# Implementation of a Web Search Engine for Restaurants using Lucene

**Sujoung Oh[1], Minsoo Lee[2]**

[1]Computer Science and Engineering, Ewha Womans University, Korea, crystal7862@gmail.com
[2]Computer Science and Engineering, Ewha Womans University, Korea, mlee@ewha.ac.kr

## ABSTRACT

Recently, many people use Internet blogs or Website related to famous restaurants. People are interested in information found from various sources and want to efficiently search them. These search engines should be easy to develop and effectively enhanced using open source search engines to provide for various international language specific techniques to be experimented for search. Therefore, we developed a search engine to find popular restaurants based on Lucene and used information retrieval techniques. The proposed search engine can extract restaurant lists related to keywords. The search engine, also, provides functions to be redirected to the original restaurant website to directly verify the information found.

**Key words :** Lucene, search engine, restaurant, information retrieval

## 1. INTRODUCTION

Many people visit websites providing restaurant information to find a restaurant which serves delicious food. Their search keywords are usually the menu or the location. As such there is a need for various search mechanisms for various user inputs on various data for restaurants which is a very popular topic on the web. However search engines have limited capability in terms of the data that they use or the queries that they can support. Thus various techniques need to be experimented. This is also true in an international environment where different requirements exist for specific languages and countries. These search engines should be easy to develop and effectively enhanced using open source search engines to provide for various international language specific techniques to be experimented for search. For this reason we built a general web search engine using Lucene. And if a user wants to know more information about the restaurant, the user can connect to the website.

This paper is organized as follows. Section 2 introduces related works and section 3 describes the methodology proposed in this paper and the implementation results. And finally in section 4 we describe the conclusions and future work.

## 2. RELATED WORKS

### 2.1 Lucene

Apache Lucene is a free and open source information retrieval software library, originally created in Java by Doug Cutting. It is supported by Apache Software Foundation and is released under the Apache Software License. Lucene has been ported to other programming languages including Delphi, Perl, C#, C++, Python, Ruby, and PHP. While suitable for any application which requires full text indexing and searching capability, Lucene has been widely recognized for its utility in the implementation of Internet search engines and local, single-site searching. At the core of Lucene's logical architecture is the idea of a document containing fields of text. This flexibility allows Lucene's API to be independent of the file format. Text from PDFs, HTML, Microsoft Word, and OpenDocument documents, as well as many others (except images), can all be indexed as long as their textual information can be extracted[1,2].

### 2.2 Web Search Engines

A web search engine is a software system for searching information on the World Wide Web. The search results are presented in a line of results often referred to as SERPs(Search Engine result pages). The information can be web pages, images, and other type of files. And search engines can be used for mining data in databases[3].

A search engine operates via the following three steps. The first step is Web crawling. The second step is indexing. The third step is searching. A web crawler is an Internet bot that browses the World Wide Web, typically for the purpose of Web indexing[4]. Search engine indexing collects, parses, and stores data to facilitate fast and accurate information retrieval. Index design incorporates concepts from linguistics, cognitive psychology, mathematics, informatics, and computer science. An alternate name for the process in the context of search engines designed to find web pages on the Internet is web indexing[5]. In the search step, a web search query that a user enters into a web search engine to satisfy his or her information needs is issued. The web search queries are distinctive in that they are often plain text or hypertext with optional search-directives[6].

## 3. SYSTEM DEVELOPMENT

The system development is divided into three parts. The first part is about extracting the text from the Web site URL. The second part is about creating the index file using the information of the web site. And finally, we describe the searching for the results.

### 3.1 Extracting Text using URLs

We extract the text from web documents using the URL information. When we save the text file with the needed information, the format we use is the Json format to enable Json parsing.



**Figure 1: Restaurant name & URL**

Figure 1 illustrates the information that contains the restaurant name and the URL of the web site. The Web crawler can continuously traverse among the URLs to create a larger subset of relevant websites. From the web sites, we can extract restaurant data such as menu, location, and client comment.

Figure 2 presents the list that contains the restaurant data. Each text file indicates the data for one restaurant.

