



Citation recommendation using semantic representation of cited papers' relations and content

Jinzhu Zhang^{a,*}, Lipeng Zhu^{a,b}

^a Department of Information Management, School of Economics & Management, Nanjing University of Science and Technology, Nanjing, China

^b Department of Information, Jinhu County People's Hospital, Huai'an, China

ARTICLE INFO

Keywords:

Citation recommendation
Cited paper
Co-citation
Citation content
Semantic representation

ABSTRACT

Citation recommendation can help researchers quickly find supplementary or alternative references in massive academic resources. Current research on citation recommendation mainly focuses on the citing papers, resulting in the enormous cited papers are ignored, including the relations among cited papers and their citation context cited in citing papers. Moreover, cited paper's content is often denoted with its original title and abstract, which is hard to acquire and rarely considers different citation motivations. Furthermore, the most appropriate method for semantic representation of cited papers' relations and content is uncertain. Therefore, this paper studies citation recommendation from the perspective of semantic representation of cited papers' relations and content. Firstly, four forms of citation context are designed and extracted as cited papers' content considering citation motivations, as well as co-citation relationships are extracted as cited papers' relations. Secondly, 132 methods are designed for generating semantic vector of cited paper, including four network embedding methods, 16 methods by combining four text representation algorithms with four forms of citation content, and 112 fusion methods. Finally, similarity among cited papers is calculated for citation recommendation and a quantitative evaluation method based on link prediction is designed, to find the most appropriate form of citation content and the optimal method. The result shows that doc2vecC (Document to Vector through Corruption) with the form of CS&SS (Current Sentences and Surrounding Sentences) performs best, in which the AUC (Area Under Curve) and MAP (Macro Average Precision) reach 0.877 and 0.889 and have increased by 0.462 and 0.370 compared with the worst-performing method. This performance is slightly improved by parameters adjustment, and a case study is performed whose results have further proved the effectiveness of this method. In addition, among four forms of cited papers' content, CS&SS performs best in almost all methods. Furthermore, the fusion methods not always perform better than the single methods, where doc2vecC (CS&SS) performs better than the best fusion method GCN (Graph Convolutional Network). These results not only prove the effectiveness of citation recommendation from the perspective of cited paper, but also provide helpful and useful suggestions for method selection and citation content selection. The data and conclusions can be extended to other text mining-related tasks. Simultaneously, it is a preliminary research which needs to be further studied in other domains using emerging semantic representation methods.

1. Introduction

When researchers write papers, they need to cite relevant references to help readers and themselves comprehensively understand the background, the current status and the innovation point. However, as increasing scientific papers created, researchers have to face challenges of academic information overload (Liu, Yu, et al., 2014), and are forced to spend a lot of time on searching for related references through reading and analysis (Kobayashi et al., 2018). More important, the references

may be limited to the authors' familiar journals and fields, which may result in absence or misquote of some relevant literatures. Therefore, how to help researchers quickly find supplementary or alternative references in massive academic resources is becoming an important issue (Ali, Qi, Kefalas, et al., 2020). Citation recommendation happens to be an effective approach to address this problem, which can not only recommend relevant and suitable references automatically, but also reduce missing citations of some important literatures (Ma et al., 2020).

Current citation recommendation approaches can be divided into

* Corresponding author.

E-mail address: zhangjinzhu@njust.edu.cn (J. Zhang).

<https://doi.org/10.1016/j.eswa.2021.115826>

Received 28 May 2021; Received in revised form 29 August 2021; Accepted 29 August 2021

Available online 11 September 2021

0957-4174/© 2021 Elsevier Ltd. All rights reserved.

three categories, including collaborative filtering-based, content-based and network-based citation recommendation (Kong et al., 2018). Collaborative filtering-based approach often recommends similar literatures which have been cited by other similar authors but not cited by the current author, and the similarity among authors is often computed by co-authorship or citation relationship (Alhijawi & Kilani, 2020). However, it rarely recommends less cited references, resulting in hardly recommending the latest literature. In addition, it considers little text content of both citing papers and cited papers. The content-based citation recommendation approach considers the text content represented by keywords, titles and abstracts of the citing papers and applies text representation methods for semantic representation of these content, which can recommend similar and latest literature by computing the text similarity (Yang et al., 2018). However, this approach cannot distinguish the importance of literatures with similar content, and do not consider the text content cited in citing papers which can give more detailed descriptions of citation motivations. The network-based citation recommendation approach uses various association relationship to build network models to learn semantic representation of papers, which can calculate the importance of each paper in the network through various network-based indicators. It converts citation recommendation problem into link prediction problem in the citation network or co-citation network (Strohman et al., 2007). However, this method focuses on the citation relationships come from general datasets for validation of the improvements of the algorithms, and pays too much attention on recommending citing paper. In addition, the text content of both citing papers and cited papers, especially the content of cited papers, are ignored in most of the times (Ganguly & Pudi, 2017).

Besides above methods, researchers have proposed some fusion methods from the perspective of comprehensive utilization of network structure and text content, e.g., TADW (Text-Associated Deep Walk) model based on matrix factorization techniques (Yang et al., 2015) and CANE (Context-Aware Network Embedding) model based on mutual attention mechanism (Tu et al., 2017), which have been proven to be effective in many tasks such as classification, label prediction and link prediction (Eto, 2019; Pan et al., 2016). Among them, some of the models can be applied for semantic representation of cited papers, which is an important component for subsequent similarity computation in citation recommendation. However, the effect of fusion method is not always better than that of the single method in specific domain, for example, the single method GraRep performs better than the fusion method naïve combination in text classification, and the single method deepwalk outperforms the fusion method TADW in link prediction (Ganguly & Pudi, 2017; Zhang et al., 2016). Whether this phenomenon happens in citation recommendation is uncertain. In addition, these fusion methods rarely consider the text content cited in citing papers which gives more detailed citation motivations and may be helpful for citation recommendation. Therefore, when the fusion methods are applied to a specific task, they need to be compared with the single methods to find the most suitable method. Moreover, the text content cited in citing papers needs to be considered in these fusion methods. In summary, this paper mainly addresses the following three problems.

- (1) Current approaches mainly recommend literatures from the perspective of citing papers. However, the larger volume of relationships among cited papers and the citation context are rarely considered. Therefore, this paper makes full use of the co-citation relationships among references to denote cited paper's relations and the text content cited in papers to denote cited papers' content for citation recommendation. It could provide a new perspective of citation recommendation using cited papers' relations, content, and their combination.
- (2) Cited paper's content is often denoted with its original title and abstract in few public datasets. However, these cited papers' content is hard to acquire, and not totally accurate and comprehensive because references are cited along with multiple citation

motivations. Therefore, this paper designs multiple forms of citation content for distinguishable description of citation motivation, and finds the most appropriate form for citation recommendation through quantitative comparison. These designed forms of citation content could be not only used for citation recommendation, but also extended to other text mining-related tasks, such as topic identification, topic evolution analysis and citation sentiment analysis.

- (3) Representation learning method performs well in many semantic representation tasks, but they are rarely applied to different forms of citation content. In addition, the most appropriate method for semantic representation of cited papers' relations and content is uncertain, because the fusion method sometimes performs worse than the single method. Therefore, this paper carefully selects and applies several representative and latest representation learning algorithms for semantic representation of cited papers' relations, content, and their combinations. Through quantitative comparison, it could find the most effective method for citation recommendation and clarify the function of different forms of citation content when combining with cited papers' relations. These results could provide helpful and useful suggestions for choosing suitable method in citation recommendation.

Therefore, we study citation recommendation from the perspective of semantic representation of cited papers' relations and content. Firstly, we parse the full-text open-source publications and design multiple forms of citation context as cited papers' content considering citation motivations, as well as extract co-citation relationships as cited papers' relations. Secondly, we carefully select various semantic representation learning methods and apply them on cited papers' relations and content, which could generate semantic vector of the cited paper. Finally, similarity among cited papers is calculated for citation recommendation and a quantitative evaluation method based on link prediction is designed to find the most appropriate form of citation content and the optimal method, and important parameters of the optimal method are adjusted.

The rest of this paper is organized as follows. Section 2 reviews the related work. Section 3 provides a detailed description of the method. Section 4 describes the experiments and results. Section 5 concludes the findings and presents the future work.

2. Related works

2.1. Citation recommendation

Citation recommendation can be classified into two categories based on the length of input text (Färber & Jatowt, 2020). This paper uses both citation context and network structure, so we divide the current citation recommendation approaches into three categories based on the adopted techniques, including collaborative filtering-based, content-based and network-based citation recommendation (Kong et al., 2018).

The collaborative filtering-based approach recommends literatures with similar authors by computing the similarity among authors according to co-authorship and citation relationship (Zarrinkalam & Kahani, 2012). The traditional collaborative filtering-based approach depends too much on the occurrence frequency of papers, which is not conducive to recommending the latest literature or less cited literature (Guns & Rousseau, 2014). Therefore, some researchers have made some improvements on the collaborative filtering-based approach to solve this problem. Gipp et al. (2009) propose a content-based collaborative filtering algorithm, which can directly capture the semantic connection among papers to recommend similar and latest literature. Wang and Blei (2011) propose a collaborative filtering-based topic regression model for citation recommendation, which combines the advantages of traditional collaborative filtering-based approach and probabilistic topic models, using contents of citing papers and ratings of other users. Liu et al. (2015) use the citation context to mine co-occurrence relationships

between citing papers, and regard this co-occurrence relationship as related information, which is added to the collaborative filtering-based model as supplementary information to improve the accuracy of recommendation. However, these methods focus on the citing papers and the enormous cited papers are often ignored, in addition, these methods rarely consider the text content of both citing papers and cited papers.

