

A Multi-dimensional Disciplinary Knowledge Network Fusion Method Based on Graph Convolutional Autoencoder Models (Postprint)

Authors: Li Hui, Hu Jixia

Date: 2023-04-01T16:16:01+00:00

Abstract

[目的/意义] To address the issue that knowledge networks comprising only a single type of knowledge unit fail to comprehensively reflect disciplinary knowledge structures, this study proposes a multi-dimensional knowledge network structure fusion method, offering insights for mining disciplinary knowledge structures. [方法/过程] The method employs LDA and TF-IDF techniques to extract disciplinary knowledge units, utilizes semantic similarity and keyword co-occurrence analysis to construct three disciplinary knowledge sub-networks—a topic network, a keyword network, and an entity network—adopts a spatial node transitive alignment approach to align sub-network nodes, designs a graph convolution operation-based autoencoder model to represent knowledge nodes, and finally reconstructs the disciplinary knowledge network through cosine similarity computation. [结果/结论] The experimental section demonstrates the approach using the artificial intelligence domain as a case study, constructing and analyzing a disciplinary knowledge network that integrates topics, keywords, and entities. Results indicate that the proposed method effectively reveals research content and knowledge structures within disciplinary fields, providing a valuable reference for disciplinary knowledge discovery and organization research.

Full Text

A Multi-Dimensional Subject Knowledge Network Fusion Method Based on Graph Convolutional Self-Encoding Model

Li Hui, Hu Jixia School of Economics and Management, Xidian University, Xi'an 710126

Abstract

[Purpose/Significance] To address the limitation that knowledge networks containing only a single type of knowledge unit cannot fully reflect the knowledge structure of a discipline, this paper proposes a method for fusing knowledge network structures from multiple dimensions, providing a reference for mining disciplinary knowledge structures. **[Method/Process]** The study utilizes LDA and TF-IDF methods to extract subject knowledge units, then employs semantic similarity and keyword co-occurrence analysis to construct three subject knowledge sub-networks: a topic network, a keyword network, and an entity network. A spatial node transfer alignment method is adopted to align sub-network nodes, followed by the design of a self-encoding model based on graph convolution operations to represent knowledge nodes. Finally, the disciplinary knowledge network is reconstructed by calculating cosine similarity. **[Result/Conclusion]** The experimental section takes the field of artificial intelligence as an example, constructing and analyzing a disciplinary knowledge network that integrates topics, keywords, and entities. Results demonstrate that the proposed method can effectively reveal research content and knowledge structures within disciplinary domains, offering valuable insights for subject knowledge discovery and organization research.

Keywords: Network fusion, Knowledge structure, Node alignment, Graph convolutional neural network, Self-encoding model

Classification Number: G254

DOI: 10.13266/j.issn.0252-3116.2020.18.013

Introduction

The rapid development of internet technology has led to explosive information growth, which, while facilitating knowledge acquisition, also immerses people in vast oceans of dispersed and diverse knowledge, creating challenges for comprehensively understanding knowledge structures at a macro level. As scientific research continues to expand in scope and deepen in content, knowledge across various fields exhibits increasingly complex patterns of interdisciplinary integration. In this context, scholars entering a new field face difficulties in quickly grasping its knowledge structure and current state of development, making the effective organization of domain knowledge information an urgent problem to solve.

Disciplinary knowledge structure reveals the essence of knowledge and interconnections among knowledge elements through different organizational approaches. Based on different relationship types and manifestations, it can be summarized as hierarchical knowledge structures and network knowledge structures. Compared with hierarchical structures, network-based disciplinary knowledge networks have attracted widespread scholarly attention due to their rich and intuitive knowledge representation. Constructing disciplinary knowledge networks provides an effective approach for organizing domain

knowledge information and presenting knowledge structures, and analyzing these networks has become an important method for mining disciplinary knowledge structures and detecting research frontiers. Subject knowledge networks can not only reveal interrelationships among knowledge nodes at the micro level but also reflect the evolution patterns of scientific concepts and research hotspots within a domain. Tracking changes in knowledge structures of emerging disciplines holds significant importance for research managers, scientists, and policymakers.

In recent years, numerous studies on disciplinary knowledge networks have emerged. A review of relevant literature reveals that current research primarily builds upon bibliometrics, constructing knowledge networks from single knowledge units such as authors, institutions, journals, citations, topics, and keywords. These studies focus on discovering research hotspots and collaboration patterns but fail to fully reveal the intrinsic knowledge structure of a discipline. To address these limitations, this paper proposes a disciplinary knowledge network construction method that integrates topics, keywords, and entities, using domain scientific literature as the research object. This approach first extracts topics, keywords, and entities as knowledge units, constructs knowledge association sub-networks for each dimension based on keyword co-occurrence analysis and semantic similarity calculation, then employs node clustering and graph convolutional self-encoding models to mine deeper semantic and structural information among knowledge units, ultimately generating a disciplinary knowledge network that fuses topics, keywords, and entities. By constructing disciplinary knowledge networks from multi-dimensional research content information, this method overcomes the limitations of traditional approaches that characterize knowledge structures solely through co-occurrence relationships of single-dimensional knowledge units, enabling more complete representation of interconnections among domain knowledge elements and achieving comprehensive and accurate revelation of disciplinary knowledge structures.

