

# HYBRID APPROACH FOR CLASSIFICATION AND LEARNING USING FUZZY RANDOM FOREST FOR CLINICAL DATA

**Biradar Tanaji G<sup>1</sup>, Ghule Ajjinath M<sup>2</sup>, Gore Dattatraya B<sup>3</sup>,**

**Mandle Yashaswi A<sup>4</sup>, Asst Prof. Sneha Farkade<sup>5</sup>**

<sup>1,2,3,4</sup> Department of Computer Engineering, GSMCOE Balewadi, Pune(India)

<sup>5</sup>Prof. Department of Computer Engineering, GSMCOE Balewadi, Pune(India)

## ABSTRACT

*Random Forests are measured for classification of multisource sensing and geographic data similar to (Health care, climate data, Stock market data etc). different ensemble classification methods have been projected in recent years. These methods have been demonstrated to improve classification correctness considerably. The most widely used ensemble techniques are boosting and bagging. Boosting is support on sample re-weighting as well as bagging uses bootstrapping. The Random Forest classifier uses bagging, or bootstrap aggregating, to appearance an ensemble of categorization and regression tree like classifiers. In accumulation, it searches only a random subset of the variables for a split at each node, in order to minimize the association between the classifiers in ensemble. This method is not responsive to noise or overtraining, as the re sampling is not based on weighting. in addition, it is computationally greatly lighter than methods based on boosting and somewhat lighter than uncomplicated bagging. In the research work, the use of Fuzzy Random Forest classifier for decision tree classification is explored. Finally compare the accuracy of the Random Forest classifier to other better-known ensemble methods on health care data*

**Keywords:** *Random Forests, Classification, Decision trees, Multisource remote sensing data*

## I. INTRODUCTION

Classification has always been a challenging problem [1]. The explosion of information that is available to companies and individuals further compounds this problem. There have been many techniques and algorithms addressing the classification issue. In the last few years we have also seen an increase of multiple classifier systems based approaches, which have been shown to deliver better results than individual classifiers [2]. However, imperfect information inevitably appears in realistic domains and situations. Instrument errors or corruption from noise during experiments may give rise to information with incomplete data when measuring a specific attribute. In other cases, the extraction of exact information may be excessively costly or unviable. Moreover, it may on occasion be useful to use additional information from an expert, which is usually given

through fuzzy concepts of the type: small, more or less, near to, etc. In most real-world problems, data have a certain degree of imprecision. Sometimes, this imprecision is small enough for it to be safely ignored. On other occasions, the imprecision of the data can be modeled by a probability distribution. Lastly, there is a third kind of problem where the imprecision is significant, and a probability distribution is not a natural model. Thus, there are certain practical problems where the data are inherently fuzzy. Therefore, it becomes necessary to incorporate the handling of information with attributes which may, in turn, present missing and imprecise values in both the learning and classification phases of the classification techniques. In addition, it is desirable that such techniques be as robust as possible to noise in the data. Here, we will focus on how to start from a multiple classifier system with performance comparable to or better than the best classifiers and extend it to handle imperfect information (missing values and fuzzy values) and make it robust to noise in nominal attributes and to outliers in numerical attributes [3,4]. To build the multiple classifier system, we follow the random forest methodology [5], and for the processing of imperfect data, we construct the random forest using a fuzzy decision tree as base classifier. Therefore, we try to use the robustness of both, a tree ensemble and a fuzzy decision tree, the power of the randomness to increase the diversity of the trees in the forest, and the flexibility of fuzzy logic and fuzzy sets for imperfect data management.

## 1.1 PROBLEM STATEMENT

In this proposed research work to design and implement a system which provide the parallel processing patient request using HDFS framework which eliminate the queue base patient time prediction system. also focus on boosting and sampling techniques for better classification and prediction for high dimension (imbalance) data with maximum truthiness using Fuzzy Random Forest Classifier.

