

Random forest application on cognitive level classification of E-learning content

Benny Thomas, Chandra J.

Department of Computer Science, CHRIST (Deemed to be University), India

Article Info

Article history:

Received Dec 19, 2019

Revised Mar 3, 2020

Accepted Mar 14, 2020

Keywords:

Blooms taxonomy

Difficulty level

E-learning

Machine learning

Random forest classifier

ABSTRACT

The e-learning is the primary method of learning for most learners after the regular academics studies. The knowledge delivery through E-learning technologies increased exponentially over the years because of the advancement in internet and e-learning technologies. Knowledge delivery to some people would never have been possible without the e-learning technologies. Most of the working professional do focused studies for carrier advancement, promotion or to improve the domain knowledge. These learner can find many free e-learning web sites from the internet easily in the domain of interest. However it is quite difficult to find the best e-learning content suitable for their learning based on their domain knowledge level. User spent most of the time figuring out the right content from a plethora of available content and end up learning nothing. An intelligent framework using machine learning algorithms with random forest Classifier is proposed to address this issue, which classifies the e-learning content based on its difficulty levels and provide the learner the best content suitable based on the knowledge level. The frame work is trained with the data set collected from multiple popular e-learning web sites. The model is tested with real time e-learning web sites links and found that the e-contents in the web sites are recommended to the user based on its difficulty levels as beginner level, intermediate level and advanced level.

Copyright © 2020 Institute of Advanced Engineering and Science.
All rights reserved.

Corresponding Author:

Benny Thomas,

Department of Computer Science,

CHRIST (Deemed to be University), India.

Email: benny.thomas@res.christuniversity.in

1. INTRODUCTION

E-learning is a popular learning method with the help of internet and other e-learning technologies. It bridges the geographical gap between the learner and teacher. E-learning become popular with the advancement in e-learning technologies and the availability of world class e-learning web sites. Currently it is the primary method of learning for most of the working professional and entrepreneurs. E-learning gives us the choice and flexibility to learn from anywhere and at any time. Because of its wider usage and potentiality, the e-learning web sites increased exponentially over the years. It is easy for any learners to find multiple e-learning web sites needed for their domain. However because of the availability of many web sites, the user most of the time get overwhelmed with the magnitude of content availability and find it difficult to understand and choose the right learning content. User spent most of the time trying to figure out the content to be chosen and end up learning nothing significant to improve the knowledge. This situation can be managed by providing intelligent content recommendations based on the domain knowledge level of the user which helps to find the right learning content. Different approaches were used to address this issue. Some of the methods used are recommended systems, good learners rating, association rule mining, learner grouping, item set mining etc.

Recommended systems are the most popular method used to suggest eLearning content to the user. Ontology based and fuzzy rule based systems are popular content based recommendation system used in eLearning [1]. Commonly used recommended systems are collaborative filtering based systems, hybrid filtering systems and content based filtering systems [2]. Content based filtering uses previously studied learners' rating. It can bring additional similar learning material based on good learner's recommendations. Collaborative filtering approach brings similar learners materials to recommend to the user. A mix of content filtering and collaborative filtering methods are used in Hybrid filtering [3]. These system does not classify contents based on the knowledge levels of the user. An intelligent e-learning framework is proposed to identify the cognitive level of e-content available in various e-learning web sites. The data set needed for this is collected through web scrapping of popular e-learning web sites. The web pages are downloaded and scraped as text files to collect the dataset. Random forest classification algorithm is used to classify the text based on difficulty levels.

2. LITERATURE REVIEW

Numerous studies have conducted to recommend the best e-learning content to the learner through various text classification methods. Atorn Nuntiyagul et. al., categorize questions kept in an item bank and reestablish it based on difficulty levels by using patterned keywords and phrase with support vector machine algorithms. Inputs are given as mathematical questions in text format. The selected keyword patterns and weights are combined and applied to vector space model as a feature matrix. The methodology takes the advantage of text classification techniques in machine learning and information retrieval [4]. G. Desai et.al., proposed that the Naive Bayes approach is one of the best method for text classification and it yields good results. It uses statistical and supervised learning methods. The methodology used is random sampling of the labeled categories of text. It explains the automatic classification of research topics based on the result analysis and textual content [5]. Sankar Perumal et. al., Proposed a new content based recommender system to provide appropriate contents by filtering the frequent item patterns obtaining through pattern mining and then ordering the final contents using fuzzy logic into different levels. It has higher efficiency and accuracy compared to the other parallel methods [1]. Tarus and Niu has conducted study to understand the different ontology based e-learning recommended systems which demonstrated ontology for representing knowledge that brought enhancement in the recommendations.