The content-based approach considers the text content represented by keywords, titles and abstracts of the citing papers, and applies text representation methods for semantic representation, to recommend similar and latest literature by computing the text similarity (Yang et al., 2018). Bhagavatula et al. (2018) use a neural network model to embed the citing papers into a vector space by encoding the titles and abstracts, then use their nearest neighbors as candidate citations and re-rank the candidate citations using a discriminative model. Nogueira et al. (2020) use BERT-based ranking models to represent the titles and abstracts of citing papers, which is able to preserve topical relevance of model outputs. In addition, they provide an extensive evaluation of domain adaptation techniques for pre-trained language models, concluding that in-domain pre-training vocabulary can greatly improve the effectiveness. Sugiyama and Kan (2010) propose a general model that not only considers the titles and abstracts of citing papers, but also considers the titles and abstracts of some cited papers. The experimental results prove that text content is crucial to the improvement of citation recommendation accuracy. Liu, Chen, et al. (2014) use the current sentence where citations appear as the citation content and recommend literatures by querying citation content related to search terms, which proves that citation content could provide important information for citation recommendation. He et al. (2010) define the 50 words before and 50 words after the placeholder as the citation content, and design a novel non-parametric probabilistic model which can measure the relevance between the citation content and papers, thus recommending a set of high-quality papers for researchers. However, these methods cannot distinguish the importance of literatures with similar content, and do not consider multiple forms of citation context cited in papers which give more detailed descriptions of different citation motivations.

In recent years, an increasing attention has been paid to network-based approach, which uses various association relationship of datasets to build network models and can easily calculate the importance of each paper in the network through various network-based indicators (Jeong et al., 2020). It converts citation recommendation problem into link prediction problem in citation network and co-citation network (Strohman et al., 2007). Son and Kim (2018) combine citation analysis with network analysis and propose a multi-level citation network model that can quantify the direct, indirect, and overall relationship among papers. Cai et al. (2018) construct a three-layer network model, including author layer, dissertation layer, and publisher layer, and the effectiveness has been proven on multiple data sets, i.e., ACL anthology network, DBLP, and CiteSeer ML datasets. Zhou et al. (2008) propose a decomposition strategy network model, which use citation relationship and co-authorship to measure the similarity among papers. West et al. (2016) propose a hierarchical clustering algorithm based on citation relationship to determine the relevance among papers, and make recommendations according to their importance in these clusters. Chen et al. (2019) propose a novel neural network model called CIREC based on citation tendency, which constructs a weighted heterogeneous network and designs a biased random walk procedure to efficiently explore papers' characteristics. Ali, Qi, Muhammad, et al. (2020) present a weighted probabilistic heterogeneous network model termed as PR-HNE, which jointly learns researchers and papers' dynamics by encoding information from six information networks into a joint latent space. However, these methods focus on the citation relationships come from general datasets for validation of the improvements of the algorithms, and paying too much attention on recommending citing paper. In addition, the text content of both citing papers and cited papers, especially the content of cited papers, are ignored in most of the times.

In summary, the collaborative filtering-based approach is often mixed with other approaches to solve the problem of limitations on hardly recommending the latest or less cited literatures. The content-based approach can solve the above problem by calculating text similarity. The network-based approach can distinguish the importance of literatures with similar content by various network-based indicators. However, the common point of these studies is recommending citations from the perspective of citing paper, resulting in the enormous cited papers are ignored, including the relations among cited papers and the multiple forms of citation context of cited papers. Therefore, this paper studies citation recommendation changing from the perspective of citing paper to cited paper. Specifically, this paper intends to make use of the co-citation relationships among cited papers to denote cited paper's relations and designs distinguishable and detailed representation form of citation content to denote cited papers' content, to expand the breadth and depth of citation recommendation from the perspective of cited paper.

2.2. Extraction of citation content

The existing citation content extraction methods can be categorized into two types (Dai et al., 2020). The first type of methods extracts several fixed number sentences around manually defined citation marks and the second extracts citation content automatically by using classification methods.

The first type of methods is simple to implement and often used by most researches. Zhou and Zhang (2019) extract the current sentence where citations appear according to manually defined citation marks and defined it as the citation content, and conduct fine-grained citation content analysis automatically for book impact assessment. G et al. (2019) add the previous and the next sentence as the citation content, and propose a system for ranking journals based on citation content.

The second type of methods can gain more accurate extraction results of citation content, but the process of implementation is much more complicated. This method first divides the query text into fixed-size segments with separators, and then uses classification methods or probability model to determine whether the segment is part of citation content based on specific features. Angrosh et al. (2010) find some features of the context, such as background, theme, advantages, disadvantages, results, methods and other keyword features based on template, and then use CRF classifier for training to filter sentences that can represent citation content. Abu-Jbara and Radev (2012) first convert the sentences into words and determine whether each word belongs to the context of the target citation through classifier. Then, they convert the extraction of citation content into a sequence labeling problem and assign categories to each word. Finally, they divide the sentence into different segments and determine the scope of the citation content according to the segment information. Angrosh et al. (2013) develop an automatic citation content extraction system named CitContExt, which defines three features of citation content and applies a classification model CRF (Conditional Random Field) for citation content recognition. Jha et al. (2017) use two classes of features, i.e., features inherent from a sentence itself and features based on relationships, to determine citation content. Subsequently, MRFs (Markov Random Fields) classification model is used to distinguish the sentences related or unrelated to citation content.

In summary, the above two types of methods can achieve good results in extraction of citation content and provide data source for other applications. When citation content is used for citation recommendation, it should be represented with different combination of surrounding sentences considering different citation motivations. Therefore, this paper designs multiple forms of citation content as cited papers' content by combining the current sentences and the surrounding sentences with different strategies, to find the best form of citation content.

2.3. Representation learning

In the task of citation recommendation, the semantic representation of cited papers' relations and content usually implement by representation learning. Representation learning aims to represent the semantic information of the object as low-dimensional and dense vectors, which has been applied to many tasks such as link prediction, technology opportunity discovery and performance prediction (Esen et al., 2008; Esen et al., 2009; Ghorbanzadeh et al., 2021; Zhang & Yu, 2020). The representation learning mainly involves network representation learning, text representation learning and the fusion of text and network representation learning (Jin & Srihari, 2007; Yin et al., 2019). It also involves the representation learning of other objects such as pictures, audios, and videos (Ye et al., 2020).

Network representation learning can represent the nodes in network as low-dimensional and dense vectors. Inspired by the word2vec proposed by Mikolov, Chen, et al. (2013), Perozzi et al. (2014) treat the nodes in network as words in sentences, and generate node sequences through random walks, and then use word2vec to learn the representation of nodes. Grover and Leskovec (2016) design a flexible neighborhood sampling strategy named node2vec that improves the process of random walk by introducing breadth-first sampling and depth-first sampling. Tang, Qu, Wang, et al. (2015) propose the LINE (Large-scale Information Network Embedding) model, which solves the problem that the random walk does not have a clear objective function for the network structure, and propose the first-order proximity and second-order proximity in the network. Tang, Qu and Mei (2015) improve LINE to handle heterogeneous networks, and propose a semi-supervised PTE (Predictive Text Embedding) model for text label prediction tasks, which can transform a set of documents with known partial labels into a heterogeneous network containing three types of nodes, including documents, words, and labels. Wang et al. (2017) jointly optimize NMF-based representation learning model and modularity-based community detection model in a unified framework, which enables the learned representations of nodes to preserve both microscopic and community structures.

Text representation learning is an important branch of representation learning, thus forming a variety of text representation methods and models. Word2vec proposed by Mikolov, Sutskever, et al. (2013) sets off the wave of text representation learning, which trains the corpus through neural network model and combines the context information of each word to map word into a fixed-dimensional vector. Then, Le and Mikolov (2014) propose an unsupervised doc2vec model of which long text is introduced into the corpus as a special paragraph ID and fully combine text context, word order, and paragraph features. Tang, Qin, et al. (2015) use CNN (Convolutional Neural Networks) and LSTM (Long Short-Term Memory) to model sentences, and then use Bi-RNN (Bidirectional Recurrent Neural Networks) to model documents, which proves that GRU (Gated Recurrent Unit) has better effect than RNN on document-level emotion classification tasks. Yang et al. (2016) add attention mechanism when modeling the sentences and documents, and model the importance of words and sentences in documents. Then document vectors are generated by making weighted calculation. Chen (2017) generates document representation by average word embedding, which introduces data dependent regularization to suppress the embedding of common and non-discriminatory words.

In recent years, many researchers have begun to study the fusion learning, which jointly utilizes the network structure information and the text information to learn the representation of nodes. Yang et al. (2015) propose a TADW model, which trains the joint modeling of network structure features and text features through matrix factorization. Zhang et al. (2016) construct a HSCA (Homophily, Structure, and Content Augmented) model, which integrates homogeneity, topology structure, and node content, and then adds three information sources to an objective function to learn network representation together. Tu et al. (2017) present a CANE model, which learns context-aware embedding

for vertices with mutual attention mechanism and models the semantic relationships between vertices more precisely. Li et al. (2017) propose a PPNE (Property Preserving Network Embedding) model to effectively fuse different types of node attribute information, which transforms the learning process of representation vectors into a joint optimization problem, and solves this joint optimization problem by using an effective stochastic gradient descent algorithm. Kipf and Welling (2017) propose a scalable approach, i.e., GCN (Graph Convolutional Network), which learns hidden layer representations that encode both local graph structure and text features of nodes through graph convolution.

In summary, the representation learning method performs excellent in many semantic representation tasks and can be applied to citation recommendation. However, the most appropriate method for semantic representation of cited papers' relations and content in this task is uncertain, because the effect of fusion method is not always better than that of the single method in specific domain (Zhang et al., 2016). Therefore, it is necessary to apply and compare various semantic representation learning methods to find the optimal method for citation recommendation, which can provide helpful and useful suggestions for method selection.

3. Method

This paper takes literatures published in PLOS ONE as the experiment data. By downloading and parsing the data, the co-citation relationships are generated to denote cited paper's relations, as well as the citation content is extracted, including the current sentence and the surrounding sentences. Subsequently, four forms of the citation content are designed to denote cited papers' content, by combining the above sentences with different strategies. Then, 132 methods are designed for generating semantic vector of cited paper, including four network embedding methods, 16 methods by combining four text representation algorithms with four forms of citation content, and 112 fusion methods. Finally, cosine similarity is used for similarity calculation among cited papers for citation recommendation and a quantitative evaluation method based on link prediction is designed to find the most appropriate form of citation content and the optimal method for citation recommendation.