## 1.2 OBJECTIVES

- To provide efficient and reliable disease diagnostic decision support system to the doctors.
- Increase the prediction Accuracy.
- Runs efficiently on large data bases.
- Handles thousands of input variables without variable deletion.
- Gives estimates of what variables are important in the classification.
- Generates an internal unbiased estimate of the generalization error as the forest building progresses.
- Provides effective methods for estimating missing data.
- Maintains accuracy when a large proportion of the data are missing.
- Provides methods for balancing error in class population unbalanced data sets.
- Generated forests can be saved for future use on other data.
- Prototypes are computed that give information about the relation between the variables and the classification.
- Computes proximities between pairs of cases that can be used in clustering, locating outliers or (by scaling) give interesting views of the data.

Capabilities of the above can be extended to unlabeled data, leading to unsupervised clustering, data views and outlier detection

## II. RELATED WORK

Random Forest produce an ensemble of choice trees. To accomplish variety among base decision trees, Breiman particular the randomization method this works well with bagging or arbitrary subspace methods [10], [11]. To generate every only tree in Random Forest Breiman followed below steps:

If the number of records in the exercise set is  $N$ , then  $N$  records are tested at random but with additional, from the innovative data, this is bootstrap taster. This example will be the training set for mounting the tree. If there are  $M$  number of input variables, then a number  $m \ll M$  is selected such that at each node,  $m$  variables are specific at random out of  $M$  and the top divided on these  $m$  attributes is used to divided the node. The value of  $m$  is held invariable during forest rising. Each tree is developed to the main extent possible. There is no pruning. In this way, many trees are induced in the forest; the numbers of trees are pre-decided by the parameter  $N$  tree. The number of variables ( $m$ ) selected at every node is also referred to as  $m$  try or  $k$  in the literature. The deepness of the tree can be controlled by a given parameters node size (i.e. number of instances in the leaf node) which is usually set to one. Once the forest is accomplished or built as explained above, to classify a new occurrence, it is run across all the trees grown in the forest. Every tree provides categorization for the new occurrence which is recorded as a division. The votes from all trees are joint and the class for which full votes are counted (majority voting) is declaring as categorization of the new instance. This process is referred to as Forest RI in the literature [11]. Here ahead, Random Forest means the forest of choice trees generated using Forest RI procedure. In the forest building process, when bootstrap section set is drawn by sample with replacement for each tree, about 1/3rd of original cases are left out. This set of cases is called OOB (Out-of-bag) data. Each tree has its individual OOB data set which is used for error approximation of individual tree in the forest, called as OOB error estimation. Random Forest algorithm also has in-built facility to compute variable position and proximities [11]. The proximities are used in swapping absent values and outliers

## III. PROPOSED SYSTEM

Classification has always been a tedious problem, As it has been commented earlier, the random forests are good categorization methods and fuzzy sets (with their estimated reasoning capability) have been introduced in the decision trees in a suitable way. Continuing these two good characteristics, in this work the multi-classifiers anticipated are a forest of fuzzy decision trees generated randomly (Fuzzy Random Forest), In this section we specify the adjustments, changes and considerations needed to construct these multi-classifiers.

