



International Journal of Innovative Research in Computer and Communication Engineering

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

Website: www.ijirce.com

Vol. 7, Issue 2, February 2019

Random Forest Algorithm in Inferring User Search Goals with Web Page Recommendation

Dr.A.Gnanabaskaran M.E., Ph.D, J. Adhipudrishwaranathan, R. Gowtham, S.M. Harish Kumar

Professor, Department of Computer Science and Engineering, K.S.Rangasamy College of Technology,

Tiruchengode, Tamilnadu, India

Department of Computer Science and Engineering, K.S.Rangasamy College of Technology,

Tiruchengode, Tamilnadu, India

ABSTRACT: An increasing number of databases have become web accessible through HTML form-based search interfaces. The data units returned from the underlying database are usually encoded into the result pages dynamically for human browsing. The inference and analysis of user search goals can be very useful in improving search engine relevance and user experience. Finally, we propose an analysis on the characteristics of random forest method, presents how to realize the self-adaptation ability with random forest method in similar situations, and verified the feasibility of the new method of using the actual data, and analysis and discussion of how to further research and improve the random forest method in big data environment.

I. INTRODUCTION

Owing to the enough accumulated data over the years in this sector, big data has gained many practical application scenarios. The whole range from vast amounts of information on the Internet to supermarket shopping bills contains significant commercial value. Rapid growth in the amount of data has overstepped the bearing capacity of traditional data analysis, which accelerates the urgency in development of big data analysis tools suitable for various application areas. Classifier technology is one focus of data mining research, and the famous classification algorithms covers association rules, Baye, decision trees, neural networks, rule learning, K-means, genetic algorithms, rough sets, fuzzy logic and other directions. Random forest method that is not subject to memory limitations and featured with rapid processing speed and good parallel scalability, is an excellent classification tool to handle massive data and a typical decision tree classification algorithm.

II. LITERATURE REVIEW

2.1 COMBINING COLLABORATIVE FILTERING

A novel approach that dynamically recommends Web services that fit users' interests. Our approach is a hybrid one in the sense that it combines collaborative filtering and content-based recommendation. In particular, our approach considers simultaneously both rating data and content data of Web services using a three-way aspect model. Unobservable user preferences are represented by introducing a set of latent variables, which is statistically estimated. To verify the proposed approach, we conduct experiments using 3,693 real-world Web services (e.g., which keywords should be used, what values should be set for a QoS attribute). Another problem is that Web services that do not satisfy user's searching query are completely excluded from the recommendation list.



International Journal of Innovative Research in Computer and Communication Engineering

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

Website: www.ijirce.com

Vol. 7, Issue 2, February 2019

2.2 REQUIREMENTS IN SERVICE RECOMMENDATION

There are three main requirements in order to conduct an effective service recommendation task:

- High recommendation accuracy. A good recommendation system should recommend more favorite Web services and fewer disliked ones, particularly in the situations where available information might be not sufficient (e.g., missing QoS of some services).
- Recommendation diversity. Recommending services that are well-known to a user is often found unsatisfactory or meaningless. If the recommended services are unfamiliar to a user, the chances of finding new Web services that match the user's requirements would increase.
- Overcoming the cold-start problem. Solving this problem not only enables users to find newly-deployed Web services, but also enhances the recommendation diversity.

2.3 WEB USAGE MINING

The motive of mining is to find users' access models automatically and quickly from the vast Web log data, such as frequent access paths, frequent access page groups and user clustering.

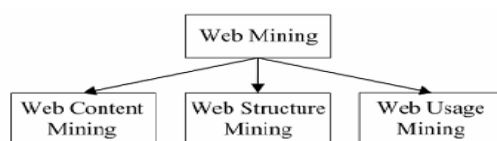


Fig no: 2.1

Through web usage mining, the server log, registration information and other relative information left by user access can be mined with the user access mode which will provide foundation for decision making of organizations. This article provides a survey and analysis of current Web usage mining systems and technologies.

