



# International Journal on Recent Researches In Science, Engineering & Technology

(Division of Computer Science and Engineering)

A Journal Established in early 2000 as National journal and upgraded to International journal in 2013 and is in existence for the last 10 years. It is run by Retired Professors from NIT, Trichy. It is an absolutely free (No processing charges, No publishing charges etc) Journal Indexed in JIR, DIIF and SJIF.

Research Paper

Available online at: [www.jrrset.com](http://www.jrrset.com)

ISSN (Print) : 2347-6729

ISSN (Online) : 2348-3105

Volume 3, Issue 11,  
November 2015.

JIR IF : 2.54

DIIF IF : 1.46

SJIF IF : 1.329

## SEMANTIC INQUIRY AND COMMON LANGUAGE QUESTIONS UTILIZING THE STRUCTURE FOR WEB CONTENT MINING

K.R. Sathishkumar<sup>1</sup>, K.S. Manojee<sup>2</sup>

<sup>1,2</sup>Assistant Professor, Department of Computer Science and Engineering, Mahendra Engineering College, Mahendhirapuri, Namakkal District, Mallasamudram, Tamilnadu, India.

**Abstract:** Searching data on World Wide Web utilizing conventional search strategies does not outcome into recovery of correct data the clients plan to get. The field of data recovery from the web is for the most part based on keyword based inquiry. However this technique misses semantic data. Recent keyword based web crawlers are not ready to search semantics in pages. In the paper, exhibit structure for semantic based web content mining (SBWCM) framework utilizing semantic ontology and SPARQL. A most important challenge in the work includes building ontology database from normal language site pages. Second challenge is mechanized interpretation of common language inquiries into more exact SPARQL questions. Proposed approach is tried on cricket sites and shows extraordinary changes over conventional keyword based indexed lists. Based on Experimental evaluations, proposed algorithm improves precision 11.71%, recall 12.73% and F-measure 11.51% of the proposed framework contrasted than previous classifiers.

**Key Words:** Semantic Inquiry, Common Language Questions, Web Content Mining, SPARQL ontology, semantic based web content mining (SBWCM).

### Introduction

Data on World Wide Web is exponentially expanding. These days' clients are more inspired by exact data on the web. Consequently, a productive method is required which mines exact data from web [1]. Semantic web systems assume vital part to convey semantic to the web mining. Execution of semantic web needs utilization of uncommon semantic innovations, for example, RDF, OWL and SPARQL.

In a present web there are billions of site pages and include data which is reasonable to be comprehended by people. In any case, this data is less helpful for programming tools in light of the fact that this data is semi organized and it doesn't contain semantics. This adds difficulties to the mechanized preparing of data. In the event of organized data like XML or RDF, mechanized handling of content and scanning data turns out to be simple for machines utilizing inquiry languages like SPARQL, however composing such question is very little agreeable to all clients [2]. Goal of this work is to diminish the gap between client's common language data and machine coherent organized data.

### Literature Survey

Aarti Singh [3] concentrates on demonstrating agent based structure for mining semantic web substance. Creator utilized clustering Techniques for giving clusters of question result which are significant. Agent based Semantic Web Mining System (SWMS) with intend to give context based and information situated outcomes to the client. Creator additionally utilized classification and clustering methods on web substance, in order to give learning based reaction to the client and

generally will point unnoticed examples. The framework includes Interface agent, accumulation agent with ontology database, clustering agent and substance mining agent. Content mining agent utilizes engaging Meta information agent and semantic Meta information agent. Creator pointed that blend of web mining systems and agent innovation will lead better outcomes.

Sharma K [4] focused learn about how to extricate the helpful data from the web and furthermore given the data and correlation about data mining. Creators present online assets for recovery of data from the web for web content mining and for distinguishing access examples of the client from web servers for web utilization mining. Creator additionally portrayed the utilization of distributed computing in web mining through distributed computing as future for web mining.

Bhatia [5] featured that looking through some data in web many time gives result that is unsatisfactory on the grounds that data came back to client is less applicable. Creators called attention to that recovering exact data from World Wide Web is as yet specialist's main area of concern. Creators additionally pointed out that semantic web are the answer for the issue. Creators recommended that utilizing the present web semantically outcome of web mining can be enhanced and which prompt working of semantic Web. The procedure of Ontology mapping is the extraction of semantics data through Grammatical Rule Extraction Technique.

