

## Semantic Extraction from Texts

Shaidah Jusoh<sup>1</sup>, Hejab M. Al Fawareh<sup>2</sup>

<sup>1</sup> College of Arts and Sciences , Universiti Utara Malaysia  
Sintok, Kedah, Malaysia  
shaidah@uum.edu.my

<sup>2</sup> College of Arts and Sciences , Universiti Utara Malaysia  
Sintok, Kedah, Malaysia  
alfwareh@gmail.com

**Abstract.** - Text documents are one of the means to store information. These documents can be found on personal desktop computers, intranets and in the Web. Thus the valuable knowledge is embedded in an unstructured form. Having an automated system that can extract information from the texts is very desirable. However, the major challenging issue in developing such an automated system is a natural language is not free from ambiguity and uncertainty problems. Thus semantic extraction remains a challenging task to researchers in this area. In this paper, a new framework to extract semantics for information extraction is proposed, where possibility theory, fuzzy sets, and knowledge about the subject and preceding sentence have been used as the key in resolving the ambiguity and uncertainty problems.

**Keywords:** Semantic Extraction, Information Extraction, Possibility Theory.

### 1. Introduction

Nowadays, the Web is considered as the world's largest repository of knowledge, and it is being constantly augmented and maintained by millions of people around the world. However, it is not in the form of a database from which records and fields are easily manipulated and understood by computers, but in natural language texts which are intended for human reading. In spite of the promise of the semantic web, the use of English language and other natural language texts will continue to be a major medium for communication, knowledge accumulation, information distribution on the Web, emails, reports, memos, blogs and etc [1]. People want to extract useful information from the texts documents quickly at a low cost. Text mining is a new area which focuses on the use of automated methods for exploiting the enormous amount of knowledge available in text documents. Text mining, sometimes alternately referred to as text data mining, refers generally to the process of deriving high quality information from texts [2].

Typical text mining tasks include text categorization, text clustering, concept/entity and fact extraction, and production of granular taxonomies, sentiment analysis, document summarization, and entity relation modeling [3]. When dealing with natural language texts, the most critical problem is ambiguity and uncertainty issues. Automated information extraction (IE) system should be able to extract correct semantics from texts. Thus the ambiguity and uncertainty issues should be resolved. In this research work, we propose a new framework for semantic extraction. The framework is based on *the knowledge of subject and relevant preceding sentence*. This paper is organized as follows. Section 2.0 will discuss information extraction; section 3.0 will present a proposed framework. The implementation and result analysis are presented in section 4.0. Section 5.0 concludes the paper.

<sup>+</sup> Corresponding author. Tel.: + 604-9284624; fax: +604-9284753.

*E-mail address:* shaidah@uum.edu.my.

In the last several years, IT practitioners have agreed that there exists a continuum of data, information and knowledge. Data is mostly structured, factual, and numeric. Data consists of facts, images, or sound. When data is combined with interpretation and meaning, information emerges. Knowledge is inferential abstract that is needed to support decision making process. Knowledge can be as simple as knowing who is the president of the United States, or it can be as complex as mathematical formula relating process variables to finish product dimensions. To distinguish between information and knowledge is not always straightforward. [1] defined knowledge as “a fluid mixed of framed experience, values, contextual information, but until people use it, it isn’t knowledge”. While [2] use knowledge definition taken from [3], that the primary elements of knowledge are concepts and relationships between concepts. Basically, [4] defined concepts as ‘perceived regularities in events or objects, or records of events or objects, designated by a label’. Knowledge exists in forms such as instinct, ideas, rules, and procedures that guide actions and decision. Most researchers agree that knowledge is a human creation. Thus we can construct new knowledge by linking new concepts/entities the knowledge that we already have [5].

## 2. Related Areas

In discussing semantic extraction, we should highlight that the most relevant application that is IE. According to [4], IE does a more limited task than full text understanding. [4] pointed that in full text understanding, all the information in the text is presented, whereas in IE, the semantic range of the output, the relations will be presented are delimited. Traditionally in IE, natural language texts are mapped to predefined, structured representation, or templates, which, when filled, represent an extract of key information from the original texts [5, 6].

