# A Review of Term Semantic Hierarchy Induction for Domain-specific Chinese Text Information Processing

Yang Yuqing

School of Software and Microelectronics

Peking University

No. 24, Jinyuan Road, Daxing District

Beijing, 102600, China

blackmud@163.com

ABSTRACT. *This paper introduces researches in the field of term semantic hierarchy extraction for domain-specific Chinese text information processing. Methods based on linguistic templates, CRFs, FCA and clustering are introduced. This paper puts emphasis on hierarchal clustering algorithms for Chinese text, and analyzes problems existed in clustering experiments.*
**Keywords:** semantic hierarchy; hyponymy; Chinese language

1. **Introduction.** In the past decades, people strived to build knowledge systems by directly editing or extracting semantic relationship in an automatic or semi-automatic way. The first one is time-consuming or a lot of people need to participate in the process to finish the massive work. The second one is to supplement the system built in the first way. It usually needs different algorithms to extract different relationships from raw texts. This paper will introduce a research direction dealing with Chinese texts that puts emphasis on automatically extracting the semantic hierarchy of domain-specific terms from texts, make a review of existing methods by comparing their features, and put an emphasis on the research status and main problems of clustering-based term semantic hierarchy induction for domain-specific Chinese text information processing.

2. **Definition.**
2.1. **Definition of Domain-specific Term.** Domain-specific Terms are words, compound

words or multi-word expressions used in a specific context. Some terms refer to different meanings in different contexts. Most terms are rarely used in daily life. [1]

2.2. **Definition of Semantic Hierarchy.** Term semantic hierarchy extraction is one of the basic parts of natural language processing and ontology building, and it has big influence on the quality of ontology [2-3]. If concept B is an object of concept A, then concept A is a *hypernym* of concept B and B is a *hyponym* of A. In the domain of computer technology, such hierarchy is often called ISA relationship, marked as ISA(B, A).
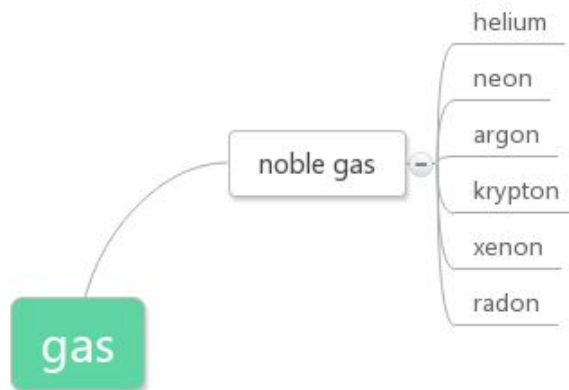


FIGURE 1. AN SIMPLE EXAMPLE OF SEMANTIC HIERARCHY

Term semantic hierarchy extraction is to automatically extract semantic relationship, including words similarity and ISA relationship from texts through different algorithms. Term semantic hierarchy extraction is a crucial step to build knowledge systems. [4-6] Methods to extract term semantic hierarchy can be categorized into templates-based methods, CRF-based methods, FCA-based methods and clustering-based methods.

3. **Main Categories of Automatic Hierarchy Acquisition Methods.**
3.1. **Templates-based Methods.** Template-based methods rely on language researchers extracting frequent rules from texts. Algorithms use human-find rules to discover terms and their hierarchies. Template-based methods are widely used.

Hearst[7] proposed six English lexical-syntactic feature, including:
- *NP such as    {NP,}\*{(or | and)} NP,*
- *Such NP as {NP,} \*{(or | and)} NP,*
- *NP {,} including {NP,}\*{(or |and)} NP,*
- *NP {, NP}\*{,} and other NP,*
- *NP {, NP}\*{,} and other NP,*
- *NP {,} especially {NP,}\*{(or | and)} NP.*

These rules can identify the ISA relationship in a sentence by detecting syntactic structure used for presenting examples. Many researchers improved this method. Iwanska[8] complemented the Hearst model, adding rules such as not only X but also Y to extract hyponyms. These algorithms are based on lexical-syntactic rules. They are easy to comprehend and can often reach high accuracy rate, so they are widely used in automatic

acquisition of term hyponymy.

