

Web People Search using Ontology Based Decision Tree

Mrunal Patil¹, Sonam Khomane², Varsha Saykar,³ Kavita Moholkar⁴

^{*}(Department of Computer Engineering, Rajarshi Shahu College of Engineering, Pune-411033, India.

^{**}(Department of Computer Engineering, Rajarshi Shahu College of Engineering, Pune-411033, India.

^{***}(Department of Computer Engineering, Rajarshi Shahu College of Engineering, Pune-411033, India.

^{****}(Lecturer, Department of Computer Engineering, Rajarshi Shahu College of Engineering, Pune-411033, India.

Abstract

Internet plays an important role in search of information. So far, most of the search done on the internet is related to person search. When we give a query for person search, it returns a set of web pages related to distinct person of given name. For such type of search the job of finding the web page of interest is left on the user. In this paper, we develop a technique for web people search which clusters the web pages based on semantic information and maps them using ontology based decision tree making the user to access the information in more easy way. This technique uses the concept of ontology thus reducing the number of inconsistencies. The result proves that ontology based decision tree and clustering helps in increasing the efficiency of the overall search.

Keywords- Clustering , Efficiency , Ontology based decision tree, Semantic information , Web people search.

1. INTRODUCTION

According to Samuel Johnson "Knowledge is of two kinds: we know a subject ourselves, or we know where we can find information about it" [6]. For years, these words remained true to those in search for information. Although not always readily available, information about person could be easily found by individuals in directories, indexed by human experts. With the advent of the Internet, however, the situation changed dramatically. Enormous amounts of data are now freely accessible to over 600 millions of users on-line. The question of how to find information of interested person on the Internet is raised by the Web People Search problem.

To solve the problem various algorithms are been adopted by the current search engines like the Google, yahoo etc. Queries are fired to such search engine. The result returned is in the form of rank list of documents along with partial content leaving the job of finding the specific document onto the user. Thus to find the relevant document user has to shift through a large set of irrelevant document.

Our approach exploits the task of searching the people in more precise manner. Web people search clusters the web pages based on the query fired. Clustering is the process of gathering similar objects in one cluster which are different to the objects

in another cluster. Web pages having semantic information are grouped in a single cluster. Later, these clusters are presented to

the user in form of ontology. Ontology is a set of concepts such as things and relations that are specified in some way in order to exchange information. Firstly input is given to the search engine in form of query. The search engine returns the top k relevant pages as soon as the query is fired. Web pages are retrieved are processed. The pre processed web pages are clusters and then presented using decision tree. The system can be made more generalized for the search of people with the help of filters. The main objective of this paper is to understand the basic concepts of ontology with particular emphasis on its application to people search problem. For this purpose, the paper contains the review of past work in part II and explains the general concepts behind ontology in Part III, followed by a general description of the system in Part IV. Part V provides the experimental results. Part VI concludes along with the future work and references.

2. Review Of The Past Work

Initially, for searching the Page Rank algorithms were used for ranking the search query results. The page rank algorithm was described by Lawrence Page and Sergey Brin. Working of the Page Rank algorithm depends upon link structure of the web pages. The Page Rank algorithm is based on the concepts that if a page contains important links towards it then the links of this page towards the other page are also to be considered as important pages. The Page Rank considers the back link in deciding the rank score. But the limitation is that if new page is inserted between two pages then the crawler should perform a large calculation to calculate the distance vector which is a time consuming process and decreases the performance. Taher Haveliwala in 2002 proposed a Topic Sensitive Page Rank as compared to the original Page Rank for improving the search-query results where a single Page Rank vector is computed using the link structure of the web to compute the relative importance of the page [2].

Hearst and Pedersen showed that relevant documents tend to be more similar to each other, thus the clustering of similar search results helps users find relevant results. In addition, Vivisimo is a real demonstration of this technique. Vivisimo was founded in 2000 by three Carnegie Mellon University scientists who decided to tackle the problem of information overload in web search. Rather than focusing just on search engine result ranking, we realized that grouping results into topics, or "clouds," made for better search and discovery. As search

became a necessity for web users, Vivisimo developed a service robust enough to handle the variety of information the everyday web user was after. The result was Clusty: an innovative way to get more out of every search. Clusty was acquired by Yippy, Inc. in May 2010. Yippy queries several top search engines, combines the results, and generates an ordered list based on comparative ranking.

