# Personalized Ontology Model for Web Information Gathering (Result Paper)

Authors

*Nilesh P.Patil[1], Dr.Santosh S.Lomte[2], Prof.Rajesh.A.Auti[3],*

Email- nileshpatil413@gmail.com,santoshlomte@hotmail.com, rajeshauti24@gmail.com

M.E.Student, Department of Computer, Dr. SeemaQuadri Institute of Tech, Aurangabad, India[1]

Professor, Department of Computer, Dr. SeemaQuadri Institute of Tech, Aurangabad, India[2]

Asst.Professor, Department of Computer, Dr. SeemaQuadri Institute of Tech, Aurangabad, India[3]

## ABSTRACT

The World Wide Web is an interlinked collection of billions of documents formatted using HTML. The amount of web based information available has increased dramatically. How to gather useful information from the web has become a challenging issue for users. Therefore, the new technology is to be introduced that that will be helpful for the web information gathering Ontology as model for knowledge description and formalization is used to represent user profile in personalized web information gathering. Ontology is the model for knowledge description and formalization. However the information of user profiles represents patterns either global or local knowledge base information, according to our analysis many models represents global knowledge. In this paper ontology system is used to recognize and reasoning over user profiles, world knowledge base and user instance repositories. This work also compares the analysis of existing system and ontology with other research areas are more efficient to represent.

*Keywords* – Local Instance Repository, Ontology, Personalization, Semantic Relations, User Profiles, Web Information gathering

## 1.1 Introduction

Today is the world of internet. The amount of the web-base information available on the internet has increased significantly. But gathering the useful information from the internet has become the most challenging job today's scenario. People are interested in the relevant and interested information from the web. The web information gathering systems before this satisfy the user requirements by capturing their information needs. For this reason user profiles are created for user background knowledge description. The user profiles represent the concepts models possessed by user while gathering the web information. A concept model is generated from user background knowledge and possessed implicitly by user. But many oncologists have observed that when user read a document they can easily determined whether or not it is of their interest or relevance to them .If the user concept model can be simulated, and then a better representation of the user profile can be build. To Simulate use concepts model, ontology's are

utilized in personalized web information gathering which are called ontological user profiles or personalized ontology's [1].In Global analysis, global knowledge bases are used for user background knowledge representation. Local analysis use local user information.

## Data Mining

The manual extraction of patterns from data has occurred for centuries. Early methods of identifying patterns in data include Bayes' theorem (1700s) and regression analysis (1800s). The proliferation, ubiquity and increasing power of computer technology has dramatically increased data collection, storage, and manipulation ability. As data sets have grown in size and complexity, direct "hands-on" data analysis has increasingly been augmented with indirect, automated data processing, aided by other discoveries in computer science, such as neural networks, cluster analysis, genetic algorithms (1950s), decision trees (1960s), and support vector machines (1990s). Data mining is the process of applying these methods with the intention of uncovering hidden patterns in large data sets. It bridges the gap from applied statistics and artificial intelligence (which usually provide the mathematical background) to database management by exploiting the way data is stored and indexed in databases to execute the actual learning and discovery algorithms more efficiently, allowing such methods to be applied to ever larger data sets.

## Personalized Ontology Construction

Personalized ontology's are a conceptualization model that formally describes and specifies user background knowledge. From observations in daily life, we found that web users might have different expectations for the same search query. For example, for the topic "New York," business travellers may demand different information from leisure travellers. Sometimes even the same user may have different expectations for the same search query if applied in a different situation. A user may become a business traveller when planning for a business trip, or a leisure traveller when planning for a family holiday. Based on this observation, an assumption is formed that web users have a personal concept model for their information needs. A user's concept model may change according to different information needs. In this section, a model constructing personalized ontology's for web users' concept models is introduced.

## World Knowledge Representation

World knowledge is important for information gathering. According to the definition provided by, world knowledge is commonsense knowledge possessed by people and acquired through their experience and education.Also, as pointed out by Nirenburg and Raskin , "world knowledge is necessary for lexical and referential disambiguation, including establishing co reference relations and resolving ellipsis as well as for establishing and maintaining connectivity of the discourse and adherence of the text to the text producer's goal and plans." In this proposed model, user background knowledge is extracted from a world knowledge base encoded from the Library of Congress Subject Headings (LCSH).