3.1. Extraction of co-citation relationships and citation content

The data related to co-citation relationship and citation content can be obtained by parsing papers with XML format. In this paper, the method of ElementTree in Python package is selected to parse the full-text papers with XML format, which performs faster and consumes less memory compared to other methods. Then co-citation relationships and citation content are extracted by locating the corresponding XML tags.

3.1.1. Extraction of co-citation relationships

The co-citation relationships are generated through projection of citation relationships between citing papers and cited papers. Therefore, this paper first locates the corresponding tags of both citing paper and cited paper, and then citation relationships could be generated and the co-citation relationships can be formed.

Upon careful inspection, this paper uses the tag "<article-title></article-title>" to locate citing papers because other tags are sometimes missing. The corresponding information of cited papers are stored in the tag "<ref-list></ref-list>". According to different types of cited papers, including paper, book and so on, the corresponding titles are usually stored in two tags, i.e., "<article-title></article-title>" and "<chapter-title></chapter-title>", respectively. In addition, some titles of a small number of papers or books are store in "<mixed-citation></mixed-citation>", which are mixed with author, publication, and other information. If none of these three tags are found, it indicates that the title of the cited papers is missing with some reasons.

3.1.2. Extraction of citation content

In the process of citation content extraction, the sentences with citation marks are found and the corresponding sentences are obtained. Then, the citation numbers in the above sentences are extracted. Finally, the citation content and the cited paper's title could be connected through citation number matching.

i. The places where the citation marks appeared are found through XML tags and the corresponding sentence are extracted as citation content. The citation mark is stored in the tag "<xref ref-type='bibr'></xref>", and the sentence including this tag is considered as the current sentence. Thus, the previous and the next sentences can be obtained according to the order of sentences.

ii. The citation numbers in the above sentence are extracted by regular expression. The citation number is usually denoted with various expressions, e.g., "[1], [1–3], [1,3–5]", which either as a single number or a combination with "-" and ",". This paper takes two regular expressions for fully extracting the citation numbers, i.e., "[\d +]" and "[\d +]-[\d +]".

iii. The corresponding title of cited paper is obtained through citation number, which occurs in both citation content and the references list. First, the citation number and the corresponding title in references list is extracted which often stored in tag "<label><label>". Then the citation number in both citation content and references list are matched, thus the citation content could correlate with the cited paper's title.

3.2. Different representation forms of citation content

In the specific task of citation recommendation, citation content should be designed with multiple forms for better description of different citation motivations. In this paper, citation content is related to the current sentence, the previous sentence and the next sentence, which can be combined into different forms of citation content with different strategies. Specifically, this paper designs four forms of citation content as follows in order to make detailed description of citation motivations and find the most appropriate form of citation content.

- (1) Current sentences (CS): Current sentences is the most direct text that related to the citation motivation, which can describe the reason why a reference is cited. Therefore, the first form of citation content is denoted by the current sentence where the citation appears. In many cases, a reference could be cited more than one time in a paper and may be cited by multiple different papers, so a reference has several related current sentences. Considering the simplicity and effectiveness, this paper combines these sentences together by simple connection to represent the content of the same reference. This strategy is also used in many citation content extraction tasks (Liu, Chen, et al., 2014).
- (2) Current sentences and surrounding sentences (CS&SS): The surrounding sentences are composed of the previous sentences and the next sentences, which may have a strong semantic correlation with the current sentences. It could be a supplement to the text that related to the motivation why a reference is cited. Therefore, the second form of citation content is denoted by the current sentences, the previous sentences and the next sentences together. This form can greatly increase the number of sentences of citation content. As with the above form, this paper combines these sentences together by simple connection.
- (3) Current sentences, and surrounding sentences that do not cite other references (CS&SS-): Some previous and next sentences may cite other references, and the motivation of these sentences may be different from the current sentences. Therefore, the third form of citation content is denoted by current sentences and the surrounding sentences that do not cite other references. This form can be compared with the second form to find out which form is more accurate. As above, these sentences are combined together by simple connection.

- (4) Automatic summarization of CS&SS (SuCS&SS): On the one hand, the number of occurrences of different references may vary greatly, e.g., one reference may appear 4 time in a paper and 10 times in other papers, and another reference may just appear 1 time in only one paper. On the other hand, as mentioned in the third form, some previous and next sentences may cite other references, and the motivation of these sentences may be different from the current sentences. Therefore, the sentences should be summarized with unified length and retain the most important sentences that can better reflect citation motivation. Therefore, the fourth form selects the most import sentences in the current sentences and the surrounding sentences by using an automatic summarization method. According to the survey, TextRank is usually used for automatic summarization because of its simplicity and effectiveness (Balcerzak et al., 2014), so this paper selects TextRank as the automatic summarization method.

3.3. Semantic representation of cited papers' relations and content

The semantic representation of cited papers' relations and content usually implement by representation learning. However, different categories of representation learning methods have different performance in specific tasks and the most effective method should be determined with comparison. Therefore, multiple representative and typical algorithms of network representation, text representation and fusion methods are carefully selected, categorized and introduced. These algorithms are then used for combining with different forms of cited papers' relations and content, thus constructing 132 methods for generating semantic vector of cited paper.

(1) Semantic representation of cited papers' relations based on network representation learning

Different network representation learning methods can generate semantic representations of network structure and are suitable for different application tasks, such as node classification and link prediction. Thus the most appropriate method for the task of citation representation may uncertain. On the basis of existing research, this paper selects two categories of network representation learning algorithms, including deepwalk, node2vec, LINE and SDNE. These four algorithms are then used for combining with co-citation network, thus constructing four methods for generating semantic vector of cited paper.

The first category is based on random walk. It first generates sequences of nodes through random walks, then use word2vec (Mikolov, Chen, et al., 2013) to learn latent representations by treating nodes in sequences as words in sentences. Deepwalk (Perozzi et al., 2014) and node2vec (Grover & Leskovec, 2016) are the two representative methods. Deepwalk is the most classic method, which select the next node in a random walk by uniformly random distribution. Node2vec improve the process of random walk by introducing breadth-first sampling and depth-first sampling. Therefore, this paper selects both deepwalk and node2vec for semantic representation of cited papers' relations in this category.

The second category preserves the first-order proximity and the second-order proximity of nodes by defining the loss function, which preserves both the local and global network structure. LINE (Tang, Qu, Wang, et al., 2015) and SDNE (Wang et al., 2016) are the two representative methods. The differences between them are the optimization ways of the first-order proximity and the second-order proximity. SNDE uses an automatic encoder structure to jointly optimize the first-order and second-order proximity, while LINE optimizes them separately. Both of these two methods perform well in many network semantic representation tasks (Yin et al., 2019). Therefore, this paper selects both LINE and SDNE for semantic representation of cited papers' relations in this category.

(2) Semantic representation of cited papers' content based on text representation learning

Different text representation methods can generate semantic

representations of cited papers' content in their own ways and are suitable for different application tasks, such as topic analysis and text classification. Thus the best text representation method for the task of citation content semantic representation may uncertain. Based on the existing research, this paper selects two main categories of text representation learning methods, including LSI, LDA, doc2vec and doc2vecC. These four algorithms are then used for combining with four forms of cited papers' content, thus constructing 16 methods for generating semantic vector of cited paper.

The first category is based on the topic model. It is often selected as an important baseline for comparison with other algorithms, which captures semantic information from the topic level and calculates probability of the text belongs to each topic. LSI (Latent Semantic Indexing) (Deerwester et al., 1990) and LDA (Latent Dirichlet Allocation) (Blei et al., 2003) are the two representative methods. LSI is based on the singular value decomposition, which can capture the grammatical features of text to extract topics. LDA is a generative probabilistic model of a corpus. The basic idea is that documents are represented as random mixtures over latent topics, where each topic is characterized by a distribution over words. These two methods perform well in text semantic representation tasks (Park et al., 2015; Poria et al., 2016). Therefore, this paper selects both LSI and LDA for semantic representation of cited papers' content in this category.

The second category is based on neural network, which generates text vector by training neural network models. Doc2vec (Document to Vector) (Le & Mikolov, 2014) and doc2vecC (Document to Vector through Corruption) (Chen, 2017) are the two representative methods. Doc2vec is the most representative method in text representation learning and is widely used in the field of natural language processing. It is concatenated with the word vectors to predict the next word in a context, and it is applicable to texts of any length (sentences, paragraphs, and documents). In addition, doc2vec performs well in text representation tasks, such as sentiment analysis and text classification (Bilgin & Şentürk, 2017; Trieu et al., 2017). Doc2VecC proves that document representation generated by average word embedding can capture good semantics of the document during the learning process. It also introduces data dependent regularization, which tends to provide informational or rare words, while suppressing the embedding of common and non-discriminatory words. More important, Doc2VecC performs well in sentiment analysis, document classification as well as other semantic representation related tasks (Chen, 2017). Therefore, this paper selects doc2vec and doc2vecC for semantic representation of cited papers' content in this category.

(3) Semantic representation of cited papers' relations and content based on fusion methods

Different fusion methods can comprehensively generate semantic representations of cited papers' relations and content, and have different performance in specific tasks. Therefore, the best fusion method for the task of citation representation may uncertain. According to the survey, this paper selects three categories of classical and effective fusion algorithms, including naïve combination, TADW, Paper2vec and GCN. Specifically, naïve combination could combine four network embedding algorithms with four text representation learning algorithms mentioned above, so it consists 16 algorithms. The other three fusion methods could only combine with their build-in network embedding method, so they consist 12 algorithms together. These 28 algorithms are then used for combining with four forms of cited papers' content, thus constructing 112 methods for generating semantic vector of cited paper.

The first category is based on vector stitching, which connects text features with network features. Naïve combination is the most representative methods, which concatenates the vectors from both text features and network features and is often considered as an important baseline for other fusion algorithms. Therefore, this paper selects naïve combination for semantic representation of cited papers in this category.