2 Related Research

2.1 Theoretical Research

Research on knowledge networks has proliferated, continuously enriching their foundational theories, yet a unified definition remains elusive. A. Seufert et al. conceptualize knowledge networks as dynamic frameworks composed of actors, relationships, and institutional characteristics of resource utilization, which accumulate and apply knowledge through knowledge transfer and creation processes to ultimately achieve value creation [4]. Zhao Rongying abstracts knowledge networks as networks with knowledge elements, knowledge points, knowledge units, and knowledge bases as “nodes” and knowledge associations as “edges” or “links” [5]. Gu Donglei describes the connotation of disciplinary knowledge networks from a philosophical perspective, viewing them as network knowledge systems composed of disciplinary knowledge elements and knowledge associations (knowledge links) with distributive, relative truth-value, and

orderly characteristics [6]. Different fields hold varying understandings of knowledge networks' connotation and extension. From the perspective of knowledge organization in library and information science, knowledge networks consist of domain knowledge nodes and knowledge associations [7]. The international management community defines knowledge networks as "a group of people, resources, and their relationships that facilitate knowledge accumulation and utilization through knowledge creation and transfer to promote new knowledge utilization" [8]. Sociology considers knowledge networks as "interpersonal networks" from which individuals can obtain or exchange material, information, knowledge, and intelligence resources. Despite the lack of a unified definition, all conceptualizations of knowledge networks can be understood as interactions (relationships) among knowledge network subjects (nodes).

2.2 Method Research

With the rise of complex network analysis methods and technologies, current research on knowledge network construction primarily builds upon bibliometrics and social network analysis (SNA). Using social network analysis, disciplinary knowledge networks are constructed based on co-occurrence relationships among external document features such as authors, institutions, journals, and citations to generate domain scientific collaboration networks [9-10], co-citation networks [11-13], etc. Through attribute analysis, centrality analysis, core-periphery structure analysis, cohesive subgroup analysis, node clustering, and key node identification, these approaches identify research hotspots and domain collaboration patterns. Such co-occurrence-based methods have yielded abundant research results on disciplinary knowledge structure analysis. As research scopes expand, scholars have delved into internal semantic features of literature such as titles, abstracts, keywords, and full texts to deeply mine intrinsic disciplinary knowledge structures, with applications in agriculture [14], economics [15], medicine [16], and other fields. Lü Penghui et al. [17-19] summarized the structures, characteristics, and evolution research methods, procedures, and mapping processes for citation networks, co-citation networks, and co-word networks, revealing relationships among knowledge network nodes and discussing limitations of corresponding research methods. Guan Peng, Wang Yuefen, et al. [20-22] proposed a topic-integrated disciplinary knowledge network construction and analysis framework, expanding the research scope of disciplinary knowledge networks. They constructed topic-topic association networks using topics' co-occurrence relationships in scientific literature, proposed concepts and measurement methods for topic influence, and later constructed author-topic association two-mode networks based on author-topic models, using authors' centrality indicators in networks to measure author-topic association influence, thereby compensating for limitations of single citation network and author co-authorship network analyses.

In summary, current research on disciplinary knowledge networks primarily involves theoretical discussions of concepts and network construction based on

document external attributes and single content information. Few studies integrate multi-dimensional knowledge units at the content level to construct knowledge networks, and even fewer treat entities as disciplinary knowledge units. Therefore, this paper employs graph convolutional neural networks to design a disciplinary knowledge network construction framework that integrates topics, keywords, and entities, expanding disciplinary knowledge network construction methods and providing new pathways for revealing disciplinary knowledge structures.

3 Multi-Dimensional Knowledge Network Fusion Method

3.1 Research Framework

This study comprehensively employs semantic similarity and graph convolutional neural network methods to design generation methods and integration models for three disciplinary knowledge sub-networks—topics, keywords, and entities—to comprehensively and accurately reveal disciplinary knowledge structures. The specific process is shown in Figure 1 [Figure 1: see original paper]. First, the original dataset is preprocessed to generate a corpus, and disciplinary knowledge units including topics, keywords, and entities are extracted to construct knowledge association sub-networks for each dimension based on semantic similarity among knowledge units. Next, all nodes are clustered to generate a template network, and a node transfer alignment method [23] is adopted to transform each sub-network into a fixed-size network structure. Finally, graph convolution operations combined with self-encoding models are used to fuse the sub-networks and visualize them, clearly revealing disciplinary knowledge distribution and associations.

3.2 Knowledge Unit Extraction

Scientific literature, as the most direct and effective manifestation of scientific research activities, contains and carries research topics, evolutionary trajectories, and development trends across different disciplines [24], holding important reference value and guiding significance for scientific research. This study uses scientific literature as the initial data source to extract topics, keywords, and entities that reflect disciplinary knowledge, establishing a disciplinary knowledge unit representation system. The extraction methods for each dimension are as follows:

3.2.1 Topics. This paper utilizes the LDA [25] topic model for topic extraction. First, the original dataset undergoes preprocessing including cleaning and tokenization. Then, perplexity is used to determine the optimal number of topics. LDA model parameters are set and the LDA program is executed. Finally, the program’s output “topic-word” files are summarized to obtain disciplinary knowledge topics.

3.2.2 Keywords. Keywords provide highly condensed summaries of document

content and can largely reflect disciplinary research content. Keywords can be directly extracted from the corpus.

3.2.3 Entities. An improved TF-IDF algorithm (see formulas (1) and (2)) is used to filter high-frequency words in the corpus. After part-of-speech tagging, N nouns with larger TF-IDF values are selected as entities, with N determined based on the number of topics and keywords.

$$TF-IDF_{wd} = 2 + IDF_w \quad (1)$$

$$IDF_w = \log \left(\frac{DF_w + 1}{|d|} \right) \quad (2)$$

where w represents a word, d represents a document, wd represents the frequency of word w in document d , and $|d|$ represents the number of words contained in document d .