Jargon Went and Jian-Yun Nie proposed a new approach to query clustering using user logs[3]. The principles are as follows. 1) If users clicked on the same documents for different queries, then the queries are similar. 2) If a set of documents is often selected for a set of queries, then the terms in these documents are related to the terms of the queries to some extent. These principles are used in combination with the traditional approaches based on query contents.

Our proposed local-cluster algorithm considers linkage structure and content generation of cluster structures to produce a ranking of the underlying clusters with respect to a user's given search query and preference. The rank of each document is then obtained through the relation of the given document with respect to its relevant clusters and the respective preference of these clusters.

3. Ontology

Ontology specifies linked concepts and terms and relations among these terms and concepts. Concepts are nothing but the entities which are language independent. Ontology shows how each concept is related what properties it has. The main purpose of ontology based decision tree is to give a more meaningful, descriptive and a readable view of concepts.

Mathematically ontology can be defined Yang *et al.*, 2008 [5] as follows:

“An ontology can be defined as an Vector $O = (C, V, P, H, \text{ROOT})$, where C is the set of concepts, V contains a set of terms and is called the vocabulary, P is the set of properties for each concept, H is the hierarchy and ROOT is the topmost concept. Concepts are taxonomically related by the directed, acyclic, transitive, reflexive relation H belongs to $C * C$. $H(c1, c2)$ shows that c1 is a subclass of c2 and for all c belongs to C it holds that $H(c, \text{ROOT})$.”

Our goal is to utilize the user context to the search results by re-ranking the results returned from the given query of search engine. The representation of the tree depends upon the user's information access behavior. Semantic information is fundamental part of user context. The best example of ontological approach has proven to be successful in the recommender system which does not consider the domain knowledge.

Ontology consists of hierarchy of classes and sub-classes for object-entity [3]. The clusters will be taken as input which is formed using the lingo algorithm, later arranged in hierarchy.

This could be achieved as follows:

- The topmost concept is the root having intermediate and leaf concepts. If a user wishes for a leaf concept then each cluster starting from the topmost concept will be traversed.

- The probability for each cluster will be calculated and depending on that clusters of user's interest will be accessed. Long with this a threshold value will be maintained.
- If the user heating ratio for a given cluster is greater and also the probability for each cluster is greater than the threshold value then parent is renamed by child name.

Mathematically it can be defined as follows:

$$P_b = (100 * HR) / \text{Total no. of clusters viewed}$$

Where, P_b = Probability

HR = hitting ratio

An example of basic ontology is shown below in the Fig 1.

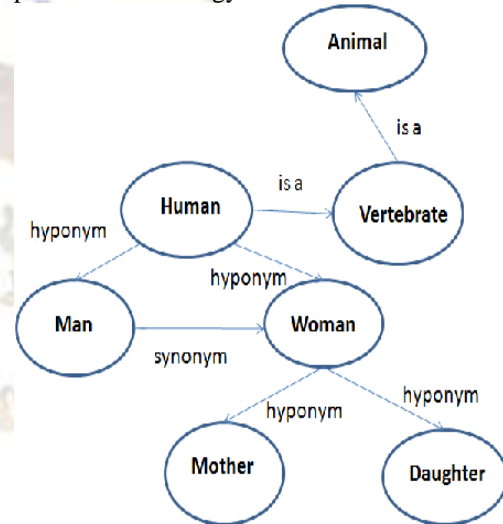


Figure 1: Example of Ontology

Ontology explains a hierarchy of classes, sub classes and their relationships. The above example describes the “HIERARCHY” for a human class. It shows that a human is a vertebrate where the man and woman are synonym and they are hyponym for class human.

4. PROPOSED SYSTEM

The system mainly works in different phases as:

4.1 Search result fetching:

The user submits the query to the search engine. The filters are used to find whether the given query is related to person search or not. Then we get the WebPages of search result lists returned by a Google web search engine. So the first search is the conventional met search based on these keywords. These WebPages are analyzed by an HTML parser and the result items are extracted.

