

An Improved Support Vector Machine Classifier Using AdaBoost and Genetic Algorithmic Approach towards Web Interaction Mining

B. Kaviyarasu

Research Scholar, PG and Research Department of Computer Computer Applications,
Hindusthan College of Arts and Science, Coimbatore – 38.

Email: bkaviyarasuphd@gmail.com

Dr. A. V. Senthil Kumar

Director, PG and Research Department of Computer Computer Applications,
Hindusthan College of Arts and Science, Coimbatore – 38.

Email: avsenthilkumar@gmail.com

ABSTRACT

Predicting the objective of internet users has divergent applications in the areas such as e-commerce, entertainment in online, and several internet-based applications. The critical part of the classifying internet queries based on obtainable features namely contextual information, keywords and their semantic relationships. This research paper presents an improved support vector machine classifier that makes use of ad boost genetic algorithmic approach towards web interaction mining. Around 31 participants are chosen and given topics to search web contents. Parameters such as precision, recall and F1 score are taken for comparing the proposed classifier with the classical support vector machine. Results proved that the proposed classifier achieves better performance than that of the conventional SVM.

Keywords: Web interaction mining, algorithm, support vector machine, classifier, ad boost, genetic algorithm.

Date of Submission: April 03, 2017

Date of Acceptance: April 13, 2017

1. INTRODUCTION

Web mining is the application of data mining methods to extract knowledge from internet information, together with internet documents, hyperlinks between records, usage logs of web sites, and many others. Web mining is the withdrawal of potentially valuable patterns and implicit understanding from pastime related to the site. This extracted knowledge will also be extra used to enhance web utilization such that prediction of subsequent page likely to accessed through consumer, crime detection and future prediction, person profiling and to recognize about person searching hobbies [Monika Dhandi, Rajesh Kumar Chakrawarti.,2016] [8].

Web Mining can be comprehensively isolated into three particular classes, as indicated by the sorts of information to be mined. The review of the three classifications of web mining [T. Srivastava et al.,2013] [11] discussed below are (1) Web Content Mining (2) Web Structure Mining (3) Web Interaction Mining.

1.1. Web Content Mining (WCM): WCM is the way toward extricating helpful data from the substance of web archives. Depicted information relates to the gathering of certainties of a web page were intended to pass on to the clients. It might comprise of content, pictures, sound, video, or organized records, for example, records and tables.

1.2. Web Structure Mining (WSM): The structure of a distinctive web comprises of Web pages as nodes,

and web link as edges associating related pages. Web Structure Mining is the way toward finding structure data from the Web. This can be further partitioned into two sorts in view of the sort of structure data utilized.

- (a) Hyperlinks: A Hyperlink is a basic unit that interfaces an area in a page to stand-out region, either inside the indistinguishable web page or on an alternate page.
- (b) Document Structure: Moreover, the substance inside a page will likewise be composed in a tree-organized structure, headquartered on the more than a couple of HTML and XML labels inside the website page. Mining endeavors right have intrigued undoubtedly by separating document object model (DOM) structures out of documents.

1.3. Web Interaction Mining (WIM): WIM is the use of data mining procedures to find intriguing utilization designs from Web information, with a specific end goal to comprehend and better serve the requirements of Web-based applications. Use of information catches the character or source of web clients alongside their perusing conduct at a webpage. WUM itself can be grouped further contingent upon the sort of use information considered:

- (a) Web Server Data: The client logs are gathered by Web server. Small range of the information incorporates IP address, page reference and get to time.

- (b) Application Server Data: Commercial application servers, for example, Web-logic, Story-Server have noteworthy components to empower E-trade applications to be based on top of them with little exertion. A key component is the capacity to track different sorts of business occasions and log them in application server logs.
- (c) Application Level Data: New sorts of occasions can be characterized in an application, and logging can be turned on for them - producing histories of these uniquely characterized occasions.

