# Model for Information Extraction and Information Retrieval Based on Cascaded Support Vector Machine and Feature Vector Optimization

Ritesh Kumar Shah<sup>1</sup>, Dr. Suresh Jain

<sup>1</sup> Ph.D Scholar, Department of Computer Science and Engineering, Mewar University, Gagnar, Raj <sup>2</sup> Director, Prestige Institute of Engineering Management & Research

## ABSTRACT

Information extraction and optimization play vital role in semantic web mining. Semantic web mining process proceeds the better way information retrieval over the ocean of internet database. In this paper design model of information optimization using feature vector optimization. The feature vector optimization process reduces the relational key constraints for the categorization of information for retrieval. Relative information categorization increases the hit ratio of retrieving data. For the better categorization of relative information used cascaded support vector machine. The cascaded support vector machine is two stages multi-kernel-based classification process. The two-stage support vector machine reduces the semantic gap between the search query and internet database. The design model removes the bottleneck of relational constraints of the post mining process. The feature vector optimization is used to find optimal attribute for the label of the cascaded support vector machine. The labeled class of support vector machine map the domain ontology database.

Keyword: - Ontology, Information Extraction, Feature Vector, CSVM, Domain

### **1. INTRODUCTION**

The rapid growth of internet database diverts the accuracy of information retrieval, for the betterment of information retrieval over the internet used semantic web mining. The semantic mining approach faced a problem of large amounts of links, tag and unstructured data on web pages. For the better retrieval of information used various data mining and attribute-based optimization technique[1, 3]. the labeling of information enhances the capacity of information retrieval. for the labeling of information used various data mining algorithms such as clustering, classification and rule mining. The classification algorithms such as support vector machine and other regression algorithms reduces the bottleneck of rules and improve the retrieval capacity [4-9]. The process of optimization reduces the semantic gap between the search query and retrieve information. For the optimization of keywords, sentence and segment of words used feature vector optimization process. The feature vector optimization process provides the optimal feature attribute of the database. It also reduces the constraints of relation mapping of words to making sentences. The cascaded support vector machine classifies the optimal words for the retrieval of information. The cascaded support vector machine used multi-kernel function for the mapping of feature attributes according to their relation[1, 10, 13]. The attribute relation mapped the domain ontologies of the database. The domain ontologies of database are not modified, just is extension of relative rules of statements[12]. The ontologies-based web mining used in various disciples of information process such as medical science, news portal, hotel industry and many more areas. In the current scenario of web-based information extraction process used regular ontologies. A regular ontology gives the predefined set of data. the regular ontologies are failed to retrieve accurate data over the internet[11]. The process of ontologies combined with optimization algorithms and achieve the better information retrieval. However, Internet archived data are growing rapidly because of the sharing of information between various systems' existing ontology-based systems can extract appropriate information to a limited extent. Now

cascaded support vector machine (CSVM) with feature vector optimization is an effective process of information retrieval over the internet[15-19]. The rest of the paper describes as in section II describe the information optimization using feature vector. In section III. Describe the process of the cascaded support vector machine. in section IV describe the proposed mode and ontologies and finally discuss the conclusion and future work in section V

#### 2. INFORMATION EXTRACTION & OPTIMIZATION

Information extraction over the internet is a very difficult task, due to large amount of links, tags, and unstructured data in web pages. For the extraction of information used various algorithms such as natural language processing and some HTMLRSS technique. these techniques extract the URLs and link and store in database[21-23]. The stored links create some patterns and extract the relevant information for the processing of queries. For the optimization of information used feature vector optimization process[19-20]. The feature vector optimization process removes the redundant URLS, links and keywords for the process of information retrieval. the process of logic describes below here[24-26].

