

Research Article

Open Access

Rongsheng Li* and Nasruddin Hassan

Information retrieval algorithm of industrial cluster based on vector space

<https://doi.org/10.1515/phys-2019-0007>

Received October 28, 2018; accepted January 28, 2019

Abstract: The current information retrieval research on industrial clusters has low precision, low recall ratio, obvious delay and high energy consumption. Thus, in this paper, a information retrieval algorithm based on vector space for industrial clusters is proposed. By optimizing the unlawful labels in the database network, dividing the web pages of the industrial cluster information database and calculating the keyword scores of the relevant information of the industrial cluster corresponding to a web page, a set of well-divided database pages is obtained, and the purification of the industrial cluster information database is realized. According to the purification of industrial cluster information database, RFD algorithm is used to extract the page data features of purified industrial cluster information database. The extracted results are substituted into the information retrieval, and the vectors composed of retrieval units are used to describe the information of various types of industrial clusters and each retrieval. The matching results of information retrieval are obtained by calculating the correlation between the information of industrial clusters and the query, and the information retrieval of industrial clusters is completed. Experimental results show that the algorithm has high precision and recall ratio, short retrieval time and low energy consumption.

Keywords: Vector space, industrial cluster, information, retrieval

PACS: 02.10.Ud; 02.90.+p; 07.05.Kf

1 Introduction

Industrial clusters refer to the collection of enterprises and related corporate bodies with geographical proximity, interrelated, and linked by virtue of mutual commonality and complementarity in a specific field [1]. The main components of industrial clusters include enterprises, governments, university research institutes, financial institutions, industry associations and intermediary institutions.

At present, the relevant platforms provide a variety of information retrieval services for the spatial database of industrial clusters, including fuzzy queries: approximate queries for enterprise names; classified queries: specific queries for production types in the enterprise database; peripheral queries: joint query for geographic coordinates and various types of industries in the enterprise tables; Site search: joint query for multiple fields in enterprise, product and industry information tables [2, 3]. Through the analysis and research of the above-mentioned retrieval methods for the platform, it is found that the retrieval module of this platform only supports the search of some data resources and some fields in the platform. The retrieval efficiency of the system is low, and it cannot meet the user's retrieval needs correctly. It lacks an efficient intelligent information retrieval algorithm or method.

The humanization of retrieval is reflected in the standardization of Web site design and the clarity of navigation. Information construction and retrieval is a hot issue in recent years. Information construction is proposed to better organize and present information on the web. An important goal is to make information understood. The organization and expression of information are essential for achieving good performance of information retrieval and acquirement [4]. Information retrieval, as an indispensable and important means of inquiry in people's daily life and work, plays an increasingly important role in today's society. The following are some widely used information retrieval methods and algorithms.

Zhang Xiaomin et al. put forward a keyword retrieval method based on temporal semantics. Temporal information was introduced to construct temporal data graph, and temporal correlation scoring mechanism was designed.

***Corresponding Author: Rongsheng Li:** School of Economic and Management, Northwest University, Xi'an, 710127 China, E-mail: jzshijian@126.com

Nasruddin Hassan: School of Mathematical Sciences, Faculty of Science and Technology, Universiti Kebangsaan Malaysia UKM Bangi, Selangor, Malaysia, E-mail: nas@ukm.edu.my

Temporal semantic constraints were introduced in the process of temporal graph search, and keyword-based temporal retrieval algorithm was designed. Experimental results showed that the retrieval time was short, but the precision was low [5]. Jiang Yu et al. proposed an information query algorithm based on Top-k. This algorithm extracted the static Top-k information of inverted index, and then calculated the initial threshold for specific query terms dynamically. On this basis, combining MaxScore and WAND algorithm, a fast-start Top-k query processing algorithm was proposed. Experimental results showed that the proposed algorithm had low computational complexity, but low recall ratio [6]. Zhao Yanni et al. proposed tree matching algorithm based on effective path weight. On the basis of maintaining the effective node and tree structure of XML document tree, the information of tree root node is the most important. With the increase of tree depth, the importance of node information is gradually weakened. The path weight was calculated automatically according to the path hierarchy, and the corresponding path was given. The matching degree of the tree was calculated according to the effective information of tree node and the effective path of tree structure. Experiments on large-scale XML document queries showed that the algorithm had a high query rate, but the delay of query process was obvious [7]. Ma Youzhong et al. proposed a similarity join query algorithm for high-dimensional data based on Chi square distribution. In order to solve the problem of dimensionality disaster and high computational cost in similarity join query of high-dimensional data, high-dimensional data was mapped to low-dimensional space based on p-stable distribution. The property of chi-square distribution proved that if the distance of low-dimensional space was greater than $k\epsilon$, the probability of that the distance of original space was greater than ϵ had a lower bound, so it could be filtered effectively in low-dimensional space at a lower computational cost. Experiments on real data sets showed that the proposed algorithm had a good recall rate, but it had the problem of high query energy consumption [8].