2. RELATED WORKS

T. Cheng et al.,2013 [9] have provided three information offerings: entity synonym information carrier, query-to-entity information service and entity tagging knowledge provider. The entity synonym service used to be an in-creation knowledge carrier that used to be presently available whilst the other two are information services presently in progress at Microsoft. Their experiments on product datasets exhibit (i) these knowledge offerings have excessive best and (ii) they've gigantic influence on consumer experiences on e-tailer web sites.

M. Nayrolles and A. Hamou-Lhadj.,2016 [7] proposed BUMPER (BUg Metarepository for dEvelopers and Researchers), a customary infrastructure for developers and researchers inquisitive about mining information from many (heterogeneous) repositories. BUMPER used to be an open supply web-founded environment that extracts information from a variety of BR repositories and variant manipulate systems. It was once equipped with a strong search engine to aid customers quickly query the repositories utilizing a single point of access. X.

Ye et al.,2015 [12] authors proposed a new studying method by means of a generalized loss function to capture the subtle relevance variations of training samples when a extra granular label constitution was once on hand. Authors have utilized it to the Xbox One's movie search mission the place session-headquartered person conduct understanding was once to be had and the granular relevance differences of coaching samples are derived from the session logs. When put next with the prevailing method, their new generalized loss function has tested sophisticated experiment efficiency measured by means of a few consumer-engagement metrics.

The purpose of T. F. Lin and Y. P. Chi.,2014 [10] was to make use of the applied sciences of TF-IDF, ok-approach clustering and indexing high-quality examination to establish the combo of key phrases to be able to advantage seo. The learn demonstrated that it might probably comfortably enhance the internet site's advancement of ranking on search engine, increase

internet site's publicity level and click on through expense.

G. Dhivya et al.,2015 [3] analyzed person conduct by using mining enriched web entry log information. The few net interaction mining approaches for extracting valuable elements used to be discussed and employ all these strategies to cluster the users of the domain to study their behaviors comprehensively. The contributions of this thesis are an information enrichment that was content and starting place situated and a treelike visualization of generic navigational sequences. This visualization makes it possible for a conveniently interpretable tree-like view of patterns with highlighted primary know-how.

Z. Liao et al.,2014 [15] introduced "task trail" to understand user search behaviors. Authors outline a mission to be an atomic person know-how want, whereas a challenge trail represents all person pursuits inside that precise project, equivalent to question reformulations, URL clicks. Previously, net search logs have been studied by and large at session or question stage the place customers may put up several queries within one venture and manage several tasks inside one session.

A. Yang et al.,2014 [2] have awarded a solution that first identifies the customers whose kNN's possibly plagued by the newly arrived content, after which replace their kNN's respectively. Authors proposed a new index constitution named HDR-tree in order to support the effective search of affected customers. HDR-tree continues dimensionality reduction through clustering and principle element evaluation (PCA) so as to make stronger the search effectiveness. To extra scale back response time, authors proposed a variant of HDR-tree, known as HDR-tree, that helps extra effective but approximate solutions.

A. U. R. Khan et al.,2015 [5] have presented a cloud carrier to explain how the status of the mass media news can be assessed utilizing users online utilization habits. Authors used knowledge from Google and Wikipedia for this comparison challenge. Google data was helpful in understanding the have an effect on of stories on web searches whereas data from Wikipedia enabled us to understand that articles related to rising information content additionally find lot of attention.

J. Jojo and N. Sugana.,2013 [4] proposed a hybrid approach which uses the ant-founded clustering and LCS classification methods to seek out and predict user's navigation behavior. As a result user profile may also be tracked in dynamic pages. Personalized search can be used to address project in the internet search community, founded on the premise that a consumer's normal choice may just aid the quest engine disambiguate the real intention of a question.

M. A. Potey et al.,2013 [6] reviewed and compared the to be had approaches to present an insight into the

discipline of query log processing for expertise retrieval.

A. Vinupriya and S. Gomathi.,2016 [1] proposed a brand new scheme named as WPP (web page Personalization) for powerful net page suggestions. WPP consist of page hit rely, complete time spent in each hyperlink, number of downloads and link separation. Founded on these parameters the personalization has been proposed. The procedure proposes a brand new implicit user feedback and event hyperlink access schemes for amazing internet web page customization together with domain ontology.

