

# Tutorial Series on Web-computing 1 Information Gathering and Searching Approaches on the Web

# Seiji YAMADA

Department of Computational Intelligence and Systems Science Interdisciplinary Graduate School of Science and Engineering Tokyo Institute of Technology 4259 Nagatsuta-cho, Midori-ku, Yokohama, 226-8502 JAPAN yamada@ymd.dis.titech.ac.jp

Hiroyuki KAWANO Department of Systems Science Graduate School of Informatics Kyoto University Sakyou-ku, Yoshida-honmachi, Kyoto, 606-8501 JAPAN kawano@i.kyoto-u.ac.jp

Received 15 November 2000

Abstract The information accessible through the Internet is increasing explosively as the Web is getting more and more widespread. In this situation, the Web is indispensable information resource for both of information gathering and information searching. Though traditional information retrieval techniques have been applied to information gathering and searching in the Web, they are insufficient for this new form of information source. Fortunately some AI techniques can be straightforwardly applicable to such tasks in the Web, and many researchers are trying this approach. In this paper, we attempt to describe the current state of information gathering and searching technologies in the Web, and the application of AI techniques in the fields. Then we point out limitations of these traditional and AI approaches and introduce two approaches: navigation planning and a Mondou search engine for overcoming them. The navigation planning system tries to collect systematic knowledge, rather than Web pages, which are only pieces of knowledge. The Mondou search engine copes with the problems of the query expansion/modification based on the techniques of text/web mining and information visualization.

Keywords: The World Wide Web, Information Gathering, Web Search Engine

# §1 Introduction

As the Internet and the Web are developing explosively, we can easily utilize them as useful information resources. Since the Web is very huge and constantly updated, information and knowledge in the Web is rich and recent. In this situation, there is a possibility that we can do information gathering and information searching in the Web anytime, anywhere to obtain our target Web pages. However there are significant problems. In general, a user does not know where his/her target Web pages exist. Thus a Web search engine is used to search for the target pages. Unfortunately, since the filtering is not sufficient, a hit list from a search engine may include a large number of non-relevant Web pages. This problem also influences information gathering. If information searching is not correct, the gathered information is also meaningless.

Though traditional information retrieval techniques<sup>28)</sup> have been applied to information gathering and search in the Web, not all of the problems are solved. For example, filtering of a hit list using semantics and the context of Web pages is hard for information retrieval methods, and natural language processing is necessary. Fortunately some AI techniques are straightforwardly applicable to such tasks in the Web, and many researchers are trying this approach.

In this paper, first we describe the current state of information gathering and searching in the Web, and survey systems developed as AI applications. Next we point out open problems in the fields. Then we introduce two approaches: navigation planning and a **Mondou** search engine for overcoming the issues.

## §2 Information gathering and searching in the Web

In this section, we survey conventional approaches and implemented systems for information gathering and searching in the Web, then we describe their limitations.

#### 2.1 Traditional Systems

#### [1] Softbot, Sage and Occam

There are studies on AI planning to generate procedures for gathering information through the Internet. The  $Softbot^{5}$  provides a planning framework which consists of operators and a planner in the similar way to STRIPS planning. The environment and the operator are described by first-order predicate calculus. An operator represents an action by UNIX commands like mkdir, find and so on. *Softbot* has user interface for a naive user to describe operators without knowledge on predicate calculus. A complete and sound partial-ordering planner is used in *Softbot*. *Occam*<sup>15)</sup> is also a planner for gathering information. It is more efficient and able to reason about the capabilities of different information sources.

Consider a command "send E-mail to Dr. Yamada in TiTech" is given to *Softbot* as a goal description. First it tries to reduce the ambiguity in the command. *Softbot* identifies a person named Dr. Yamada in TiTech. In this processing, the knowledge about the organization of TiTech is used. Next *Softbot* 

finds applicable operators and generates a plan for achieving the given goal. In this example, the operators describing finger and sendmail are selected and added into a plan. When a complete plan is obtained, *Softbot* executes it and the goal is satisfied. *Sage*<sup>14)</sup> is also developed for integrating planning, execution, replanning and sensing to gather information in distributed resources. The aim of these studies is to generate a plan as a procedure of gathering information, and a plan consists of UNIX commands or database operations.