Aiming at the problems existing in the current research results, an information retrieval algorithm for industrial clusters based on vector space is proposed. The detailed process is as follows:

The improved VIPS algorithm is used to purify the information database of industrial clusters, so as to improve the precision and recall ratio of information retrieval, and reduce the retrieval delay and the energy consumption.

RFD algorithm is used to extract the page data features of purified industrial cluster information database, which lays a foundation for information retrieval of industrial clusters.

The search space is defined, to calculate the correlation between documents and queries, and realize the information retrieval of industrial clusters by the idea that the higher the correlation between information and query words is, the more relevant the information is.

The proposed algorithm is verified.

The full text is summarized and the next research plan is proposed.

2 Material and methods

2.1 Purification of industrial cluster information database

In order to improve the precision and recall ratio of information retrieval in industrial clusters, reduce the retrieval delay and energy consumption, the information database of industrial clusters needs to be processed. Noise removal module is indispensable[9]. In this paper, the improved VIPS algorithm is used to debase information blocks. Through a large number of statistics and analysis, noise semantic blocks are identified by using the number of text and links, the relative position of the page blocks and the content attributes of the page blocks. The origin coordinates of the database web page window are defined as the top left corner of the web page, the abscissa coordinate of the web page block center X is the abscissa coordinate of the center point of the web page block in the window, the ordinate coordinate of the web page block center is Y , the width of the web page is W , and the height of the web page is H . The spatial position of the web page is defined by the relative space position K of the web block, and the expression of K is:

$$K = \begin{cases} u & Y/H \leq R_1 \\ d & Y/H \geq R_2 \\ l & X/W \leq R_3 \\ r & X/W \leq R_4 \\ mid & else \end{cases} \quad (1)$$

Where u , d , l , r , and mid represent the top, bottom, left, right and middle positions of the database pages.

According to the definition of equation (1), VIPS algorithm is used to partition the web pages, and the rules are used to purify the web pages. The optimized sorting algorithm is as follows:

Input: database web page set P and the keyword set Q of industry cluster related information.

Output: a good set S_P of database pages.

The detailed process is as follows:

Optimize the illegal labels in the database network.

Use VIPS algorithm to segment web pages.

Calculate the scores of key information related to industrial clusters in a web page:

$$S_j = \frac{c_j \sum (b \times f_i \times t_{ij} \times r)}{\sqrt{\sum (b^2 \times f_i^2)} l_j} Q \cdot K \quad (2)$$

Where S_j represents the score of web page j corresponding to relevant information keywords of industrial clusters, c_j represents the number of entries containing relevant information keywords of industrial clusters, t_{ij} represents the occurrence frequency of relevant information keywords of industrial clusters in web page j , f_i represents the frequency of inverted words in web pages of relevant information keyword i of industrial clusters. b represents the field parameter, and l_j represents the length of page j .

According to the equation (2), it can get a good set of database page S_P :

$$S_P = \frac{S}{P} \cdot b \quad (3)$$

Where S represents the set of original database page. The result of equation (3) is the result of purifying industrial cluster information database.

2.2 Feature extraction of industry cluster information

According to the purification of industrial cluster information database in Section 2.1, RFD algorithm is used to extract the page data features of the purified industrial cluster information database.

Usually, if a feature item becomes a representative feature of a category, most samples of that category have this feature; if a feature item becomes a discriminant feature of a category, then most samples of other categories do not have this feature. In feature extraction, the representative and discriminant features should be selected as vector representations of a class [10, 11].

Supposing that $p(x'|c')$ can approximate the ratio of the number of information containing feature item x' in training set category c' to the total number of information containing c' in training set, then $p(x'|\bar{c}')$ can approximate the ratio of the number of information not belonging to category c' and containing feature item x' to the total number of information not belonging to category c' in training set. The similarity measure of characteristic $RFD(x', c')$ can be expressed as:

$$RFD(x', c') = S_P \left(\frac{A}{M} - \frac{B}{N-M} \right)^2 = \frac{(A \times D - B \times C)^2}{M^2(N-M)} \quad (4)$$

Where A represents the number of training information belonging to category c' and containing characteristic item x' . B represents the number of training information that does not fall into category c' and contains characteristic item x' . C represents the number of training information data that belong to category c' and does not contain characteristic item x' . D represents the number of training information that does not fall into category c' and does not contain characteristic item x' . M represents the number of information belonging to category c' . N represents the total number of training data.

