



# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 6, Issue 2, February 2018

## Enhancing Information Extraction Performance for E-Commerce Systems

Sudhanshu Chourey<sup>1</sup>, Prof. Prashant Kumar Koshta<sup>2</sup>

Research Scholar, Department of Computer Science, Gyan Ganga College of Technology, India<sup>1</sup>

Professor, Department of Computer Science, Gyan Ganga College of Technology, India<sup>2</sup>

**ABSTRACT:** Web based business sites depend vigorously on outlining and breaking down the conduct of clients, making a push to impact client activities towards the advancement of achievement measurements, for example, Click through Rate, Cost per Conversion, Basket and Lifetime Value and User Engagement. Information extraction from the current web based business sites datasets, utilizing information mining and machine learning procedures, has been significantly impacting the Web promoting exercises. At the point when looked with another web based business site, the machine learning expert begins a web mining process by gathering chronicled and ongoing information of the site and breaking down/changing this information in request to be equipped for extricating data about the site structure and substance and its clients' conduct. As it were after this procedure the information researchers can manufacture pertinent models and calculations to upgrade advertising exercises. This is a costly procedure in assets and time since it will dependably rely upon the condition. We may not know a priori that a visit to a Delivery Conditions page is applicable to the expectation of a client's ability to purchase and hence would not empower following on those pages.

### I. INTRODUCTION

E-commerce is one of the most disruptive innovations in trading. Marketing and advertising techniques are used to influence costumers' behaviour, trying to increase sales and profits. Recommendation systems are one of the used techniques. Data mining and machine learning techniques had been applied to e-commerce as a way to improve e-metrics. Customer retention and engagement, click-through rate, conversion rate, shopping cart abandonment rate, customer lifetime value.

### II. LITERATURE SURVEYS

The Paper [1] proposed with the objective of demonstrating the feasibility and applicability of the developed process, In this way, the implementation consists of a chain of processes arranged such as the output of each element of the chain is the input of the next one., we firstly have our data sources, consisting of the e-commerce website itself and the usage logs associated with it. On one hand, we have the website that is generally built using common web technologies and other resources. On the other hand, we got the logs associated with the website. These logs contain information about the requests or events associated with the interaction of the user with the website. If the data is captured at the server layer, this logs contains information about the HTTP requests of the user as he navigates between website's pages. If the data is captured at the application layer, this logs contains richer data about the interaction of the user, containing not only information about the pages he visits but also information about user interactions (i.e. clicks).

# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 6, Issue 2, February 2018

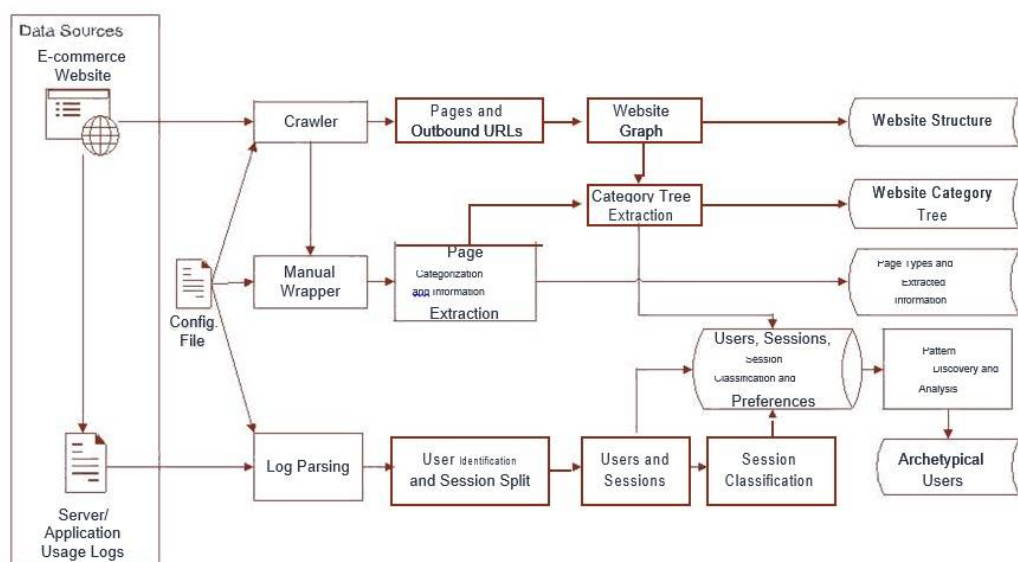


Fig 1.0 System overview

In This paper [2], The deep web contents are the information content that cannot be indexed by search engines as they stay behind searchable web interfaces. Current deep web directories mostly have less coverage of relevant web resources which degrade their ability. A Crawler goes across a variety of web pages during a crawling process. Hence to achieve efficient crawling and wide coverage, ranking and prioritizing links of different sites is necessary. The objective of this system is to extract deep web information with wide coverage for hidden web resource and uphold efficient crawling for focused web crawler.

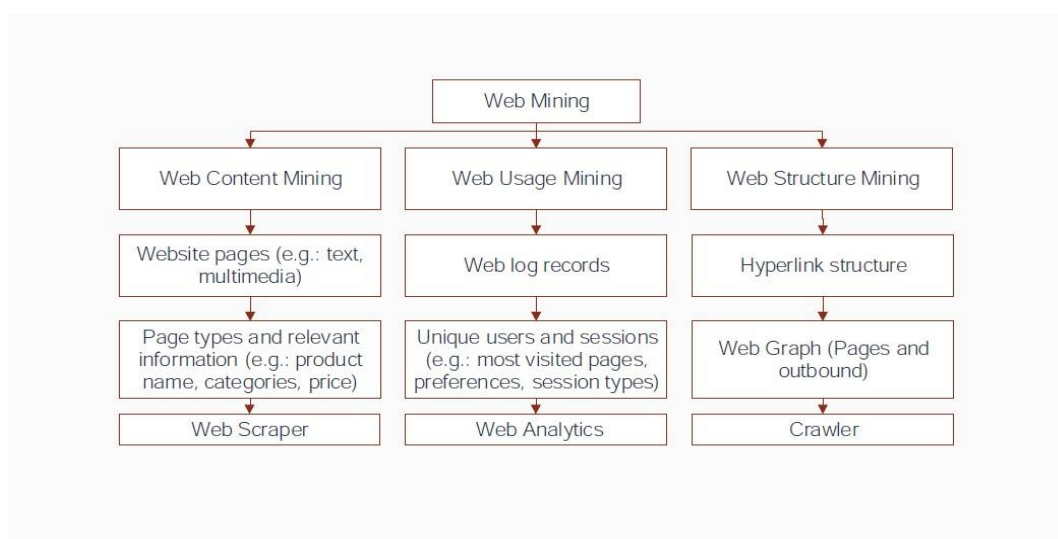


Fig 2.0



# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 6, Issue 2, February 2018

In This Paper [3] The more the number of search engines accessing a website, the more will be its visibility when searching for a particular web site. The observed results show that all search engine crawlers are not visiting all the websites. In this experiment the data set 1 was accessed by more number of search engines compared to data set 2. Certain search engines were consistent in the number of visits and number of pages crawled while a few were not consistent or irregular in their visits and pages crawled. It is found that data set 1 is more visible to search engine crawlers as it is crawled by more number of search engines compared to data set 2. The results also showed a positive correlation between the number of visits and number of pages crawled. A better search engine optimization policy can be followed to make the websites visible to different search engines so that the websites will be listed top in the search engine rankings.

In This Paper [4] The k-means clustering algorithm is one of the most commonly used data partitioning Algorithms. Despite its wide use the algorithm suffers from serious drawbacks. In this paper, an improved k-means partitioning algorithm, named SkM, is proposed dealing with the selection of the initial cluster centers. In the improved algorithm, the initial centers are selected based on a factor that utilizes the standard deviation and the max value of all the data points found in the data set, as compared to the traditional k-means that performs a random selection.

In This Paper[5] Based on the thought of K-means algorithm, the object sets of e-commerce transaction data of 300 phones can be deemed as input to be clustered, in order to get clustering center and object sets of clustering data. objects can be randomly selected from data sets as the center of initial group, then assign each object to the most similar group according to the object's mean value of the group, and update the group's mean value, calculating the object's mean value of each group. Repeat the above steps until there is no change for the number of group.

