

Knowledge Graph Analysis of Internal Control Field in Colleges

Jun WANG

Abstract: Knowledge graph is a new method to describe concepts, instances and their relationships in the objective world. In recent years, it has attracted people's wide attention. Knowledge graph can effectively expand the breadth of search results. At present, in the field of internal control of colleges and universities, the key word search technology is mainly used to retrieve relevant knowledge content. It is difficult to retrieve information by using the relation between objects. Therefore, this paper first proposes a knowledge graph construction method for internal control in universities, through which the knowledge graph of internal control policy in universities can be constructed. Then, this paper proposes a knowledge inference method based on inference rules, which USES the knowledge hidden in the knowledge graph of internal control policy to realize intelligent data retrieval. Finally, the *CiteSpace* software is used to realize the visualization display in the field of internal control in universities and realize the construction of internal control knowledge graph and visualization research. The system can effectively utilize the relationship between knowledge objects of internal control and give full play to the value of information resources of internal control.

Keywords: CiteSpace; Construction Methods; Internal Control; Knowledge Graph; Visualization

1 INTRODUCTION

As an important part of administrative institutions, colleges and universities shoulder important responsibilities such as scientific and technological innovation, personnel training, social services and cultural inheritance. At present, colleges and universities have entered a period of rapid development. Colleges and universities are not only expanding in size, but also diversifying their revenues and expenditures. With the increasing complexity of economic activities, colleges and universities have become increasingly large public welfare socialized organizations, and the corresponding management risks are also increasing. However, the level and mode of internal control of colleges and universities are far behind the expansion of the scale of colleges and universities. Although the internal control of enterprises has achieved long-term development, the internal control of administrative institutions, especially colleges and universities, is neglected to study and promote. Therefore, administrative institutions have not been given enough attention. Among them, the internal control knowledge and the internal control policy involve a lot of knowledge, with a wide range of aspects and a deep content, which leads to the university administrators' insufficient in-depth and thorough grasp of the internal control knowledge of universities, so the implementation of internal control lacks theoretical support [1]. This is an important reason why the internal control of administrative institutions is not valued. How to enable university managers to learn internal control knowledge efficiently has become the focus of researchers in this field.

In 2012, Google first proposed the concept of knowledge graph and successfully applied it to semantic search, which greatly improved the quality of search results. A knowledge graph is essentially a knowledge base called the semantic web. It can be expressed as a knowledge base with a directed graph structure, where the nodes of the graph represent instances or concepts, and the edges of the graph represent various semantic relationships between instances/concepts. Knowledge graph provides a way for people to better organize, manage, and understand huge amounts of information. At the same time, knowledge

graph has also become the foundation of knowledge retrieval, intelligent question answering, personalized recommendation and other applications. Knowledge graph technology is applied to internal control in universities. The knowledge graph of university internal control constructed by using the object data of internal control policy documents can make use of the correlation between the objects of university internal control to expand the scope of the retrieval results. Aiming at the data of internal control objects in colleges and universities, this paper proposes a method to construct the knowledge graph of internal control information in colleges and universities. On this basis, knowledge inference based on inference rules is implemented to further explore the knowledge hidden in the knowledge graph of internal control information in colleges and universities. Finally, the above technology is applied to the knowledge graph construction and visualization application of internal control information in colleges and universities to realize the intelligent retrieval and recommendation of internal control information in colleges and universities.

As an important research method of information science and scientific metrics, graph of scientific knowledge has been extended to other fields due to its unique advantages. Scientific knowledge graph is mainly studied in the domain of knowledge, with dual attributes of graph and spectrum. Therefore, scientific knowledge graph cannot only realize the visualization of knowledge graph, but also realize the serialization of knowledge graph. This paper intends to use *CiteSpace* analysis software to make a visual analysis of the research results of the current university internal control knowledge, so as to grasp the overall context of the current university internal control knowledge and further explore the current research situation in this field.

