

No.	Meta-path	Description
1.	<i>common_author</i>	Core venues share an author
2.	<i>common_term</i>	Core venues share a term
3.	<i>direct_cites</i>	Core venue cites core venue
4.	<i>direct_cited_by</i>	Core venues cited by core venues
5.	<i>citation_paper</i>	Core venues share a reference (ref)
6.	<i>co_citation_paper</i>	Core venues co-cited together (cite)

If we had applied a simple vector addition without any weights, we would have got a vector  $A_i^{v'}$  as:

$$A_i^{v'} = [1.3, 1.1, 1.2, 0.8, 0.6] \quad (24)$$

We can clearly see the difference between  $A_i^v$  and  $A_i^{v'}$ . It clearly indicates the influence of topic distribution vector of recent year 2012 in the calculation of  $A_i^v$  where as in  $A_i^{v'}$ , all the year wise vectors contribute equally. Furthermore, venue-venue similarity is done among venues exploiting their corresponding weighted vector  $A_i^v$  and  $T_i^v$  respectively.

### 8.2.2. Venue2Vec edge weighting

Using  $A_i$  and  $T_i$  for a venue  $v_i$ , we compute cosine similarity between any two venues. We get two cosine similarities,  $Sim_a(v_i, v_j)$  and  $Sim_t(v_i, v_j)$ , for a pair of venues,  $v_i$  and  $v_j$ , using  $(A_i, A_j)$  and  $(T_i, T_j)$  respectively.

$$Sim_a(v_i, v_j) = \frac{A_i^v \cdot A_j^v}{|A_i^v| |A_j^v|} = \frac{\sum_{b=1}^k (a_{b,i} * a_{b,j})}{\sqrt{\sum_{b=1}^k a_{b,i}^2} * \sqrt{\sum_{b=1}^k a_{b,j}^2}} \quad (25)$$

$$Sim_t(v_i, v_j) = \frac{T_i^v \cdot T_j^v}{|T_i^v| |T_j^v|} = \frac{\sum_{b=1}^k (t_{b,i} * t_{b,j})}{\sqrt{\sum_{b=1}^k t_{b,i}^2} * \sqrt{\sum_{b=1}^k t_{b,j}^2}} \quad (26)$$

Now we utilize these two similarity metrics to get one final metric,  $Sim(v_i, v_j)$  with the help of an adjustment parameter  $m$  as:

$$Sim(v_i, v_j) = m * Sim_a(v_i, v_j) + (1 - m) * Sim_t(v_i, v_j) \quad (27)$$

where  $m \in [0, 1]$ .

We consider this similarity score as contextual similarity features (CSF). We are using this CSF score in Section 8.2.3 to generate a weighted VVG (venue-venue) graph and also to compute the edge-weight among venues.

### 8.2.3. Generation of venue-venue graph (VVG)

In this section, we will create a homogeneous undirected venue-venue graph (VVG) from the HIN graph to recommend relevant venues to the input seed paper. We define this graph as an undirected graph,  $VVG = (B, D)$  with a vertex type mapping function  $\omega: B \rightarrow B$  and an edge type mapping function  $\pi: D \rightarrow D$ . Here, we have one type of vertex  $B$  for each venue.

$$B = \{\text{set of venues where only } P\_main \text{ papers published}\} \quad (28)$$

The type of edge  $D$  is defined as  $b_1 \xrightarrow{\text{connects}} b_2: \omega(b_1) \in \{P\_main\}, \omega(b_2) \in \{P\_main\}, b_1, b_2 \in B$ .

It joins two venues using only one type of edge such that  $b_1 \xrightarrow{d_1} b_2$ , where  $\pi(d_1) \in D$ . Table 6 lists all types of meta-paths defined in our model. We are extracting the venue of  $P\_main$  and considering as a core venue to maintain a homogeneous VVG graph. Initially, the CSF score as computed in Section 8.2.2 among venues is used to create the VVG graph. The average CSF score is used as a threshold to create the edge between venues. No edge exists with less than average CSF score found among venues.

### 8.2.4. Combining meta-path features into VVG

Since meta-paths are mostly composite relations of various edge types in a HIN graph, they can capture the distinct relationship between a pair of HIN vertices [55]. We assume that a meta-path connects two different  $P\_main$  papers  $x, y$  that belong to two disjoint core venues  $v_i$  and  $v_j$  respectively.

We observed that meta-path features with more than two degree<sup>16</sup> are not much meaningful in our work and even not able to create much difference to compute the similarity among venues. To reduce the time complexity and to obtain a tightly coupled relationship among venues, only one-degree and two-degree meta-path features are incorporated into this VVPN model, and a homogeneous VVG graph is exploited to recommend academic venues. We believe that research papers that share many similar references may use a common set of background knowledge. By using this hypothesis, this information could be used to compute the possible associations among papers.

### 8.2.5. Computing meta-path edge weights as features

To discover the latent association between venues, we have divided the above six meta-paths as depicted in Table 6 into 3 categories of edge weighting.

- (i) *Common\_Features* (CF): Common author and common term meta-path belong to this category. Common author similarity and common term similarity between two venues  $v_i$  and  $v_j$  are represented by  $Sim_A(v_i, v_j)$  and  $Sim_T(v_i, v_j)$  respectively. Term appearing in titles or abstracts of a  $P\_main$  paper after stop word removal and stemming are consider for similarity computation. We use snowball stemmer to get the root words [69]. Jaccard similarity coefficient is used to calculate both  $Sim_A(v_i, v_j)$  and  $Sim_T(v_i, v_j)$  Eq. (29). In case of computation of  $Sim_A(v_i, v_j)$ , sets  $E$  and  $F$  denote list of authors associated with venue  $v_i$  and  $v_j$  respectively.

$$J(E, F) = \frac{|E \cap F|}{|E \cup F|} \quad (29)$$

where  $0 \leq J(E, F) \leq 1$ .

Similarly during  $Sim_T(v_i, v_j)$  computation, sets  $E$  and  $F$  denote sample terms occur in venue  $v_i$  and  $v_j$  respectively [70]. Then we are combining the above two similarity scores to obtain CF score (Common\_Features) between two venues  $v_i$  and  $v_j$  respectively. The computation of CF edge weighting between  $v_i$  and  $v_j$  is defined below.

$$CF(v_i, v_j) = Sim_A(v_i, v_j) + Sim_T(v_i, v_j) \quad (30)$$

Generally, none of the CF similarity scores among two venues will get a perfect score of 1, and also random walk is sensitive to a higher probability score. Normalization of data within a uniform range (e.g., (0–1)) is essential to prevent larger applies to the output variables. This representation numbers from overriding smaller ones. One way is to scale input and output variables ( $z$ ) in the interval  $[\rho_1, \rho_1]$  corresponding to the range of the transfer function [71]. Before adding this meta-path CF score into the model, we are individually applying the normalization to be in the range of [0.1–0.9] as shown in Eq. (31).

$$z_i = \rho_1 + (\rho_2 - \rho_1) \frac{(x_i - x_i^{min})}{(x_i^{max} - x_i^{min})} \quad (31)$$

After applying this normalization, we will get a normalized CF score  $CF'(v_i, v_j)$  among two venues  $v_i$  and  $v_j$ .

<sup>16</sup> The degree of a meta-path indicates its length and the distance between two main papers.

- (ii) *Direct – Citation\_Features* (DCF): The meta-paths such as *direct\_cites* and *direct-cited-by* are included in this group. The computation of edge weighting of DCF is defined below.

$$DCF(v_i, v_j) = |P_{ij}| + |P_{ji}| \quad (32)$$

Where  $P_{ij}$  denotes set of papers published at venue  $v_i$  and referring to papers published at venue  $v_j$ . After applying the normalization defined in Eq. (31), we will get a normalized DCF score  $DCF'(v_i, v_j)$  among two venues  $v_i$  and  $v_j$ .

- (iii) *Co – Citation\_Features* (CCF): The remaining meta-paths such as *citation\_paper* and *co-citation\_paper* are within this group. The computation of edge weighting of DCF is defined below.

$$CCF(v_i, v_j) = \sum_{\substack{k \neq i, \\ k \neq j}} |P_{ik} \cap P_{jk}| + \sum_{\substack{k \neq i, \\ k \neq j}} |P_{ki} \cap P_{kj}| \quad (33)$$

where  $P_{ik}$  is the set of papers published at venue  $v_i$  and referring to papers published at venue  $v_k$ . After applying the normalization defined in Eq. (31), we will get a normalized CCF score  $CCF'(v_i, v_j)$  among two venues  $v_i$  and  $v_j$ .

We add each normalized meta-path scores into the model to analyze their effect on the recommendation quality. We already have initial edge weighting score CSF, which is computed based on the age-discounted scheme (inverse log weighting scheme) based abstract and title similarity as calculated in Section 8.2.2. It was purely based on the contextual similarity to be in the range of (0–1). So after applying normalization defined in Eq. (31), we will get a normalized CF score  $CSF'(v_i, v_j)$  among two venues  $v_i$  and  $v_j$ .

$$CSF'(v_i, v_j) = Sim(v_i, v_j) \quad (34)$$

Initially, the recommendation will be provided based on the normalized  $CSF'$  matching score.

$$CWS(v_i, v_j) = CSF'(v_i, v_j) \quad (35)$$

We need to combine individual normalized meta-path scores into the model, and we call it a combined weighted score (CWS). In addition to normalized CSF score all normalized scores obtained from Eqs. (30), (32) and (33) are added to obtain the  $CWS(v_i, v_j)$  to increase the probability of recommending relevant venues during recommendation. The CWS score can be used as a probability score between venues in VVG graph as computed using Eq. (38) to apply random walk with restart (RWR).

$$CWS(v_i, v_j) = CSF'(v_i, v_j) + CF'(v_i, v_j) + DCF'(v_i, v_j) + CCF'(v_i, v_j) \quad (36)$$

## 9. Fusion model: CNAVER (Layer-4)

To be more specific, the predictions resulting from the PPPN model and VVPN model are first produced separately, allowing us to leverage the individual strengths of both approaches since there is no interdependency between them.

### 9.1. Top venues recommendation (PPPN model)

We apply LDA on abstract and Doc2Vec on the title for Set-II papers (Section 6) and the top  $t_2$  similar papers are chosen. Abstract and title similarity is computed as discussed in Section 8.1.1. We have four assumptions regarding the inclusion of these  $t_2$  papers obtained from Set-II paper for abstract similarity.

- There may be few papers which are recently got published without having any citations (Set-II), may be involved with many reputed venues.
- The seed paper's title and keywords are matching with some papers in Set-II so there is a possibility that the seed paper may get accepted at similar venues as that of Set-II papers.
- Generally the papers published in reputed venues get a high number of citations. Chances of getting acceptance in a new venue are relatively easier than reputed venues.
- New venues should get an equal chance of inclusion in the final recommendation to reasonable address the new venue cold-start issue.

### 9.2. Top venues recommendation (VVPN model)

To exploit collaboration network information along with publication content, we employ a popular network-based approach known as a random walk with restart (RWR). RWR provides an excellent way to measure how closely related two nodes are in a graph [72]. The core equation of the RWR model is shown in Eq. (37).

$$R^{(t+1)} = \alpha S R^{(t)} + (1 - \alpha) Q \quad (37)$$

where  $S$  is the transfer matrix, representing the probability for each node to jump to other nodes.  $R^{(t)}$  is the rank score vector at step  $t$  and  $Q$  is the initial vector of the form  $(0, \dots, 1, \dots, 0)$ . Initially, the rank score of the target node is 1, while others are 0.  $\alpha$  is the damping coefficient. With probability  $(1 - \alpha)$ , walker restarts from the start node. We use the transfer matrix  $S$  to bias our walker's behavior.

We use the weighted combined score (CWS) found after aggregating various meta-paths features in Eq. (36), to bias the walker towards nodes with a higher content as well as semantical similarity. Edge weight  $w_{v_i, v_j}$  for an edge from  $v_i$  to  $v_j$  is given by the equation below:

$$w_{v_i, v_j} = \frac{CWS(v_i, v_j)}{\sum_{x \in N(v_i)} CWS(v_i, x)} \quad (38)$$

where  $N(v_i)$  is set of nodes which have incoming links from  $v_i$ . RWR is an iterative process. After certain iterations,  $R^{(t)}$  converges to a steady-state probability vector. We use  $R^{(t+1)}$  venue-rank score vector to give our final top N recommendation.