Since both N and $N - M$ in the equation (4) are constants, the equation (4) can be simplified to:

$$RFD(x', c') = (A \times D - B \times C)^2 \quad (5)$$

In order to reduce the error of feature extraction, the equation (5) is improved.

$$RFD = \begin{cases} (A \times D - B \times C)^2 & A \times D - B \times C > 0 \\ 0 & A \times D - B \times C \leq 0 \end{cases} \quad (6)$$

Based on the above considerations, the calculated feature items of equation (6) have more classification discrimination ability, which is mainly to remove the information data features of industrial clusters which do not have the classification ability.

The main idea of improved feature extraction based on RFD is that for a feature to become a representative feature of a certain category, it must have the following two characteristics: representative and discriminant [12]. The absolute value of the sum of the representativeness measure of feature item x' and the discriminability measure of feature item x' are used to measure the correlation between features and categories, which is called ARFD. Conditional probability $p(x'|c')$ is a representative measure of characteristic item, while $-p(x'|\bar{c}')$ is a discriminant measure of characteristic item x' . The improved feature extraction is to calculate by using equation (7):

$$ARFD(x', c') = \left[p(x'|c') - p(x'|\bar{c}') \right] S_P \quad (7)$$

Equation (7) can be approximated to:

$$ARFD(x', c') = \left| \frac{A}{M} - \frac{B}{N-M} \right| S_P \quad (8)$$

Where the greater the value of $ARFD(x', c')$ is, the more the relevant information of feature item x' and class c' is.

2.3 Information retrieval algorithm based on vector space for industrial clusters

In order to improve the precision and recall ratio of information retrieval in industrial clusters, information retrieval is realized on the basis of information feature extraction of industrial clusters. The retrieval pattern of vector space is a relatively easy to understand retrieval pattern, and is a widely used information retrieval algorithm model in the field of information retrieval[13, 14]. The basic idea is that information and query are made up of words, and each query can be described by a vector composed of retrieval units. When searching, the correlation between information and query is calculated, and the higher the correlation with a specific query is considered the more relevant information.

The common way to describe information and retrieval vectors is that the retrieval space is composed of all retrieval units contained in information and retrieval, and the information and retrieval are represented as vectors in this space.

It is assumed that the information retrieval space of industrial clusters is $\Omega = \langle t'_1, t'_2, \dots, t'_{n'} \rangle$. Among them, $t'_i = (i' = 1, 2, \dots, n')$ is the different retrieval units contained in information and query, n' is the size of the whole retrieval space Ω , that is, the total number of different retrieval units contained in information query.

In retrieval space Ω , all information can be represented by vectors: $d' = \langle \omega_{d'1}, \omega_{d'2}, \dots, \omega_{d'n'} \rangle$. Among them, $\omega_{d'n'}(i' = 1, 2, \dots, n')$ is a series of descriptions of the information meaning, when the retrieval unit t'_i appears in the information type, $\omega_{d'n'}$ is 1, conversely, when the retrieval unit t'_i does not appear in the information, $\omega_{d'n'}$ is 0. Usually, most of the items in $\omega_{d'n'}$ are zero because the size of search space is much larger than the length of each industry information file.

Combining the above information, we can approximately understand that in search space, all queries can also be represented by vectors: $q = \langle \omega_{q1}, \omega_{q2}, \dots, \omega_{qn'} \rangle$. Among them, $\omega_{qn'}(i' = 1, 2, \dots, n')$ is a series of descriptions of the query meaning, when retrieval unit t'_i appears in the query, $\omega_{qn'}$ is 1, conversely, when retrieval unit t'_i does not appear in the query, $\omega_{qn'}$ is 0. In general, because the length of queries is shorter than that of industrial clusters, more entries will be zero in $\omega_{qn'}$.

According to the above analysis, not every retrieval unit is equally important in information retrieval of industrial clusters (for example, keywords should be more important than non-keywords). So, how to embody such information in vectors needs to be solved urgently [15–21]. One of the feasible schemes is to adjust the weight of vec-

tors manually, which enlarges the weight of retrieval units that users care about. However, manual intervention is difficult to achieve because of the huge workload. Therefore, another method is more commonly used in information retrieval: the weights based on the statistical frequency of the information file set, also known as TF-IDF weights.

TF-IDF weights consist of two parts, one is the frequency of the retrieval unit appearing in the information file, that is, TF, the other is called inverted file frequency, that is, IDF. TF-IDF weight is usually the product of TF and IDF for a given retrieval unit.

For the convenience of illustrating the problem, the following definition is made: $TF_{ij'}$ represents the frequency of the retrieval unit t'_i appearing in the industrial cluster information database, DF_j represents the amount of information containing the retrieval unit t'_i in the entire industrial cluster information database.