The methodology used is survey of the e-learning recommended systems and compared and analyzed the results of various ontology based recommended systems. The ontology-based systems uses hybrid recommendation systems and knowledge-based techniques and other approaches such as content-based filtering, collaborative filtering, fuzzy-based context-aware and trust-based techniques [6]. Shuai et. al., provide extensive review of the researches in deep learning recommended systems and devised a nomenclature for deep learning based recommended models. Also proposed an organization scheme for classifying and clustering existing works with advantages and disadvantages of using deep learning based recommended systems [7]. Dragi Kocev et.al has done an experimental evaluation of Multi label learning methods by selecting competent methods and using data set from different domains and based on the previous usage statistics. Multi label learning is learning by examples. The results are experimentally analyzed and found that the best performing methods are random forests of predictive clustering trees and hierarchy of multi-label classifier [8]. Salem et. al., suggested an automatic meta-learning recommendation model which extract learning contents from knowledge units as teaching notes using Natural Language Processing techniques. NLP helps to extract the verbs based on the cognitive levels. It describes the use of Blooms taxonomy for computer science domain [9]. Amal et. al., provided a graph based Tringularity system for knowledge unit identification and classification using Bloom's taxonomy levels. When a knowledge unit is given the system finds out the hidden association in the knowledge unit using Tringularity graph. This model can be used to give the new sequential ordering of knowledge units in a textbook [10]. Colin et. al., suggest that the existing learning taxonomies are not useful for practical oriented subjects such as computer programming and devised a new taxonomy for programming.

A two dimensional approach of Bloom's taxonomy called Matrix Taxonomy is proposed to provide more practical framework to assess the learner. The two dimensions are divided as producing and interpreting which removes the strict ordering and retain the Blooms concepts [11]. Othman et. al., uses Naïve base classifier method to identify Bloom's taxonomy levels in text based on rule set in training data. The concept can be used to order a text book using Blooms Taxonomy cognitive levels and give a new sequential ordering to the book. The results shows that the several parts of the book which are described as intermediate become advanced and some advanced topics become intermediate [12]. Yahyaa et. al., analyzed the difficulty levels of questions asked in class during teaching, using Bloom Taxonomy verbs. The classification is done based on its difficulty levels using Bloom levels with K-NN, Naïve Bayes and SVM algorithms using term

frequency as the selection criteria. It uses four stages such as text representation, selection, classifier construction and testing the classifier [13]. Yamaguchi et. al., has developed personalized English teaching material for beginner level learners which identifies the cognitive level of a text document content. The difficulty level is identified by the personalized vocabulary of the learner. The learning materials are recommended based on the vocabulary knowledge of the students. The number of unknown words estimate the difficulty level of the content. The research has relevance today because of the exponential growth of e-learning content availability in the web [14]. Yahya and Osman identified difficulty levels using support vector machine algorithms for question paper categorization based on difficulty levels. Already categorized questions are gathered and support vector machine classifier is applied on these questions for classification [15]. Yang et. al., proposed a method to apply teaching and learning to learners from diverse background effectively using item response theory. The approach uses web document classification by introducing the concept of knowledge unit obtained from the subject. The questions are developed to measure and analyze them using item response theory. The learners are assessed through self-assessment based on the subject and user's knowledge level the sub topics of the subject. Based on this, the students are divided in to different groups to learn from the web [16].

3. RESEARCH METHOD

Dataset is collected from popular e-learning websites through web scraping. The contents were divided in to three different difficulty levels namely beginner, intermediate and advanced. The difficulty level is identified before downloading the e-content by reading each page of the e-content in the website. Each page is parsed to produce the text file for further processing. The dataset with varying dimensions are created to check the robustness of the algorithm at different dimensionalities. The data set is divided in to three sizes, 600 files, 2100 files and 4000 files to check the performance as shown in Table 1.