2.4 PRICE DYNAMIC TREND ANALYSIS

Commodities price of others e-supermarkets or online shopping systems is the most important data for the shopkeepers of shop online. This requirement becomes actuality because of the Web mining developing very fast. The Web mining algorithm from extracting directory tree of different Website, the commodities name on the Webpage and commodities price based on participle are described in detailed.

2.5 CONCLUSIONS FROM THE LITERATURE REVIEW

Web services recommendation and selection is a fundamental issue in service-oriented computing. Existing Web services discovery and recommendation approaches focus on either perishing UDDI registries, or keyword-dominant, QoS-based Web service search engines.

A new web recommendation system based on the proposed Weighted Association Rule (WAR) model. We extend the association rule mining by assigning a significant weight to the pages based on time spent by each user on each page and visiting frequency of each page. The proposed weighting measure can be used to judge the importance of a page to a user, and try to give more consideration to pages which are more useful to the user.

System performance was evaluated under different settings and in comparison, with traditional Association Rule based model. The experimental results show that our method is better in precision and coverage rates than the conventional association rule-based recommendation.

Web usage mining model is a kind of mining to server logs. Web Usage Mining plays an important role in realizing enhancing the usability of the website design, the improvement of customers' relations and improving the requirement of system performance and so on. Web usage mining provides the support for the web site design, providing personalization server and other business making decision, etc.

The market of mobile phone change very fast. So, the price dynamic trend analysis is pay more attention by the shopkeepers. All data mining is as same importance as individual analysts from enterprise standpoint. We believe



International Journal of Innovative Research in Computer and Communication Engineering

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

Website: www.ijirce.com

Vol. 7, Issue 2, February 2019

that in the future, dynamic trend analysis system using the data mining technology can have greater development. Although our work achievement the very useful effect that touched the shopkeepers mind. But future interesting work is waiting on our hard work such as the decision support system.

III. PROPOSED SYSTEM

3.1 DATA VALIDATION

Data sets used in the test are originated from real customer data of financial industry. The data amount accounts for 200,000 pieces with around 10,000 pieces of data from each quarter, which was sampled from a larger original 5-year data set. The data set contains a target category and 16 feature attributes, which includes both the continuous numerical attributes and discrete attributes. Random forest kit in R language version was used in the test.

3.2 RESTRUCTURING WEB SEARCH RESULTS

Since search engines always return millions of search results, it is necessary to organize them to make it easier for users to find out what they want. Restructuring web search results is an application of inferring user search goals. We will introduce how to restructure web search results by inferred user search goals at first. Then, the evaluation based on restructuring web search results will be described. The inferred user search goals are represented by the vectors and the feature representation of each URL in the search results can be computed result. Then, we can categorize each URL into a cluster centred by the inferred search goals. In this paper, we perform categorization by choosing the smallest distance between the URL vector and user-search-goal vectors. By this way, the search results can be restructured according to the inferred user search goals.

3.2.1 Presentation Style (PS)

This feature describes how a data unit is displayed on a webpage. It consists of six style features: font face, font size, font color, font weight, text decoration (underline, strike, etc.), and whether it is italic. Data units of the same concept in different SRRs are usually displayed in the same style.

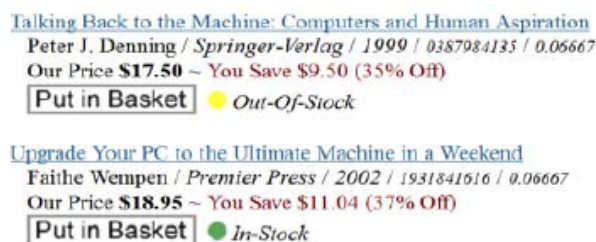


Fig no: 3.1

3.3 RANDOM FORESTS

Random forests are an idea of the general technique of random decision forests that are an ensemble learning technique for classification, regression and other tasks, that control by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests accurate for decision trees' habit of overfitting to their training set.

A Random Forest is a classifier consisting of collection of tree-structured classifiers where independent random vectors are distributed identically and each tree cast a unit vote for the most popular class at input x.



