

GCO Based Web Page Retrieval Through Semantic Techniques

A.Baskar, C.Gomathi

Abstract-- Now-a-days the online user face the several problems like unwanted data over the unwanted web pages, security issues, retrieval time consumption, unwanted advertisement, etc., To overcome the web page retrieval problem, the new idea has been proposed to retrieve the web pages through calculated value of relevance relationship between user queries and keywords by using web page weight calculation and find the distance value between web pages using Semantic concept like Geometric Interpretation(G), Cross web page (C), Optimality Properties(O). Finally, the user satisfactory value has been calculated to rank the web pages through relationship value.

Index Terms- Query Mining, Annotation based Retrieval, Relevant Web Content, Semantic concept.

I. INTRODUCTION

The user can retrieve the information from the web browser by issuing the query. The browser retrieve the information from different web pages. Some web pages deliver the user content based on how much time the user have been seen the particular content. The web browser count the value, how much time the user have seen the particular content. The most visited web content opened as first link in the web browser. By applying this method, the new uploaded useful content cannot open within first few links. The users are searching more time to get the relevant document. Annotation based retrieval systems satisfies the user query with partial manner. Based on text, the web page has been retrieved based on annotation based retrieval. Through this, the user should search more time to get the relevant document. To overcome this problem, we calculate the relationship value for retrieved web pages and arrange it in decreasing order the web pages to the user. The web page weight has been calculated based on the user query. We calculate the distance between the user query using the distance formula. Then we calculate the similarity and dissimilarity keyword weight value in retrieved web pages. Finally we add these the three

values that web page weight value, web page distance value, similar and dissimilar keywords weight value to find the relationship value. The web page, which one have more relationship value that will open as first web page to the user. We apply the user query to the stop word removal to remove the unnecessary words.

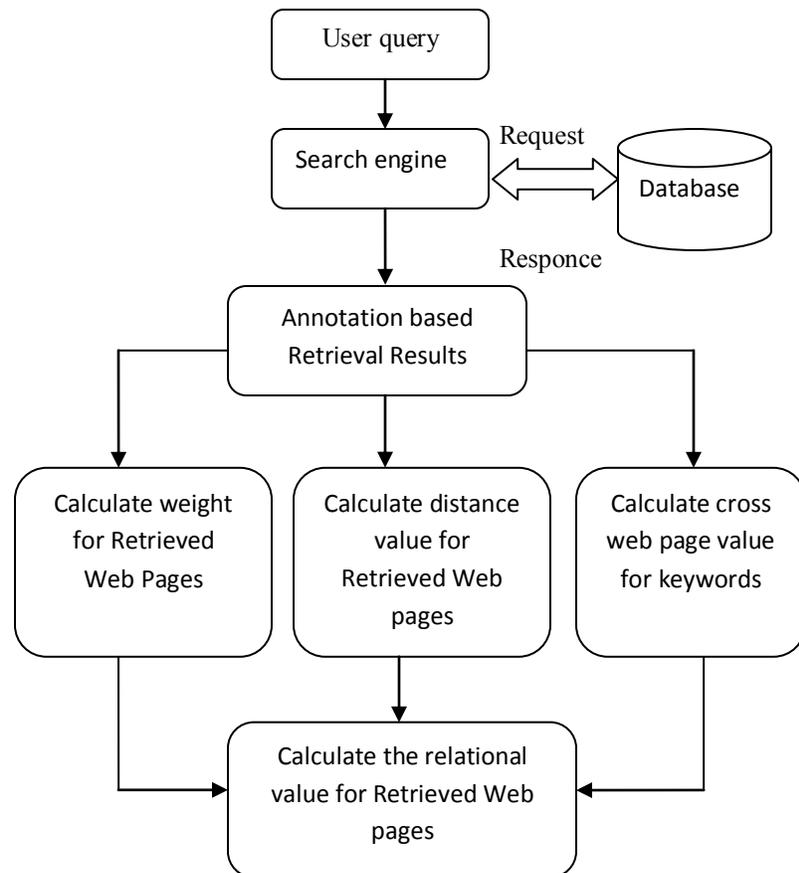


Fig 1. SYSTEM ARCHITECTURE MODEL

Fig 1 notifies that the user enter the query in search engine. The search engine gather the related results from the web data base. The results must be based on the user query. Calculate the web page weight value to retrieved web pages. Calculate the distance value to retrieved web pages. Finally, calculate the cross web page

