# CONSTRUCTION OF SEMANTIC LEXICON OF DERMATOLOGY

Yang Zhou[1], Nan Xiang[1], Ruixiang Wang[1], Yao Liu[2]
Xingliang Qi[1] and Zhenguo Wang[1]

[1]Institute of Chinese Medical History and Literature
Shandong University of Traditional Chinese Medicine
No. 4655, Daxue Road, Jinan 250355, P. R. China
zhenguow@163.com

[2]Institute of Scientific and Technical Information of China
No. 15, Fuxing Road, Haidian District, Beijing 100038, P. R. China
liuy@istic.ac.cn

Abstract. *This paper introduces the construction of the Semantic Lexicon of Dermatology by using the theory and technology of Natural Language Processing (NLP) which can provide the database, such as automatic semantic analysis, word sense disambiguation, for NLP. This paper analyzes some problems of the Semantic Lexicon of Dermatology including system, standard and code of classification, range of vocabulary, etc. It also introduces the main function of the editing system of the Semantic Lexicon of Dermatology.*
**Keywords:** Semantic Lexicon, Dermatology, Natural Language Processing (NLP)

1. **Introduction.** Construction of the Semantic Lexicon is the basic work of Natural Language Processing (NLP). The Semantic Lexicon occupies a key position in the study of NLP. Concerned people around the world have developed many large-scale Semantic Dictionaries. Semantic Dictionaries in medical field have also been constructing gradually [1]. With the development of Precision Medicine, the study of NLP in medical field has become a new research hot topic. Accordingly, the construction of medical professional Semantic Lexicon is imperative.

Usually, there are two ways to construct the Semantic Lexicon. One way is based on the subjective judgment of experts in concerned fields. The other way is based on Automatic Clustering [2]. Construction of the Semantic Lexicon of Dermatology combines the two ways mentioned above. On one hand, construction of the professional Semantic Lexicon must take the existing research results in the field as a reference which means that we should make full use of the existing dermatology dictionaries [3]. On the other hand, we should introduce new terms of dermatology by using Dermatology Lexicon Editing System based on corpus of works and thesis of dermatology. The Semantic Lexicon of Dermatology will be used in NLP widely.

This paper introduces the principle for construction of the Semantic Lexicon of Dermatology, the main function and framework of the Editing System of the Semantic Lexicon of Dermatology.

2. **Principles for the Construction of the Semantic Lexicon of Dermatology.**

2.1. **Principles for collecting terms.** The terms have wide coverage. It contains not only dermatology, venereology, cosmetology, dermatological surgery, but also clinical dermatology and basic medical knowledge, such as molecular biology, immunology, biochemistry, pathology, microbiology, pharmacology, and biology.

We collect new words, including disease names, drug names and other professional terms.

We choose the vocabulary with high frequency use for translating English to Chinese. If there are multiple Chinese words which express the same conceptual mean, we will get the word which is used most commonly.

2.2. **Classification of the Semantic Lexicon of Dermatology.** Skin disorders are not only the diseases of skin, but also involve many human organs which are closely related with social and natural environment. Accordingly, the terms of dermatology not only contains dermatosis itself, but also involves many related medical specialties, such as internal medicine, surgery, pharmacy, physiology and pathology. Therefore, the classification of the Semantic Lexicon of Dermatology actually covers the whole medical system. Moreover, it is also related to the non-medical factors, such as natural science, and geography.

The classification of the Semantic Lexicon of Dermatology uses MESH classification system for reference and it has been extended and modified [4-6]. In order to collect the terms maximally, we use the existing dermatology dictionary, the latest clinical knowledge and the corpus of dermatology. We collect 300,000 terms and classify them into twenty-four first-rate categories.

First-rate category contains Anatomy, Organisms, Diseases, Pathogeny, Disease Location, Chemicals and Drugs, Analytical, Diagnostic and Therapeutic Techniques and Equipment, Psychiatry and Psychology, Biological Sciences, Physical Sciences, Anthropology, Education, Sociology and Social Phenomena, Technology, Industry, Agriculture, Humanities, Information Science, Persons, Health Care, Geographicals, Institutions, and Others (Figure 1).

Each category is divided into the lower categories; for example, the 'Pathogeny' category is divided into categories of Acute Geographical Factors, Neural Factors, Aggravated Factors, etc. And the category 'Aggravated Factors' is divided into categories of Insolation Weathering, Improper Diet, Improper Use of Drugs, etc. The lowest category is the fourth level directory.

```
Anatomy, an
Organisms, or
Diseases, di
        Skin and Connective Tissue Diseases, sctd
        Nutritional and Metabolic Diseases, nmd
        Endocrine Diseases, edi
        Immunologic Diseases, id
        Injuries，Poisonings，Occupational Diseases, ipod
         Animal Diseases, ad
        Symptoms and General Pathology, sgp
        ......
Chemicals and Drugs, cd
Analytical，Diagnostic and Therapeutic Techniques and Equipment, adtte
Psychiatry and Psychology, pp
Biological Sciences, bs
Physical Sciences, ps
Anthropology,Education,Sociology and Social Phenomena, aessp
Technology，Industry，Agriculture, tia
Humanities, hu
Information science, is
Persons, pe
Health Care, hc
Geographicals, ge
Pathogeny, pa
Disease Location, dl
Institutions, inst
Other, oth
```

FIGURE 1. Classification of the Semantic Lexicon of Dermatology

2.3. **Principles for marking terms.** We use the first letter of English words as the tab of the term. For example, we mark 'biological sciences' as 'bs'.