International Journal of Innovative Research in Computer and Communication Engineering

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

Website: www.ijirce.com

Vol. 7, Issue 2, February 2019

A random vector is generated which is independent of the past random vectors with same distribution and a tree is generated by using the training test [Brieman (2001)]. For random forests, an upper bound is derived to obtain the generalization error in terms of two parameters that are given below:

- The accuracy of individual classifiers
- The dependency between the individual classifiers
- The generalization of error for random forest includes two segments. These segments are defined below:
 - The strength of the individual classifiers in the forest.
 - The correlation between them in terms of raw margin function

IV. EXPERIMENTAL SETUP AND PROCEDURE

4.1 EXPERIMENTAL DESIGN METHODOLOGY

In the Expt, we use Weka data mining tool to conduct the experiment. We compared the classification performance of the Decision tree i.e. Random Forest and Random Tree models employing attribute selection filter. The Breast Cancer DataSet from UCI repository is used in this expt. We use 10-fold cross validation as the test mode to record classification accuracy. This approach is suitable to avoid biased results and provide robustness to the classification. Also, the parameters of a classification algorithm are chosen to their default values. The following steps have been applied to generate experimental data in order to draw inference:

- Find classification performance of the classifiers in the dataset.
- Find classification performance using Attribute selection filter.

4.2 ACCURACY AND PRUNING OF RANDOM TREES

In order to meet the needs in the big data environment, improved algorithm should have the following characteristics:

1. It can quickly generate a classifier on a given data training set;
2. The resulting classifier can quickly classify new streaming data;
3. An algorithm should be of adaptability so as to respond to the changes in data modes and guarantee its accuracy;
4. It should limit the scale, namely, the number of trees, which will, on the one hand, ensure the efficiency of the algorithm, while on the other hand, will also guarantee its accuracy.

4.3 APPLICATIONS OF RANDOM FOREST

Online Learning and Tracking: Incremental Extremely Random Forest algorithm was introduced in 2009 [Wang et al. (2009)]. This algorithm is used for online learning classification and in video tracking problems. It deals with the small steaming labeled data. Whenever the examples are arrived at leaf node, Gini Index is calculated to determine the splitting node of the tree. This technique also increases the capability as compared with other co-training framework. The examples are stored in the memory so that they can be reuse again to perform the split test with small number and avoid the calculation of Hoeffding bounds for large number of attributes.

4.3.1 Semi-Supervised Random Forest:

Random Forest does not require various binary classifiers for the evaluation of multi-class problem. But, it has some limitation that it needs large amount of labeled data to move to the best level of performance. Therefore, to overcome this problem Semi-Supervised Random Forest algorithm was proposed. This algorithm processes both labeled and unlabeled training data. It is based on Deterministic Annealing. Using this approach, unlabeled data consist of labels that can effectively treated as additional optimization variables. The other advantage of using this algorithm is that by estimating the out-bag-error monitoring of unlabeled data is performed.



International Journal of Innovative Research in Computer and Communication Engineering

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

Website: www.ijirce.com

Vol. 7, Issue 2, February 2019

4.4 RANDOM FOREST

This method is proposed for generating the base classifiers ensembles on the basis of feature extraction. The training set for base classifier is created by splitting the feature set into K subsets and Principal Component Analysis (PCA) is applied to each subset of features. All principal components are retained so that the variability information in the data could be preserved. Therefore, K axis of rotation is performed to develop the new features for base classifiers. This approach is used to maintain the accuracy and diversity within the ensemble.

V. RESULT AND DISCUSSION

5.1 EXPERIMENTAL RESULT

It is performed several runs in Weka tool and gathered the data for the inference. Table-1 summarizes the classification accuracy in percentage of all the classifiers across the datasets with original features while

Table no: 5.1

the classification performance after a kNN algorithm-based feature selection. It is observed in the tabulated data that the performance of the Random Forest and Random Tree classifiers with attribute selection filter.

