

Query based Personalization in Semantic Web Mining

Mahendra Thakur
Department of CSE
Samrat Ashok Technological
Institute
Vidisha, M.P., India

Yogendra Kumar Jain
Department of CSE
Samrat Ashok Technological
Institute
Vidisha, M.P., India

Geetika Silakari
Department of CSE
Samrat Ashok Technological
Institute
Vidisha, M.P., India

Abstract— To provide personalized support in on-line course resources system, a semantic web-based personalized learning service is proposed to enhance the learner's learning efficiency. When a personalization system relies solely on usage-based results, however, valuable information conceptually related to what is finally recommended may be missed. Moreover, the structural properties of the web site are often disregarded. In this Paper, we present a personalize Web search system, which can helps users to get the relevant web pages based on their selection from the domain list. In the first part of our work we present Semantic Web Personalization, a personalization system that integrates usage data with content semantics, expressed in ontology terms, in order to compute semantically enhanced navigational patterns and effectively generate useful recommendations. To the best of our knowledge, our proposed technique is the only semantic web personalization system that may be used by non-semantic web sites. In the second part of our work, we present a novel approach for enhancing the quality of recommendations based on the underlying structure of a web site. We introduce UPR (Usage-based Page Rank), a Page Rank-style algorithm that relies on the recorded usage data and link analysis techniques based on user interested domains and user query.

Keywords-Semantic Web Mining; Personalized Recommendation; Recommended System

I. INTRODUCTION

Comparing with the traditional face-to-face learning style, e-learning is indeed a revolutionary way to provide education in the life-long term. However, different learners have different learning styles, goals, previous knowledge and other preferences; the traditional "one-size-fits-all" learning method is no longer enough to satisfy the needs of learners. Nowadays more and more personalized systems have been developed and are trying to find a solution to the personalization of the learning process, which affect the learning function outcome. The Semantic

Web is not a separate web but an extension of the current one, in which information is given well-defined meaning, and better enabling computers and people to work in cooperation [1]. Under the conditions of Semantic Web-based learning system the learning information is well-defined, and the machine can understand and deal with the semantics for the learning contents to provide adaptable learning services with a powerful technical support.

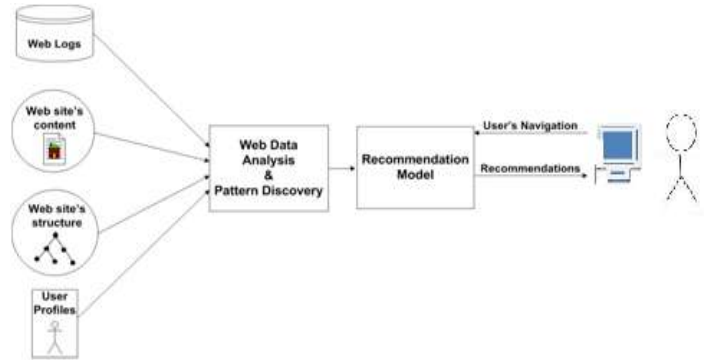


Figure: 1 the web personalization process

The problem of providing recommendations to the visitors of a web site has received a significant amount of attention in the related literature. Most of the research efforts in web personalization correspond to the evolution of extensive research in web usage mining, taking into consideration only the navigational behavior of the (anonymous or registered) visitors of the web site. Pure usage-based personalization, however, presents certain shortcomings. This may happen when, for instance, there is not enough usage data available in order to extract patterns related to certain navigational actions, or when the web site's content changes and new pages are added but are not yet included in the web logs. Moreover, taking into consideration the temporal characteristics of the web in terms of its usage, such systems are very vulnerable to the training data used to construct the predictive model. As a result, a number of research approaches integrate other sources of information, such as the web content or the web structure in order to enhance the web personalization process [1] and [2].