In IE, there are two levels of extractions; entity extraction and fact extractions. Extracting entity/concepts from the texts require a person to read them. Fact extraction is a process of spreading out the facts from entities. This is very time consuming. It can become a challenging task if the person does not have enough background related to the texts. Having an automated system that can extract required information from the texts is becoming an urgent need. However, this desire is not easy to achieve. Natural language texts are not free from the ambiguity problems. It is not only many words may refer to one meaning and one word may have more than one meaning, but also a structure of the sentence can be interpreted into more than one meaning.

On the other hand Singh [7] and Hale [8] addressed information extraction is based on understanding of the structure and meaning of the natural language in which documents are written and the goal of information extraction is to accumulate semantic information from text. Technically extracting information from texts requires lexical knowledge, grammars describing the specific syntax of the texts to be analyzed as well as semantic [9].

Today, most of the IE systems that involve semantic analysis exploit the simplest part of the whole spectrum of domain and task knowledge, that is to say, named entities. However, the growing need for IE application to domains such as functional genomics that require more text understanding. For example, in biomedical domain, entities would be gene, protein names and drugs. NER often forms the starting point in a text mining system, meaning that when the correct entities are identified, the search for patterns and relations between entities can begin. [10] also claim that one of the major problems in NER is ambiguous protein names; one protein name may refer to multiple gene products.

Although [11] have put effort to resolve ambiguous terms using sense-tagged corpora and UMLS, the ambiguity is still the major “world problem” [10] in IE. In fact [11] work focus only on biomedical terms only. Recognizing and classifying named entities in texts require knowledge on the domain entities. List entities are used to tag text entities, with the relevant semantic information; however exact character strings are often not reliable enough for precise entity identification [8].

Recent applications in information extraction include apartment rental ads [12], job announcements [13], geographic web documents [14], and medical abstracts [10]. [15] point out that much published work on IE reports on closed experiments; systems are built and evaluations are conducted based on carefully annotated

training and test corpora. Although IE has been implemented for varieties of applications as mentioned above, up to date, automated IE has not yet involved with semantic extraction.

### 3. Proposed framework

Our proposed framework solves the ambiguity and uncertainty problems in semantic extraction for IE at two levels of extraction. The first is at the entity extraction level and the second is at the fact extraction level as shown in Figure 3.1. The whole process of extracting entity and facts from texts can be condensed into 3 steps as illustrated in Figure 3.1.

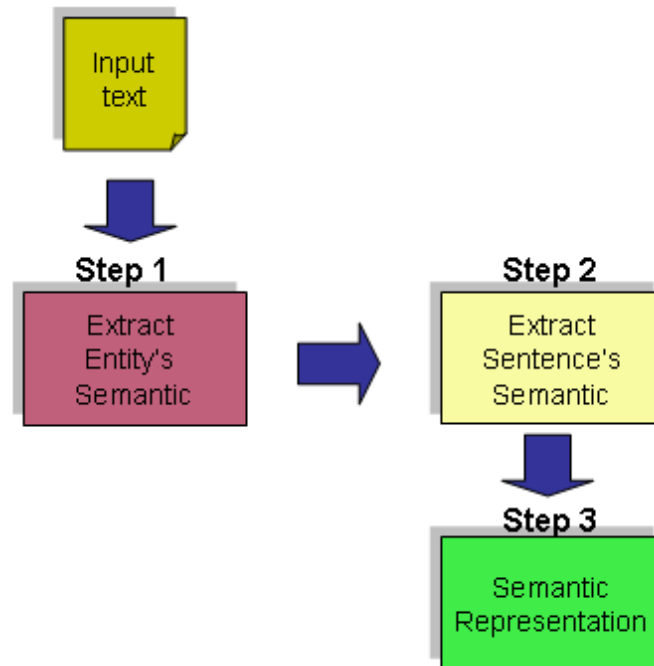


Figure 1: The steps for semantic extraction

#### 3.1. Step 1

In this step, the text input is segmented into sentences. Each sentence will be processed syntactically to recognize its part of speech. The word that is belong to a verb or a noun part-of-speech category is defined as an entity. Let us consider the following sentences as examples:

I put the baby in the pen

She runs the company

From syntactic processing, the system would be able to determine that the word pen is belong to part-of-speech for a noun category. The syntactic processor also can determine that “runs” is a verb. However, when the system needs to extract semantic of the word, the system would face ambiguity and uncertainty problem. For example, a word ‘pen’ can be interpreted as a writing tool, or an enclosure, in which babies may be left to play. While the word ‘run’ can be interpreted as an activity of controlling or as a physical action. In information extraction, semantic of the texts should be correctly interpreted.

