

Association for Information Systems AIS Electronic Library (AISeL)

AMCIS 2002 Proceedings

Americas Conference on Information Systems
(AMCIS)

December 2002

ANALYSIS OF AN AGENT-BASED APPROACH FOR DISCOVERING TERM SEMANTIC RELATIONSHIP

Lina Zhou

University of Maryland, Baltimore County

Dongsong Zhang

University of Maryland, Baltimore County

Follow this and additional works at: <http://aisel.aisnet.org/amcis2002>

Recommended Citation

Zhou, Lina and Zhang, Dongsong, "ANALYSIS OF AN AGENT-BASED APPROACH FOR DISCOVERING TERM SEMANTIC RELATIONSHIP" (2002). *AMCIS 2002 Proceedings*. 207.

<http://aisel.aisnet.org/amcis2002/207>

This material is brought to you by the Americas Conference on Information Systems (AMCIS) at AIS Electronic Library (AISeL). It has been accepted for inclusion in AMCIS 2002 Proceedings by an authorized administrator of AIS Electronic Library (AISeL). For more information, please contact elibrary@aisnet.org.

ANALYSIS OF AN AGENT-BASED APPROACH FOR DISCOVERING TERM SEMANTIC RELATIONSHIP

Lina Zhou and Dongsong Zhang
University of Maryland, Baltimore County
zhoul@umbc.edu zhangd@umbc.edu

Abstract

Despite exponential accumulation of electronic texts in the Internet age, effective search and reuse is not available for most texts due to the lack of association among them. Terms and the semantic relationship between terms have great potential in organizing and sharing textual information. We develop an agent-based approach for discovering term relationship in texts collected from the Web. The goal of the paper is to investigate which of the five generic relationships in text, including synonym, hypernym, hyponym, holonym and meronym, can be effectively discovered with the proposed approach.

Introduction

With the ever-growing expansion of computer and Internet technology into many disciplines, electronic texts are accumulated exponentially at an unprecedented speed. However, effective search and reuse is not available for most texts due to the lack of association among them. Terms and the semantic relationship between terms have great potential in organizing and sharing text information. If we know how terms within a domain are semantically related, we can reorganize a large number of texts by linking them with semantic relations. It enables us to navigate and search for required information in an intelligent and efficient manner. Therefore, it has of both theoretical and practical significance to identify term semantic relations from texts.

With computers getting more powerful and information increasing distributed, agent-based (AB) approach has been adopted for many purposes (Jennings and Wooldridge, 1998). However, the application of AB approach to discovering term relationship from large number of texts has rarely been explored. We are going to analyze an AB approach in discovering term semantic relationship in this paper. The approach automatically identifies terms and their related terms in texts collected from the Web. The results are expected to help people better understand the terminology in a domain in an interrelated way and provide a coarse framework for building domain knowledge or ontology.

The discussion on automatic discovery of similar words has been mainly devoted to finding words that belong to the same semantic class for general purpose and word pairs that have a specific relationship (Hearst, 1992; Jean-David, 1998; Li and Abe, 1998; Lin, 1998; Riloff and Shepherd, 1999; Maedche and Staab, 2000). The prior approaches were mostly tested on established corpus, which was not easily adaptable to different domains. Therefore, the issue of what kinds of relations are more likely to be discovered by an AB approach is still left open. If an AB approach is found to be better at detecting some relation types than others, we can customize the approach to the former relation types and leave the latter to other approaches. The types of identified relationship could also potentially become a dimension for comparing and integrating AB approaches. In this paper, we focus on five types of the popular generic term relationship in texts, including synonym, hypernym, hyponym, holonym and meronym. The goal is to find out which ones of the above-mentioned relationships can be effectively discovered with an AB approach.

The rest of the paper is organized as follows. In Section 2, we introduce the concept of semantic relationship between terms, review the related research on automatically learning term relations, and propose research questions. Then, we describe STARTER, an AB approach we developed for acquiring term relations in Section 3. Next, we present metrics for measuring the effectiveness of the STARTER in detecting the five selected relations, and analyze the evaluation results in details in Section 4. Finally, we conclude our study and discuss potential applications of the findings.

