



WEB INTERACTION MINING USING IMPROVED EXTREME LEARNING MACHINE CLASSIFIER

¹ B. Kaviyarasu, ² Dr. A. V. Senthil Kumar

¹ Research Scholar, ² Director

^{1,2} PG and Research Department of Computer Computer Applications,

^{1,2} Hindusthan College of Arts and Science,

^{1,2} Coimbatore – 38.

ABSTRACT: Predicting the intention of internet users contains different applications in the areas such as e-commerce, entertainment in online, and several internet-based applications. The integral section of classifying the internet queries based on accessible features such as contextual information, keywords and their semantic relationships. This research article aims in proposing improved extreme learning machine classifier for 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 ELM [16]. Results proved that the proposed classifier attains better performance than that of the conventional ELM.

Keywords: [Web interaction mining, algorithm, extreme learning machine, classifier, precision, recall, F1-Score.]

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.

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.

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.

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.

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.

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:

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.

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.

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 grading 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 web patterns 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.

3. PROPOSED WORK

3.1. Improved Feature Collection using Grading Method

As far as improved approach is concerned, a class of web patterns are obtained as the basic unit or context for computing meaning scores for words. Meaning measure fundamentally portrays how expected a particular words' frequency is in a class of web patterns when compare to the other classes of web patterns. When there are unexpected words in the tweet are present then meaning measure results in

high meaning scores. In this portion it is analogous to the Multinomial Naive Bayes in which the all the web patterns in a class is merged into a single tweet and then the probabilities are estimated from this one large class web patterns. In improved meaning measure, parameter c_j represents web patterns that fit in to class j and S represents the complete training set. It is presumed that a feature w appears k times in the dataset S , and m times in the web patterns of class c_j . The length of dataset (i.e. training set) S and class c_j measured by the total term frequencies is L and B respectively. N is the rate of the length of the dataset and the class which calculate in (3). Based on these the number of false alarms (NFA) is defined in (4)

$$L = \sum_{d \in S} \sum_{w \in d} tf_w \quad \dots (1)$$

$$B = \sum_{d \in c_j} \sum_{w \in d} tf_w \quad \dots (2)$$

$$N = \frac{L}{B} \quad \dots (3)$$

$$NFA(w, c_j, S) = \binom{k}{m} \cdot \frac{1}{N^{m-1}} \quad \dots (4)$$

The meaning score of the word w in a class c_j is defined as:

$$meaning(w, c_j) = -\frac{1}{m} \log NFA(w, c_j, S) \quad \dots (5)$$

In order to simplify the calculations meaning formula can be re-written as:

$$meaning(w, c_j) = -\frac{1}{m} \log \binom{k}{m} - [(m-1) \log N] \quad \dots (6)$$

The larger the meaning score of a word w in a class c_j can be perceived as that the given word w is further meaningful, important or edifying for that class. It is firmly to mention that, the words with larger meaning scores be in contact to more meaningful, significant or informative words for that particular class. On the other hand,

for feature collection we need a way to combine these class-based scores into one and select top R features. In order to do this grading method is applied. Grading perform sort the features by using their meaning scores for each class. For example, the rank of the first element on each sorted list will be 1 and the last element will be the dictionary size. We use rank of the features in each class instead of their meaning scores. When combining these class based lists into a single feature list, for each feature we pick the highest rank among all classes as in (7).

$$score(w) = \max_{c_j \in C} (Rank(w, c_j)) \quad \dots (7)$$

3.2. Extreme Learning Machine Classifier

Once when the feature collection task is completed, IELM is employed for performing classification task. Given a set of N training samples (x_i, t_i) and $2L$ hidden neurons in total (that is, each of the two hidden layer has L hidden neurons) with the activation function $g(x)$. At first randomly initialize the connection cost matrix between the input layer and the first hidden layer W and the bias matrix of the first hidden layer B , and then calculate the cost matrix between the second hidden layer and the output layer.

$$g(W_H H + B_1) = H_1 \quad \dots (8)$$

where W_H denotes the cost matrix between the first hidden layer and the second hidden layer. It is presumed that the first and second hidden layers have the same number of neurons, and thus W_H is a square matrix. The notation H denotes the output between the first hidden layer with respect to all N training samples. The matrices B_1 and H_1 respectively represent the bias and the expected output of the second hidden layer.

