`RESEARCH ARTICLE`                                                                                    `OPEN ACCESS`

# Distinct Web Search Engine with Reducing Ambiguity Word Complexity

## Ragavi.R[1], Sowmiya.S.R[2]

*[1] Department of Computer Science and Engineering,Dhanalakshmi Srinivasan Engineering College,Perambalur,India*

Email−[1]ragavi616@gmail.com

*[2] Department of Computer Science and Engineering,Dhanalakshmi Srinivasan Engineering College,Perambalur,India*

Email−[2]jitgiri10@gmail.com

--------------------------------------\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*-------------------------------

## Abstract:

Searching is one of the common mission accomplished at the Internet. Search engines are the primary tool of the net, from in which you'll be able to acquire associated information and searched in step with the required key-word given by way of the person. The statistics at the internet is developing dramatically. The consumer has to spend more time inside the internet a good way to locate the precise data they may be inquisitive about. Existing internet system do not remember particular needs of user and serve every person similarly. For this ambiguous query, a number of documents on awesome topics are returned by way of engines like Google. Hence it turns into difficult for the consumer to get the specified content material. Moreover it also takes extra time in searching a pertinent content. Privacy based totally Personalized Web Search Engine is considered as a promising answer to address these problems, considering the fact that distinct search outcomes can be furnished depending upon the choice and records needs of users. It exploits consumer records and seek context to learning in which experience a question refer. In order to perform Personalized Web seek it is critical to version User's hobby. User profiles are built to version person's need based on his/her web utilization information. This Enhanced User Profile will help the person to retrieve concentrated information. This paper proposes structure for constructing consumer profile and enhances the person profile the usage of historical past know-how. It can be used for suggesting suitable net pages to the person based totally on his search question and background expertise. And also put into effect the pruning set of rules to put off the consumer info from anonymous person for maintain the important thing phrase privacy. On other side we want to hide the privacy contents present in the person profile to place the privateness hazard under manage. User privateness may be provided in form of safety like without compromising the customized seek satisfactory.

Index Terms— ***Search engine, Personalized web search, User profile, Greedy approach, Privacy***

--------------------------------------\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*-------------------------------

## I.   INTRODUCTION

Web mining is the software of statistics mining strategies to discover styles from the World Wide Web. As the call proposes, this is information collected via mining the web. It makes utilization of automatic apparatuses to show and extricate data from servers and web2 reviews, and it lets in businesses to get to both organized and unstructured

statistics from browser activities, server logs, internet site and hyperlink structure, web page content material and exceptional sources.

The aim of Web structure mining is to generate structural summary about the Web website and Web web page. Technically, Web content material mining mainly specializes in the shape of internal-record, whilst Web structure mining attempts to discover the link shape of the links on the inter-record stage. Based on the topology of the hyperlinks, Web shape mining will categorize the Web pages and generate the information, which includes the similarity and courting among specific Web sites.

Web structure mining can also have some other route -- coming across the structure of Web report itself. This type of structure mining may be used to reveal the structure (schema) of Web pages, this would be proper for navigation reason and make it possible to examine/combine Web page schemes. This type of structure mining will facilitate introducing database strategies for gaining access to data in Web pages by using providing a reference schema. Web mining may be divided into three different sorts – Web utilization mining, Web content material mining and Web structure mining.

Current net search engines like Google and yahoo are built to serve all users, impartial of the special desires of any man or woman user. With the exponential evolution of the available records on the World Wide Web, a traditional hunt engine, even though based totally on sophisticated report indexing algorithms, has trouble meeting efficiency and effectiveness overall performance demanded via customers trying to find applicable facts. Personalization of net search is to carry out retrieval for every user incorporating his/her pursuits. Personalized net search differs from time-honored web seek, which returns same consequences to all customers for same queries, regardless of various person hobbies and facts needs. When queries are issued to search engine, maximum go back the