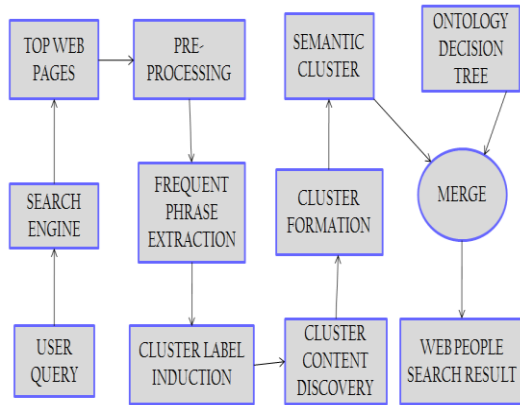


Figure 2: Proposed system

Generally, there are only titles and query-dependent snippets available in each result item. We assume these contents are informative enough because most search engines are well designed to facilitate users' relevance judgment only by the title and snippet, thus it is able to present the most relevant contents for a given query. Each extracted phrase is in fact the name of a candidate cluster, which corresponds to a set of documents that contain the phrase.

4.2 Ontology Based Web People Search:

When designing a Cluster Based Web Search, special attention must be paid to ensuring that both content and description (labels) of the resulting groups are meaningful to humans. There are various algorithms such as K means, K-medoid, Suffix tree clustering, Hierarchical clustering. There are various drawbacks of these algorithms and to overcome these we use the Lingo algorithm. Lingo reverses the process—first attempt to ensure that we can create a human-perceivable cluster label and only then assign documents to it. Specifically, extract frequent phrases from the input documents, hoping they are the most informative source of human-readable topic descriptions. Finally, match group descriptions with the extracted topics and assign relevant documents to them.

Our novel algorithm, Lingo clusters the search results. Unlike other algorithms lingo first discovers the name of the clusters and then evaluate each cluster with appropriate name. Basically lingo consists of five phases. The first phase deals with preprocessing to data. It detects the word boundaries and applies the stemming and stop word removal. The second phase deals with the frequent phrase extraction here a term appearing certain number of times is discovered and combined in all set of documents. Phase three is the cluster label induction. SVD is used to extract orthogonal vectors of the term document matrix, believed to represent distinct topic in the input data [6]. The fourth phase is the cluster content discovery phase. In this phase a Vector Space Model is applied to put the input documents under the cluster label discovered in the previous phase. Highest scoring documents for each cluster are assigned as that cluster's content [7]. Last phase is the i.e. the fifth phase is final cluster formation. Later, ontology based decision tree is referred for

rearranging the clusters formed. The steps that will be followed are:

- The ontology tree is formed according to the hitting ratio. Each object forms here the separate cluster.
- After the cluster formation, weight is given to each node of the tree and threshold value is maintained.
- Rearrange the clusters considering the weight given to each cluster and then display them to the users. Procedure that calculates the similarity between objects and clusters that estimate the similarity between clusters and ontology objects are used for this purpose.
- If desired number of clusters are obtained, then stop else goto step a.

5. Experiment Results

This section deals with experimental results when applied on web people search using ontology based decision tree. We have searched for the different persons with their relations like Abdul Kalam as the president of India, and listed the results obtained from various search engines. The results obtained from our search engine where compared with traditional search engine likes Google, yahoo etc. Then these results were used to find the overall efficiency and accuracy of the search engines to give relevant results. We found that the accuracy for our search engine i.e. Web People Search Engine is good than other search engines.