As already implied, the users' navigation is largely driven by semantics. In other words, in each visit, the user usually aims at finding information concerning a particular subject. Therefore, the underlying content semantics should be a dominant factor in the process of web personalization. The web site's content characterization process involves the feature extraction from the web pages. Usually these features are keywords subsequently used to retrieve similarly characterized content. Several methods for extracting keywords that characterize web content have been proposed. The similarity between documents is usually based on exact matching between these terms. This way, however, only a binary matching between documents is achieved, whereas no actual *semantic* similarity is taken into consideration. The need for a

more abstract representation that will enable a uniform and more flexible document matching process imposes the use of semantic web structures, such as ontology's. By mapping the keywords to the concepts of an ontology, or topic hierarchy, the problem of binary matching can be surpassed through the use of the hierarchical relationships and/or the *semantic similarities* among the ontology terms, and therefore, the documents. Finally, we should take into consideration that the web is not just a collection of documents browsed by its users. The web is a directed labeled graph, including a plethora of hyperlinks that interconnect its web pages. Both the structural characteristics of the web graph, as well as the web pages' and hyper links' underlying semantics are important and determinative factors in the users' navigational process. The main contribution of this paper is a set of novel techniques and algorithms aimed at improving the overall effectiveness of the web personalization process through the integration of the content and the structure of the web site with the users' navigational patterns. In the first part of our work we present the semantic web personalization system for Semantic Web Personalization that integrates usage data with content semantics in order to compute semantically enhanced navigational patterns and effectively generate useful recommendations. Similar to previously proposed approaches, the proposed personalization framework uses ontology terms to annotate the web content and the users' navigational patterns. The key departure from earlier approaches, however, is that Semantic Web Personalization is the only web personalization framework that employs automated keyword-to-ontology mapping techniques, while exploiting the underlying semantic similarities between ontology terms. Apart from the novel recommendation algorithms we propose, we also emphasize on a hybrid structure-enhanced method for annotating web content. To the best of our knowledge, Semantic Web Personalization is the only semantic web personalization system that can be used by any web site, given only its web usage logs and a domain-specific ontology [3] and [4].

II. BACKGROUND

The main data source in the web usage mining and personalization process is the information residing on the web site's logs. Web logs record every visit to a page of the web server hosting it. The entries of a web log file consist of several fields which represent the date and the time of the request, the IP number of the visitor's computer (client), the URI requested, the HTTP status code returned to the client, and so on. The web logs' file format is based on the so called "extended" log format.

Prior to processing the usage data using web mining or personalization algorithms, the information residing in the web logs should be preprocessed. The web log data preprocessing is an essential phase in the web usage mining and personalization process. An extensive description of this process can be found. In the sequel, we provide a brief overview of the most important pre-processing techniques, providing in parallel the related terminology. The first issue in the *pre-processing* phase is *data preparation*. Depending on the application, the web log data may need to be cleaned from entries involving page accesses that returned, for example, an error or graphics file accesses. Furthermore, crawler activity usually should be

filtered out, because such entries do not provide useful information about the site's usability. A very common problem to be dealt with has to do with web pages' *caching*. When a web client accesses an already cached page, this access is not recorded in the web site's log. Therefore, important information concerning web path visits is missed. Caching is heavily dependent on the client-side technologies used and therefore cannot be dealt with easily. In such cases, cached pages can usually be inferred using the referring information from the logs and certain heuristics, in order to re-construct the user paths, filling out the missing pages. After all page accesses are identified, the *page view identification* should be performed. A *page view* is defined as "the visual rendering of a web page in a specific environment at a specific point in time". In other words, a page view consists of several items, such as frames, text, graphics and scripts that construct a single web page. Therefore, the page view identification process involves the determination of the distinct log file accesses that contribute to a single page view. Again such a decision is application-oriented. In order to personalize a web site, the system should be able to distinguish between different users or groups of users. This process is called *user profiling*. In case no other information than what is recorded in the web logs is available, this process results in the creation of aggregate, anonymous user profiles since it is not feasible to distinguish among individual visitors. However, if the user's registration is required by the web site, the information residing on the web log data can be combined with the users' demographic data, as well as with their individual ratings or purchases. The final stage of log data pre-processing is the partition of the web log into distinct *user* and *server sessions*. A user session is defined as "a delimited set of user clicks across one or more web servers", whereas a *server session*, also called a *visit*, is defined as "a collection of user clicks to a single web server during a user session". If no other means of session identification, such as cookies or session ids is used, *session identification* is performed using time heuristics, such as setting a minimum timeout and assumes that consecutive accesses within it belong to the same session, or a maximum timeout, assuming that two consecutive accesses that exceed it belong to different sessions [1] and [5] and [6].