By defining the above definition, the frequency of inversion information can be defined as:

$$IDF_j = \log \left(\frac{d'}{DF_j} \right) \quad (9)$$

Where, IDF_j represents the frequency of reversal information.

For a given information file, the vector describing the information file is composed of n' elements, which correspond to n' retrieval units in the information file set. The weights of each element are determined by the frequency of the corresponding retrieval unit appearing in the industrial cluster information database and the frequency of the retrieval unit appearing in the entire industrial cluster information database, as shown in Eq. (10):

$$\omega_{ij'} = TF_{ij'} \times IDF_j \quad (10)$$

Using as the weight of each element in the vector, the vector of information and retrieval is further adjusted. When the value range is [0.25, 0.30], the vector of information and retrieval can be adjusted best. This vector can describe the information and query more accurately.

For vector space retrieval model, it not only needs to define vectors to represent information and retrieval, but also needs to choose an appropriate method to calculate the relevance of information and query to determine whether information and query are related. The cosine of vector angle is used as the basis for judging the relevance of industrial cluster information.

According to the above, the similarity between information d' and retrieval q is defined in retrieval space Ω . The retrieval matching process can be expressed as fol-

lows:

$$SC(d', q) = \frac{\sum_{i'=1}^{n'} \omega_{d'_{i'}} \times \omega_{q_{i'}}}{\left[\sum_{i'=1}^{n'} \omega_{d'_{i'}}^2 \sum_{i'=1}^{n'} \omega_{q_{i'}}^2 \right]^{1/2}} \cdot ARFD \quad (11)$$

Where $SC(d', q)$ calculated by equation (11) is the result of information retrieval based on vector space.

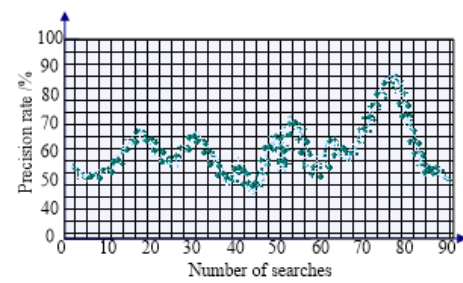
3 Results

In order to verify the validity of the vector space based information retrieval algorithm for industrial clusters, a correlation experiment is conducted. In the experiment, two industrial gathering information of education and entertainment in a province are selected as the source of experimental data. Experimental environment: Intel Pentium Dual E2140@1.60GHz; operating system: Microsoft Windows XP; hard disk: 160 GB; memory: 1 GB; development tools: Eclipse 3.2. The experimental indicators are: Retrieval precision; Retrieval recall ratio; Retrieval delay; Network energy consumption of retrieval. The results are as follows:

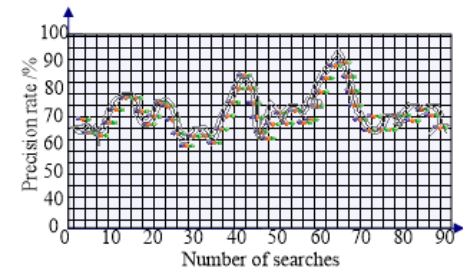
Figures 1 and 2 show that the information retrieval algorithm based on vector space for industrial clusters has higher precision and recall ratio, and is more robust than the current research results. RFD algorithm is used to extract the page data features of purified industrial cluster information database, and preliminarily determine the characteristics of industrial cluster information data, which provides support for information retrieval. Based on feature extraction, the retrieval space is set, and the cosine of vector angle is used to judge the correlation between information and retrieval in industrial clusters. The results of information retrieval in industrial clusters are obtained, which effectively improves the precision and recall ratio of information retrieval.

As can be seen from Figures 3 and 4, compared with the current research, the information retrieval algorithm based on vector space has a great advantage in terms of retrieval delay and energy consumption. Before searching for industrial clusters, this algorithm uses improved VIPS algorithm to purify the information database of industrial clusters. The noisy semantic blocks are identified and removed by the number of text and links, the relative position of web blocks and the content attributes of web blocks, thus greatly reducing the information retrieval delay and reduce retrieval energy consumption.