The second category treats text content as special information to extend node's attributes in the network. TADW (Yang et al., 2015) and

Paper2vec (Ganguly & Pudi, 2017) are two representative methods. They have differences in the way of adding text information to the network. TADW introduces text information into the network based on matrix factorization, which is the first time that researchers have tried to fuse network information with text information. Paper2vec creates artificial text-based links through text features and add them into network for training to jointly utilize the network structure and text content. They perform well in many semantic representation tasks (Ganguly & Pudi, 2017; Kim et al., 2017). Therefore, this paper selects both TADW and Paper2vec for semantic representation of cited papers in this category.

The third category is based on convolutional neural network, which not only considers the links between nodes, but also combines the content information of each node. It learns semantic representation of node by convolving the feature information in the network, which is also considered as a fusion method. GCN (Kipf & Welling, 2017) is the most representative method. Firstly, each node sends its transformed characteristic information to its neighbor nodes, and then each node aggregates the characteristics of these neighbor nodes, and finally the final node representation is obtained through a non-linear transformation. GCN has been proven to perform well in many semantic representation tasks, such as classification and recommendation (Jeong et al., 2020; Zhang et al., 2021). Therefore, this paper selects GCN for semantic representation of cited papers in this category.

3.4. Similarity calculation among cited papers for citation recommendation

After semantic representation of cited papers' relations and content, each cited paper could be denoted as a semantic vector. Then, multiple vector similarity indicators can be used for similarity calculation among cited papers. Cited papers pairs with higher similarity and have not been co-cited in citing papers currently will be recommended as the candidates for citation recommendation.

This paper selects the commonly used cosine similarity to calculate the similarity among cited papers. It has been proved effective in multiple tasks, such as semantic extraction and scientific community detection (Carusi & Bianchi, 2019; Greiner-Petter et al., 2020). For the convenience of explanation, the learned semantic vectors of cited papers are denoted as $\vec{x} = (x_1, x_2, x_3, \dots, x_i, \dots, x_n)^T$ and $\vec{y} = (y_1, y_2, y_3, \dots, y_i, \dots, y_n)^T$, the similarity among all cited papers except the links in the training set can be calculated by formula (1).

$$\text{sim}(\vec{x}, \vec{y}) = \cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \|\vec{y}\|} = \frac{\sum_{i=1}^n x_i \times y_i}{\sqrt{\sum_{i=1}^n (x_i)^2} \times \sqrt{\sum_{i=1}^n (y_i)^2}} \quad (1)$$

3.5. Quantitative evaluation based on link prediction

The method based on link prediction is designed to quantitatively evaluate the performance of citation recommendation, to find the optimal method and most appropriate form of citation content. Suppose the corresponding recommendation of reference i is j , which means they are similar and have not been co-cited currently. If they are co-cited in the future, it can be considered as a successful prediction of citation recommendation.

Specifically, the data in the previous period is considered as the training set, and the data in the next period is considered as testing set. The testing set is regarded as the positive samples and the same number of negative samples is generated for fair evaluation. Then the AUC (Area Under Curve) and MAP (Macro Average Precision) could be computed to quantitatively evaluate the results of citation recommendation (Zhang, 2017).

AUC measures the overall accuracy of link prediction results. The similarity score of a link in positive sample is compared with that of a link in negative sample. The comparison is conducted independently n

times. Suppose there are n' times when the link in positive sample has a higher similarity score than the link in negative sample, and n'' times when the links in positive sample have the same similarity score with links in negative sample. The AUC value can be calculated with formula (2).

$$AUC = (n' + 0.5n'')/n \quad (2)$$

MAP is the macro average precision, which denotes as a weighted mean of precisions at a precision-recall curve. MAP is a comprehensive indicator of precision and recall which can better evaluate the overall performance of the recommendation algorithm. It is computed from the similarity score of the links in the positive and negative sample. It is defined as formula (3), in which P_n and R_n are the precision and recall at the n^{th} threshold.

$$MAP = \sum_n (R_n - R_{n-1})P_n \quad (3)$$

4. Results

Firstly, this paper evaluates the effects of four network representation methods on cited papers' relations to find the best one. Then, the effects of 16 text representation methods are evaluated to find the best form of citation content, the most appropriate text representation algorithm, and the best combination of them. Thirdly, 112 fusion methods which joint use cited papers' relations and multiple forms of citation content are evaluated to find the most appropriate one. Finally, some important parameters in the optimal method are adjusted to improve the performance of citation recommendation, and the optimal method is chosen to recommend top k candidates for case study.

4.1. Dataset description

In the process of data preparation, full-text papers with XML format are downloaded and parsed, and then stored in MySQL database. After that, the data are preprocessed, and the co-citation relationships and the citation content are extracted. Finally, the data are divided into training set and testing set for fair comparison.

(1) Data collecting and parsing

PLOS ONE accepts research papers in over 200 disciplines in the fields of science, engineering, medicine, and related social and human sciences. Specifically, all papers in PLOS ONE are open access to all researchers and their full-text data with XML format could be freely downloaded. We collected a research dataset in PLOS ONE under the subject area of AI (artificial intelligence) in Nov 25, 2018, where 1,675 full text papers are included and 5 correction papers are eliminated.

Then, the data related to co-citation relationship and citation content are extracted through data parsing. The parsed data is then stored in the MySQL database for searching and processing. A total of 115,653 citation relationships are obtained, in which part of citation relationships may be repeated because a reference may be cited multiple time by the same citing paper. In addition, a total of 346,927 sentences are obtained, including 115,653 current sentences, 115,653 previous sentences and 115,621 next sentences. The number of the next sentences is less than the number of the current sentences and previous sentences because the current sentences may appear at the end of the paper.

(2) Data preprocessing

In the preprocessing of co-citation relationship generation, the citing papers and cited papers with missing title are removed first. Then, one-to-many relationships between citing papers and cited papers are generated through SQL query among tables. Finally, a bipartite network between citing papers and cited papers is formed by using networkx package in Python, and the co-citation relationship is generated through projection of the bipartite network.

In the preprocessing of citation content, the cited papers with missing title are removed first. Then, the mapping relationship between

cited papers' title and citation content is formed through SQL query among tables. Due to subsequent need for text representation of the citation content, word segmentation is performed and stop words are removed. In addition, case conversion and word stemming are performed for unifying the same word. Finally, the sentences with a character length of less than 25 are removed, because these sentences usually contain few semantic information through manual observation.

After preprocessing, a total of 71,677 citation relationships were remained, as well as 109,837 current sentences, 109,837 previous sentences, and 109,805 next sentences.

(3) Data division

The data set is divided into training set and testing set according to the publication year of the citing papers. The papers published from 2007 to 2016 are regarded as training set and the papers published from 2017 to 2018 are regarded as testing set. According to the papers in training and testing set, the corresponding cited papers' relations and content are treated as training and testing set respectively. The details of divided datasets are shown in Table 1.

The co-citation relationships in training set are obtained by the projection on citation network. The final testing set is obtained through following steps. First, the testing set is formed, including 542,554 co-citation relationships projected from the citation network. Then, all unlinked edges in the co-citation network generated from training set are obtained. Finally, the intersection of all unlinked edges with the testing set is considered the final testing set.

4.2. Citation recommendation using semantic representation of cited papers' relations

This paper takes co-citation relationships among references to denote cited papers' relations as the model input and selects four different network representation methods for citation recommendation. For fair comparison, two important parameters are set uniformly in the four methods, i.e., "window = 10" and "representation-size = 128", where "window = 10" means taking ten citations around the current position as context information and "representation-size = 128" means each citation is represented as a point in the 128-dimensional space. The rest parameters are assigned with the default values of the models (Grover & Leskovec, 2016; Perozzi et al., 2014; Tang, Qu, Wang, et al., 2015; Wang et al., 2016). Table 2 shows the AUC and MAP of four network representation learning methods for citation recommendation.

The result shows that deepwalk performs best in all methods whose accuracy is marked in bold in Table 2. Compared with LINE which performs worst, the AUC and MAP of deepwalk have increased by 0.26 and 0.194. The result indicates that though deepwalk does not have a clear loss function, it can make better use of network structure information through random walks. It also shows that the best suitable method should be selected through quantitative evaluation in a new task or domain. For example, the algorithms like node2vec, LINE, and SDNE are all appeared after deepwalk, and researchers have stated that these methods perform better than deepwalk in their experiment. However, deepwalk performs best in this special task of citation recommendation which is very different from the previous conclusions. In summary, deepwalk is considered as the optimal representation learning method for semantic representation of cited papers' relations in this situation.

Table 1
Description of training and testing datasets.

	Training set (2007–2016)	Testing set (2017–2018)
Number of cited papers	42,352	20,402
Number of co-citation relationships	1,222,897	542,554
Number of current sentences	76,661	33,176
Number of previous sentences	76,661	33,176
Number of next sentences	76,646	33,159

Table 2

AUC and MAP of four network representation learning methods for citation recommendation.

Method	Deepwalk	Node2vec	LINE	SDNE
AUC MAP	0.685 0.725	0.569 0.625	0.425 0.531	0.536 0.542

The algorithm node2vec has designed a flexible neighborhood sampling strategy that improves the process of random walks by using breadth-first sampling and depth-first sampling, but the performance is much worse than deepwalk. Compared to deepwalk, the AUC and MAP of node2vec have decreased by 0.116 and 0.1. The possible reason is that the important parameters in node2vec, i.e., p and q , may significantly affect the experimental results. Parameter p controls the likelihood of immediately revisiting a node in the walk and parameter q allows the search to differentiate between “inward” and “outward” nodes. This paper sets “ $p = 0.25$ ” and “ $q = 0.25$ ” because such parameter settings have achieved good results in many tasks (De Winter et al., 2018). However, these parameter setting may not be suitable for this citation recommendation task.

Unlike the deepwalk and node2vec, LINE and SDNE consider both first-order proximity and second-order proximity in the process of learning semantic representation of cited papers’ relations, but the performance is not good. Compare with deepwalk which perform best, the AUC and MAP of LINE have decreased by 0.26 and 0.194 and SNDE has decreased by 0.149 and 0.183. It shows that for citation recommendation task, excessively considering the information of common neighbors will weaken the effect.