value to retrieved web pages. The relationship value will be calculate by adding these three values. Arrange the relationship value in decreasing order with all web pages and allocates the rank to all web pages. Then, send the web pages as results to the user with rank order. The relationship value will be calculated to web pages using web page weight, web page distance and cross web page value. The user query, the single user query may have multiple words itself. If the user query have 'and' , 'to', 'or', 'space' in the statements, that will be remove using the stop word removal algorithm. Remaining each keyword considered as the each component. At first, find the web page weight value for each component separately. If one user query have more than one keywords, calculate the web page weight value by entering the keywords and apply these values into probability methods. We calculate the distance value to retrieved web pages. Here, retrieved web pages means it notifies the retrieved documents. We calculate the web page distance value in retrieved documents. We should analyze the maximum query length in search engine. The cross web page value means, some words have same meaning but different names. Based on this method, we calculate the cross web page value to retrieved web page documents. Finally, add these three values and finding the relationship value for retrieved documents. Allocate the rank to all web pages by arranging the relationship value in decreasing order then send the web page documents to the user.

The rest of the paper arranged as follows: In Section 2, we present the related work. Section 3 contains proposed methods based on our approach. Section 4 contains results and discussion, Section 5 contains conclusion about this paper, Section 7 contains references.

II. RELATED WORK

The web page mining tool gather the most used web paths and retrieve the web pages based on the document through web usage mining tool. Some web pages form based web pages. These web pages does not constructed in properly. Even though the mining tool search the content in the form based web pages which is based on user content[1] Bettina Bernendt. The web browser analyze the user query and it store web page access system's ip address, access time and date. When the user enter the query in the system, it will automatically annotated under the search. Based on query, the browser stored some web pages automatically. They apply the cloning method in double order transitions probability to retrieve the web pages[2] Jose Borges.

Page ranking system satisfies the user search partially and improve the user efficiency. Multiple documents are stored in the web data base. The entire database has been monitored using the indexing methods. The document count has increased sequentially.

The page ranking algorithm specifies that the most visited document considered as the high level document to retrieve to the user. In decreasing order, the document has been opened through the web pages by page ranking system[3] Lawrence Page. The first order, second order up to kth order markov model predict the web user behavior which is depends upon the user search. Smoothing techniques has been proposed for the user search that will be used to link from first web page to fourth web page. The second order markov model reduce the time complexity. The third order markov model measure that how much time it has been applied to normal prediction in the user search[4] Mukund Deshpande.

The conceptual log integrate the web documents with indexing mechanism. It displays the invisible documents in the web personalization. Semantic annotation combines with conceptual log that characterize the content using the vocabulary. Association rule retrieve the relevant content from the largest data set. Web pages has been designed by the user approach with taxonomy. The web page user profiles has been classified based on the real time environment. The individual web personalization systems improved through the web content profile. The original web usage logs leads the conceptual log and developed itself. Conceptual log delivers the patterns, sequential rules with the uniform resource identifier[5] M.Eirinaki. The online users retrieve the document from the various web pages as global model and collaborate the entire content to develop the mixture model. The global model develop the mixture model through the behavior model. Data cleaning method clear the unwanted data from web page while user browsing. The behavior model decides the way to design the future web pages and it represent the personalized probabilistic sequential models[6] Eren Manavoglu.