2 LITERATURE REVIEW

2.1 Internal Control of Colleges and Universities

The audit procedure committee of AICPA (American Institute of Certified Public Accountants) (1949) [3] gave the first authoritative definition of internal control. COSO committee (2013) issued a new version of internal control

framework, which strengthened the concept of comprehensive internal control. In addition to being more general, the new framework can be widely applied to entities such as government agencies, non-profit and for-profit organizations. In [5] the authors believe that it is the essence of effective internal control to collect information extensively and conduct sorting and analysis, because the authenticity of financial information depends on the reliability of information. Lack of internal control will have a significant impact on enterprises and units. In [6] the authors investigated and studied the companies before the implementation of SOX act, and found that the companies with insufficient investment in internal control had larger internal control defects, such as the company's complex operation and frequent change of auditors. In [6] and [7] the authors in order to analyze the influencing factors of internal control, they tested the position of the company's internal defects. Through the study, it is found that the earlier the company is established and listed, the more comprehensive its internal system will be. In [9] the authors pointed out that the company's performance is positively correlated with the quality of internal control, and only when the company achieves the established financial goals can managers attach importance to internal control. In [10] the authors found that when a company complied with SOX404, the characteristics of the audit committee were affected by the quality of internal control. While the authors in [8] pointed out that whether the company made mistakes or significant defects was not affected by the size of the audit committee of the company. The authors in [9] constructed a knowledge model and text model for knowledge maps.

2.2 Knowledge Graph Application

Generic knowledge graphs cover all areas. The graph not only contains a lot of common sense knowledge, but also emphasizes the breadth of the knowledge graph. At present, many achievements have been made in the field of general knowledge mapping. For example, *DBpedia* (2016) [10] built with knowledge extracted from Wikipedia; The integration of Wikipedia, WordNet and GeoNames (2009) knowledge formed by *Yago* [11]; *NELL* (2008) [12] and *Concept Graph* (2010) [13], *Freebase* (2015) [14], *Wikidata* (2017) [15] were constructed by using knowledge extracted from the Internet. In recent years, many general knowledge graphs have appeared in China. For example, *CNDBpedia* (2017) [16] published by Fudan University integrates encyclopedia data and knowledge of some fields; *Zhishi.me* published by Shanghai Jiao Tong University, includes data of Chinese Wikipedia, interactive encyclopedia and Baidu encyclopedia. Based on Chinese and English encyclopedia *XLore*, Baidu bosom friend, *Sogou's* dog cube, etc. These generic knowledge graphs cover a wide range of areas. They cannot only provide ordinary users with intelligent questions and answers, but also provide better services in personalized recommendation.

Vertical domain knowledge graphs are usually oriented towards a specific domain. The knowledge granularity of the graph is more detailed, and more attention is paid to the depth and completeness of knowledge. For example, by

extracting medical facts from EMR information, the knowledge graph of breast cancer was constructed. The construction of cultural relics ontology is instantiated with the information of cultural relics, and the cultural relics knowledge graph is formed extracting information from different software resources, constructing software knowledge instances and forming software knowledge atlas; Triples are extracted from the data of carbon trading field, and then converted into associated data to build a knowledge graph of carbon trading field. In addition, you can build GeoNames (geographic domain), IMDB (movie & TV domain), MusicBrainz [17]. These knowledge graphs oriented to a particular domain not only help to give full play to the value of domain data, but also lay a foundation for the subsequent intelligent application research. Although there are some domain knowledge graphs, knowledge graph technology in the field of university internal control has not been widely concerned. At the same time, these domain knowledge graphs are mostly in the aspect of construction, and the application of knowledge graphs is not considered.

3 METHODS OF CONSTRUCTING KNOWLEDGE GRAPH OF INTERNAL CONTROL IN COLLEGES AND UNIVERSITIES

The knowledge graph in this paper adopts the top-down construction method, which constructs the concept layer and the instance layer of the knowledge graph successively. As shown in Fig. 1, the concrete method of constructing conceptual layer and instance layer of knowledge graph of internal control in colleges and universities is given.

3.1 Construct Concept Layer

(1) Construct concept node C . According to *General Rules for Classification and Coding of Internal Control Objects*, with the help of domain experts, concept node C , such as financial concept node, can be determined by combining internal control object data. Each concept node corresponds to an object directory table and an object base information table, such as the financial concept node corresponding to the financial object directory table and the financial base information table.

(2) Construct attribute side P_C and attribute value type node N_C . First, determine the fields and field types in the object base information table corresponding to the concept node. On this basis, the field is extracted as the attribute side P_C of concept node C , and the field type is extracted as the attribute value type node N_C . For example, the engineering equivalent field in the basic financial information table is extracted as the attribute edge of the financial concept, and the corresponding attribute value type node is the type of the engineering equivalent field, that is, the integer type.