4.3. Citation recommendation using semantic representation of different forms of citation content

In this experiment, four text learning algorithms are combined with four forms of citation content. Doc2vec is implemented by using doc2vec module in gensim package. Some important parameters in doc2vec are set as follows: “window = 10”, “vector_size = 128” and “dm = 0” (DBOW model is chosen). The implementation of doc2vecC uses the source code provided by Chen (2017), where the embedding dimension is also set to 128. The rest parameters in doc2vec and doc2vecC are assigned with the default values (Chen, 2017; Le & Mikolov, 2014). LSI and LDA are implemented by using lsimodel and ldamodel module in gensim package. The “num_topics = 128” is set for each topic model to ensure that the embedding dimension is the same as other algorithms. The rest parameters in LDA and LSI are assigned with the default values (Gupta & Varma, 2017).

Table 3 shows the AUC and MAP of four text representation algorithms combined with four forms of citation content for citation recommendation, in which the numbers emphasized in bold correspond to the highest AUC and MAP in each row, and the numbers emphasized in italic correspond to the highest value in each column.

It can be found that doc2vecC performs much better than other three algorithms in all forms of citation content, whose AUC and MAP are emphasized italic. Compared with doc2vec, the accuracy within all forms of citation content has been significantly improved, i.e., the AUC and MAP of CS have increased by 0.264 and 0.29, the AUC and MAP of

Table 3

AUC and MAP of four text representation algorithms combined with four forms of citation content.

Method	Form			
	CS	CS&SS	CS&SS-	SuCS&SS
LSI	0.745 0.735	0.798 0.798	0.788 0.787	0.722 0.718
LDA	0.740 0.745	0.777 0.788	0.758 0.769	0.712 0.713
doc2vec	0.566 0.561	0.712 0.713	0.647 0.642	0.692 0.692
doc2vecC	<i>0.830 0.851</i>	<i>0.877 0.889</i>	<i>0.867 0.882</i>	<i>0.826 0.835</i>

CS&SS have increased by 0.165 and 0.176, the AUC and MAP of CS&SS- have increased by 0.22 and 0.24, and the AUC and MAP of SuCS&SS have increased by 0.134 and 0.143. Similarly, doc2vecC performs much better than LSI and LDA in all forms of citation content. The reason may be that doc2vecC generates document representation through average word embedding, which can better capture the semantics of the document during the learning process. On the other hand, each form of citation content only consists of an average of 3 sentences which is often treated as short text, and doc2vecC may be more suitable for short text than other three methods. This result could provide useful suggestions in selecting text representation learning methods for short text.

From the perspective of four forms of citation content, no matter which algorithm is used, CS&SS performs better than all other forms of citation content, whose AUC and MAP are emphasized in bold. For LSI, the AUC and MAP of CS&SS are both 0.798 which has increased by 0.053 and 0.063 compared with CS, 0.01 and 0.011 compared with CS&SS- and 0.076 and 0.08 compared with SuCS&SS, respectively. Similarly, CS&SS is also the best form of citation content for LDA, doc2vec and doc2vecC. The possible reason is that although the previous sentence and the next sentence may cite other references, these contents still have strong semantic relevance to the citation, which is conducive to citation recommendation. This result can be borrowed on other similar studies related to citation content analysis and may draw more accurate and sufficient conclusions, because the combination of current, previous and next sentences performs best in almost all situations.

In summary, doc2vecC is the best text representation method and CS&SS is the best form of citation content as cited papers’ content. In addition, the combination of doc2vecC and CS&SS perform best whose AUC and MAP reach 0.877 and 0.889.

4.4. Citation recommendation using semantic fusion of cited papers’ relations and content

Naïve combination is the first fusion method to combine the vectors generated by network and text representation learning. The dimension of the vectors generated by the naïve combination is 256, because the dimension of the vectors generated by both single methods is 128. Table 4 shows the AUC and MAP of naïve combination for citation recommendation, in which the numbers emphasized in bold correspond to the highest AUC and MAP in each method.

The result shows that the combination of deepwalk and doc2vec(CS) performs best, where the AUC and MAP reach 0.830 and 0.846. Compared with the combination of LINE and LSI(CS) which performs worst, the AUC and MAP have increased by 0.415 and 0.327. The reason may be that deepwalk performs better than other three network representation methods, thus the combination of deepwalk with text representation methods may perform better than other combinations.

When node2vec is used as network representation method, the combination of node2vec and LSI (CS&SS) performs best, where the AUC and MAP reach 0.702 and 0.746 and have decreased by 0.128 and 0.1 compared with the best one. When LINE is selected, the combination of LINE and LSI (CS&SS) performs best, where the AUC and MAP reach 0.619 and 0.688 and have decreased by 0.211 and 0.158 compared with the best one. The possible reason is that the performance of node2vec and LINE is not as good as deepwalk, resulting in poor effect with combination.

It should be noted that the AUC and MAP are always the same in four forms of citation content when SDNE is selected. The possible reason is that each digit of the vector generated by SDNE is relatively large, some of them even reach 100, whereas each digit of the vector generated by other text representation methods is between -1 and 1 . Therefore, when they are directly joined, the text feature hardly affects the result after stitching.

Another finding is that the best fusion method may not be the combination of the best single methods of cited papers’ relations and content respectively. In this situation, the best single methods are

Table 4
AUC and MAP of naïve combination for citation recommendation.

Method		Form			
		CS	CS&SS	CS&SS-	SuCS&SS
Deepwalk	LSI	0.708	0.816	0.793	0.743
		0.751	0.838	0.822	0.770
	LDA	0.686	0.686	0.686	0.686
		0.726	0.725	0.725	0.725
	Doc2vec	0.830 	0.825	0.826	0.822
		0.846	0.840	0.840	0.836
	Doc2vecC	0.750	0.732	0.734	0.738
		0.776	0.761	0.763	0.765
	Node2vec	0.504	0.702 	0.643	0.608
		0.578	0.746	0.697	0.657
LINE	LSI	0.570	0.569	0.569	0.569
		0.625	0.625	0.622	0.624
	LDA	0.668	0.684	0.669	0.670
		0.698	0.710	0.698	0.700
	Doc2vec	0.669	0.595	0.598	0.594
		0.698	0.646	0.647	0.644
	Doc2vecC	0.415	0.619 	0.548	0.510
		0.519	0.688	0.629	0.591
	SDNE	0.426	0.425	0.425	0.425
		0.532	0.531	0.530	0.531
SDNE	LSI	0.479	0.501	0.489	0.480
		0.568	0.588	0.577	0.573
	LDA	0.441	0.436	0.437	0.435
		0.545	0.541	0.542	0.540
	Doc2vec	0.536	0.536	0.536	0.536
		0.542	0.542	0.542	0.542
	Doc2vecC	0.536	0.536	0.536	0.536
		0.542	0.542	0.542	0.542
	Doc2vecC	0.536	0.536	0.536	0.536
		0.542	0.542	0.542	0.542

deepwalk and doc2vecC(CS&SS), and their corresponding fusion method's AUC and MAP are 0.732 and 0.761. However, the truly best fusion method is deepwalk combined with doc2vec(CS), where doc2vec(CS) is not the best single method. This best fusion method's AUC and MAP reach to 0.830 and 0.846 which have increased by 0.098 and 0.085 than the fusion of two best single methods. Therefore, it is necessary to try different combinations of network and text representation methods with multiple forms of citation content, to find the best fusion method for specific task.

Like the above result, the fusion method is not always performing better than the single method in citation recommendation. For example, the single method deepwalk that learning cited papers' relations performs best with AUC 0.685 and MAP 0.725, and doc2vec(CS&SS) that learning cited papers' content performs best with AUC 0.877 and MAP 0.889. However, after the fusion the AUC and MAP are only 0.732 and 0.761, respectively. There are some other similar examples, which means that both semantic vectors play important role for citation recommendation and the fusion method should be improved for this specific task. In summary, the best representation learning method for citation recommendation should be determined by quantitative evaluation and more fusion methods should be applied and improved.

The other three fusion algorithms are TADW, Paper2vec and GCN. The dimension or hidden layer is set to the same value of 128, and the rest parameters are assigned with the default values that perform best in their experiment. These three algorithms have the unchanged build-in network embedding method, so we can only change text representation learning methods for fusion and comparison. Specifically, in TADW, the combination of deepwalk with LSI, LDA, doc2vec and doc2vecC are named as TADW1, TADW2, TADW3, and TADW4, respectively. Similarly, the Paper2vec and GCN has the same four combinations. Table 5 shows the AUC and MAP of TADW, Paper2vec and GCN for citation recommendation, in which the numbers emphasized in bold correspond to the highest AUC and MAP for each of them.

Table 5
AUC and MAP of TADW, Paper2vec and GCN for citation recommendation.

Method		Form			
		CS	CS&SS	CS&SS-	SuCS&SS
TADW	TADW1	0.718 0.753	0.752 0.777	0.729 0.770	0.748 0.767
	TADW2	0.802 0.815	0.770 0.784	0.778 0.794	0.774 0.783
	TADW3	0.784 0.794	0.791 0.811	0.774 0.789	0.770 0.784
	TADW4	0.806 0.827	0.800 0.821	0.803 0.824	0.790 0.810
Paper2vec	Paper2vec1	0.720 0.737	0.731 0.749	0.729 0.748	0.728 0.741
	Paper2vec2	0.701 0.722	0.728 0.742	0.725 0.735	0.707 0.723
	Paper2vec3	0.732 0.737	0.737 0.749	0.740 0.748	0.736 0.743
	Paper2vec4	0.743 0.760	0.738 0.763	0.747 0.767	0.729 0.750
GCN	GCN1	0.501 0.502	0.511 0.521	0.537 0.545	0.542 0.548
	GCN2	0.839 0.866	0.853 0.876	0.843 0.868	0.850 0.873
	GCN3	0.625 0.666	0.560 0.591	0.545 0.594	0.577 0.605
	GCN4	0.529 0.532	0.557 0.600	0.565 0.565	0.546 0.580

The result shows that GCN performs better than other three fusion methods, with the AUC of 0.853 and the MAP of 0.876. Compared with the optimal result of naïve combination, the AUC and MAP have increased by 0.023 and 0.03. Compared with the optimal results of TADW, the AUC and MAP have increased by 0.047 and 0.049. Compared with the optimal results of Paper2vec, the AUC and MAP have increased by 0.106 and 0.109. The possible reason is that GCN learns semantic representation of objects by convolving multiple-level features automatically, which makes better use of cited papers' relations and content compared with other fusion methods. In summary, GCN is considered as the optimal fusion method in this situation.