To resolve the problem we have applied *subject context knowledge* during the semantic processing. Figure 3.2 illustrates the process.

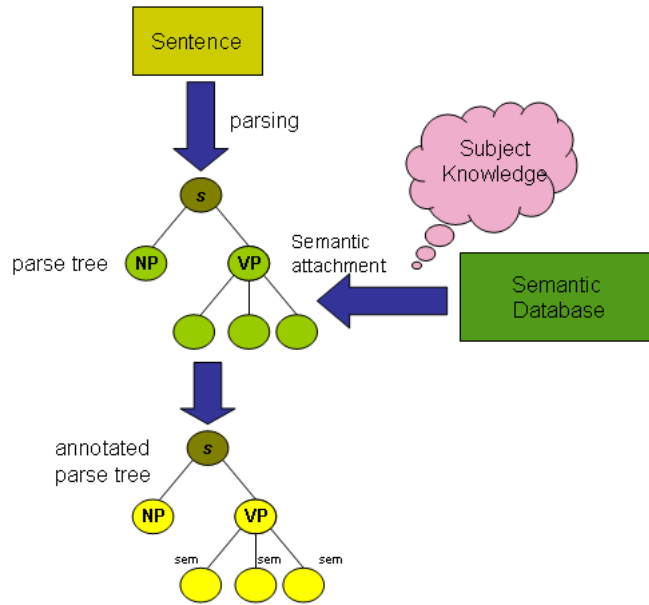


Figure.2: The process of Step 1

As previously mentioned, the structure of a sentence (parse tree) is obtained through a parsing/syntactic process. Using the possibility theory, the possibility value is assigned to each meaning of the words. The value is determined by the subject context knowledge. Let us consider pen as a word ( $w$ ) and its meanings; a tool for writing ( $m_1$ ) and an enclosure ( $m_2$ ). The possibility ( $\rho$ ) of  $w = m_1$  or  $w = m_2$ , is determined by subject context knowledge (SK), which can be formulated as follows

$$w = (m_1, m_2, m_3, \dots, m_n) \quad (1)$$

where  $m_1, \dots, m_n$ , represent the possible meaning of the word  $w$ , and  $n$  is a finite number of the meaning.

$$\rho = (\rho_1 = m_1, \rho_2 = m_2, \rho_3 = m_3, \rho_n = m_n) \quad (2)$$

The possible meanings of  $w$  is represented by  $\rho_1, \rho_2, \dots, \rho_n$ . The value of  $\rho_1, \rho_2, \dots, \rho_n$  is decided based on the SK as represented in the Table 1.

Word ( $w$ )	Semantic ( $m$ )	Possibility Value ( $\rho$ )
Pen	A writing tool with a point from which ink flows	0.5
Pen	An enclosure for confining livestock	0.1
Pen	An enclosure in which babies may be left to play	0.9
Pen	A correctional institutions for those convicted crime	0.2
Pen	Female swan	0.4

Table 1: A semantic database for “baby” context.

In Table 1, the context of the word pen is “baby”. In this work, fuzzy operator  $max$  is used to select the most possible meaning of the pen as formulated in Eq. 3

$$\rho = \max(\rho_1, \rho_2, \rho_3, \dots, m_n) \quad (3)$$

Thus, by applying Eq. (3), the syntactic processor is able to decide the most possible meaning of the word 'pen', which is an enclosure in which babies are left to play. Therefore, if the subject knowledge is "writing" the values of the possibility in Table 1 would be different. Once the ambiguity and uncertainty problems, a correct semantic is attached to the parse tree. The annotated parse tree would be used for the process in the step 2.