(3) Construct the inter-concept relationship R_C . Construct the bipartite graph of concept relation, and determine the relationship R_C between concept nodes C through the graph of bipartite graph. The definition of the concept relation bi-graph is as follows:

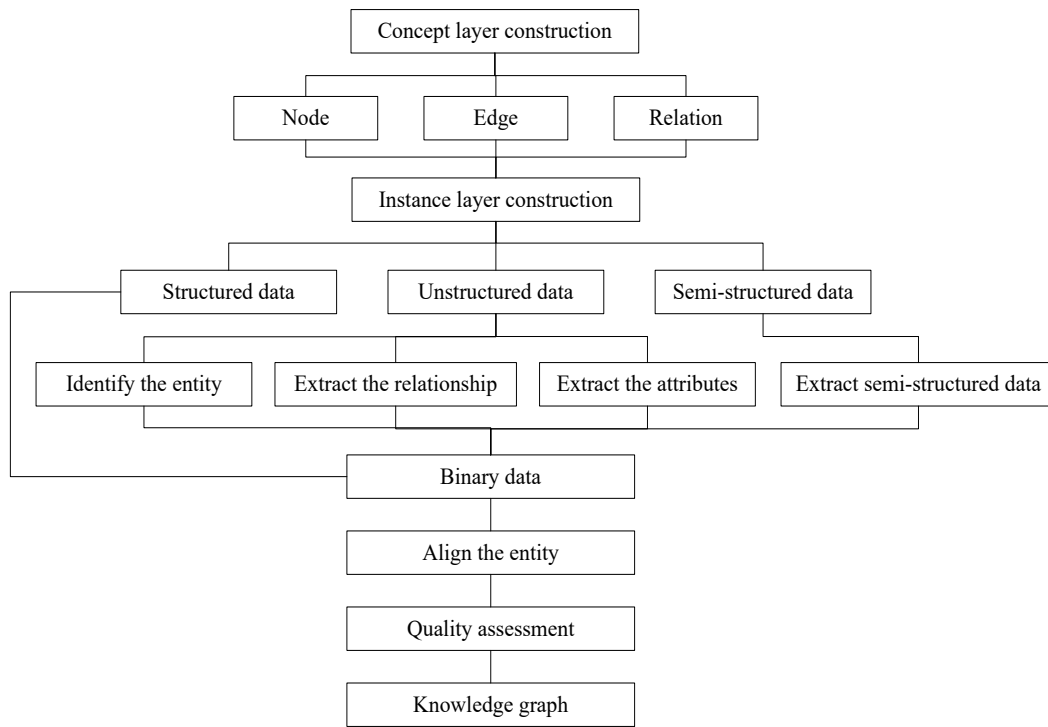


Figure 1 Knowledge graph construction process

Definition 1 (concept layer G_C)

The concept layer is the core of the knowledge graph, which describes the data pattern of the knowledge graph and regulates the facts in the instance layer. In the knowledge graph of internal control information, the conceptual layer $G_C = (C, P_C, N_C, R_C)$, where:

C represents the concept nodes in the graph, such as financial concept and audit concept;

P_C represents the attribute side of the concept node, such as the attribute side of the financial concept including the engineering class and the financial type;

N_C represents the node of attribute value type. For example, the node of other attribute value types such as engineering is integer, and the node of financial attribute value type is character type;

R_C represents the relationship between concept nodes and concept nodes, such as the inflow relationship between financial concepts and audit concepts.

Definition 2 (instance layer G_E)

The instance layer consists of a set of facts. In the knowledge graph of internal control information, instance layer $G_E = (E, P_E, N_E, R_E)$, where:

E represents instance nodes in the graph, such as finance and audit;

P_E represents the attribute edge of the instance node, such as internal control and financial type;

N_E represents attribute value node, such as internal control and other attribute value nodes are 1;

R_E represents the relationship between instance nodes and instance nodes, such as the inflow relationship between finance and audit.

Definition 3 (knowledge graph G of internal control information)

Knowledge graph of internal control information $G = (G_C, G_E, R)$, where:

G_C represents the conceptual layer of knowledge graph of internal control information;

G_E represents the instance layer of knowledge graph of internal control information;

R represents the relationship between concept nodes in G_C and instance nodes in G_E , such as the relationship between instance concepts between financial concepts and financial operations.