---

#### Algorithm 3: Fusion of PPPN and VVPN models

---

**Input:** shortlisted papers after Okapi BM25+ ( $t_1$ ) and shortlisted Set-II papers ( $t_2$ ),  
Venue of interest ( $Z$ ) for a given seed paper  $p_m$

**Output:** Top N recommended list of venues for  $p_m$

**Initialization:** Let

$T = t_1 + t_2$  be the set of candidate papers

$\mathcal{L}$  = Ordered list of unique venues from top-ranked papers based on abstract similarity scores (Section 8.1.4)

$\mathcal{A} = \{a_1, a_2, \dots, a_N\}$

Borda Count  $B_c(a_i) \leftarrow N - i + 1$

$\mathcal{M}$  = Ordered list of unique venues in decreasing order (Section 9.2)

$\mathcal{B} = \{b_1, b_2, \dots, b_N\}$

Borda Count  $B_c(b_i) \leftarrow N - i + 1$

$\mathcal{V}$  = Final list of unique venues

$\mathcal{V} = \{v_1, v_2, \dots, v_{|\mathcal{V}|}\}$  where  $|\mathcal{V}| \leq 2N$

```

for  $i \leftarrow 0$  to  $|\mathcal{L}|-1$  do
  for  $j \leftarrow 0$  to  $|\mathcal{M}|-1$  do
    if ( $a_i == b_j$ ) then
      Borda Count  $B_c(v_i) \leftarrow B_c(a_i) + B_c(b_j)$  /*they are same venue*/
    else
      individually consider Borda Count  $B_c(v_i) \leftarrow B_c(a_i)$  and  $B_c(v_j) \leftarrow B_c(b_j)$ 
    end
  end
end

```

Sort venues in the decreasing order of Borda Count ( $B_c(v_i)$ )

Prepare the final list of top N venues recommendation

---



### 9.3. Final venues recommendation (Fusion model)

Although the social network analysis (SNA), content-based filtering (CBF) and random walk with restart (RWR) are widely used for making venue recommendations, they may not provide the best recommendation results due to their limitations. After getting the individual top N recommendations from both the PPPN model and VVPN model, we need to apply some rank-based fusion because the fusion can provide better recommendation than a single approach and the disadvantages of one approach can be overcome by the other.

Fusion has been widely investigated in the recommendation community. They were often divided into two categories: score-based and ranking-based. Score-based combination methods require similarity information to conduct ranking list aggregation, such as CombSum, CombMNZ, and weight combination [73,74]. Ranking-based combination methods need rank or position information to integrate different candidate's ranking lists, such as Borda fusion, Condorcet fusion, and MAPFuse [75]. In this research, the Borda fusion technique is applied to incorporate the existing prediction lists generated by the PPPN model, and the VVPN model as PPPN provides scores for each venue while VVPN provides ranks of them [76]. The complete steps are quoted in Algorithm 3.

## 10. Experiments

In this section, we present the experiments of the proposed fusion model "CNAVER" to evaluate the effectiveness of it. In this section, we present the experimental datasets, evaluation strategy, evaluation metrics, experimental setting, parameter tuning, and baseline methods. All experiments are performed on a laptop with 64 bit Windows 10 operating system, Intel i7-3540M, CPU@3.00 GHz, and 8 GB memory. All the programs are implemented in python.

### 10.1. Dataset used

We use a real-world dataset DBLP-citation-network V10,<sup>17</sup> the citation data extracted from DBLP, ACM, MAG (Microsoft academic graph), and other sources [77] to demonstrate the effectiveness of our proposed method. The tenth version contains 3,079,007 papers and 25,166,994 citations. Each paper is associated with abstract, authors, title, publishing year, venue, and references list. After removing duplicate papers, papers with missing fields, and inconsistent entries in the database, we are left with 2,236,968 papers.

We performed our experiments on a subset of the dataset (data collected in between the year 2000–2017) from Tang et al. [77]. Due to hardware constraints, only a subset of the original dataset is used for the experimentation. Since papers coming from various fields not only have varied research interests but also may have interdisciplinary collaborations resulting in a diverse number of published venues. Moreover, it is common to practice with sampling with many studies. In this paper, we, divided the dataset into two parts according to the year of publication: data during the years 2000–2012 as the training set, and the rest as a testing set.

### 10.2. Evaluation strategy

We adopt the following two kinds of evaluation to measure the performances of CNAVER against other state-of-the-art methods.

- (a) *Coarse-level or offline evaluation*: As the name suggests, it provides some raw-level quick notion of how the proposed

CNAVER fares vis-a-vis other systems. We focus on the prediction accuracy to see whether the original publication venue for the test paper is predicted or not, and if yes, at what rank within some top N recommendations. Accuracy, MRR, and  $F - measure_{macro}$  evaluation metrics are used during the evaluation (detailed below). We call this scenario *offline* because we can evaluate a system this way only when we have annotated test data.

- (b) *Fine-level or online evaluation*: This evaluation-scenario is more realistic (and, that is why we call *online*) as a researcher needs to have more than one venue recommendation from a system for her paper-in-writing that she wants to communicate. Here we go a little deeper and aim to see the relevance, usefulness, and quality of the recommended results. The system recommends an ordered list of venues that are assessed by experts in terms of graded relevance Eq. (52). Precision, nDCG, and average venue quality are used as evaluation metrics in the evaluation.

### 10.3. Evaluation metrics

We employed eight metrics such as accuracy, MRR,  $F - measure_{macro}$ , precision@k, nDCG@k, average venue-quality (Ave-quality), diversity, and stability to evaluate the performance of CNAVER. Detailed information about these metrics has been discussed [78].

- (a) Accuracy@N: It is the ratio of no. of times our system correctly predicts the real publication venue within top N recommendations for a set of papers [20,36]. Here N varies among 3, 6, 9, 12 and 15 respectively.

$$Accuracy@N = \frac{\# \text{ correctly predicts venues within top } N}{\text{Total number of test papers}} \quad (39)$$

- (b) Mean Reciprocal Rank (MRR): MRR is the arithmetic mean of reciprocal rank (RR) which is the inverse of the first rank where the correct venue is recommended in the ranked result [79].

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_{rel_i}} \quad (40)$$

where  $rank_{rel_i}$  denotes the rank position of the first relevant document for the  $i$ th query in a query set Q.

- (c)  $F - measure_{macro}$  ( $F_1$ ): The macro-average is the average of the same measures calculated for all classes. It treats all classes equally. For an individual class  $C_i$  (number of venues), the assessment is defined by  $tp_i$  (true positives),  $tn_i$  (true negatives),  $fp_i$  (false positives), and  $fn_i$  (false negatives) [80,81].

$$Precision_{macro} = \frac{\sum_{i=1}^N \frac{tp_i}{tp_i + fp_i}}{N} \quad (41)$$

$$Recall_{macro} = \frac{\sum_{i=1}^N \frac{tp_i}{tp_i + fn_i}}{N} \quad (42)$$

$$F - measure_{macro} = \frac{2Precision_{macro} \times Recall_{macro}}{Precision_{macro} + Recall_{macro}} \quad (43)$$

- (d) Precision: Precision is the fraction of retrieved items that are relevant. In our context, it is the fraction of recommended venues that are relevant, as shown in Eq. (44).

$$Precision = \frac{|\text{relevant venues} \cap \text{recommended venues}|}{\text{total recommended venues}} \quad (44)$$

Precision@k means when k venues are recommended, i.e.,

$$Precision@k = \frac{|\text{relevant venues} \cap \text{recommended venues}|}{k} \quad (45)$$

<sup>17</sup> <https://aminer.org/citation>

- (e) Normalized discounted cumulative gain (nDCG): It represents the ratio of discounted system gain and discounted ideal gain accumulated at a particular rank  $p$ , where gain at a rank  $p$  is the sum of relevance values from rank 1 to rank  $p$  [79]. Relevance value in our system ( $rel_{sj}$ ) is a score (0, 1 or 2) assigned by a researcher to the venue at position  $j$ . Ideal vector is constructed hypothetically where all relevance scores ( $rel_{ij}$ ) are ordered in decreasing order to ensure the highest gain at any rank.

$$DCG_{sp} = rel_{s1} + \sum_{j=2}^p \frac{rel_{sj}}{\log_2(j)} \quad (46)$$

$$IDCG_p = rel_{i1} + \sum_{j=2}^p \frac{rel_{ij}}{\log_2(j)} \quad (47)$$

$$nDCG_p = \frac{DCG_p}{IDCG_p} \quad (48)$$

- (f) Diversity (D): It is defined as the average dissimilarity (opposite of similarity) between all pairs of items in a result set [82].

$$D = 2 * \frac{\sum_{i=1}^N \sum_{j=1}^N (1 - \text{Similarity}(v_i, v_j))}{N(N-1)} \quad (49)$$

where  $N$  is the length of the recommendation,  $v_i$  and  $v_j$  are the venues appearing in the recommendation lists and  $\text{Similarity}(v_i, v_j)$  denotes the content (abstract, keywords) similarity among venues  $v_i$  and  $v_j$ .

- (g) Stability: A recommender system is stable if the predictions do not change strongly over a short period of time [34,83]. It is also called mean absolute shift (MAS), designed to capture the internal consistency among predictions made by a given recommendation algorithm [3]. It is also defined through a set of training data  $R_1$  and a set of prediction (ranking of original venue) of seed paper,  $P_1$ . For an interval of time (addition of new data into the training data), the recommender system can now make prediction,  $P_2$ . MAS is defined as

$$\text{Stability} = \text{MAS} = \frac{1}{|P_2|} \sum_{(u,i) \in P_2} |P_2(u, i) - P_1(u, i)| \quad (50)$$

where  $P_1, P_2$  are the predictions made in phase 1 and phase 2, respectively.

- (h) Average-Venue Quality (Ave-quality): It evaluates the quality of the venues recommended by CNAVER based on Google's h5-index [22].

$$\text{Average-venue quality} = \frac{\sum_{v \in V} H5_v}{|V|} \quad (51)$$

where  $V$  is the set of recommended venues and  $H5_v$  is the h5-index of venue  $v$ . Higher the Ave-quality, we can claim, the better is the recommendation.

#### 10.4. Experimental setting

While preparing the test dataset, we consider two scenarios. Firstly, due to operational constraints, 20 sub-domains of computer science were selected as a testing dataset in our experiment. A total of 120 seed papers (6 from each sub-domains) are chosen manually from 20 sub-domains: information retrieval (IR), image processing (IP), security (SC), wireless sensor network (WSN), machine learning (ML), software engineering (SE), computer vision (CV), artificial intelligence (AI), data mining (DM), theory of computation (TC), databases (DB), human-computer interaction (HCI), algorithms and theory (AT), natural language

processing (NLP), parallel and distributed systems (PDS), world Wide Web (WWW), web semantics (WS), computer architecture (CO), compiler design (CD) and multimedia (MM).

Secondly, while identifying seed papers following conditions are taken into consideration to measure the effectiveness of CNAVER to handle cold start issues like a new venue and new researcher.

- (i) *Category 1* ( $2 \leq v_c < 8$ ): Select papers whose associated venues have publications greater than or equal to 2 but less than 8.
- (ii) *Category 2* ( $8 \leq v_c < 15$ ): Select papers whose associated venues have publications greater than or equal to 8 but less than 15.
- (iii) *Category 3* ( $15 \leq v_c$ ): Select papers whose associated venues have publications greater than or equal to 15.
- (iv) *Category 4* ( $2 \leq p_c < 8$ ): Select papers whose associated authors have publications greater than or equal to 2 but less than 8.
- (v) *Category 5* ( $8 \leq p_c < 15$ ): Select papers whose associated authors have publications greater than or equal to 8 but less than 15.
- (vi) *Category 6* ( $15 \leq p_c$ ): Select papers whose associated authors have publications greater than or equal to 15.

There are two major categories, i.e., venue count ( $v_c$ ) and publication count ( $p_c$ ). Generally  $v_c$  denotes the number of published papers of individual venue and  $p_c$  denotes the number of publications of a researcher. It is ensured that each category is well represented in the seed papers.