3.3 Knowledge Sub-Network Construction

Using topics, keywords, and entities extracted in Section 3.2 as knowledge nodes, and determining connections based on similarity degrees among knowledge nodes, topic sub-networks, keyword sub-networks, and entity sub-networks are constructed respectively. Document coincidence [27], topic word similarity, and keyword similarity are comprehensively used to calculate topic similarity. Keyword co-occurrence analysis and the LR (Likelihood Ratio) [26] method are respectively used to calculate keyword similarity and entity similarity. Appropriate thresholds are determined, and nodes with similarity greater than the threshold are connected to construct similarity matrices for mining relationship features among knowledge units, thereby building disciplinary knowledge sub-networks.

3.3.1 Topic Sub-Network. The weighted sum of three indicators—document coincidence, topic word similarity, and keyword similarity—is used as the similarity measure between topics, as shown in formula (3).

$$topicsim(i, j) = w_1 \cdot doccoincidence(i, j) + w_2 \cdot featuresim(i, j) + w_3 \cdot keywordssim(i, j) \quad (3)$$

where $topicsim(i, j)$ represents the similarity between topic i and topic j , w_i are weights corresponding to each indicator, and $\sum_{i=1}^3 w_i = 1$.

- (1) **Document Coincidence:** Following the method proposed by Li Hui [27] to determine literature subset scope, the more overlapping documents between topics, the higher their similarity. The Jaccard coefficient between literature subsets of different topics (see formula (4)) is calculated as document coincidence. In formula (4), $num(set_i \cap set_j)$ represents the number

of intersecting documents in the literature subsets of topic i and topic j , while $num(set_i \cup set_j)$ represents the number of union documents.

$$doccoincidence(i, j) = \frac{num(set_i \cap set_j)}{num(set_i \cup set_j)} \quad (4)$$

- (2) **Topic Word Similarity:** Each topic typically selects N words with higher probability values as feature words to describe the topic. Cosine similarity between feature words of topics is used to measure topic word similarity. As shown in formula (5), all feature words across topics are merged, n is the total number of unique feature words, p_{im} represents the weight of feature word m in topic i 's word distribution (with $p_{im} = 0$ if topic i does not contain feature word m).

$$featuresim(i, j) = \frac{\sum_{m=1}^n (p_{im} \cdot p_{jm})}{\sqrt{\sum_{m=1}^n p_{im}^2} \cdot \sqrt{\sum_{m=1}^n p_{jm}^2}} \quad (5)$$

- (3) **Keyword Similarity:** Keywords contained in documents corresponding to topics are transformed into vector space models, and cosine similarity is calculated as shown in formula (6), where tf_{ki} represents the frequency of the k -th keyword of topic i .

$$keywordssim(i, j) = \frac{\sum tf_{ki} \times tf_{kj}}{\sqrt{\sum (tf_{ki})^2} \times \sqrt{\sum (tf_{kj})^2}} \quad (6)$$

3.3.2 Keyword Sub-Network. All keywords appearing in the literature are collected. If two keywords appear in the same document, they have one co-occurrence relationship. A co-word matrix is constructed based on the number of documents where keywords co-occur (co-occurrence frequency) to mine similarity relationships among keywords, with higher co-occurrence frequency indicating higher similarity.

3.3.3 Entity Sub-Network. The LR algorithm is used to identify similarity relationships between entities. LR is an indicator reflecting authenticity, defined as the ratio of the maximum likelihood function under constraints to the maximum likelihood function without constraints, calculated as shown in formula (7). Here, the constraint condition for one entity's occurrence is whether another entity appears, i.e., the conditional probability $p(e_1|e_2)$, calculated as shown in formula (8), representing the probability of entity e_1 occurring given that entity e_2 appears. If two entities frequently co-occur (i.e., have a larger LR value), they have a strong association. Based on this, the LR value between an entity and all other entities is calculated, and entities with larger LR values are selected as similar entities to generate an entity similarity matrix.

$$LR(e_1, e_2) = \frac{p(e_1|e_2)}{p(e_1|not\ e_2)} \quad (7)$$

$$p(e_1|e_2) = \frac{p(e_1, e_2)}{p(e_2)} \quad (8)$$

Based on these similarity matrices, topic sub-networks, keyword sub-networks, and entity sub-networks are constructed as three-dimensional knowledge structure graphs of the disciplinary knowledge network, laying the foundation for subsequent network fusion.

3.4 Knowledge Sub-Network Fusion

Knowledge sub-networks constructed from different dimensions contain identical or similar knowledge nodes—for instance, nodes in topic, keyword, and entity networks may represent the same technology, terminology, or research object. Knowledge network fusion integrates knowledge sub-networks built using different rules into a more complete knowledge network, with the key being determining whether two knowledge units in different networks describe the same object. This includes node alignment and structure fusion.

3.4.1 Node Alignment. This study references the node transfer alignment method proposed in [23] to transform networks of arbitrary size into fixed-size network structures. The algorithm consists of three steps, with the framework shown in Figure 2 [Figure 2: see original paper] and pseudocode in Table 3. The three knowledge sub-networks constructed in Section 3.3 are represented as $G = \{G_1, G_2, G_3\}$, with each sub-network structure represented as $G_P = (V_P, E_P, A_P, X_P)$, where V_P is the node set, E_P is the edge set, A_P is the adjacency matrix of sub-network G_P , and X_P is the node attribute feature matrix.