equal consequences to users. In truth, the giant majority of queries to search engines are short and ambiguous. Different users may additionally have absolutely one-of-a-kind records desires and desires whilst using exactly the equal question. Personalized net search can be completed by means of checking content similarity among net pages and user profiles. Some paintings has represented person pursuits with topical categories. User's topical pastimes are both explicitly specified by means of customers themselves, or can be robotically learned by using classifying implicit person facts. Search outcomes are filtered or re-ranked by means of checking the similarity of subjects between search outcomes and consumer profile. The web mining approach is shown in fig 1
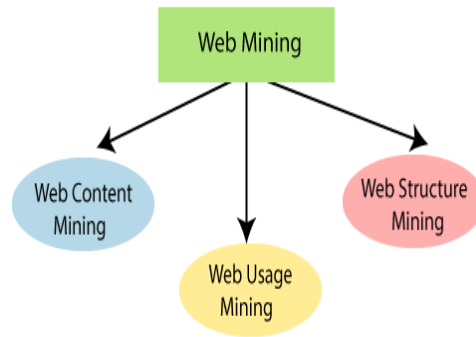


**Fig 1: Types of Web-mining**

## II. RELATEDWORK

P. Yin, et.al,…[1] applied the look at of consumer perceptions approximately websites discloses that the most important layout functions for distinct website domains consist of navigations, timeliness, clarity, visualization, accuracy, and safety. The clean-to-navigate feature is ranked a few of the top three for all domains. Web customers look forward to extra at ease surfing reviews which require the WWW surroundings to be each powerful and efficient. Effective browsing method that the users can easily search the most interesting internet site by using specifying applicable key phrases, while efficient browsing indicates the users can reach the target web site in a website with just

few clicks. Both necessities may be facilitated by the usage of the web mining techniques in the design segment. In this observe we advise a new system for the internet site structure optimization (WSO) trouble based on a complete survey of present works and practice considerations. An more suitable tabu search (ETS) algorithm is proposed with superior seek capabilities of more than one neighborhoods, adaptive tabu lists, dynamic tabu tenure, and multi-degree aspiration standards.

M. Chen, et.al,…[2] addressed the query of a way to improve the shape of a internet site in place of reorganize it appreciably. Specifically, we develop a mathematical programming (MP) version that helps consumer navigation on a internet site with minimum modifications to its current structure. Our version is particularly suitable for informational web sites whose contents are static and comparatively stable through the years. Examples of companies that have informational websites are universities, vacationer attractions, hospitals, federal businesses, and sports agencies. Our version, but, won't be appropriate for websites that basically use dynamic pages or have risky contents. While various strategies were proposed to relink webpages to improve navigability the usage of consumer navigation records, the absolutely reorganized new shape may be enormously unpredictable, and the price of disorienting users after the changes stays unanalyzed. This paper addresses the way to enhance a website with out introducing great modifications. Specifically, we endorse a mathematical programming version to enhance the user navigation on a website while minimizing alterations to its current structure. Results from massive checks conducted on a publicly available actual facts set suggest that our model now not best significantly improves the person navigation with very few changes, however also can be efficiently solved. We have additionally tested the model on big synthetic data units to demonstrate that it scales up thoroughly. In addition, we define two assessment metrics and use them to evaluate the performance of the progressed internet site the use of the real facts set. Evaluation

outcomes confirm that the person navigation on the progressed structure is indeed substantially improved. More apparently, we find that closely disoriented users are much more likely to enjoy the progressed shape than the less disoriented users.

C. Kim, et.al,…[3] carried out the framework for performed excessive productiveness of publishing, the webpages in lots of web sites are automatically populated by way of the usage of not unusual templates with contents. For human beings, the templates offer readers clean get admission to to the contents guided through constant systems despite the fact that the templates aren't explicitly introduced. However, for machines, the unknown templates are considered dangerous due to the fact they degrade the accuracy and overall performance because of the irrelevant phrases in templates. Thus, template detection and extraction techniques have obtained loads of interest these days to improve the overall performance of internet programs, such as information integration, search engines like google and yahoo, category of internet documents, and so on. Thus, template detection techniques have received a variety of attention recently to improve the overall performance of search engines like google, clustering, and category of net documents. In this paper, we present novel algorithms for extracting templates from a big variety of web files which are generated from heterogeneous templates. We cluster the net documents based totally on the similarity of underlying template structures inside the documents so that the template for each cluster is extracted concurrently. We increase a singular goodness measure with its fast approximation for clustering and offer complete evaluation of our algorithm. Our experimental outcomes with actual-life information units affirm the effectiveness and robustness of our algorithm in comparison to the state of the art for template detection algorithms.