##### 10.4.1. Procedure of online evaluation

For this evaluation, we did not have the ready annotation, but we need one. The annotation or relevance assessment is collected from the volunteers through crowdsourcing in the best effort basis. There are 57 researchers with expertise in the subjects of the papers provided with input and output of our recommender system where for each paper, 15 venues recommended. Out of 57 researchers, 23 evaluated 3 papers each, 17 researchers evaluated 2 each and the rest 17 were evaluated by 17 researchers.

All the experts were identified from academia with a minimum of 3 years of research experience. Most were having a Ph.D. except few research students and research assistants who were pursuing a Ph.D. with bachelors' or masters' degree in science or technology. The experts or researchers were so chosen that their active areas of research perfectly match with the topics of seed papers. Among 57 researchers, there were 8 professors, 11 associate professors, 19 assistant professors, 12 senior research students, and the remaining 7 were research assistants.

The experts check the titles, abstracts, authors, year of publication, and venue recommendations of the papers and determine the relevance-level of the recommendations. In this experiment the relevance value  $r$  is ternary, i.e.,  $r \in \{0, 1, 2\}$ .

$$\text{Relevance}(r) = \begin{cases} 2 & \text{perfectly matching} \\ 1 & \text{partial matching} \\ 0 & \text{otherwise} \end{cases} \quad (52)$$

It is set to 2 if the expert agrees that the research paper is completely matching with the scope of the journal, set to 1 if there is a partial matching or set to 0 otherwise. But while computing precision, we have assumed the partial relevance as not relevant, i.e., relevance 1 is substituted with a relevance value of 0.

To comprehensively evaluate our proposed method and more specifically, to address the research questions (RQs) discussed in Section 1, we prefer to examine the following sub-queries SQs:

**Table 7**

Experimental parameter settings.

Parameter	Range	Default
Vector dimension ( $k$ )	(10–200)	100
Adjustment parameter ( $c$ and $m$ )	(0.1–0.95)	0.7
Similarity threshold ( $T$ )	(0.2–0.55)	0.35
Number of neighbor ( $F$ )	(5–50)	10
Top similar paper ( $R$ )	(5–20)	10
Number of Set-II papers ( $t_2$ )	(5–45)	15
Damping constant ( $\alpha$ )	(0.1–0.95)	0.65
Number of recommended nodes	(5–50)	15

**SQ1:** How effective is CNAVER in comparison to other state-of-the-art methods?

**SQ2:** How is the quality of venues recommended by CNAVER as compared to other state-of-the-art methods?

**SQ3:** How does CNAVER handle cold-start issues and other issues like data sparsity, diversity, and stability?

### 10.5. Parameter tuning and optimization

In this section, we demonstrate the impact of various experimental parameter settings including dimensions of vectors ( $k$ ) for  $A_i$  and  $T_i$  calculations, adjustment parameter ( $c$ ), threshold ( $T$ ), minimum number of neighbor ( $F$ ), top similar papers ( $R$ ), number of Set-II papers ( $t_2$ ) to perform PPPN recommendation and adjustment parameter ( $m$ ), and damping constant ( $\alpha$ ) to perform VVPN model respectively.

The ranges and default values of the parameters are depicted in Table 7. When the effect of the parameter is under examination, the other parameters are set to default values. These experimentations are performed in the training phase contains known output, and the model learns on this data in order to be generalized to other data later on. The ranges of values of various parameters for which the model achieves higher performance are identified as optimal parameters. During this experimentations, the best results are marked by the 'bold-face' in each position.

#### 10.5.1. Influence of vector dimension ( $k$ )

In order to find the ideal dimension ( $k$  = no. of topics) for LDA, we conduct experiments on four values for vector dimension, i.e., {10, 50, 100, 200}. To find the ideal dimension for vectors  $A_i$  and  $T_i$ , the value of the adjustment parameter is set to be 0.7, and  $\alpha$  is set to be 0.65. We extracted the venues of identified seed papers and selected them as a target node to run the VVPN model respectively, then, observed the average performance of the VVPN model in terms of MRR upon various categories. We repetitively performed such experiments with varying recommendation lists in length to evaluate the influence of the vector dimension on the results. We conduct experiments on four values for vector dimension ( $k$ ), i.e. {10, 50, 100, 200}. It is observed that the model performs best when the value of the vector dimension is 100.

Table 8 represents the performance of the VVPN model on various vector dimensions. As we can see, MRR score keeps on increasing and behaving a consistent performance while incrementing of vector dimension. From the whole point of view, the model performs a little better with vector dimension 100. Although with vector dimension 200 its results with the best performance ever. There is no significant improvement in MRR while changing the size of  $k$  from 100 to 200. As we know, it is computationally costly as compared to vector dimension 100. So we have considered the vector dimension ( $k$ ) as 100 in our experiment.

#### 10.5.2. Influence of adjustment parameter ( $c$ and $m$ )

In order to find the ideal value of  $m$  to get the efficient combined score of vectors  $A_i$  and  $T_i$ , we conduct experiments

on 10 possible values for adjustment parameter, i.e. {0.5, 0.45, 0.4, 0.35, 0.3, 0.25, 0.2, 0.15, 0.1, 0.05}. The value of the vector dimension is set to 100, and  $\alpha$  is set to be 0.65. We have followed a similar procedure, as explained in Section 10.5.1.

In Table 9, we can observe that the variation tendency of MRR score performs roughly consistent. We can see that the MRR shows a downward trend with the decreasing value of adjustment parameter  $1 - m$ . The model performs the best while the value of the adjustment parameter is 0.3. This is due to the case that, in most of the cases, the abstract is giving a better clarity of topic similarity while in some instances, the title resulting better. So considering a similar nature, in this experiment, the value of  $(1 - m)$  and  $(1 - c)$  has been taken as 0.3.

#### 10.5.3. Influence of $T$ , and $F$ in intra-graph clustering

Similarly, during Jarvis-Patrick, by varying the value of  $T$ , we observed the effect on sub-clusters size and similarity among papers belong to each sub-clusters. We investigate the performance of sub-clusters found after varying threshold  $T$  as 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5 and 0.55, and  $F$  as 5, 7, 10, 12, 15, 20, 30, and 50 while performing intra-graph clustering. We observed that, while considering  $T$  as 0.5 or more than that, it results in a larger number of clusters with less number of papers in each sub-clusters. We also saw that due to high  $T$ , chances of forming singleton clusters are more. Due to the low value of  $T$ , there is a high chance of forming less number of clusters with a larger number of papers in each cluster.

The average similarity among papers in each cluster is less due to their loosely coupled relationship. A similar pattern is shown by the value  $F$  during intra-graph clustering. We observed that with the value of  $T$  as 0.35, resulting in desired clusters with high average similarity among papers and also showing a tightly coupled relationship. The modularity value observed as 0.69. The best result is found with a value of  $F$  as 10 during the experimentation.

#### 10.5.4. Influence of $R$ in recommendation of PPPN

We first examine whether increasing the number of papers can produce desired recommendation performance. We gradually changed the value of  $R$  as 5, 8, 10, 12, 15, and 20, respectively. We observed that, after 10 papers, no such changes are occurring in the recommendation order. This is because most similar papers occur in the list of top 10 and papers that are strongly coupled to those 10 papers are also exhibit a high contextual similarity as well as tightly coupled semantic relationship. After applying Jarvis-Patrick, we observed that 10 papers are sufficient to capture most similar papers to recommend appropriate venues as discussed in Algorithm 2. We found that there are no such changes on the final recommendation after increasing the value of  $R$ . So finally, the value of  $R$  is considered as 10 during the experimentation.

#### 10.5.5. Influence of $t_2$ in recommendation of PPPN

We have also experimentally tested the effect of the number of papers ( $t_2$ ) selected from Set-II to perform abstract similarity. We first examine whether increasing the number of papers can produce desired recommendation performance. We gradually changed the value of  $t_2$  from 5 to 45 and noticed that, after 15 papers, there are no such changes occurring in the recommendation order. The upper limit is taken as 45 to offer equal opportunity to Set-II path along with Set-I path (papers found after intra-graph clustering) as on an average the intra-graph clustering results in 45–85 number of papers. Our proposed model is assumed to recommend a maximum of 15 venues.

#### 10.5.6. Influence of damping constant ( $\alpha$ )

Also, we measured the performance of CNAVER on damping constant  $\alpha$ . The damping constant  $\alpha$  is an important parameter in RWR. In order to find the ideal value of  $\alpha$  to perform the

**Table 8**  
Influence of vector dimension ( $k$ ) on MRR.

Topic dimension	MRR					
	$2 \leq v_c < 8$	$8 \leq v_c < 15$	$15 \leq v_c$	$2 \leq p_c < 8$	$8 \leq p_c < 15$	$15 \leq p_c$
10	0.0573	0.0881	0.0923	0.0495	0.0761	0.0982
50	0.0619	0.0893	0.1044	0.0562	0.0892	0.1016
100	0.0932	<b>0.0993</b>	0.1097	0.0838	0.1038	<b>0.1134</b>
200	<b>0.0946</b>	0.0989	<b>0.1104</b>	<b>0.0866</b>	<b>0.1039</b>	0.1128

**Table 9**  
Influence of adjustment parameter ( $m$ ) on MRR.

Adjustment prob.( $1 - m$ )	MRR					
	$2 \leq v_c < 8$	$8 \leq v_c < 15$	$15 \leq v_c$	$2 \leq p_c < 8$	$8 \leq p_c < 15$	$15 \leq p_c$
0.5	0.0793	0.0849	0.0853	0.0854	0.0861	0.0864
0.45	0.0798	0.0879	0.0851	0.0853	0.0893	0.0847
0.4	0.0867	0.0905	0.0893	0.0915	0.0841	0.0859
0.35	0.0972	0.0895	0.0949	0.0858	0.0903	0.0885
0.3	<b>0.1093</b>	<b>0.1197</b>	<b>0.1267</b>	<b>0.1134</b>	<b>0.1127</b>	<b>0.1185</b>
0.25	0.0972	0.1014	0.1258	0.1016	0.1039	0.1073
0.2	0.0668	0.0848	0.1132	0.0894	0.0917	0.0995
0.15	0.0526	0.0773	0.0866	0.0739	0.0877	0.0914
0.1	0.0473	0.0637	0.0725	0.0683	0.0746	0.0828
0.05	0.0437	0.0591	0.0683	0.0565	0.0677	0.0769

**Table 10**  
Influence of restart probability on MRR.

Restart prob. ( $1 - \alpha$ )	MRR					
	$2 \leq v_c < 8$	$8 \leq v_c < 15$	$15 \leq v_c$	$2 \leq p_c$	$8 \leq p_c < 15$	$15 \leq p_c$
0.5	0.0633	0.0729	0.1134	0.0635	0.0862	0.1067
0.45	0.0765	0.0914	0.1048	0.0714	0.0975	0.1095
0.4	0.0933	0.0975	0.1086	0.0837	0.1037	0.1135
0.35	<b>0.1357</b>	<b>0.1432</b>	<b>0.1791</b>	<b>0.1137</b>	<b>0.1174</b>	<b>0.1248</b>
0.3	0.1141	0.1265	0.1464	0.1089	0.1146	0.1195
0.25	0.0973	0.1012	0.1296	0.1019	0.1034	0.1078
0.2	0.0763	0.0847	0.1134	0.0896	0.0917	0.0995
0.15	0.0525	0.0772	0.0865	0.0734	0.0877	0.0918
0.1	0.0472	0.0636	0.0724	0.0687	0.0745	0.0823
0.05	0.0432	0.0594	0.0688	0.0563	0.0671	0.0766

random walk on venue-venue graph, we conduct experiments on ten possible values for damping constant, i.e. {0.5, 0.45, 0.4, 0.35, 0.3, 0.25, 0.2, 0.15, 0.1, 0.05}. The value of the vector dimension is set to 100, and the adjustment parameter is set to be 0.7. With higher values of  $\alpha$ , the probability of random walker reaching far away as the number of nodes increases. Hence, the chances of getting new venues will be more, but it may result in irrelevant venues.