A. Web Usage Mining and Personalization:

Web usage mining is the process of identifying representative trends and browsing patterns describing the activity in the web site, by analyzing the users' behavior. Web site administrators can then use this information to redesign or customize the web site according to the interests and behavior of its visitors, or improve the performance of their systems. Moreover, the managers of e-commerce sites can acquire valuable business intelligence, creating consumer profiles and achieving market segmentation. There exist various methods for analyzing the web log data. Some research studies use well known data mining techniques such as association rules discovery, sequential pattern analysis, clustering, probabilistic models, or a combination of them. Since web usage mining analysis was initially strongly correlated to data warehousing, there also exist some research studies based on OLAP cube models. Finally some proposed web usage mining approaches that require registered user profiles, or combine the usage data

with semantic meta-tags incorporated in the web site's content. Furthermore, this knowledge can be used to automatically or semi-automatically adjust the content of the site to the needs of specific groups of users, i.e. to personalize the site. As already mentioned, web personalization may include the provision of recommendations to the users, the creation of new index pages, or the generation of targeted advertisements or product promotions. The usage-based personalization systems use association rules and sequential pattern discovery, clustering, Markov models, machine learning algorithms, or are based on collaborative filtering in order to generate recommendations. Some research studies also combine two or more of the aforementioned techniques [2] and [4].

B. Integrating Content Semantics in Web Personalization:

Several frameworks supporting the claim that the incorporation of information related to the web site's content enhances the web personalization process have been proposed prior or subsequent to our work. In this Section we overview in detail the ones that are more similar to ours, in terms of using a domain-ontology to represent the web site's content.

Dai and Mobasher proposed a web personalization framework that uses ontologies to characterize the usage profiles used by a collaborative filtering system. These profiles are transformed to "domain-level" aggregate profiles by representing each page with a set of related ontology objects. In this work, the mapping of content features to ontology terms is assumed to be performed either manually, or using supervised learning methods. The defined ontology includes classes and their instances therefore the aggregation is performed by grouping together different instances that belong to the same class. The recommendations generated by the proposed collaborative system are in turn derived by binary matching of the current user visit, expressed as ontology instances, to the derived domain-level aggregate profiles, and no semantic similarity measure is used. The idea of semantically enhancing the web logs using ontology concepts is independently described in recent. This framework is based on a semantic web site built on an underlying ontology. The authors present a general framework where data mining can then be performed on these semantic web logs to extract knowledge about groups of users, users' preferences, and rules. Since the proposed framework is built on a semantic web knowledge portal, the web content is already semantically annotated focuses solely on web mining and thus does not perform any further processing in order to support web personalization.

In recent (through the existing RDF annotations), and no further automation is provided. Moreover, the proposed framework t work also proposes a general personalization framework based on the conceptual modeling of the users' navigational behavior. The proposed methodology involves mapping each visited page to a topic or concept, imposing a concept hierarchy (taxonomy) on these topics, and then estimating the parameters of a semi-Markov process defined on this tree based on the observed user paths. In this Markov models-based work, the semantic characterization of the content is performed manually. Moreover, no semantic similarity measure is exploited for enhancing the prediction process, except for generalizations/specializations of the

ontology terms. Finally, in a subsequent work, explore the use of ontologies in the user profiling process within collaborative filtering systems. This work focuses on recommending academic research papers to academic staff of a University. The authors represent the acquired user profiles using terms of research paper ontology (is-a hierarchy). Research papers are also classified using ontological classes. In this hybrid recommender system which is based on collaborative and content-based recommendation techniques, the content is characterized with ontology terms, using document classifiers (therefore a manual labeling of the training set is needed) and the ontology is again used for making generalizations/specializations of the user profiles [7] and [8] and [9].