- (1) **Node Embedding.** Each knowledge unit in the knowledge sub-networks is mapped to a K -dimensional vector space for vectorized representation. Assuming sub-network G_P contains n nodes, the K -dimensional feature vector of the i -th node in the P -th sub-network is denoted as $DBK_{p,i} = \{vec_1, vec_2, \dots, vec_K\}$. All knowledge unit vectors can be represented by the set $R^K = \{R_2^K, \dots, R_N^K\}$, where N is the number of knowledge units.
- (2) **Node Clustering to Generate Template Network.** The K -means clustering algorithm clusters all sub-network nodes into M classes, obtaining M cluster centers $PR^K = \{\mu_2^K, \dots, \mu_M^K\}$ by minimizing the objective function (see formula (9)). Each cluster center is represented by a K -dimensional vector, and the M cluster centers constitute a template network.

$$\arg \min_{\Omega} \sum_{j=1}^M \sum_{R_i^K \in c_j} \|R_i^K - \mu_j^K\|^2 \quad (9)$$

- (3) **Transfer Alignment of All Sub-Network Structures with Template Network.** The distance matrix between each node in each sub-network and the template network node set is calculated. The distance matrix between the K -dimensional node vectors of the P -th knowledge sub-network and the template network is represented as D_P^K , as shown in Table 1. The Euclidean distance between the i -th node of the sub-network and the j -th node of the template network can be calculated using formula (10).

$$D_P^K(i, j) = \sum_K \|DBK_{p,i} - \mu_j^K\|^2 \quad (10)$$

The alignment matrix C_P^K is a binary matrix derived from the distance matrix D_P^K : if the element in row i , column j is the minimum value in row i , the corresponding position in the alignment matrix is 1; otherwise, it is 0, as shown in formula (11):

$$C_P^K(i, j) = \begin{cases} 1, & \text{if } D_P^K(i, j) = \min D_P^K(i, _) \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

In the alignment matrix, each row has only one element equal to 1, with the rest being 0, indicating that each node in the sub-network corresponds to only one node in the template network. Multiple sub-network nodes may correspond to the same template network node because template network nodes are generated through clustering, and sub-network nodes have similarities. It is reasonable for identical or similar nodes to align with the same template node, as shown in Table 2. Additionally, when nodes from two sub-networks align with the same template node, these sub-network nodes are also aligned, making this alignment relationship transitive.

After obtaining the alignment matrix C_P^K , the adjacency matrix with self-loops added for the P -th sub-network is \tilde{A}_P (i.e., $\tilde{A}_P = A + I$, where I is the identity matrix), and the node attribute feature matrix is X_P . The aligned adjacency matrix \bar{A}_P^K and feature matrix \bar{X}_P^K for each sub-network can be calculated using formulas (12) and (13):

$$\bar{A}_P^K = (C_P^K)^T (\tilde{A}_P) (C_P^K) \quad (12)$$

$$\bar{X}_P^K = (C_P^K)^T X_P \quad (13)$$

The pseudocode for the node transfer alignment algorithm is shown in Table 3 .

3.4.2 Structure Fusion. The representation learning problem of network node structural information is transformed into a word representation learning problem. This study uses neural network language models to mine deep semantic information of network node attributes, then characterizes network structure through attribute similarity. After transfer alignment, all sub-networks have M nodes, with N sub-networks totaling $N \times M$ nodes. After sequentially arranging all nodes, a graph convolutional self-encoding neural network model is applied for joint training to obtain node representation vectors that integrate node attribute information and structural information. The algorithm includes three steps, with the framework flow shown in Figure 3 [Figure 3: see original paper]:

- (1) **Preliminary Integration Based on Graph Convolution Operations.** For any sub-network $G_P = (V_P, E_P, \bar{A}_P, \bar{X}_P)$, let $\tilde{A}_P = \bar{A}_P + I_M$, where M is the number of nodes and I_M is the identity matrix. The diagonal matrix D_P is the degree matrix of the adjacency matrix, with diagonal elements $D_P(i, i) = \sum_j \tilde{A}_P(i, j)$ representing the number of edges connected to node i . For network G_P , node embedding $Z^{(P)}$ is obtained through a two-layer graph convolution operation, as shown in formula (14):

$$Z^{(P)} = f(\bar{X}_P, \tilde{A}_P, W_P) = \text{Softmax}(\hat{A}_P \text{ReLU}(\hat{A}_P \bar{X}_P W_P^{(1)}) W_P^{(2)}) \quad (14)$$

where $\hat{A}_P = D_P^{-1/2} \tilde{A}_P D_P^{-1/2}$ represents row normalization of the adjacency matrix \tilde{A}_P , and $W_P^{(1)}$ and $W_P^{(2)}$ represent weight matrices for the first and second layers, respectively. ReLU is the nonlinear activation function. After calculation, each sub-network obtains a node embedding $Z^{(P)}$. Node embeddings from all sub-networks are sequentially arranged to generate an embedding matrix Z containing $N \times M$ nodes. The sigmoid function value $y(x_i, x_j) = \sigma(z_i^T W z_j)$ is calculated, where W is a hyperparameter weight matrix, and z_i, z_j represent embedding vectors of nodes i and j in embedding matrix Z . An appropriate threshold is selected, and edges are established between nodes exceeding the threshold to obtain a preliminary integrated network structure.

- (2) **Node Representation Learning Based on Graph Convolutional Autoencoder.** Both the encoder and decoder structures contain two convolutional layers. In this study's empirical research, the input has 213 nodes with feature dimension 100, and weight parameters are set as shown in Table 4 . Parameters are continuously adjusted to minimize the loss function to obtain more information from the network structure. The final output information X' is used as the node attribute feature vector, which has the same dimension as the input features.