Related Work and Research Questions

A term can be described as a meaningful unit that consists of content word(s) and has distinct attributes in texts within a certain domain. Every term is potentially related to a group of other terms through some kind of semantic relations. The relations could be either domain-specific or generic. The generic relations include hyponym/hypernym (IS-A/HAS-A), meronym/holonym (Part-of/Has-Part), and synonym, etc. In this paper, we focus on the above five types of generic relations, for they are most frequently discussed in the literature.

We found two lines of research that are related to the discovery of term relations. One is concerned with improving the effectiveness of finding similar words on general basis. Lin (1998) automatically extracted similar words from texts and created a tree structure for words based on semantic similarities. Rioloff and Shepherd (1999) proposed a bootstrapping algorithm for creating semantic lexicon from text corpus supplemented with a set of seed words. Both research only dealt with single words rather than terms. Zhou et al. (2002) made an effort to automatically acquire relationship between terms, which are better candidates for domain terminology than words. Most of the above researches adopted statistical approach, but none of them investigated the specific types of relations that are embedded in the discovered similar words or terms. The other line of research focuses on extracting specific relations from texts, such as hyponym (Berland and Charniak, 1999; Hearst, 1992), which normally uses manually compiled heuristics rather than AB approach.

If we know the types of relations that are more likely to be discovered with an AB approach, we can continue to refine algorithms for detecting those relations, and develop different approaches for other relation types. Therefore, we are mainly interested in the following research question: among the five selected relations, which one(s) can an AB approach discover most effectively? As secondary objectives, we also look into another two related questions:

- (1) Can an AB approach effectively identify term relations from texts?
- (2) What kinds of relations are more likely to be discovered simultaneously?

STARTER: An Agent-Based Approach

In this section, we briefly introduce an agent-based approach, STARTER, which was developed to address the above-mentioned research questions. STARTER has the following desirable features: 1) it is adaptable to different domains; 2) it recognizes terms in addition to words; 3) it adopts multi-featured criteria for selecting terms; and 4) it is semi-autonomous.

There are four main stages involved in the STARTER: document collection, term identification, key term selection, and related term extraction and ranking.

Document Collection

The first stage in STARTER is collecting documents from the Web. Internet spiders are software applications that collect Internet pages by following outgoing links in each page recursively. The input to spiders could include one or more addresses of Internet pages, one or more keywords related to a concerning domain, and so on. The fetched pages are unformatted by a parser, which extracts textual content from the pages. The transformed pages whose sizes are greater than 50 bytes are collected into the set of source documents for studying terms relationship.

Term Identification

STARTER uses an NLP tool, EngLite (Voutilainen, 2000), to recognize term candidates and partial terms from documents. Based on the linguistic information generated by the tool, such as lemma, part-of-speech (POS), morphology, and light syntax, we perform post-processing to automatically construct multi-word terms.

Key Term Selection

A collection of documents usually contains many term candidates. Motivated by keyword extraction technologies in IR, we filter candidates by removing those with overly high and low occurrence frequencies. In addition, we propose multi-featured filtering criteria, which are based on a term's POS and weights. During the POS filtering, we eliminate term candidates with certain parts-

of-speech, such as articles and propositions, for such candidates rarely become key terms. The weight filtering attempts to remove terms whose weights are lower than a certain threshold, for it is not meaningful to examine the relations between low-weight terms and other collocated terms. Among the variety of possible weighting algorithms, we choose the adapted *tfidf* scheme, commonly used in question-answering system (Singhal, et al., 1998). This scheme weights a term by a factor dependent upon both its importance within a document and its importance within a whole collection. The raw term frequency is normalized using logarithmic function, thus dampening effects of large differences in term frequencies. The weight of *kth* term in *ith* document, w_{ik} , is defined as follows:

$$w_{ik} = f_{ik} \times \ln(N_D / n_k) \tag{1}$$

$$f_{ik} = \begin{cases} 0, & \text{if } f_{ik} = 0; \\ \ln(f_{ik}) + 1, & \text{otherwise} \end{cases}$$

where N_D is the total number of documents in a collection C ,
 n_k is the number of documents in C that contain term k ,
 f_{ik} is the frequency of term k in document i .

