# Using Search Session Context for Named Entity Recognition in Query

Junwu Du[1], Zhimin Zhang[2], Jun Yan[2], Yan Cui[1], Zheng Chen[2]

[1]School of Computer Science
Beijing Institute of Technology
Beijing, 100081, P.R. China

du.junwu@gmail.com; cui.yan@live.cn

[2]Microsoft Research Asia
Sigma Center, 49 Zhichun Road
Beijing, 100080, P.R. China

{zhzha, junyan, zhengc}@microsoft.com

## ABSTRACT

Recently, the problem of Named Entity Recognition in Query (NERQ) is attracting increasingly attention in the field of information retrieval. However, the lack of context information in short queries makes some classical named entity recognition (NER) algorithms fail. In this paper, we propose to utilize the search session information before a query as its context to address this limitation. We propose to improve two classical NER solutions by utilizing the search session context, which are known as Conditional Random Field (CRF) based solution and Topic Model based solution respectively. In both approaches, the relationship between current focused query and previous queries in the same session are used to extract novel context aware features. Experimental results on real user search session data show that the NERQ algorithms using search session context performs significantly better than the algorithms using only information of the short queries.

## Categories and Subject Descriptors

H.3.3 [**Information Storage and Retrieval**]: Information Search and Retrieval – *search process*

## General Terms

Algorithms, Experimentation

## Keywords

Search Session, Named Entity Recognition, CRF, Topic Model

## 1. INTRODUCTION

Nowadays, the Named Entity Recognition problem in Query (NERQ) has attracted increasingly attention from information retrieval community since it can be used to fix a number of failure cases in the relevance based search engines. Figure 1 gives an example of failure case in a commonly used commercial search engine. If we can identify "*Eagle Vision*" as a named entity in query, this irrelevant result will not be returned to end users with high rank. To our best knowledge, there are only a few previous studies that have tried to recognize named entities in queries. As one of the earliest related works, Guo et al [3] has tried to use a topic model to identify named entities in queries and they showed that around 70% of the real search queries contain named entities. However, since the queries are always very short (i.e., 2-3 words on average) and many of them do not satisfy the natural language grammar [3], both the classical Named Entity Recognition (NER)

algorithms [2] and the algorithms particularly designed for query [3] may fail when the queries are short with limited context.



**Figure 1: A failure case of a commercial search engine**

In this paper we propose to address the NERQ problem by utilizing the context information in search sessions, once used to perform QC [1]. We extract two new types of features from the search sessions, namely *class feature* and *overlap feature* to help performing NERQ.

## 2. CONTEXT AWARE NAMED ENTITY RECOGNITION

In this Section, we propose two novel search session features for NERQ. Suppose the predefined classes are $C = \{c_1, c_2, \dots\}$ and the queries, which are in the same search session, are $Q = \{q_1, q_2 \dots\}$. Let the words, which make up of a query $q_i$ be $W_i = \{w_{i1}, w_{i2}, \dots\}$.

### 2.1 Class Feature and Overlap Feature

In a search session, a sequence of queries can more accurately identify the class than just the current focused query.

**Definition-1:** the class feature extracted from queries before the current focused query $q_i(i \geq 1)$ in the same search session is defined as $F_{class} = \{f_1, f_2, \dots f_k, \dots f_{|C|}\}$, and the $f_k$ is defined as:

$$f_k = \begin{cases} \sum_{j=1}^{i-1} e^{j-i+1} \Pr(c_k|q_j) & , i > 1 \\ 0 & , i \leq 1 \end{cases} \quad (1)$$

**Definition-2:** the overlap feature extracted from the current focused query $q_i(i \geq 1)$ and its previous queries in the same search session is based on single word. This feature is defined as $F_{overlap} = \{g_1, g_2, \dots g_k, \dots g_{|W_i|}\}$, and $g_k$ is defined as:

$$g_k = \begin{cases} \sum_{j=1}^{i-1} e^{j-i+1} O_{i,j,k} & , i > 1 \\ 0 & , i \leq 1 \end{cases} \quad (2)$$

In Eqn. (2), $O_{i,j,k}$ is defined as:

$$O_{i,j,k} = \begin{cases} 1 & , w_{ik} \in W_j \\ 0 & , others \end{cases} \quad (3)$$

## 2.2 Applications of New Features

In the CRF [4] based method, two new types of search session features are used by the classifier. Following are the details of the new features:

- Class feature: using every element of $F_{class}$ in Definition-1 as a CRF feature.
- Overlap feature: using every element of $F_{overlap}$ in Definition-2 as a CRF feature.