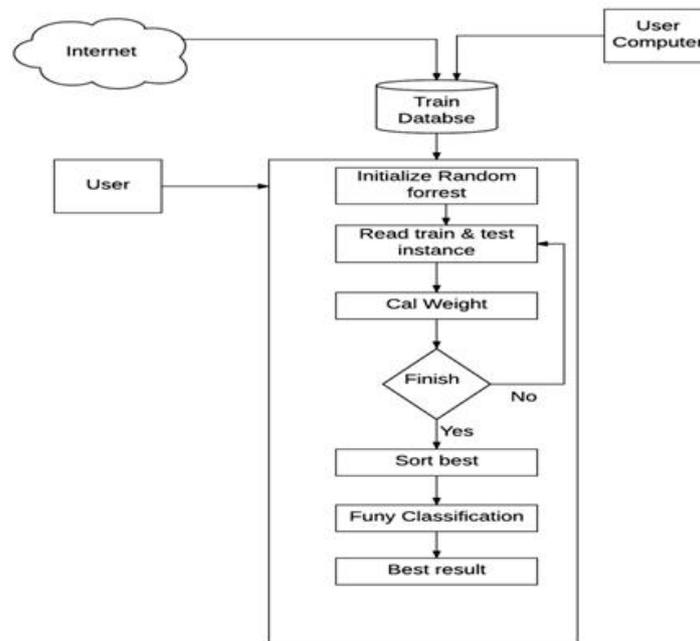


Figure 1 Proposed System architecture

a) Forest structure: at first a forest is constructed from ten trees. For that classical random forest is joint with performance measurement criteria's like Relief and numerous estimators. The forest construction is revealed below. At first, the forest started with ten trees and select a greatest fit is selected from the residual dataset and the construction is made. The similar process is continual up to the 10 trees.

b) Polynomial fitting process for feature selection: Forest construction is an iterative process. Here we have to find similarity of two Nodes  $\vec{a} = (a_1, a_2, a_3, \dots)$  and  $\vec{b} = (b_1, b_2, b_3, \dots)$ , where  $a_n$  and  $b_n$  are the components of the nodes (features of the train node, or values for each features of the node ) and the  $n$  is the dimension of the node:

$$\vec{a} \cdot \vec{b} = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

c) Fuzzy classification: the fuzzy classification executes on probability basis. The RF best nodes get input to fuzzy classifier, it will first analyze the each attribute values base on probability basis then finally classify with label.

**Algorithm**

In this work we propose to use Algorithm 1 to generate a random forest whose trees are fuzzy decision trees, proposing, therefore, a basic algorithm to generate a Fuzzy Random Forest (FRF). Each tree in the forest will be a fuzzy tree generated following the guidelines of [6], adapting it where is necessary.

**Algorithm 1: Random Forest Initialization**

**Step 1 :** Start with examples set of entry, having the weights of the examples (in the root node) equal to 1.

**Step 2 :** At any node N still to be expanded, compute the number of examples of each class. The examples are distributed in part or in whole by branches. The distributed amount of each example to a branch is obtained as the product of its current weight and the membership degree to the node.

**Step 3:** Compute the standard information content.

**Step 4:** At each node search the set of remaining attributes to split the node.

- Select with any criteria, the candidate attributes set to split the node.
- Compute the standard information content to each child node obtained from each candidate attribute.
- Select the candidate attribute such that information gain is maximal.

**Step 5 :** Divide N in sub-nodes according to possible outputs of the attribute selected in the previous step.

**Step 6 :** Repeat steps 2-5 to stop criteria is satisfied in all nodes.

## **Algorithm 2. Fuzzy Decision Tree Learning The fuzzy trees random generator methodology:**

- 1 Read each entry as subset from database according to current weight
2. for each subset Find the best form n partition nodes base on all attributes.
- 3: Display all best attributes from best subset.
- 4: Repeat this step when step 2 return null

## **Algorithm 3. Fuzzy Random Forest ensemble Learning**

**Input:** E, Fuzzy Partition;

**Output:** Fuzzy Random Forest

**Step 1.** Take a random sample of X examples with replacement from the dataset E

**Step 2.** Apply Algorithm 2 to the subset of examples obtained in the previous step to construct a fuzzy tree, using the fuzzy partition

**Step 3.** Repeat steps 1 and 2 until all fuzzy trees are built to constitute the FRF ensemble.

## **Algorithm 4. Fuzzy Decision Tree approach**

**Input:** E, Fuzzy Partition;

**Output:** Classified Fuzzy Tree

**Step 1:** Start with the examples in E with values Fuzzy tree nodes

**Step 2:** Let M be the set of attributes where all numeric attributes are partitioned according to the Fuzzy Partition

**Step 3:** Choose an attribute to do the split at the node N

**3.1:** Make a random selection of attributes from the set of attributes M

**3.2:** Compute the information gain for each selected attribute using the values of fuzzy tree each e in node N

**3.3.** Choose the attribute such that information gain is maximal

**Step 4.** Divide N in children nodes according to possible outputs of the attribute selected in the previous step and remove it from the set M. Let  $E_n$  be the dataset of each child node