It is evident from Table 10 that there is a drastic increase in MRR while decreasing the probability ( $1 - \alpha$ ) till 0.35. Afterward, it exhibits a downtrend with the decreasing value of damping constant. The MRR score performs the upper convex curve, rapidly rising with the value of ( $1 - \alpha$ ) as 0.35 and then shows a decline and downtrend in performance. So based on the above statistics, we have considered the value of damping constant ( $1 - \alpha$ ) as 0.35 in the rest of the experiment.

#### 10.6. Baseline methods

To measure the effectiveness of the proposed venue recommendation, we compare our results with eight state-of-the-art methods. Detailed settings of the systems are presented below.

- (a) Friend based model (FB): The basic idea of the friend-based model is to recommend venues as per the number of neighbors such as a researcher's co-author and co-author's co-author. If more frequently a venue is related to neighbors, the venue is recommended to the researcher [84].

- (b) Collaborative filtering model (CF): It is based on the memory-based collaborative filtering with a given paper-venue matrix. The underlying assumption is that there is a high probability for a paper to get published in venues where other similar papers have been published [36].
- (c) Co-authorship network-based model (CN): This model used a new approach that builds a social network for each author and then recommends venues based on the reputation of the author's social network and other information such as venue name, venues sub-domain, number of publications [21].
- (d) Content-based filtering model (CBF): The main idea behind this approach is to compute the similarity between researchers and venues. Here we have taken researcher's publication and venues publications content as feature vectors respectively, which are calculated by LDA model [20].
- (e) Random walk with restart model (RWR): It runs a random walk with restart model on a co-publication network with two types of nodes: authors and venues. This model is similar to AVER, but the probability of skipping to the next neighbor node is equal in RWR [28].
- (f) Publication recommender system (PRS): It is based on a new content-based filtering (CBF) recommendation model using chi-square and softmax regression. It mainly consists of two modules, such as feature selection module and softmax regression module [30].
- (g) Personal venue rating-based collaborative filtering model (PVR): It is based on the implicit rating given to individual venues, created from references of a researcher's



**Table 11**

PPPN and VVPN recommendation performance in terms of accuracy and MRR.

Approach	Acc@3	Acc@6	Acc@9	Acc@12	Acc@15	MRR
FB	0.0555	0.0972	0.1250	0.1666	0.1944	0.0338
CF	0.0972	0.1111	0.1527	0.1805	0.2361	0.0451
CN	0.1111	0.1388	0.1805	0.2222	0.2500	0.0516
CBF	0.1527	0.1805	0.2083	0.2361	0.2916	0.0648
RWR	0.1944	0.2222	0.2500	0.2916	0.3194	0.0775
PVR	0.2083	0.2361	0.2368	0.3194	0.3472	0.0863
PRS	0.2063	0.2291	0.2486	0.2793	0.3419	0.0875
PAVE	0.2500 <sup>+</sup>	0.2916 <sup>+</sup>	0.3055 <sup>+</sup>	0.3611 <sup>+</sup>	0.4305 <sup>+</sup>	0.0906 <sup>+</sup>
PPPN	<b>0.3334</b>	<b>0.3611</b>	<b>0.4027</b>	0.4722	0.6805	0.1150
VVPN	0.3055	0.3457	0.3888	<b>0.5138</b>	<b>0.7361</b>	<b>0.1169</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

publications and the papers which cited researcher's past publications [19].

- (h) Personalized academic venue recommendation model (PAVE): It is similar to the popular random walk model except for the definition of transfer matrix with bias. The probability of skipping to the next neighbor node is biased using co-publication frequency, relation weight, and researcher's academic level in PAVE [22].

Among these eight methods CF and PVR are based on collaborative filtering approach, PAVE, RWR is based on random walk with restart algorithm exploiting co-authorship networks, CN, FB is based on co-authorship network, and CBF and PRS are based on content-based filtering method.

## 11. Results and discussion

In this section, we evaluate the effectiveness of CNAVER against existing state-of-the-art methods. Before assessing the performance of the proposed fusion model CNAVER individual performance analysis of PPPN model and VVPN model are analyzed in two phases such as Offline or Coarse-level evaluation and Online or Finer-level evaluation. During the assessment, best results and the second-best are marked by 'bold-face' and '+' symbol respectively.

### 11.1. Offline evaluation of PPPN model

The complete results of accuracy and MRR are presented in Table 11 during the position 3, 6, 9, 12, and 15 respectively. We can see that the PPPN model reveals to a consistent accuracy over all other state-of-the-art strategies. More than 33% of the time (Acc@3 = 0.3334), it can predict the original venue of the seed paper within top 3 recommendations. The PPPN approach shows an accuracy of 0.6805 while recommending top 15 recommendations. FB strategy exhibits bad performance with an accuracy of 0.1944 while recommending 15 recommendations. More than 68% time, PPPN model can predict the original venue of the seed paper within top 15 recommendation.

For MRR, PPPN performs excellent behavior (MRR 0.1150). The proposed approach can predict the original venue at early ranks compared to all other methods. In the case of MRR also, the least performance is demonstrated by the FB method.

### 11.2. Online evaluation of PPPN model

In this section, we analyze the performance of the PPPN model against other state-of-the-art methods. The evaluation metrics, including precision, nDCG, and average venue quality (H5-Index), are taken into consideration throughout this assessment.

**Table 12**

PPPN and VVPN performance in terms of precision.

Methods	P@3	P@6	P@9	P@12	P@15
FB	0.5646	0.5493	0.5347	0.5116	0.5392
CF	0.5656	0.5887	0.5889	0.5998	0.5885
CN	0.6114	0.5994	0.6028	0.6003	0.5904
CBF	0.6117	0.6113	0.6108	0.6008	0.6001
RWR	0.6551	0.6317	0.6254	0.6273	0.6299
PRS	0.6675	0.6976 <sup>+</sup>	0.6533 <sup>+</sup>	0.6205	0.6234
PVR	0.6559	0.6248	0.6229	0.6318	0.6301
PAVE	0.7005 <sup>+</sup>	0.6835	0.6492	0.6659 <sup>+</sup>	0.6678 <sup>+</sup>
PPPN	<b>0.7243</b>	<b>0.7307</b>	0.6948	0.6651	0.6509
VVPN	0.6992	0.6998	<b>0.7207</b>	<b>0.7114</b>	<b>0.7219</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

**Table 13**

PPPN and VVPN performance in terms of nDCG.

Methods	nDCG@3	nDCG@6	nDCG@9	nDCG@12	nDCG@15
FB	0.5913	0.5734	0.5743	0.5748	0.5786
CF	0.6109	0.6198	0.6213	0.6231	0.6217
CN	0.6255	0.6339	0.6296	0.6319	0.6325
CBF	0.6671	0.6511	0.6517	0.6587	0.6639
RWR	0.6589	0.6584	0.6527	0.6592	0.6657
PRS	0.6549	0.6585	0.6691	0.6604	0.6695
PVR	0.6612	0.6672	0.6632	0.6637	0.6543
PAVE	0.6759 <sup>+</sup>	0.6691 <sup>+</sup>	0.6701 <sup>+</sup>	0.6749 <sup>+</sup>	0.6771 <sup>+</sup>
PPPN	<b>0.6939</b>	<b>0.7013</b>	<b>0.6939</b>	0.6689	0.6706
VVPN	0.6584	0.6699	0.6892	<b>0.7209</b>	<b>0.7425</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

#### 11.2.1. Precision@k

In Fig. 4, we can see the significance of the PPPN model as far as precision@k and nDCG@k over all other standard approaches. The PPPN model exhibits the highest precision of 0.7348 at position 4, and after that, it marginally downgrades and furthermore achieves a precision about 0.6775 at position 11 as depicted in Fig. 4(a). The PPPN model performs superior until the initial 11 recommendations. Afterward, it shows a descending pattern because of which it is unable to maintain consistency as depicted in Table 12. This model demonstrates a lower precision of 0.6531 at position 13.

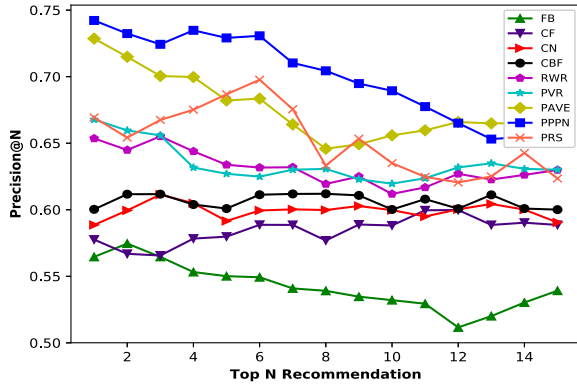
For the first venue, PPPN accomplishes the highest precision among all other methods. Later on, those precision continues diminishing and furthermore achieves a precision about 0.6509 at position 15. PAVE method indicates higher performance over PPPN model at position 12, 13, 14, and 15 respectively. The worst performance among all methods is demonstrated by the FB method.

#### 11.2.2. nDCG@k

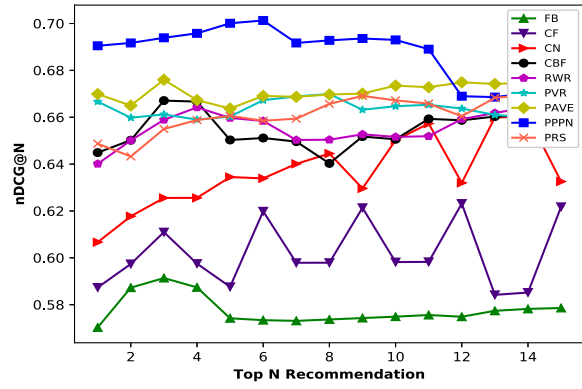
The overall nDCG@k of all methods are shown in Table 13. Throughout nDCG@k evaluation, PPPN model demonstrates superior scores over all other state-of-the-art methods. The PPPN model performs an upward trend and furthermore achieves the most astounding nDCG 0.7013 during position 6, also subsequently again, it reveals to a descending pattern and shows nDCG of 0.6685 at position 13. Afterward, it gradually increases and accomplishes a decent nDCG 0.6706 at position 15, as depicted in Fig. 4(b).

#### 11.2.3. Average venue quality (H5-Index) analysis

We have additionally assessed the quality of venues recommended by PPPN model as compared to other existing methodologies. The PPPN model outperforms other methods in terms of average H5-Index of recommended venues, as illustrated in Fig. 5(a). While assessing average venue quality, the PPPN model performs an upward trend from the beginning and shows an

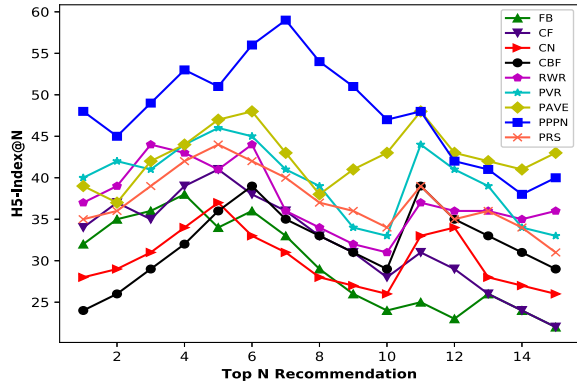


(a) PPPN performance in terms of precision@k

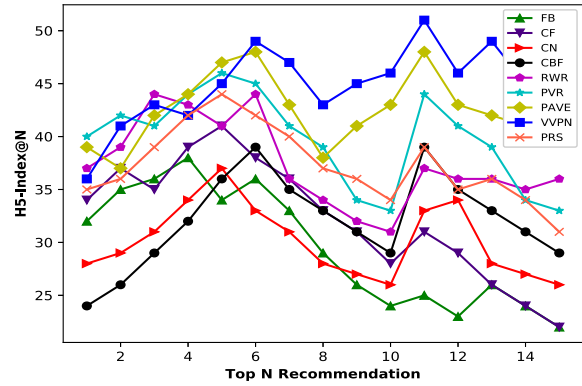


(b) PPPN performance in terms of nDCG@k

Fig. 4. PPPN recommendation performance in terms of precision@k and nDCG@k.



(a) PPPN Average venue quality



(b) VVPN Average venue quality

Fig. 5. PPPN and VVPN performance in terms of average venue quality.

overall average H5-Index about 49. The top-quality venues recommended by PPPN is in position 7 with the highest H5-Index of 59. Then it indicates a descending pattern furthermore achieves an H5-Index of value 40 at position 15, as shown in Fig. 5(a). The lowest quality of venues recommended by the FB method with an average H5-Index of 30, whereas the second-highest quality venues recommended by PAVE model with an average H5-Index about 43.