Meanwhile, each instance node E has a relationship R with only one concept node C . If the two instance nodes E_1 and E_2 have an inter-instance relationship R_E . Moreover, E_1 and E_2 have the relations between instance concepts R_1 and R_2 respectively with concept nodes C_1 and C_2 , so there must be R_C between C_1 and C_2 , and the relations R_C and R_E have the same name. For example, there is an inflow relationship between finance and financial practice, there is an inflow relationship between financial practice and financial concept, and there is an inflow relationship between audit practice and audit concept. Fig. 1 is the structure of partial internal control information knowledge atlas.

Definition 4 (concept relation bipartite diagram G_B)

$G_B = (V_C, V_T, E_B)$, where:

V_C represents the concept node;

V_T represents the relational table node;

E_B represents the edge between the conceptual node V_C and the relational table node V_T .

Firstly, the concept relation bipartite graph is constructed, and then the relationship between different object directory tables and object relation tables is obtained by using the relational schema extracted from the database. The object directory table and object relational table are then mapped to nodes in the graph. The concept nodes are also treated as nodes in the diagram. The concept node is replaced with the corresponding object directory node, and the concept relation bipartite graph is obtained. By using the mapping algorithm of concept relation bisection graph, the bisection graph is mapped to the graph that only has V_C of concept node set. Thus, the relationship R_C between

concept nodes can be obtained, and the conceptual layer of knowledge graph of internal control information can also be constructed.

Key algorithms are shown in Tab. 1.

Table 1 Mapping Algorithm of Concept Relation Bipartite Graph

Input: $G_B = (V_C, V_T, E_B)$ Output: conceptual layer of knowledge graph of internal control in universities while($V_T \neq null$) Take A relational table node A from V_T ; If(A is connected to 2 concept nodes) Obtain two concept nodes B and C connected to A; Add a relationship between B and C; Delete the edges between A and B and C; else Get concept node B connected to A; Add a relation R_C pointing from B to B itself; Delete the edge between A and B; End if Specify the relationship trigger word description R_C ; End while
--

3.2 Construct Instance Layer

(1) Construct instance node E and the relationship R between concepts and instances.

According to the object directory table corresponding to concept node C in G_C , the corresponding instance node E can be constructed by extracting the name of internal control object, such as finance and audit. Meanwhile, the relationship R between concept and instance is added between concept node C in G_C and corresponding instance node E in G_E . For example, add relationship R between the financial node and the audit concept node.

(2) Construct attribute side P_E and attribute value node N_E .

According to the property side P_C of C and the property value type node N_C in G_C , the same field as P_C is extracted as the property side P_E of instance node E , and the corresponding field value is extracted as the property value node N_E of its property side P_E . For example, for finance, fields with the same name are extracted from the basic information table of internal control according to the attribute edges of the previously constructed concept node of internal control. The audit type and audit level are the attributes of the internal control node. The corresponding node with field value 1 can be extracted as the attribute value node of the attribute side.

(3) Construct the relationship R_E between instances.

Build the inter-instance relationship R_E from the object relation table. According to the name of R_C , determine the name of R_E , and complete the construction of knowledge graph of internal control information. According to the audit and report relationship table, there is a relationship between the finance and the financial statements, and a relationship is added between the two instances. At the same time, there is an inflow relationship between finance and audit, as there is an instance relationship between finance and audit.

4 MODEL CONSTRUCTION

4.1 Knowledge Model

We define the conditional probability of events in a knowledge graph as (h, r, t) , as shown in formula (1).

$$Pr(h | r, t) = \frac{\exp\{z(h, r, t)\}}{\sum_{\tilde{h} \in \tau} \exp\{z(\tilde{h}, r, t)\}} \quad (1)$$

We have named the model $pTransE$ (probabilistic Trans E). We define $Pr(r | h, t)$ and $Pr(t | h, r)$ in the same way by selecting the corresponding normalization terms respectively. We define the possibility of observing triple events as formula (2):

$$T_f(h, r, t) = \log Pr(h | r, t) + \log Pr(t | h, r) + \log Pr(r | h, t) \quad (2)$$

The goal of the knowledge model is to maximize the conditional likelihood probability of triple events in the knowledge graph as formula (3):

$$T_K = \sum_{(h, r, t) \in \Delta} T_f(h, r, t) \quad (3)$$

4.2 Text Model

We propose the following key assumptions for text modeling that link word embedding with knowledge embedding. Although we do not know what words are, there is a relationship between words.