Related Term Extraction and Ranking

Term Representation

Before measuring term similarity, we need first create a model for representing terms. We select feature set to model a term, for it is flexible in accommodating different kind of information. Elements of the feature set should be independent of domain knowledge base, for STARTER is expected to be adaptable to different domains. Among the limited information that is produced from the previous steps, we choose documents in which a term occurs and the associated term weights in those documents as elements of the feature set. As a result, we develop a term vector model to represent terms, consisting of related document numbers and corresponding term weights. The maximum size of each vector is the total number of documents in the collection. The vectors of all terms constitute a term vector space.

Similarity Measurement

A similarity metrics needs to be established for measuring term similarity. We expect that the higher the similarity value between a pair of terms, the more closely they are related, and vice versa. Thus, the similarity value reflects the strength of relatedness. It is common to formulate a metrics as a function of co-occurrence frequencies of two terms in documents. High co-occurrence frequency of a pair of terms may indicate strong relation between them, although it does not always hold true in some cases. One typical exception refers to very high-frequency terms, which appear in the same documents as many other terms, but their frequent co-occurrence with other terms may happen by chance. Aiming at differentiating significant and non-significant co-occurrences, STARTER selects cosine (Rijsbergen, 1979) as similarity metrics.

Cosine similarity measure (Rijsbergen, 1979) has been widely applied in information retrieval to match a query to documents, whose value ranges from 1 for perfectly relevant 0 for perfectly irrelevant. We extend *cosine* function to measuring term similarity in the term vector space. The similarity between term t_i and term t_j is defined as:

$$sim(t_i, t_j) = \frac{\sum_{k=1}^L w_{ik} * w_{jk}}{\sqrt{\sum_{k=1}^N (w_{ik})^2 * \sum_{k=1}^M (w_{jk})^2}} \tag{2}$$

where w_{ik} and w_{jk} are the weights of t_i and t_j in the document k respectively,
 N and M are the number of documents where t_i and t_j occur respectively,
 L is the number of documents in a collection where both t_i and t_j occur.

As shown in formula (2), the similarity between two terms is measured by the angle between their vector representations in a multi-dimensional space. The similarity value is used to estimate the relatedness of one term with respect to the other.

Related Term Set (RTSET) extraction

According to the *cosine* measure, a term may be related to many other terms to different degrees, thus we need to develop criteria for automatically extracting those related terms that are more reliable. A Related Term SET (RTSET) of a term is defined by a list of extracted related terms and their corresponding similarity values with the term. We set up two thresholds for generating RTSET based on similarity calculation results. One is g , minimum similarity value. The other is f , maximum number of related terms that could be extracted for a given term. Both can be manually or automatically adjusted based on the amount of output and the expected accuracy. The extraction of RTSET proceeds in three steps. First, all the term candidates are extracted and filtered according to g . Second, the remaining candidates are ranked in descending order of their similarity values. Third, the top f candidates, if any, are extracted into the final RTSET.

Evaluation

Test Domain and Data

We chose community development (CD) as the test domain. It was of practical advantage and significance for us to make such a choice. First, we have experienced domain experts available, who can help us in collecting data and evaluating results. Second, community development is such a domain that is in great need of domain dictionary or lexicon. We can help build knowledge resource while attaining our research goals. Third, common terms take a significantly large percentage of texts in. This enables us to compare the experimental results with generic lexical resource such as WordNet (Miller, 1990).

STARTER first collected 3,725 documents on community development. Then, it identified 6,032 unique terms and their associated RTSET from the document collection. In order to make the evaluation more efficient, we randomly selected about 10 percent of terms and their RTSET as the test data. Finally, 576 terms were extracted for evaluation.