### 11.3. Offline evaluation of VVPN model

VVPN model shows a consistent accuracy over all other standard approaches (Table 11). More than 30% time it can predict the original venue of the seed paper within the top 3 recommendations. Initially, the VVPN model shows an accuracy of 0.3457 at position 6. Then slowly it shows an upward trend and exhibits an excellent performance with an accuracy of 0.7361 at position 15.

VVPN also shows excellent performance over other standard approaches in terms of MRR. VVPN model exhibits the overall MRR of 0.1169. The second-best performance is shown by the PAVE model with MRR 0.0906. The proposed approach could predict the original venue at early recommendations as compared to all other methods. In the case of MRR also, the least performance is exhibited by the FB method.

### 11.4. Online evaluation of VVPN model

We evaluate at a finer level, the effect of VVPN recommendations and compare it against other state-of-the-art methods. We use various metrics such as precision, nDCG, and average venue quality (H5-Index), respectively.

#### 11.4.1. Precision@k

The compared results are shown in Fig. 6. It can be easily observed that the proposed approach VVPN model has made a significant improvement of precision@k over the standard approaches as depicted in Fig. 6(a). Initially, the VVPN model shows a precision of 0.7132 at a position 2. Then it slowly indicates a downward trend and reaches a precision of 0.6998 at position 6 as depicted in Table 12.

But afterward, it shows an upward trend and shows the highest precision of 0.7219 at position 15 and least precision value of 0.6804 after recommending 5 recommendations. The PRS model shows the second-best performance at position 5, 6, and 7 respectively. PAVE exhibits excellent performance at position 1, 2, 3 and 4 respectively. The worst performance among all methods is shown by the FB method.

#### 11.4.2. nDCG@k

The overall nDCG scores are shown in Table 13. Initially, the VVPN model shows a lower nDCG 0.6587 at position 2. Then slowly it shows a downward trend and reaches the nDCG 0.6578 at position 5. Afterward, it shows an upward trend and is able to show consistency at other positions of the recommendations. It is clearly shown in Fig. 6(b) that the graph of VVPN is consistent after recommending 5 venues and shows a nDCG 0.7209 at position 12. But method PAVE shows higher nDCG than VVPN model at position 1, 2, 3, 4 and 5 respectively. It consistently shows the second-best performance throughout. The FB model exhibits the worst performance.

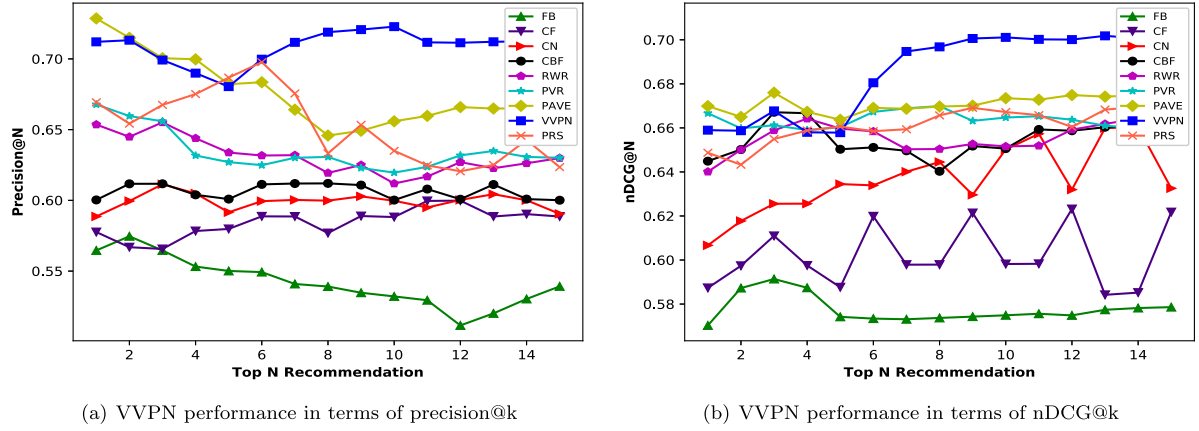


Fig. 6. VVPN recommendation performance in terms of precision@k and nDCG@k.

Table 14

Accuracy and MRR results of CNAVER and other compared approaches.

Approach	Acc@3	Acc@6	Acc@9	Acc@12	Acc@15	MRR
FB	0.0555	0.0972	0.1250	0.1666	0.1944	0.0338
CF	0.0972	0.1111	0.1527	0.1805	0.2361	0.0451
CN	0.1111	0.1388	0.1805	0.2222	0.2500	0.0516
CBF	0.1527	0.1805	0.2083	0.2361	0.2916	0.0648
RWR	0.1944	0.2222	0.2500	0.2916	0.3194	0.0775
PVR	0.2083	0.2361	0.2368	0.3194	0.3472	0.0863
PRS	0.2063	0.2291	0.2486	0.2793	0.3419	0.0875
PAVE	0.2500 <sup>+</sup>	0.2916 <sup>+</sup>	0.3055 <sup>+</sup>	0.3611 <sup>+</sup>	0.4305 <sup>+</sup>	0.0906 <sup>+</sup>
<b>CNAVER</b>	<b>0.3572</b>	<b>0.3888</b>	<b>0.4583</b>	<b>0.5833</b>	<b>0.7916</b>	<b>0.1402</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

#### 11.4.3. Average venue quality (H5-Index) analysis

We investigate the performance of venue quality recommended by VVPN as compared to other existing approaches. VVPN model outperforms other methods in terms of average H5-Index of recommended venues. Overall, the average H5-Index of venues recommended by the VVPN model is 45. The top-quality venues recommended by VVPN are at position 11 with the highest H5-Index of 51 as displayed in Fig. 5(b).

#### 11.5. Offline evaluation of fusion model: CNAVER

The complete results of accuracy and MRR after fusion are depicted in Table 14. It is evident from the overall results of accuracy and MRR that the proposed approach CNAVER shows a consistent performance over all other standard approaches. More than 35% of the time it can predict the original venue of the seed paper within the top 3 recommendations.

The proposed approach shows an accuracy of 0.7916 after recommending the top 15 recommendations. Similarly, during the evaluation of MRR, we can see that CNAVER outperforms all other state-of-the-art methods and shows excellent behavior with a MRR 0.1402. The proposed approach can predict the original venue at early recommendations better than all other methods. The second-best performance is exhibited by the PAVE, whereas the FB performs the worst among all different standard approaches.

We have also investigated the efficacy of the proposed model CNAVER in terms of  $F - measure_{macro}$  ( $F_1$ ) against other state-of-the-art methods. The complete results of  $F - measure_{macro}$  are shown in Table 15.  $F_1$  scores are generally seen to increase with rank up to a certain point (around 9–12) and drop thereafter. This is possibly due to the fact that precision and recall both increase till that point until the original venues are retrieved, causing an increase in  $F_1$  score. However, with further increase in ranks,

Table 15

Macro-average analysis in terms of  $F - measure_{macro}$  ( $F_1$ ).

Approach	$F_1@3$	$F_1@6$	$F_1@9$	$F_1@12$	$F_1@15$
FB	0.0128	0.0231	0.0458	0.0412	0.0408
CF	0.0351	0.0437	0.0759	0.0694	0.0621
CN	0.0561	0.0672	0.1045	0.1004	0.0938
CBF	0.0894	0.1025	0.1289	0.1167	0.1125
RWR	0.1148	0.1413	0.2141	0.1894	0.1663
PVR	0.1297	0.1568	0.1854	0.1945	0.1867
PRS	0.1231	0.1456	0.1959	0.1889	0.1826
PAVE	0.1631 <sup>+</sup>	0.2012 <sup>+</sup>	0.2674 <sup>+</sup>	0.2248 <sup>+</sup>	0.2179 <sup>+</sup>
<b>CNAVER</b>	<b>0.2769</b>	<b>0.3179</b>	<b>0.3627</b>	<b>0.3561</b>	<b>0.3524</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

precision drops sharply without much increase in recall leading to an overall drop in  $F_1$  scores. Here also CNAVER demonstrates the efficacy in comparison to other state-of-the-art methods. The second-best performance is exhibited by PAVE, whereas FB performs the worst.

#### 11.6. Online evaluation of fusion model: CNAVER

In this section, the performance of CNAVER against other state-of-the-art methods is discussed. We demonstrate the performance of the proposed system considering various evaluation metrics such as precision, nDCG, and average venue quality (H5-Index), respectively.

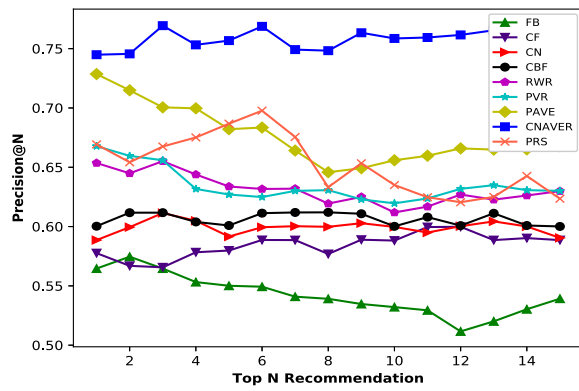
##### 11.6.1. Precision@k

The overall results precision and nDCG evaluations are shown in Fig. 7. In Fig. 7(a), we can see the significance of CNAVER in terms of precision over all other standard approaches. Initially, the proposed CNAVER exhibits a precision of 0.7456 at position 2, and after that, it slightly shows an upward trend and shows a precision of 0.7634 at position 9 (Table 16).

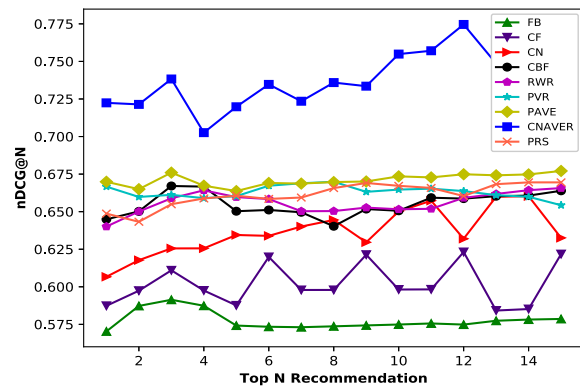
The proposed model CNAVER exhibits the highest precision of 0.7704 after recommending 15 recommendations. It shows a lower precision of 0.7449 at position 1. Similarly, the PAVE method performs the second-best at positions 1, 2, 3 and 4 respectively. PRS method exhibits slightly higher precision than the PAVE method at position 5, 6, 7 and 9 respectively. The worst performance among all methods is shown by the FB method.

##### 11.6.2. nDCG@k

The nDCG@k evaluations of all methods are shown in Table 17. Proposed CNAVER also exhibits better nDCG scores over all other state-of-the-art methods. The CNAVER model performs



(a) CNAVER results in terms of precision@k



(b) CNAVER results in terms of nDCG@k

Fig. 7. CNAVER performance in terms of precision@k and nDCG@k.

Table 16

Precision of CNAVER and other compared approaches.

Methods	P@3	P@6	P@9	P@12	P@15
FB	0.5646	0.5493	0.5347	0.5116	0.5392
CF	0.5656	0.5887	0.5889	0.5998	0.5885
CN	0.6114	0.5994	0.6028	0.6003	0.5904
CBF	0.6117	0.6113	0.6108	0.6008	0.6001
RWR	0.6551	0.6317	0.6254	0.6273	0.6299
PRS	0.6675	0.6976 <sup>+</sup>	0.6533 <sup>+</sup>	0.6205	0.6234
PVR	0.6559	0.6248	0.6229	0.6318	0.6301
PAVE	0.7005 <sup>+</sup>	0.6835	0.6492	0.6659 <sup>+</sup>	0.6678 <sup>+</sup>
<b>CNAVER</b>	<b>0.7694</b>	<b>0.7687</b>	<b>0.7634</b>	<b>0.7629</b>	<b>0.7704</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

Table 17

nDCG of CNAVER and other compared approaches.