In the Topic Model based method, a single-named-entity query $q$ is represented as triples$(e, t, c)$, where $e$ denotes named entity, $t$ denotes the context of $e$ in $q$, and $c$ denotes the class of $e$ [3]. In our experiment, for simplicity, only one query before the current focused query is considered. In online prediction, we try to segment both the current focused query and the previous query in all possible ways together, and the triples with the highest $P_{QueryPair}$ value are output results for NERQ. $P_{QueryPair}$ is defined as:

$$P_{QueryPair} = \mathrm{Pr}_{pre}(e, t, c) + \mathrm{Pr}_{cur}(e, t, c) + P_{class} + P_{overlap}(4)$$

where $\mathrm{Pr}_{pre}(e, t, c)$ and $\mathrm{Pr}_{cur}(e, t, c)$ are the joint probabilities with respect to previous query and current query in the same search session. $P_{class}$ represents the class distribution similarity between the two sequential queries, which is defined as :

$$P_{class} = \frac{1}{|C|} \sum_{k=1}^{|C|} \mathrm{Pr}(c_k | e_{pre}) \mathrm{Pr}(c_k | e_{cur}) \qquad (5)$$

In Eqn. (5), $|C|$ is the number of predefined classes. $P_{overlap}$ represents the effectiveness of the overlap feature, which is defined as:

$$P_{overlap} = \frac{\lambda}{l} \sum_{k=\alpha}^{\alpha+l-1} g_k \qquad (6)$$

In Eqn. (6), $l$ is number of words in $e_{cur}$, $\alpha$ is the entity's first word's index in the current focused query, and $g_k$ is defined in Eqn. (2). $\lambda$ is the weight parameter of overlap feature, and in our experiment we define $\lambda$ as $\frac{\mathrm{Pr}_{cur}(e,t,c)}{|C|}$.

## 3. EXPERIMENTS

In the experiments of this work, we only use the car model domain queries for demonstration. We have two sets of experiments by CRF and topic model respectively to compare the NERQ results with and without context information in search sessions.

In the experiment of using CRF model for NERQ, we start with a small car name dictionary which contains 1,663 standard car names to build a dataset for evaluation. Using the car name dictionary, we firstly scan one month's click-through log data of a commercial web search engine. If a query which contains car names in the car name dictionary is found by exact match, then all the queries in the same search session are added to our dataset with the query order in session unchanged. In this experiment, we defined 5 named entity (NE) tags, namely 0 (Not entity), 1 (Single word entity), 2 (Beginning of the entity), 3 (Middle of the entity), 4 (Ending of the entity). We have asked 10 human labelers to manually label car domain named entities in queries of our dataset. From the labeled data, 5000 queries and their corresponding search sessions are used as training set and 1000 queries in remaining search sessions are used as testing set. In this

experiment, the CRF classifier without considering the search session context is used as baseline, which is represented by CRF. On the other hand, our approach is the CRF algorithm that uses our proposed session level context features. We name it as the Context aware CRF (CCRF). The evaluation metrics used for the two classifiers are precision, recall and F-measure. Results of evaluation comparison of the two classifiers are shown in Table 1.

**Table 1: With/without session context by CRF classifiers**

| NE tags | Precision | | Recall | | F1 | |
|---|---|---|---|---|---|---|
| | CRF | CCRF | CRF | CCRF | CRF | CCRF |
| 0 | 45.45% | 49.35% | 57.38% | 69.09% | 50.72 | 57.58 |
| 1 | 62.69% | 79.10% | 66.67% | 71.62% | 64.62 | 75.18 |
| 2 | 85.71% | 83.93% | 63.16% | 67.14% | 72.73 | 74.60 |
| 3 | 47.37% | 63.16% | 47.37% | 50.00% | 47.37 | 55.81 |
| 4 | 69.64% | 71.43% | 69.64% | 76.92% | 69.64 | 74.07 |
| **Overall** | **62.91%** | **69.09%** | **62.91%** | **69.09%** | **62.91** | **69.09** |

In the experiment by using topic model, 160 car names are selected as seed named entities. We constructed the topic model in the same way with the work [3] using the seeds and search engine log. 1,000 query pairs in the same search session were randomly selected from the dataset, generated in the experiment of using CRF, as test set. In this experiment, the Topic Model without considering search session context is used as baseline. In contrast, our approach is the Topic Model algorithm that has used the search session context features. We name it as the Context aware Topic Model (Ca Topic Model). We use two word indexes of the query to label the beginning and the ending of the named entity in the query. The evaluation metrics used for the two topic models are precision. The experiment results are shown in Table 2.

**Table 2: With/without session context by Topic Models**

| | Precision | | |
|---|---|---|---|
| | NE Beginning | NE Ending | Total NE |
| Topic Model | 65.9% | 59.3% | **54.2%** |
| Ca Topic Model | 74.8% | 65.4% | **60.6%** |

## 4. CONCLUSION

In this extended abstract, we proposed two context aware features extracted from user search sessions to tackle the problem of named entity recognition in query. We applied the two features to improve two classical NER algorithms, which are known as Conditional Random Field (CRF) and Topic Model. The experimental results on real user search session data show that the NERQ algorithms using search session context through our proposed features can perform significantly better than the classical NER algorithms using the short query information only.

## 5. REFERENCES

[1] Cao, H., Hu, D., Shen, D., Jiang, D., Sun , J., Chen, E., Yang, Q. 2009. Context-Aware Query Classification. In SIGIR '09

[2] Ekbal, A., Haque, R., Das, A., Poka, V., Bandyopadhyay, S. 2008. Language Independent Named Entity Recognition in Indian Languages. In Proc.IJCNLP-08.

[3] Guo, J., Xu, G., Cheng, X., Li, H. 2009. Named Entity Recognition in Query. In SIGIR '09, pages 268.

[4] Lafferty, J., McCallum, A., Pereira, F. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. In Proc. ICML.