Although CS&SS is not always performing best in all representation learning algorithms, but it is still stable and has a very little difference compared with other forms. In GCN, the most appropriate form of citation content is still CS&SS and GCN2(CS&SS) performs best in all fusion methods. In TADW, the best form is CS, whose best AUC and MAP are 0.806 and 0.827. They have only increased by 0.006 compared with CS&SS. In Paper2vec, the best form is CS&SS-, whose best AUC and MAP are 0.747 and 0.767. They have increased by 0.009 and 0.004 compared with CS&SS. These results prove that CS&SS performs best in almost all situations considering all single and fusion methods.

In summary, doc2vecC(CS&SS) is the best method for citation recommendation which performs better than the best network embedding method and the best fusion method. Its AUC and MAP reach 0.877 and 0.889. Compare with the best fusion method, i.e., GCN, its AUC and MAP have increased by 0.024 and 0.013. Compare with the best method in network representation learning, i.e., deepwalk, its AUC and MAP have increased by 0.192 and 0.164. Therefore, doc2vecC(CS&SS) could be treated as the optimal method for citation recommendation.

4.5. Parameter adjustment of optimal method

The performance of citation recommendation can be improved by adjusting the important parameters in specific domain (Liu et al., 2020; Zhang et al., 2016). The best method in this paper for citation recommendation is the doc2vecC(CS&SS), so we select doc2vecC for parameter adjustment. According to investigation, the number of iterations and the size of embedding dimension have a great influence on doc2vecC (Chen, 2017), so this paper adjusts these two parameters. In this experiment, the parameter "epochs" is set to 1, 5, 10, 15, 20, 25 and 30, and "embed_dims" is set to 80, 100, 120, 140, and 160, respectively. Fig. 1 shows the AUC and MAP with varying "epochs" and "embed_dims" in doc2vecC(CS&SS).

The result shows that the trends of AUC and MAP are similar with the adjustment of the parameters and the value of MAP is slightly higher than that of AUC overall, as shown in Fig. 1. Under the same dimension, when "epochs" varies from 1 to 5, both AUC and MAP have increased rapidly. It means that the smaller "epochs" cannot make the model learn all the information. On the contrary, when "epochs" varies from 5 to 30,

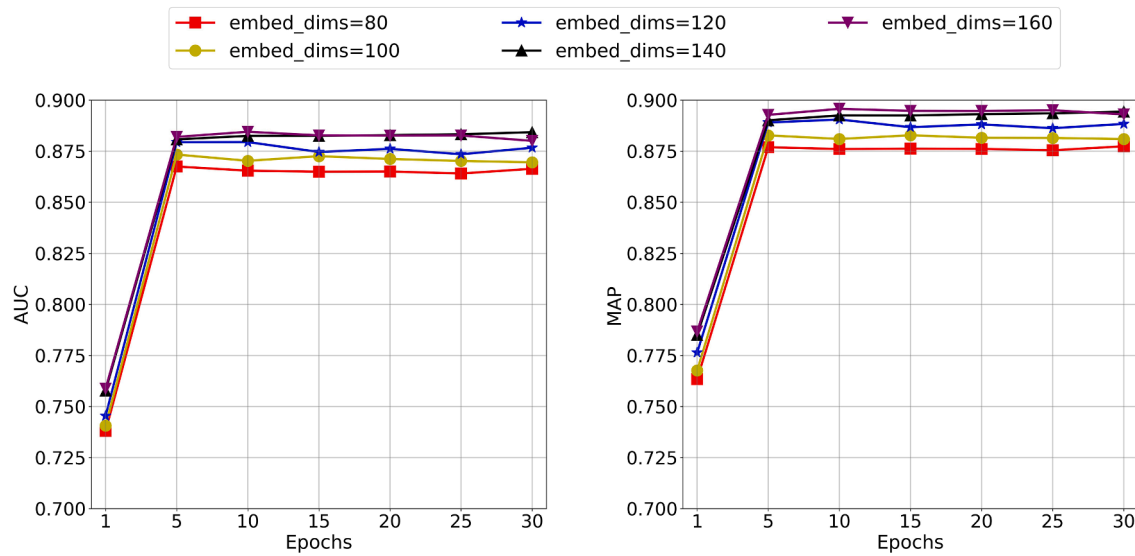


Fig. 1. AUC and MAP of doc2vecC(CS&SS) with varying “epochs” and “embed_dims”.

the AUC and MAP are relatively stable which have not increased according to epochs. The possible reason is that the model has learned all the semantic information of cited papers’ content when “epochs” is set as 5, and when the epoch exceeds 5, there is little semantics added.

Another finding is that both AUC and MAP are always increasing as the dimension increases except few cases. The possible reason is that the high embedding dimension can encode more useful text information into the vector and can achieve more accurate semantic representation of cited papers’ content. When the two important parameters are set as “epochs = 10” and “embed_dims = 160”, the method performs best with the AUC and MAP reach 0.885 and 0.896. Compared with the result that the parameters are assigned with the default value, the AUC and MAP have slightly improved by 0.008 and 0.007. It indicates that specific parameter adjustments can improve the performance of citation recommendation, but in this specific task, the improvement is small compared with other methods, i.e., the accuracy of TADW has increased by 0.02 for multi-class classification with parameter adjustment, and the F1 of Grarep has increased by about 0.06 for clustering with parameter adjustment (Cao et al., 2015; Yang et al., 2015).

4.6. Case study

We choose the optimal method doc2vecC(CS&SS) after parameters adjustment to make citation recommendation for case study. Two representative cited papers are selected and top k ($k = 5$) citation recommendations are presented for result analysis. In addition, two representative sentences are selected from current sentences and surrounding sentences of each recommendation to analyze the relevance of recommendation and the selected cited paper.

The first selected cited paper is a classic algorithm entitled “Random forests”, and the top 5 citation recommendations are shown in Table 6. It shows that the recommendation results are highly correlated with the selected cited paper, which can be used to help researchers supplement or update references. For example, recommendation “Rotation forest: a new classifier ensemble method” is the most relevant to the selected cited paper, and the similarity reaches 0.877. Researchers can cite this recommendation to explain why random forests performs better as an ensemble method. In addition, the fourth ranking recommendation is “Support-vector networks” and the similarity reaches 0.815. Researchers can cite this recommendation to show that besides random forests, SVM is also a relatively advanced classifier, and some researches has integrated them into a framework.

The second selected cited paper is an application paper entitled

“DNA binding protein identification by combining pseudo amino acid composition and profile-based protein representation”, and the top 5 citation recommendations are shown in Table 7. It shows that the recommendation results are all related to DNA and protein, and the 1st, 3rd and 5th ranking recommendations are indeed co-cited in the future. This indicates that the proposed method is effective. Specifically, the first ranking recommendation can be cited to illustrate that classification algorithms are often considered as conventional methods for predicting DNA binding protein based on sequence. In addition, the fourth ranking recommendation can be cited to illustrate the general characteristics of sequence-based methods.

5. Conclusions

Existing citation recommendation works mostly focus on citing papers, resulting in the enormous cited papers are ignored. Therefore, this paper designs multiple forms of citation content to denote cited paper’s content, and compares several semantic representation methods to find the most appropriate form of citation content and the optimal method for citation recommendation. The result of quantitative comparisons shows that deepwalk, doc2vecC(CS&SS) and GCN(CS&SS) is the optimal method respectively, corresponding to network representation, text representation and the fusion method. Besides, a case study is performed for qualitative analysis whose results have further proved the effectiveness of the method. These jointly prove the effectiveness of this framework and could provide useful suggestions for citation recommendation.

Doc2vecC(CS&SS) is the optimal method with the highest AUC and MAP of 0.885 and 0.896, which have increased by 0.462 and 0.37 compared with the worst-performing method. Compared with other similar tasks, the values of these two evaluation indicators perform equal or better. For example, in the news recommendation task, the best AUC in the 6 datasets just reaches 0.8116, and in the link prediction task of heterogeneous social networks, the best AUC reaches about 0.89 (Hu et al., 2020; Ozcan & Oguducu, 2019). Moreover, the best MAP of entity recommendation task in social networks of five datasets reaches about 0.85, and in cultural event recognition task, the MAP reaches only 0.767 (Salvador et al., 2015; Yu et al., 2017). Although these datasets are for the specific tasks and come from different domains, they can indirectly prove our method is effective and the conclusions are credible. Simultaneously, embedding vector of cited papers generated in this way can be used as auxiliary information in other tasks. On the whole, this method is effective to find supplementary or alternative references, and

Table 6

Top 5 recommendations of the first selected cited paper.

Rank	Title	Similarity	Representative sentences
1	Rotation forest: a new classifier ensemble method	0.877	Random Forest or Rotation Forest significantly outperform classical decision tree classification models in terms of classification accuracy or other classification performance metrics, they are not suitable for knowledge discovery process. When forming ensembles, creating diverse classifiers (but maintaining their consistency with the training set) is a key factor to make them accurate.
2	The elements of statistical learning	0.837	Thus, the tree analysis is run on multiple similar datasets, and the aggregate results are used to construct the final model and relevant statistics. Friedman recommends bagging with 50% of the database.
3	Estimating continuous distributions in Bayesian classifiers	0.817	Furthermore, different prediction models were established using the (III) naïve Bayes algorithm, calculating prior-, conditional- and posterior-probabilities, (IV) the logistic regression classifier, characterized by membership function for each class, and (V) the multi-layer perceptron (an artificial neural network), combining various linear models for non-linear classification. We applied different such methods for classifying the feature based data into classes (TB, not TB), as desired.
4	Support-vector networks	0.815	Support vector machines are supervised machine learning algorithms which build models based on 'training' data and search for similar patterns in 'test' data. In this section, we summarize four state of the art classifiers in the literature, namely SVM, Random Forest (RF), Naive Bayes, and K-nearest Neighbor (KNN) classifier.
5	Which classification algorithm works best with stylistic features of Portuguese in order to classify web texts according to users' needs?	0.812	We utilized the RCRT approach that constructed an ensemble of classifiers with Random Tree as the base classifier. The random committee algorithm generated predictions by averaging probability estimates over the generated classification trees. The final prediction was a straight average of the predictions generated by the individual base classifiers.