Y. C. Fan et al.,2016 [14] proposed an information cleansing and understanding enrichment framework for enabling consumer alternative working out by way of Wi-Fi logs, and introduces a sequence of filters for cleansing, correcting, and refining Wi-Fi logs.

Y. Kiyota et al.,2015 described learn how to construct a property search habits corpus derived from micro blogging timelines, in which tweets concerning property search are annotated. Authors applied micro task-established crowd sourcing to tweet knowledge, and construct a corpus which contains timelines of special customers that are annotated with property search phases.

In our previous works extreme learning machine classifier [16] and penta layered artificial neural networks [17] are used for web interaction pattern mining.

3. PROPOSED WORK

3.1. Background of Support Vector Machine

The support vector machines (SVM) classifier is one among the machine learning algorithms which is based on the structural risk minimization principle and statistical learning theory. The basic idea of SVM is to transform the data into a higher dimensional space and find a classification hyper-plane that separates the data with the maximum margin. The standard SVM model is as follows:

$$\min \frac{1}{2} \| \omega \|^2 + C \sum_{i=1}^l \xi_i$$

$$s.t \ y_i ((\omega \cdot \phi(x_i)) + b) \geq 1 - \xi_i \quad \dots (1)$$

$$\xi_i \geq 0, i = 1, 2, \dots, l$$

where $n \ x_i \in R^n$ and $y_i \in \{ -1, +1 \}$ are the training samples and the corresponding class label respectively, ϕ is a nonlinear map that transforms the data to the high dimensional feature space, ω is the normal vector to the bounding plane, b is a bias value, ξ_i ($i= 1, 2, \dots, l$) are the slack variables, and C is a penalty parameter.

Instead of solving the above said optimization problem, it is feasible to solve the dual problem:

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l y_i y_j \alpha_i \alpha_j k(x_i, x_j) - \sum_{j=1}^l \alpha_j$$

$$s.t \ \sum_{i=1}^l y_i \alpha_i = 0$$

$$0 \leq \alpha_i \leq C, i = 1, 2, \dots, l$$

... (2)

where $k(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$ is called kernel function. The binary classifier $g(x)$ and the decision function $f(x)$ could be calculated as follows:

$$f(x) = \text{sign}(g(x)) = \text{sign} \left(\sum_{x_i \in SV} y_i \alpha_i^* k(x, x_i) + b \right)$$

(3)

is α_i^* an optimal solution of the problem (2), SV is the set of support vectors, $b^* = \sum_{x_i \in SV} y_i \alpha_i^* k(x, x_i)$. if $0 < \alpha_i^* < C$.

3.2. Support Vector Machine Classifier with Adaboost Genetic Algorithmic Approach

The AdaBoost (adaptive boosting) algorithm is one of the most popular ensemble methods. It creates a collection of moderate classifiers by maintaining a set of weights over training data and adjusting these weights after each learning cycle adaptively. The weights of the training samples which are correctly classified by current classifier will decrease while the weights of the samples which are misclassified will increase. Since the proposed work has the scope to incorporate AdaBoost algorithm, the standard SVM needs to be extended to the SVM for which each training sample has different weights. Now the SVM model is transformed to

$$\min \frac{1}{2} \| \omega \|^2 + C \sum_{i=1}^l \mu_i \xi_i$$

$$s.t \ y_i ((\omega \cdot \phi(x_i)) + b) \geq 1 - \xi_i \quad \dots (4)$$

$$\xi_i \geq 0, i = 1, 2, \dots, l$$

where (μ_1, \dots, μ_l) indicates the weight of the sample x_i . And the dual problem is as follow:

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^i \sum_{j=1}^i y_i y_j \alpha_i \alpha_j k(x_i, x_j) - \sum_{j=1}^i \alpha_j$$

$$s.t \sum_{i=1}^i y_i \alpha_i = 0$$

$$0 \leq \alpha_i \leq C, i = 1, 2, \dots, l$$

... (5)