In Chinese domain, Liu Lei[9] proposed a method of hyponym acquisition based on "is a" ("是一个") pattern, which divides the sentences by judging whether the sentence has commas or the character "的" and use different rules to deal with sentences with different types. In our project, researchers also use template-based method to discover term hyponymy in the terminology dictionary of electronics and chemical industries, such as extracting hypernym by the rule of "*belong to*" ("属于……的一种").

Although these template-based methods have high accuracy rate, but they depend more on the effectiveness of rules constructed by researchers. Chinese language put more emphasis on expressing the true meaning of texts no matter in what structure, so it can be exhaustive to find all the rules; in addition, there are many implicit hyponymy relations in Chinese sentences, so texts with explicit hyponymy structures tend to be sparse. Sparseness however can be made up by the big data on the Internet. For example, the Probase Concept Library of the Microsoft Corp [10], which is based partially on the Hearst's theory, extracts the word pairs from the one billion and six hundred million Web pages and finds ISA relationship according to the term probability.

In Chinese areas, there are also studies extracting ISA relationships through structured online Encyclopedia pages. Song Wenjie et al.[11] extract hyponym based on dictionaries and online Encyclopedia. They construct regular expressions, match the to-be-extracted hyponyms in web pages, and put the target word and their lower word collection into Hash tables. The concept mining method based on large corpus in online Encyclopedia can be effective when dealing with everyday words, but the domain terms are sparse in the common text, so it is difficult to find their hyponymy through this method.

**3.2. CRF-based Methods.** To build a conditional random field model (CRF) [12] need to solve two problems: one is the estimation of the parameters; the other is the feature selection. It needs to learn the training data set and then obtain the weight parameters of each feature. Feature selection is to pre-process the text and select features which are useful in the establishment of CRF model.

Let X be the random variable of the observation sequence, and let Y be the random variable of the labeled sequence. $X = (X_1, X_2, ... X_n)$, $Y = (Y_1, Y_2, ... Y_n)$, where $Y_i$ belongs to T, (i=1, 2,..., n). T is a limited set of tags. For example, X can be the words or parts of speech in the sentence; Y can be the corresponding sentence syntax elements, such as noun phrases, verb phrases, etc. The random variable Y and X constitute a joint distribution, and a conditional probability model P (Y|X) is constructed according to the observation and labeled sequences. The model definition is as follows:

$$P(y|x) = \frac{1}{Z(x)} \exp \sum_i \sum_k \lambda_k f_k(y_{i-1}, y_i, x, i) + \sum_i \sum_k \mu_k g_k(y_i, x, i)$$

(1)

$Z(x)$ is the normalization factor; $f_k(y_{i-1}, y_i, x, i)$ is the state transfer function; $G_K (Y_i, x, )$ is the state characteristic function.

Researchers has already used CRFs to identify the hyponymy of words from the general

field, such as Deschacht[13] use Wordnet to label synonyms, hypernyms and hyponyms, then construct CRF model and train it. But most words in Wordnet are from everyday life, so it has little use in labeling domain-specific terms and extract hyponymy.

Chinese researchers also use CRF to find the relationship between words. Huang Yi [14] finds the context information of domain terms from online Encyclopedia, summarize and extract relevant pattern. Characteristics include words and part of speech; characteristic dictionary and punctuation information is also added as characteristics. It trains and learns the model on the basis of CRF model, and eventually establishes the classification model. But the experiment only considers words and part of speech features, but the sentence usually contains more than one noun or noun phrase, so it will influence the accuracy of the results.

In addition, the ICTCLAS system is used to find the domain terms, but the result is not accurate enough，so the researchers assumed that CRF model can be improved into double layers while the low level model is used to identify the domain terms. In addition, the results are calculated by the accuracy rate and recall rate, which is improper as the semantic level is complex and the semantic distance is often subjective. The article did not describe whether experts' opinions are considered, so the result can only be considered as a reference.

