

ADVANCED BOOST BRAIN TUMOR CLASSIFICATION WITH RANDOM TREE & KNN SEGMENTATION

Ujwala Sane¹

*Student, Department of Computer Science and Engineering, CSMSS, Chh. Shahu College of Engineering, Aurangabad,
Maharashtra India¹*

Saneujwala1991@gmail.com

Abstract:- Advanced Boost Brain Tumor Classification based on two main components of their system: discrimination and randomization. Discrimination refers to the use of SVM to learn the splits at each node, whereas randomized as a form of features to learn the splits at each node. There are several problems that may arise from this randomization procedure. First upon, we considered image patches of size 50×50 in a 500×500 images, sampling space may contain thousand of patches, which makes for the images categorization. In this, randomly selected samples are more likely to over-lap with each other, which would cause redundancy. Therefore in this project, find out new ways for selecting image patch selection should result in higher quality splits at each tree node, which in turn should increase overall accuracy of the classifier.

Keywords: SVM, KNN, RK-RF

I INTRODUCTION

Attributed their accomplishment to the two main components of their system: discrimination and randomization. Discrimination make reference to the use of SVM to acquire a knowledge of the splits at each node, whereas randomization select to a random selection of image patches, which are used as a form of features to learn the break at each node. There are several problems that may causes from this randomization procedure. Firstly, if we take up image patches of size 50X50 in an 500X500 image, sampling space may accommodate thousands of patches, which assemble it less likely that a randomly selected patch will accommodate an object of interest for the image categorization. In addition, randomly selected samples are more likely to overlay with each other, which would rise a redundancy. For that reason, in this project, find out new ways for selecting image patches. In theory, more informative patch selection should produce in higher quality splits at each tree node, which in turn should increase overall accuracy of the classifier.

II LITERATURE SURVEY

FCM and KNN Based Automatic Brain Tumor Detection

A brain tumor is produce when abnormal cells get acquire within the brain. These cells multiply in an uncontrolled manner and harm the brain tissues. Magnetic Resonance Image scans are commonly used to recognition brain tumors.

However, segmenting and detecting the brain tumor manually is a annoying task for the radiologists. Hence, there is a need for automatic systems which submit correct results. A fully automatic method is introduced to determine brain tumors. It contains the five stages Image Acquisition, Pre-processing, Segmentation, using Fuzzy C-means technique; Harris Corner Detection based feature extraction and classification using K-NN. Performance metrics such as accuracy, precision, sensitivity and specificity are used to find out the performance.

A schematic overview of the proposed approach is demonstrated in. A random forest classifier was applied to the characteristics of data from each modality independently, not only to obtain single-modality classification results for comparison, but also to derive the similarities required for manifold learning. The resulting similarity matrices were combined, and classical was applied to generate a joint embedding for multi-modality classification.

Automated Diagnosis of Diseases Based on Classification: Dynamic Determination of the Number of Trees in Random Forests Algorithm

Most important task of any diagnostic system is the process of take out to determine and identify a possible disease or disorder and the decision reach out by this process. Therefore, machine learning algorithms are popularly working [1], [2]. Therefore machine learning techniques to

be useful in medical illusions problems, they should be distinguish by high performance, the capability to carry out with missing data and noisy data, the clarity of diagnostic knowledge, and the ability to describe decisions. In this paper, the betterment of the random forests classification algorithm, which meets the above mention characteristics, is solve it. This is achieved by to find out the only tuning parameter of the algorithm, which is the number of base classifiers that create the altogether and affects its performance. Random forests are a sub sequential modification of KNN [3]–[6]. It constructs a large number of unproved and decorrelated. trees. The generation of the trees is depend on the combination of two sources of changeability. First, each tree is constructed on a bootstrap reflect of the original data set, as in KNN, and second a random characteristics subset, of fixed predefined size, is studied for splitting each node of the tree. Gini index is used as the feature evaluation measure that determines the best split. The decision tree is built to the maximum size without pruning. The random forests classify each new instance by the majority vote of the full set of trees.