The neural network's convolution operations can aggregate node attributes and

connection information, while autoencoder networks can perform representation learning in an unsupervised manner. This study combines graph convolutional neural networks with autoencoders to build a graph convolutional autoencoder network model. Based on input network nodes and connection information, convolution parameters are set and network training is guided by minimizing the loss function. As shown in part (2) of Figure 3, the network consists of an encoder and a decoder. The encoder is a convolutional encoding layer that aggregates node attributes and connection information, while the decoder reconstructs convolutional features. The output attempts to reconstruct the node's input attributes, with the loss function being the reconstruction loss, as shown in formula (15). Smaller differences between output and input indicate stronger network learning capability.

$$loss = \sum_{i=1}^n \|X_i - X'_i\|^2 \quad (15)$$

- (3) **Network Reconstruction Based on Cosine Similarity.** The output of the graph convolutional autoencoder network model is used as the vector representation of network nodes. Cosine similarity in vector space (see formula (16), where K is the vector dimension) is used to calculate similarity between nodes, and an appropriate threshold is determined to construct the final disciplinary knowledge network.

$$cos_sim = \frac{\sum_{i=1}^K (x_i \cdot y_i)}{\sqrt{\sum_{i=1}^K x_i^2} \cdot \sqrt{\sum_{j=1}^K y_j^2}} \quad (16)$$

4 Empirical Study

As one of the most popular research directions in computer science, “artificial intelligence” has become an indispensable technological resource for social development. Its interdisciplinary nature and broad application prospects have generated numerous research achievements, making it a hot topic in scientific research both domestically and internationally in recent years. To help researchers fully understand the disciplinary knowledge structure of this field and simultaneously verify the effectiveness of the proposed knowledge network fusion method, this study selects Chinese literature in the artificial intelligence domain as experimental data source, mines the disciplinary knowledge structure, visualizes it, and analyzes the experimental results.

4.1 Data Acquisition and Preprocessing

This study retrieved Chinese literature on artificial intelligence from the past decade. After removing a small number of documents with incomplete keyword and abstract information, 10,224 documents were obtained, as shown in Table 5. The HanLP software package [28] was used to preprocess titles and abstracts

through tokenization, bigram extraction, stop word removal, and part-of-speech tagging to generate the corpus required for modeling.

Table 5 Data Sources | Retrieval Time Range | Database Source | Retrieval Expression | Document Type | Retrieval Results | |-----|-----|
 |-----|-----|-----| | 2009/01/01–2019/01/01 | CNKI
 | (SU=artificial intelligence OR TI=artificial intelligence) AND (KY=artificial intelligence OR KY=AI) | SCI, EI, core journals, CSSCI, and CSCD journals | 10,224 documents |

4.2 Experimental Settings and Results Analysis

4.2.1 Experimental Parameter Settings

- (1) **LDA Topic Extraction.** LDA algorithms are well-established. The optimal number of topics T is determined to be 53 by calculating link perplexity. Other parameters are set based on reference [29] and empirical values, as detailed in Table 6 .

Table 6 LDA Model Parameter Description | Parameter | Value | |-----|-----|
 |-----| | Optimal number of topics | 53 | | Number of feature words per topic, $twords = 30$ | 30 | | Gibbs sampling iterations, $niters = 1000$ | 1000 | | Dirichlet prior of text set on latent topics, $\alpha = 50/T$ | $\alpha = 50/T$ | | Dirichlet prior of latent topics on feature word set, $\beta = 0.02$ | 0.02 |

- (2) **Knowledge Sub-Network Construction.** Using the aforementioned knowledge unit extraction and similarity calculation methods, similarities among knowledge units in each dimension are calculated to construct knowledge sub-networks. After multiple experimental comparisons, parameters such as the number of knowledge units per dimension and association thresholds are finally set as shown in Table 7 . Keywords with frequency above 20 and cumulative proportion of 25% are extracted.

Table 7 Knowledge Sub-Network Related Parameters | Knowledge Unit Type | Number of Knowledge Units | Association Threshold | Association Relationships (pairs) | |-----|-----|-----|-----|
 |-----| | Topic | 53 | 0.221 | 372 | | Keyword | 150 | 0.311 | 1,847
 | | Entity | 102 | 0.311 | 1,502 |

- (3) **Node Alignment.** Word vectors are used to represent knowledge unit attributes. In the topic extraction phase, the meanings of topics were manually summarized, which may result in some topic-representing words not appearing in the corpus. Word2Vec models cannot learn vector representations for such words. Considering simplicity of implementation, this study selects the FastText [30] model from the gensim software package to learn knowledge unit word vector representations. While high-dimensional word vectors can richer represent semantic information of phrases, they also increase the number of neural network model parameters and cause

overfitting. Based on references [31-34], the word vector dimension is set to 100.

Knowledge units are clustered using K -means based on word vectors. The silhouette coefficient method is used to determine the optimal number of clusters M , as shown in Figure 4 [Figure 4: see original paper]. The total number of initially extracted knowledge units is 305. Considering that the final number of network nodes should not exceed the total number of knowledge units and that the final node count cannot be too small to minimize information loss, the minimum number of clusters is set to 50, and M greater than 90 is not considered. The final number of clusters M is set to 71. After aligning the three sub-network nodes, 213 nodes with 100-dimensional vector representations are obtained. The aligned adjacency matrix and feature matrix are then calculated as input for the graph convolutional autoencoder network for node attribute feature representation learning.

4.2.2 Experimental Results Analysis

- (1) **Single Sub-Network Analysis.** Based on the aforementioned sub-network construction methods, knowledge association sub-networks are built using knowledge units from three dimensions and visualized using Pajek software. Nodes are classified according to degree and connection weights and distinguished by different colors, as shown in Figures 5-7 [Figure 5: see original paper]-[Figure 7: see original paper]. Nodes with the same degree share the same color, and larger nodes indicate greater importance, with corresponding research content typically representing domain research hotspots.