Mo Yuanyuan [15] improved Huang Yi's method with CRFs of double layers to extracting the semantic relationship. Low layer CRF model is constructed by words, with each word as a unit, taking long-distance dependencies between words into consideration. When extracting the domain-specific concept, they combine the word according to the template, and put the training texts into training files, each line consisting of a pair of domain-specific concept, then use the high layer CRF model label the domain-specific concept pairs.

Algorithms based on conditional random fields are usually used for texts parsing and named entity recognition, and it is respectively uncommon to use it in extracting the hyponymy. This kind of algorithm has good performance in term extraction, but it is time-consuming and is a heavy task to label the training texts.

3.3. **FCA Algorithm.** Harris [16] proposed distributional hypothesis theory, which means that if two words are in similar contexts, then these two words is approximate. FCA algorithm is based on Harris hypothesis.

FCA algorithm is derived from the mathematical theory, using the contexts of the concept to show its properties and generating the concept hierarchy by the concept lattice construction algorithm. OBITKO[17] gets a syntax tree through the syntax parser, gets the subject-verb relationship and the verb-object relationship, and then change the verb and noun into the original form to show the concept property. But Chinese does not have a general syntax parser, so some researchers use qualifiers as properties.

Wen Chun [18] compares the method based on VSM hierarchical clustering and FCA to extract Chinese terminology hyponymy. They found VSM-based hierarchical clustering method works better, but it still need to label the class by HowNet, a Chinese knowledge system like WordNet. But words in HowNet are mostly from everyday life, so the result is

influenced by the fact, while the performance of FCA is limited on account of Chinese characteristics, so its effect is poorer.

3.4. **Algorithms Based on Hierarchal Clustering.** Hierarchal Clustering is based on Harris hypothesis like FCA method. Hierarchal clustering algorithm first calculates term similarity, then put similar terms into same cluster. [19-23]

Hierarchical classification algorithm is divided into two kinds. One is the bottom-up hierarchical clustering algorithm, also known as the Agglomerative NESting. Top-down hierarchical clustering, also known as DIANA clustering, is divisive from the top clustering. The following figure shows the different processing process of the two levels of classification.
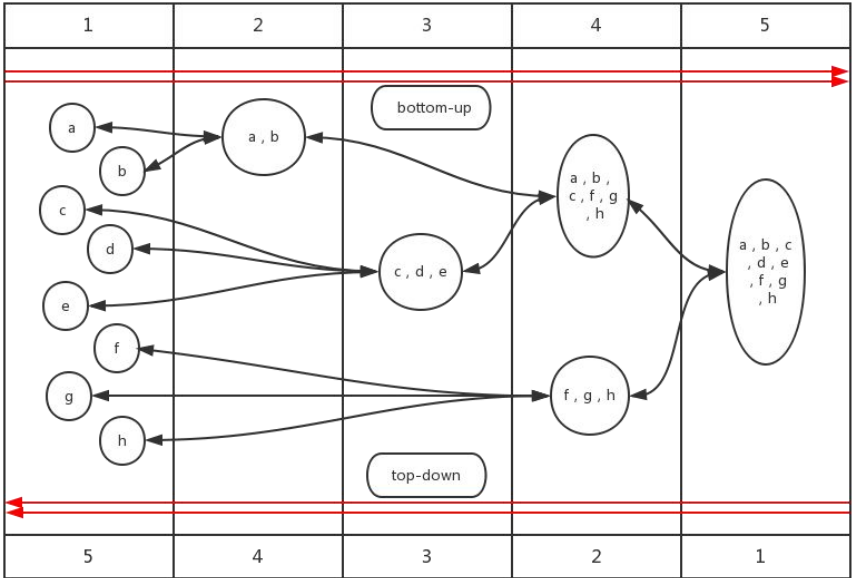


FIGURE 2. TWO DIRECTIONS OF HIERARCHAL CLUSTERING METHODS

Chen Yuanhao [24] proposed WAC algorithm that is based on bottom-up hierarchal clustering to deal with common words. WAC improves the traditional hierarchical clustering method, cluster the words from the bottom level, and select a representative point set to form upper set of nodes. The similarity calculation between two nodes uses the feature of spectral clustering, which allows a node own multiple father nodes. The algorithm aims to generate undirected graphs that have different number of nodes and the relationship between adjacent graphs. These undirected graphs constitute a pyramid, and the top of the pyramid is an undirected graph with only one node. However, if the input that the undirected graph is not precise enough, which means that the words similarity calculation [25] before clustering calculation is not precise enough, then the clustering precision cannot be improved. It is necessary to better describe the similarity between words.