One of the most important issues in the creation of an ensemble classifier, such as random forests, is the size of the at once, the number of classifiers composing the at once, and how the unneeded classifiers are removed from the ensemble. The factors that may affect the size of the ensemble are:

- 1) The required accuracy,
- 2) The calculating cost,
- 3) The nature of the classification problem, and
- 4) The number of available processors.

The methods reported in the literature, dealing with this problem, can be grouped into three categories:

- 1) Methods that preselect the ensemble size,
- 2) Methods that post select the ensemble size (pruning of the ensemble) and
- 3) Methods that select the ensemble size during training.

Preselection methods are the simplest way to find the ensemble size. More precisely, the number of the base classifiers is a tuning parameter of the algorithm, which can be set by the user. Pruning methods consist of precombining and postcombining methods [8]. In the first case, pruning is performed before combining the classifiers. The classifiers that seem to perform well are added in the ensemble. The predictive strength of a classifier is determined using

different estimation of measures. In postcombining pruning methods, the classifiers are deleted from the ensemble based on their contribution to the collective. More accurately, most of the postcombining pruning methods are based on the overproduce-and-choose strategy, which consists of two phases. The overproduction phase aims to produce a large initial pool of candidate classifiers, while the selection phase aims to select adequate classifiers from the pool of classifiers so that the selected group of classifiers can achieve optimum positive predictive rate. In the second phase (selection phase), different approaches are used. More specifically, at once selection methods can be grouped into the following categories:

- 1) Weighted voting methods
- 2) search-based methods
- 3) clustering-based methods
- 4) Ranking methods
- 5) Optimization of a measure or function methods.

III PROBLEM STATEMENT

Attributed their success to the two main components of their system: discrimination and randomization. Discrimination refers to the use of SVM to pursue the splits at each node, whereas randomization mention to a random selection of image patches, which are used as a form of features to pursue the splits at each node. There are several problems that may rise from this randomization procedure. Firstly, if we examine image patches of size 50X50 in an 500X500 image, sampling space may contain thousands of patches, which make out it less likely that a randomly selected patch will contain an object of interest for the image categorization. In addition, randomly selected samples are more likely to overlap with each other, which would cause redundancy. In this project, find out new ways for selecting image patches. In theory, more informative patch selection should produced higher quality splits at each tree node, which in turn should increase overall accuracy of the classifier.

3.1 GOALS AND OBJECTIVE

Before starting Random Forest procedure, I standardize each image by recycling them to the same size and then apply Selective Search Segmentation to derive important regions from each image. Each region is described by 4 coordinates in the image (points in the bottom left and top right corners of the region). Then, SVM is activated to all the regions that were returned by Selective Search Segmentation and its

centroids are chosen as then all candidate regions. In this particular case, I used 1024 centroids

3.2 STATEMENT SCOPE

To fix the problems related to random patch selection I combined a selective search segmentation algorithm into the original random forest framework. Image patches selected using selective search segmentation is more likely to contain the objects of interest. In addition, segmentation should eliminate redundant overlapping between the image patches,

which will make our feature space more diverse. Solving these two problems should result in an increased particular power of random forest.

3.3 OUTCOMES

- ^ Effective recommendation system
- ^ Highly Scalable
- ^ Less time consuming

IV PROJECT IMPLEMENTATION

4.1 ARCHITECTURE OF PROPOSED SCHEDULER

The performance in find out biomarkers for premalignant pancreatic Brain tumor could be improved by using the decision tree ensemble techniques instead of a single algorithm counterpart. These techniques had show more likely to accurately distinguish disease class from normal class as indicated by a larger area under the Receiver Operating Characteristic curve. Moreover, they achieved comparatively lower root mean squared errors.

As stated in to their method, the peptide mass-spectroscopy data were processed first to improve data integrity and over come variations among data due to the differences in sample loading conditions. The pre-processing steps include baseline adjustment using group median, smoothing to delete noise using a Gaussian kernel, and normalization to make all the data similar. After that, the data were one by one sampled such that 90% formed a training set and the remaining 10% formed a test set.