C. Integrating Structure in Web Personalization:

Although the connectivity features of the web graph have been extensively used for personalizing web search results, only a few approaches exist that take them into consideration in the web site personalization process. To use citation and coupling network analysis techniques in order to conceptually cluster the pages of a web site. The proposed recommendation system is based on Markov models. In previous, use the degree of connectivity between the pages of a web site as the determinant factor for switching among recommendation models based on either frequent item set mining or sequential pattern discovery. Nevertheless, none of the aforementioned approaches fully integrates link analysis techniques in the web personalization process by exploiting the notion of the *authority* or *importance* of a web page in the web graph.

In a very recent work, address the data sparsity problem of collaborative filtering systems by creating a bipartite graph and calculating linkage measures between unconnected pairs for selecting candidates and make recommendations. In this study the graph nodes represent both users and rated/purchased items.

Finally, subsequent work, proposed independently two link analysis ranking methods, *Site Rank* and *Popularity Rank* which are in essence very much like the proposed variations of our *UPR* algorithm (*PR* and *SUPR* respectively). This work focuses on the comparison of the distributions and the rankings of the two methods rather than proposing a web personalization algorithm [9] and [10].

III. PROPOSED TECHNIQUE

In this paper, we present Semantic Enhancement for Web Personalization, a web personalization framework that integrates content semantics with the users' navigational patterns, using ontologies to represent both the content and the usage of the web site. In our proposed framework we employ web content mining techniques to derive semantics from the web site's pages. These semantics, expressed in ontology terms, are used to create semantically enhanced web logs, called C-logs (concept logs). Additionally, the site is organized into thematic document clusters. The C-logs and the document clusters are in turn used as input to the web mining process, resulting in the creation of a broader, semantically enhanced set of recommendations. The whole process bridges the gap between Semantic Web and Web Personalization areas, to create a Semantic Web Personalization system.

A. Semantic Enhancement for Web Personalization System Architecture:

Semantic Enhancement for Web Personalization uses a combination of web mining techniques to personalize a web site. In short, the web site's content is processed and characterized by a set of ontology terms (categories). The Web personalization process include (a) The collection of Web data, (b) The modeling and categorization of these data (preprocessing phase), (c) The analysis of the collected data, and (d) The determination of the actions that should be performed. When a user sends a query to a search engine, the search engine returns the URLs of documents matching all or one of the terms, depending on both the query operator and the algorithm used by the search engine. *Ranking* is the process of ordering the returned documents in decreasing order of relevance, that is, so that the "best" answers are on the top. When the user enters the query, the query is first analyzed. The Query is given as input to the semantic search algorithm for separation of nouns, verbs, adjectives and negations and assigning weights respectively. The processed data is then given to the personalized URL Rank algorithm for personalizing the results according to the user domain, interest and need. The sorted results are those results in which the user is interested. The personalization can be enhanced by categorizing the results according to the types. Thus after building the knowledge base, the system can give use recommendation based on the similarity of the user interested domain and the user query. The recommendation procedure of the System has two steps:

- The system gives user a list of interested domains .Detect user's current interested domain.
- Based on user's current interested domain and combined his or her profile, the system will give him or her set of URLs with ranking scores.