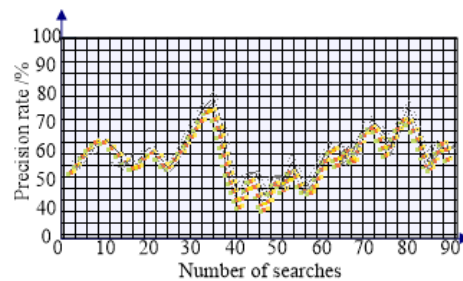
Experiment 1:



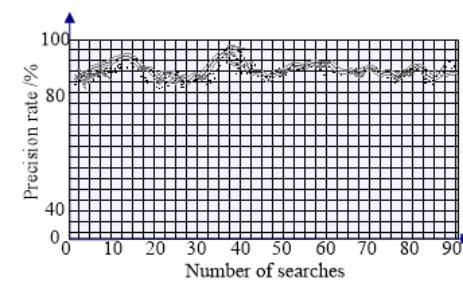
(a) Precision of keyword retrieval based on temporal semantics



(b) Precision of information query based on Top-k



(c) Precision of tree matching algorithm based on effective path weight



(d) Precision of information retrieval in industrial clusters based on vector space

Figure 1: Comparison of precision of different information retrieval methods

Experiment 2:

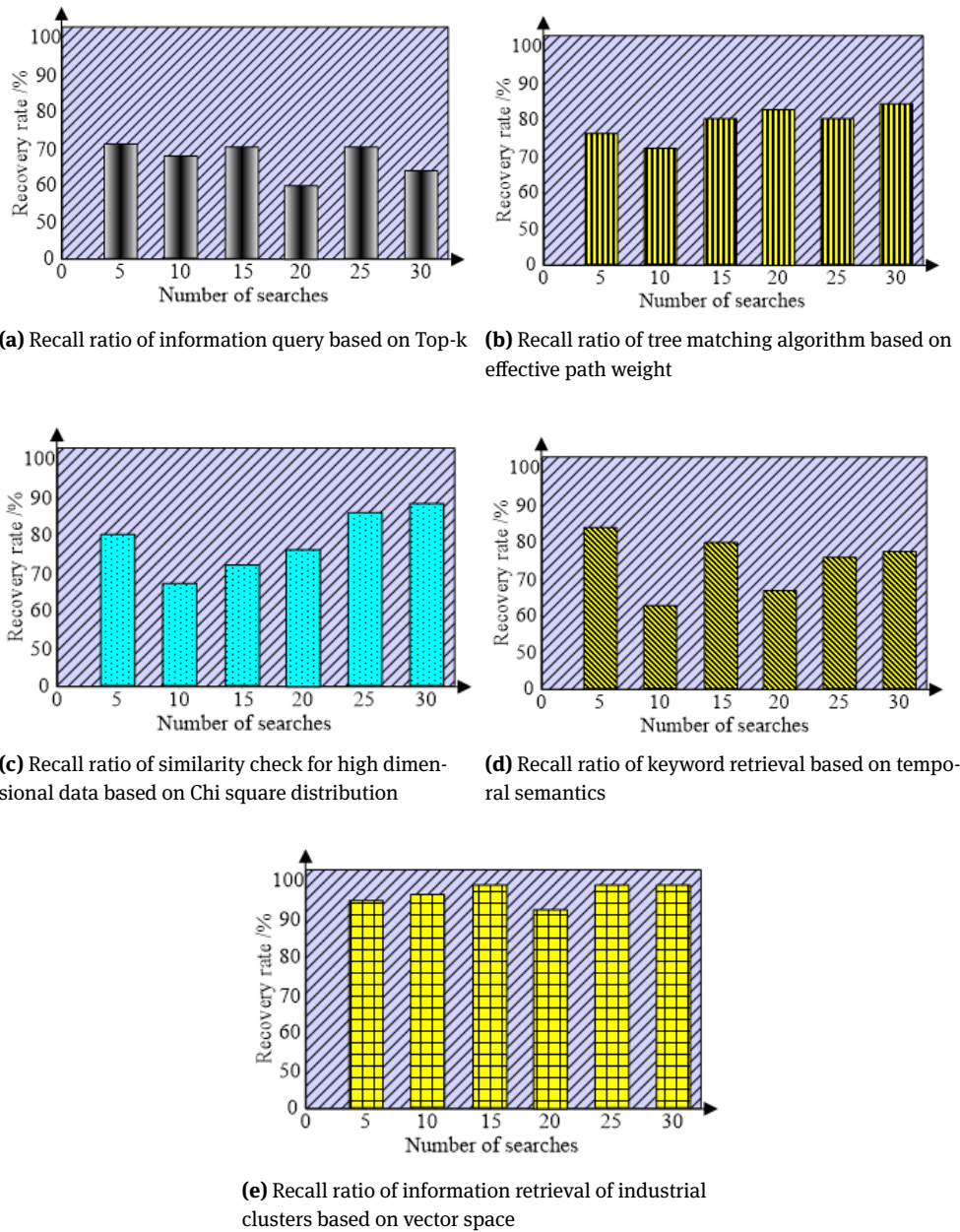


Figure 2: Comparison of recall ratio with different information retrieval methods