Table 1. Description of the datasets used to test the model

Total Files in each dataset	Beginner files	Intermediate files	Advanced files	Training files	Testing files
600	200	200	200	150	50
2100	700	700	700	525	175
4000	1800	900	1300	3000	1000

3.1. Block diagram

Figure 1 shows the work flow of the proposed framework. Figure 2 show generic architecture of the proposed framework. The content obtained through web crawling is loaded in to the framework. After preprocessing and data reduction, bloom taxonomy verbs and its synonyms were added to the selected feature set to improve the accuracy of the classifier. The random forest classifier is used to train and test the model. The model is validated using the e-learning web pages from different e-learning web sites. Finally, the eLearning content is recommended to the learner based the domain knowledge of the learner.

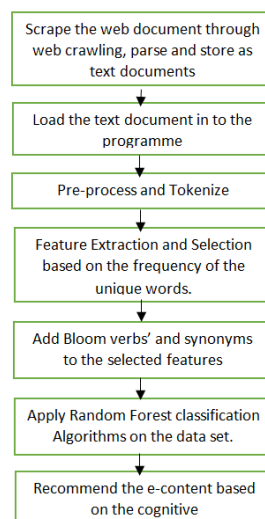


Figure 1. Work flow diagram of the classification

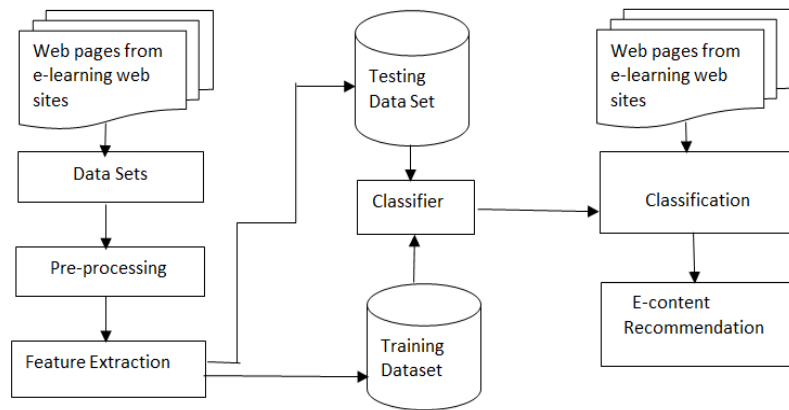


Figure 2. Generic architecture of the proposed framework

3.2. Preprocessing

Pre-processing minimize noise in the dataset by removing the unrelated and irrelevant data from the document. It removes all the tab spaces, punctuations, one letter words, two letter words, numeric strings and stop words. All the letters in the documents are converted to lower case and all words with two letters are removed from the document. The document is divided into flat array of words after removing metadata elements.

3.3. Feature extraction

The size of the document obtained after preprocessing is further reduced through feature extraction methods. The number and frequency of unique words in the feature set is obtained and all the duplicate words are removed to reduce the size of the data. The feature size is further reduced using feature selection methods by removing less important features and taking only relevant percentage of the total feature set. This percentage is calculated using N-Fold cross validation as follows. The N-Fold cross validation is done by taking 15 to 75 percentage of the total feature set.

$$n = (\text{Total No.of Features} * \text{percentage}) / 100 \quad (1)$$

Feature = words [0: n], Where percentage is an integer value less than 100 .The different percentage value of the total feature set is calculated to find the best percentage of feature selection. The data size percentage between 15 and 25 is giving the best accuracy and it reduced the data size further by 75 to 85 percentage. The documents after preprocessing is divided into training and testing. 75 percentage of the data is used for training and remaining 25 is used for testing. The best training and testing percentage is calculated using N-Fold validation with varying training and testing data size percentage such as 65-35, 70-30, 75-25 ,80-20 and 85-15.

3.4. Bloom's taxonomy

Bloom taxonomy helps to divide a learning content in to different cognitive levels based on the difficulty level of the learning content [17]. Bloom Taxonomy is hierarchical, learning at the higher level depends on having attained sufficient knowledge in the lower level. Before understanding a concept, one must remember it. In order to apply a concept, one must first understand it. The concept should be evaluated before analyzing it. To create a concept/product, it must be thoroughly evaluated. Different levels are:

- Remembering-lowest level
- Understanding-lowest level
- Applying-Intermediate level
- Analyzing-Intermediate level
- Evaluating-Advanced level
- Creating-Advanced level

Bloom's taxonomy helps to identify the difficulty level of the content with the help of different verbs used in each context [18]. The Bloom's taxonomy verbs are added to the feature set obtained after data reduction. The synonyms of each of these verbs are extracted using WordNet from NLTK tool kit. The Bloom's synonyms are also added to the feature set. This helps to improve the performance of the classification algorithm and to predict the classes in which the document belong [19].