## [2] ARACHNID

ARACHNID<sup>17)</sup> (or *InfoSpider*<sup>18)</sup>) is an information gathering system which manipulates multiple search agents in the Web. ARACHNID applies an A-Life like approach for controlling the agents. Each agent has energy for describing its activity. If an agent finds a relevant Web page, its energy increases and otherwise the energy decreases. When an agent gets fully active, it produces children, and if it significantly looses energy, it will be dead. In ARACHNID, since information gathering is done on-line, it can obtain more recent information than a search engine using off-line search.

## [3] Meta Search Engine

While research on a single search engine is active, meta search engines have also been developed. A meta search engine utilizes other search engines for information search in the Web. First it sends a query inputted by a user to several search engines. Next it receives a hit lists from them. Finally it integrates the results and indicates them to the user. MetaClower<sup>23,24)</sup> is a typical meta search engine. SavvySearch<sup>8)</sup> is able to characterize search engines using user's feedback.

#### [4] Web Robot

Programs called *Spider*, *Clower*, *Web robot* are utilized to gather Web pages for construction of a search engine database.<sup>3,9)</sup> When start points are given to a Web robot, it begins to trace the links from the current Web pages and gather pages using breadth-first search. In general, the purpose of a Web robot is to collect Web pages as many as possible, and it has no control mechanism for selective information gathering.

#### [5] Recommender and Shopping Agent

 $WebWatcher^{10}$  is a recommender agent which guides an user in Web browsing. In the similar way to a human guide in a museum, WebWatcher can infer and recommend the Web page candidates which a user is likely to select next by using his/her browsing history. The knowledge for prediction of the Web pages to which a user goes next is acquired by machine learning methods. Reinforcement learning and instance-based learning are applied to obtain the knowledge. *Letizia*<sup>16</sup> is also able to recommend to a user the next interesting Web pages.

Some learning systems have been developed for information gathering and browsing in the Web.  $ShopBot^{4)}$  learns the text pattern indicating the price of

CD-ROMs, and searches for the cheapest one more efficiently than a human. Fist *ShopBot* collects both of positive and negative training examples by actively investigating on-line shop Web sites. It tries to give different types of input to the site and the returned Web pages are dealt as training examples. Using them, *ShopBot* can utilize an inductive learning method and obtain a description of a vender site.

## 2.2 Open Problems

As mentioned earlier, various AI applications to information gathering and searching in the Web have been developed thus far. Some systems successfully overcome the limitation of traditional information retrieval approach like TFIDF indexing, statistical information and so on. However we can point out open problems in the following.

- Gathering pieces of knowledge: In all of conventional systems, the gathered information is only a set of pieces of knowledge. For example, any search engine collects Web pages which mention only surface knowledge on a query. We need more systematic and deep knowledge on a target information described in a query.
- Appropriate combination of keywords: We need frequent patterns or associative rules to determine a much more suitable combination of keywords in a query. For instance, when we describe a query, we usually consider the distance of keywords in the meaning networks, and evaluate the coverage of expressions.

In the following sections, we introduce two approaches to address the above limitations. Navigation planning and a Mondou search engine are described.

# **§3** Navigation Planning

The Web is very useful for a user who wants to understand a *target concept*. He/she can browse helpful Web pages to understand a target concept. However, in general, this task is very hard because he/she may not know where such Web pages are, and has to search for them over the vast Web search space. A solution to this problem is to use a search engine with the target concept as a query. However, since the retrieved Web pages are not filtered sufficiently, a user has to select useful ones from them. Furthermore, since in most cases the gathered Web pages are pieces of surface knowledge which directly explain the target concept and include concepts that a user does not understand, he/she must search the useful Web pages for them using a search engine again. This task is repeated until a user understands the target concept, and hence wastes time. We consider the task as planning, and propose *navigation planning*<sup>25,26)</sup> to automatically generate a sequence of Web pages which can guide a user to understand a target concept deeply.