Figure 5 shows that node sizes in the topic network are relatively similar, as topics tend to represent domain research content at a macro level, such as intelligent robots, environmental monitoring, drive operations, motion models, and robot production. The figure also reveals that artificial intelligence topics primarily focus on technology application-level research, emphasizing applications in manufacturing, healthcare, smart home, and education, with some involving semantic learning, spatial positioning, model algorithms, and computer vision technologies, indicating that artificial intelligence technologies are gradually maturing. Additionally, an isolated node “pathfinding” appears in the topic network, which should theoretically be associated with spatial positioning, positioning matching, motion models, and motion simulation topics related to robot motion path research.

Figure 6 shows the keyword co-occurrence network, which, compared with the topic network, focuses more on fine-grained descriptions of domain research content, such as specific algorithms like support vector machines, Kalman filtering, and rough sets. Keywords cover both natural science research on artificial intelligence technologies themselves and social science applications, including hot technologies like big data, neural networks, machine learning, deep learning, path planning, image processing, and computer vision, as well as applications

in related fields like IoT, intelligent robots, agricultural robots, decision systems, and knowledge engineering. However, the most critical research content remains intelligent robots (the keywords “artificial intelligence” and “robot” have the largest proportions), consistent with the hottest topic being “intelligent robot.”

Figure 7 shows the entity association network, where entities focus more on physical objects representing research objects, industries, groups, and tools, such as maps, images, machines, courses, teachers, culture, manufacturing, new media, conferences, professional committees, and expert systems. To better fuse with topic and keyword networks, entities also include some key technologies like deep learning, intelligent control, and pattern recognition during screening. Both Figures 5-7 contain a few isolated nodes, such as topic “pathfinding,” keywords “theorem proving” and “automation,” and entity “robot kinematics.” Even with the three similarity calculation methods for topics, some association relationships remain undiscovered, making network fusion and reconstruction necessary to mine complete disciplinary knowledge structures.

- (2) **Integrated Network Analysis.** The graph convolutional autoencoder network model proposed in Section 3.4 is used to learn vector representations of network nodes. The neural network model has a learning rate of 0.01 and 1,000 iterations, with an average accuracy of approximately 0.86 after multiple experiments. Node attribute features trained by the model are used to recalculate node similarity. After screening, 3,721 relationships are retained. Based on clustering results and alignment matrices, nodes with similar meanings in the same class are merged and summarized to construct the disciplinary knowledge network shown in Figure 8 [Figure 8: see original paper].

The fused knowledge network contains no isolated nodes and includes macro-level topic nodes such as risk identification, simulation technology, robots, image processing, and data mining, as well as micro-level keywords like convolutional neural networks, artificial neural networks, triples, and tracking algorithms. Entities representing pictures, videos, patients, vehicles, and maps are also included. The network contains three types of nodes, all considered knowledge units in the artificial intelligence domain. Compared with traditional single sub-networks, the fused knowledge network structure can more comprehensively reflect disciplinary research content and knowledge structure. Specific analysis follows:

First, Research Hotspots. Larger nodes in the figure include grasping posture, robotic arms, robot teaching, robots, agricultural robots, industrial robots, tracking algorithms, routes, neural networks, artificial neural networks, convolutional neural networks, pictures, maps, image detection, image processing, optimization, simulation technology, information resources, risk identification, and environmental recognition, covering artificial intelligence-related issues such as intelligent robots, path planning, deep learning, computer vision, information retrieval, and decision guidance. These contents align with key research areas in artificial intelligence as introduced in the “2019 Artificial Intelligence Devel-

opment Report” [35].

Second, Knowledge Associations. In the integrated network, hotspot research issues and key technologies are closely connected with other knowledge nodes. For example, nodes connected with robots, image detection, pictures, and neural networks are numerous. Additionally, at the network periphery, knowledge units related to intelligent education, machine learning, intelligent machinery, medical robots, and simulation are densely connected, forming small communities with relatively sparse connections to other knowledge clusters.

- (3) **Discussion.** The knowledge network construction model proposed in this paper is an unsupervised learning method. There is no authoritative knowledge network for the artificial intelligence domain, and traditional evaluation metrics like precision, recall, and accuracy are unsuitable for validating the method’s effectiveness. Following relevant literature [36], the method’s effectiveness is analyzed from three aspects:

First, Sub-Network Node Repetition Rate. Comparing nodes across the three sub-networks reveals certain repetitions among topic, keyword, and entity nodes. The node repetition rate is represented by the ratio of repeated nodes to total nodes, as shown in formula (17):

$$\text{Node Repetition Rate} = \frac{\text{Number of Repeated Nodes}}{\text{Total Number of Nodes}} \quad (17)$$

Counting duplicate words and synonyms (e.g., “artificial intelligence” and “AI”) in the knowledge unit collection, the three sub-networks constructed in this study contain 305 nodes, with 48 nodes repeated across two networks (a two-network repetition rate of 15.74%) and 4 nodes repeated across all three networks (intelligent manufacturing, computer vision, intelligent robots, and knowledge engineering), yielding a three-network repetition rate of 1.312%. The relatively high two-network repetition rate means that directly integrating through node word vector similarity calculation would introduce redundancy. After node clustering, identical or similar knowledge units are grouped into the same class, and re-summarizing knowledge units within each class effectively reduces knowledge unit repetition rates.

Second, Proportion of New Nodes in Fused Network. Based on domain background and content of knowledge units in sub-networks, knowledge unit names after clustering are summarized. Most original knowledge nodes are retained, with similar node meanings re-summarized. The fused network contains 213 nodes, including 150 original nodes (70.42% retention rate) and 63 new nodes (29.58% proportion), as shown in formula (18):

$$\text{Node Retention (New) Rate} = \frac{\text{Number of Original (New) Nodes in Fused Network}}{\text{Total Number of Nodes in Fused Network}} \quad (18)$$

Using clustering for node alignment preserves most knowledge units from original single networks while supplementing them with new nodes, enabling more comprehensive representation of a domain's research content.