3.5. Machine learning

The machine learning algorithm is used to train the machine with the model created with the help of different e-learning documents. In Machine learning, training is done using the fit method, the system is able to automatically understand difficulty levels of contents passed through using classification algorithms [20, 21]. The random forest machine learning algorithms is used for training the model after comparing many other algorithms on the data set [22]. The accuracy of the random forest classifier is found to be the highest in comparison with other algorithms.

3.6. Random forest classifier

Random forest belongs to ensemble family of classifiers as shown in Figure 3. It consist of number of random decision trees. It recursively generate many binary decision trees from a bagged random set of data. Each tree is independent from the other and it is constructed from a bootstrap sample of training data [23]. The trees in the random forest is constructed by an addition of randomness and therefore algorithm is named as random forest [24]. The data which is left without any decision tree is called the out of bag data and it is used for testing the performance of the decision tree. It uses feature selection method and feature ranking based on the importance of a feature in the overall dataset. The accuracy and robustness of the RF is very high. Its main advantages are its power to handle over-fitting and missing data [25] as well as its capacity to handle large datasets without eliminating the variables in the feature selection and resilience to high dimensionality data, insensitivity to noise, and resistance to over fitting. The trees in forest gives out a prediction result and the class which has more voting is taken as the model prediction value.

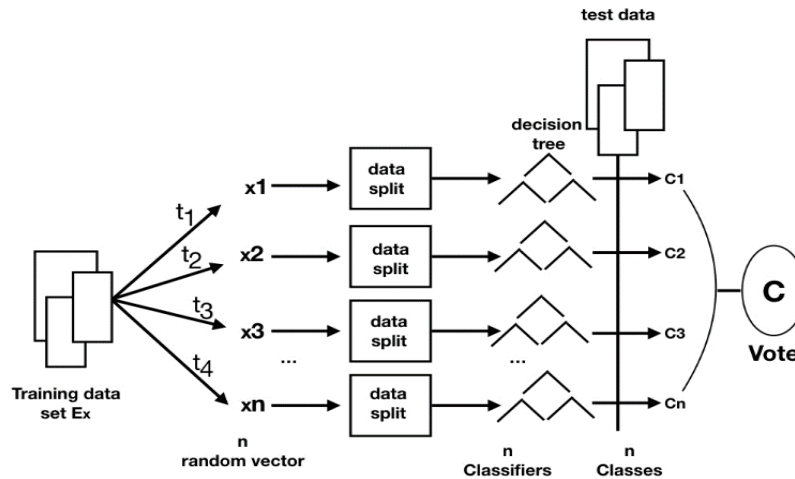


Figure 3. Generic architecture of random forest classifier

Random forest consist of a collection of randomized regression trees

$$R_{N(X, K_M, D_N), M \geq 1} \tag{2}$$

Where K_1, k_2 , are independently distributed random variable [24]. The Random trees are combined to form aggregate estimate of regression.

$$R_N(X, D_N) = E_K[R_N(X, K, D_N)] \tag{3}$$

Where E_K is the expectation with respect to the random parameter on X and the data set D_N . the RF get the most important features from the feature set to construct the decision trees. RF chooses the best fit using the gini index.

$$gini(ar) = 1 - \sum [P_j]^2 \tag{4}$$

$$gini_{split} = \sum_{ar=1}^n \frac{n_{ar}}{n} gini(ar) \tag{5}$$

Where P_j is the frequency of the feature set (ar) at class J, n_{ar} is the number of randomly selected training records. Each tree in the random forest output prediction value and the class with more vote is taken as the prediction value. The six Blooms' taxonomy classes are condensed into three classes with three difficulty levels. The Random forest classifier performance is improved with the addition of Blooms verbs and its synonyms into the feature set as RF internally use Gini index to select the best features and Gini index look for the most frequently occurring features from the feature set [25]. As the Blooms features are abundant in e-learning content, it improves the performance of the classifier.