Table 7

Top 5 recommendations of the second selected cited paper.

Rank	Title	Similarity	Representative sentences
1	Efficient prediction of nucleic acid binding function from low-resolution protein structures	0.781	As a summary, sequence based prediction methods for DNA-binding proteins have been investigated with several classifiers such as logistic regression, random forest, support vector machine. For instance, Szilágyi and Skolnick used logistic regression to predict the DNA-binding proteins from the amino acid composition. Recently, the concept of PseAAC was further extended to represent the feature vectors of DNA and nucleotides. SVM has been widely used in the realm of bioinformatics. Later, the same authors introduced the grouped weight to represent protein samples for predicting DNA-binding proteins. These methods can be divided into two categories, structure based modeling and sequence based prediction.
2	iRSpot-PseDNC: identify recombination spots with pseudo dinucleotide composition nucleic acids research	0.780	Sequence-based methods use amino-acid compositions, sequence homology, and sorting signals as features. In 2009, Gao et al reconstructed the benchmark dataset for NRs and introduced the pseudo amino acid composition (PseAAC) to represent the protein samples in hope to improve the prediction quality.
3	An ensemble of reduced alphabets with protein encoding based on grouped weight for predicting DNA-binding proteins	0.777	Gao and Skolnick proposed a threading-based method which required only the target protein sequence to identify the DNA-binding domains based on a template library composed of DNA-protein complex structure. The 1st type is actually using both the structural of proteins and their sequences information for identifying the DNA-binding proteins.
4	Prediction of protein cellular attributes using pseudo amino acid composition	0.772	
5	A threading-based method for the prediction of DNA-binding proteins with application to the human genome	0.766	

provides helpful and useful suggestions for representation learning method and citation content selection in citation recommendation.

Through in-depth analysis of the experimental results, it can be found that it is necessary to design multiple forms of citation content for comparison, which is conducive to improving the performance of citation recommendation. Among four forms of citation content, CS&SS performs best in almost all methods, which shows that the surrounding sentences have a strong semantic correlation with the current sentences and could be a good supplement to describe motivation why the reference is cited. The multiple forms of citation content designed in this paper can be easily extended to other text mining related tasks, such as topic analysis and citation sentiment analysis, etc., which could enrich the experimental data and improve the experimental performance. Certainly, the optimal form of citation content in other tasks may be different from that in this paper and the optimal form in other tasks need

to be determined by quantitative evaluation.

In addition, the results show that among all the semantic representation learning methods, doc2vecC performs best within any form of citation content. The reason may be that doc2vecC generates document representation through average word embedding, which can better capture the semantics of the document during the learning process. Simultaneously, this single method of text representation method performs better than fusion method. We attribute this to two reasons. On the one hand, this paper does not assign weight to co-citation relationships according to their location, e.g., the relation that two papers co-cited in different section is less important than in the same sentence or paragraph. On the other hand, parameters are assigned with default value may influence the result.

Furthermore, the results show that the best fusion method is not the combination of the best single methods of network and text respectively. It is also worth mentioning that the fusion method sometimes performs worse than the single method in citation recommendation. These two findings may not consistent with normal cognition before our experiment. They indicate that each method has its own scope of application and the conclusion may be changed when an algorithm is applied to a different task. Therefore, it is necessary to apply as many representation learning methods as possible for comparisons to get the best semantic vector of the cited paper. This conclusion is also applicable to other task when using representation learning methods, such as clustering, recommendation and classification.

It should be noted that the effect of citation recommendation has not been greatly improved through the adjustment of specific parameters, where the AUC and MAP have only increased by 0.008 and 0.007. However, in some fields, parameter adjustments can significantly improve the performance, such as multi-label classification and phrase extraction in specific tasks. Therefore, subsequent experiments in other fields are necessary to be conducted for exploring the influence of these parameters.

Finally, this study is a preliminary research and still need to be further studied and improved in following directions. Firstly, although this paper has selected several typical and representative algorithms to perform semantic representation of cited papers' relations and content, other latest methods is still emerging and may have better semantic representation effect on them. Therefore, it may necessary to apply other methods for comparative studies to find the optimal method for citation recommendation in the future. Secondly, this paper only conducts experiment in the field of artificial intelligence, but the results may be different in other fields because authors may have different citation behaviors, styles and motivations in different fields. Therefore, it is necessary to expand research field to others, e.g., social science, data management and so on, to further prove the feasibility of citation recommendation from the perspective of cited papers.

Funding

This work was supported by the National Natural Science Foundation of China (Grant No. 71974095) and the Social Science Fund Youth Project of Jiangsu Province (No. 17TQC003).

CRediT authorship contribution statement

Jinzhu Zhang: Conceptualization, Methodology, Formal analysis, Writing – original draft, Writing – review & editing, Funding acquisition.
Lipeng Zhu: Software, Validation, Data curation, Writing – original draft, Writing – review & editing.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- Abu-Jbara, A., & Radev, D. (2012). Reference scope identification in citing sentences. Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Montreal, Canada.
- Alhijawi, B., & Kilani, Y. (2020). A collaborative filtering recommender system using genetic algorithm. *Information Processing & Management*, 57(6), Article 102310.
- Ali, Z., Qi, G., Kefalas, P., Abro, W. A., & Ali, B. (2020). A graph-based taxonomy of citation recommendation models. *Artificial Intelligence Review*, 53(2), 1–44.
- Ali, Z., Qi, G., Muhammad, K., Ali, B., & Abro, W. A. (2020). Paper recommendation based on heterogeneous network embedding. *Knowledge-Based Systems*, 210, Article 106438.
- Angrosh, M., Cranefield, S., & Stanger, N. (2010). Context identification of sentences in related work sections using a conditional random field: towards intelligent digital libraries. Proceedings of the 10th Annual Joint Conference on Digital libraries, New York, United States.
- Angrosh, M., Cranefield, S., & Stanger, N. (2013). Conditional random field based sentence context identification: enhancing citation services for the research community. Proceedings of the First Australasian Web Conference, Sydney, Australia.
- Balcerzak, B., Jaworski, W., & Wierzbicki, A. (2014). Application of TextRank algorithm for credibility assessment. 2014 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT), Warsaw, Poland.
- Bhagavatula, C., Feldman, S., Power, R., & Ammar, W. (2018). Content-based citation recommendation. Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics, New Orleans, Louisiana.
- Bilgin, M., & Şentürk, İ. F. (2017). Sentiment analysis on Twitter data with semi-supervised Doc2Vec. 2017 International Conference on Computer Science and Engineering, Antalya, Turkey.
- Blei, D. M., Ng, A. Y., Jordan, M. I., & Lafferty, J. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- Cai, X., Han, J., Li, W., Zhang, R., Pan, S., & Yang, L. (2018). A three-layered mutually reinforced model for personalized citation recommendation. *IEEE transactions on neural networks and learning systems*, 29(12), 6026–6037.
- Cao, S., Lu, W., & Xu, Q. (2015). GraRep: learning graph representations with global structural information. Proceedings of the 24th ACM International on Conference on Information and Knowledge Management, Melbourne, Australia.
- Carusi, C., & Bianchi, G. (2019). Scientific community detection via bipartite scholar/journal graph co-clustering. *Journal of Informetrics*, 13(1), 354–386.
- Chen, M. (2017). Efficient vector representation for documents through corruption. Proceedings of the 5th International Conference on Learning Representations, Toulon, France.
- Chen, X., Zhao, H.-J., Zhao, S., Chen, J., & Zhang, Y.-P. (2019). Citation recommendation based on citation tendency. *Scientometrics*, 121(2), 937–956.
- Dai, T., Zhu, L., Wang, Y., & Carley, K. M. (2020). Attentive stacked denoising autoencoder with Bi-LSTM for personalized context-aware citation recommendation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28, 553–568.
- De Winter, S., Decuyper, T., Mitrović, S., Baesens, B., & De Weerd, J. (2018). Combining temporal aspects of dynamic networks with Node2Vec for a more efficient dynamic link prediction. 2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), Barcelona, Spain.
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6), 391–407.
- Esen, H., Inalli, M., Sengur, A., & Esen, M. (2008). Performance prediction of a ground-coupled heat pump system using artificial neural networks. *Expert Systems with Applications*, 35(4), 1940–1948.
- Esen, H., Ozgen, F., Esen, M., & Sengur, A. (2009). Artificial neural network and wavelet neural network approaches for modelling of a solar air heater. *Expert Systems with Applications*, 36(8), 11240–11248.
- Eto, M. (2019). Extended co-citation search: Graph-based document retrieval on a co-citation network containing citation context information. *Information Processing & Management*, 56(6), Article 102046.
- Färber, M., & Jatowt, A. (2020). Citation recommendation: Approaches and datasets. *International Journal on Digital Libraries*, 21(4), 375–405.
- Lakshmanan, G. P., & Ramanathan, L. (2019). Using citation context to improve the retrieval of research article from cancer research journals. *Asian Pacific Journal of Cancer Prevention*, 20(3), 951–960.
- Ganguly, S., & Pudi, V. (2017). Paper2vec: Combining graph and text information for scientific paper representation. European Conference on Information Retrieval, Aberdeen, Scotland.
- Ghorbanzadeh, H., Sheikhhahmadi, A., Jalili, M., & Sulaimany, S. (2021). A hybrid method of link prediction in directed graphs. *Expert Systems with Applications*, 165, Article 113896.
- Gipp, B., Beel, J., & Hentschel, C. (2009). Scienstein: A research paper recommender system. Proceedings of the International Conference on Emerging Trends in Computing, Tamil Nadu, India.
- Greiner-Petter, A., Youssef, A., Ruas, T., Miller, B. R., & Gipp, B. (2020). Math-word embedding in math search and semantic extraction. *Scientometrics*, 125(3), 3017–3046.
- Grover, A., & Leskovec, J. (2016). node2vec: Scalable feature learning for networks. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, United States.
- Guns, R., & Rousseau, R. (2014). Recommending research collaborations using link prediction and random forest classifiers. *Scientometrics*, 101(2), 1461–1473.