If there is only one word, we use the first two letters, such as marking 'diseases' as 'di'.

If two tabs are the same, we use the first three letters, such as marking "genetics" as "gen". If there is still a repetition, then we use the Arabia number to distinguish (Figure 2).

Biological Sciences生物科学,bs
Genetics遗传学,gen
aberrant不正常的(异常的), abe1
Aberration畸变(失常), abe2
abnormal异常的, abn1
abnormality异常(失常、畸形、畸变), abn2

FIGURE 2. An example of marking terms



FIGURE 3. System framework

3. **Editing System of the Lexicon of Dermatology.** The system is based on the B/S architecture of the Internet. It can be accessed by the browser without the client. Professionals can use the system to modify and check the lexicon according to the reference statements and their own experience, and it supports multi-collaborative editing. It can also derive lexicon and word attributes.

There are three modules in the system structure. They are presentation layer, application layer and data layer. This system based on database is developed by using web application framework (based on Strusts2 + Hibernate + j2ee + Spring) (Figures 3 and 4). The function of 'User Management' module includes registration, login, modifying password. This part will not be detailedly described.

3.1. **Proofreading and editing terms.** It can add, delete, modify, query and upload the new domain lexicon. We can get the most relevant information from Internet sites in accordance with the data itself by using the system. There are four functions in the system of proofreading and editing terms, Batch Input, Terms Retrieval, Terms Editing, Terms Outputting (Figure 5).
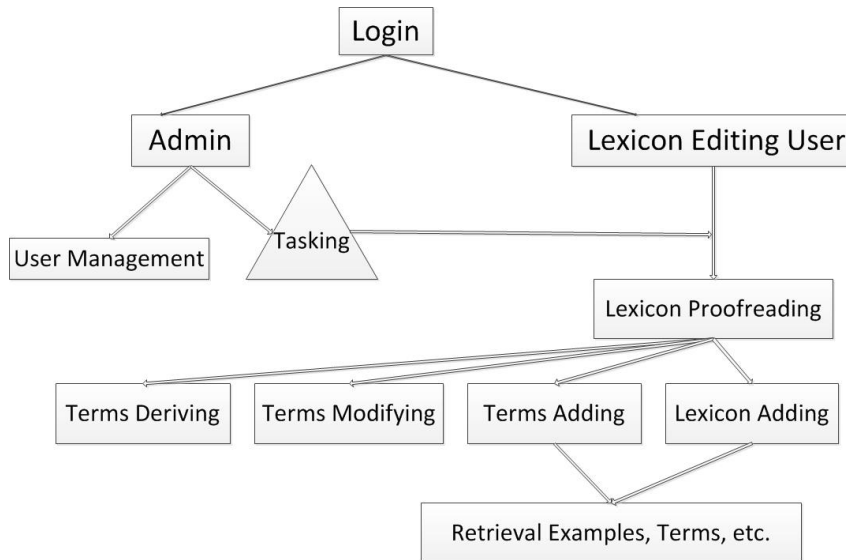
FIGURE 4. The procession of the system



FIGURE 5. Proofreading and editing terms

3.1.1. *Batch input.* It can import the data into the database by using batch processing method. Select the TXT or Excel file, and click 'Upload' to import the data into the database. Then, it shows a word in a row in the TXT file, and it indicates word, English name, definition, another name in columns in the Excel file.

3.1.2. *Terms retrieval.* This system provides two search ways and search functions that base on combination of multiple choices. It is word characteristic retrieval, keywords searching and domain retrieval.

Users choose the characteristic of words from the word character tabulations in the search area. The system will provide the relative words tabulations of the word characters selected. We input keywords in the relative input box, and click 'search' button. We will get the word character selected and words information which contains the key words. The key words should be fuzzy matching which means the words contain the key words.

If we check the 'yes' choice of whether searching entry that has no match with upper and lower entry only, the query range will be limited in the entry that has no match with upper and lower entry.

3.1.3. *Terms editing.* It can edit the related properties of terms, such as Chinese names. It can also add new terms with name, characteristic, related words, etc. At the same time, Chinese name, English name, another name and examples will be added into the corresponding properties by using retrieval module automatically.