Define some term for process of algorithms URLS U links L Keyword K, pattern P, Redundant R

Input:

FVO (U, L, K)

**Output:** 

Optimal (U, L.K)

 $U_i \leftarrow \emptyset, L_o \leftarrow \emptyset, K \leftarrow \emptyset$ 

Start fetching of page documents

For all  $u_i \in u_v$  do

$$U_i \leftarrow \phi, U_n \leftarrow \phi, U_d \leftarrow \phi, U_o \leftarrow \phi$$

$$v_v \leftarrow R(u_v), U_v \leftarrow V_r$$

If  $u_v$  is-A Relation

Create feature vector of attribute of U

else if Vr is with P then

create an optimal set of Ui

end if

$$Lp \leftarrow R(L_v), L_v \leftarrow L_r$$

if Lv is IS-A then

then creates pattern  $L_p$ 

end if

 $k_p \leftarrow R(k_v), k_v \leftarrow k_r$ 

# if *kv* is IS-A then then creates pattern $k_p$ for all Patterns of P create a feature vector UV(u1,u2,....un), RV(r1,r2,...rn)KV(k1,k2,k3,....kn)creates the patterns weight of U, R, K as 1 end if for all patterns *P* do create an optimal value of U create optimal value of R create optimal value of K as all weight value is 1. end for end for for all input patterns P do for all mapping relation do creates segments of patterns end for $Ui \leftarrow u_i \cup, R_o \leftarrow R_o \cup UK_o \leftarrow Ko$ do end for for all optimal value of pattern. end for



Figure 1: process block diagram of feature vector optimization

### 3. CASCADED SUPPORT VECTOR MACHINE

The cascaded support vector machine is two stage classifiers. The two-stage classifier used for the labeling of domain ontology relation data. A query vector is then evaluated by every function in the cascade in turn and if at any point it is classified negative the evaluation stops[30]:

$$f_c(x) = sgn(f_1(x)) sgn(f_2(x)) \dots \dots \dots \dots \dots (1)$$

Where,  $f_c(x)$  is the cascade evaluation function. In other words, the decision functions in the cascade can be biased in such a way that their negative classification is very confident while the positive decisions are passed on to the next, more complex function.

In practice, the original full SVM (form ) can run on all the queries that pass through the cascade:

$$F(x) = \begin{cases} -1 & \text{if } f_c(x) < 0\\ sgn(f_{svm}(x)) & \text{if } f_c(x) \ge 0 \end{cases}$$
(2)

Biasing of the functions is done by setting the offset parameter b to achieve a desired accuracy of the function of an evaluation set. Romdhani introduces a method to model a desired receiver operator curve, although we settled for a

simpler approach, requiring that all positive objects in the original training set that are correctly classified by the full SVM have to be classified correctly by every level of the cascade as well and setting the offsets accordingly [27-29].

### 4. PROPOSED METHODOLOGY

The process of model is combination of the cascaded support vector machine, feature vector optimization and domain ontology. The process of feature vector optimization reduces irrelevant link, URLs and keywords. The optimal value of FVO passes through the cascaded support vector machine. the cascaded support vector machine(CSVM) mapped the relation with domain ontologies(DO). The process of descriptions given below. Mapped the feature vector data pattern  $P_r \in \Re^D$  in the cascaded support vector machine in a two-stage process. The mapped data of CSVM define the class level of domain ontology(DO). The labeling of class to DM creates semantic accuracy of query retrieval(UR).

1. Input: FVO data patterns

Output: *Query* Retrieval

- 2. Compute  $DM_{(p,k)}$  and k similar(RD)
- 3. for all  $FVO \in DM_{(p_r,k)}$  do
- 4. estimate-Relation (pt, FVO)
- 5. end for
- 6.  $CSVM \leftarrow DM (p_t,k)$  {label of class map
- 7. for all  $FVO \in CSVM$  and  $DM \in UR$  do
- 8. iterate k pattern (RD) and rmap pattern (DM, UR)
- 9. if UR<sub>(Label)</sub> then
- 10.  $CSVM \leftarrow CSVM \cup \{DM\}$
- 11. end if
- 12. end for
- 13. for all  $FVO \in CSVM$  do
- Update fvo(Label) and class({Relation})
- 15. end for
- 16. retrieve user query result
- 17. return DM