- Gupta, S., & Varma, V. (2017). Scientific article recommendation by using distributed representations of text and graph. Proceedings of the 26th International Conference on World Wide Web Companion, Perth, Australia.
- He, Q., Pei, J., Kifer, D., Mitra, P., & Giles, L. (2010). Context-aware citation recommendation. Proceedings of the 19th International Conference on World Wide Web, Beijing, China.
- Hu, L., Li, C., Shi, C., Yang, C., & Shao, C. (2020). Graph neural news recommendation with long-term and short-term interest modeling. *Information Processing & Management*, 57(2), Article 102142.
- Jeong, C., Jang, S., Park, E., & Choi, S. (2020). A context-aware citation recommendation model with BERT and graph convolutional networks. *Scientometrics*, 124(3), 1907–1922.
- Jha, R., Jbara, A. A., Qazvinian, V., & Radev, D. R. (2017). NLP-driven citation analysis for scientometrics. *Natural Language Engineering*, 23(1), 93–130.
- Jin, W., & Srihari, R. K. (2007). Graph-based text representation and knowledge discovery. *Proceedings of the 2007 ACM Symposium on Applied Computing*.
- Kim, S. M., Paris, C., Power, R., & Wan, S. (2017). Distinguishing individuals from organisations on Twitter. Proceedings of the 26th International Conference on World Wide Web Companion, Perth, Australia.
- Kipf, T. N., & Welling, M. (2017). Semi-supervised classification with graph convolutional networks. Proceedings of the 5th International Conference on Learning Representations (ICLR), Toulon, France.
- Kobayashi, Y., Shimbo, M., & Matsumoto, Y. (2018). Citation recommendation using distributed representation of discourse facets in scientific articles. Proceedings of the 18th ACM/IEEE on Joint Conference on Digital Libraries, Fort Worth Texas, USA.
- Kong, X., Mao, M., Wang, W., Liu, J., & Xu, B. (2018). VOPRec: Vector representation learning of papers with text information and structural identity for recommendation. *IEEE Transactions on Emerging Topics in Computing*, 1–1.
- Le, Q., & Mikolov, T. (2014). Distributed representations of sentences and documents. Proceedings of the 31st International Conference on Machine Learning, Beijing, China.
- Li, C., Wang, S., Yang, D., Li, Z., Yang, Y., Zhang, X., & Zhou, J. (2017). PPNE: property preserving network embedding. International Conference on Database Systems for Advanced Applications, Suzhou, China.
- Liu, H., Chen, Z., Tang, J., Zhou, Y., & Liu, S. (2020). Mapping the technology evolution path: A novel model for dynamic topic detection and tracking. *Scientometrics*, 125(2), 2043–2090.
- Liu, H., Kong, X., Bai, X., Wang, W., Bekele, T. M., & Xia, F. (2015). Context-based collaborative filtering for citation recommendation. *IEEE Access*, 3, 1695–1703.
- Liu, S., Chen, C., Ding, K., Wang, B., Xu, K., & Lin, Y. (2014). Literature retrieval based on citation context. *Scientometrics*, 101(2), 1293–1307.
- Liu, X., Yu, Y., Guo, C., & Sun, Y. (2014). Meta-path-based ranking with pseudo relevance feedback on heterogeneous graph for citation recommendation. Proceedings of the 23rd ACM International Conference on Information and Knowledge Management, New York, United States.
- Ma, S., Zhang, C., & Liu, X. (2020). A review of citation recommendation: From textual content to enriched context. *Scientometrics*, 122(4), 1445–1472.
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. Proceedings of the 1st International Conference on Learning Representations, Scottsdale, Arizona, USA.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *Advances in Neural Information Processing Systems*, Lake Tahoe, Nevada, USA.
- Nogueira, R., Jiang, Z., Cho, K., & Lin, J. (2020). Navigation-based candidate expansion and pretrained language models for citation recommendation. *Scientometrics*, 125(3), 3001–3016.
- Ozcan, A., & Oğuducu, S. G. (2019). Multivariate time series link prediction for evolving heterogeneous network. *International Journal of Information Technology & Decision Making*, 18(1), 241–286.
- Pan, S., Wu, J., Zhu, X., Zhang, C., & Wang, Y. (2016). Tri-party deep network representation. Proceedings of the 25th International Joint Conference on Artificial Intelligence, New York, United States.
- Park, H., Kwon, K., Khiati, A. I. Z., Lee, J., & Chung, I. J. (2015). Agglomerative hierarchical clustering for information retrieval using latent semantic index. 2015 IEEE International Conference on Smart City/SocialCom/SustainCom (SmartCity), Chengdu, China.
- Perozzi, B., Al-Rfou, R., & Skiena, S. (2014). Deepwalk: Online learning of social representations. Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, New York, United States.
- Poria, S., Chaturvedi, I., Cambria, E., & Bisio, F. (2016). Sentic LDA: Improving on LDA with semantic similarity for aspect-based sentiment analysis. 2016 International Joint Conference on Neural Networks (IJCNN), Ancouver, Canada.
- Salvador, A., Zeppelzauer, M., Manchon-Vizuete, D., Calafell, A., & Giro-I-Nieto, X. (2015). Cultural event recognition with visual ConvNets and temporal models. 2015 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), Boston, United States.
- Son, J., & Kim, S. B. (2018). Academic paper recommender system using multilevel simultaneous citation networks. *Decision Support Systems*, 105, 24–33.
- Strohman, T., Croft, W. B., & Jensen, D. (2007). *Recommending citations for academic papers*. Amsterdam, Netherlands: SIGIR.
- Sugiyama, K., & Kan, M.-Y. (2010). Scholarly paper recommendation via user's recent research interests. Proceedings of the 10th Annual Joint Conference on Digital Libraries, New York, United States.
- Tang, D., Qin, B., & Liu, T. (2015). Document modeling with gated recurrent neural network for sentiment classification. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Lisbon, Portugal.
- Tang, J., Qu, M., & Mei, Q. (2015). PTE: Predictive text embedding through large-scale heterogeneous text networks. Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, New York, United States.
- Tang, J., Qu, M., Wang, M., Zhang, M., Yan, J., & Mei, Q. (2015). LINE: Large-scale information network embedding. *Proceedings of the 24th International Conference on World Wide Web*.
- Trieu, L. Q., Tran, H. Q., & Tran, M.-T. (2017). *News classification from social media using twitter-based doc2vec model and automatic query expansion*. Nha Trang, Vietnam: Information and Communication Technology.
- Tu, C., Liu, H., Liu, Z., & Sun, M. (2017). Cane: Context-aware network embedding for relation modeling. Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Vancouver, Canada.
- Wang, C., & Blei, D. M. (2011). Collaborative topic modeling for recommending scientific articles. Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, New York, United States.
- Wang, D., Cui, P., & Zhu, W. (2016). Structural deep network embedding. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, United States.
- Wang, X., Cui, P., Wang, J., Pei, J., Zhu, W., & Yang, S. (2017). Community preserving network embedding. Proceedings of 31st the AAAI Conference on Artificial Intelligence, California, United States.
- West, J. D., Wesley-Smith, I., & Bergstrom, C. T. (2016). A recommendation system based on hierarchical clustering of an article-level citation network. *IEEE Transactions on Big Data*, 2(2), 113–123.
- Yang, C., Liu, Z., Zhao, D., Sun, M., & Chang, E. (2015). Network representation learning with rich text information. Twenty-Fourth International Joint Conference on Artificial Intelligence, Buenos Aires, Argentina.
- Yang, L., Zheng, Y., Cai, X., Dai, H., Mu, D., Guo, L., & Dai, T. (2018). A LSTM based model for personalized context-aware citation recommendation. *IEEE Access*, 6, 59618–59627.
- Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., & Hovy, E. (2016). Hierarchical attention networks for document classification. Proceedings of the 15th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego, California.
- Ye, Y., Zhang, S., Li, Y., Qian, X., Tang, S., Pu, S., & Xiao, J. (2020). Video question answering via grounded cross-attention network learning. *Information Processing & Management*, 57(4), Article 102265.
- Yin, Y., Ji, L., Huang, R., & Du, L. (2019). Research and development of network representation learning. *Chinese Journal of Network and Information Security*, 5(2), 81–91.
- Yu, C., Zhao, X., An, L., & Lin, X. (2017). Similarity-based link prediction in social networks: A path and node combined approach. *Journal of Information Science*, 43(5), 683–695.
- Zarrinkalam, F., & Kahani, M. (2012). A multi-criteria hybrid citation recommendation system based on linked data. 2nd International eConference on Computer and Knowledge Engineering (ICCKE), Chongqing, China.
- Zhang, D., Yin, J., Zhu, X., & Zhang, C. (2016). Homophily, structure, and content augmented network representation learning. International Conference on Data Mining, Barcelona, Spain.
- Zhang, J. (2017). Uncovering mechanisms of co-authorship evolution by multirelations-based link prediction. *Information Processing & Management*, 53(1), 42–51.
- Zhang, J., & Yu, W. (2020). Early detection of technology opportunity based on analogy design and phrase semantic representation. *Scientometrics*, 125(1), 551–576.
- Zhang, Y., Satapathy, S. C., Guttery, D. S., Górriz, J. M., & Wang, S. (2021). Improved breast cancer classification through combining graph convolutional network and convolutional neural network. *Information Processing & Management*, 58(2), Article 102439.
- Zhou, D., Zhu, S., Yu, K., Song, X., Tseng, B. L., Zha, H., & Giles, C. L. (2008). Learning multiple graphs for document recommendations. Proceedings of the 17th International Conference on World Wide Web, Beijing, China.
- Zhou, Q., & Zhang, C. (2019). Using citation contexts to evaluate impact of books. Proceedings of the 17th International Conference on Scientometrics and Informetrics, Rome, Italy.