Evaluation Approach

Baseline

As stated before, there does not exist any dictionary or thesaurus in CD that can be used as a perfect baseline for our evaluation. Instead, we took advantage of a general-purpose thesaurus – WordNet (WN). It implied that the evaluation focuses on generic rather than domain-specific terms and relations. WN captures many types of generic relations between words. The generic relations selected for this study include synonym, hypernym, hyponym, holonym, and meronym. Even with such a compromise, our research goals can still be reasonably achieved for the following reasons: 1) the selected five relations are the most popular generic relation types; 2) all the relations are subject to the same restriction of being evaluated based on general terms. The impact of the restriction can be eliminated in comparison between different relations; and 3) CD is characterized with common word usage and broad content coverage, which establishes a large base for comparison with WN. This is demonstrated by the fact that 421 out of 576 randomly selected terms are found in WN.

WN is organized into Synonym sets. A SYNonym SET (SYNSET) in WN is similar to a RTSET in the acquisition results of STARTER (ST). For example, the SYNSET and RTSET of term *plan* are:

plan 3 4 program programme design architectural *plan* (SYNSET)
plan activity 0.39 program 0.39 management 0.37 change 0.35 (RTSET)

According to the overlap between terms in ST and WN and between RTSET and SYNSET of the same term, we classified terms in the ST into seven categories, as shown in Table 1. If a term is found in WN (or ST), and the size of the SYNSET (or RTSET) is at least one, we call the term *exist*. We label a term with *empty* if the corresponding SYNSET (or RTSET) is empty. If a term is not included in WN, or a term is included in both ST and WN but none of the terms in the SYNSET appears in the same document as the original term, it is considered as *null*. Even though SYNSET is the basic unit of WN, we can literally call SYNSET of hypernym as HYPESET. Likewise, we call SYNSET of hyponym, holonym, and meronym as HYPOSET,

HOLOSET, and MEROSSET respectively. In addition to SYNSET, term classification in Table 1 can be applied to the intersection between RTSET from ST and any of the other four types of SET from WN such as HYPESET.

Table 1. Classification of Terms

WN	ST		
	exist	empty	null
exist	ws	wn	wnu
empty	ns	nn	n/a
null	nus	un	n/a

Evaluation Metrics

Two metrics were selected to evaluate the performance of the AB approach: precision and recall. The *precision* of a term is the percentage of terms in the RTSET that are truly related to it, while the *recall* is the percentage of truly similar terms that are included in the RTSET. The classification of terms in Table 1 helped us tailor the evaluation metrics to each type. Since *nus* and *nun* terms were not available in WN, we ignored them in the evaluation. For *wnu* type, there was no way for an AB approach to find them due to the data-driven nature, so we excluded them in the evaluation. RTSET and SYNSET of *wn* and *ns* types had no intersection, so both of their precisions and recalls were 0. On the contrary, the precision and recall of *nn* type were both 1, for terms that do not have the specific relationship were correctly identified in this case. Finally, we only left the precision and recall of *ws* type for definition as follows:

$$Precision_t = \frac{|SYNSE_T^{WN} \cap RTSET_t^{ST}|}{|RTSET_t^{ST}|} \tag{3}$$

$$Recall_t = \frac{|SYNSE_T^{WN} \cap RTSET_t^{ST}|}{|SYNSE_T^{WN} \cap \forall k (D(t_k) \cap D(t_i) \neq \Phi)|} \tag{4}$$

where $D(t_i)$ is all the documents in which term t_i appears.

Even though only SYNSET is listed in the above formula, we can replace SYNSET with any of the other four relation sets, such as HYPOSET, in order to get the precision and recall for that relation.

The overall precision and recall are the average precision and recall of all the terms. Since we only considered four types of terms from the previous analysis, namely *nn*, *ns*, *wn*, and *ws*, the formula for overall precision was simplified as:

$$OverallPrecision = \frac{\sum_{i=1}^T Precision_t}{T} = \frac{|nn| + \sum_{i=1}^{|ws|} Precision_t}{T} \tag{5}$$

where $T = |nn| + |ns| + |wn| + |ws|$

The formula for the overall recall was simplified in the same way.

Results and Analyses

Overall Result Analysis

In this section, we briefly show the overall effectiveness of an AB approach in detecting term semantic relationships and the relation types that can be more effectively discovered by the approach.

For each of the seven term types in Table 1, the number of terms for each of the five relations was listed in Table 2. The five relation types are short noted with the first four characters in the rest of this paper. For example, hype stands for hypernym.