Y. Yang, et.al,…[4] applied the facts extraction (IE) scheme which performs an crucial function in internet information discovery and management. The maximum essential obligations in facts extraction from the Web are web site structure expertise and herbal language sentences processing.

However, little work has been done closer to an included statistical version for knowledge webpage structures and processing herbal language sentences within the HTML elements. Our latest work on web site expertise introduces a joint model of Hierarchical Conditional Random Fields (i.e. HCRF) and extended SemiMarkov Conditional Random Fields (i.e. Semi-CRF) to leverage the page structure information outcomes in loose text segmentation and labeling. In this pinnacle-down integration version, the selection of the HCRF version ought to manual the decision-making of the Semi-CRF model. However, the drawback of the top-down integration strategy is likewise obvious, i.E., the choice of the Semi-CRF version couldn't be used by the HCRF model to manual its selection-making. This paper proposed a novel framework known as WebNLP, which allows bidirectional integration of page shape expertise and textual content information in an iterative manner. We have applied the proposed framework to neighborhood commercial enterprise entity extraction and Chinese character and agency call extraction. Experiments display that the WebNLP framework carried out significantly higher performance than existing methods.

J. Hou, et.al,…[5] advocate two algorithms that use web page similarity to discover applicable pages. The new page source, based totally on which the algorithms are hooked up, is constructed with required properties. The web page similarity evaluation and definition are based on link facts the various Web pages. The first set of rules, Extended Cocitation set of rules, is a cocitation algorithm that extends the conventional cocitation concepts. It is intuitive and concise. The second one, named Latent Linkage Information (LLI) set of rules, reveals applicable pages more efficiently and precisely by means of the usage of linear algebra theories, specially the singular price decomposition of matrix, to reveal deeper relationships most of the pages. This paper gives two link analysis-based totally algorithms to locate relevant pages for a given Web page (URL). The first set of rules comes from the prolonged cocitation analysis of the Web pages. It is intuitive and easy to put into effect. The 2nd one takes benefit of linear algebra theories to expose deeper relationships some of the Web pages and to perceive relevant pages more exactly and efficaciously. The experimental outcomes display the feasibility and effectiveness of the algorithms. These algorithms might be used for various Web programs, such as enhancing Web seek. The ideas and strategies on these pictures might be useful to other Web-associated researches.

## III. EXISTING METHODOLOGY

The current machine put into effect personalized internet seek for enhancing common experience and folksonomy primarily based smart seek systems. A large department of the cutting-edge web is characterized by user generated content material classified the use of collaborative tagging or folksonomy. It makes very tricky to look for suitable content material due to ambiguity in lexical illustration of principles and variances in preferences of users. A past effort to use this approach has shown encouraging results in obtaining applicable content however it does no longer deal with the difficulty of noise in seek outcomes. In existing machine, the put in force the system that employ the customized web seek method of traditional internet seek systems to pay attention on the issue of inappropriate seek effects in commonplace feel the usage of K-Means and Page Rank algorithm.

### 3.1 K-Means algorithm:

The K means a set of rules is easy to implement, requiring an easy records shape to hold some facts in every iteration to be used inside the next new release. The idea makes K-means more efficient, especially for dataset containing big number of clusters. Since, in every generation, the k-means set of rules computes the distances between data point and all facilities; this is computationally very high-priced particularly for huge datasets. Therefore, we do can use from preceding new release of k-approach algorithm. K-Means is one of the top ten clustering algorithms that are extensively utilized in actual global

packages. It is a very simple unsupervised mastering algorithm that discovers actionable information by means of grouping comparable gadgets into numerous clusters. However, it wishes the range of clusters to be known priori. We can calculate the distance for each records factor to nearest cluster. At the next generation, we compute the gap to the preceding nearest cluster. The point stays in its cluster, if the new distance is less than or same to the previous distance, and it isn't required to compute its distances to the alternative cluster centers. This saves the time required to compute distances to k-1 cluster centers. "K-means algorithm is one in all first which a statistics analyst will use to investigate a new records set due to the fact it's far algorithmically easy, enormously robust and gives "true enough" solutions over a huge type of records sets." The K-means set of rules is the most normally used partitioned clustering algorithm because it could be effortlessly carried out and is the maximum efficient one in terms of the execution time.