In this way, the system could help the user to retrieve his or her potential interested domains. Besides, a user can change his or her current interested domain by clicking the interested domain list on the same page but with more convenience. In the beginning, if the user does not have a profile in the database, the system displays the user available domains, and then keeps a track of the user's selections. The user's selections is used to construct a table that uses URL weight calculation. The current interested domains recommendation is based on last selections. The figure 2 shows the complete process.

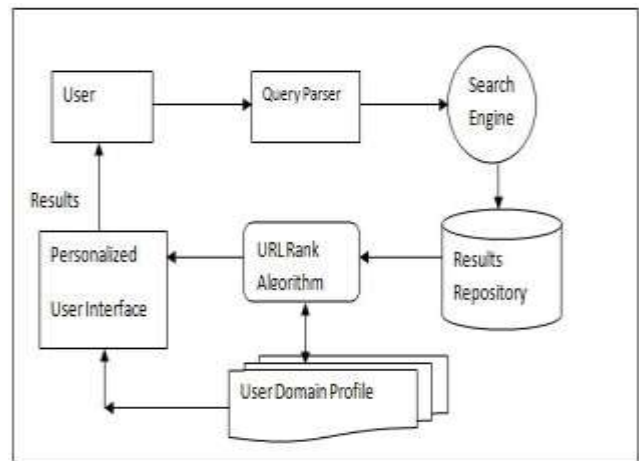


Figure 2: Web Personalization architecture

B. Recommendation process:

The learner's implicit query defined previously under both of its shapes constitutes the input of the recommendation phase. The recommendation process task is accomplished using basically: content based filtering (CBF) and collaborative filtering (CF) approaches (Figure 3). First, we apply the (CBF) approach alone using the search functionalities of the search engine. We submit the term vector to the search engine in order to compute recommendation links. Results are ranked according to the cosine similarity of their content (vector of *TF-IDF* weighted terms) with the submitted term vector. Second, we apply the collaborative approach (CF) alone by comparing, first, the sliding window pages to clusters (groups of learners obtained in the offline phase by applying two-level model based collaborative filtering approach) in order to classify the active learner in one of the learner's group. Then, we use the *ARs* of the corresponding group to give personalized recommendations. The current session window is matched against the "condition" or left side of each rule.

It is worth noting that several recommendation strategies using these approaches have been investigated in our work. After applying a CF and CBF approaches alone, we included next the possibility to combine both of the recommendation approaches (CBF and CF) in order to improve the recommendation quality and generate the most relevant learning objects to learners. Hence, two approaches are to be considered: Hybrid content via profile based collaborative filtering with cascaded/feature augmentation combination, which performs collaborative recommendation followed by content recommendation (the reverse order could also be considered); and Hybrid content and profile based collaborative filtering with weighted combination, where the collaborative filtering and content based filtering recommendations are performed simultaneously, then the results of both techniques are combined together to produce a single recommendation set. In the Hybrid content via profile based collaborative filtering with cascaded/feature augmentation combination approach, we apply first CF approach giving as output a set of recommended links, then we apply CBF approach on these links. In fact, recommended links are mapped to a set of content terms in

order to compose a term vector (top k frequent terms), a parser tool must be used for this task. Finally, these terms are submitted to the search engine which returns the final recommended links.

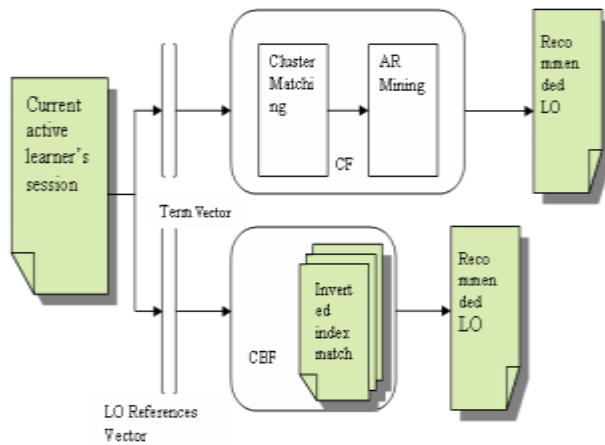


Figure 3: Recommendation process

In the Hybrid content and profile based collaborative filtering with weighted combination approach, the collaborative filtering and content based filtering are performed separately, then the results of both techniques are combined together to produce a single recommendation set.