The web user efficiency improved through the Eigen vector value with web mining tools. The particular web domain concert problem identified in structural analysis. Hyperlink can be helped to develop the web pages with millions of document. The web mining tool make more response time to retrieve the relevant content. The different types queries are used to classify the query. Filter operation remove the unwanted consecutive web pages over the retrieval documents. Navigational methods displays the more information to the user with hyperlinked environment. General browser considered as the authoritative web pages to retrieve the more relevant document[7] Jon M.Kleinburg. The relevant web document retrieved through the unsupervised learning methods. Stochastic process allocates the probability value for all documents and assign the probability value to each words in the document. Markov chain property considered to calculate the probability value. The authorship information has included in probabilistic topic model using the latent dirichlet allocation techniques. The author identification methods used to analyze the author stylistic features. Graph and network model analyze the relationship

between the authors. The scientific model reduce the author topic model by discover the topic with authors[8] Mark Steyvers.

The web server log has given the permission to access data to the user. The first order markov model validate the user query and retrieve the web pages to the user using Bayesian approach and dirichlet prior. Bayesian approach and dirichlet prior used to calculate the cross validation between the web pages. The web page login entry avoids the problems that multiple user using the web page in the same computer. The user query classified in to java files, post script documents, HTML document. The user query classification comes through the grouping visit[9] Rituparna Sen. The web pages are improved automatically using the adaptive web page techniques. The improved process based on the user query. Page gathering algorithm used to implement the web pages. The web navigation process integrate the multiple information using the hyperlink. The web page improving process eliminates the unwanted retrieval web pages using the hyperlink environment[10] Mike Perkowitz.

III. PROPOSED METHODS

We calculate the relationship value to retrieved web documents by using three methods such that Geometric Interpretation(G), Optimality properties of proposed distance(o), Cross web page value(C). These three methods are applied to find the web page relationship value using the probability and dividing formulas.

A. Annotation based Retrieval

Annotation based retrieval specifies that automatically web pages retrieved based on the user typed keywords. The search engine analyze the user keyword and it compare with the stored documents in the web data base. Here, the web pages specifies that web documents. The web documents has been retrieved dynamically from the web data base, based on the user typed query. The search engine check the query with the data base using keywords. The keywords has been generated in the user typed query. Annotation based retrieval does not contain any algorithm or formula to retrieve the web documents itself. By having the keywords only, the search engine search the document in the entire data base and retrieve the results to the user.

B. Geometric Interpretation (G)

Geometric interpretation specifies the web page weight value to retrieve the web documents which is based on the user query. A single user query may have more than one keyword in itself. The keyword considered as the component. A single user query may have multiple components. We calculate the web page weight

value separately to each component. We apply the formula to calculate the web pages

$$W_p = \frac{W_r}{W_t} \quad (1)$$

In equation (1), W_p specifies web page weight value, W_r specifies relevant web pages based on the user query. W_t specifies total number of web pages on the particular search engine. The above equation may suited if the user query have single component. If the user query have two or more components, the above equation does not suited directly.

$$P(A, B) = \frac{P(A \cap B)}{P(A \cup B)} \quad (2)$$

Equation (2) describes, if the user query have two components itself, we may apply the equation (2) to find web page weight value between the components. In Equation (2), A and B are two different components. $P(A \cap B)$ specifies that intermediary relationship between the two components. $P(A \cup B)$ specifies that entire relationship between the two components. To find the $P(A \cup B)$ value, we can use the formula

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \quad (3)$$

we apply the equation (3) in equation (2), then we may write the equation

$$P(A, B) = \frac{P(A \cap B)}{P(A) + P(B) - P(A \cap B)} \quad (4)$$

Equation (4) used to find the web page weight value between the two components. We may use the same formula to find web page weight value between multiple components.

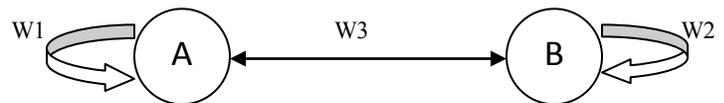


Fig 2. SIMPLE DIAGRAM FOR CALCULATING GEOMETRIC INTERPRETATION

In Fig 2, we have two components A & B. Two components specifies two different keywords. The two different key words are occurred in the single statement as user typed query. We should calculate the web page weight value for these two different queries. At first, we find the web page weight value W_1 for A

component, then we find the web page weight value W_2 for B component in separately. Next, we find the web page weight value W_3 for between the two components. Finally we find the entire web page weight value for both components by applying the equation (4).