## 1.2 Necessity

Compared with the TREC model, the Ontology model had better recall but relatively weaker precision performance. The Ontology model discovered user background knowledge from user local instance repositories, rather than documents read and judged by users. Thus, the Ontology user profiles were not as precise as the TREC user profiles; The Ontology profiles had broad topic coverage. The substantial coverage of Possibly-related topics were gained from the use of the WKB and the large number of training documents. Compared to the web data used by the web model, the LIRs used by the Ontology model were controlled and contained less uncertainties. Additionally, a large number of uncertainties were eliminated when user background knowledge was discovered. As a result, the user profiles acquired by the Ontology model performed better than the web model.

## 1.3 Objective

Data mining, which aims at extracting interesting information from large collections of data, has been widely used as an active decision making tool. Real world applications of data mining require a dynamic and resilient model that is aware of a wide variety of diverse and unpredictable contexts. Contexts consist of circumstantial aspects of the user and domain that may act the data mining process. The underlying motivation is mining datasets in the presence of context factors may improve performance and easy of data mining as identifying the factors, which are not easily detectable with typical data mining techniques. Ontology's
are the structural frameworks for organizing information and are used in arterial intelligence, the Semantic Web, systems engineering, software engineering, biomedical informatics, library science, enterprise bookmarking and information architecture as a form of knowledge representation about the world or some part of it. The creation of domain ontology's is also fundamental to the definition and use of an enterprise architecture framework. Sharing common understanding of the structure of information among people or software agents is one of the more common goals in developing ontology's.

Applying data mining for ontology building:

Ontology represents the concepts and the relationship between them for specialized domain. Building ontology is a complex work, in order to build ontology you need a domain expert to help you to declare all domain concepts and the relationship between them. In this work we propose a methodology for building ontology based on the output of data mining result.

Ontology in Software Engineering:Modelling ontology is a tedious task always important to demonstrate can gain by applying ontologies in software engineering, the current advent of logic based formalisms in the context of the semantic web effort is an important factor. Activities by the W3C and others have helped to flesh out standards like RDF or OWL receive increasing attention by tool builders and users. Important factor is the flexibility of ontologies with information integration as a major use case, ontologies are well to combine information from various sources and infer new factors and also the flexibility. Further promoted by the"web"- focus of current ontology approaches due to the fact that software systems also get increasingly web-enabled and must thuscope with data

from heterogeneous sources that may not be known at developmenttime, software engineers seek technologies that can help in this situation. Thus,experts in the field like Grady Booch are expecting semantic web technology to beone of the next big things in the architecture of web-based applications [4]. Also, theweb makes it easier to share knowledge. Having URIs as globally unique identifiers,it is easy to relate one's ontology to someone else's conceptualization. This in turnencourages interoperability and reuse.Regarding more Software Engineering-specific advantages, ontologies makedomain models first order citizens. While domain models are clearly driving the coreof every software system, their importance in current Software Engineering processesdecreases after the analysis phase. The core purpose of ontologies is by definition theformal descriptions of a domain and thus encourages a broader usage throughout thewhole Software Engineering lifecycle.

## RELATED WORK

The LGSM (Local Global search methodology) it is used to calculate the hit/miss rate. For calculating hit ratio,

$$\text{Hit Ratio} = \frac{\text{Number of Hits}}{(\text{Number of Hits} + \text{Number of Miss})}$$

The performance of memory is frequency measured in terms of quantity is called hit ratio. When cpu needs to find the word in cache, if word is found in cache then it produces a hit. If the word is not found in the cache, it is in main memory it is counted as miss. If it retrieves information from the local repository it is considered as hit. If it retrieves data directly from global it is considered as miss [6].

Our project refer ontology model and the proposed ontology model aims to discover user background knowledge and learns personalized ontology's to represent user profiles see the figure.
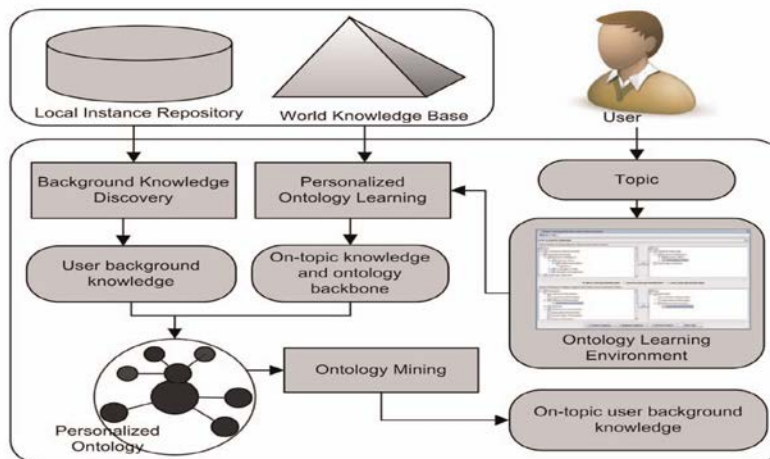