4. IMPLEMENTATION

The algorithm is made to run with different data size to understand the best performance and to identify the optimum data size percentage required. This is accomplished with the help of N-Cap cross validation. The optimum data performance is obtained at 8 percentage to 15 percentage of data size. Table 2 and Figure 4 shows the accuracy, training time and testing time of Random forest classifier on the data set with different dimensionality of the data set. 8 to 14 percentage of the total data is used for testing the model after N-fold cross validation.

Table 2. Random forest classifier run with 600 files with 8 to 14% feature selection

Dimensionality in Percentage	Time needed to Train	Time needed to Test	Accuracy
8	0,171s	0,017s	0,986
9	0,171s	0,017s	0,990
10	0,155s	0,017s	0,986
11	0,1712s	0,017s	0,986
12	0,171s	0,017s	0,986
13	0,186s	0,017s	0,990
14	0,187s	0,017s	0,918

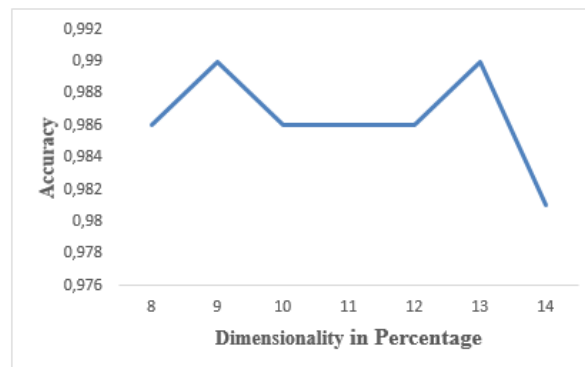


Figure 4. Result of random forest classifier run with different data size

The algorithm is made to run again after applying Bloom's Taxonomy verbs and all its synonyms into the feature set using WordNet from NLTK tool kit. The accuracy of the algorithm is improved further with the addition of Bloom's Taxonomy verbs and synonyms. The accuracy, training time and testing time of the classifier with bloom taxonomy verbs in the feature set is shown in the Table 3 and Figure 5. These values increased with the addition of Blooms verbs in the feature set.

Table 3. Random forest classifier run using 600 files with 8 to 14% feature selection using bloom's taxonomy verbs and synonyms

Dimensionality in Percentage	Time needed to Train	Time needed to Test	Accuracy
8	0,282s	0,000s	0,983
9	0,264s	0,017s	0,992
10	0,266s	0,017s	0,992
11	0,281s	0,017s	0,992
12	0,295s	0,017s	0,992
13	0,344s	0,017s	0,992
14	0,297s	0,017s	0,992

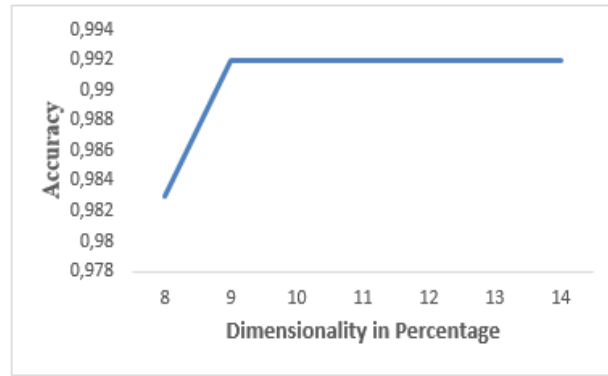


Figure 5. Random forest classifier run with different data size with bloom’s taxonomy verbs and synonyms

4.1. Recommendation

The trained model is saved and validated on real data using many web sites links from popular e-learning web sites, The web pages from e-learning web sites are passed to the Framework. The classifier divides the e-learning content from multiple websites in to three different difficulty levels and recommended to the user as shown in Figure 6.

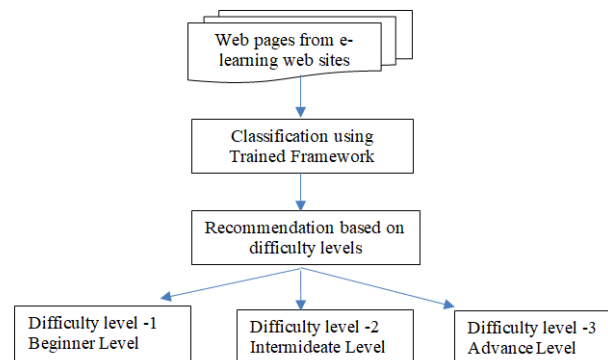


Figure 6. Recommendation of e-learning content using trained Framework based on difficulty levels