Except for the weights of samples, the classification performance of the proposed Adaboost-SVM is equally affected by its model parameters. During the AdaBoost iterations, if parameters make the classification accuracy of ISVM less than 50%, the requirement on a component classifier in AdaBoost cannot be satisfied. In contrast, if the accuracy is too high, boosting classifiers may become inefficient because the errors of these component classifiers are highly correlated. Hence, how to select appropriate model parameters is important. There are many evolutionary algorithms for searching the suitable solution in real-valued spaces. With the advantages consisting of parallel search, solving complex problems, and large search space, the genetic algorithm (GA) is applied to perform the model parameters selection in the k-fold cross-validation set.

However, the process is time consuming and may cause the overfitting. Therefore, we adjust the model selection procedure so that the GA could be stopped when the cross-validation accuracy is over 0.5. In this result, the component classifiers conform to the condition of AdaBoost and the computational time could be saved. Moreover, the randomness of result produced by GA would decrease significantly after many times of Adaboost iterations. As a result, the outcomes corresponding to several independent runs of the hybrid method are similar to each other. So, the process of the proposed ISVM mechanism is relatively stable.

AdaBoost-GA-ISVM algorithm

Step1: Input - Training samples $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $x_i \in R^u$ and $y_i \in \{1, \dots, K\}$; Moderate classifier ISVM; The number of classes K ; The total number of the iterations T .

Step 2: Initialize: The weights of training samples: $w_i^1 = 1/n, i = 1, \dots, n$; The GA parameters include size of population N ; Maximum number of generations $MaxI$; Length of chromosome of C and kernel parameters l ; Crossover rate p_c and mutation rate p_m .

Step 3: For $t = 1, 2, \dots, T$.

a) Select appropriate parameters

i) Encode the parameter C and kernel parameters as an l-bit string which consists of l_1 bits standing for C and l_2 bits standing for kernel parameters, here $l = l_1 + l_2$; Generate an initial population consisted of N strings of

binary bit. To avoid trapping into same local optimum in GA process, the parameters are estimated starting from a completely new initial population for each t .

ii) For $j = 1, 2, \dots, N$, obtain the j th group of parameters by decoding the string j and train a component multi-classifier ISVM g_i^j using these parameters on the k -fold cross-validation data set.

iii) Calculate the average cross-validation error:

$$E_t = \frac{1}{N} \frac{1}{n} \sum_{j=1}^N \left(\sum_{i=1}^n I(Y_i, g_i^j(x_i)) \right),$$

where the

indicator function I produces 0 if the arguments are equal and 1 if they are different.

iv) If $E_t < 0.5$, do step v) else the parameters satisfy the requirement of AdaBoost

v) Perform reproduction: selection, crossover and mutation. Then,

a) Generate new offspring population. If the number of generation exceeds $MaxI$, it means that the t -th moderate classifier is invalid, then stop the GA iteration

b) Train a multi-classifier ISVM G_t using the suitable parameters from a) and obtain a probabilistic output vector

c) Compute the training error of G_t .

d) Set weight for the current classifier

$$G_t : \alpha_t = 0.5 \ln \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$

e) Update the weights:

$$w_i^{t+1} = w_i^t \exp \{-\alpha_t y_i G_t(x_i)\}$$

Step 4: Output: When the number of valid classifier reaches T the proposed algorithm is completed.

4. Experimental Results

31 participants are taken in order to build the dataset for evaluating the proposed model. The people that are chosen belong to heterogeneous age groups and web experience; similar considerations apply for education, even though the majority of them have a computer science or technical background. All participants were requested to perform ten search sessions organized as follows:

- Four guided search sessions;
- Three search sessions in which the participants know the possible destination web sites;
- Three free search sessions in which the participants do not know the destination web sites.

This led to 129 sessions and 353 web searches, which were recorded and successively analyzed in order to manually classify the intent of the user according to the two-level taxonomy. Starting from web searches, 490 web pages and 2136 sub pages were visited. The interaction features were logged by the inbuilt YAR plug-in that is present in Google Chrome web browser.