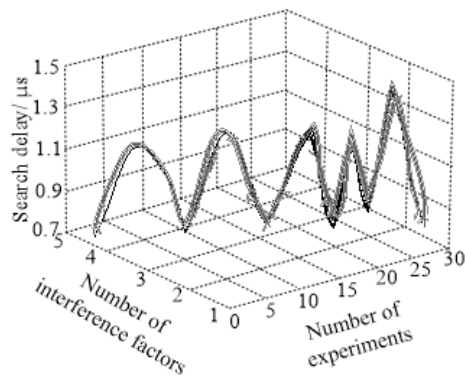
4 Discussion

This paper discusses the effect of adjusting the weight ω_{ij} of each element in the information data vectors of industrial clusters to the information and retrieval vectors. The value of ω_{ij} is defined in $[0.19, 0.24]$, $[0.25, 0.30]$ and $[0.31, 0.36]$ respectively, and the effect of weight ω_{ij} on information and retrieval vectors is observed. The larger the adjustment coefficient is, the more accurate the vector

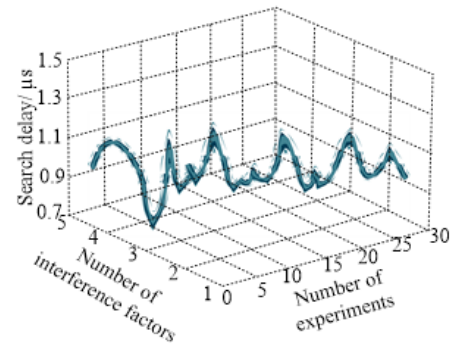
can describe the information and the content of the query. The simulation results are as follows:

In Figure 5, when the value of ω_{ij} is $[0.19, 0.24]$ and the vector adjustment coefficients of information and retrieval fluctuate greatly, which indicates that the accuracy of vector description information and query content will also be affected to varying degrees, and then the effect of information retrieval in industrial clusters will be affected. When the value of ω_{ij} is $[0.25, 0.30]$, the vector adjust-

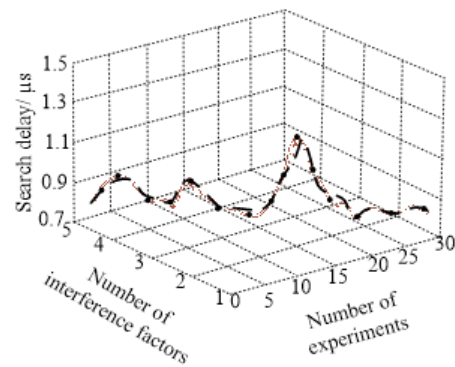
Experiment 3:



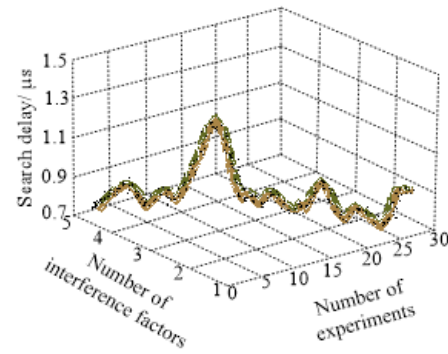
(a) Query delay of tree matching algorithm based on effective path weight



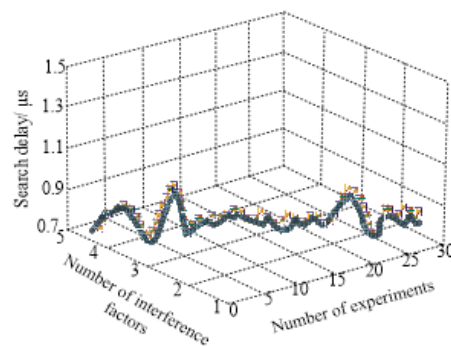
(b) Query delay of keyword retrieval based on temporal semantics



(c) Information query delay based on Top-k



(d) Similarity query delay for high dimensional data based on Chi square distribution



(e) Query delay of industrial clusters information retrieval based on vector space

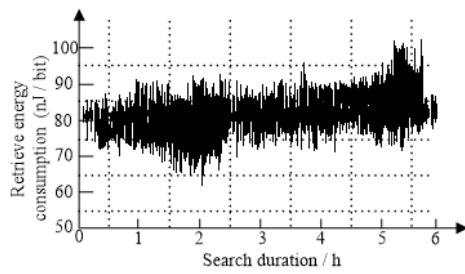
Figure 3: Comparison of retrieval time for different information retrieval methods

ment coefficient of information and retrieval is the largest, which means that the vector can describe information and query content to the greatest extent.