**Step 5.** Repeat steps 3 and 4 with each tree until the stopping criteria is satisfied.

## **IV. SYSTEM APPLICATION**

- Patient disease recognition system.
- Health care recommendation hospitalized system
- Queue recommendation base system

## V. HARDWARE & SOFTWARE COMPONENTS

### 5.1 SOFTWARE REQUIREMENT:

#### Front End

- Jdk 1.7.0
- Hadoop 1.2
- Internet Explorer 6.0/above

#### Back-End

- MongoDB

### 5.2 HARDWARE REQUIREMENT:

- Processor:- Intel Pentium 4 or above
- Memory:- 512 MB or above
- Other peripheral:- Printer
- Hard Disk:- 10gb

## VI. CONCLUSION

In this research work, we presented approaches for improving performance of Random Forest classifier using fuzzy logic in terms of accuracy, and time for learning and classification. In case of accuracy improvement, research is done using different attribute evaluation measures and combine functions. A fuzzy decision tree model along with weighted voting is suggested which improves the accuracy using sub classification . Improvement in learning time mainly concerns on reducing number of base decision trees in Random Forest so that learning and in turn, classification is faster. The approaches suggested in this direction are different partitions of training datasets to learn the base decision trees, and ranking of training bootstrap samples on the basis of diversity. Both these approaches are leading to efficient learning of Random Forest classifier. An attempt is made to find optimal subset of Random Forest classifier using Fuzzy classifier. Random Forest has inherent parallelism and can be easily parallelized for scalability and efficiency. A new parallel approach is proposed in which both, individual tree as well as entire forest is generated in parallel. The new approaches presented here are leading to effective learning and classification using Fuzzy Random Forest algorithm

## VII. ACKNOWLEDGEMENT

We are deeply indebted to our Project Guide, Asst. Prof. Sneha Farkade for her valuable guidance and support for completion of this seminar. We are thankful to all our teachers and professors of our department for giving us their expertise in the related topic. We would also like to thank our library staff, internet staff and laboratory assistants for providing us cordial support and necessary facilities which were of great help for preparing this report.

## REFERENCES

- [1] H. Ahn, H. Moon, J. Fazzari, N. Lim, J. Chen, R. Kodell, Classification by ensembles from random partitions of high dimensional data, Computational Statistics and Data Analysis 51 (2007) 6166–6179.

- [2] D. Opitz, R. Maclin, Popular Ensemble Methods: an empirical study, *Journal of Artificial Intelligence* 11 (1999) 169–198.
- [3] P.P. Bonissone, J.M. Cadenas, M.C. Garrido, R.A. Díaz-Valladares, A fuzzy random forest: fundamental for design and construction, in: *Information Processing and Management of Uncertainty in Knowledge-Based Systems International Conference (IPMU2008)*, Málaga, Spain, 2008, pp. 1231–1238.
- [4] J.M. Cadenas, M.C. Garrido, R.A. Díaz-Valladares, Hacia el Diseño y Construcción de un fuzzy random forest, in: *Proceedings in II Simposio sobre Lógica Fuzzy y Soft Computing*, Zaragoza, Spain, pp. 41–48.
- [5] L. Breiman, Random forests, *Machine Learning* 45 (1) (2001) 5–32.
- [6] C.Z. Janikow (1998). Fuzzy decision trees: issues and methods. *Trans. Syst., Man and Cybernetics B* 28(1), 1–15.
- [7] Boinee P, Angelis A, Foresti G, Meta Random Forest, *International Journal of Computational Intelligence* 2, (2006)
- [8] Bonissone P, Cadenas J, Garrido M, Diaz R, A Fuzzy Random Forest: Fundamental for Design and Construction, *Studies in Fuzziness and Soft Computing*, Vol 249, 23-42, (2010)
- [9] Bonissone P, Cadenas J, Garrido M, DiazValladares R, A Fuzzy Random Forest, *International Journal of Approximate Reasoning*, 51, 729-747, (2010)
- [10] Breiman L, Bagging Predictors , Technical report No 421, (1994)
- [11] Brieman L, Random Forests, *Machine Learning*, 45, 5-32,(2001)