For performing query classification, the proposed ISVM presumes that the queries in a user session are independent; Conditional Random Field (CRF) considers the sequential information between queries, whereas Latent Dynamic Conditional Random Fields (LDCRF) models the sub-structure of user sessions by

assigning a disjoint set of hidden state variables to each class label.

In order to evaluate the effectiveness of the proposed model, we adopted the classical evaluation metrics of Information Retrieval: precision, recall, and F1-measure. In order to simulate an operating environment, 60% of user queries were used for training the classifiers, whereas the remaining 40% were used for testing them. The values of the precision, recall and F1-Score of the participants are given in Annexure 1.

Precision: It is the fraction of retrieved documents that are relevant to the query which is calculated using (6).

$$precision = \frac{|\{relevant\ documents\} \cap \{retrieved\ documents\}|}{|\{retrieved\ documents\}|} \dots (6)$$

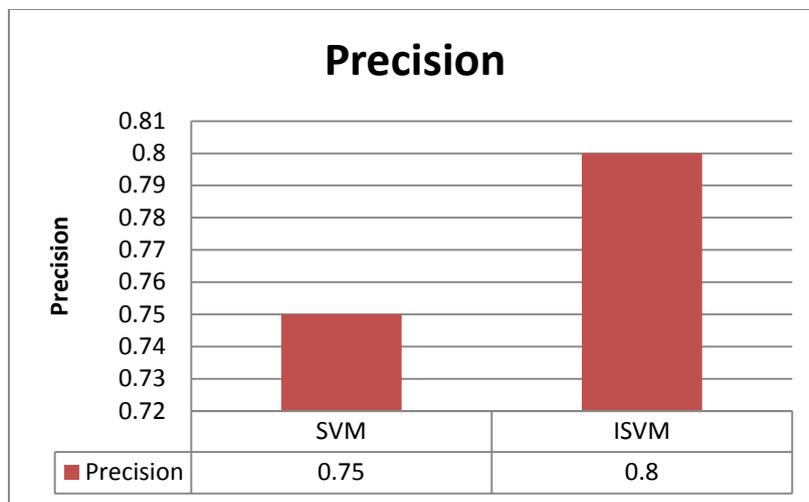


Figure 1. Comparison of Precision

Recall: It is the fraction of the documents that are relevant to the query that are successfully retrieved. The

recall value of the SVM and ISVM algorithm is 1 and is calculated using (7).

$$precision = \frac{|\{relevant\ documents\} \cap \{retrieved\ documents\}|}{|\{relevant\ documents\}|} \dots (7)$$

F1 – Measure: F1 score is a measure of a test's accuracy. It considers both the precision p and the recall

r of the test to compute the score. The F-1 measure is calculated using (8).