The basic algorithm pseudo code as follows:

Input: X = {x1, x2, x3,…..,xn} be the set of data points , Y= {y1,y2,y3…yn} be the set of data points and V = {v1,v2,v3,….,vn} be the set of centers

Step 1: Select 'c' cluster centers arbitrarily

Step 2: Calculate the distance between each pixels and cluster centers using the Euclidean Distance metric as follows

$$Dist(X,Y) = \sqrt{\sum_{j=1}^{n}(X_{ij} - Y_{ij})^2} \qquad \text{------- Eqn(1)}$$

X, Y are the set of data points

Step 3: Pixel is assigned to the cluster center whose distance from the cluster center is minimum of all cluster centers

Step 4: New cluster center is calculated using

$$V_i = \frac{1}{c_i}\sum_{1}^{ci} x_i \qquad \text{------Eqn(2)}$$

Where Vi denotes the cluster center, ci denotes the number of pixels in the cluster

Step 5: The distance among every pixel and new obtained cluster facilities is recalculated

Step 6: If no pixels were reassigned then stop. Otherwise repeat steps from 3 to 5

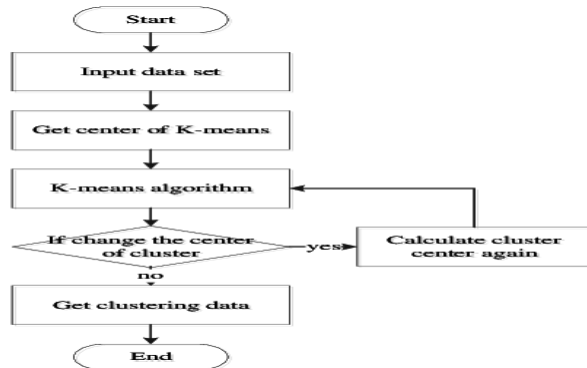The flowchart of the algorithm is shown in fig 3.1



**Fig 2 :Flow chart of K-Means clustering**
**Limitations:**

- The most important trouble with this set of rules is that it's far touchy to the choice of the preliminary partition and might converge to nearby optima.
- It is computationally highly-priced and requires time proportional to the made from the quantity of information gadgets, variety of clusters and the wide variety of iterations.
- The first-class of the ensuing clusters closely relies upon on the choice of initial centroids which reasons it to converge at nearby finest.
- Empty clusters hassle, which arise to defined fixed cluster in staring of the set of rules.

**3.2 PAGE RANK ALGORITHM**

PageRank (PR) is a set of rules utilized by Google Search to rank web sites of their search engine results. PageRank changed into named after Larry Page, one of the founders of Google. It isn't always the handiest algorithm used by Google to order search engine outcomes, however it is the primary algorithm that was utilized by the agency, and it's miles the pleasant-acknowledged. The above centrality measure is not carried out for the multi-graphs. The PageRank algorithm outputs a chance distribution used to symbolize the probability that a person randomly clicking on

hyperlinks will arrive at any specific page. PageRank can be calculated for collections of documents of any length. It is believed in several studies papers that the distribution is lightly divided amongst all files within the collection at the beginning of the computational technique. The PageRank computations require numerous passes, referred to as "iterations", thru the gathering to regulate approximate PageRank values to greater carefully reflect the theoretical true cost. The size of every query is proportional to the full size of the opposite faces which can be pointing to it.

The pseudo code for the algorithm is:

Given a web graph with n nodes, where the nodes are pages and edges are hyperlinks

- Assign each node an initial page rank
- Repeat until convergence calculate the page rank of each node (using the equation in the previous slide)

$$PR(A) == (1-d) + d * (PR(T1)/C(T1)+\ldots+ (PR(Tn)/C(Tn))$$

After all, the sum of the weighted page ranks of all pages Ti is multiplied with a damping factor d which can be set between 0 and 1. Thereby, the extend of page rank benefit for a page by another page linking to it is reduced