C. Optimal Properties of Proposed Distance (O)

Optimal properties derived from the Markov chain properties. These properties are used to calculate the web page distance value. Optimality properties measure distance value for retrieved web documents. Here, web pages notifies the web documents.

1. **Removal of Stop Word:** Removal of Stop Word removes the unwanted characters from the user query. Unwanted characters like `space`, `?`, `;`, `!`, `-', `_` These characters may be occurred in the user query. These characters can be removed by the stop word removal techniques. We calculate the distance value after eliminating the unwanted characters.
2. **Distance Calculation:** We measure the distance value for retrieved web documents. We have considered the user query length and the retrieved web document length. We have used the formula to calculate the distance value

$$Wd = \frac{D(A) + D(B)}{Td} \quad (5)$$

In Equation (5), Wd notifies retrieved web document's distance total value. $D(A)$ notifies user query total length, $D(B)$ notifies retrieved web content total length. Td notifies that maximum query length in the search engine. $D(B)$ value will be changed based on the retrieved web document.

D. Cross web page calculation (C)

The user enter the query in the search engine. The search engine have given the results to the user based on the user query. Lot of keywords are available in the user query. The search engine search the web contents using the keywords. Basically keywords are classified into different category. Some of keywords have same meaning but in name, the key words have different names. The search engine should consider this problem. Because, the search engine search the content based on keywords only. We have find some solution to overcome this problem. We have find some words which is having two names in same meaning. We have used the formula

$$Wc = \left(\frac{Ks}{T} + 1\right) - (1 - \mu)\left(1 - \frac{Kd}{T}\right) \quad (6)$$

In equation (6), Wc - specifies the total cross web page weight value and Ks - specifies the total number of similar keywords weight, Kd - specifies the total number of dissimilar keywords weight, here dissimilar keywords means another alternate word for same keyword, T - specifies total number of keywords in the search engine.

E. Calculating the Relationship value

By having the above three values, we may find the relationship value for each an every retrieved content by using the formula

$$Vr = (\alpha * Wp) + (\beta * Wd) + (\gamma * Wc) \quad (7)$$

In Equation (7), Vr notify the relationship value for retrieved web content.

IV. EXPERIMENTAL RESULTS

In this section, we have proved the results in Geometric Interpretation and the optimality properties of proposed distance. For example, if the user enter the query "How to cook biriyani" in search engine, we can retrieve some results through the annotation based retrieval. At first, we apply the stop word removal techniques to remove unwanted character. In the query, `to` is the unwanted character, we can remove this word.

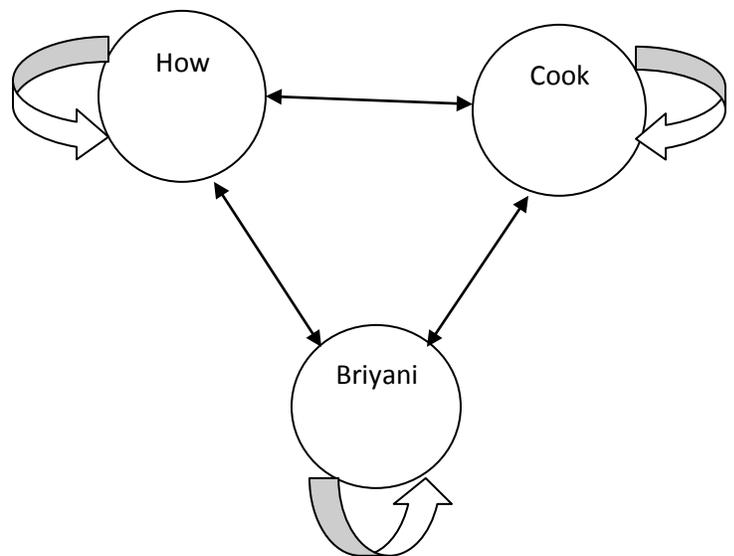


Fig 3. SAMPLE QUERY DIAGRAM FOR CALCULATING WEB PAGE WEIGHT

In Fig 3, based on the user query, we have three different keywords such that how, cook and briyani. These three keywords consider as the three different components. First, we calculate the web page weight value separately for each component. Then, we combined the keywords like How-briyani, cook-briyani, how-cook to find the web page weight value. Finally, we calculate the web page weight value for the full user query. Now we have some web page weight value. These values will apply in to the Equation(4). We find the web page weight value for user query. We find the web page weight value for retrieved document. The documents retrieved based on the user query. We have given some results for retrieved documents.