Ji Peipei [26] proposed a top-down hierarchical clustering method, which uses the most

widely used VSM (vector space model) to calculate the similarity of words, and then uses the improved K-means method to cluster. When re-calculate the cluster center, according to the different distribution of each class, it would select $n$ terms, and then those the one which is closest to the average vector as the new central word, so as to speed up the convergence process; after each layer of clustering, calculate the comprehensive similarity of all the terms. If the term has higher comprehensive similarity value, it may have a wider semantic range in this layer, so it will be chosen as the class tag and the other terms will be into the next clustering process.

The specific steps of the algorithm are: determine the initial center words manually; calculate the similarity between each term and each center word, and label the term with the closest center; re-calculate the center, calculate the average similarity of each type, choose $n$ terms which is nearest to the center term and calculate their average vector, and choose the term which is closest to the average vector as the new center; calculate the Squared-Error, and if the center list is not stable, then recalculate the center; if the center list is stable, then the algorithms is over.

The drawback of the approach is that the hypernym is always improper. In Ji Peipei's experiment, *Aspirin* became the hypernym of *chemical drugs*, *patent medicine*, *medicine*, *prescription drug*, but aspirin is a high-frequency word, so it has a higher score in the integrated similarity calculation and becomes the wrong hypernym. We selected one hundred FINANCE terms in Sougou FINANCE news corpus to extract the hyponymy by this algorithm. But as word similarity calculation is not precise enough, clustering effect is not good enough, and it is difficult to obtain a reasonable hypernym.

Peng Cheng, Ji Peipei [27] and other researchers later proposed hierarchical clustering algorithm based on the *deterministic annealing*. The algorithm processes is similar to literature [26]'s hierarchical clustering process, but when doing the hierarchical clustering annealing algorithm is used. It uses the principle of maximum entropy in information theory. By keeping the entropy within a certain range, algorithm has a chance to jump out of local optimum, and algorithms can get a better result with less cost. Deterministic annealing is different with K-means in that different initial words have little influence on it. In addition, the algorithm adds the control to the tree depth. If the current class of entry is less than the depth of the entire tree, then clustering downward is stopped. Thereby making the results more reasonable as it avoids tree branches has same depth. By comparing he results generated by the algorithm and experts in the field, they found basic semantic structure. And as the corpus increases, the similarity increases though in a slow speed. The slow speed may due terms are sparse in the raw texts.

Gu Jun, Zhu Ziyang [28] proposed hierarchical clustering method. It also calculates similarity between the terms by the vector space model. The first clustering uses the ant colony algorithm, and get some initial classes by its robustness and excellent performance in distributed computing, then finish the multiple clustering by K-means. Ant colony clustering algorithm [29] to simulate the process that ants pick up the body of its companions and forms the corpse stack, that is to say the entire array is divided into grids, while ants move from the starting point, pick up objects, move the object, and places the object in

reasonable place to form clusters. The ants' act to pick up objects, move objects and drop the object is judge based on the average similarity of the object and the surroundings. In Equation (2), d ($t_i$, $t_j$) is the similarity between the terms calculated in advance. If the average similarity between the objects and the surrounding is high, then the probability that ants drop the object is large; If the average similarity between the objects and the surrounding is low, then the probability that ants pick up the object is large; so that after a certain time, similar objects will gather together and achieve the effect of clustering. The algorithm uses the advantage of ant colony algorithm that it does not need to determine initial numbers of clusters. However, due to the ant colony algorithm requires a certain number of iterations, so it is rather time-consuming.

$$f(t_i) = \max\{0, \frac{1}{s^2} * \sum_{t_j \in Neights* s(r)} [1 - \frac{d(t_i, t_j)}{\alpha(1 + (v-1)/v_{max})}]\}$$

(2)