C. This process uses the following steps:

I. Step 1 is performed in the same way as in CF approach; the result is called Recommended Set 1;

II. Step 2 maps each LO references in the sliding window to a set of content terms (top k frequent terms). Then these terms are submitted to the search engine which returns recommended links. This result is called Recommended Set 2;

III. Final collaborative and content based filtering recommendation combination: both recommended sets obtained previously are combined together to form a coherent list of related recommendation links, which are ranked based on their overlap ratio.

IV. METHODOLOGY

- **Data Set** The two key advantages of using this data set are that the web site contains web pages in several formats (such as pdf, html, ppt, doc, etc.), written both in Greek and English and a domain-specific concept hierarchy is available (the web administrator created a concept-hierarchy of 150 categories that describe the site's content). On the other hand, its context is rather narrow, as opposed to web portals, and its visitors are divided into two main groups: students and researchers. Therefore, the subsequent analysis (e.g. association rules) uncovers these trends: visits to course material, or visits to publications and researcher details. It is essential to point out that the need for processing online (up-to-date) content, made it impossible for us to use other publicly available web log

sets, since all of them were collected many years ago and the relevant sites' content is no longer available. Moreover, the web logs of popular web sites or portals, which would be ideal for our experiments, are considered to be personal data and are not disclosed by their owners. To overcome these problems, we collected web logs over a 1-year period (01/01/10 – 31/12/10). After preprocessing, the total web logs' size was approximately 105 hits including a set of over 67.700 distinct anonymous user sessions on a total of 360 web pages. The sessionizing was performed using distinct IP & time limit considerations (setting 20 minutes as the maximum time between consecutive hits from the same user).

- **Keyword Extraction: Category Mapping:** We extracted up to 7 keywords from each web page using a combination of all three methods (raw term frequency, inlinks, outlinks). We then mapped these keywords to ontology categories and kept at most 5 for each page.
- **Document Clustering:** We used the clustering scheme described in recent, i.e. the DBSCAN clustering algorithm and the similarity measure for sets of keywords. However, other web document clustering schemes (algorithm & similarity measure) may be employed as well.
- **Association Rules Mining:** We created both URI-based and category-based frequent item sets and association rules. We subsequently used the ones over a 40% confidence threshold.

V. RESULTS

In our paper work we compare the performance of the three ranking methods based on pure similarity, plain Page Rank and weighted (personalized) URL Rank.

The personalization accuracy was found to be 75%; the random search accuracy is 74.6 %. The average of personalization accuracy is 74.7%. Because the interested domains personalization is done considering the user selected domain, the accuracy is higher than the random recommendation in our experiment. Above Fig. 4 is a comparison of the interested domains personalization accuracy based on random selection and based on our personalization method. Figure 4 shows Relevance Query Results vs. Random & Personalization Selection graph.

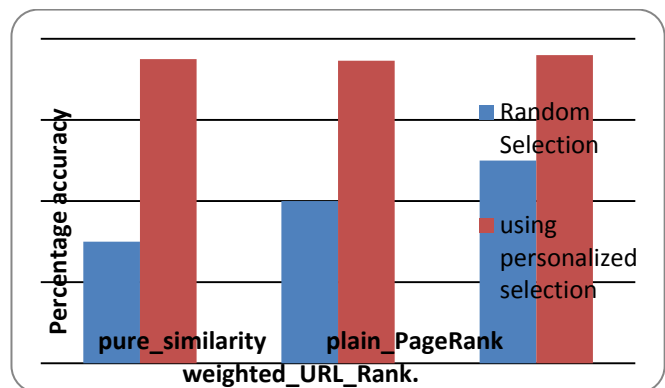


Figure 4 – Random Selection accuracy

The URL personalization accuracy based on the interested domains selection is 71.3%; and the URL personalization accuracy without the interested domains selection assistance is 31.9 % in Fig. 5. From this result, we can see that the interested domains recommendation help the system to filter lots of URLs that the user might not be interested in. Moreover, the system could focus on the domains that users are interested in to select the relevant URL. Figure 5 shows Relevance Query Results vs. Random & Personalization Selection graph.