Figure 2: process block diagram of user query retrieval with cascaded support vector machine

## 4. CONCLUSIONS

In this paper presents the new model for the retrieval of web-based information using Domain ontology and cascaded support vector machine. the cascaded support vector machine proceeds the optimal data of feature vector optimization. The feature vector optimization proceeds the data of information extraction using an information extractor. The extracted information contains URLs, links and keywords. The optimal value of FVO reduces the irrelevant relation of patterns for the mapping of the cascaded support vector machine. The cascaded support vector machine mapped the relation label of domain ontology of any data. the process of model reduces the semantic gap of user query and information retrieval. In future the proposed model implements in news portal, hotel industry and many more ontology-based information retrieval system

### **5. REFERENCE**

[1] Farman Ali, Daehan Kwak, Pervez Khan, Shaker Hassan A. Ei-Sappagh, S. M. Riazul Islam, Daeyoung Park and Kyung-Sup Kwak "Merged Ontology and SVM-Based Information Extraction and Recommendation System for Social Robots", IEEE, 2017, Pp 12364-12379.

[2] Nandhini M, Janani M and Dr.Sivanandham S.N "Association Rule Mining Using Swarm Intelligence and Domain Ontology", IEEE, 2012, Pp 537-541.

[3] Sanghamitra Bandyopadhyay and Koushik Mallick "A New Feature Vector Based on Gene Ontology Terms for Protein-Protein Interaction Prediction", IEEE, 2017, Pp 762-770.

[4] Dr Kehinde K. Agbele, Eniafe F. Ayetiran, Kehinde D. Aruleba and Daniel O. Ekong "Algorithm for Information Retrieval Optimization", IEEE, 2016, Pp 1-8.

[5] Belainine Billal, Alexsandro Fonseca and Fatiha Sadat "Efficient Natural Language Pre-processing for Analyzing Large Data Sets", IEEE, 2016, Pp 3864-3871.

[6] Muhammad Rio Bastian and Ayu Purwarianti "Information Extraction in Statistics Indicator Tables using Rule Generalizations and Ontology", ICITSI, 2016, Pp 1-6.

[7] Sonia Haiduc, Venera Arnaoudova, Andrian Marcus and Giuliano Antoniol "The Use of Text Retrieval and Natural Language Processing in Software Engineering" IEEE, 2016, Pp 898-899.

[8] Sharon Gower Small and Larry Medsker "Review of Information Extraction Technologies & Applications", Springer, 2014, Pp 1-29.

[9] Christopher Baechle, Ankur Agarwal, Ravi Behara and Xingquan Zhu "Co-Occurring Evidence Discovery for COPD Patients using Natural Language Processing", IEEE, 2017, Pp 321-324.

[10] S.S. Dhenakaran and S.Yasodha "Semantic Web Mining - A Critical Review", International Journal of Computer Science and Information Technologies, 2011, Pp 2258 – 2261.

[11] G. Babu and Dr. T. Bhuvaneswari "association rule mining for identifying optimal customers using MAA algorithm", Journal of Theoretical and Applied Information Technology, 2014, Pp 829-838.

[12] Farman Ali, Kyung-Sup Kwak and Yong-Gi Kim "Opinion mining based on fuzzy domain ontology and Support Vector Machine: A proposal to automate online review classification", Applied Soft Computing, 2016, Pp 235–250.

[13] Wafa Damak, Imene Khanfir Kallel and Issam Rebai "Semantic object recognition by merging decision tree with object ontology", Advanced Technologies for Signal and Image Processing, 2014, Pp 65-70.

[14] Hongsheng Xu, Ruiling Zhang, Chunjie Lin and Wenli Gan "Semantic Annotation of Ontology by Using Rough Concept Lattice Isomorphic Model", International Journal of Hybrid Information Technology, 2015, Pp 93-108.