4. **Main Problems in Current Researches for Hierarchal Clustering.** The processes of top-down hierarchical clustering algorithm in most Chinese researches to build the semantic hyponymy tree are similar. They can be summarized in the following flow chart:
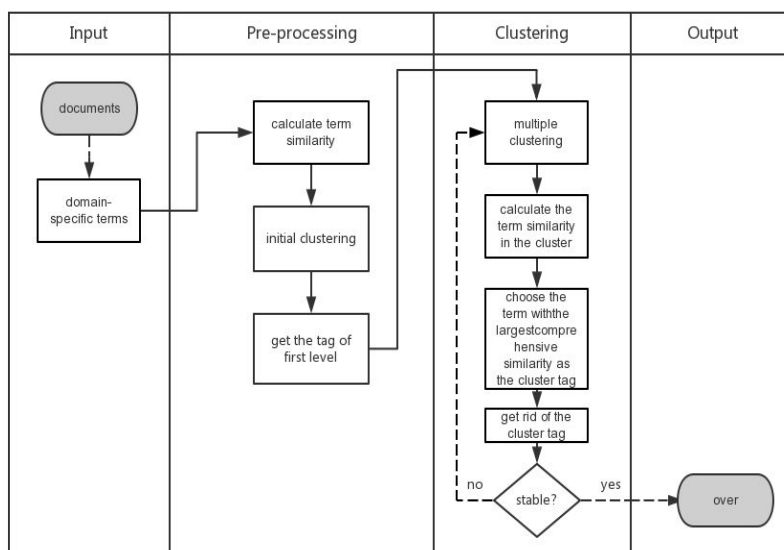


FIGURE 3. THE PROCESS OF TOP-DOWN HIERARCHAL CLUSTERING

As can be seen in the Figure, the third step multiple clustering is calculated on the basis of term extraction and term similarity; and as semantic relationship is subjective and terminology is strongly professional, so there are no uniform evaluation criteria.

4.1. **Term Extraction.** Term extraction from raw texts is always a hot research spot. First of all, there is no space between Chinese words, so term extraction result will be based on the performance of parsing algorithms; second is extraction methods have their own problem. But in recent years, machine-learning methods has got some progress in Chinese

56

terminology extraction, particularly the CRFs algorithm. Liu Bao and his partners [30] combine CRF and the method based on templates to extract scientific terms automatically, and the F value reaches 84.4% in open test; Liu Lei [31] extract Chinese and English terms from patent by means CRFs, and the F value of Chinese term extraction reaches 88.43%. Conditional random fields describe the model by defining the conditional probability P (Y | X), rather than the joint probability distribution P (X, Y). CRF model does not have independence assumptions as in hidden Markov model, so more text features can be added; and CRF model calculates the probability that is a global optimum rather than a local optimum, so it solves the problem of maximum entropy model, but training texts annotation requires a lot of labor and template selection is also difficult. [1]

4.2. **Term Similarity.** From the foregoing flow chart summarized we can see that the basic part of clustering is to calculate the semantic similarity between words, such as the k-means algorithm first chooses *n* objects as the initial cluster centers, then compute the similarity distance between the cluster center and other objects respectively, and assign them to the most similar cluster; In ant colony algorithm, ant decide whether to pick up or drop the object upon the average similarity between the object and the things around it, while the average similarity formula need to know the similarity among objects beforehand. Consequently, the reasonable similarity calculation is the basis of clustering.

Most Chinese researchers calculate similarity of words by cosine value, namely establish VSM model by using TF-IDF weighting and calculating the cosine similarity. Wen Chun and his partners [18] compared the document vector model established by the method of the window and TF-IDF method. The method based on the window put *n* words around the key word within the window as its properties, and build word vector space according to the frequency of properties. Another method establishes a concept-document model, using TF-IDF to calculate weights. No matter how big the corpus size is, the latter method, the traditional VSM model weighted by TF-IDF always has better performance.

4.3. **Judgment.** There is no uniform standard to calculate the accuracy of the semantic tree. Judgment method is divided into two categories. One is to judge manually the accuracy of clustering and the reasonability of the hypernym; the other is to calculate the similarity of trees built by algorithms and the tree built by experts with certain function.