According to Eq. (1), we defined the same conditional probability $Pr(\mathbf{w} | r_{wv}, \mathbf{v})$ to construct a model to evaluate the occurrence of two words at the same time. As opposed to knowledge embedding, variable r_{wv} is a hidden variable, not an explicit variable.

The challenge now is to deal with the hidden variable r_{wv} . Obviously, without any more assumptions, the number of different r_{wv} is set to $|\mathbf{v}| \times \bar{N}$, where N is the average number of unique words that each word has in common. The number N is very large. Therefore, it is almost impossible to estimate one vector for each r_{wv} . We have to restrict the degree of freedom of r_{wv} . Here, we use auxiliary variables to reduce the size of the variable we need to estimate: let $\mathbf{w}' = \mathbf{w} + r_{wv}$, then Eqs. (4) and (5) can be obtained:

$$z(\mathbf{w}, r_{wv}, \mathbf{v}) \triangleq z(\mathbf{w}', \mathbf{v}) = b - \frac{1}{2} \|\mathbf{w}' - \mathbf{v}\|^2 \quad (4)$$

$$Pr(\mathbf{w} | r_{wv}, \mathbf{v}) \triangleq Pr(\mathbf{w} | \mathbf{v}) = \frac{\exp\{z(\mathbf{w}', \mathbf{v})\}}{\sum_{\mathbf{w} \in \mathbf{v}} \exp\{z(\mathbf{w}', \mathbf{v})\}} \quad (5)$$

So we need to estimate the vectors \mathbf{w} and \mathbf{w}' for each term, and there are $2 \times |\mathbf{v}|$ vectors. The goal of the text model is to maximize the likelihood of word pairs appearing together.

$$\Gamma_T = \sum_{(\mathbf{w}, \mathbf{v}) \in C} n_{wv} \log Pr(\mathbf{w} | \mathbf{v}) \quad (6)$$

In formula (6), C is all word pairs that appear simultaneously in a fixed-size text window. n_{wv} is the co-

occurrence number (w, v) of the pair. Interestingly, this text model is almost equivalent to Skip-Gram.

5 EXAMPLE

5.1 Sample Collection

In order to achieve better results of visual analysis, targeted topic retrieval and data collection of a certain scale are needed. In view of the consideration of time and efficiency, this paper selects the core journals in CNKI database (including CSSCI and Peking University Chinese core journal) from 2010 to 2019 as the retrieval source, and selects relevant literature on *Xi Jinping's* thought on ecological civilization as the data source. In order to ensure the integrity of sample collection, the retrieval expression set in this paper is: "subject = (TI= 'university internal control policy')". The search time range of this article is from January 1, 2010 to August 31, 2019, and the matching method is accurate. After the advanced search of the literature by setting the above conditions, a total of 2,450 journal samples were obtained. In order to focus on the research topic, this paper excludes non-study samples such as book reviews, conference notices, and news reports. The number of documents that were ultimately related to the topic was determined to be 1,500, and this sample was used as a sample library for the next analysis. Although these data are not completely covered, they are still representative.

5.2 Knowledge Graph Analysis of Internal Control in Colleges and Universities

Key words have a high-level generalization of the research content and theme of an article. A node in the graph represents a keyword, and the size of the node is proportional to the co-occurrence frequency of the keyword. The larger the node, the higher the co-occurrence frequency, the stronger the research hotspot of the word. The obtained sample data is imported into *Citespace*. After the preliminary test, the clustering result is set reasonably. The time interval is set to 2010-2019, the time slice defaults to 1, and the data extraction threshold is "TOP50".

The keyword knowledge graph is obtained (see Fig. 2). For example, in Fig. 3, the node where the university is located has the most annual rings, indicating that it appeared earlier in the research field. The yellow squares outside the nodes such as internal control, risk prevention, risk assessment, risk-oriented, and risk management reflect that this is still the current research hotspot. After MTS graph pruning, the number of nodes in the keyword co-occurrence graph is 172, the number of nodes is 165, and the network density is 0.0112.