[15] S. Chitra and G. Aghila "A Survey on Tools and Algorithms of Ontology Operations", Research Journal of Engineering Sciences, 2014, Pp 12-25.

[16] Ismat Ara Reshma, Md Zia Ullah and Masaki Aono "Ontology based Classification for Multi-Label Image Annotation", ICAICTA, 2014, Pp 1-7.

[17] Dr. S. Chitra "A Novel Approach to Analyse User Satisfaction Level on Web pages using Ontologies", IRJET, 2017, Pp 294-203.

[18] V'aclav Jirkovsky, Petr Kadera and Nestor Rychtyckyj "Semi-Automatic Ontology Matching Approach for Integration of Various Data Models in Automotive", Springer, 2017, Pp 1-14.

[19] Vindula Jayawardana, Dimuthu Lakmal, Nisansa de Silva, Amal Shehan Perera, Keet Sugathadasa and Buddhi Ayesha "Deriving a Representative Vector for Ontology Classes with Instance Word Vector Embeddings", arXiv, 2017, Pp 1-6.

[20] Zhen-Shu Mi, Ahmad C. Bukhari and Yong-Gi Kim "An Obstacle Recognizing Mechanism for Autonomous Underwater Vehicles Powered by Fuzzy Domain Ontology and Support Vector Machine", Hindawi Publishing Corporation, 2014, Pp 1-11.

[21] O. Abuomar, S. Nouranian, R. King, T.M. Ricks and T.E. Lacy "Comprehensive mechanical property classification of vapor-grown carbon nanofiber/vinyl ester nanocomposites using support vector machines", Computational Materials Science, 2015, Pp 316–325.

[22] Maciej Zieba, Jakub M. Tomczak, Marek Lubicz and Jerzy Swiatek "Boosted SVM for extracting rules from imbalanced data in application to prediction of the post-operative life expectancy in the lung cancer patients", Applied Soft Computing, 2014, Pp 99–108.

[23] Kyle Williams, Jian Wu, Sagnik Ray Choudhury, Madian Khabsa and C. Lee Giles "Scholarly Big Data Information Extraction and Integration in the CiteSeer Digital Library", IEEE, 2014, Pp 1-6.

[24] Jaime Zabalza, Jinchang Ren, Jiangbin Zheng, Junwei Han, Huimin Zhao, Shutao Li and Stephen Marshall "Novel Two-Dimensional Singular Spectrum Analysis for Effective Feature Extraction and Data Classification in Hyperspectral Imaging", IEEE, 2015, Pp 1-32.

[25] Jie Tao, Amit V. Deokar and Omar F. El-Gayar "An Ontology-based Information Extraction (OBIE) Framework for Analyzing Initial Public Offering (IPO) Prospectus", Hawaii International Conference on System Science, 2014, Pp 769-778.

[26] Zehra , Camlica1, H.R. Tizhoosh and Farzad Khalvati "Medical Image Classification via SVM

using LBP Features from Saliency-Based Folded Data", IEEE, 2015, Pp 1-5.

[27] Sharon Gower Small and Larry Medsker "Review of Information Extraction Technologies & Applications", Neural Computing and Applications, 2013, Pp 1-28.

[28] Esraa Elhariri, Nashwa El-Bendary, Mohamed Mostafa M. Fouad, Jan Platos, Aboul Ella Hassanien and Ahmed M.M.Hussein "Multi-class SVM Based Classification Approach for Tomato Ripeness\*", Springer, 2014, Pp 175-188.

[29] Wenjia Li, Jigang Ge and Guqian Dai "Detecting Malware for Android Platform: An SVM-based Approach", IEEE, 2015, Pp 464-469.

[30] Ignas Kukenys and Brendan McCane "Classifier Cascades for Support Vector Machines", IEEE, 2008, Pp 1-6.