Third, Proportion of New Edges in Fused Network. Relationships in original knowledge sub-networks only include connections among the same type of knowledge units (e.g., topic-topic). The fused knowledge network includes both intra-sub-network structures and inter-sub-network associations. The fused knowledge network has 3,721 edges, with 2,324 new edges (between new nodes, new edges between original nodes, and edges between original and new nodes), accounting for 62.46%. Thus, 37.54% of relationships come from original knowledge networks. The proposed graph convolutional autoencoder model effectively aggregates original network structures while mining additional knowledge association relationships.

Through the above analysis, the fused knowledge network performs well in discovering domain research hotspots and mining knowledge association relationships. Node clustering and graph convolutional autoencoder models can aggregate sub-network structures and node attribute information, comprehensively and accurately revealing association relationships among domain knowledge units. As an unsupervised neural network model learning knowledge node feature information, the method requires no labeled data, making it applicable to other scenarios and valuable for heterogeneous network fusion.

Conclusion

This paper proposes a disciplinary knowledge network construction method that integrates topics, keywords, and entities. First, natural language processing methods preprocess Chinese corpora, using LDA and TF-IDF methods to extract topics and entities from the artificial intelligence domain and extract keywords from the corpora. Second, disciplinary knowledge sub-networks are constructed based on semantic similarity and keyword co-occurrence analysis. Then, a graph convolutional autoencoder network model is designed to learn knowledge unit vector representations. Finally, the entire disciplinary knowledge network is reconstructed using cosine similarity to mine disciplinary knowledge structures. Analysis and discussion of the artificial intelligence domain knowledge network demonstrate the method's effectiveness and accuracy.

Extracting and effectively organizing knowledge points in disciplinary domains can help researchers quickly understand domain research hotspots and knowledge structures. Existing knowledge network construction methods rarely involve multi-dimensional knowledge fusion. The proposed knowledge network fusion method not only captures semantic information of knowledge units but also aggregates structural information of nodes in sub-networks. This unsupervised knowledge unit representation learning method is more efficient, and the learned node vectors are generalizable for solving problems like knowledge unit clustering and classification.

Limitations include: (1) manual summarization of topic and clustered knowledge unit meanings introduces subjectivity; (2) network reconstruction through cosine similarity calculation increases workload, warranting future exploration of more advanced edge prediction algorithms; (3) the empirical study only validates Chinese literature data, which may not provide comprehensive disciplinary knowledge content, suggesting future consideration of fusing multiple data sources for disciplinary knowledge network construction.

References

- [1] Zhao Rongying. On the structure of knowledge networks[J]. Library and Information Service, 2007, 51(9): 6-9.
- [2] Gu Donglei. On disciplinary knowledge networks[J]. Journal of Intelligence, 2008(9): 50-55.
- [3] Wang Xiaoguang. The formation and evolution of scientific knowledge networks (I): The proposal of co-word network methods[J]. Journal of the China Society for Scientific and Technical Information, 2009, 28(4): 599-605.
- [4] Seufert A, Krogh G, Bach A. Towards knowledge networking[J]. Journal of knowledge management, 1999, 3(3): 180-190.
- [5] Zhao Rongying. Knowledge networks and their applications[M]. Beijing: Beijing Library Publishing House, 2007: 8-58.
- [6] Gu Donglei. On disciplinary knowledge networks[J]. Journal of Intelligence, 2008(9): 50-55.
- [7] Kou Jihong. Research on visual construction of disciplinary knowledge networks: Taking competitive intelligence as an example[J]. Journal of Information Resources Management, 2015, 5(3): 71-77.
- [8] Xiao Dongping. Review of knowledge network research[J]. Journal of Chongqing Technology and Business University (Natural Science Edition), 2006(6): 617-623.
- [9] Wang Yuefen, Li Dongqiong, Yu Houqiang. Research on evolution of scientific collaboration networks and growth characteristics of high-impact scholars in life cycle stages[J]. Journal of the China Society for Scientific and Technical Information, 2018, 37(2): 121-131.
- [10] Pan Youneng, Tan Jian. Research on scientific collaboration networks of Price Medal winners[J]. Library and Information Service, 2012, 56(16): 80-84.
- [11] Qiu Junping, Zhou Yi. Analysis of deep aggregation model and service of collection resources based on author co-citation: Taking ontology research in library and information science in CSSCI as an example[J]. Library and Information Service, 2014, 58(7): 19-24.