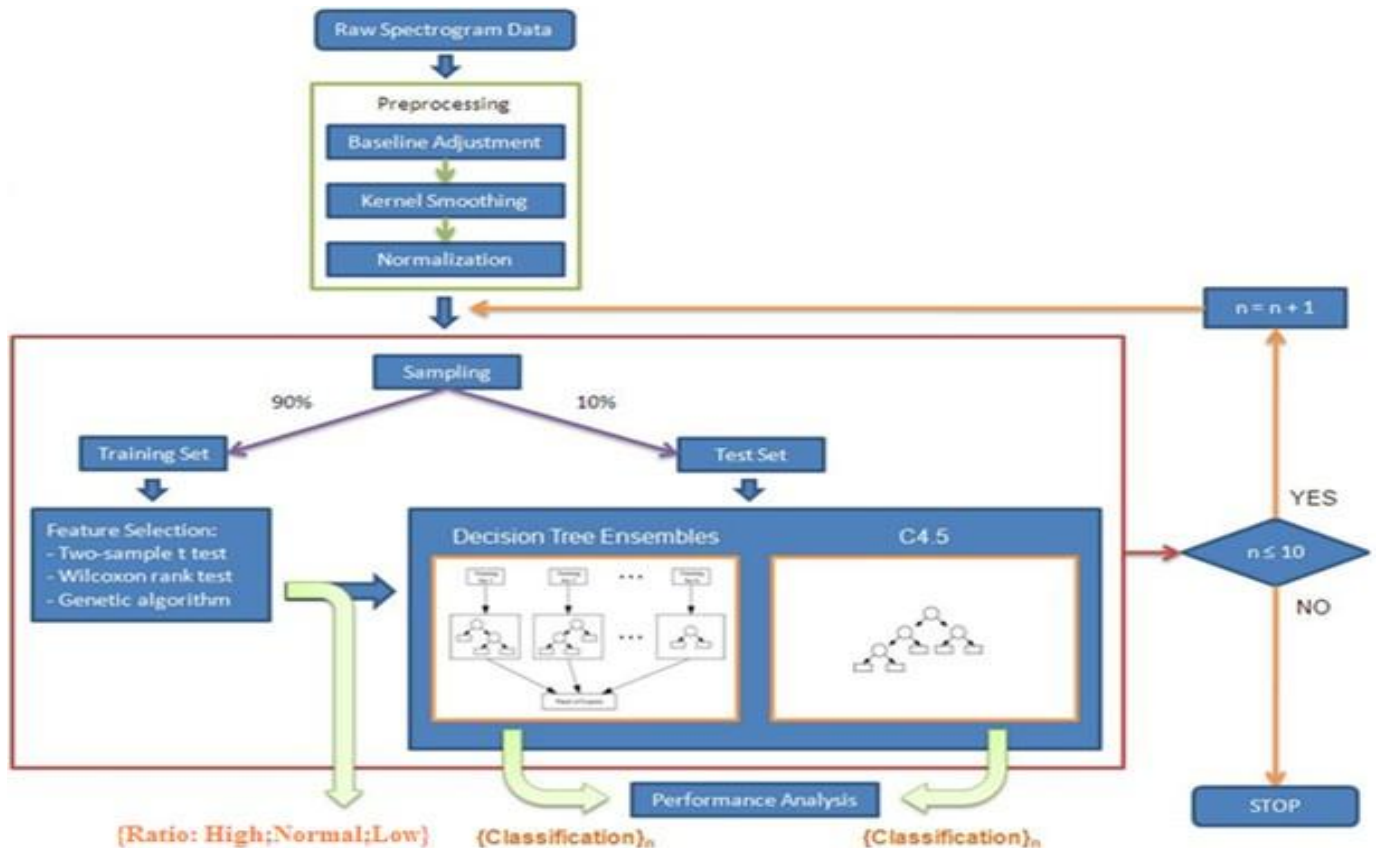


Figure 4.1: Architecture of Proposed System

The training set was used in feature selection. In the study, the authors considered three different feature selection methods. The first method was a two-sample homoscedastic t test, which was used under the assumption that all the features from either normal or disease class had normal distribution. Unlike the first method, the second method based on ANDI rank test considered that the features had no distribution. The last feature selection method was a genetic algorithm.