**Figure 2.2 Architecture of Ontology Model**

## IMPLEMENTATION

### Tf-Idf and Cosine similarity Algoritham's

In the year 1998 Google handled 9800 average search queries every day. In 2012 this number shot up to 5.13 billion average searches per day. The graph given below shows this astronomical growth.

Step 1: Term Frequency (TF)

Term Frequency also known as TF measures the number of times a term (word) occurs in a document. Given below are the terms and their frequency on each of the document.

TF for Document 1

| Document1 | the | game | of | life | is | a | Everlasting | learning |
|---|---|---|---|---|---|---|---|---|
| Term Frequency | 1 | 2 | 2 | 1 | 1 | 1 | 1 | 1 |

In reality each document will be of different size. On a large document the frequency of the terms will be much higher than the smaller ones. Hence we need to normalize the document based on its size. A simple trick is to divide the term frequency by the total number of terms. For example in Document 1 the term game occurs two times. The total number of terms in the document is 10. Hence the normalized term frequency is 2 / 10 = 0.2. Given below are the normalized term frequency for all the documents.

Normalized TF for Document 1

| Document1 | the | game | of | life | is | a | everlasting | learning |
|---|---|---|---|---|---|---|---|---|
| Normalized TF | 0.1 | 0.2 | 0.2 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |

Step 2: Inverse Document Frequency (IDF)

The main purpose of doing a search is to find out relevant documents matching the query. In the first step all terms are considered equally important. In fact certain terms that occur too frequently have little power in determining the relevance. We need a way to weigh down the effects of too frequently occurring terms. Also the terms that occur less in the document can be more relevant. We need a way to weigh up the effects of less frequently occurring terms. Logarithmshelps us to solve this problem.

Let us compute IDF for the term game

$IDF(game) = 1 + \log_e(\text{Total Number Of Documents} / \text{Number Of Documents with term game in it})$

There are 3 documents in all = Document1, Document2, Document3

The term game appears in Document1

IDF(game) = 1 + $\log_e$(3 / 1)

= 1 + 1.098726209

= 2.098726209

Step 3: TF * IDF

Remember we are trying to find out relevant documents for the query: life learning For each term in the query multiply its normalized term frequency with its IDF on each document. In Document1 for the term life the normalized term frequency is 0.1 and its IDF is 1.405507153. Multiplying them together we get 0.140550715(0.1 * 1.405507153). Given below is TF * IDF calculations for life and learning in all the documents.

Step 4: Vector Space Model – Cosine Similarity

From each document we derive a vector. The set of documents in a collection then is viewed as a set of vectors in a vector space. Each term will have its own axis. Using the formula given below we can find out the similarity between any two documents.

Cosine Similarity (d1, d2) = Dot product(d1, d2) / ||d1|| * ||d2||

Dot product (d1,d2) = d1[0] * d2[0] + d1[1] * d2[1] * … * d1[n] * d2[n]

$||d1||$ = square root(d1[0]$^2$ + d1[1]$^2$ + ... + d1[n]$^2$)

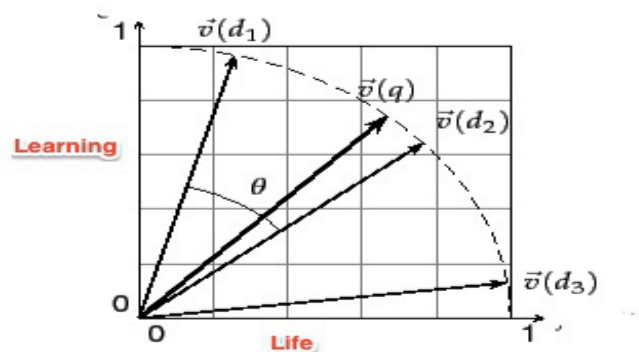$||d2||$ = square root(d2[0]$^2$ + d2[1]$^2$ + ... + d2[n]$^2$)



Figure.3.2.1 Graph for vector medel

## RESULTS AND ANALYSIS

We have done experimentation of users for the results and analysis. The Following parameters come in to the picture.