## 3.1 Web Browsing as Planning

In this research, *navigation* means a task that indicates a sequence of useful Web pages which guides a user to understand a concept. The sequence of Web pages is called a *plan*, and *navigation planning* means the automatic generation of the plan. We can summarize the Web browsing task for concept understanding in the following. This procedure is iterated until terminated by the user.

- 1. Search Web pages using a search engine.
- 2. Read and understand the pages retrieved by the search engine.
- 3. Select unknown concepts in the Web pages.
- 4. Go to Step 1 with unknown concepts as target concepts.

The procedure above is considered as *planning*<sup>7</sup> under the following correspondence. Using this formalization, we can apply a classical planning framework to navigation planning.

- Action: Understanding concepts on a Web page.
- *State*: A user's knowledge state described with a set of words describing concepts which he/she knows.
- Initial state: A user's initial knowledge state.
- *Goal state*: A target concept described with a set of words which a user wants to understand.
- Operator: U-Op(URL) defined by the followings.
  - Label: URL of the Web page
  - Condition:  $C = \{c_1, \dots, c_i\}$ , where C means the condition words which are necessary to understand the pages.
  - Effect:  $E = \{e_1, \dots, e_j\}$ , where E is effect words which a user obtains by understanding the page.

This navigation planning contains a significant problem which has not been presented in planning. It is that the U-Op(URL) operators are not given in advance. This is because it is impossible for a human designer to generate the operators from all the Web pages in the Web. Hence the operators need to be automatically generated from Web pages when they are necessary.

## 3.2 Generating Operators from Web Pages

## [1] Using Tag Structure in a HTML File

Various methods to extract keywords from text have been studied.<sup>22)</sup> Though most methods are based on the frequency of words, one of the most effective methods is to utilize the structure in text. Since a Web page is described in a HTML format, we can utilize tag structures.

The prime candidates for condition words are those linked to other Web pages, i.e. the words between <A HREF=URL> and </A>, because this tag is a sign of reference to relevant topics, which are important for understanding the current Web page in many cases.

S. Yamada and H. Kawano



(b) Expansion of a node-i

Fig. 1 Navigation Planning

Since the title of the Web page describes words which a user may acquire by reading the page, the words between <TITLE> and </TITLE> are candidates for the effect words. In the same way, headings describe knowledge which a user may obtain by reading the section. Thus the words between <Hn> and </Hn> are also candidates for the effect words.

## [2] KeyGraph: A Keyword Extraction Method

The extraction of condition and effect words using the tag structure is not sufficient. All the linked words are not necessarily candidates for condition words, and all the condition words are not necessarily linked. Thus we look for another method to assist it, and *KeyGraph* is found to be appropriate for this task.

KeyGraph is a fast extraction method of keywords representing the asserted core idea in a document.<sup>19)</sup> KeyGraph composes clusters of terms, based on co-occurrences between terms in the document. Each cluster represents a concept on which the document is based (i.e. condition words), and terms connecting clusters tightly are obtained as author's assertion (i.e. effect words). Furthermore the likelihood for condition and effect words can be computed by KeyGraph, and used for weight of an operator. Another merit of KeyGraph is that it does not employ a corpus.

The extraction of condition and effect words using tag structure and *Key-Graph* is integrated. First condition and effect words are obtained with the tag structure. Next *KeyGraph* also extracts condition and effect words and both of



Fig. 2 Plan for "Concept Formation"

them are integrated.