Methods	nDCG@3	nDCG@6	nDCG@9	nDCG@12	nDCG@15
FB	0.5913	0.5734	0.5743	0.5748	0.5786
CF	0.6109	0.6198	0.6213	0.6231	0.6217
CN	0.6554	0.6339	0.6296	0.6319	0.6325
CBF	0.6671	0.6511	0.6517	0.6587	0.6639
RWR	0.6589	0.6584	0.6527	0.6592	0.6657
PRS	0.6549	0.6585	0.6691	0.6604	0.6695
PVR	0.6612 <sup>+</sup>	0.6672	0.6632	0.6637	0.6643
PAVE	0.6599	0.6691 <sup>+</sup>	0.6701 <sup>+</sup>	0.6749 <sup>+</sup>	0.6771 <sup>+</sup>
<b>CNAVER</b>	<b>0.7283</b>	<b>0.7247</b>	<b>0.7235</b>	<b>0.7467</b>	<b>0.7511</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

an upward trend and reaches a nDCG 0.7359 at position 8, and afterward, it shows an upward trend and reaches nDCG 0.7611 at position 15 (Fig. 7(b)). The performance of CNAVER is consistent and shows the highest nDCG 0.7746 at position 12. The PAVE model demonstrates the second-best performance. The FB model shows the worst performance among all other standard approaches.

### 11.6.3. Average venue quality (H5-Index) analysis

We also investigate the performance of venue quality recommended by CNAVER as compared to other existing approaches. CNAVER outperforms other methods in terms of average H5-Index of recommended venues as depicted in Table 18. Overall, the average H5-Index of venues recommended by CNAVER is 54. The top-quality venues recommended by CNAVER are at position 7 with the highest H5-Index of 63 as displayed in Fig. 8.

### 11.6.4. Evaluation of diversity

Diversity is defined in terms of content dissimilarity. We group all papers published at a particular venue and extract their corresponding keywords. We apply the similarity score to define

Table 18

H5-Index of CNAVER and other compared approaches.

Approach	HI@3	HI@6	HI@9	HI@12	HI@15
FB	36	36	26	23	22
CF	35	38	31	29	22
CN	31	33	27	34	26
CBF	29	39	31	35	29
PRS	39	42	36	35	31
RWR	44 <sup>+</sup>	44	32	36	36
PVR	41	45	34	41	33
PAVE	42	48 <sup>+</sup>	41 <sup>+</sup>	43 <sup>+</sup>	43 <sup>+</sup>
<b>CNAVER</b>	<b>51</b>	<b>58</b>	<b>52</b>	<b>52</b>	<b>53</b>

Best results are highlighted in bold, and 2nd best are marked by ('+').

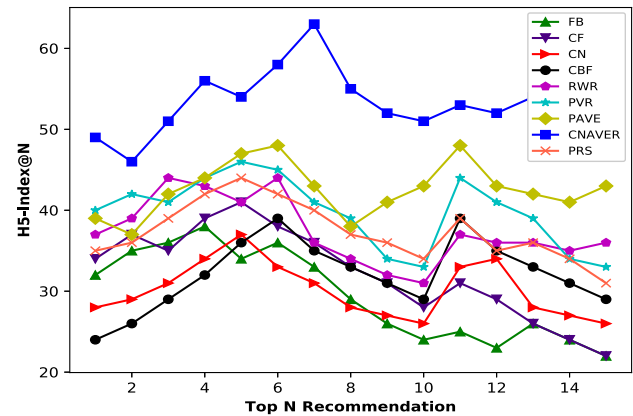


Fig. 8. Average venue quality of CNAVER and other approaches.

diversity in Eq. (49), and the scores are in Table 19. CNAVER is seen to show the best diversity, whereas the second-best performer is the method PVR.

### 11.6.5. Evaluation of stability

We have also provided a comprehensive investigation of the stability of CNAVER as defined in Eq. (50). CNAVER shows lower MAS than all other standard approaches (Table 20). It shows a MAS of 4.359 on the DBLP dataset, meaning that on an average every predicted venue will shift by a position of 4.359 after adding new data into the training data of the system. We have considered the average MAS-score as a threshold to decide whether a particular method provides stability or not.

### 11.6.6. Statistical significance test

To see the difference between the performance of CNAVER and the second best is significant or not, we did paired-samples t-test.



**Table 19**

Diversity (D) of CNAVER and other approaches.

Methods	Diversity (D)
FB	0.219
CF	0.338
CN	0.273
CBF	0.204
RWR	0.309
PVR	0.394 <sup>+</sup>
PRS	0.273
PAVE	0.316
<b>CNAVER</b>	<b>0.497</b>

Best results and 2nd best are marked by bold, and ('+').

**Table 20**

Stability (MAS) of CNAVER and other approaches.

Methods	MAS
FB	9.821
CF	8.757
CN	9.452
CBF	5.769
RWR	7.884
PVR	8.236
PRS	5.351 <sup>+</sup>
PAVE	8.349
<b>CNAVER</b>	<b>4.359*</b>

Best results and 2nd best are marked by bold, and ('+').

It determines whether there is statistical evidence that the mean difference between paired observations on a particular outcome is significantly different from zero. This is a parametric test done using t-statistic as defined below.

$$t = \frac{m}{s/\sqrt{n}} \quad (53)$$

Where  $m$  and  $s$  denote the mean and standard deviation of the differences between all pairs. Here,  $n$  denotes the size of the sample space.

For a two-sided test the null and alternative hypothesis are considered the following way.

$H_0 : \mu_d = 0$  (True mean difference is equal to zero)

$H_1 : \mu_d \neq 0$  (True mean difference is not equal to zero)

To claim superiority (or inferiority) of a system over another, one-sided test is done where  $H_0$  remains the same, but  $H_1$  is modified as follows.  $H_1 : \mu_d > (or <) 0$

Under the null hypothesis, this statistic is supposed to follow a t-distribution with  $n - 1$  degrees of freedom (df). If the  $p$ -value is much less than  $\alpha$  (here  $\alpha = 0.01, 0.05, 0.1$ ) then we can reject the null hypothesis ( $H_0$ ) in favor of alternative hypothesis ( $H_1$ ).

We performed one-sided test ( $\mu_d = \mu_1 - \mu_2 > 0$ ) on overall precision, nDCG, accuracy, MRR,  $F - measure_{macro}$ , average venue quality, and diversity and one-sided test ( $\mu_d = \mu_1 - \mu_2 < 0$ ) in terms of stability to validate the significance of proposed CNAVER against other state-of-the-art methods. During the test, we have compared the paired samples t-tests results at positions 3, 6, 9, 12, and 15, respectively. It required 120 pairwise comparisons ( $df = 119$ ) for each position during the evaluation. The statistically significant results are highlighted in the star annotated symbol at each position. The proposed model CNAVER consistently outperforms other approaches, and the differences are statistically significant, as depicted in Tables 22–28.

## 12. Study of the proposed approach

The main findings concerning our SQs as introduced in Section 10.4.1 are summarized below:

### 12.1. SQ1: How effective is CNAVER in comparison to other state-of-the-art methods?

The overall results of CNAVER and other state-of-the-art methods are displayed in Tables 14, 15, 16, and 17 respectively. It demonstrates the best performance in terms of precision@k, nDCG@k, accuracy, MRR, and  $F - measure_{macro}$ , respectively. Also, the difference with the second-best is statistically significant even at 1% level of significance.

### 12.2. SQ2: How is the quality of venues recommended by CNAVER as compared to state-of-the-art methods?

The complete results of venue quality in terms of H5-Index is depicted in Table 18. The venues recommended by CNAVER are of high quality when contrasted with other cutting edge techniques as portrayed in Fig. 8. The average H5-index of CNAVER is 54 after recommending 15 venues. PAVE recommends venues having the second-best H5-index. The least quality of recommendation performed by the FB model. The most elevated H5-index recommended by CNAVER is 63, and the least is 46, whereas the most noteworthy H5-index suggested by the second-best PAVE is 48 and the least is 37.

### 12.3. SQ3: How does CNAVER handle cold-start issues and other issues like data sparsity, diversity, and stability

- Cold-start issues:** To specifically address “cold-start” issues like a new researcher and a new venue, PPPN and VVPN are fused to give venue recommender framework customizable for personalization. Examination of Tables 9 and 10 reveal that, regardless of whether the seed paper identified with the new researcher and new venue, CNAVER can predict the original venue at an early stage of recommendations. It does not require past publication records or co-authorship networks for the recommendations. Rather it only focuses on the work at hand. It considers only the current area of interest along with the title, and abstract as inputs to recommend the same.
- Data sparsity:** To explicitly address the data sparsity issue, both importance and relevance parameters are considered at the beginning phase of the proposed method. Social network analysis through different centrality measures and content features like abstract and title were used to capture the quality of essentiality, relevance, and importance separately. To extract only related papers, the entire citation network is apportioned, and later on, intra-graph clustering is performed. So given the seed paper, just essential and relevant papers are filtered from the sub-clusters. It has been noticed that the number of papers found after centrality measures are around 32,069 out of total 2, 236, 968 papers as input. The average number of papers involved after Intra-graph clustering for abstract similarity is in the range of 80–120. After the initial step, we are left with important papers for further computation, which is close to the area of interest. Hence there is no data sparsity issue in our proposed approach, as indicated in Table 21.
- Diversity:** To resolve the issue of diversity, both connection and contextual similarity-based relevance parameters are taken into consideration. Mainly age-discounted Venue2Vec, meta-path features, and biased random walk are incorporated in VVPN to recommend venues from diverse publishers. 1-degree and 2-degree meta-paths capture different rich latent information in VVPN model. In the PPPN model, topic modeling alongside intra-graph clustering captures both contexts as well as links to suggest

**Table 21**

Cold-start and other issues available in CNAVER and other approaches.

Methods	Cold-start	Sparsity	Diversity	Stability
FB	yes (new researcher)	no	yes	yes
CF	yes (researcher and venue)	yes	no	yes
CN	yes (new venue)	no	yes	yes
CBF	yes(new venue)	no	yes	no
RWR	yes (new researcher)	no	yes	yes
PRS	yes(new venue)	no	yes	no
PVR	yes (researcher and venue)	yes	no	yes
PAVE	yes(new researcher)	no	yes	yes
<b>CNAVER</b>	no	no	no	no

**Table 22**

p values of paired t-test with CNAVER (Precision).

Methods	P@3	P@6	P@9	P@12	P@15
FB	0.00011*	0.00008*	0.00004*	0.00002*	0.00006*
CF	0.00013	0.00016*	0.00015*	0.00019*	0.00015*
CN	0.00023*	0.00016*	0.00019*	0.00021*	0.00014*
CBF	0.00023*	0.00021*	0.00020*	0.00017*	0.00013*
RWR	0.00035*	0.00029*	0.00024*	0.00026*	0.00021*
PRS	0.00039*	0.00057*	0.00038*	0.00027*	0.00022*
PVR	0.00036*	0.00026*	0.00023*	0.00029*	0.00027*
PAVE	0.00053*	0.00046*	0.00031*	0.00039*	0.00037*

\*p values are statistically significant at  $\alpha = 0.05$ .**Table 23**

p values of paired t-test with CNAVER (nDCG).

Methods	nDCG@3	nDCG@6	nDCG@9	nDCG@12	nDCG@15
FB	0.00023*	0.00012*	0.00009*	0.00003*	0.00005*
CF	0.00025*	0.00027*	0.00031*	0.00020*	0.00017*
CN	0.00031*	0.00039*	0.00037*	0.00024*	0.00017*
CBF	0.00052*	0.00045*	0.00041*	0.00032*	0.00038*
RWR	0.00049*	0.00053*	0.00048*	0.00052*	0.00036*
PRS	0.00045*	0.00043*	0.00048*	0.00029*	0.00033*
PVR	0.00051*	0.00049*	0.00053*	0.00032*	0.00031*
PAVE	0.00065*	0.00057*	0.00062*	0.00041*	0.00049*

\*p values are statistically significant at  $\alpha = 0.05$ .

relevant papers from diverse publishers. CNAVER, therefore shows the highest value of  $D$  (diversity) as compared to all other approaches (Table 19).