The test set was used to generate a single decision tree including the decision tree ensembles. The ensemble methods being studied were Random Forest, Random Tree, KNN, boost, Stacking, Adaboost, and Multiboost. Their performances were measured in terms of accuracy and error in the classification of the features, selected by each selection method. Then, they were compared against the performance of a single decision tree generated by C4.5 algorithm. The process repeated ten times to validate the resulting performance consistency. According to the results reported, the decision tree ensembles achieved higher accuracy up to 70% regardless of the feature selection methods used. In terms of biomarker identification, both the t test and the ANDI rank test had similarly impressive performance by consistently selecting the same biomarker-suspect features. Unlike the first two methods, the performance of the genetic algorithm was considerably poor. also noted that 70% accuracy was still lower than expected. This could be as a result from a naturally low concentration of the biomarkers at the premalignant stage of the Brain tumor. In addition, it was also possible that one dataset might not be suitable for all algorithms, thus underestimating the accuracy.

V RESULTS

The proposed procedure was evaluated using eight biomedical datasets () and five versions of the KNN-random forests classification technique (the classical KNN-random forests, RF with KnnRF, RF with me, RK-RF, and RK-RF with me). The classification problem is various depending on the dataset that is used. More specifically, in the scenario of Brain tumor, Parkinson, breast cancer, Pima Indians Diabetes, and SPECTF heart dataset, a two-class classification problem is addressed, if a patient suffers from the corresponding disease (Brain tumors, Parkinson, breast cancer, diabetes, heart) or not. In the scenario of Haberman's survival and Ecoli dataset, a prediction concerning the survival of a patient that has undergone a surgery and the protein localization sites, respectively, must be made.

Finally, in the scenario of breast tissue dataset, the classification problem is reduced to a tissue characterization problem. A breast tissue must be classified into one of the following categories: 1) carcinoma, 2) fibro adenoma, 3) mastopathy, 4) glandular, 5) connective, of the two curves (original and fitted one) using measures of similarity or dissimilarity. When a new tree is added in the forest, the graph of correctness is updated. For this curve, the eight polynomial fits are process and the best one is selected. The two curves are segmented in sliding parts of five points and are compared using: 1) the correlation coefficient (cc), 2) the mean square error (mse) (the average of the squares of the differences between the predicted and actual values), 3) the mean absolute relative error (mare) (the average of the absolute differences between the predicted and actual values divided by the true values), and 4) the mean absolute percentage error (mape) (the average of the absolute errors, as a percentage of the actual values). If there are consecutive parts in the curves, where the above measures are below given thresholds, the procedure stops and returns the size of the collective, otherwise, it continues until and 6) dispose. The criterion for the selection of the specific databases is the diversity they present as far as it concerns the number of samples, the number and nature of the predicted feature data s, the number of classes, and the medical problems they cover. The utilization of diverse datasets in conjunction with the variants of the KNN-random forests technique s aims to indicate the invariance of the procedure to the predicted factors.

The experiment we conducted in order to measure the performance of the proposed procedure is the following. First, we run each one of the variations of the KNN-random forests technique for each dataset for a number of times. In each iteration, the number of trees is increased by one. The procedure is terminated when a forest of 100 trees is created. From the 100 various forests that are created, the one with the best correctness is selected.

The number of the trees of the selected forest, which expresses the optimal size of the collective, is reported in Fig. 2. In the figure, the number of trees understand by the proposed procedure, using each one of the criteria described in Section IIB, is depicted. The analysis of the bar charts reveals that the utilization of the first criterion results to 10 disagreements between the optimal and the proposed collective size. The number of disagreements, when the second criterion is process, using either the curve of

correctness or the curve of correlation, depends on the comparative measure, which is used. More specifically, the mse fails to understand the optimal size of the collective into 12 cases, while the correlation coefficient (cc), the mare, and the mape fail in 6 cases. On the contrary, the combination of correctness and correlation in the stopping criterion (third criterion) leads only to four disagreements.

The comparative results, in scenario the third criterion is used, are reported in Table I. More specifically, the correctness (Acc), the Brier score (BS), and the correlation (Cor) achieved by the best forest and the forest created using the proposed procedure are reported only for the cases, where a disagreement exists between the optimal and the proposed number of trees of the forest.