# 3.3 Planning Procedure

We have developed a navigation planning procedure. Figure 1 shows the overview of our navigation planning system, called *NaviPlan*. *NaviPlan* automatically generates a plan for given target concepts. It uses *backward beam* search<sup>21</sup> from a goal state (Fig. 1(a)). The expansion of a node (Fig. 1(b)) includes the search for relevant Web pages with a search engine and the generation of operators.

We fully implemented *NaviPlan* using Perl. Figure 2 shows a plan (depth = 4) generated by *NaviPlan* with the target concept "concept formation." Four useful Web pages are indicated with hyperlinked titles and URLs. By reading the Web pages in this order, a user is able to understand "concept formation" deeply.

# 3.4 Experimental Evaluations

Now let us evaluate *NaviPlan* according to the usage of the system as a navigation planner. In this experiment, we employed the search engine MetaCrawler, for expanding nodes in the planning procedure of *NaviPlan*. The system compared with *NaviPlan* here was also MetaCrawler.

For each of 10 subjects (8 graduate school students and 2 university staffs), we investigated the number of Web pages he/she read for understanding a target concept and the brief explanation on the goal. As results, we found out *NaviPlan* achieved efficient navigations to the goal, due to the strong intension of operators to reach the goal from user's initial state.

In order to investigate the effect of having a plan instead of only one page to read, we also made experiments using more difficult target concepts for the same subjects to understand. Consequently we conclude that planning in *NaviPlan* contributes to increasing user's knowledge until he/she understands a goal concept, by reading multiple but small number of pages.



Fig. 3 The Structure of Mondou

# §4 Intelligent Search System

In 1996, We have developped "Mondou" web search engine, which belongs to the first generation of research oriented Japanese web search engines.  $*^{i \ 12}$  We have applied the rapid emerging technologies of data mining to the mechanism of query modification and expansion. We also implemented Java applets in order to support the search steps. In this section, we explain major technical aspects of Mondou systems, such as data cleaning, data mining and information visualization.

#### 4.1 Architecture of Mondou Search Engine

The Mondou consists of three primary modules, such as *web robot*, *text database* and *query server*, which are shown in Fig. 3.

The fundamental functions of the *web robot* are described in Section 2.1 [4]. In addition to the basic mechanisms of the web robot, our web robot parses HTML tags and extracts the important keywords from the *title, headings and sub-headings, anchor strings*, and other tagged attributes. Moreover it gives appropriate weight values to derived keywords in order to measure the importance of web pages using the heuristic functions.

Our text database handles not only the characters in the web pages, but also the connectivities of hyperlinks. Therefore, the primary database table consists of the following attributes, "keywords, URLs, hyper links, http servers, IP addresses", and the system also keeps the control/management data tables for operating the web robots.

The *query server* is the most important program in the **Mondou** system. We implement typical data mining algorithms to derive appropriate association rules. The associative keywords with some attribute values are visualized on the java applets. Furthermore, in order to adjust the demand of searchers, we tried to apply agent technologies and the integrating mechanism of heterogeneous

<sup>\*&#</sup>x27; The URL of Mondou web search engine is http://www.kuamp.kyoto-u.ac.jp/labs/ infocom/ mondou/index.e.html. The name of Mondou comes from the term of Zen, which is a sophisticated discussion style in Japan, so the discovery is in the processes of questions and answers.

databases to Mondou systems.

#### 4.2 Data Cleaning by Web Robots

As we mentioned in Section 2.1, the simple web robot visits web servers sequentially in the breadth-first manner. However, we have to consider the performance of gathering the web pages more effectively, then we try to collect meaningful or useful web pages as fast as possible.

It is so hard to discover the meaningful knowledge from a web graph, then we try to extract the interesting features of web pages using the following two types of hyperlinks.

- Inside link: the link to other pages on the same web server.
- Outside link: the link to other pages on the other web servers.

Actually, when we browse web pages, we consider the connectivities and distance of the web page from/to other web pages. Therefore, it may be helpful to analyze and characterize the connectivities of hyperlinks, in order to fast collect important web pages in the same and different web servers.