Chen Yuanhao and his partners [24] compare the similarity between tags pairs and the one marked manually in order to evaluate the semantic tree. Equation (3), the Leacock Equation [32], calculates the relationship between a pair of semantic tags in the tree.

$$sim(t_{i}, t_{j}) = -\log[\frac{ShortestDist(t_i, t_j)}{2 * D}]$$

(3)

*ShortestDist* is the shortest distance between tag $t_i$ and $t_j$ in the semantic tree. D is height of the semantic tree. If the semantic tree is correct, then the more relevant tags are, the shorter the semantic distance will be in the tree. This evaluation method, however, requires manual annotation. The task would be rather huge if all pairs of tags were annotated, so the authors chose to mark fifty tags, which are relevant to a certain tags, but the total task still

reached 34k.

Peng Cheng, Ji Peipei and their partners [27] proposes a method to compare the tree built by algorithms and built by experts in certain field. They compare semantic tree matrix value $M_A$ and the standard semantic tree matrix $M_B$, as shown in Equation (4). Equation (5) is the unrolled one.

$$M_{AB} = \sum_{i=1}^{n} \sum_{j=1}^{n} [M_A(i,j) - M_B(i,j)]^2$$
(4)

$$M_{AB} = \sum_{i=1}^{n} \sum_{j=1}^{n} [M_A(i,j)]^2 - 2\sum_{i=1}^{n} \sum_{j=1}^{n} [M_A(i-j)M_B(i-j)] + \sum_{i=1}^{n}\sum_{j=1}^{n} [M_B(i,j)]^2$$
(5)

Third term in Equation 4 is a constant, which is a semantic tree matrix, which has been obtained. The first term is also a semantic tree matrix value, which varies in account of different clustering methods. The second term is the relationship between $M_A$ and $M_B$, the more similar the greater the value. This method requires the participation of experts in certain field. They need to generate standard semantic tree beforehand.

5. **Conclusions.** This paper introduces researches in the field of term semantic hierarchy extraction for domain-specific Chinese text information processing. Methods based on linguistic templates and clustering are introduced. This paper puts emphasis on hierarchal clustering algorithms for Chinese text, and analyzed problems existed in clustering experiments.

**REFERENCES**

[1] Ji Peipei, Yan XiaoYan, Ceng Yonghua. A Survey of Term Recognition and Extraction for Domain-specific Chinese Text Information Processing [J]. Library and Information Service, 2010, 54(16): 124-129.

[2] Lin D, Pantel P. Induction of semantic classes from natural language text[C]//Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 2001: 317-322.

[3] Maedche A, Staab S. Mining ontology from text[C]//12th International Workshop on Knowledge Engineering and Knowledge Management. 2000.

[4] Wen C, Shi Z, Zhang X. A Survey on Ontology Concept Hierarchy Acquisition [J]. Computer Applications and Software, 2010, 9: 032.

[5] He T T, Zhang X P. Approach to Automatical Construction of Domain Ontology [J]. Computer

Engineering, 2007, 22: 081.

[6]   Jiang Yong, Hypernymy Extraction With Hybrid Text Kernel [D]. 2013.

[7]   Hearst M A. Automatic acquisition of hyponyms from large text corpora[C]//Proceedings of the 14th conference on Computational linguistics-Volume 2. Association for Computational Linguistics, 1992: 539-545.

[8]   Iwanska L, Mata N, Kruger K. Fully automatic acquisition of taxonomic knowledge from large corpora of texts: limited-syntax knowledge representation system based on natural language[C]//In LM Iwanksa and SC Shapiro, editors, Natural Language Processing and Knowledge Processing. 2000.

[9]   Liu L, Cao C, Wang H, et al. A method of hyponym acquisition based on" isa" pattern[J]. Journal of Computer Science, 2006: 146-151.

[10]  Wu W, Li H, Wang H, et al. Probase: A probabilistic taxonomy for text understanding[C]//Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data. ACM, 2012: 481-492.