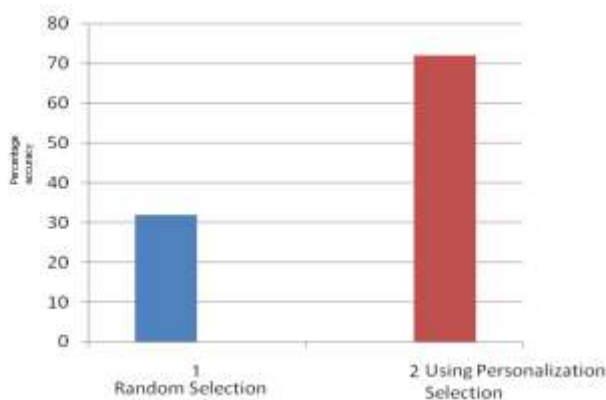


Figure 5- Personal accuracy in interested domain

VI. CONCLUSION

In this paper contribution is a core technology and reusable software engine for rapid design of a broad range of applications in the field of personalized recommendation systems and more. We present a web personalization system for web search, which not only gives user a set of personalized pages, but also gives user a list of domains the user may be interested in. Thus, user can switch to different interests when he or she is surfing on the web for information. Besides, the system focuses on the domains that the user is interested in, and won't waste lots of time on searching the information in the irrelevant domains. Moreover, the recommendation won't be affected by the irrelevant domains, and the accuracy of the recommendation is increased.

REFERENCES

- [1] Changqin Huang, Ying Ji, Rulin Duan, "A semantic web-based personalized learning service supported by on-line course resources", 6th IEEE International Conference on Networked Computing (INC), 2010.
- [2] V. Gorodetsky, V. Samoylov, S. Serebryakov, "Ontology-based context-dependent personalization technology", IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, vol. 3, pp. 278-283, 2010.
- [3] Pasi Gabriella, "Issues on preference-modelling and personalization in information retrieval", IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, pp. 4, 2010.
- [4] Wei Wenshan and Li Haihua, "Base on rough set of clustering algorithm in network education application", IEEE International Conference on Computer Application and System Modeling (ICCASM 2010), vol. 3, pp. V3-481 - V3-483, 2010.
- [5] Shuchih Ernest Chang, and Chia-Wei Wang, "Effectively generating and delivering personalized product information: Adopting the web 2.0 approach", 24th IEEE International Conference on Advanced Information Networking and Applications Workshops (WAINA), pp. 401-406, 2010.