Therefore, we proposed a focused crawling algorithm in our paper,<sup>12)</sup> which is based on the assumption that the important web pages are referred many times from other servers. And it is also helpful to reduce the number of noisy/meaningless web pages from the view point of data cleaning.

Furthermore, in order to reduce the load of the network and the number of web page requests, we are trying to implement cooperative web robots based on a distributed scheduling algorithm.<sup>27)</sup>

#### 4.3 Data Mining for Query Modification

The research of data mining<sup>6</sup> consists of databases, statistics and various kinds of AI research, such as machine learning, inductive learning, knowledge representation, information visualization and others.

The data mining plays an important role in order to support the discovery procedures in the Mondou system. We improve the typical data mining algorithms of association rules,<sup>1)</sup> and apply them to the set of web pages.

Shortly speaking, given a set of web documents, where each document consists of various keywords, an association rule is an expression  $X \Rightarrow Y$ , where X and Y are sets of keywords. The intuitive meaning of such a rule is that documents which contain the keywords X tend to also contain the keywords Y. The support threshold value of the rule  $X \Rightarrow Y$  is the percentage of documents that contain both X and Y.

Table 1 shows typical examples of derived keywords, the number of URLs and associative keywords by our algorithm.  $^{\ast 2}$ 

In order to refine these rules much more by Mondou, it is effective to classify several attributes and use the relationships of different attributes by using semi-structured data. Moreover, our Mondou system derives the different

<sup>&</sup>lt;sup>+2</sup> A derived keyword, "rcaau", stands for "Retrieval loCation by weighted AssociAtion rUle," which appears in the web page of Mondou.

keywords	No. of URLs	related keywords
Applications	975	windows, speech, computer, internet, helper, ···
Informatics	462	genome, medical, department, center, school, ···
Japan	20,366	japan country, webec, japan, japanese, back
Engine	2,836	lycos, yahoo, dragon, search, creative, japan
Mining	134	akita, geology, college, net, data
Mondou	28	rcaau, search

 $\begin{tabular}{ll} Table 1 & Examples of Associative Keywords \end{tabular}$ 

set of association rules from the different databases, such as netnews archives, electronic journals and other digital libralies. Finally our system deduces much more effective combination of keywords from heterogeneous databases by using the difference of derived rules. We also implement the agent oriented communication interface, so that different **Mondou** systems can communicate each other by high level query languages like KQML.

## 4.4 Information Visualization

In order to visualize the characteristics of association rules and recognize the features of document clusters, we constructed interactive visual interfaces<sup>13)</sup> as Java applets.

In Fig. 4, this 3D graph shows the distributions of derived keywords based on ROC graphs.<sup>11,20)</sup> We also implemented fisheye graph drawing techniques.<sup>2)</sup> By the simple mouse operations, users can determine a combination of keywords easily.

For example, in Fig. 4, by a query with a keyword "mediator", Mondou derives associative keywords, such as "information", "data", "distributed", "supporting", "connections", "personal" and others. Figure 5 shows an example of visualizing results for a selected query using some of above derived keywords. In Fig. 5, numerical numbers mean the length of queries, and the color of a keyword changes by pointing the number. In the right window, this applet displays the search results which satisfy the combination of black keywords.

Consequently, users can easily select the interesting combination of keywords by data mining and information visualization, which are provided by our Mondou system.

# §5 Conclusion

We surveyed the conventional systems of information gathering and searching in the Web. Information gathering systems, *Softbot*, *Sage*, ARACHNID and so on were introduced. Though some of them have successfully overcome the problems which traditional IR systems hardly cope with, there are still open problems. We also pointed out two significant problems, the gathering pieces of knowledge and the appropriate combination of keywords in the fields.

As an attempt to introduce a promising approach to gather systematic information, not pieces of knowledge, we described navigation planning that generates a plan guiding a user to understand a concept in the Web. It has the



Fig. 4 Visual Interface for Keyword Selection

ability to generate operators during planning from Web pages using keyword extraction methods. The search for useful Web pages for a user to understand goal concepts was formalized using an AI planning framework, and an operator corresponding to the understanding of a Web page was defined with condition and effect knowledge. Then we described the whole planning procedure, and verified the effectiveness experimentally.