Table I. Web page weight value for retrieved content

Retrieved Content	Wp Value
How to cook briyani	0.4840
How to make a chicken briyani	0.3020
How to cook vegetarian briyani	0.1660
How to prepare briyani	0.0579
How to make muslim briyani	0.0479
How to cook briyani in tamil	0.0221
How to cook mutton briyani	0.3080
How to cook briyani rice	0.2779
Chicken briyani recipe	0.3240
How to cook veg briyani	0.1430
How to cook egg briyani	0.1250
How to cook briyani in hindi	0.0250
Easy chicken briyani	0.2950
How to cook basmathi rice for briyani	0.0182
Secrets of making a perfect briyani	0.0172
Vegetable briyani recipe	0.3170
Vegetable briyani video by Tarla Dalal	0.0247
How to make chicken briyani at home	0.2760
Hyderabad chicken dum briyani	0.0239
Hara masala briyani	0.0214
10 minutes easy briyani recipe	0.2210
Cooking food and stories: Kolkatta briyani	0.2540
Recipe for the perfect Eid briyani	0.0121
Easy aromatic chicken briyani	0.0059
Cooking briyani in oven	0.0481
Kerala briyani recipe and cooking	0.0339
Easy chicken dum briyani recipe with stepwise pictures	0.0011
Kitchen secrets and snippets: vegetable briyani	0.0064
Small potato briyani meal	0.0928
How to make shahi korma briyani without oil	0.0147

The Table I specifies the values that web page weight for relevant web document. The user query “ How to cook briyani”. We have given 30 relevant document based on the user query. We calculate the web page weight for 30 documents. We have assume the maximum documents in the search engine 100000 laksh documents.

Table II. Web page distance value for retrieved content

Retrieved Content	Wd Value
How to cook briyani	0.0078
How to make a chicken briyani	0.0117
How to cook vegetarian briyani	0.0127
How to prepare briyani	0.0092
How to make muslim briyani	0.0107
How to cook briyani in tamil	0.0112
How to cook mutton briyani	0.0107
How to cook briyani rice	0.0097
Chicken briyani recipe	0.0097
How to cook veg briyani	0.0092
How to cook egg briyani	0.0092
How to cook briyani in hindi	0.0112
Easy chicken briyani	0.0087
How to cook basmathi rice for briyani	0.0151
Secrets of making a perfect briyani	0.0146
Vegetable briyani recipe	0.0107
Vegetable briyani video by Tarla Dalal	0.0161
How to make chicken briyani at home	0.0142
Hyderabad chicken dum briyani	0.0127
Hara masala briyani	0.0083
10 minutes easy briyani recipe	0.0127
Cooking food and stories: Kolkatta briyani	0.0180
Recipe for the perfect Eid briyani	0.0142
Easy aromatic chicken briyani	0.0127
Cooking briyani in oven	0.0097
Kerala briyani recipe and cooking	0.0142
Easy chicken dum briyani recipe with stepwise pictures	0.0229
Kitchen secrets and snippets: vegetable briyani	0.0205
Small potato briyani meal	0.0107
How to make shahi korma briyani without oil	0.0176

Table II describes that, the results for web document distance value. We have given 30 document results for the query “How to cook briyani”. We remove some unwanted characters and space in the document name to calculate the web page distance. Equation (5) are applied to measure the web page distance.

We find the similar and dissimilar keywords weight value by applying the equation (7). First, we should analyze that what are the keywords have same meaning and different name. we have identified some keywords have different name in same meaning. We have given that name in Table III.