## IV.    PROPOSED METHODOLOGIES

Searching is one of the usually used venture on the Internet. Search engines are the simple tool of the net, from which related statistics may be gathered consistent with the desired query or key-word given by way of the consumer, and are extraordinarily famous for recurrently used websites. With the first rate development of the World Wide Web (WWW), the records search has grown to be a primary commercial enterprise phase of a worldwide, competitive and money-making market. An ideal search engine is the one which need to journey via all the web pages inside the WWW and have to listing the related statistics based totally at the given user key-word. In spite of the current developments on internet seek technologies, there are nevertheless many conditions wherein seek engine users obtains the non-relevant seek effects from the search engines. A customized Web seek has numerous tiers of efficiency for distinctive users, queries, and seek contexts. This technique can help in offering greater applicable information for a particular consumer with the aid of reorganizing the search results from Web search. Hence it permits users to acquire the proper records in accordance with their hobby primarily based on Greedy algorithm.

Based on the hardness of the hassle, we use a greedy algorithm. Implicit statistics includes past sports as recorded in Web server logs through cookies or consultation tracking modules. Explicit facts generally come from registration bureaucracy and score questionnaires. Additional facts including demographic and application records (for example, e-trade transactions) also can be used. In a few cases, Web content material, shape, and alertness data can be introduced as extra assets of statistics, to shed extra light on the subsequent tiers. Data is often pre-processed to put it right into a format that is well matched with the evaluation approach to be used inside the subsequent step. Preprocessing can also include cleansing statistics of inconsistencies, filtering out irrelevant facts consistent with the aim of evaluation (instance: automatically generated requests to embedded photos can be recorded in internet server logs, despite the fact that they upload little records approximately consumer pursuits), and completing the lacking hyperlinks (because of caching) in incomplete click on thru paths. Most importantly, specific periods want to be diagnosed from the exceptional requests, based totally on a heuristic, which include requests originating from an same IP cope with within a given time period. Analysis of Web statistics - Also called Web Usage Mining, this step applies device mastering or Data Mining strategies to discover thrilling usage patterns and statistical correlations among net pages and user organizations. This step frequently effects in automated user profiling, and is normally implemented offline, so that it does no longer add a burden on the net server. The ultimate section in personalization uses the effects of the preceding

analysis step to deliver recommendations to the person. The advice technique typically entails producing dynamic Web content on the fly, along with adding hyperlinks to the ultimate web page asked by the person. In the start, a person profile is randomly selected as the seed of a brand new cluster. The closest user profile is continuously selected and combined with the seed until the cluster satisfies p-linkability or the size of the cluster |Gi| satisfies the constraint |Gi| ≥ |U|avg p . At next step, a user profile with the longest distance to the previous seed is selected as the seed of the new cluster.

result ← ∅
 C ← ∅
seed ← a randomly picked user profile from S
 while |S| > 0 do
 seed ← the furthest user profile(with the min similarity value) to seed
while C does NOT satisfy p-linkability AND |S|>0 do add the closest user profile (with the max similarity value) to C
 end while
if C does satisfy p-linkability then
 result ← result ∪ C;
C ← ∅
end if
end while
for each user profile in C do  assign it to the closest cluster end for

The component to protect privacy is generating an online profile that is put into effect on a search proxy running on a client machine itself. This proxy will have the hierarchical user profile and customized privacy requirements. Phases in this Architecture consists both online and offline phase. Hierarchical generation of user profile on client side and customized privacy requirements specified by the user are handled. The proposed work can be described as follows:

**4.1 User enrollment**

In this module, can construct user information system and includes admin and data users. Admin can be responsibility for maintain all records with secured manner. Admin provides approval system. Data users search information from server. User can be register with their details such as name, age, gender and other details.

**4.2 User profile**

This module introduces an approach to personalize digital multimedia content based onuses profile information. For this two main mechanisms were developed a profile generator that automatically creates user profiles representing the user preferences. Profile generator can be used to predict the query whether it is common or personalized. Common query can consider as new search. And personalized search can be referred as already search query.