$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall} \dots (8)$$

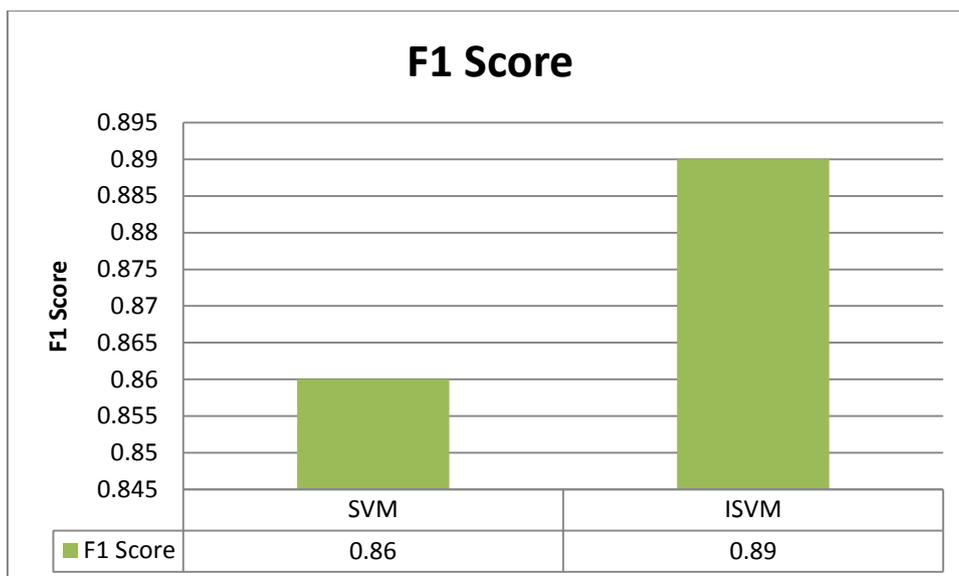


Figure2. Comparison of F-1 Score

Conclusions

This research work aims in design and development of an improved support vector machine classifier that makes use of adaboost genetic algorithmic approach towards web interaction mining. Appropriate parameters are chosen with the help of the genetic algorithm. In order to improve the performance of the SVM classifier Adaboost algorithm is employed. Performance metrics such as precision, recall and F-1 score are chosen. From the results it is evident that the proposed ISVM algorithm outperforms SVM classifier. As a future work, the proposed algorithm will be redesigned for mining mouse movements in the web content.

References

- [1] A.Vinupriya and S. Gomathi, "Web Page Personalization and link prediction using generalized inverted index and flame clustering," 2016 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, 2016, pp. 1-8.
- [2] A.Yang, X. Yu and Y. Liu, "Continuous KNN Join Processing for Real-Time Recommendation," 2014 IEEE International Conference on Data Mining, Shenzhen, 2014, pp. 640-649.
- [3] G. Dhivya, K. Deepika, J. Kavitha and V. N. Kumari, "Enriched content mining for web applications," Innovations in Information, Embedded and Communication Systems (ICIIECS), 2015 International Conference on, Coimbatore, 2015, pp. 1-5.
- [4] J. Jojo and N. Sugana, "User profile creation based on navigation pattern for modeling user behaviour with personalised search," Current Trends in Engineering and Technology (ICCTET), 2013 International Conference on, Coimbatore, 2013, pp. 371-374.
- [5] A.U. R. Khan, M. B. Khan and K. Mahmood, "Cloud service for assessment of news' Popularity in internet based on Google and Wikipedia indicators," Information Technology: Towards New Smart World (NSITNSW), 2015 5th National Symposium on, Riyadh, 2015, pp. 1-8.
- [6] M. A. Potey, D. A. Patel and P. K. Sinha, "A survey of query log processing techniques and evaluation of web query intent identification," Advance Computing Conference (IACC), 2013 IEEE 3rd International, Ghaziabad, 2013, pp. 1330-1335.
- [7] M. Nayrolles and A. Hamou-Lhadj, "BUMPER: A Tool for Coping with Natural Language Searches of Millions of Bugs and Fixes," 2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER), Suita, 2016, pp. 649-652.
- [8] Monika Dhandi, Rajesh Kumar Chakrawarti, "A Comprehensive Study of Web Usage Mining", 2016 Symposium on Colossal Data Analysis and Networking (CDAN), INDORE, India, 2016, Pages: 1 - 5.
- [9] T. Cheng, K. Chakrabarti, S. Chaudhuri, V. Narasayya and M. Syamala, "Data services for E-tailers leveraging web search engine assets," Data Engineering (ICDE), 2013 IEEE 29th International Conference on, Brisbane, QLD, 2013, pp. 1153-1164.
- [10] T. F. Lin and Y. P. Chi, "Application of Webpage Optimization for Clustering System on Search Engine V Google Study," Computer, Consumer and Control (IS3C), 2014

International Symposium on, Taichung, 2014, pp. 698-701.