- (iv) **Stability:** To deal with the stability issue, a series of techniques are used, broadly divided into two classes: PPPN and VVPN-which are finally fused together to provide a single ranked list of recommendations. Both PPPN and VVPN are a pipeline of techniques where papers are filtered and/or added to enrich the pool of shortlisted papers. Techniques like content-similarity, various centrality measures, meta-path, random walk are used to invite diversity as well as robustness to the system. Each of these techniques participates in a co-operative manner where the contribution of any single technique is not immensely decisive. Rather, we have some amount of redundancy such that a paper is potentially shortlisted by several techniques. To counter the destabilizing nature of network-based approaches, content-based approaches are incorporated at several places in the pipeline. In all, these batteries of techniques together provide stability to the recommendations. CNAVER shows the minimum MAS than all other standard approaches (Table 20).

#### 12.4. More insightful discussion on the results

The overall performance results obtained and discussed in Section 11 showcase the efficacy of the proposed CNAVER. The

**Table 24**

p values of paired t-test with CNAVER (Accuracy and MRR).

Approach	Acc@3	Acc@6	Acc@9	Acc@12	Acc@15	MRR
FB	0.00009*	0.00013*	0.00010*	0.00008*	0.00001*	0.00103*
CF	0.00011*	0.00010*	0.00011*	0.00008*	0.00003*	0.00123*
CN	0.00013*	0.00012*	0.00014*	0.00008*	0.00005*	0.00218*
CBF	0.00012*	0.00013*	0.00011*	0.00009*	0.00006*	0.00276*
RWR	0.00016*	0.00017*	0.00015*	0.00011*	0.00007*	0.00318*
PRS	0.00022*	0.00017*	0.00013*	0.00011*	0.00009*	0.00047*
PVR	0.00021*	0.00023*	0.00013*	0.00016*	0.00007*	0.00446*
PAVE	0.00057*	0.00065*	0.00043*	0.00039*	0.00012*	0.00594*

\*p values are statistically significant at  $\alpha = 0.05$ .**Table 25**p values of paired t-test with CNAVER ( $F_1$ ).

Approach	$F_1$ @3	$F_1$ @6	$F_1$ @9	$F_1$ @12	$F_1$ @15
FB	0.00009*	0.00013*	0.00010*	0.00008*	0.00001*
CF	0.00011*	0.00010*	0.00011*	0.00008*	0.00003*
CN	0.00013*	0.00012*	0.00014*	0.00008*	0.00005*
CBF	0.00012*	0.00013*	0.00011*	0.00009*	0.00006*
RWR	0.00016*	0.00017*	0.00015*	0.00011*	0.00007*
PRS	0.00022*	0.00017*	0.00013*	0.00011*	0.00009*
PVR	0.00021*	0.00023*	0.00013*	0.00016*	0.00007*
PAVE	0.00061*	0.00063*	0.00057*	0.00051*	0.00047*

\*p values are statistically significant at  $\alpha = 0.05$ .**Table 26**

p values of paired t-test with CNAVER (H1:H5-Index).

Approach	H1@3	H1@6	H1@9	H1@12	H1@15
FB	0.00037*	0.00012*	0.00014*	0.00011*	0.00009*
CF	0.00032*	0.00024*	0.00019*	0.00016*	0.00015*
CN	0.00034*	0.00017*	0.00014*	0.00011*	0.00007*
CBF	0.00017*	0.00023*	0.00019*	0.00025*	0.00012*
RWR	0.00059*	0.00034*	0.00029*	0.00039*	0.00034*
PRS	0.00031*	0.00026*	0.00027*	0.00023*	0.00016*
PVR	0.00036*	0.00027*	0.00019*	0.00016*	0.00007*
PAVE	0.00042*	0.00038*	0.00034*	0.00043*	0.00036*

\*p values are statistically significant at  $\alpha = 0.05$ .

excellent overall precision implies that the models can effectively recommend the relevant venues. However, there are a few limitations of our work.

- The proposed system has multiple parameters involved in both PPPN and VVPN models. Most of the steps involved in the PPPN model are purely based on empirical assumptions but backed by observations from rigorous experimentation.
- If the topmost  $R$  papers similar to a given seed paper are loosely coupled in the bibliographic citation network, Jarvis-Patrick may create clusters with less number of related papers. Hence, the proposed model CNAVER may fail to capture the relevant papers resulting in possibly irrelevant venue recommendations.
- In the VVPN model, while choosing the venue of interest ( $Z$ ) if the topmost paper is contextually similar with the seed paper but its corresponding venue associated with entirely different domains then few of the topmost recommendations by random walk algorithm may not be relevant.
- If the original venue of a seed paper is comparatively new and the venue does not have sufficient number papers, the system may perform poorly. Although the venue of interest  $Z$  is contextually similar in content but due to meta-paths features aggregation, other venues may be recommended in the VVPN model, but the original venue may not appear at the top of the recommendation list. Hence, it may results in low accuracy and low MRR during on-line evaluation.

**Table 27**

p values of paired t-test with CNAVER (Diversity).

Methods	Diversity (D)
FB	0.00008*
CF	0.00019*
CN	0.00013*
CBF	0.00009*
RWR	0.00021*
PVR	0.00047*
PRS	0.00018*
PAVE	0.00025*

\*p values are statistically significant at  $\alpha = 0.05$ **Table 28**

p values of paired t-test with CNAVER (Stability).

Methods	Stability (MAS)
FB	0.00024*
CF	0.00046*
CN	0.00032*
CBF	0.00089*
RWR	0.00073*
PVR	0.00053*
PRS	0.00095*
PAVE	0.00046*

\*p values are statistically significant at  $\alpha = 0.05$ 

### 13. Conclusion and future work

Academic venue recommendation is an emerging area of research in recommendation systems. The prevalent techniques are few in numbers and suffer from various limitations. One of the major issues is that of cold-start having two sub-parts: a new venue and a new researcher. Additionally, there exist problems of sparsity, diversity, and stability in venue recommender systems that are not adequately addressed in existing state-of-the-art methods.

We proposed a fusion-based scholarly venue recommender system CNAVER incorporating paper-paper peer network (PPPN) model and venue-venue peer network (VVPN) model that reasonably addresses the above-mentioned issues. Several techniques like topic modeling based contextual similarity, link analysis, and topic-oriented intra-graph clustering, abstract similarity using Okapi BM25+ algorithm are used to reinforce the PPPN model. To identify relevant venues, age-discounted based Venue2Vec, different meta-paths features, and biased random walk with restart (RWR) algorithm are incorporated into the VVPN model.

We conducted an extensive set of experiments on a real dataset DBLP and showed that CNAVER consistently outperforms state-of-the-art methods. It shows substantially higher scores of precision@k, nDCG@k, accuracy, MRR,  $F - measure_{macro}$ , diversity and stability than other best in class techniques. CNAVER proposes top-notch venues when contrasted with cutting edge techniques as far as H5-index.

Nonetheless, there is scope for continuous update of the model. Considering the fast development of digital information technology, we would like to employ a web crawler to update the training dataset and the learning model continuously. This crawler will automatically extract and collect the relevant data to generate the training dataset. To continually enhance the quality of the recommendation of CNAVER, we plan to collect feedback from users through a web-based application. We plan to adopt some information retrieval techniques like relevance feedback, or pseudo relevance feedback to improve the relevance of final recommendations. In the future, we would like to incorporate advanced machine learning techniques such as gradient descent

optimization, etc., in such a way that it will enforce the random walker not to go too far from the initial venue of interest (Z).

We intend to explore with different datasets and to broaden it for various controls with the objective of enhancing precision, accuracy, diversity, novelty, coverage, and serendipity. We would like to explore the same with the assistance of heterogeneous bibliographic data coordinate with all conceivable meta-paths (till degree four) highlights to recommend scholarly venues.

### 14. Compliance with ethical standards

The authors declare no conflicts of interest. The article utilizes a rank based fusion model CNAVER exploiting paper-paper peer network (PPPN) model and venue-venue peer network (VVPN) model to recommend publication venues for a researcher. The article does not contain any examinations with human or creature subjects.

### References

- [1] D. Liang, L. Charlin, J. McInerney, D.M. Blei, Modeling user exposure in recommendation, in: Proceedings of the 25th International Conference on World Wide Web, International World Wide Web Conferences Steering Committee, 2016, pp. 951–961.
- [2] G. Adomavicius, A. Tuzhilin, Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions, IEEE Trans. Knowl. Data Eng. 17 (6) (2005) 734–749.
- [3] J. Bobadilla, F. Ortega, A. Hernando, A. Gutiérrez, Recommender systems survey, Knowl.-Based Syst. 46 (2013) 109–132.
- [4] J. Tang, S. Wu, J. Sun, H. Su, Cross-domain collaboration recommendation, in: Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, 2012, pp. 1285–1293.
- [5] S. Cohen, L. Ebel, Recommending collaborators using keywords, in: Proceedings of the 22nd International Conference on World Wide Web, ACM, 2013, pp. 959–962.
- [6] P. Chaiwananorn, C. Lursinsap, Collaborator recommendation in interdisciplinary computer science using degrees of collaborative forces, temporal evolution of research interest, and comparative seniority status, Knowl.-Based Syst. 75 (2015) 161–172.
- [7] J. Son, S.B. Kim, Academic paper recommender system using multilevel simultaneous citation networks, Decis. Support Syst. 105 (2018) 24–33.
- [8] Y. Sebastian, E.-G. Siew, S.O. Orimaye, Learning the heterogeneous bibliographic information network for literature-based discovery, Knowl.-Based Syst. 115 (2017) 66–79.
- [9] G. Wang, X. He, C.I. Ishuga, Har-si: A novel hybrid article recommendation approach integrating with social information in scientific social network, Knowl.-Based Syst. 148 (2018) 85–99.
- [10] A.S. Raamkumar, S. Foo, N. Pang, Using author-specified keywords in building an initial reading list of research papers in scientific paper retrieval and recommender systems, Inf. Process. Manage. 53 (3) (2017) 577–594.
- [11] W. Zhao, R. Wu, H. Liu, Paper recommendation based on the knowledge gap between a researcher's background knowledge and research target, Inf. Process. Manage. 52 (5) (2016) 976–988.
- [12] W. Huang, Z. Wu, L. Chen, P. Mitra, C.L. Giles, A neural probabilistic model for context based citation recommendation, in: AAAI, 2015, pp. 2404–2410.
- [13] X. Liu, Y. Yu, C. Guo, Y. Sun, Meta-path-based ranking with pseudo relevance feedback on heterogeneous graph for citation recommendation, in: Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, ACM, 2014, pp. 121–130.
- [14] Q. He, D. Kifer, J. Pei, P. Mitra, C.L. Giles, Citation recommendation without author supervision, in: Proceedings of the Fourth ACM International Conference on Web Search and Data Mining, ACM, 2011, pp. 755–764.
- [15] Z. Yang, D. Yin, B.D. Davison, Recommendation in academia: A joint multi-relational model, in: Advances in Social Networks Analysis and Mining, ASONAM, 2014 IEEE/ACM International Conference on, IEEE, 2014, pp. 566–571.
- [16] X. Tang, X. Wan, X. Zhang, Cross-language context-aware citation recommendation in scientific articles, in: Proceedings of the 37th International ACM SIGIR Conference on Research & Development in Information Retrieval, ACM, 2014, pp. 817–826.
- [17] J. Beel, B. Gipp, S. Langer, C. Breiteringer, Paper recommender systems: a literature survey, Int. J. Digit. Lib. 17 (4) (2016) 305–338.
- [18] F. Xia, N.Y. Asabere, J.J. Rodrigues, F. Basso, N. Deonauth, W. Wang, Socially-aware venue recommendation for conference participants, in: Ubiquitous Intelligence and Computing, 2013 IEEE 10th International Conference on and 10th International Conference on Autonomic and Trusted Computing, UIC/ATC, IEEE, 2013, pp. 134–141.