Table III. Different words in same meaning

Keywords	Another Word for Keyword
buy	Purchase
big	Large
quickly	Speedily
middle	Between

look	See
mountain	Hills
mobile	Cell phone
happy	Pleasure
Rich	Spicy
At present	Current time
Glass	Mirror
Lady	Women
Silent	Pease
pause	break

Next we have find the keywords weight value to calculate the cross web page value. Table III have more keywords itself which specifies same meaning but in different keywords.

Table IV. Different Keywords Weight Value

Keywords	(Ks/T) Value	Another Word for Keyword	(Kd/T) Value
Buy	0.1250	Purchase	0.4220
Big	0.1350	Large	0.0127
Quickly	0.2230	Speedily	0.0020
Middle	0.4260	Between	0.0085
Look	0.1080	See	0.0247
Mountain	0.2860	Hills	0.1880
Mobile	0.2080	Cell phone	0.0002
Happy	0.7140	Pleasure	0.1160
Rich	0.2260	Spicy	0.0004
At present	0.3370	Current time	0.0318
Glass	0.3310	Mirror	0.0014
Lady	0.3060	Women	0.0072
Silent	0.0886	Pease	0.0018
Pause	0.0771	Break	0.0032

Table IV shows that different keywords and each keyword have one alternate and those keyword weight value. In Table IV, fourth column notify that dissimilar keywords weight value. We may apply these two values in in Equation (6), we may get the cross web page value.

V. CONCLUSION

This paper hybrid the features of three different methods like Geometric interpretation, Cross web page calculation, Optimality properties of proposed distance. This GCO used for web page retrieval through the relational value of web pages. Experimental results provide the detail of web page weight, cross web page calculation for keywords and web page distance value for optimal retrieval of web page. In future this work has extended to image and on line video retrieval.

VI. ACKNOWLEDGEMENT

The authors like to given the thanks to reviewers for their comments that will help to improve this paper significantly.

VII. REFERENCES

- [1] Bettina Berendt, Myra Spiliopoulou, "ANALYSIS OF NAVIGATION BEHAVIOUR IN WEB PAGES INTEGRATING MULTIPLE INFORMATION SYSTEMS" The VLDB Journal-2009.
- [2] Jose Borges, Mark Levene, " A DYNAMIC CLUSTERING BASED MARKOV MODEL FOR WEB USAGE MINING", May 26, 2004.
- [3] Sergey Brin, Lawrence Page, " THE ANATOMY OF A LARGE SCALE HYPERTEXTUAL WEB SEARCH ENGINE", computer networks and ISDN systems, 1998.
- [4] Mukund Despande, George Karypis, " SELECTIVE MARKOV MODELS FOR PREDICTING WEB PAGE ACCESS", ACM Transaction and internet technology, Vol 4 , May 2004.
- [5] M. Erinaki, M. Vazirgiannis, I. Varlamis, " USING SITE SEMANTICS AND A TAXONOMY TO ENHANCE THE WEB PERSONALIZATION PROCESS".
- [6] Eren Manavoglu, Dmitry Pavlov, C. Lee Giles, " PROBABILISTIC USER BEHAVIOUR MODEL".
- [7] Jon M Kleinberg, " AUTHORITATIVE SOURCES INA HYPERLINKED ENVIRONMENT", ACM – SIAM symposium, 1998.
- [8] Mark Steyvers, Padhraic Smyth, Thomas Griffiths, "PROBABILISTIC AUTHOR – TOPIC MODELS FOR INFORMATION DISCOVERY" 10th ACM SigKDD and Data Mining Scattle, 2004.
- [9] Rituparna Sen and Mark H. Hamsen, "PREDICTING WEB USERS NEXT ACCESS BASED ON LOG DATA" Journal of Computational and Graphical Statistical, Vol-12, 2003.
- [10] Mike Perkowitz, Oren Etzioni, "TOWARDS ADAPTIVE WEB SITES : CONCEPTUAL FRAMEWORK AND CASE STUDY".