Figure 1 shows the home page of the Ontology model , in that we see that we insert the queiry in search engine. And we select by with algoritham we can find the result of the inserted quiry.
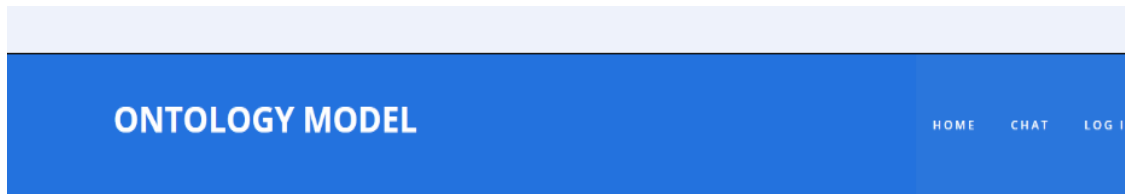


**Figure 1. Home Page**

In this figure we see that the home page of the Ontology Model ,on that we insert the quiey that will be search we also see that by which algoritham(Tf-IDF and Cosine Similarity Algoritham) we can find the result of the inserted queiry.



| Document | Index | TF | TWC | IDF |
|----------|-------|----|----|-----|
| Index | 0.00000 | 0 | 16 | 2.29928298413026 |
| Index | 0.00000 | 0 | 23 | 2.29928298413026 |
| Index | 0.00000 | 0 | 36 | 2.29928298413026 |
| Index | 0.00000 | 0 | 34 | 2.29928298413026 |
| Index | 0.00000 | 0 | 35 | 2.29928298413026 |
| Index | 0.00000 | 0 | 37 | 2.29928298413026 |

**Figure 2 Normalization Form**

In figure 3 we see that the chating page of the model , for the chating we need the username and password if we are allrady register if user is not register than register first and than enter, if we forget password that click on forget password? Then we can change the password.
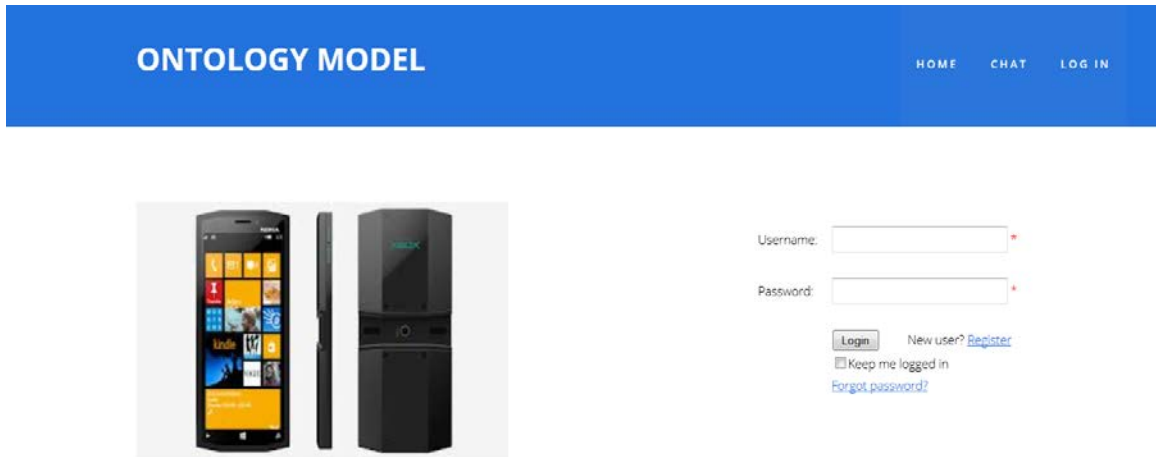