Jayatilaka [6] pointed out that Semantic web has its underlying foundations on ontology's and the greater part of the ontology databases are physically fabricate which is repetitive task and in addition consuming and huge domain information is required for architect. Physically fabricating ontology has tested the development of semantic web improvement. An issues in extricating data from huge number of site pages keeping in mind the end goal to construct ontology database and proposed strategy that joins web content mining with web information utilization mining in the data extraction process. Creators have considered both the web clients and web creator's viewpoints as for web content which prompted the extraction of more objective data. The assessed outcome also demonstrated the viability of the proposed approach. Proposed method by creators will be helpful for transformation of expansive set of common language site pages to semantic web database and furthermore can be utilized to fabricate cross domain ontology database.

WANG Yong-gui [7] focused to that semantic ontology based web mining cab be valuable to enhance web services. Semantic ontology based web information mining is a mix of the semantic web and web mining. Use of Semantic web gets outcomes because of web effortlessly and in addition enhances the adequacy of web mining. The related learning of Semantic Web and Web mining and afterward creator has examined the semantic based Web Mining which proposed to assemble a semantic based Web mining model utilizing the system of Agent. Be that as it may, because of the immaturity of the important innovations and different limitations, implementation is left for future work.

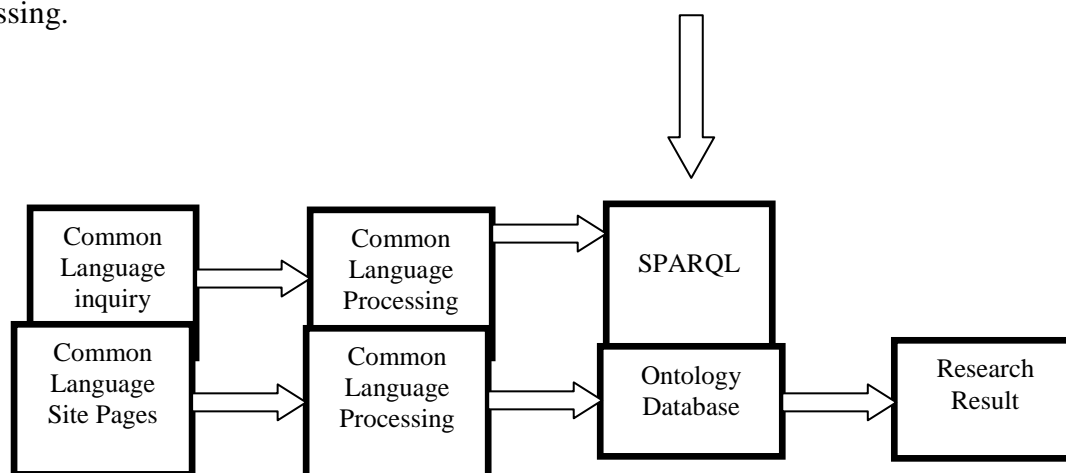
ZhusongLiu, Yuqin Zhang [8] has suggested that conventional grammar level based search has prompted low quality of outcomes because of absence of semantics in it. Creators acquainted methods of semantic web with internet business area and outlined a recovery framework with semantic ontology network framework and featured the essential methods in web based business search with semantics. Creators analyzed semantic structure and semantic recovery algorithm with conventional keyword based inquiry and reasoned that semantic search recovers more applicable data to clients search question.

Soner Kara [9] and others designed semantic ontology based data extraction from web and inquiry framework. Creators exhibited its application to soccer domain. They featured critical issues like search ease of use, scalability of framework and recovery performance in ontology based semantic inquiry. Creators proposed a keyword based semantic inquiry approach. Creators likewise featured the way that execution of the framework can be impressively enhanced utilizing area particular data extraction from site pages, domain particular derivation of extra data and applying area particular principles. Performance in proposed framework is enhanced by utilizing ordering semantic philosophy database. They implement the framework OWL and thought about the performance of framework against conventional keyword based search strategy.

Brijendra Singh [10] featured significance of web mining in a region of Data Mining and particularly significance of the extraction of data from the World Wide Web. Creators classified web mining into three unique classes in light of information they mine to be specific web data mining, web structure mining and web uses mining. Creators gave audit of past, present and future of web data mining, web structure mining and web utilizations mining. Creators had also provided future directions for inquire about in the field. Creators additionally displayed the relative investigation and synopsis of different methods utilized for web data mining with their applications and a few critical research issues.

### Proposed System

Fundamental thought behind proposed system is to decrease gap between common language pages and questions for human and proposed a semantic based web content mining (SBWCM) framework for organized semantic website pages and exact inquiry languages like SPARQL for machine processing.