5. RESULTS AND DISCUSSION

Random forest classifier is used for developing the intelligent framework to assess the difficulty level of the e-content in multiple websites. The framework can be used to recommend right learning content to a learner based on the domain knowledge level. The model is trained with data sets of varying size using different percentage of feature selection. The model is test with N-Fold cross validation to obtain the best accuracy. The trained model is tested with different percentage of test data set to ensure the accuracy before saving the final model. The saved model is tested with real time web site links and found that the values obtained are correct. The output is validated using the link and check the actual data in the web sites and found that the difficulty levels are predicted with accuracy. The output of the classifier show in Figure 7.

URL	Subject	Topic	Difficulty Level
https://www.w3schools.com/java/java_booleans.asp	java	Booleans	beginner
https://www.w3schools.com/java/java_arraylist.asp	java	ArrayList	intermediate
https://www.w3schools.com/java/java_break.asp	java	Break and Continue	beginner
https://www.w3schools.com/java/java_comments.asp	java	Java Comments	beginner
https://www.w3schools.com/java/java_constructors.asp	java	Constructors	beginner
https://www.w3schools.com/java/java_data_types.asp	java	Data Types	beginner
https://www.w3schools.com/java/java_date.asp	java	Date and Time	intermediate
https://www.javatpoint.com/InetAddress-class	java	Java InetAddress	advanced

Figure 7. The output of the classifier using e-learning content from websites with different difficulty levels

6. CONCLUSION

The number of e-learning web sites and the users of these web sites increased dramatically over the years because of the advancement in internet and e-learning web sites. The e-content available is so much that the user often found it hard to figure out the right content. This leads to the necessity of an intelligent solution to recommend the right e-learning content to the user. The proposed framework uses a Random forest Classifier to develop a trained model. The model is used to classify the e-learning content in web sites based on its difficulty levels. The framework is useful to find the exact learning content from a large collection of content from the web site based on the knowledge level of the user.