Table 2. Frequency Matrix of Term and Relation Types

Term Type	Term Relation				
	holo	mero	hype	hypo	syno
nn	44	43	0	12	0
ns	262	278	0	99	0
wn	2	2	9	8	9
ws	2	2	30	20	31
wnu	111	96	382	282	381
nus	148	148	148	148	148
nun	7	7	7	7	7

It was shown from Table 2 that 155 (*nus* + *nun*) out of 576 terms were not available in WN, and the number of *wnu* terms was different from one relation to another. The overall average precisions and recalls of the remaining term types for each of the five relations were displayed in Figure 1 and Figure 2 respectively. It revealed that the precisions and recalls of holo and mero relations were the highest, while those of hype and syno were the lowest. If we only considered *ws* terms, the precisions and recalls of holo and mero were the lowest, for there were only 2 terms fall in this category and neither of them was correct. For the remaining three relations, the average recalls ranged from 7.5 to 8.89 percent, and precisions from .77 to 1.03 percent. Since by default there exists at most one type of relationship between a pair of terms, so we derived the overall average precision of *ws* terms, 2.79 percent, by aggregating individual precisions of five relations. It indicated that about 2.8 percent of terms in the extracted RTSET of *ws* terms really had the corresponding relations with the original terms.

We conclude from the above analysis that STARTER is better at discovering terms that do not have holo and mero relationship, and those that have hype, hypo and syno relationship. It is partly because statistically fewer terms (*ws+wn+wnu*) have holo and mero relationship than the other relationship, partly because less percentage ($ws/(ws+wnu)$) of mero and holo relationship occurs in the same documents than other relationship, as shown in aggregated measures on term types in Table 3.

Table 3. Aggregated Measures on Term Types

Measurement	holo	mero	hype	hypo	syno
<i>ws+wn+wnu</i>	115	100	421	310	421
$ws/(ws+wnu)$	2.91%	3.35%	8.09%	7.63%	8.33%

Hype and syno relations displayed lower precisions and recalls than others. It is not difficult to find from Table 2 that the precisions and recalls of hype and syno were exclusively contributed by *ws* terms. It was because there were 0 *nn* terms and by definition the values of *wn* and *ns* terms were 0. It can be inferred from Table 2 that 26.91 percent ($(nun+nus)/576$) of extracted terms for evaluation were not collected in WN. Even though these two types of terms were not considered in the evaluation, they were potentially useful for the discovery of new terms falling in one of the following potential categories:

- (1) Word combinations that have been frequently used together and have their special meanings in CD, such as *program coordinator*, *downtown revitalization*, *economic education*, *child program*, and *outcome measurement*.
- (2) New acronyms that are recently created and adopted by people in certain areas of CD, such as *CDC*, and *CSBG*.
- (3) Proper nouns related to CD, such as *Empowerment Zone*, and *Community Development Society*.

It is important to capture all these new usages when creating domain resources.

Detailed Result Analysis

As shown in the previous section, there were differences on precision and recall between different relation types. In this section, we further look into whether the differences are significant. Moreover, we are going to explore the types of relations that are more likely to be discovered simultaneously by analyzing the correlation between the results of individual relation types.

