

Title	Automatic Extraction of Named Entity Related Relations for Searching
Author(s)	Nguyen, Thanh Tri
Citation	
Issue Date	2008-03-04
Type	Conference Paper
Text version	publisher
URL	<a href="http://hdl.handle.net/10119/8226">http://hdl.handle.net/10119/8226</a>
Rights	
Description	JAIST 21世紀COEシンポジウム2008「検証進化可能電子社会」= JAIST 21st Century COE Symposium 2008 Verifiable and Evolvable e-Society, 開催：2008年3月3日～4日, 開催場所：北陸先端科学技術大学院大学, GRP研究員発表会 セッションA-2発表資料

# Automatic Extraction of Named Entity Related Relations for Searching

Nguyen Thanh Tri  
 School of Information Science  
 t-thanh@jaist.ac.jp

## 1. The aim of research

Named entities (NEs) play an important role in many Natural Language Processing applications including semantic search. Discovering the NE-related relations may be beneficial to these applications. Our study proposes a to extract relations of NEs in the form of  $\langle ne, category, related-to, object \rangle$  quadruples which describe that the named entity *ne* ISA *category*, and the *category* IS-RELATED-TO *object*. We extended the Person Category Extraction (PCE) algorithm to extract these tuples, and experiments on Wall Street Journal (WSJ) corpus give promising results.

## 2. The approach

Text documents often contain valuable relations of entities. For example, in the sentence taken from the WSJ corpus: *There's a generally more positive attitude toward the economy, said Bette Raptapoulos, analyst for Prudential-Bache Securities Inc., ...* (1) there are relations: "Bette Raptapoulos" is-a analyst, and analyst for "Prudential-Bache Securities Inc." Such relations may be beneficial in many NLP applications, such as for answering *Who* and *List* questions, e.g., "Who is Bette Raptapoulos?" or "Give me the list of analyst for Prudential-Bache Securities Inc."

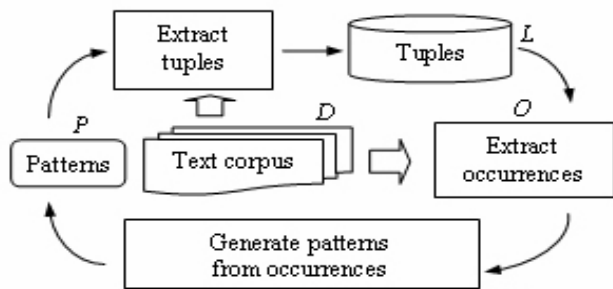


Figure 1. PCE's model

## Person Category Extraction:

The purpose of NECE is to extract  $\langle named\_entity, category \rangle$  tuples, in which *category* is the fine-grained categories of *named\_entity*, so the set of named entity classes can be expanded by automatically extracting from texts. When we extract the tuple  $\langle \text{"Bette Raptapoulos"}, \text{"analyst"} \rangle$  from (1), we only have information: "Bette Raptapoulos" ISA "analyst". If we can extract the relation: "analyst" for "Prudential-Bache Securities Inc.", we will have complete information about "Bette Raptapoulos". The PCE algorithm is depicted in Figure 1. Starting with two seed patterns, PCE extracted  $\langle person, category \rangle$  tuples. The extracted tuples were used to extract the occurrences of  $\langle person, category \rangle$  tuples in texts for generating new patterns. Again, new patterns were used to extract new tuples. The process terminated when no more patterns were produced. A pattern is defined as a 4-tuple:  $(order, person\_slot, middle, category\_pattern)$ , where *order* indicates the occurrence order of person and category in a sentence; *person\_slot* is a slot which will be replaced with a person named entity; *middle* is the string surrounded by person and category; *category\_pattern* is defined as: *category\_pattern* := *noun\_phrase1* (and *noun\_phrase2*)? where *noun\_phrase<sub>i</sub>* is a regular expression that matches a noun phrase with added POS tags. In order to extract new  $\langle person, category \rangle$  tuples, for every sentence *s*, from a pattern, and for each person NE *named\_entity* in *s*, we construct a regular expression:

$* named\_entity middle category\_pattern *$

If *s* matches the above regular expression, the  $\langle person, category \rangle$  tuple is extracted. For improving the performance of the algorithm, at each match if the *category* is not valid (i.e., the *category* is not a type of person) the match is ignored.

