

# A Simple Approach to NTCIR-10 MedNLP Task

Yuka Tateisi  
National Institute of Informatics  
yucca@nii.ac.jp

Takashi Okumura  
National Institute of Public Health  
taka@niph.go.jp

## ABSTRACT

For MedNLP complaint and diagnosis subtask we tried a simple dictionary-matching method using MeCab, and achieved 61.10% F-score in official evaluation. Our method can evaluate the coverage of terminology data in a simple and inexpensive way.

## Team Name

NIPH

## Subtasks

Complaint and diagnosis (Japanese)

## Keywords

## 1. INTRODUCTION

Recently, the volume of medical records written in electronic format is increasing. As a result, information processing technique in medical fields is in demand. Medical records contain information in natural language such as description of the medical status of patients, which needs to be processed with NLP techniques. However, few NLP resources such as ontology and electronic dictionary are publicly available in the Japanese medical community, in contrast to bio- or genome- science community. Computerized medical records systems have such resources built-in, but they differ from system to system and not for use outside the system.

In such circumstances some doctors develop their own terminology lists and dictionaries privately, which are customized to their own purposes, and as such, they may have a limited scope of coverage. Our attempt here is to use one of such resources against the NTCIR-MedNLP task [2] data to evaluate the coverage of such customized data.

We take one of such dictionaries, developed by Dr. Keijiro Torigoe and used for developing a diagnosis support system National Institute of Public Health [3]. The dictionaries, called Disease Master Data and Symptom Master Data, consist of records of about 1900 disease names and 800 symptom names respectively, with English translations and short descriptions<sup>1</sup>. We try to find complaint information by simply matching the text with these Master data, with the method described in the following section.

## 2. METHODS

<sup>1</sup>The data are available from I'ROM Holdings Co. Ltd. (<http://www.irom-hd.co.jp/>) under BSD license.

The Japanese names of diseases and symptoms are extracted from the Master Data. In addition to these, we added the names of diseases, syndromes, symptoms, deficiencies, and side effects extracted from a laboratory test manual [4] to increase the volume of candidate names. We used Japanese morphological analyzer MeCab [1] as a dictionary matcher. That is, we built a user dictionary from these data and analyzed the text with MeCab with the user dictionary. The approach taken here is very simple so that everyone can employ with minimal training. We hope that, if this simple approach is successful, non-experts in NLP can test their own resources casually and in the end help enrich public NLP resources.

### 2.1 Data

The dictionary that we built for MeCab is based on the following three datasets. The numbers of dictionary entries taken from each dataset are shown in parentheses.

- Disease Master Data: List of names of diseases and syndromes, compiled by Dr. Keijiro Torigoe (1912).
- Symptom Master Data: List of names of symptoms and findings, compiled by Dr. Keijiro Torigoe (861).
- Lab Test Data: Names of diseases, syndromes, symptoms, deficiencies, and side effects extracted from laboratory test manual (2492).

XML character entity references in these master data, representing Greek characters, are mapped to corresponding characters in JIS Zenkaku. Alphanumeric characters in the records are also mapped to corresponding characters in Zenkaku. With these master datasets combined as a user dictionary, MeCab ver 0.995 with IPA dictionary was used for matching. All the extract names were treated as common noun, with smaller costs than those in the system dictionary, to give higher preferences. All the remaining fields for the grammatical attributes were set to the default value for common nouns. The costs for those names from the Disease Master Data and the Symptom Master Data were set to 100, and the costs for names from the Lab Test Data were set to 200. These are significantly smaller than most of words in the system (IPA) dictionary, which are assigned the value in the range of 4000–6000.

To signify that the words are from the user dictionary, the “reading” field and the “pronunciation” field of the MeCab record are utilized. For each entry, the reading field was set to the name of the dataset, and the pronunciation field