The expected output of the second hidden layer can be calculated as

$$H_1 = TS^\dagger \quad \dots (9)$$

where \dagger is the MP generalized inverse of the matrix S . The calculating method of \dagger is the same as previously discussed for H^\dagger ,

namely $S^\dagger = (S^T S)^{-1} S^T$ if $S^T S$ is nonsingular, or alternatively

$S^\dagger = S^T (S^T S)^{-1}$ if SS^T is nonsingular. Consequently it is defined the augmented matrix $W_{HE} = [B_1 W_H]$, and calculate it as

$$W_{HE} = g^{-1}(H_1) H_E^\dagger \dots (10)$$

where H_E^\dagger is the generalized inverse of $H_E = [1 H]^T$, 1 denotes a one-column vector of size N whose elements are the scalar unit 1, where the notation $g(x)$ indicates the inverse of The calculation of H^\dagger proceeds in the fashion described before. The experiments conducted to test the performance of the ELM algorithm. In order to perform the classification task extensively used logistic sigmoid function $g(x) = 1/(1 + e^{-x})$ is used. The actual output of the second hidden layer is calculated as

$$H_2 = g(W_{HE} H_E) \dots (11)$$

and finally, the cost matrix S_{new} between the second hidden layer and the output layer is calculated as

$$S_{new} = H_2^\dagger T \dots (12)$$

where H_2^\dagger is the MP generalized inverse of H_2 , obtained using the approach discussed before. The ELM output after training can be expressed as

$$f(x) = H_2 S_{new} \dots (13)$$

Algorithm 1. IELM Algorithm

Input: N training samples $X = [x_1, x_2, \dots, x_N]^T$, $T = [t_1, t_2, \dots, t_N]^T$ and 2L hidden neurons in total with activation function $g(x)$

1: Randomly generate the connection cost matrix between the input layer and the first hidden layer W and the bias matrix of the first hidden layer B and for simplicity, W_{IE} is defined as $[B W]$ and similarly, X_E is defined as $[1 X]^T$.

2: Calculate $H = g(W_{IE} X_E)$;

3: Obtain cost matrix between the second hidden layer and the output layer $= H^\dagger T$

4: Calculate the expected output of the second hidden layer $H_1 = T S^\dagger$

5: Determine the parameters of the second hidden layer (connection cost matrix between the first and second hidden layer and the bias of the second hidden layer)

$$W_{HE} = g^{-1}(H_1) H_E^\dagger$$

6: Obtain the actual output of the second hidden layer $H_2 = g(W_{HE} H_E)$

7: Recalculate the cost matrix between the second hidden layer and the output layer

$$S_{new} = H_2^\dagger T$$

Output: The final output of IELM is $f(x) = \{[W_H g(W X + B) + B_1]\} S_{new}$

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 IELM 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.

Precision: It is the fraction of retrieved documents that are relevant to the query which is calculated using (14).

$$precision = \frac{|\{relevant\ documents\} \cap \{retrieved\ documents\}|}{|\{retrieved\ documents\}|}$$

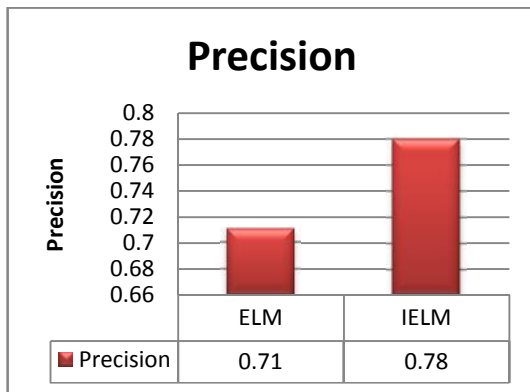


Figure - 1 Comparison of Precision

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 (15).

$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall} \dots (15)$$

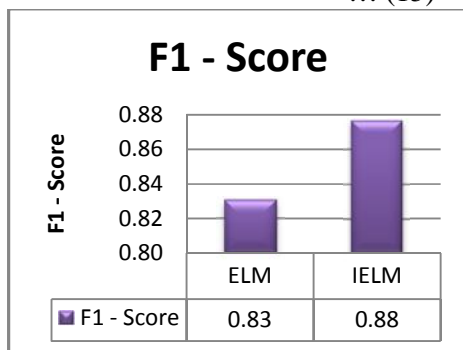


Figure - 2 Comparison of F-1 Score

CONCLUSIONS

This research work aims in design and development of improved extreme learning machine classifier in order to perform web interaction mining.

Modification is made in conventional extreme learning machine classifier with the help of improved feature collection using grading method strategy. Performance metrics such as precision, recall and F-1 score are chosen. From the results it is evident that the proposed IELM algorithm outperforms ELM classifier.

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]. Guang-Bin Huang, Qin-Yu Zhu, Chee-Kheong Siew, "Extreme learning machine: Theory and applications," Neurocomputing, vol. 70, pp. 489-501, 2006.