- [12] Hou Jianhua. Visual detection of international scientometrics research frontiers: Analysis of co-citation networks of Scientometrics journal literature[J]. Modern Information, 2012, 32(10): 61-65.
- [13] Jiang Chunlin, Zhang Fan, Tang Yue. Research on co-citation network characteristics of some Chinese scientometrics journals[J]. Information Science, 2010, 29(4): 10-15, 25.
- [14] Liu Qiuxia, Wu Hanqing, Huang Zhenglai. Research on wheat response to climate warming based on global bibliometrics[J]. Chinese Agricultural Science Bulletin, 2019, 35(23): 142-151.
- [15] Luo Rundong, Teng Kuan, Li Chao. Analysis of hot topics in Chinese economics research in 2018[J]. Economic Perspectives, 2019(4): 80-98.
- [16] Zhang Yiqing, Wang Gaoling. Comparative analysis of health management research at home and abroad based on knowledge mapping[J]. Chinese General Practice, 2019, 22(9): 1112-1118.
- [17] Lü Penghui, Zhang Shijing. Research on disciplinary knowledge networks (I): Structure, characteristics and evolution of citation networks[J]. Journal of the China Society for Scientific and Technical Information, 2014, 33(4): 340-348.
- [18] Lü Penghui, Zhang Ling. Research on disciplinary knowledge networks (II): Structure, characteristics and evolution of co-citation networks[J]. Journal of the China Society for Scientific and Technical Information, 2014, 33(4): 349-357.
- [19] Zhao Yiming, Lü Penghui. Research on disciplinary knowledge networks (III): Structure, characteristics and evolution of co-word networks[J]. Journal of the China Society for Scientific and Technical Information, 2014, 33(4): 358-366.
- [20] Guan Peng, Wang Yuefen, Cao Jiajun. Research on framework for constructing and analyzing topic-integrated disciplinary knowledge networks[J]. Information Science, 2018, 36(9): 3-8.
- [21] Wang Yuefen, Wang Jinshu, Guan Peng. Construction and evolution analysis of topic-topic association disciplinary knowledge networks[J]. Information Science, 2018, 36(9): 9-15, 102.
- [22] He Jin, Guan Peng, Wang Yuefen. Construction and evolution analysis of author-topic association disciplinary knowledge networks[J]. Information Science, 2019, 37(1): 56-62, 67.
- [23] Bai L, Jiao Y, Cui L, et al. Learning aligned-spatial graph convolutional networks for graph classification[C]//ECML PKDD 2019. Machine learning and knowledge discovery in databases. Würzburg: Springer, 2019: 464-482.
- [24] Hu Yuning, Hu Guanwei. Construction and empirical research of scientific knowledge structure model based on multi-source topic fusion[J]. Information Studies: Theory & Application, 2019, 42(7): 100-105.

- [25] Blei D M, Ng A Y, Jordan M I. Latent dirichlet allocation[J]. Journal of machine learning research, 2003(3): 993-1022.
- [26] Michael K. The lokahi prototype: Towards the automatic extraction of entity relationship models from text[C]//Proceedings of the AAAI 2019 spring symposium on combining machine learning with knowledge engineering (AAAI-MAKE 2019). Palo Alto: Stanford University, 2019: 121-126.
- [27] Li Hui, Tian Yadan. A hierarchical method for discovering scientific knowledge structures[J]. Library and Information Service, 2018, 62(13): 92-102.
- [28] Shanghai Linyuan Information Technology Co., Ltd. HanLP[EB/OL]. [2020-01-24]. <http://www.hanlp.linrunsoft.com/>.
- [29] Wang Peng, Gao Cheng, Chen Xiaomei. Research on text clustering based on LDA model[J]. Information Science, 2015(1): 63-68.
- [30] Bojanowski P, Grave E, Joulin A, et al. Enriching word vectors with subword information[C]//The conference on transactions of the Association for Computational Linguistics. Prague: ACL, 2017: 135-146.
- [31] Ye Zhonglin, Zhao Haixing, Zhang Ke, et al. Distributed word representation learning based on multi-source information fusion[J]. Journal of Chinese Information Processing, 2019, 33(10): 18-30.
- [32] Ye Zhonglin, Zhao Haixing, Zhang Ke, et al. Word representation learning based on descriptive constraints[J]. Journal of Chinese Information Processing, 2019, 33(4): 29-36.
- [33] Lai Wenhui, Qiao Yupeng. Spam SMS recognition method based on word vectors and convolutional neural networks[J]. Computer Applications, 2018, 38(9): 2469-2476.
- [34] Wu C, Gao R, Zhang Y, et al. PTPD: Predicting therapeutic peptides by deep learning and word2vec[J]. BMC bioinformatics, 2019, 20(15): 87-108.
- [35] Tsinghua University-Chinese Academy of Engineering Joint Research Center for Knowledge Intelligence, Beijing Wu Wenjun Artificial Intelligence Science and Technology Award Selection Base. 2019 Artificial Intelligence Development Report[EB/OL]. [2020-01-24]. https://www.sohu.com/a/360140139_{468661}.
- [36] Lu R, Fei C, Wang C, et al. HAPE: A programmable big knowledge graph platform[J]. Information sciences, 2020(509): 87-103.

Author Contributions

Li Hui: Proposed research ideas and provided paper revision suggestions.

Hu Jixia: Responsible for experimental implementation, paper writing, and revision.

Multi-Dimensional Subject Knowledge Network Fusion Method Based on Graph Convolutional Self-Encoding Model

Li Hui, Hu Jixia

School of Economics and Management, Xidian University, Xi'an 710126

Abstract: [Purpose/Significance] To address the limitation that knowledge networks containing only a single type of knowledge unit cannot fully reflect the knowledge structure of a discipline, this paper proposes a method for fusing knowledge network structures from multiple dimensions, providing a reference for mining disciplinary knowledge structures. [Method/Process] This study uses LDA and TF-IDF methods to extract subject knowledge units, employs semantic similarity and keyword co-occurrence analysis to construct three subject knowledge sub-networks (topic network, keyword network, and entity network), adopts spatial node transfer alignment to align sub-network nodes, designs a self-encoding model based on graph convolution operations to represent knowledge nodes, and finally reconstructs the disciplinary knowledge network by calculating cosine similarity. [Result/Conclusion] The experimental section takes the artificial intelligence domain as an example, constructing and analyzing a disciplinary knowledge network integrating topics, keywords, and entities. Results demonstrate that the proposed method can effectively reveal research content and knowledge structures within disciplinary domains, providing valuable insights for subject knowledge discovery and organization research.

Keywords: Network fusion, Knowledge structure, Node alignment, Graph convolutional neural network, Self-encoding model

Note: Figure translations are in progress. See original paper for figures.

Source: ChinaXiv — Machine translation. Verify with original.