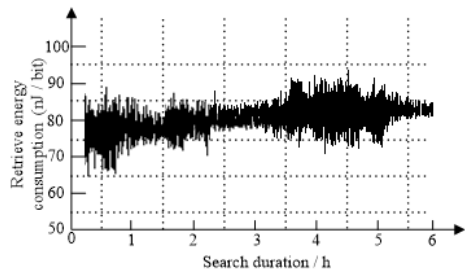
5 Conclusions

As a hot social content, industrial clusters play a positive role in social development. Information retrieval of industrial clusters is conducive to understanding the develop-

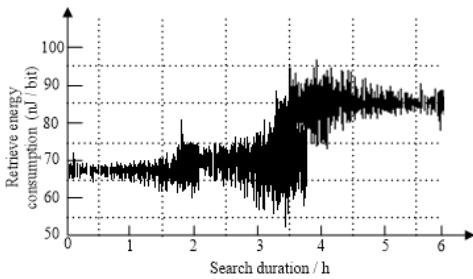
Experiment 4:



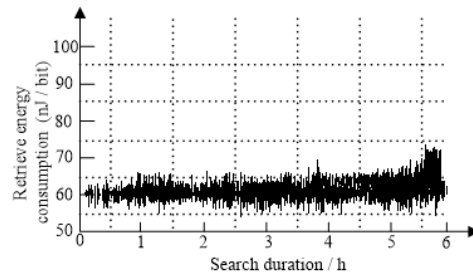
(a) Energy consumption of high-dimensional data similarity based on Chi square distribution



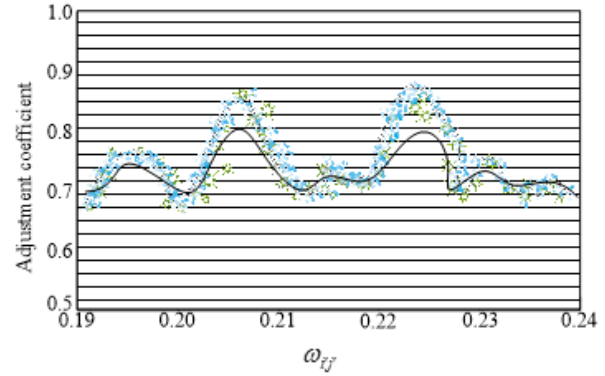
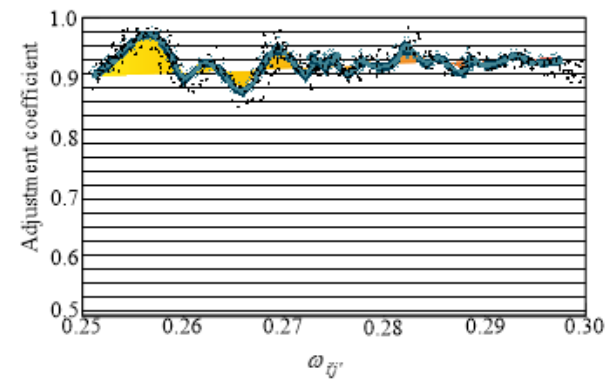
(b) Energy consumption of information query based on Top-k



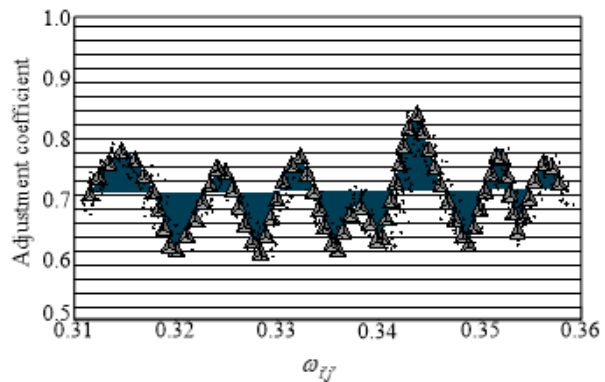
(c) Energy consumption of keyword retrieval based on temporal semantics



(d) Energy consumption of industrial cluster information retrieval based on vector space

Figure 4: Comparison of energy consumption in different information retrieval methods(a) Vector adjustment coefficients of information and retrieval when the value of ω_{ij} is [0.19, 0.24]

(b) Vector adjustment coefficients of information and retrieval when the value of is [0.25, 0.30]

(c) Vector adjustment coefficients of information and retrieval when the value of ω_{ij} is [0.31, 0.36]**Figure 5:** Influence of different values of ω_{ij} on vector adjustment coefficients of information and retrieval

ment status of industrial clusters, which is of great significance to the regulation and progress in this field. Therefore, an information retrieval algorithm of industrial clusters based on vector space is proposed. The information

retrieval of industrial clusters is completed by purifying the information database of industrial clusters, extracting the information features of industrial clusters and matching the information similarity of industrial clusters. Experimental results show that the proposed algorithm has a high retrieval rate and retrieval efficiency, and has absolute advantages over the current research results.