[11]  Song Wenjie, Zhou Junsheng, Qu Weiguang, Chinese Hyponymy Extraction Based on Dictionary and Encyclopedia Resources [J]. Journal of Data Acquisition and Processing 2014(5)

[12]  Lafferty J, McCallum A, Pereira F C N. Conditional random fields: Probabilistic models for segmenting and labeling sequence data [J]. 2001.

[13]  Deschacht K, Moens M F. Efficient hierarchical entity classifier using conditional random fields[C]//Proceedings of the 2nd Workshop on Ontology Learning and Population. 2006: 33-40.

[14]  HUANG Y, WANG Q, LIU Y. An acquisition method of domain-specific terminological hyponymy based on CRF [J]. Journal of Central South University (Science and Technology), 2013: S2.

[15]  MO Y, GUO J, YU Z, JIANG N, XIAN Y, Hyponymy Extraction of Domain Ontology Concept Based on CCRF [J]. Computer Engeneering, 2014.

[16]  Harris Z S. Mathematical structures of language [J]. 1968.

[17]  Obitko M, Snasel V, Smid J, et al. Ontology Design with Formal Concept Analysis[C]//CLA. 2004, 110.

[18]  WEN C, SHI Z, ZHANG L. Contrast research of Chinese domain ontology concept hierarchy induction methods [J]. Application Research of Computers, 2009, 8: 011.

[19]  Duan M, Research and Application of Hierarchal Clustering [D], Zhongnan University, 2009

[20]  Chuang S L, Chien L F. Towards automatic generation of query taxonomy: A hierarchical query clustering approach[C]//Data Mining, 2002. ICDM 2003. Proceedings. 2002 IEEE International Conference on. IEEE, 2002: 75-82.

[21]  de Mantaras R L, Saitia L. Comparing Conceptual, Divisive and Agglomerative Clustering for Learning Taxonomies from Text[C]//ECAI 2004: 16th European Conference on Artificial Intelligence, August 22-27, 2004, Valencia, Spain: Including Prestigious Applications of Intelligent Systems (PAIS 2004): Proceedings. IOS Press, 2004, 110: 435.

[22]  Brooks C H, Montanez N. Improved annotation of the blogosphere via auto-tagging and hierarchical clustering[C]//Proceedings of the 15th international conference on World Wide Web. ACM, 2006: 625-632.

[23]  Ferragina P, Gulli A. The anatomy of a hierarchical clustering engine for Web-page, news and book snippets[C]//Data Mining, 2004. ICDM'04. Fourth IEEE International Conference on. IEEE, 2004: 395-398.

[24]  CHEN Yuan-hao, ZHANG Ben-yu. ZHANG Hong-jiang, Weighted Aggregation Based Clustering Algorithm for Blog Tag Taxonomy Construction [J]// JOURNAL OF CHINESE COMPUTER

SYSTEMS 2009(7)

[25] Baeza-Yates R, Ribeiro-Neto B. Modern information retrieval[M]. New York: ACM press, 1999.

[26] Ji Peipei, Yan Xiaoyan, Ceng Yonghua, Wang Lingyan, Research of Term Semantic Hierarchy Induction for Domain-specific Chinese Text Information Processing，2010(9)：37-41

[27] Cheng P, Pei-pei J I. Research of term semantic hierarchy correlations based on deterministic annealing [J]. Application Research of Computers, 2011, 9: 009.

[28] Gu Jun，Zhu Ziyang, Study on Ontology Hierarchy Relation Induction on Clustering Algorithm [J]//New Technology of Library and Information Service，2011(12)：46-51

[29] Dorigo M, Blum C. Ant colony optimization theory: A survey[J]. Theoretical computer science, 2005, 344(2): 243-278.

[30] Liu B, Zhang G, Cai D. Technical term automatic extraction research based on statistics and rules[J]. Computer Engineering and Applications, 2008, 44(23): 147-150.

[31] Liu L, Research on Automatic Bilingual Term Extraction Technology For Patents [D]. Shenyang, 2009.

[32] Leacock C, Chodorow M. Combining local context and WordNet similarity for word sense identification [J]. WordNet: An electronic lexical database, 1998, 49(2): 265-283.