**Figure 3 Chatting Form**

In figure 4 we see that verify the serch query by Tittle if it is alrady in the database than display it if not than search the query.
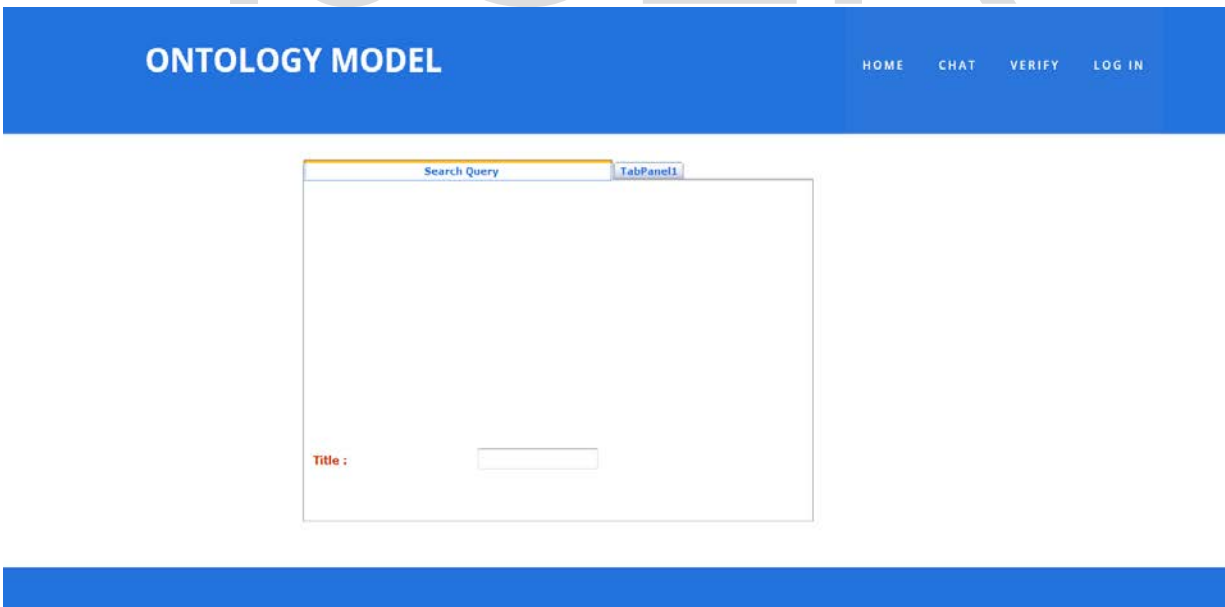


**Figure 4 Query Upload Form**

In figure 5 we see the user online form in that we see the which user are online on this model on this form we gives the send button with the help of this button we sed the message to other online use. And it also gives the color with our conveneace.
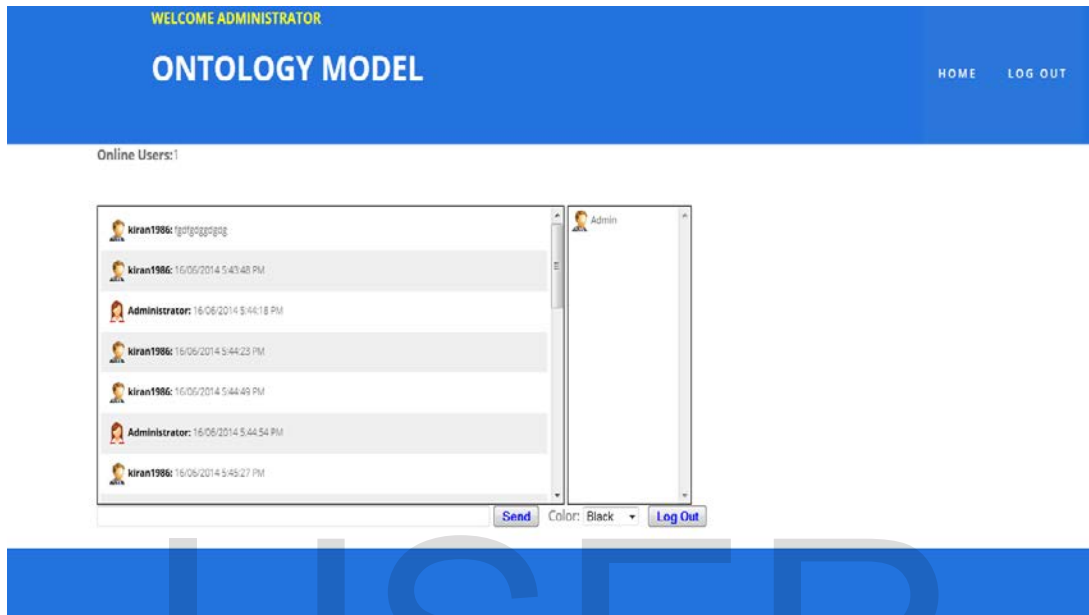