In This Paper[6] does an improved K-means clustering algorithm for identifying internet user behavior. Web data analysis includes the transformation and interpretation of web log data find out the information, patterns and knowledge discovery. The efficiency of the algorithm is analyzed by considering certain parameters. The parameters are date, time, S\_id, CS\_method, C\_IP, User agent and time taken. The research done by using more than 2 years of real data set collected from two different group of institutions web server .this dataset provides a better analysis of Log data to identify internet user behavior.

IN This Paper[7] In this experiment checked on E-commerce application where database used for non-hypertext links i.e. Deep Web. Non hypertext content indexed quicker by Google crawler with the help of sitemap metadata and robot.txt. Most of the Web's information is buried far down on dynamically generated sites, and general search engines do not find it. These traditional search engines cannot "see" or retrieve content in the Deep Web. Today wealth of information that is great source in Deep Web and therefore missed. General Search engine is crawl but not showing indexed resources in result web search engine Page. It is a great achievement in www if every Deep Web on all major sectors likes - Education, Business, Governance, Media, Career, Multimedia etc. implements aforesaid technique for the hidden public access content. This sitemap and metadata also similar useful for other general search engine like Yahoo, Bing, Ask etc. but only Robot.txt page change according targeted search engine.

In This Paper[8] we have presented vector analysis and KMeans based algorithms for mining user clusters. We have also applied the proposed algorithms to the real world data and our experimental results show the proposed algorithm is feasible, and have scalability.

### III. PROBLEM STATEMENT

Starting with the crawling step, some kind of relevance index could be calculated and associated with each link, that can become useful when combining the user's navigational path with the relevance associated with each link. A manual wrapper approach was used, but there is space to experiment with other approaches. Also, there was made two data crossings, but exists space to experiment with the crossing of data from other sources.

# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 6, Issue 2, February 2018

## IV. PROPOSED WORK

- An all-in-one approach for collecting and processing information present on a given e-commerce website and usage log files;
- A consistent and adaptable model that represents the website structure, content and users, establishing connections and relationships between the data (that could otherwise pass unnoticed);
- A reduction of the need of developing and applying a different approach for each website, trying to reduce costs and resources.

## V. IMPLEMENTATION

To find the archetypical users from our user profile database, we proceed to apply the k-means clustering algorithm.

- Keyword-based user profiles clustering;
- Session type based clustering.

The result of the application of this algorithm gives us a set of clusters that contains users with similarities between them (e.g. preferences and session types). From this set of similar groups of users we can get a grasp of the archetypical website users.