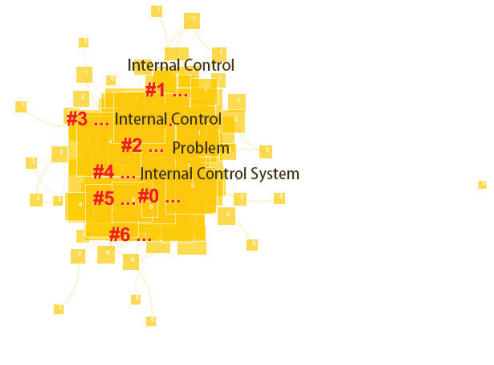


Figure 2 Key words knowledge graph

By cluster analysis of keywords, the structural features between clusters can be found, so as to highlight the key nodes and their important connections in the network. In *Citespace*, two important indexes to evaluate the mapping effect are module value and average contour value respectively. The former is Q value, which is generally taken 0.3 as the standard, and greater than 0.3 means that the classified community structure is significant; the latter is S value, whose value is above 0.7, and the clustering effect is convincing. According to Fig. 3, *Modularity Q* of the co-occurrence graph is 0.8238 and *MeanSilhouette* is 0.6252, indicating that the knowledge structure presented by the knowledge graph is significant and clustering is efficient and ideal.

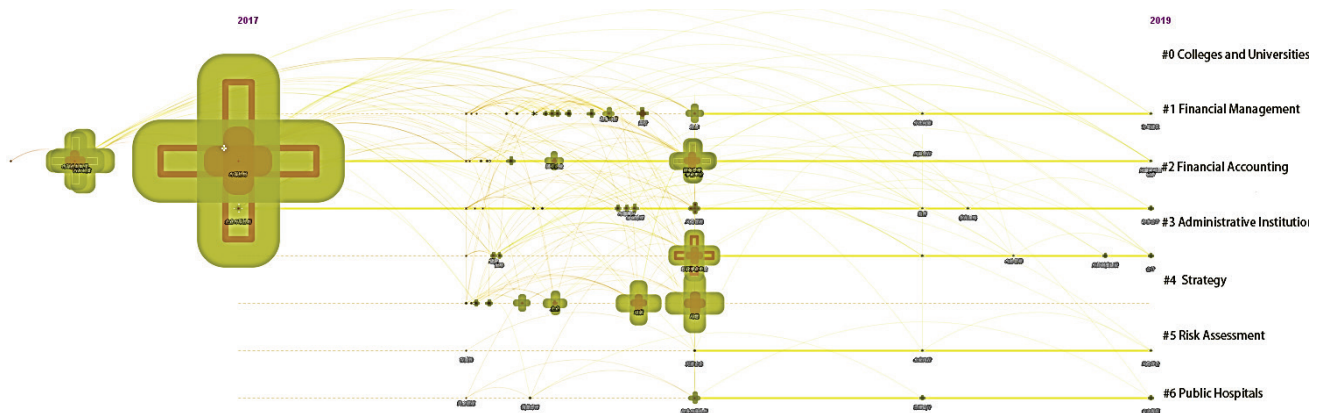


Figure 3 Clustering results

The generated keyword graphs are automatically clustered to generate a keyword clustering graph. Among them, different colors mean different keyword categories. According to the color and classification results, the

clustering results of the keywords are analyzed for feasibility, rationality and other factors, and technical processing is carried out to obtain a comprehensive analysis graph. After repeated analysis and comparison of

keyword nodes and overall structure, the clustering results can be divided into four categories. These four categories are named clockwise: risk prevention, risk assessment, risk orientation, and risk management. As can be seen from Figure 3, the clustering graph shows that the internal control group of the university is in the center, and other groups surround the surrounding pattern. The intertwined links between keywords indicate that the keywords are still interrelated.

6 CONCLUSION

Based on the data of internal control knowledge objects in colleges and universities, this paper proposes a method for constructing internal control knowledge graphs in colleges and universities. A knowledge reasoning method based on inference rules is implemented. The above technology can be applied to the internal control knowledge graph construction and retrieval system in colleges and universities. We cannot only build a knowledge graph of internal control in colleges, but also realize intelligent data retrieval and recommendation. This paper explores the construction and application of knowledge graphs in the field of internal control in colleges and universities. However, this paper mainly considers the use of internal control knowledge object data to construct the internal control information knowledge graph, and the use of unstructured text in policy documents is not sufficient. In the future work, the internal control knowledge contained in the policy text will be further explored. And it is added to the knowledge graph of internal control information of colleges and universities to enrich the knowledge graph of internal control information in colleges and universities.