- [6] Xiangwei Mu, Van Chen, and Shuyong Liu, "Improvement of similarity algorithm in collaborative filtering based on stability degree", 3rd International Conference on Advanced Computer Theory and Engineering (ICACTE), vol. 4, pp. V4-106 - V4-110, 2010
- [7] Dario Vuljani, Lidia Rovani, and Mirta Baranovi, "Semantically enhanced web personalization approaches and techniques", 32nd IEEE International Conference on Information Technology Interfaces (ITA), pp. 217-222, 2010.
- [8] Raymond Y. K. Lau, "Inferential language modeling for selective web search personalization and contextualization", 3rd IEEE International Conference on Advanced Computer Theory and Engineering (ICACTE), vol. 1, pp. V1-540 - V1-544, 2010.
- [9] Esteban Robles Luna, Irene Garrigos, and Gustavo Rossi, "Capturing and validating personalization requirements in web applications", 1st IEEE International Workshop on Web and Requirements Engineering (WeRE), pp. 13-20, 2010.
- [10] B. Annappa, K. Chandrasekaran, K. C. Shet, "Meta-Level constructs in content personalization of a web application", IEEE International conference on Computer & Communication Technology-ICCT'10, pp. 569 - 574, 2010.
- [11] F. Murtagh, "A survey of recent advances in hierarchical clustering algorithms", Computer Journal, vol. 26, no. 4, pp. 354 -359, 1983.
- [12] Wang Jicheng, Huang Yuan, Wu Gangshan, and Zhang Fuyan, "Web mining: Knowledge discovery on the web system", IEEE International Conference systems, Man and cybernatics, vol.2, pp. 137 - 141, 1999.
- [13] B. Mobasher, "Web usage mining and personalization in practical handbook of internet computing", M.P. Singh, Editor. 2004, CRC Press, pp. 15.1-37.
- [14] T. Maier, "A formal model of the ETL process for OLAP-based web usage analysis", 6th WEBKDD- workshop on Web Mining and Web Usage Analysis", part of the ACM KDD: Knowledge Discovery and Data Mining Conference, pp. 23-34, Aug. 2004
- [15] R. Meo, P. Lanzi, M. Matera, R. Esposito, "Integrating web conceptual modeling", WebKDD, vol. 3932, pp. 135-148, 2006.
- [16] B. Mobasher, R. Cooley, and J. Srivastava, "Automatic personalization based on web usage mining", Communications of the ACM, vol. 43, no. 8, pp. 142-151, Aug 2000.
- [17] B. Mobasher, H. Dai, T. Luo, and M. Nakagawa, "Effective personalization based on association rule discovery from web usage data", 3rd international ACM Workshop on Web information and data management, 2001.
- [18] O. Nasraoui, R. Krishnapuram, and A. Joshi, "Mining web access logs using a relational clustering algorithm based on a robust estimator", 8th International World Wide Web Conference, pp. 40-41, 1999.
- [19] D. Pierrakos, G. Paliouras, C. Papatheodorou, V. Karkaletsis, and M. Dikaiakos, "Web community directories: A new approach to web personalization", 1st European Web Mining Forum (EWMF'03), vol. 3209, pp. 113-129, 2003.
- [20] Schafer J. B., Konstan J., and Reidel J., "Recommender systems in e-commerce", 1st ACM Conference on Electronic commerce, pp. 158-166. 1999

AUTHORS PROFILE



Dr. Yogendra Kumar Jain presently working as head of the department, Computer Science & Engineering at Samrat Ashok Technological Institute Vidisha M.P India. The degree of B.E. (Hons) secured in E&I from SATI Vidisha in 1991, M.E. (Hons) in Digital Tech. & Instrumentation from SGSITS, DAVV Indore(M.P), India in 1999. The Ph. D. degree has been awarded from Rajiv Gandhi Technical University, Bhopal (M.P.) India in 2010.

Research Interest includes Image Processing, Image compression, Network Security, Watermarking, Data Mining. Published more than 40 Research papers in various Journals/Conferences, which include 10 research papers in International Journals. Tel:+91-7592-250408, E-mail: ykjain_p@yahoo.co.in.

Geetika Silakari presently working as Asst. Professor in Computer Science & Engineering at Samrat Ashok Technological Institute Vidisha M.P India. The degree of B.E. (Hons) secured in Computer Science & engineering. She secured M.Tech in Computer science and Engineering from Vanasthali University. She is currently pursuing PHD in Computer science and engineering. E-mail:geetika.silakari@gmail.com



Mr. Mahendra Thakur is a research scholar pursuing M.Tech in Computer Science & Engineering from Samrat Ashok Technological Institute Vidisha M.P India. He secured degree of B.E. in IT from Rajiv Gandhi Technical University, Bhopal (M.P.) India in 2007.

E-mail-mahendrasati2010@gmail.com