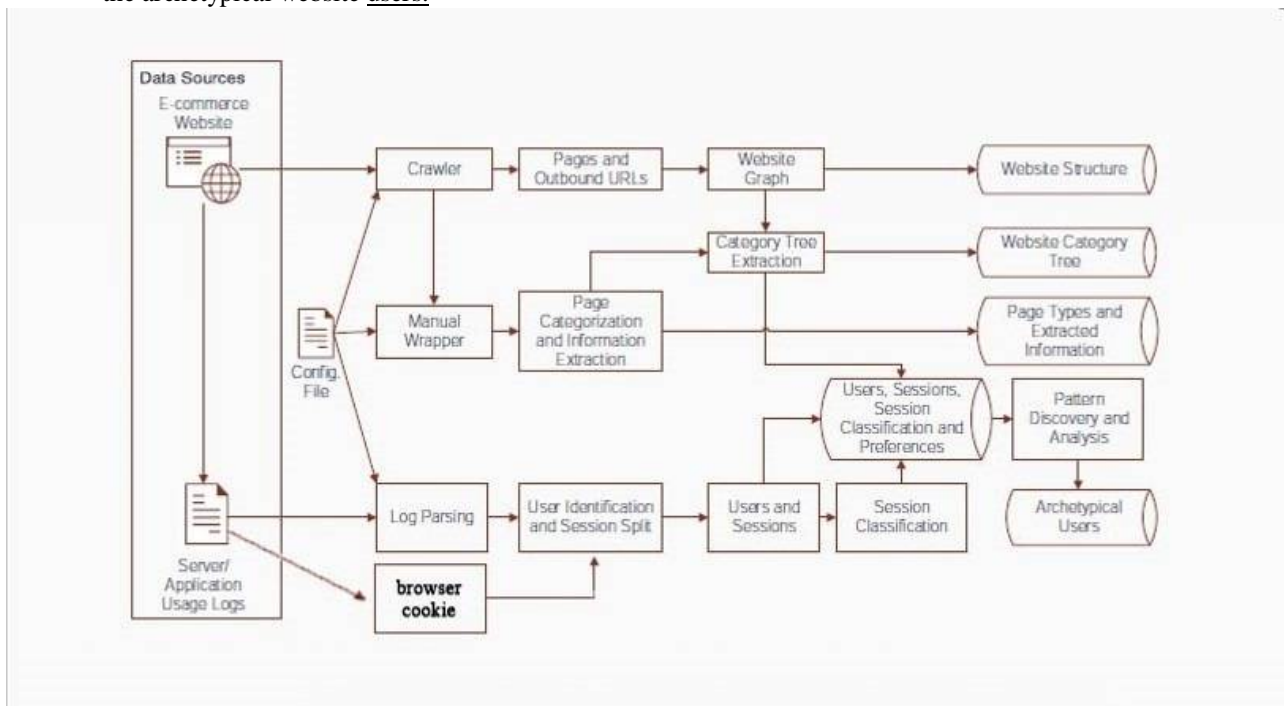


Fig. System Overview



# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 6, Issue 2, February 2018

Website's Category Tree	Keyword-based Proles
Sources: *Web Graph; *Information extracted from Pages (Categories).	Sources: *User sessions; *Information extracted from pages (Categories and Page Types).

## VI. COMPARISON AND RESULT

we consider the following Comparison:

- Improve the crawler implementing parallelism and/or prioritisation of the frontier ;
- Identify and differentiate static from dynamic hyperlinks;
- Carry experiments with another kind of web scrapers (e.g. wrapper induction);
- Increase the data crossings (e.g. cross content and usage data to get to know the favourite user brands or range of prices);
- Apply different algorithms to finding and understand the archetypical website's users (e.g. other cluster algorithms or other pattern discovery techniques);
- Analyse the possibility of expanding this methodology beyond e-commerce websites, finding other user cases

Website's Category Tree	Keyword-based Proles
<b>Output:</b> <b>Tree structure with categories and sub-categories present in the website product catalogue.</b>	<b>Output:</b> <b>Information about pages visited by category and by type in user proles.</b>

## VII. CONCLUSION AND FUTURE WORK

Since we were dealing with the web, the data sources are so large and diverse that a new field was born, called Web Mining. The research on this area becomes crucial when retrieving and dealing with the data coming from the e-commerce websites. There are two main sources of information on the web, the websites, and the usage logs. This applies to e-commerce too, and, different techniques exist to retrieve the website's pages and hyperlinks with the use of crawlers, extract the pages content through different wrappers, and, finally, identify users and sessions through different methods.

In this case, we applied the k-means clustering algorithm to find similar user groups. Other techniques may be applied to, to get other patterns or user models, and easily integrated into the information model developed.

## REFERENCES

- [1] Automating the Extraction of Static Content and Dynamic Behaviour from e-Commerce Websites By Joao Pedro Diasa., Hugo Sereno Ferreira
- [2] Enhancing Crawler Performance for Deep Web Information Extraction by Parigha V. Suryawanshi, D. V. Patil
- [3] Mining Web Logs to Identify Search Engine Behaviour at Websites by Jeeva Jose (Informatica)
- [4] Extending the k-means Clustering Algorithm to Improve the Compactness of the Clusters by Antonia Nasiakou (JPRR)
- [5] Clustering Analysis on E-commerce Transaction Based on K-means Clustering by Xuan HUANG 2014 ACADEMY PUBLISHER
- [6] Clustering of User Behaviour based on Web Log data using Improved K-Means Clustering Algorithm by S. Padmaja (IJET)
- [7] Deep Web Performance Enhance on Search Engine by Deepak Kumar (ICSCTI)
- [8] Web User Clustering Analysis based on KMeans Algorithm by linHuaXu (ICINA)