- [19] H. Alhoori, R. Furuta, Recommendation of scholarly venues based on dynamic user interests, *J. Infometrics* 11 (2) (2017) 553–563.
- [20] E. Medvet, A. Bartoli, G. Piccinin, Publication venue recommendation based on paper abstract, in: Tools with Artificial Intelligence, ICTAI, 2014 IEEE 26th International Conference on, IEEE, 2014, pp. 1004–1010.
- [21] H. Luong, T. Huynh, S. Gauch, L. Do, K. Hoang, Publication venue recommendation using author network's publication history, *Intell. Inf. Database Syst.* (2012) 426–435.
- [22] S. Yu, J. Liu, Z. Yang, Z. Chen, H. Jiang, A. Tolba, F. Xia, PAVE: Personalized academic venue recommendation exploiting co-publication networks, *J. Netw. Comput. Appl.* 104 (2018) 38–47.
- [23] F. Xia, W. Wang, T.M. Bekele, H. Liu, Big scholarly data: A survey, *IEEE Trans. Big Data* 3 (1) (2017) 18–35.
- [24] J. Lu, D. Wu, M. Mao, W. Wang, G. Zhang, Recommender system application developments: a survey, *Decis. Support Syst.* 74 (2015) 12–32.
- [25] N.M. Villegas, C. Sánchez, J. Díaz-Cely, G. Tamura, Characterizing context-aware recommender systems: A systematic literature review, *Knowl.-Based Syst.* 140 (2018) 173–200.
- [26] M.C. Pham, D. Kovachev, Y. Cao, G.M. Mbogor, R. Klamma, Enhancing academic event participation with context-aware and social recommendations, in: Advances in Social Networks Analysis and Mining, ASONAM, 2012 IEEE/ACM International Conference on, IEEE, 2012, pp. 464–471.
- [27] K. Sugiyama, M.-Y. Kan, Towards higher relevance and serendipity in scholarly paper recommendation by Kazunari Sugiyama and Min-Yen Kan with Martin Vesely as coordinator, *ACM SIGWEB Newsl.* (Winter) (2015) 4.
- [28] Z. Chen, F. Xia, H. Jiang, H. Liu, J. Zhang, AVER: random walk based academic venue recommendation, in: Proceedings of the 24th International Conference on World Wide Web, ACM, 2015, pp. 579–584.
- [29] H. Alhoori, How to identify specialized research communities related to a researcher's changing interests, in: Digital Libraries, JCDL, 2016 IEEE/ACM Joint Conference on, IEEE, 2016, pp. 239–240.
- [30] D. Wang, Y. Liang, D. Xu, X. Feng, R. Guan, A content-based recommender system for computer science publications, *Knowl.-Based Syst.* 157 (2018) 1–9.
- [31] P. Lops, M. De Gemmis, G. Semeraro, Content-based recommender systems: State of the art and trends, in: Recommender Systems Handbook, Springer, 2011, pp. 73–105.
- [32] R. Klamma, P.M. Cuong, Y. Cao, You never walk alone: Recommending academic events based on social network analysis, in: International Conference on Complex Sciences, Springer, 2009, pp. 657–670.
- [33] M. Hornick, P. Tamayo, Extending recommender systems for disjoint user/item sets: The conference recommendation problem, *IEEE Trans. Knowl. Data Eng.* 24 (8) (2012) 1478–1490.
- [34] G. Adomavicius, J. Zhang, Stability of recommendation algorithms, *ACM Trans. Inf. Syst.* 30 (4) (2012) 23.
- [35] Z. Yang, B.D. Davison, Distinguishing venues by writing styles, in: Proceedings of the 12th ACM/IEEE-CS Joint Conference on Digital Libraries, ACM, 2012, pp. 371–372.
- [36] Z. Yang, B.D. Davison, Venue recommendation: Submitting your paper with style, in: Machine Learning and Applications, ICMLA, 2012 11th International Conference on, Vol. 1, IEEE, 2012, pp. 681–686.
- [37] T. Huynh, K. Hoang, Modeling collaborative knowledge of publishing activities for research recommendation, *Comput. Collect. Intell. Technol. Appl.* (2012) 41–50.
- [38] J. Yu, K. Xie, H. Zhao, F. Liu, Prediction of user interest based on collaborative filtering for personalized academic recommendation, in: Computer Science and Network Technology, ICCSNT, 2012 2nd International Conference on, IEEE, 2012, pp. 584–588.
- [39] A.J. Trappey, C.V. Trappey, C.-Y. Wu, C.Y. Fan, Y.-L. Lin, Intelligent patent recommendation system for innovative design collaboration, *J. Netw. Comput. Appl.* 36 (6) (2013) 1441–1450.
- [40] M. Kochen, R. Tagliacozzo, Matching authors and readers of scientific papers, *Inf. Storage Retr.* 10 (5–6) (1974) 197–210.
- [41] M. Errami, J.D. Wren, J.M. Hicks, H.R. Garner, eTBLAST: a web server to identify expert reviewers, appropriate journals and similar publications, *Nucleic Acids Res.* 35 (suppl\_2) (2007) W12–W15.
- [42] M.J. Schuemie, J.A. Kors, Jane: suggesting journals, finding experts, *Bioinformatics* 24 (5) (2008) 727–728.
- [43] N. Kang, M.A. Doornenbal, R.J. Schijvenaars, Elsevier journal finder: recommending journals for your paper, in: Proceedings of the 9th ACM Conference on Recommender Systems, ACM, 2015, pp. 261–264.
- [44] W.H. Hsu, A.L. King, M.S. Paradesi, T. Pydimarri, T. Weninger, Collaborative and structural recommendation of friends using weblog-based social network analysis, in: AAAI Spring Symposium: Computational Approaches to Analyzing Weblogs, vol. 6, 2006, pp. 55–60.
- [45] T. Silva, J. Ma, C. Yang, H. Liang, A profile-boosted research analytics framework to recommend journals for manuscripts, *J. Assoc. Inf. Sci. Technol.* 66 (1) (2015) 180–200.
- [46] M.C. Pham, Y. Cao, R. Klamma, Clustering technique for collaborative filtering and the application to venue recommendation, in: Proc. of I-KNOW, Citeseer, 2010.
- [47] M.C. Pham, Y. Cao, R. Klamma, M. Jarke, A clustering approach for collaborative filtering recommendation using social network analysis, *J. UCS* 17 (4) (2011) 583–604.
- [48] H.P. Luong, T. Huynh, S. Gauch, K. Hoang, Exploiting social networks for publication venue recommendations, in: KDIR, 2012, pp. 239–245.
- [49] I. Boukhris, R. Ayachi, A novel personalized academic venue hybrid recommender, in: Computational Intelligence and Informatics, CINTI, 2014 IEEE 15th International Symposium on, IEEE, 2014, pp. 465–470.
- [50] E. Minkov, B. Charrow, J. Ledlie, S. Teller, T. Jaakkola, Collaborative future event recommendation, in: Proceedings of the 19th ACM International Conference on Information and Knowledge Management, ACM, 2010, pp. 819–828.
- [51] Z. Lu, N. Xie, W. Wilbur, Identifying related journals through log analysis, *Bioinformatics* 25 (22) (2009) 3038–3039.
- [52] A. Singhal, C. Buckley, M. Mitra, Pivoted document length normalization, in: Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, 1996, pp. 21–29.
- [53] Y. Sun, B. Norick, J. Han, X. Yan, P.S. Yu, X. Yu, Pathselclus: Integrating meta-path selection with user-guided object clustering in heterogeneous information networks, *ACM Trans. Knowl. Discov. Data* 7 (3) (2013) 11.
- [54] Y. Sun, J. Han, Mining heterogeneous information networks: a structural analysis approach, *ACM SIGKDD Explor. Newsl.* 14 (2) (2013) 20–28.
- [55] Y. Sun, J. Han, X. Yan, P.S. Yu, T. Wu, Pathsim: Meta path-based top-k similarity search in heterogeneous information networks, *Proc. VLDB Endow.* 4 (11) (2011) 992–1003.
- [56] A. Grover, J. Leskovec, Node2vec: Scalable feature learning for networks, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, 2016, pp. 855–864.
- [57] B. Zhu, S. Watts, H. Chen, Visualizing social network concepts, *Decis. Support Syst.* 49 (2) (2010) 151–161.
- [58] Y. Liang, Q. Li, T. Qian, Finding relevant papers based on citation relations, in: International Conference on Web-Age Information Management, Springer, 2011, pp. 403–414.
- [59] T. Opsahl, F. Agneessens, J. Skvoretz, Node centrality in weighted networks: Generalizing degree and shortest paths, *Social Networks* 32 (3) (2010) 245–251.
- [60] L.C. Freeman, Centrality in social networks conceptual clarification, *Social Networks* 1 (3) (1978) 215–239.
- [61] D.M. Blei, A.Y. Ng, M.I. Jordan, Latent dirichlet allocation, *J. Mach. Learn. Res.* 3 (Jan) (2003) 993–1022.
- [62] J.H. Lau, T. Baldwin, An empirical evaluation of doc2vec with practical insights into document embedding generation, 2016, arXiv preprint arXiv:1607.05368.
- [63] Q. Le, T. Mikolov, Distributed representations of sentences and documents, in: International Conference on Machine Learning, 2014, pp. 1188–1196.
- [64] V.D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, Fast unfolding of communities in large networks, *J. Stat. Mech. Theory Exp.* 2008 (10) (2008) P10008.
- [65] M.E. Newman, M. Girvan, Finding and evaluating community structure in networks, *Phys. Rev. E* 69 (2) (2004) 026113.
- [66] M.E. Newman, Analysis of weighted networks, *Phys. Rev. E* 70 (5) (2004) 056131.
- [67] R.A. Jarvis, E.A. Patrick, Clustering using a similarity measure based on shared near neighbors, *IEEE Trans. Comput.* 100 (11) (1973) 1025–1034.
- [68] K. Jones, S. Walker, S.E. Robertson, A probabilistic model of information retrieval: development and comparative experiments: Part 2, *Inf. Process. Manage.* 36 (6) (2000) 809–840.
- [69] M.F. Porter, Snowball: A Language for Stemming Algorithms, 2001.
- [70] R. Real, J.M. Vargas, The probabilistic basis of jaccard's index of similarity, *Syst. Biol.* 45 (3) (1996) 380–385.
- [71] I.A. Basheer, M. Hajmeer, Artificial neural networks: fundamentals, computing, design, and application, *J. Microbiol. Meth.* 43 (1) (2000) 3–31.
- [72] H. Tong, C. Faloutsos, J.-Y. Pan, Fast Random Walk with Restart and its Applications, IEEE, 2006.
- [73] S. Wu, C. Huang, L. Li, F. Crestani, Fusion-based methods for result diversification in web search, *Inf. Fusion* 45 (2019) 16–26.
- [74] S. Wu, Applying the data fusion technique to blog opinion retrieval, *Expert Syst. Appl.* 39 (1) (2012) 1346–1353.
- [75] D. Lillis, L. Zhang, F. Toolan, R.W. Collier, D. Leonard, J. Dunnion, Estimating probabilities for effective data fusion, in: Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, 2010, pp. 347–354.
- [76] J.A. Aslam, M. Montague, Models for metasearch, in: Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, 2001, pp. 276–284.



- [77] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, Z. Su, Arnetminer: extraction and mining of academic social networks, in: *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, 2008, pp. 990–998.
- [78] M. Deshpande, G. Karypis, Item-based top-n recommendation algorithms, *ACM Trans. Inf. Syst.* 22 (1) (2004) 143–177.
- [79] H. Stuckenschmidt, Approximate information filtering on the semantic web, in: *Annual Conference on Artificial Intelligence*, Springer, 2002, pp. 114–128.
- [80] M. Sokolova, G. Lapalme, A systematic analysis of performance measures for classification tasks, *Inf. Process. Manage.* 45 (4) (2009) 427–437.
- [81] E. Gibaja, S. Ventura, A tutorial on multilabel learning, *ACM Comput. Surv.* 47 (3) (2015) 52:1–52:38, <http://dx.doi.org/10.1145/2716262>, URL <http://doi.acm.org/10.1145/2716262>.
- [82] M. Kunaver, T. Požrl, Diversity in recommender systems—A survey, *Knowl.-Based Syst.* 123 (2017) 154–162.
- [83] G. Adomavicius, J. Zhang, Improving stability of recommender systems: a meta-algorithmic approach, *IEEE Trans. Knowl. Data Eng.* 27 (6) (2014) 1573–1587.
- [84] C. Desrosiers, G. Karypis, A comprehensive survey of neighborhood-based recommendation methods, in: *Recommender Systems Handbook*, Springer, 2011, pp. 107–144.