**4.3 Repository creation**

In this module used to create a repository for personalized search user can select a personalized web search they can accurate data from search. Once a user has entered a query, the query can be compared against the contextual information available to determine if the query can be refined to include other terms. If the query is on a topic the user has previously seen, the system can reinforce the query with similar terms, or suggest results from prior searches. Generically refers to a central place where data is stored and maintained. A repository can be a place where multiple databases or files are located for distribution over a network, or a repository can be a location that is directly accessible to the user without having to travel across a network. If it is a new topic, chances are the system should not augment the query, or if it does, it can help define what the topic is not about by providing a diverse set of results to the user. The final output of query augmentation is a more precise query that can be shown to the user and submitted to a search engine for processing.

**4.4 Search module**

In this module used to search data in this page. In this page provide search engine for user to search data. User customizable privacy preserving search, it's the framework assumes that the queries do not contain any sensitive information, and aims at protecting the privacy in individual user profiles

while retaining their usefulness for pws. In this section, the procedures carried out for each user during two different execution phases, namely the offline and online phases. Generally, the offline phase constructs the original user profile and then performs privacy requirement. Customization according to user specified topic sensitivity can be determined. The subsequent online phase finds the optimal risk generalization solution in the search space determined by the customized user profile

## 4.5 Results

The study of the efficiency of the proposed generalization algorithms is quite realistic as it clearly seen from the output. Here we implement mp model on the profiles, which has an edge over other search engines. The queries are randomly selected from their respective query log. The profile-based personalization contributes little or even reduces the search quality, while exposing the profile to a server would for sure risk the user's privacy. To address this problem, we develop an online mechanism to decide whether to personalize a query. The basic idea is straightforward. If a distinct query is identified during generalization, the entire runtime profiling will be aborted and the query will be sent to the server without a user profile. In this page provide search result to user. User views their search result in this page. The proposed architecture is shown in fig 3.
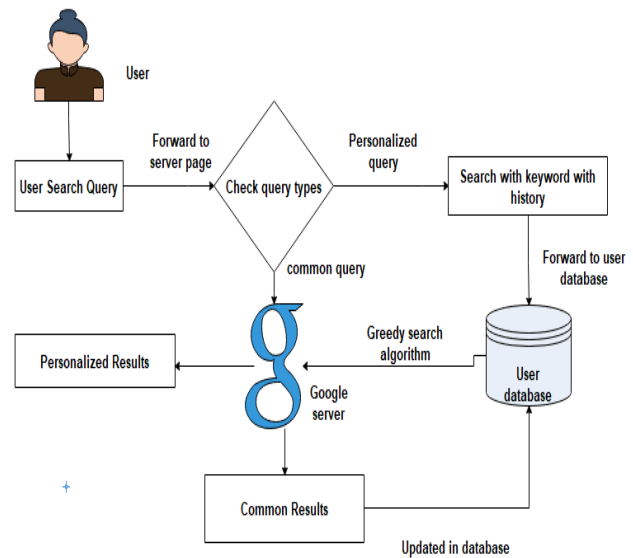


**Fig 3: Proposed framework**

## V. RESULTS AND DISCUSSION

We can evaluate the performance using accuracy metrics. The accuracy metric is evaluated as

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} *100$$

The proposed algorithm provide improved accuracy rate than the machine learning algorithms.

| Algorithm | True positive | True negative | False positive | False negative |
|---|---|---|---|---|
| K-means | 5 | 10 | 20 | 30 |
| Page rank | 10 | 8 | 15 | 20 |
| Greedy algorithm | 20 | 5 | 10 | 10 |

Table (1) Performance measurement

Accuracy table shown in table 2.

| Algorithm | Accuracy (%) |
|---|---|
| K-means | 23 |
| Page rank | 34 |
| Greedy algorithm | 55 |