**Figure.1 Work Flow Diagram of Semantic Based Web Content Mining (SBWCM)**

Figure 1 indicates essential thought behind the proposed system. Propose to keep web and search requests in common language for better ease of use. Common language site pages crawled from sites is handled utilizing domain particular formats and particular data is removed and mapped to domain particular ontology. Following is the brief explanation about every module in proposed framework.

**Crawling:** Crawler module is in responsible for gathering website pages like the conventional keyword based search engine. Crawler module executes intermittently to get refreshed website pages. Furthermore domain particular conditions can be determined in crawler module like which website pages should be crawled and which site pages should not be.

**Data extraction:** This module utilizes pages crawled from World Wide Web. Domain particular format based characteristic language processing is utilized for removing helpful data. Structure of format differs from domain to domain in the meantime also changes from site to site. This template essentially abuses weak structures in website pages like, formatting of specific data, stream of data, CSS classes utilized to show data, tables structures, page titles and metadata in pages and so forth. Different layout is wanted for each objective site. Precision of this module depends on consistency in organizing website pages specifically site, which the greater part of the expert sites already possess. This module gives machine readable values.

**Semantic Mapping:** in the module is in charge of mapping extricated data to domain particular ontology. Mapping is as far as ontology triplet in the types of <object1, connection, object2>. Data extracted from many pages may map to single question utilizing such triplets. E.g. In our illustration usage of cricket web crawler, data about players role will be taken from group squad page where as data about batting will be accessible from score summary and data about his

particular hit will be accessible from commentary page. This data is mapped to single player object utilizing distinctive triplet. Deduction is utilized to derive extra triplets from existing ones.

**Inquiry Interface:** Search interface takes common language search questions and changes over into more exact SPARQL question. Recognizing Keywords and assets is a domain particular process. While Quantifiers might be resemble, equivalent to more prominent than, not as much as, like and so forth. SPARQL questions are executed on ontology database and result is shown to the client. In conventional web index outcomes are indicating an important website page, instead of showing precise response to search inquiry. For instance “Most elevated score in ICC title trophy”, will restore a list of pages where keyword in question appears. In proposed framework result is only a solitary value.

**Result and Discussion**

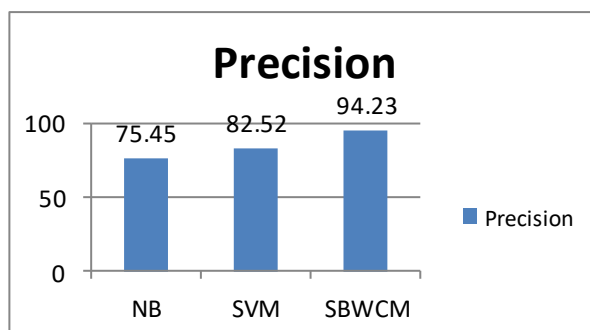
Proposed semantic based web content mining (SBWCM) framework is executed for cricket domain for ICC World Twenty20. OWL API is utilized for semantic mapping. Performance of keyword based search is contrasted and SPARQL question search. SPARQL inquiry search gives more exact outcomes than keyword based search. Common language questions are changed over into SPARQL through NLP inquiry preparing module and performance is again contrasted and direct SPARQL question input.

The proposed SBWCM strategy determines the estimation features such as precision, recall and F-measure to compute effectiveness of the proposed SBWCM framework and overcome the earlier frameworks in web site data set. In the framework decreases the gap between common language pages and questions for human.

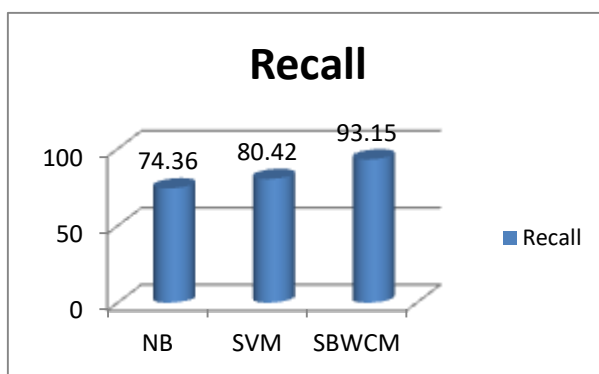
Table 1 demonstrates the precision, recall and F-measure for input features with previous frameworks. Table 1 exhibits the average value of all estimation features with input constraints. The proposed SBWCM strategy is computed with following existing methodologies namely Naïve Bayes (NB), Support Vector Machine (SVM) classifiers. According to Table 1, it noticed that SBWCM framework has the best score on each specify features for classifiers.