The results are listed in the following table:

Search Engine Query:	Google Results Obtained	Yahoo Results Obtained	Web People Search Results obtained (Clusters)
A.P. J. Abdul Kalam as president	<ol style="list-style-type: none"> 1. A.P.J. Abdul Kalam - Wikipedia, the free encyclopedia 2. Dr APJ Abdul Kalam - ASITE FOR INSPIRATION AND NATION BUILDING 3. Images for Abdul Kalam 4. Abdul Kalam Abdul Kalam Biography Abdul Kalam Quotes 5. Profile: Dr. A.P.J. Abdul Kalam: Form er President of India: Speeches 6. Dr. A.P.J. Abdul Kalam: Former President of India 7. Abdul Kalam Quotes - BrainyQuote 8. Dr. A.P.J. Abdul Kalam 9. Abdul Kalam: Latest News, Photos and Videos 10. President of India : Smt. Pratibha Devi Singh Patil: Rashtrapati Bhavan 	<ol style="list-style-type: none"> 1. A.P.J. Abdul Kalam - Wikipedia, the free encyclopedia 2. A.P.J. Abdul Kalam: Information from Answers.com 3. Abdul Kalam - Image Results 4. Dr APJ Abdul Kalam - A SITE FOR INSPIRATION AND NATION BUILDING 5. Abdul Kalam - News Results 6. Dr. A.P.J. Abdul Kalam: Former President of India 7. Abdul Kalam - Video Results 8. Abdul Kalam - Wikipedia 9. Abdul Kalam: Facts, Discussion Forum, and Encyclopedia Article. 10. Abdul Kalam - Simple English Wikipedia, the free encyclopedia. 	<ol style="list-style-type: none"> 1. A.P.J. Abdul Kalam (15) 2. Abdul Kalam Quotes (13) 3. President APJ Abdul Kalam (10) 4. Amul Pabir Jambhdeen Abdul Kalam Tamil (9) 5. Tamil Hindi (9) 6. A. P. J. Abdul Kalam (7) 7. Latest News (5) 8. Abdul Kalam's Death (3) 9. Security (3) 10. Abdul Kalam Horoscope (2) 11. Madras Institute (2) 12. Notable Scientist (2) 13. Linked Kingdom LinkedIn (2) 14. Other Topics (34)
Narayana Murthy as CEO of Infosys	<ol style="list-style-type: none"> 1. N.R. Narayana Murthy - Wikipedia, the free encyclopedia 2. Infosys - N.R. Narayana Murthy Management Profiles About Us 3. N.R. Narayana Murthy Infosys Catamaran Fund Personalities 4. NR Narayana Murthy's last day as Infosys Chairman - Times Of India 5. Narayana Murthy's last day as Infosys chairman 6. Narayana Murthy Steps Down as Infosys Chairman - Outlook India 7. Narayana Murthy Steps Down as Infosys Chairman - Outlook India 8. News for Narayana murthy as ceo of infosys 9. Narayana Murthy quits CEO forum - Indian Express 10. An article on N.R. Narayana Murthy chairman Infosys Technologies Ltd. 	<ol style="list-style-type: none"> 1. N.R. Narayana Murthy - Wikipedia, the free encyclopedia 2. Infosys - N.R. Narayana Murthy Management Profiles About Us 3. Narayana Murthy's last day as Infosys Chairman - India - DNA 4. Narayana Murthy's last day as Infosys Chairman on Friday 5. Narayana Murthy's last day as Infosys chairman The Asian Age 6. Soul of Infosys will always be Narayana Murthy: iGate CEO 7. Narayan Murthy, chairman and CEO 8. Narayan Murthy, chairman and CEO 9. N.R. Narayana Murthy Infosys Catamaran Fund 10. Narayana Murthy: The man behind Infosys 	<ol style="list-style-type: none"> 1. N.R. Narayana Murthy (18) 2. Chief Mentor (15) 3. N.R. Narayana Murthy (13) 4. KV Kamath (11) 5. Narayana Murthy will Step (9) 6. New CEO (9) 7. Kris Gopalakrishnan (5) 8. Software Professionals (5) 9. Emeritus (4) 10. Asia Society (2) 11. CIOL News Reports (2) 12. Economic Times Corporate Dossier List (2) 13. Infosys with an Equity Capital of Rs.10 (2) 14. MITN Jadhav (2) 15. Other Topics (30)