Table (2) Accuracy table

## Accuracy (%)



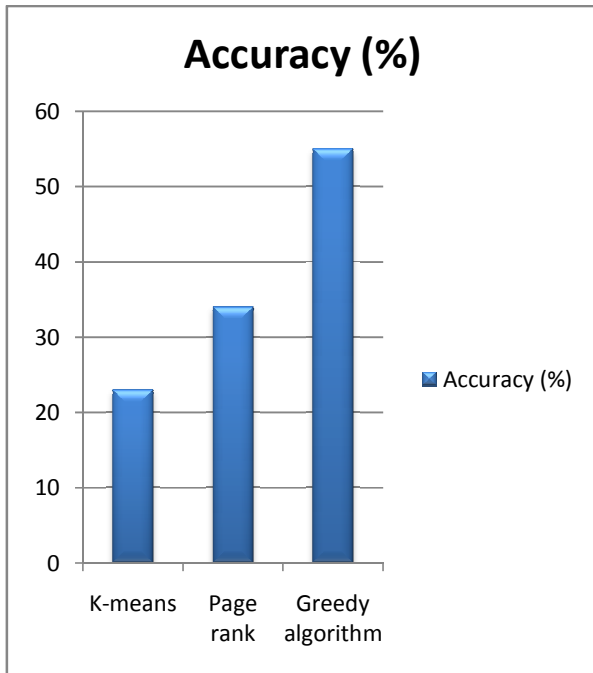**Fig4 : Performance chart**

From the performance chart, Greedy algorithm provide high level accuracy than the existing machine learning algorithms

## VI.     CONCLUSION

Personalized web search customizes the quest consequences to enhance the search first-rate for net customers. However, user's private information is probably uncovered within the consumer profile that's the premise in personalized net search. In this challenge, proposed a grouping approach for anonymizing person profiles with likability notion to bound the possibility of linking a doubtlessly touchy term to a consumer by way of p. We offered a greedy clustering technique with novel semantic similarity metric based on augmented consumer profiles in an effort to deal with the sparsity of user profiles and consider semantic relationships between person profiles. Personalized search is a promising manner to improve seeks high-quality. However, this technique calls for customers to furnish the server full access to non-public facts on Internet, which violates customers' privacy. In this work, we investigated the feasibility of reaching stability

among users' privateness and search excellent. First, a set of rules is provided to the person for gathering, summarizing, and organizing their personal facts right into a hierarchical user profile, wherein widespread phrases are ranked to better stages than precise phrases.

### REFERENCES

[1]   P. Yin and Y. Guo, "Optimization of multi-criteria website structure based on enhanced tabu search and web usage mining," Applied Mathematics and Computation, vol. 219, no. 24, pp. 11082-11095, 2013.

[2] M. Chen and Y. Ryu, "Facilitating Effective User Navigation through Website Structure Improvement," IEEE Transactions on Knowledge and Data Engineering, vol. 25, no. 3, pp. 571-588, 2013

[3] C. Kim and K. Shim, "TEXT: Automatic Template Extraction from Heterogeneous Web Pages," IEEE Transactions on Knowledge and Data Engineering, vol. 23, no. 4, pp. 612-626, 2011.

[4] Y. Yang, Y. Cao, Z. Nie, J. Zhou and J. Wen, "Closing the Loop in Webpage Understanding," IEEE Transactions on Knowledge and Data Engineering, vol. 22, no. 5, pp. 639-650, 2010

[5] J. Hou and Y. Zhang, "Effectively Finding Relevant Web Pages from Linkage Information," IEEE Transactions on Knowledge and Data Engineering, vol. 15, no. 4, pp. 940-951, 2003.

[6] A. Paranjape, R. West, L. Zia and J. Leskovec, "Improving Website Hyperlink Structure Using Server Logs," in Proceedings of the Ninth ACM International Conference on Web Search and Data Mining, 2016.

[7] H. Kao, J. Ho and M. Chen, "WISDOM: Web Intrapage Informative Structure Mining Based on Document Object Model," IEEE Transactions on Know-ledge and Data Engineering, vol. 17, no. 5, pp. 614-627, 2005.

[8] H. Kao, S. Lin, J. Ho and M. Chen, "Mining Web Informative Structures and Contents Based on Entropy Analysis," IEEE Transactions on

Knowledge and Data Engineering, vol. 16, no. 1, pp. 41-55, 2004

[9] P. Loyola, G. Martínez, Muñoz, V. J. D. K., Maldonado and C. A. P., "Combining eye tracking and pupillary dilation analysis to identify website key objects," Neurocomputing, vol. 168, pp. 179-189, 2015.

[10] M. Butkiewicz, H. Madhyastha and V. Sekar, "Characterizing Web Page Complexity and Its Impact," IEEE/ACM TRANSACTIONS ON NETWORKING, vol. 22, no. 3, pp. 943-956, 2014