VI. CONCLUSION

In this analytical study, we consider 2 popular decision tree classifiers with publicly available microarray datasets and one attribute selection filter for classification. Following a methodical approach, we gathered experimental data using Weka, a popular data mining tool. Both the techniques compete with each other to provide classification accuracy across the dataset So, we lean towards statistical inference and find that attribute selection filter based kNN classification selection can produce slightly better classification accuracy than that of Random Forest.

Classifier's Name	Accuracy	Time Taken to Build the Model (Seconds)	ROC Area	F-Measure	Root Mean Squared Error
Random Forest	70.24	0.21	0.543	0.53	0.34
kNN	68.35	0.25	0.432	0.43	0.46

ACKNOWLEDGEMENT

“We acknowledge DST-File No.368. DST -FIST (SR/FIST/College -235/2014 dated 21-11-2014) for financial support and DBT-STAR-College-Scheme-ref.no: BT/HRD/11/09/2018 for providing infrastructure support.”

REFERENCES

- [1] Adeniyi, Z. Wei, Y. Yongquan (2015), “Automated web usage data mining and recommendation system using K-Nearest Neighbor (KNN) classification method”, Saudi Computer Society, King Saud University, Applied Computing and Informatics, Production and hosting by Elsevier B.V, pp.111-119.
- [2] Bussa V. R. R. Nagarjuna, Akula Ratna babu and Miriyala Markandeyulu, A. S. K. Ratnam (2016), “Web Mining: Methodologies, Algorithms and Applications”, International Journal of Soft Computing and Engineering (IJSCE), ISSN: 2231-2307, Volume-2, Issue-3.
- [3] Haidong Zhong, Shaozhong Zhang, Yanling Wang, Shifeng Weng and Yonggang Shu (2017), “Mining Users' Similarity from Moving Trajectories for Mobile Ecommerce Recommendation”, International Journal of Hybrid Information Technology Vol.7, No.4, pp.309-320



ISSN(Online): 2320-9801
ISSN (Print) : 2320-9798

International Journal of Innovative Research in Computer and Communication Engineering

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

Website: www.ijirce.com

Vol. 7, Issue 2, February 2019

- [4]Kavita Sharma, Gulshan Shrivastava, Vikas Kumar (2017), "Web Mining: Today and Tomorrow", 2017 6th International Conference on Electronics Computer Technology , 2017 IEEE,pp.444-458.
- [5]Lina Yao and Quan Z. Sheng, Aviv Segev, Jian Yu (2015) "Recommending WebServices via Combining Collaborative Filtering with Content-based Features", 2015 IEEE 24th International Conference on Web Services, 2015 IEEE,pp.236-242.
- [6]Neha Sharma & Pawan Makhija (2015), "Web usage Mining: A Novel Approach for Web user Session Construction", Global Journal of Computer Science and Technology: E Network, Web & Security Volume 15 Issue 3,pp.329-341.
- [7]Petrović, P. Perković and I. Štajduhar (2015), "A Profile- and Community-Driven Book Recommender System", 2015 38th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO),pp.240-253.
- [8]Quanyin Zhu, Hong Zhou, Yunyang Yan, Jin Qian and Pei Zhou (2016), "Commodities Price Dynamic Trend Analysis Based on Web Mining", 2016 Third International Conference on Multimedia Information Networking and Security. IEEE,pp.341-352.
- [9]Rana Forsati, Mohammad Reza Meybodi, Afsaneh Rahbar (2015), "An Efficient Algorithm for Web Recommendation Systems", 2015 IEEE/ACS International Conference on Computer Systems and Applications,pp.356-373.
- [10]Renuka Mahajan, J. S. Sodhi, Vishal Mahajan (2016), "Web Usage Mining for Building an Adaptive e-Learning Site:A Case Study", International Journal of e-Education, e-Business, e-Management and e-Learning, Manuscript,pp.800-818.
- [11]Ricardo Terra, Marco Tulio Valente, Krzysztof Czarnecki and Roberto S. Bigonha (2015), "A recommendation system for repairing violations detected by static architecture conformance checking", SOFTWARE – PRACTICE AND EXPERIENCE, Softw. Pract. Exper. 2015; pp.342-354.