REFERENCES

- [1] S. Pariserm Perumal, G. Sannasi, and K. Arputharaj, "An intelligent fuzzy rule-based e-learning recommendation system for dynamic user interests," *J. Supercomput.*, vol. 75, pp. 5145-5160, 2019.
- [2] G. Geetha, M. Safa, C. Fancy, and D. Saranya, "A Hybrid Approach using Collaborative filtering and Content based Filtering for Recommender System," *J. Phys. Conf. Ser.*, vol. 1000, no. 1, 2018.
- [3] K. I. Ghauth and N. A. Abdullah, "Learning materials recommendation using good learners' ratings and content-based filtering," *Educ. Technol. Res. Dev.*, vol. 58, no. 6, pp. 711-727, 2010.
- [4] A. Nuntiyagul, K. Naruedomkul, N. Cercone, and D. Wongsawang, "Adaptable learning assistant for item bank management," *Comput. Educ.*, vol. 50, no. 1, pp. 357-370, 2008.
- [5] M. K. M. S. D. H. P. G. Desai, and N. Chiplunkar, "Text Mining Approach to Classify Technical Research Documents using Naïve Bayes," vol. 4, no. 7, pp. 386-391, 2015.
- [6] J. K. Tarus, Z. Niu, and G. Mustafa, "Knowledge-based recommendation: a review of ontology-based recommender systems for e-learning," *Artif. Intell. Rev.*, vol. 50, no. 1, pp. 21-48, 2018.
- [7] S. Zhang, L. Yao, A. Sun, and Y. I. Tay, "Deep Learning Based Recommender System : A Survey," *ACM Journals*, vol. 52, no. 1, 2019.
- [8] G. Madjarov, D. Kocev, D. Gjorgjevikj, and S. Džeroski, "An extensive experimental comparison of methods for multi-label learning," *Pattern Recognit.*, vol. 45, no. 9, pp. 3084-3104, 2012.
- [9] F. Nafa, J. I. Khan, and S. Othman, "Extending cognitive skill classification of common verbs in the domain of computer science for algorithms knowledge units," *CSEDU 2017 - Proc. 9th Int. Conf. Comput. Support. Educ.*, vol. 1, no. Csedu, pp. 173-183, 2017.
- [10] F. Nafa, J. I. Khan, S. Othman, and A. Babour, "Mining cognitive skills levels of knowledge units in text using graph tringularity mining," *Proc. - 2016 IEEE/WIC/ACM Int. Conf. Web Intell. Work. WIW 2016*, pp. 1-4, 2017.
- [11] U. Fuller *et al.*, "Developing a computer science-specific learning taxonomy," *ITiCSE-WGR 2007 - Work. Gr. Reports ITiCSE Innov. Technol. Comput. Sci. Educ.*, pp. 152-170, 2007.
- [12] F. Nafa, S. Othman, and J. Khan, "Automatic concepts classification based on bloom's taxonomy using text analysis and the Naïve Bayes Classifier Method," *CSEDU 2016 - Proc. 8th Int. Conf. Comput. Support. Educ.*, vol. 1, no. Csedu, pp. 391-396, 2016.
- [13] A. A. Yahya, A. Osman, A. Taleb, and A. A. Alattab, "Analyzing the Cognitive Level of Classroom Questions Using Machine Learning Techniques," *Procedia - Soc. Behav. Sci.*, vol. 97, pp. 587-595, 2013.
- [14] I. Horie, K. Yamaguchi, K. Kashiwabara and Y. Matsuda, "Improvement of difficulty estimation of personalized teaching material generator by JACET," *2014 Information Technology Based Higher Education and Training (ITHET)*, York, pp. 1-8, 2014.
- [15] A. A. Yahya and A. Osman, "Automatic Classification of Questions Into Bloom's Cognitive Levels Using Support Vector Machine," *Proc. Int. Arab Conf. Inf. Technol.*, pp. 1-6, 2011.
- [16] K. M. Yang, R. J. Ross, and S. B. Kim, "Constructing different learning paths through e-learning," *Int. Conf. Inf. Technol. Coding Comput. ITCC*, vol. 1, pp. 447-452, 2005.
- [17] F. Nafa and J. Khan, "Conceptualize the domain knowledge space in the light of cognitive skills," *CSEDU 2015 - 7th Int. Conf. Comput. Support. Educ. Proc.*, vol. 1, pp. 285-295, 2015.
- [18] Anderson, L. W, & Krathwohl, D. R. "Bloom ' s Taxonomy Action Verbs," 2001.
- [19] B. Thomas, "The Effect of Bloom's Taxonomy on Random Forest Classifier for cognitive level identification of E-content," *EasyChair The world for scientists*, pp. 1-6, 2020.
- [20] K. Kowsari, K. J. Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text classification algorithms: A survey-Good," *Inf.*, vol. 10, no. 4, pp. 1-68, 2019.
- [21] C. C. Aggarwal and C. X. Zhai, "A survey of text classification algorithms," *Min. Text Data*, vol. 9781461432, pp. 163-222, 2012.
- [22] D. C. J. Mr. Benny Thomas, "Machine learning based Text document classification for e-learning," *Int. J. Recent Technol. Eng.*, vol. 8, no. 4, pp. 1-9, 2019.
- [23] D. Grissa, *et al.*, "Feature selection methods for early predictive biomarker discovery using untargeted metabolomic data," *Front. Mol. Biosci.*, vol. 3, no. 30, pp. 1-15, 2016.
- [24] P. Biau, "Analysis of a Random Forests Model G' erard," *Anaesthesiol. Intensive Ther.*, vol. 49, no. 5, pp. 373-381, 2017.
- [25] B. H. Menze *et al.*, "A comparison of random forest and its Gini importance with standard chemometric methods for the feature selection and classification of spectral data," *BMC Bioinformatics*, vol. 10, pp. 1-16, 2009

BIOGRAPHIES OF AUTHORS

Mr. Benny Thomas research scholar from the Christ (Deemed to be University) Bangalore, has completed his MPhil in Computer science from Madurai Kamaraj University. Worked in industry in different roles such as software developer, trainer, and training manager. His research interests are Artificial Intelligence, Big data, Data Analytics and Deep Learning Neural Networks.



Dr. Chandra J. Associate Professor from Department of Computer Science; CHRIST (Deemed to be University) holds a Masters in Computer Applications from Bharathidasan University. PhD from Hindustan University. Her research interests include Artificial Neural Network, Data Mining, Genetic algorithm, Big data analytics, Deep learning, and Convolutional Neural Networks Predictive analytics. And Medical Image Processing.