**Table.1 Comparison of Precision, Recall and F-measure**

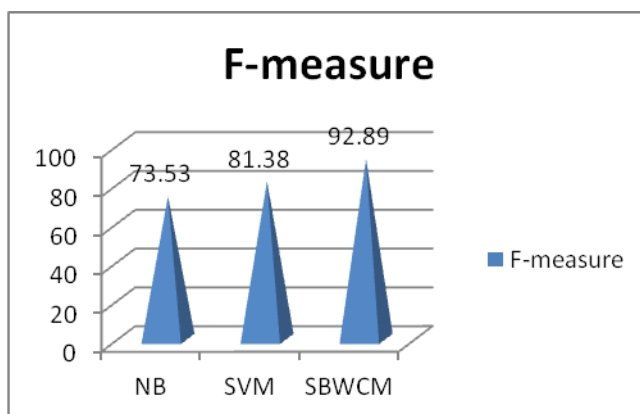
Algorithm	Precision	Recall	F-measure
Naïve Bayes	75.45	74.36	73.53
SVM	82.52	80.42	81.38
SBWCM	94.23	93.15	92.89



**Figure.2 Comparison of Precision**



**Figure.3 Comparison of Recall**



**Figure.4 Comparison of F-measure**

According to Figure 2 to 4 clarifications, it clarified the proposed SBWCM strategy is computed based on precision, recall and F-measure. Proposed SBWCM is estimated with Naïve Bayes (NB), Support Vector Machine (SVM) classifiers behalf of precision, recall and F-measure. SVM is the nearest challenger. It is classified huge number of web pages. However, SVM is not decrease the human common language web site and questions. An SBWCM framework decreases the decreases the gap between common language pages and questions for human and it improves precision 11.71%, recall 12.73% and F-measure 11.51%. Finally, the paper announces the proposed SBWCM framework is best on all several variables.

### Conclusion

In the paper we proposed a system for web content mining utilizing semantic web ideas. We likewise proposed method for changing over common language inquiries into SPARQL questions. Implementation of structure for cricket domain indicates extraordinary change over conventional keyword based searching. Composing SPARQL inquiries is troublesome for clients, yet our proposed strategy for change of common language questions into SPARQL utilized by naïve clients indicated confident outcomes. Our future work incorporates outlining area independent algorithms to change over common language website pages into ontology database.

### References

- [1] Sheila Kinsella, Aidan, Hogan, Andreas Harth, Jürgen Umbrich, Axel Polleres, and Stefan Decker. "Searching and browsing linked data with SWSE: the semantic web search engine." Web semantics: science, services and agents on the world wide web 9, no. 4 (2011): 365-401.
- [2] Danica, Damljanovic, Milan Agatonovic, and Hamish Cunningham. "Natural language interfaces to ontologies: Combining syntactic analysis and ontology based lookup through the user

interaction." In *The Semantic Web: Research and Applications*, pp. 106-120. Springer Berlin Heidelberg, 2010S.

[3] Aarti Singh, "Agent Based Framework for Semantic Web Content Mining", *International Journal of Advancement in Technology*, April 2012.

[4] Sharma K and Kumar and Shrivastava V, "Web Mining: Today and Tomorrow", in proceedings of the IEEE 3rd International Conference on Electronics Computer Technology, 2011.

[5] Jain S, and Bhatia C.S., "Semantic Web Mining: Using Ontology Learning and Grammatical Rule Interface Technique", In IEEE 2011.

[6] Wimalarathne G.D.S.P, and Jayatilaka A.D.S "Knowledge Extraction for Semantic Web Using Web Mining", *The International Conference on Advances in ICT for Emerging Regions - ICTer2011*, IEEE, 2011.

[7] JIA Zhen, WANG Yong-gui, "Research on Semantic Web Mining", *International Conference on Computer Design and Applications*, IEEE, 2010.

[8] Yuqin Zhang, ZhusongLiu, "Research and Design of E-commerce Semantic Search", *3rd International Conference on Information Management, Innovation Management and Industrial Engineering*, IEEE, 2010.

[9] Nihan K. C, zgur Alan, Soner Kara, OrkuntSabuncu, SametAkpınar, Ferda N. Alpaslan, "An Ontology-Based Retrieval System Using Semantic indexing", IEEE 2010.

[10] Hemant Kumar Singh, Brijendra Singh, "Web data mining research: a survey", IEEE, 2010.