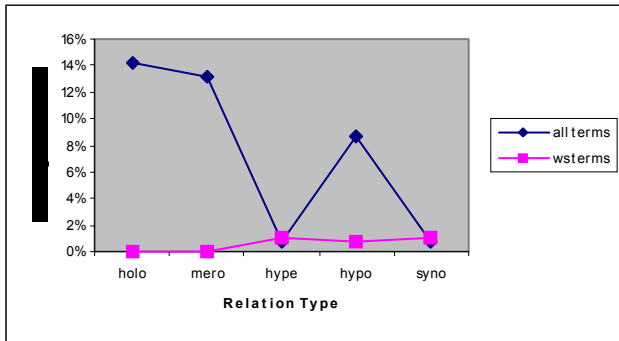


Figure 1. Results of Average Precision

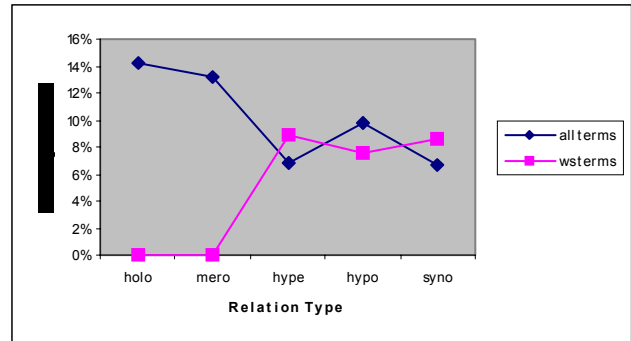


Figure 2. Results of Average Recall

Taking each term as a sample and each relation type as a dependent variable, we conducted correlation analysis and paired T-test. Before the statistics analysis, we checked term type patterns in order to make sure the data were clean. The most frequent term patterns across five relations were listed in Table 4. The first pattern told us that 134 terms did not have holo and mero relations but STARTER mistakenly generated such relations for them, and they had hype, hypo, and syno relations but none of the relationship occurred in the same documents. The 8th row showed that 11 terms did not have holo, mero, and hypo relationship and STARTER correctly recognized them, and they had hype and syno relationship but none of the related terms was found in the same documents. Other patterns can be interpreted in a similar way. It can be inferred from Table 4 that statistically fewest terms had holo and mero relationship, relatively more terms had hypo relationship, and most terms had hype and syno relationship.

Table 4. Most Frequent Patterns of Term Types across Five Relations

Pattern	holo	mero	hype	hypo	syno	Count
1	ns	ns	wnu	wnu	wnu	134
2	ns	ns	wnu	ns	wnu	70
3	wnu	wnu	wnu	wnu	wnu	37
4	wnu	ns	wnu	wnu	wnu	33
5	ns	wnu	wnu	wnu	wnu	31
6	nn	nn	wnu	wnu	wnu	19
7	wnu	ns	wnu	ns	wnu	17
8	nn	nn	wnu	nn	wnu	11
9	ns	ns	ws	wnu	ws	10

wnu terms did not have any related term from WN occurring in the same documents in the collection, but we could not guarantee that the related terms would not co-occur in other documents. However, we noticed from Table 4 that wnu frequently appeared in most frequent patterns. It cautioned us to cleaning the sample data by removing wnu terms. Since it rarely happened that all five relations had wnu terms at the same time, so we had to filter wnu terms for each pair of relations in comparison separately. As a result, the sample data included terms belonging to such types as ws, nn, ns, and wn.

Based on our analysis, the correlations between the recalls of five relation types were consistent with those between the precisions. It demonstrated that holo & mero (N=264) and syno & hype (N=39) were perfectly correlated, hypo & mero (N=118, p=.000) and hypo & holo (N=106, p=.000) relations were positively correlated at $\alpha=0.01$ level, and other pairs of relations were negatively

correlated but not significant. The positive correlation implied that if we found one relation correctly, we were likely to detect the other relation correctly as well; and the more likely we identified one relation, the higher probability we found the other relation. For example, if we discovered a holo relation of a term, we were likely to find mero relation as well, and vice versa. We figured that the nearly perfect correlation between mero and holo was attributed to the large number of *nm* cases, and the correlation between syno and hype may be due to the fact that neither of them had *nm* or *ns* cases.

It was shown from the Paired T-test results that there was significant difference on the average precision between hype & holo ($p=.007$), hypo & holo ($p=.009$), and syno & holo ($p=.007$) relations at the confidence level $\alpha = 0.01$, and significant difference between hype & mero ($p=.016$), hypo & mero ($p=.015$), and syno & mero ($p=.016$) at $\alpha=0.025$. In addition, there was significant difference on the average recall between hype & holo ($p=.037$), and syno & holo ($p=.037$) at $\alpha=0.05$.

By combining the results of correlation analysis and T-test, we can see that precisions of holo and mero were significantly higher than those of hype, hypo, and syno, and the recall of holo was significantly higher than that of hype and syno. It was interesting to see that hype & mero and hypo & holo were both significantly correlated despite their significantly different precisions.