Sachin Tendulkar as Cricketer	<ol style="list-style-type: none"> 1. Sachin Tendulkar - Wikipedia, the free encyclopedia 2. Sachin Tendulkar India Cricket Cricket Players and Officials 3. News for sachin tendulkar as cricketer 4. All about Sachin Tendulkar 5. Profile of Sachin Tendulkar 6. Sachin Tendulkar 1st runs in One Day Cricket -- 36 vs NZ 4th ODI 7. More videos for sachin tendulkar as cricketer 8. Sachin Tendulkar God Of Cricket Photos Pictures Records Videos 9. Sachin Tendulkar Sachin Tendulkar Photos ... - Cricket - Yahoo 10. Famous quotes on Sachin Tendulkar by other cricketers ... 	<ol style="list-style-type: none"> 1. Sachin Tendulkar - Wikipedia, the free encyclopedia 2. Kapil Dev wants Sachin Tendulkar to retire immediately 3. Sachin Tendulkar As Cricketer - Video Results 4. Sachin Tendulkar is a cricket god: Hussey - The Times of India 5. Sachin Tendulkar God of cricket Features Latest News 6. Why Kapil Dev's comments on Sachin Tendulkar are uncharitable 7. Achievements of Sachin Tendulkar - Wikipedia 8. ESPNCricinfo Awards 2011: Sachin Tendulkar 9. Drop Sehwag, Sachin from India ODI team: Cricketnet 10. Sachin Tendulkar Sachin Tendulkar Records 	<ol style="list-style-type: none"> 1. Sachin Tendulkar - Wikipedia, the free encyclopedia (15) 2. Sachin Ramesh Tendulkar is an Indian Cricketer (15) 3. Batsman Sachin (13) 4. Sachin Ramesh Tendulkar Born (11) 5. History of Cricket (10) 6. World Cup (7) 7. Year (7) 8. Sachin Tendulkar Biography (6) 9. Sachin Tendulkar Profile (6) 10. Latest News (4) 11. Sachin Tendulkar in 2011 Cricket World Cup (3) 12. Sydney (3) 13. Australian Batsman Michael Hussey Feels (2) 14. BBC News (2) 15. Cricinfo (2) 16. Other Topics (46)
Shahrukh Khan as Film Star	<ol style="list-style-type: none"> 1. Shahrukh Khan - Wikipedia, the free encyclopedia 2. Shahrukh Khan photos, videos, latest news, Shahrukh Khan pictures 3. Shahrukh Khan Actor Shahrukh Khan, Shahrudin film star Shah Rukh Khan in scuffle row - Yahoo! News 4. News for shahrukh kha n as film star 5. Indian filmstar Shahrukh Khan's detention at a US Airport 6. Ra. One challenges Shah Rukh Khan, the star - Times Of India 7. Shah Rukh Khan SRK the Don Ra. One Don 2 Shah Rukh Khan 8. Shahrukh Khan. Watch Shahrukh Khan Videos Movies Songs Online. 9. Indian film star Shah Rukh Khan in scuffle row - Yahoo! News 10. Bollywood Gossip Indian Movie Star Gossip Movie Masala 	<ol style="list-style-type: none"> 1. Indian film star Shah Rukh Khan in scuffle row - Yahoo! News 2. First Celebrity: Shahrukh Khan Indian film star 3. Bonding with co-stars Shahrukh Khan - Shahrudin film star Shah Rukh Khan in scuffle row - Yahoo! News 4. BBC News - Top Bollywood star Shah Rukh Khan in party bust-up 5. International Film Star Shahrukh Khan Bullied by Angry 6. Bollywood Star Shah Rukh Khan Bragges Extremist Party 7. Indian film star Shah Rukh Khan in scuffle row 8. Shah Rukh Khan - IMDb 9. Shahrukh Khan to star with Katrina Kaif in Yash Chopra's movie 10. Latest Bollywood Gossip 	<ol style="list-style-type: none"> 1. Shah Rukh Khan (13) 2. King Khan (9) 3. Born 2 November 1965 (8) 4. International (6) 5. Shahrukh Khan Picture (6) 6. BBC News (5) 7. Embroided in an Ugly Row on Monday (5) 8. Don 2 SRK (3) 9. Videos (3) 10. Dubaifor New Year by Nazia Khan (2) 11. Kareena Kapoor (2) 12. Lady Gaga (2) 13. Leonardo DiCaprio (2) 14. Other Topics (37)