Technique	Correctness (%)	TP Rate	FP Rate	TN Rate	FN Rate	Sensitivity	Specificity	Precision	Measure	RMSE
Random Forest	0.6500	0.79	0.53	0.48	0.21	0.79	0.48	0.65	0.71	0.4569
KNN	0.6833	0.78	0.44	0.56	0.22	0.78	0.56	0.69	0.73	0.4285
Logitboost	0.6889	0.83	0.49	0.51	0.17	0.83	0.51	0.69	0.75	0.4402
Stacking	0.6444	0.99	0.79	0.21	0.01	0.99	0.21	0.61	0.76	0.4761
Multiboost	0.6889	0.81	0.46	0.5	0.19	0.81	0.54	0.70	0.74	0.5175
Logistic	0.7500	0.79	0.30	0.70	0.21	0.79	0.70	0.78	0.78	0.4224
Naivebayes	0.6833	0.64	0.26	0.74	0.36	0.64	0.74	0.76	0.68	0.5289
Bayesnet	0.6722	0.63	0.28	0.73	0.37	0.63	0.73	0.74	0.67	0.5308
Neural Network	0.7000	0.70	0.30	0.70	0.30	0.70	0.70	0.75	0.72	0.4517
RBFnet	0.6722	0.76	0.44	0.56	0.24	0.76	0.56	0.69	0.71	0.4632
KNN_RF	0.9644	0.71	0.33	0.68	0.29	0.71	0.68	0.74	0.71	0.5489

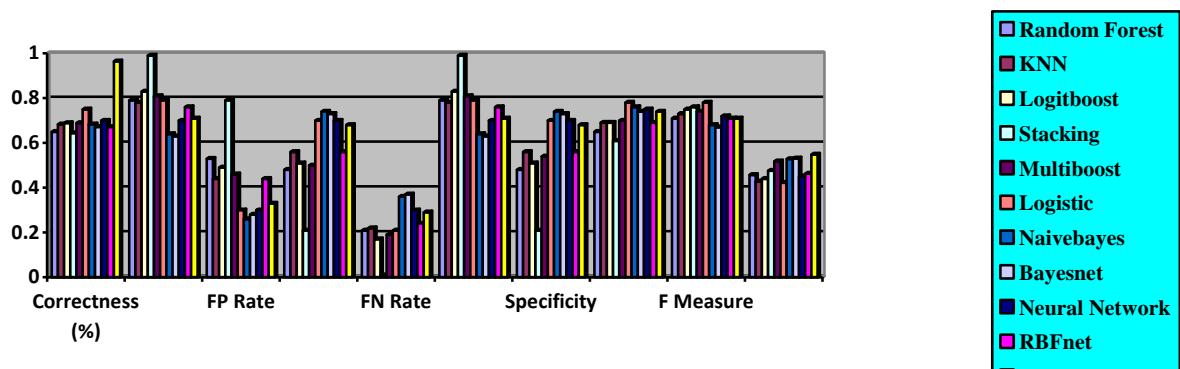


Fig5.1: Graph indicating Classification results using predicted feature data s selected by Student test.

The comparative results for the first and the second criterion are omitted due to the large number of disagreements they present. Table I indicates that the differences in accuracies range from 0% to 1.9% while the differences in brier score and correlation range from 0.002 to 0.009, and 0.001 to 0.036, respectively. In scenario of the first and the second criteria, the differences in the evaluation measures (Acc, BS, and Cor) belong, on average, in the following intervals: Acc $\in [0.45, \dots, 0.88]$, BS $\in [0.004, \dots, 0.01]$, and Cor $\in [0.02, \dots, 0.023]$. According to the description of the proposed procedure, the termination is achieved when for a consecutive number of points the criterion that is used is fulfilled. Thus, an interesting issue that should be examined is the stopping point of the proposed procedure. The number of trees that should be constructed for the stopping criterion to be satisfied is shown on Table II. We must pay attention to the cases, where the stopping point equals to 100, which is the high amount of number of trees that can participate in the collective. This is observed in the following cases: 1) breast cancer with RK-RF and RK-RF with me, 2) Ecoli with classical RF and RF with KnnRF, 3) SPECTF heart with RF with KnnRF, 4) Brain tumor with RK-RF with me and mare, 5) Pima Indians diabetes with classical RF and mare, and 6) This study explores the utility of three various predicted feature data selection schemas to reduce the high dimensionality of a pancreatic Brain tumor proteomic dataset. Using the top predicted feature data s selected from each procedure, we compared the prediction performances of a single decision tree technique C4.5 with six various decision-tree based classifier collectives (Random forest, Stacked generalization, KNN, Adaboost, Logitboost and Multiboost). We show that collective classifiers always out process single decision tree classifier in having greater accuracies and smaller prediction errors when applied to a pancreatic Brain tumor dataset.