Discussion

The overall precision and recall do not seem high. However, in view of the large frequency of term occurrences in the documents and little input of domain knowledge to the STARTER approach, the results were reasonably good. We can improve precision or recall by adjusting the thresholds for RTSET extraction. If we reduce the value of *g* and/or increase the value of *f*, the recall could be improved. The precision could be improved the other way around. Furthermore, WordNet is a general-purpose resource, which does not reflect the domain specificity of community development knowledge. We severely punished the mismatched terms between WN and ST by assigning 0 to both precision and recall of *wn* and *ns* terms. Since all five relations were subject to the same impact, it did not prevent us from comparing different types of relations consistently.

Conclusions and Future Research

In this paper, we investigated the potential of AB approach in discovering term relationship by testing the STARTER. Based on comprehensive result analyses, we found that the AB approach was better at discovering terms that do not have holo and/or mero relations, and terms that have hype, hypo and/or syno relations. It was shown from the results of correlation analysis and T-test that precisions of holo and mero were significantly higher than those of hype, hypo, and syno, and the recall of holo was significantly higher than those of hype and syno. In addition, *nus* and *nun* terms could be useful for discovering new domain-specific terms, *ns* for discovering new relations, and *wnu* for removing non-domain-specific relations.

Domain specificity is an important characteristic of knowledge. With growing interest in domain ontology in recent years, the AB approach can be extended to automatically acquiring ontological knowledge. It could become an efficient way of capturing and updating knowledge as new terms and relations are constantly created and old ones deprecated. After analyzing semantic relations that are discovered by the proposed AB approach, we can target the approach to the relation types that it is suitable for and develop different agents or approaches for identifying other relation types. In the long term, we will be able to build a multi-agent system that can discover various relation types with coordinated efforts.

References

- Berland, M., and Charniak, E. "Finding Parts in Very Large Corpora," *Proceedings of the ACL-99*, 1999, pp. 57-64.
- Hearst, M.A. "Automatic Acquisition of Hyponyms from Large Text Corpora," *Proceedings of the Proceedings of the Fourteenth International Conference on Computational Linguistics*, Nantes, France, 1992, pp. 539-545.
- Jean-David, S. "Automatic acquisition of terminological relations from a corpus for query expansion," *Proceedings of the ACM SIGIR'98*, Melbourne, Australia, 1998, pp. 371-372.
- Jennings, N.R., and Wooldridge, M. (eds.). *Agent Technology: Foundations, Applications, and Markets*, Springer Verlag, 1998.
- Li, H., and Abe, N. "Word Clustering and Disambiguation Based on Co-occurrence Data," *Proceedings of the COLING-ACL'98*, 1998, pp. 749-755.
- Lin, D. "Automatic retrieval and clustering of similar words," *Proceedings of the COLING-ACL'98*, Montreal, Canada, 1998, pp. 768-774.

- Maedche, A., and Staab, S. "Mining ontologies from text," *Proceedings of the 12th International Workshop on Knowledge Engineering and Knowledge Management*, French Riviera, 2000, pp. 189-202.
- Miller, G.A. "WORDNET: An On-Line Lexical Database," *International Journal of Lexicography* (3-4), 1990, pp. 235-312.
- Rijsbergen, C.J.V. *Information Retrieval*, London: Butterworths, 1979.
- Riloff, E., and Shepherd, J. "A Corpus-Based Bootstrapping Algorithm for Semi-Automated Semantic Lexicon Construction," *Journal of Natural Language Engineering* (5:2), 1999, pp. 147-156.
- Singhal, A., Choi, J., Hindle, D., Lewis, D.D., and Pereira, F. "AT&T at TREC-7," *Proceedings of the Seventh Text Retrieval Conference (TREC-7)*, Gaithersburg, Maryland, 1998, pp. 239-252.
- Voutilainen, A. "Helsinki taggers and parsers for English," In *Analyses and Techniques in Describing English*, J. M. Kirk (ed.) Rodopi: Amsterdam & Atlanta, 2000.
- Zhou, L., Booker, Q.E., and Zhang, D. "ROD - Toward Rapid Ontology Development for Underdeveloped Domains," *Proceedings of the 35th Hawaii International Conference on System Sciences*, Big Island, Hawaii, 2002.