Table1. Results obtained from different search engines.

Based on the above results obtained we can find the overall accuracy of the various search engines which is shown by following graph.

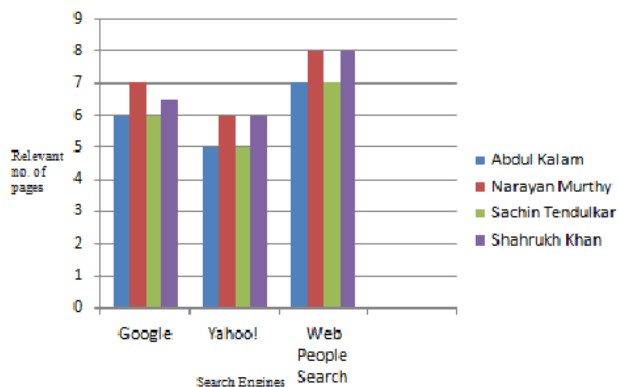


Figure 3: Result Graph

5. CONCLUSION

Our method exploits a web mining tool, knowledge driven cluster based search system that helps the user to find web information based on individual preferences. Our system also shows that, while it is possible to improve the efficiency of search through ontology method discussed above, it infact works best when operated in conjunction with one another and provide better search result.

References

Proceeding Papers:

- [1] 'Lingo: Search Results Clustering Algorithm Based on Singular Value Decomposition' By Stanis law Osinski, Jerzy Stefanowski, and Dawid Weiss, 2004.
- [2] Topic Sensitive PageRank by Haveliwala.

Journal:

- [3] International Applications Journal of Computer (0975 – 8887) Volume 22– No.2, May 2011.
- [4] 'Web People Search via Connection Analysis' by Dmitri V. Kalashnicov, Zhaoqi (Stella) Chen, Sharad Mehrotra, Member, IEEE, and Rabia Nuray-Turan. IEEE TRANCTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 20, NO. 11, NOVEMBER 2008.

Proceeding Papers:

- [5] Yang, X., Guo, D., Cao, X. and Zhou, J. (2008) Research on Ontology-Based Text Clustering, Proceedings of the 2008 Third International Workshop on Semantic Media Adaptation and Personalization, IEEE Computer Society Washington, DC, USA.
- [6] Michael W. Berry, Susan T. Dumais, and Gavin W. O'Brien. Using linear algebra for intelligent information retrieval. Technical Report UT-CS-94-270, 1994.

- [7] Stanis law Osinski and Dawid Weiss , "Conceptual Clustering Using Lingo Algorithm Evaluation on Open Directory Project Data", Institute of Computing Science, Pozna'n University of Technology,2003.

- [8] Stanis law Osinski, Jerzy Stefanowski, and Dawid Weiss,"Lingo: Search Results Clustering Algorithm Based on Singular Value Decomposition", Institute of Computing Science, Pozna'n University of Technology, 2003.

- [9] R. Guha and A. Garg, Disambiguating People in Search. Stanford Univ., 2004.

- [10] An algorithm for clustering of web search result. By Stanislaw Osinski,2003.

- [11] 'Ontology-Driven Induction Trees at Multiple Levels of Abstraction' by Jun Zhang, Adrian Selvescu and Vasant Honavar. Artificial Intelligence Research Laboratory Department of Computer Science, Iowa State University Ames, Iowa-50011-1040 USA.

- [12] An Empirical Evaluation on Semantic Search Performance of Keyword- Based and Semantic Search Engines: Google, Yahoo, Msn and Hokia Duygu Tümer1, Mohammad Ahmed Shah2, Yltan Bitirim1 2009 Fourth International Conference on Internet Monitoring and Protection IEEE .

- [13] D. V. Kalashnikov, R. Nuray-Turan, and S. Mehrotra. Towards breaking the quality curse. A web-querying approach to Web People Search. In Proc. of Annual International ACM SIGIR Conference, Singapore, July 20–24 2008.