### 3.2. Step 2

In Step 2, annotated parse tree is used to determine the semantic meaning of the sentence. Let us consider the sentence "I put the baby in the pen". Although, step 1 has resolved that the ambiguity problem for the word pen, during the parsing process, the syntactic processor would also generate more than one parse tree. This happens because of the ambiguity in the grammar itself. The sentence can be parsed in two ways; the first parse tree is parsed through production grammar rules in 1, and the second parse tree through production grammar rules in 2, as illustrated below.

$$1. \left. \begin{array}{l} s \rightarrow \text{Pronoun, VP} \\ \text{VP} \rightarrow \text{Verb, NP, PP} \end{array} \right\} \text{parse tree 1}$$

$$2. \left. \begin{array}{l} s \rightarrow \text{Pronoun, VP} \\ \text{VP} \rightarrow \text{Verb, NP} \\ \text{NP} \rightarrow \text{Det, Noun, PP} \end{array} \right\} \text{parse tree 2}$$

When the sentence can be parsed in two ways, there will be two possible meanings of the sentence. The first parsing could be interpreted as the "the person put the baby who is located at some place into the pen" and the second parsing could be interpreted as "the baby is already in the pen, and the person put him/her into some place". To extract semantic from the sentence, the processor should be able to determine the most possible meaning.

To resolve the problem, the processor refers to the previous related sentence and uses its semantic to determine most possible meaning of the current sentence. As example, the preceding sentence of the sentence "I put the baby in the pen" is "A baby is left alone on the floor". By using the knowledge about the most relevant preceding sentence, a possible value ( $\sigma$ ) is attached to the derived production rules. Thus the production rule of grammar can be represented as  $\alpha \xrightarrow{\sigma} \beta$  where  $\sigma$  is a *plausibility function* in each grammar rule, and  $\sigma \in [0,1]$  indicates the plausibility for substituting  $\alpha$  with  $\beta$  in a parsing process. A string  $S$  of symbols in  $V_T$  is said to be in the language  $L(G)$  if and only if  $s \rightarrow S$ , i.e.  $S$  is derivable from  $s$ . When  $Tr$  is a parse tree generating  $S$ , the plausibility of  $Tr$  is

$$\min\{\mu(s \rightarrow \alpha_1), \dots, (\alpha_n \rightarrow S)\} > 0 \quad (4)$$

where  $s \rightarrow \alpha_1, \alpha_1 \rightarrow \alpha_2, \dots, \alpha_m \rightarrow S$  is the derivation chain from which  $Tr$  is constructed, and  $\mu(\alpha_i \rightarrow \alpha_{i+1})$  is the non-zero  $\sigma_{(i+1)}$ . The restricting fuzzy set  $F_s$  is defined as

$$F_s = \{Tr\} \quad (5)$$

and its membership function is

$$\mu_{F_s}(Tr) = \left\{ \begin{array}{ll} \min\{\mu(s \rightarrow \alpha_1), \dots, (\alpha_m \rightarrow S)\} & \text{if } s \rightarrow Tr S \\ 0 & \text{otherwise} \end{array} \right\}$$

where  $\rightarrow_{Tr}$  is the chain  $s \rightarrow \alpha_1, \alpha_1 \rightarrow \alpha_2, \dots, \alpha_m \rightarrow S$  from which  $Tr$  is constructed. When a sentence is ambiguous, the fuzzy max operator is used to select the most possible parse tree, which is formulated in Eq. (6).

$$\mu F_s(G_{Tr}) = \{\max(Tr_1, \dots, Tr_n)\} \quad (6)$$

Semantically, the sentence “I put the baby in the pen” is resolved to the meaning “the person put the baby who is located at some where into a pen”.

### 3.3. Step 3

For further computation, predicate calculus is used for semantic representations. For example, the semantic for a sentence “I put the baby in the pen” is represented in the form of `put (baby, pen)`.