An occurrence of a person, category tuple is defined as a 4-tuple:  $\langle order, person, middle, category \rangle$ , where *middle* is a string surrounded by person and category. An occurrence of a  $\langle person, category \rangle$  tuple is extracted if a sentence *s* matches the regular expression:  $* person middle category *$  or  $* category middle person *$  After having been extracted, occurrences are used to generate new patterns. A *middle* of an occurrence is not necessarily reliable,

some constraints were proposed to retain reliable ones: Repetition of a middle ( $\text{repetition}(\text{middle})$ ) is the number of times the middle appears between the person and category of  $\langle \text{person}, \text{category} \rangle$  tuples of same person. Diversity of a middle ( $\text{diversity}(\text{middle})$ ) is the number of times the middle appears between the person and category of  $\langle \text{person}, \text{category} \rangle$  tuples of different persons. A middle that has  $\text{repetition}(\text{middle}) > \text{threshold}_R$  seems reliable and is kept. A pattern seems specific if it is generated based on tuples of a person, so only middles that have  $\text{diversity}(\text{middle}) > \text{threshold}_D$  are kept to make the generated patterns general (Condition 1). If a middle contains a verb phrase, the verb phrase should express the relation person ISA category (Condition 2).

### Named-Entity-Related Relations Extraction:

We extend the PCE to extract  $\langle \text{named\_entity}, \text{category}, \text{related-to}, \text{object} \rangle$  quadruples describing the relations: *named\_entity* ISA *category* (or ISA relations for short), and *category* related-to *object* (or *related-to* relations for short). From our observations, the related-to relations can be expressed in the following ways:

- a) The *category* and *object* are linked by a preposition: “*category preposition object*”, e.g., “analyst for Prudential-Bache Securities Inc.”
- b) The *category* and *object* are connected by a possessive apostrophe, e.g., “Semi-Tech’s chief executive officer”. This can be interpreted as “category of object”, e.g., “chief executive officer of Semi-Tech”.
- c) The *object* and *named\_entity* are linked by a preposition, e.g., “. . . said economist David Littmann of Manufacturers National Bank . . .”, from which an expected quadruple is  $\langle \text{“David Littmann”}, \text{“economist”}, \text{“of”}, \text{“Manufacturers National Bank”} \rangle$ .
- d) The *object* is embedded in *category*, e.g., “IBM president”. This can also be interpreted as “category of *object*”, e.g., “president of IBM”.
- e) The *related-to* relation is implicitly expressed, e.g., “Mr. Baird, who heads the Manhattan U.S. attorney’s securities-fraud unit, denied the quote . . .”, from which an expected quadruple is  $\langle \text{“Baird”}, \text{“head”}, \text{“of”}, \text{“securities-fraud unit”} \rangle$ .

Since case e) does not have fixed expressions, we do not treat such cases. In case d), because *object* is already embedded in *category*, we do not need to extract the *object*. For cases a), b) and c), we construct regular expressions to extract the *object* and *related-to*.

### Utilization of Named-Entity-Related Relations for Semantic Search:

The extracted  $(\text{named\_entity}, \text{category}, \text{related-to}, \text{object})$  quadruples are valuable for NLP applications. In this section, we use them for answering some types of questions. If *named\_entity* in a quadruple is a person, the quadruple helps answer the query: “Who is *named\_entity*?”, e.g., “Who is Bette Raptapoulos?” If *named\_entity* is of another type, such as *organization* or *location*, the quadruple helps answer the query: “What is *named\_entity*?”, e.g., “What is IBM?” For answering the question, we just search for quadruples having the same *named\_entity* as that of the question. If a quadruple is found, then the answer is: “*named\_entity* is a(n) *category* related-to *object*”.

The extracted quadruples also help answer *list* questions, e.g., “Give me the list of analyst for Prudential-Bache Securities Inc.” The general form of this question type is “Give me the list of *category* [*related-to object*]”, where the part in square brackets is optional. For answering this question type, we search for the list *L* of quadruples having the same *category*, [*related-to*, and *object*] as those of the question. The answer is the list of *named\_entity* of quadruples in *L*.

## 3 The progress of the search

Our work was published in the “Proceedings of the 21st Pacific Asia Conference on Language, Information and Computation” (PACLIC21).

## 4 Future direction

There are still other named-entity-related relations in documents, and discovering such relations may help answer other question types. In the future, we intend to study to extract other relations in documents for searching.

## 5 Publication

1. **T. T. Nguyen** and A. Shimazu, “Automatic Extraction of the Fine Category of Person Named Entities from Text Corpora”, *IEICE Transactions on Information and Systems, Special section on Knowledge, Information and Creativity Support System*, Vol. E90-D, No. 10, pp. 1542-1549, 2007.
2. **T. T. Nguyen** and A. Shimazu, “Acquisition of Named-Entity-Related Relations for Searching”, *Proceedings of the 21st Pacific Asia Conference on Language, Information and Computation*, (PACLIC21), pp. 349-357, 2007.
3. **T. T. Nguyen**, L. M. Nguyen, and A. Shimazu, “Using Semi-supervised Learning for Question Classification”, *Journal of Natural Language Processing*, Vol. 15, No. 1, 2008. (to appear)