7 REFERENCES

- [1] Gong, D., Liu, S., Liu, J. & Ren, L. (2019). Who benefits from online financing? A sharing economy E-tailing platform perspective. *International Journal of Production Economics*. <https://doi.org/10.1016/j.ijpe.2019.09.011>
- [2] Amit, S. (2012). *Official Google Blog: Introducing the Knowledge Graph*. <https://www.mendeley.com/catalogue/official-google-blog-introducing-knowledge-graph-things-not-strings/>
- [3] AICPA. (1949). Internal Control, A special Report by Committee on Auditing Procedure of AICPA, New York.
- [4] Control framework. <http://www.controlsframework.com/index.php>
- [5] Kinney, W. (2000). Research on opportunities in internal control quality and quality assurance. *Auditing, A journal of practice and theory*, 19, 83-90. <https://doi.org/10.2308/aud.2000.19.supplement.83>
- [6] Ashbaugh-Skaife, H., Collins, D. W., & Kinney, W. (2007). The discovery and reporting of internal control deficiencies prior to SOX-mandated audits. *Journal of Accounting and Economics*, 44. <https://doi.org/10.1016/j.jacceco.2006.10.001>
- [7] Doyle, J., Ge, W., & McVay, S. (2007). Determinants of Weaknesses in Internal Control over Financial Reporting. *Journal of Accounting and Economics*, 44. <https://doi.org/10.1016/j.jacceco.2006.10.003>
- [8] Ashbaugh-Skaife, H., Collins, D. W., & Kinney, W. R. (2007). The discovery and reporting of internal control deficiencies prior to SOX-mandated audits. *Journal of Accounting and Economics*, 44(1-2), 0-192. <https://doi.org/10.1016/j.jacceco.2006.10.001>
- [9] LaFond, R. & You, H F. (2010). The federal deposit insurance corporation improvement act, bank internal controls and financial reporting quality. *Journal of Accounting and Economics*, 49(1-2), 75-83. <https://doi.org/10.1016/j.jacceco.2009.09.007>
- [10] Hoitash, U., Hoitash, R., & Bedard, J. C. (2009). Corporate Governance and Internal Control over Financial Reporting: A Comparison of Regulatory Regimes. *The Accounting Review*, 84(3), 839-867. <https://doi.org/10.2308/accr.2009.84.3.839>
- [11] Lu, C., Laublet, P., & Stankovic, M. (2016). Travel Attractions Recommendation with Knowledge Graphs. *European Knowledge Acquisition Workshop*. Springer International Publishing. https://doi.org/10.1007/978-3-319-49004-5_27
- [12] Bizer, C., Lehmann, J., Kobilarov, G., Auer, Sören, Becker, C., Cyganiak, R., et al. (2009). Dbpedia - a crystallization point for the web of data. *Social Science Electronic Publishing*, 7(3), 154-165. <https://doi.org/10.1016/j.websem.2009.07.002>
- [13] Suchanek, F. M., Kasneci, G., & Weikum, A. G. (2008). Yago - a large ontology from Wikipedia and WordNet. *Web Semantics Science Services & Agents on the World Wide Web*, 6(3), 203-217. <https://doi.org/10.1016/j.websem.2008.06.001>
- [14] Carlson, A., Betteridge, J., Kisiel, B., Settles, B., & Mitchell, T. M. (2010). Toward an Architecture for Never-Ending Language Learning. *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2010, Atlanta, Georgia, USA, July 11-15, 2010*. AAAI Press.
- [15] Wang, Z., Wang, H., Wen, J. R., & Xiao, Y. (2015). An Inference Approach to Basic Level of Categorization. *ACM International on Conference on Information & Knowledge Management*. ACM. <https://doi.org/10.1145/2806416.2806533>
- [16] Yue, B., Gui, M., Guo, J., Yang, Z., & You, S. (2017). An Effective Framework for Question Answering over Freebase via Reconstructing Natural Sequences. *The 26th International Conference. International World Wide Web Conferences Steering Committee*. <https://doi.org/10.1145/3041021.3054240>
- [17] Bo, X., Yong, X., Liang, J., Xie, C., Liang, B., Cui, W., et al. (2017). Cn-dbpedia: a never-ending Chinese knowledge extraction system.

Contact information:

Jun WANG
 School of Economics and Management,
 Beijing Jiaotong University
 No. 3, Shangyuancun, Haidian District, Beijing, China
 E-mail: 15113125@bjtu.edu.cn