- [11] T. Srivastava, P. Desikan, V. Kumar, "Web Mining – Concepts, Applications and Research Directions", Studies in Fuzziness and Soft Computing Foundations and Advances in Data Mining, Springer Berlin Heidelberg, 2013, pp 275-307.
- [12] X. Ye, Z. Qi, X. Song, X. He and D. Massey, "Generalized Learning of Neural Network Based Semantic Similarity Models and Its Application in Movie Search," 2015 IEEE International Conference on Data Mining Workshop (ICDMW), Atlantic City, NJ, 2015, pp. 86-93.
- [13] Y. C. Fan, Y. C. Chen, K. C. Tung, K. C. Wu and A. L. P. Chen, "A framework for enabling user preference profiling through Wi-Fi logs," 2016 IEEE 32nd International Conference on Data Engineering (ICDE), Helsinki, Finland, 2016, pp. 1550-1551.
- [14] Y. Kiyota, Y. Nirei, K. Shinoda, S. Kurihara and H. Suwa, "Mining User Experience through Crowdsourcing: A Property Search Behavior Corpus Derived from Microblogging Timelines," 2015 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT), Singapore, 2015, pp. 17-21.
- [15] Z. Liao, Y. Song, Y. Huang, L. w. He and Q. He, "Task Trail: An Effective Segmentation of User Search Behavior," in IEEE Transactions on Knowledge and Data Engineering, vol. 26, no. 12, pp. 3090-3102, Dec. 1 2014.
- [16] B. Kaviyarasu, Dr. A. V. Senthil Kumar, "WEB INTERACTION MINING USING IMPROVED EXTREME LEARNING MACHINE CLASSIFIER", International Journal for Research in Science Engineering and Technology, vol. 3, no.12, pp. 45-51, 2016.
- [17] B. Kaviyarasu, Dr. A. V. Senthil Kumar, "WEB INTERACTION MINING USING PENTA LAYERED ARTIFICIAL NEURAL NETWORK CLASSIFIER", International Journal of Computer Science Engineering and Technology, vol.3, no.1, pp. 64-70, 2017.

Appendix – 1

	Retrieved Documents	Precision for ISVM	Precision for SVM	Recall for ISVM	Recall for SVM	F1 Score for SVM	F1 Score for ISVM
P1	353	0.80	0.74	1	1	0.85	0.89
P2	350	0.81	0.77	1	1	0.87	0.90
P3	360	0.81	0.73	1	1	0.84	0.89
P4	365	0.79	0.73	1	1	0.84	0.88
P5	352	0.82	0.77	1	1	0.87	0.90
P6	358	0.79	0.76	1	1	0.87	0.88
P7	363	0.79	0.75	1	1	0.86	0.88
P8	350	0.83	0.77	1	1	0.87	0.91
P9	352	0.83	0.77	1	1	0.87	0.91
P10	365	0.82	0.75	1	1	0.86	0.90
P11	352	0.83	0.76	1	1	0.86	0.91
P12	357	0.80	0.75	1	1	0.86	0.89
P13	351	0.81	0.74	1	1	0.85	0.89
P14	360	0.80	0.76	1	1	0.87	0.89
P15	363	0.80	0.75	1	1	0.86	0.89
P16	364	0.79	0.72	1	1	0.84	0.88
P17	354	0.82	0.74	1	1	0.85	0.90
P18	353	0.80	0.75	1	1	0.86	0.89
P19	363	0.77	0.74	1	1	0.85	0.87
P20	352	0.80	0.76	1	1	0.87	0.89
P21	357	0.78	0.76	1	1	0.86	0.88
P22	359	0.83	0.76	1	1	0.86	0.91
P23	360	0.79	0.76	1	1	0.87	0.88
P24	358	0.79	0.73	1	1	0.84	0.88
P25	361	0.80	0.73	1	1	0.84	0.89
P26	350	0.84	0.75	1	1	0.86	0.91
P27	360	0.79	0.74	1	1	0.85	0.88
P28	358	0.78	0.75	1	1	0.86	0.88
P29	365	0.80	0.74	1	1	0.85	0.89
P30	367	0.80	0.74	1	1	0.85	0.89
P31	360	0.79	0.76	1	1	0.86	0.88