In a field of a Web search engine, many researchers are trying to automatically generate the appropriate combination of keywords in a query. A **Mondou** web search engine was introduced as one of web search engines in the first generation, which was based on the emerging technologies of data mining, text mining, web mining and information visualization. In particular, we focused on the query modification, which is a typical and important problem in the research of information retrieval. Our implemented functions, the derivation and visualization of associative keywords in order to modify the initial query, are empirically found very effective for novice users.



Fig. 5 Visual Interface of Document Clustering

# References

- Agrawal, R. and Srikant, R., "Fast Algorithms for Mining Association Rules," in *Proc. of the 20th International Conference on Very Large Data Bases*, pp. 487–499, Santiago, Chile, 1994.
- Chaomei, C., Information Visualisation and Virtual Environments, Springer-Verlag, 1999.
- 3) Cheong, F. C., Internet Agents: Spiders, Wanderers, Brokers, and Bots, New Riders, 1996.
- 4) Doorenbos, R. B., Etzioni, O. and Weld, D. S., "a Scalable Comparison-shopping Agent for the World-Wide Web," in *Proc. of the First International Conference on Autonomous Agent*, pp. 39–48, 1997.
- 5) Etzioni, O. and Weld, D., "A SoftBot-based Interface to the Internet," Communication of the ACM, 37, 7, pp. 72-76, 1994.
- 6) Feldman, R., "Practical Text Mining," in Second Symposium on Principles of Data Mining and Knwoledge Discovery (PKDD-97), Nantes, France, 1998.

- Fikes, R. E. and Nilsson, N. J., "STRIPS: a New Approach to the Application of Theorem Proving to Problem Solving," *Artificial Intelligence*, 2, pp. 189–208, 1971.
- 8) Howe, A. E. and Dreilinger, D., "Savvy Search: a Metasearch Engine that Learns Which Search Engines to Query," *AI Magazine*, 18, 2, pp. 19–25, 1997.
- 9) Jamsa, K., Lalani, S. and Weakley, S., Web Programming, Jamsa Press, 1996.
- 10) Joachims, T., Freitag, D. and Mitchell, T., "Webwatcher: a Tour Guide for the World Wide Web," in Proc. of the Fifteenth International Joint Conference on Artificial Intelligence, pp. 770-775, 1997.
- 11) Kawahara, M. and Kawano, H., "Performance Evaluation of Bibliographic Navigation System with Association Rules from Roc Convex Hull Method," *Transactions of the IPSJ:Database, 40*(SIG3(TOD1)), pp. 105–113, 1999.
- 12) Kawano, H., "Mondou: Web Search Engine with Textual Data Mining," in Proc. of IEEE Pacific Rim Conference on Communications, Computers and Signal Processing, pp. 402-405, 1997.
- 13) Kawano, H. and Kawahara, M., "Mondou: Information Navigator with Visual Interface," in *Data Warehousing and Knowledge Discovery, Second International Conference, DaWaK 2000*, pp. 425–430, London, UK, Sep. 2000.
- 14) Knoblock, C. A., "Planning, Executing, Sensing, and Replanning for Information Gathering," in *Proc. of the Fourteenth International Joint Conference on Artificial Intelligence*, pp. 1686–1693, 1995.
- 15) Kwok, C. T. and Weld, D. S., "Planning to Gather Information," in Proc. of the Thirteenth National Conference on Artificial Intelligence, pp. 32–39, 1996.
- 16) Lieberman, H., "Letizia: a Agent that Assists Web Browsing," in Proc. of the Fourteenth International Joint Conference on Artificial Intelligence, pp. 924-929, 1995.
- 17) Menczer, F., "ARACHNID: Adaptive Retrieval Agents Choosing Heuristic Neighborhoods for Information Discovery," in *Proc. of the Fourteenth International Conference on Machine Learning*, pp. 227–235, 1997.
- 18) Menczer, F. and Monge, A. E., "Scalable Web Search by Adaptive Online Agents: an Inforspiders Case Study," in *Intelligent Information Agents*, pp. 323– 347. Springer, 1999.
- 19) Ohsawa, Y., Benson, N. E. and Yachida, M., "KeyGraph: Automatic Indexing by Co-occurrence Graph Based on Building Construction Metaphor," in Proc. of IEEE Advanced Digital Library Conference, pp. 12–18, 1998.
- 20) Provost, F. and Fawcett, T., "Analysis and Visualization of Classifier Performance: Comparison Under Imprecise Class and Cost Distributions," in Proc. of 3rd Int'l Conference on Knowledge Discovery and Data Mining (KDD-97), pp. 43-48, 1997.
- Russell, S. and Norvig, P., Artificial Intelligence -A Modern Approach-, Prentice-Hall, 1995.
- 22) Salton, G. and Buckley, C., "Term-weighting Approaches in Automatic Text Retrieval," in *Readings in Information Retrieval* (Jones, K. S. and Willet, P., eds.), *Morgan Kaufmann*, pp. 323–328. Morgan Kaufmann, 1997.
- 23) Selberg, E. and Etzioni, O., "Multi-service Search and Comparison Using the Metacrawler," in *the 1995 World Wide Web Conference*, 1995.