3.1.4. *Terms outputting.* This system provides three output formats of terms (Figures 6-8). One format is to export the words list of the current retrieval conditions to TXT. The export format is to display a word in line; the default file name is 'user ID-current time'. Another way is to export the words list of the current retrieval conditions to an Excel table, the table head for the entry name, Chinese name, English name, etc. The default file name is 'current time-user ID-property'. The last format is to export structured lexicon.
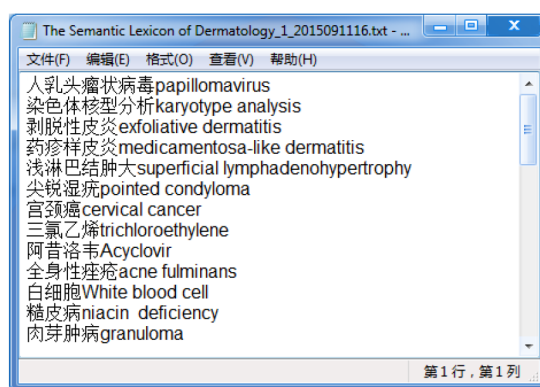


FIGURE 6. Format1 – TXT



FIGURE 7. Format2 – Excel table

3.2. **Displaying details of upload file.** The file of lexicon type can create hyponymy. The system can save it automatically and create hyponymy and synonyms.

The file of corpus type can analyze the corpus documents. The system resolves the corpus document to a corresponding lexicon file of Excel format. Then save it and create hyponymy automatically.

3.3. **Editing lexicon tree structure.** The lexicon tree structure can be modified by adding, deleting, modifying, sequencing, and checking. We can also drag the nodes to other nodes directly.
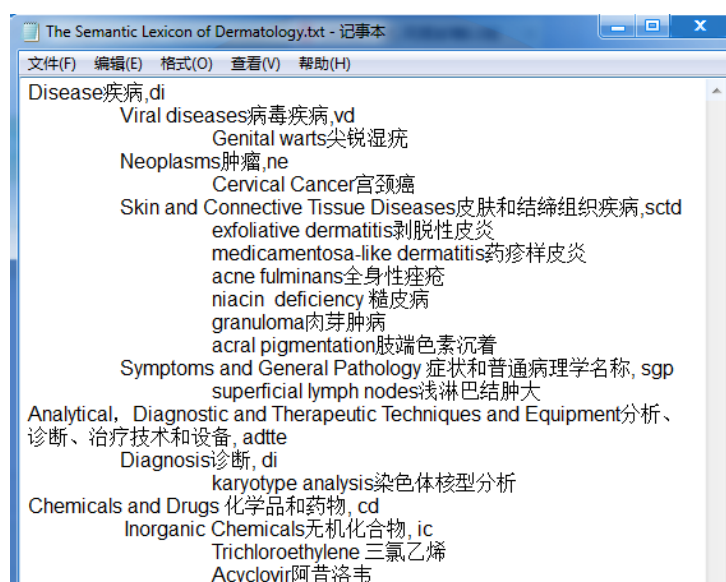
FIGURE 8. Format3 – Structured lexicon

4. **Application of the Semantic Lexicon of Dermatology.** Semantic Lexicon of Dermatology is a mechanical dictionary in the professional field which serves for the language information processing. We constructed lexical category system for information processing based on characteristics of skin diseases in the lexicon. It fulfills the words classification, and describes the different word attributes by classification to achieve the structuring in a way. It will excavate implicit knowledge from massive text. Semantic analysis plays an important role in the NLP. Semantic Lexicon of Dermatology will provide comprehensive semantic knowledge for the sentence meaning analysis, lexical ambiguity elimination which will improve accuracy of interpretation of documents.

5. **Conclusion.** In this paper, firstly, we obtain the terms from the textbook of dermatology and the clinical literature. Secondly, we classify the terms of dermatology by using MESH classification system. Finally, we introduce the Editing System of the Lexicon of Dermatology.

With the development of dermatology, there will be many new professional vocabularies. In the future work, we will continue to collect vocabulary by using the Editing System of the Lexicon of Dermatology, so as to optimize the Dermatology Lexicon.

**REFERENCES**

[1] P. Ruch, R. H. Baud, A. M. Rassinoux, P. Bouillon and A. G. Robert, Medical document anonymization with a semantic lexicon, *Proc. of AMIA Symp.*, vol.7, no.1, pp.729-733, 2000.
[2] T. Chen and M. Sun, Automated construction of Chinese thesaurus based on self-organizing map, *Journal of the China Society for Scientific and Technical Information*, vol.26, no.1, pp.77-83, 2007.
[3] W. Xu, *English Chinese Dictionary of Dermatovenereology*, Phoenix Science Press, 2007.
[4] *https://www.nlm.nih.gov/class//terms_cond.html [EB]*, 2014.
[5] B. Zhao, *China Clinical Dermatology*, Phoenix Science Press, 2010.
[6] X. Zhang, *Textbook Series of New Century – Dermatovenereology*, People's Medical Publishing House, 2002.