### 3.2 Index Creation

We read the text files, and then create the index for the restaurant data. It contains the restaurant name, main text, and URL. Figure 3 illustrates the files in the index folder. Lucene is used for the creation of index structures.
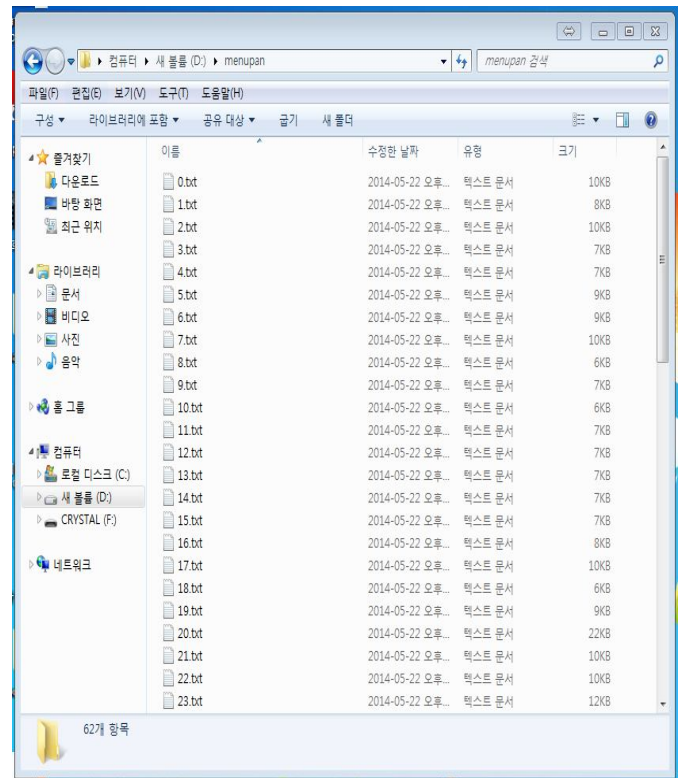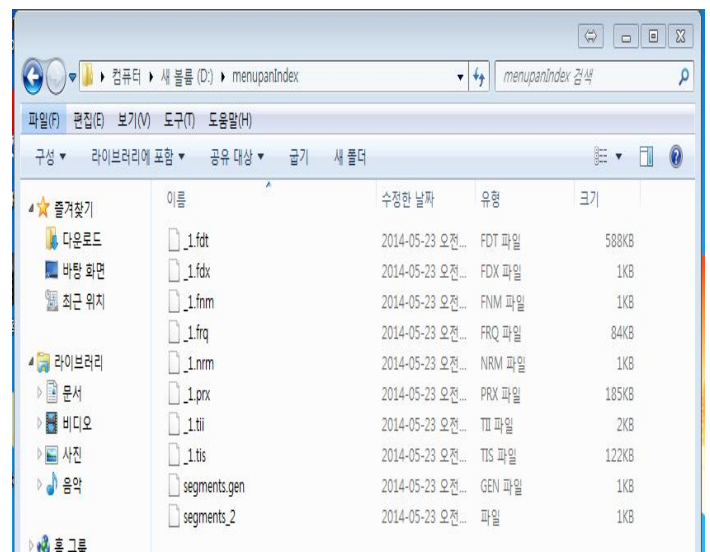


**Figure 2: Extracted restaurant data files**



**Figure 3: Index file creation**

### 3.3 Search

Our engine is based on Lucene and will return the search results based on the keyword. The results include the input query, restaurant lists about the input query, address of the restaurant, and URL of the restaurant.

**Figure 4: Search results**

Figure 4 shows the search results. In this example, we used 'sushi' as the keyword. After searching, we got three restaurant data related to sushi.
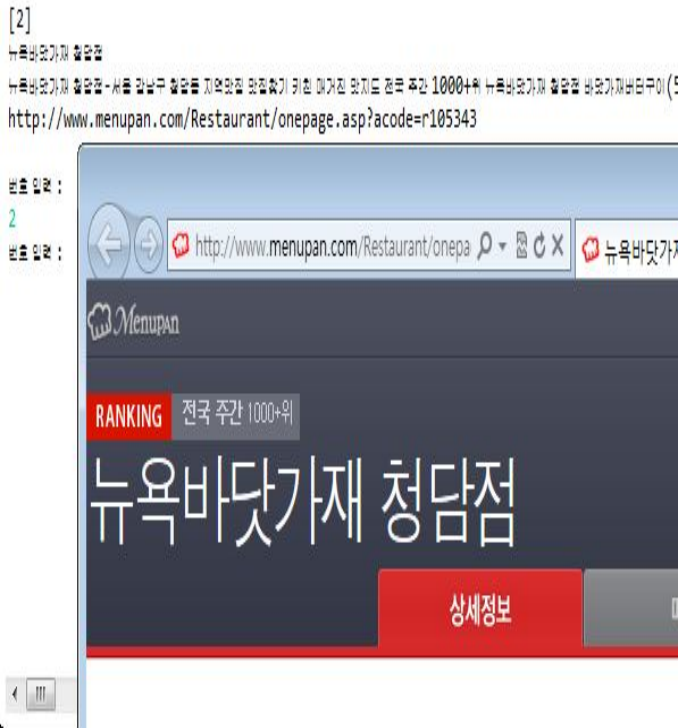


**Figure 5: Connecting to the Web site**

And our engine can connect to the web site of the restaurant when the information on the web site needs to be verified. The web site is provided when a user selects the restaurant index number. Figure 5 shows the restaurant web site which is accessed through the developed search engine.

## 4. CONCLUSION & FUTURE WORK

This paper only search the documents contained search keyword. For this reason, our engine is hard to find popular restaurant. To solve the problem, in the future work, we develop search engine that use ranking by adding the heuristic values such as TV program (that recommend popular restaurant) and positive word such as delicious, excellent, wonderful, and "visit again".

## REFERENCES

1. Lucene, Wikipedia, [Online] Available: http://en.wikipedia.org/wiki/Lucene
2. M. McCandless, E. Hatcher, O. Gospodnetic, Lucene IN ACTION, 2nd ed., Manning Publications, 2010.
3. Web Search Engine, Wikipedia, [Online] Available: http://en.wikipedia.org/wiki/Web_search_engine
4. Web Crawler, Wikipedia, [Online] Available: http://en.wikipedia.org/wiki/Web_crawler
5. Indexing, Wikipedia, [Online] Available: http://en.wikipedia.org/wiki/Searcg_engine_indexing
6. Web Search Query, Wikipedia, [Online] Available: http://en.wikipedia.org/wiki/Web_search_query