## 4. Implementation Issues

The proposed framework has been implemented in C language. Dynamic programming technique has been used to create a parser for syntactic processing, where [16] has been applied. The semantic attachment has been conducted by using *lambda reduction technique* [17]. In this work, seventy fuzzy grammar rules have been used. Fifteen data sets have been used for the framework. Each data set consists of ambiguous and unambiguous sentences. Each sentence may contain ambiguous and unambiguous words. The length of data set is between five to seven sentences. The process is conducted at a sentence level. The obtained results have been compared to human judgment, and the results indicate the proposed framework is successful.

## 5. Summary

This paper proposes a new framework for extracting semantics from texts. The novelty of this framework is the knowledge of about subject and the most relevant preceding sentence have been used to resolve ambiguity in extracting semantics for information extraction. Possibility theory and fuzzy sets have been used to extract the most possible semantics from the texts based on the knowledge about subject and preceding sentence. Experimental results indicate that the proposed framework is successful.

## 6. Acknowledgment

This research is supported by the Ministry of Higher Education, Malaysia for FRGS grant.

## 7. References

- [1] A. McCallum, (2005), Information extraction: distilling structured data from unstructured text, *Queue* 3(9), 2005, pp. 48–57.
- [2] H. M. AlFawareh, S. Jusoh and W. R. S. Osman (2008). Ambiguity in text mining. In: *Proceedings of the International Conference on Computer and Communication Engineering (ICCCCE'08)*. 2008, pp. 1172–1176.
- [3] J. Redfean, Text mining, JISC, 2006, pp. 1–2.
- [4] R. Grishman, Information extraction: Techniques and challenges. In: *SCIE*, 1997, pp. 10–27.
- [5] H. Karanikas, C. Tjortjis, and B. Theodoulidis, An approach to text mining using information extraction. In: *Workshop of Knowledge Management: Theory and Applications in Principles of Data Mining and Knowledge Discovery 4<sup>th</sup> European Conference*, 2000.
- [6] R. Rao, From unstructured data to actionable intelligence. *IEEE Computer Society*, 2003.
- [7] N. Singh, The use of syntactic structure in relationship extraction. *Master's thesis*. Department of Electrical Engineering and Computer Science, MIT, 2004.
- [8] R. Hale, Text mining: Getting more value from literature resources. *Drug Discovery Today* 10(6), 2005, 377–379.
- [9] C. Nédellec and A. Nazarenko, Ontologies and information extraction: A necessary symbiosis. In: *Ontology Learning from Text: Methods, Evaluation and Applications* (P. Buitelaar, P. Comiano and B. Magnini, Eds.), IOS Press Publication, 2005.
- [10] R. Malik, CONAN: Text Mining in Biomedical domain. PhD thesis. Utrecht University, Austria, 2006.
- [11] Q. Liu, H. Yu, X. Cheng, and S. Bai, Chinese Named Entity Recognition Using Role Model, *Computational Linguistics and Chinese Language Processing*, Vol. 8, No. 2, 2003, pp. 1-31
- [12] S. Soderland, Learning information extraction rules for semi-structured and free text. *Machine Learning* 34, 1999, pp. 233–272.
- [13] M.E. Calliff and R.J. Mooney, Relational learning of pattern-match rules for information extraction, In: *Proceedings of the Sixteenth National Conference on Artificial Intelligence*, 1999, pp. 328–334.
- [14] O. Etzioni, M. Cafarella, S. Downey, A-M. Kok, T. Popescu, T. S. Shaked, D.S. Soderland, A. Weld and A. Yates, Web-scale information extraction in know it all. In: *Proceedings of the Thirteenth International World Wide Web Conference*, 2004.
- [15] R. Yangarber, R. Grishman, P. Tapanainen and S. Huttunen, Automatic acquisition of domain knowledge for information extraction. In: *Proceedings of the 18th International Conference on Computational Linguistics (COLING 2000)*, 2000, pp. 940-946.
- [16] J. Earley, An efficient context-free parsing algorithm. *Communication of the ACM*, Vol. 13, No. 2, 1970, pp. 94-102.
- [17] D. Jurafsky and J. H. Martin, *Speech and Language Processing*, United States of America, Prentice-Hall, 2000.