Acknowledgement: Major project of applied research of philosophy and social science in Henan higher schools “research on development strategy of higher career education in the context of higher education globalization” (2016-yyzd-21);

Research project of decision-making of Henan provincial government “study on occupational education promoting industrial upgrading mechanism of Henan province” (2016B090).

References

- [1] Wu Z., Gao K., Wang Z., Wei C., Wali F., Zan G. et al., Direct information retrieval after 3D reconstruction in grating-based X-ray phase-contrast computed tomography, *J. Synchr. Radiat.*, 2018, 25(Pt 4), 1222-1228.
- [2] Li G.L., Chen J.L., Liu B., Yin Y., Zhang H.B., Cross-media retrieval of online product based on tag-rank and CCA, *Sci. Technol. Eng.*, 2016, 16(14), 222-227.
- [3] Lu N., Gao Q.M., An algorithm for retrieving internet tourism resources based on mixed feature threshold, *B Sci. Technol.*, 2017, 33(8), 162-165.
- [4] Hu J.H., Qin Z.C., Shi L., Lu Z., Zhou B., Research on spatial information cloud service platform and application, *J. China Acad. Electron. Inf. Tech.*, 2016, 11(1), 51-58.
- [5] Zhang X.M., Qi W., Zhang J., Gui X.Q., T-STAR: Keywords-based temporal information retrieval method over relational databases, *Appl. Res. Comput.*, 2017, 34(10), 3051-3056.
- [6] Jiang Y., Song X.S., Yang Y.X., Jiang K., Rapid start top-k query based on threshold, *J. Chinese Inf. Process.*, 2017, 31(5), 163-170.
- [7] Zhao Y.N., Guo H.L., XML tree matching algorithm based on effective path weight, *Comput. Eng. Des.*, 2016, 37(4), 949-953.
- [8] Ma Y.Z., Jia S.J., Zhang Y.X., Chi-square distribution based similarity join query algorithm on high-dimensional data, *J. Comp. Appl.*, 2016, 36(7), 1993-1997.
- [9] Li Y.X., The simulation research on the optimization management of mass library information retrieval, *Comput. Simulat.*, 2017, 34(5), 389-392.
- [10] Ren Y., Design and implementation of a fault searching system combined with semantic web, *Comput. Meas. Control*, 2017, 25(5), 35-37.
- [11] Wei Y.P., Banawan K., Ulukus S., Cache-aided private information retrieval with partially known uncoded prefetching: fundamental limits, *IEEE J. Sel. Area Comm.*, 2017, (99), 1-1.
- [12] Huang X.X., He Y., Application of cloud computing technology in library group resource retrieval, *Automat. Instrum.*, 2017, (2), 139-142.
- [13] Rocha V., Kon F., Cobe R., Wassermann R., A hybrid cloud-P2P architecture for multimedia information retrieval on VoD services, *Comput.*, 2016, 98(1-2), 73-92.
- [14] Jiang Y., Zhang J., Zhu L.X., Ontology based knowledge graph model of genealogical record and retrieval system, *Electron. Des. Eng.*, 2017, 25(12), 161-165.
- [15] Metwally O.N., Sinha S.R., Sa2026 a novel voice-activated web application for rapid knowledge generation and information retrieval through semantic parsing of verbal communication, *Gastroenterol.*, 2016, 150(4), S433-S433.
- [16] Zandebasiri M., Soosani J., Pourhashemi M., Evaluating existing strategies in environmental crisis of Zagros Forests of Iran, *Appl. Ecol. Env. Res.*, 2017, 15(3), 621-632.
- [17] Jezewska-Frackowiak J., Seroczynska K., Banaszczyk J., Wozniak D., Skowron M., Ozog A. et al., Detection of endospore producing bacillus species from commercial probiotics and their preliminary microbiological characterization, *J. Environ. Biol.*, 2017, 38(6), 1435-1440.
- [18] Delgado J., Peña J.M., Monotonicity preserving representations of curves and surfaces, *Appl. Math. Nonlin. Sci.*, 2016, 1(2), 517-528.
- [19] Li D., Wang L., Peng W., Ge S., Li N., Furuta Y., Chemical structure of hemicellulosic polymers isolated from bamboo biocomposite during mold pressing, *Polym. Compos.*, 2017, 38(9), 2009-2015.
- [20] Brown T., Du S., Eruslu H., Sayas F.J., Analysis of models for viscoelastic wave propagation, *Appl. Math. Nonlin. Sci.*, 2018, 3, 55-96.
- [21] Gao W., Zhu L., Guo Y., Wang K., Ontology learning algorithm for similarity measuring and ontology mapping using linear programming, *J. Intell. Fuzzy Syst.*, 2017, 33(5), 3153-3163.