- 24) Selberg, E. and Etzioni, O., "the Metacrawler Architecture for Resource Aggregation on the Web," in *IEEE Expert*, IEEE, January-February, pp. 11–14, 1997.
- 25) Yamada, S. and Ohsawa, Y., "Navigation Planning to Guide Concept Understanding in the World Wide Web," in *Proc. of the Fourth International Conference* on Autonomous Agent, pp. 114-115, 2000.
- 26) Yamada, S. and Osawa, Y., "Planning to Guide Concept Understanding in the WWW," in AAAI 1998 Workshop on AI and Information Integration, pp. 121–126, 1998.
- 27) Yamana, H., Tamura, K., Kawano, H., Kamei, S., Harada, M., Nishimura, H., Asai, I., Kusumoto, H., Shinoda, Y. and Muraoka, Y., "Expreriments of Collecting WWW Information Using Distributed WWW Robots," in *Proc. of SIGIR'98*, pp. 379–380, Melbourne, Australia, 1998.
- 28) Yates, R. B. and Neto, B. R., Modern Information Retrieval, Addison Wesley, 1999.



Seiji Yamada, Dr.Eng.: He received the B.S., M.S. and Ph.S. degrees in control engineering and artificial intelligence from Osaka University, Osaka, Japan, in 1984, 1986 and 1989, respectively. From 1989 to 1991, he served as a Research Associate in the Department of Control Engineering at Osaka University. From 1991 to 1996, he served as a Lecturer in the Institute of Scientific and Industrial Research at Osaka University. In 1996, he joined the Department of Computational Intelligence and Systems Science at Tokyo Institute of Technology, Yokohama, Japan, as an Associate Professor. His research interests include artificial intelligence, planning, machine learning for a robotics, intelligent information retrieval in the WWW, human computer interaction. He is a member of AAAI, IEEE, JSAI, RSJ and IEICE.



**Hiroyuki Kawano, Dr.Eng.:** He is an Associate Professor at the Department of Systems Science, Graduate School of Informatics, Kyoto University, Japan. He obtained his B.Eng. and M.Eng. degrees in Applied Mathematics and Physics, and his Dr.Eng. degree in Applied Systems Science from Kyoto University. His research interests are in advanced database technologies, such as data mining, data warehousing, knowledge discovery and web search engine (Mondou). He has served on the program committees of several conferences in the areas of Data Base Systems, and technical committees of advanced information systems.