Classification results using predicted feature data s selected by Student test.

Breast tissue with RK-RF with me and mape. The results in the first four cases are expected since the proposed procedure is terminated after a small number of trees, 20 to 25, from the optimal point, which in those case are 1) 91, 92, 2) 80, 80, 75, 3) 71, and 4) 69, respectively. In the scenario of the Pima Indians diabetes with classical RF and mare, the criterion is not satisfied, while in the last scenario (Breast Tissue with RK-RF with me and mape), the criterion is not satisfied only for the curve of correctness . When the criterion is not

fulfilled using either one or both curves, the procedure returns as the best collective size the point that was selected as candidate important of are obtained using the third criterion. point is the same for both curves, it is suggested, by the Breast tissue with RK-RF with me and mape. The results in the first four cases are expected since the proposed procedure is terminated after a small number of trees, 20 to 25, from the optimal point, which in those case are 1) 91, 92, 2) 80, 80, 75, 3) 71, and 4) 69, respectively. In the scenario of the Pima Indians diabetes with classical RF and mare, the criterion is not satisfied, while in the last scenario (Breast Tissue with RK-RF with me and mape), the criterion is not satisfied only for the curve of correctness procedure, as the best collectivesize, otherwise, the procedure selects the one with the highest correctness . However, the termination point is 100 trees since one or both curves should be fully grown. It must be mentioned that the curves depicted in Fig. 3 are obtained using the third criterion, which examines both the curve of correctness and correlation. Although, someone should expect, by observing the curve of correctness, the procedure to stop earlier (e.g., when the numbers of trees is equal to 35) the fact that the curve of correlation continues to decline leads the procedure to terminate at a various point, a point with high correctness and low correlation.

Another observation that arises from Table II is that using the mare as a comparative measure, between the original and the fitted curve, smaller number of trees is needed to be constructed in order to terminate the procedure. On average, the mare uses 62 Fewer trees compared to the mape. Thus, the mare performs better than the mape both in time and determination of the best collective size.

VI CONCLUSION

Our proposed system implements a novel classification mechanism for efficiently analyze the brain tumor images using RDTNN classifier. We utilized ROI (Region of Interest) segmentation method for CT image. Using DWT, the key features are extracted; the extracted features are taken as input for RDT to reduce the dimensionality of features. Then the images were trained with KNN classifier. Finally, the proposed algorithm is significantly efficient for classification of the human brain image is benign and malignant with high sensitivity, specificity and accuracy rates. The performance of this study appears some advantages of this technique: it is accurate, robust easy to operate, noninvasive and inexpensive. In future work, we have a plan to explore different types of medicinal images as



well as some other application domains and study some formal properties of image features.

REFERENCES

- [1] I. Kononenko, Machine learning for medical diagnosis: History, state of the art and perspective, *Artif. Intell. Med.*, vol. 23, no. 1, pp. 89109, 2001.
- [2] G. D. Magoulas and A. Prentza, Machine learning in medical applications, *Mach. Learning Appl. (Lecture Notes Comput. Sci.)*, Berlin/Heidelberg, Germany: Springer, vol. 2049, pp. 300307, 2001.
- [3] L. Breiman, KNN predictors, *Mach. Learning*, vol. 24, no. 2, pp. 123140, 1996.
- [4] Y. Freund and R. E. Schapire, A decision-theoretic generalization of online learning and an application to boosting, *J. Comput. Syst. Sci.*, vol. 55, no. 1, pp. 119139, 1997.
- [5] T.K.Ho, The random subspace method for constructing decision forests, *IEEE Trans.*