**Figure 5 Online user Form**

In figure 6 we see the performance evaluation of the screen with the help algoritham and with is the accurate result of the search query is display,with the help of graph.
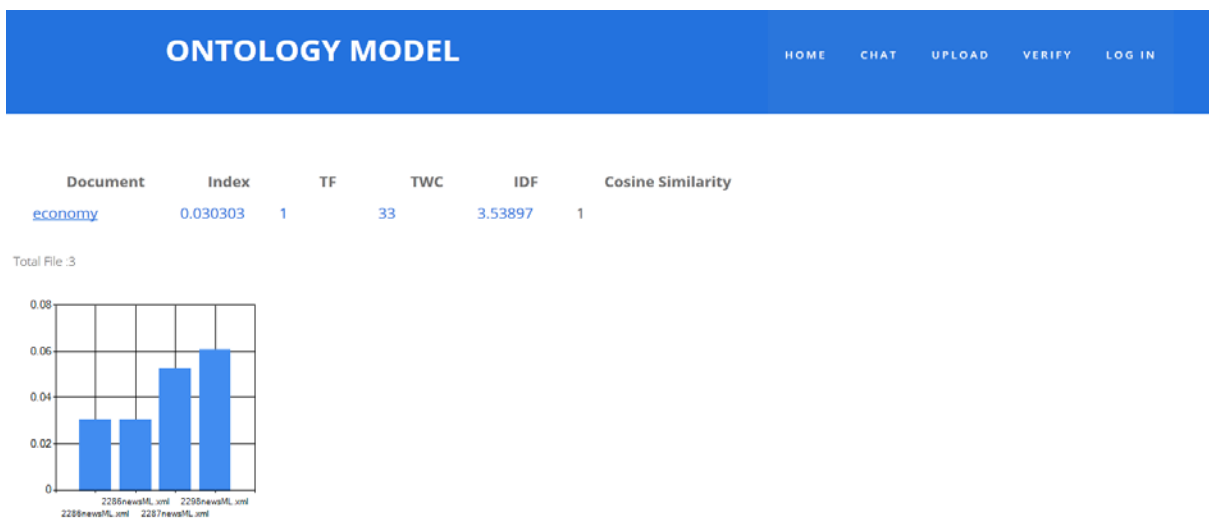


**Figure 6 Performance evaluation (Accuracy plot)**

## CONCLUSION

In this way we study an ontology model for personalized web information gathering. The model constructs user personalized ontology's by extracting world knowledge from the LCSH system and discovering user background knowledge from user local instance repositories. The ontology model in this paper provides a solution to emphasizing global and local knowledge in a single computational model. The findings in this paper can be applied to the design of web information gathering systems. The model also has extensive contributions to the fields of Information Retrieval, web Intelligence, Recommendation Systems, and Information Systems.

## ACKNOWLEDGEMENT

## REFERENCES

1. S.Gauch, J.Chaffee, and A.Pretschner,"Ontology-Based Personalized Search and Browsing" Web Intelligence and Agent Systems, vol. 1, nos. 3/4, pp.219-234, 2003.
2. Y.Li and N. Zhong, "Web Mining Model and Its Applications for information Gathering" Knowledge-Based Systems, vol. 17, pp. 207-217, 2004.
3. Y. Li and N. Zhong, "Mining Ontology for Automatically Acquiring Web User Information Needs," IEEE Trans. Knowledge and Data Eng., vol. 18, no. 4, pp.554-568, Apr. 2006.
4. X. Tao, Y. Li, N. Zhong, and R. Nayak, "Ontology Mining for Personalized Web Information Gathering," Proc. IEEE/WIC/ACM Int"l Conf. Web Intelligence, pp. 351-358, 2007.
5. A. Sieg, B. Mobasher, and R. Burke, "Web Search Personalization with Ontological User Profiles," Proc. 16th ACM Conf. Information and knowledge Management(CIKM"07),pp.525-534,2007.
6. J.D. King, Y. Li, X. Tao, and R. Nayak, "Mining World Knowledge for Analysis of Search Engine Content," Web Intelligence and Agent Systems, vol. 5, no. 3, pp. 233-253, 2007
7 .R.Baeza-Yates and B. Ribeiro-Neto, Modern Information Retrieval. Addison Wesley, 1999.
8. L.M. Chan, Library of Congress Subject Headings: Principle and Application. Libraries Unlimited, 2005.
9. P.A. Chirita, C.S. Firan, and W. Nejdl, "Personalized Query Expansion for the Web," Proc. ACM SIGIR ('07), pp. 7-14, 2007.
10. R.M. Colomb, Information Spaces: The Architecture of Cyberspace. Springer, 2002.