



Journal of Selected Areas in Microelectronics (JSAM)

Singaporean Journal of Scientific Research(SJSR)

Vol 6.No.1 2014 Pp. 43-47

available at: www.iaaet.org/sjsr

Paper Received : 05-11-2013

Paper Accepted: 18-12-2013

Paper Reviewed by: 1Prof.. S.Soundera Valli 2. Chai Cheng Yue

Editor : Dr. Binod Kumar

Ensemble Classification Techniques for the Diagnosis of Lung Cancer

VIJAYA GURUMANI

RESEARCH SCHOLAR ANNAMALAI UNIVERSITY

Email: viji.pooshan@gmail.com

Dr.SUHASINI ANNAMALAI

ASSOCIATE PROFESSOR ANNAMALAI UNIVERISTY

Email: suha_babu@yahoo.com

ABSTRACT

Throughout the years, Classification has been assuming an essential part in the fields of information mining and in the investigations of machine learning, statistics, neural networks and and numerous master frameworks. Diverse Classification calculations have been effectively executed in different provisions. Lately, restorative information characterization particularly tumor information order has gotten an enormous engage around the specialists. Decision tree classifiers are utilized broadly for diverse sorts of tumor cases. In this study troupe strategies, for example boosting, bagging and random forest have been acknowledged for the examination of execution of precision and time unpredictability for the characterization of tumor datasets. At last random forest beats the other three group techniques.

Keywords: Boosting, Bagging, Random Forest, machine learning, neural networks, decision tree.

1. INTRODUCTION

Lung malignancy is the most well-known growth which prompts demise for both ladies and men, so the unanticipated location of lung tumor increments the treatment victory. Distinctive methods are utilized to give the unanticipated identification, for example Computer Aided Detection (CAD) framework. Lung malignancy could

be seen on conventional x-ray and computed tomography (CT scan). Medication and visualization hinge on upon the histological sort of disease, the stage (level of spread), and the patient's execution status, yet generally just 14% of individuals diagnosed with lung growth survive five years after the diagnosis[1].Drug is an age-old field which holds higher complexities and information instead of whatever possible field. Data mining on

therapeutic information can cause in straightforward characterization to exceptionally might be to get generally thought of the information dependent upon different characteristics, with the intention that the multifaceted nature could be lessened and recognition of oddities could get simpler. Cancer is one such sickness that has more extensive run of spread in India. Measurably, India is discovered to have higher rate of increment in malignancy patients.

The primary excuse for why of growth is tumor. Tumor is unusual development that might be either generous or harmful. Benign tumors are non obtrusive while malignant tumors are destructive and spread to other part of the figure. With the quick headway in data innovation, numerous diverse information mining systems and approaches have been connected to integral prescriptions for tumors [2]. Malignancy information has higher complexities because of different sorts of disease and different techniques for conclusion. The association of the paper is as accompanies: Section 2 arrangements with identified work, with the popoular ensemble frameworks. In Section 3 test effects are looked at and Section 4 presents the conclusion.

2.BACKGROUND

With the gigantic measure of information saved in documents, databases, and different vaults, it is progressively significant, if not fundamental, to advance influential means for examination and maybe elucidation of such information and for the extraction of fascinating learning that could help in choice making. Data Mining [3], likewise ubiquitously reputed to be Knowledge Discovery in Databases (KDD), alludes to the nontrivial extraction of implied, formerly obscure and possibly convenient data from information in databases. While information mining and learning revelation in databases (or KDD) are habitually treated as equivalent words, information mining is really part of the learning finding process.

The classifiers that make up the group are called base models and the studying frameworks that prepared these models the base learners . Both hypothetical and observational examination has demonstrated that an exceptional group is one where the base models are correct and assorted. An exact classifier is one that predicts more than half of the new cases rightly. Two classifiers are differing assuming that they make autonomous

precise expectations. The focal point over utilizing characterization on medicinal information blunders on new information.

Numerous strategies have been proposed through the years to generate different classifiers [10,11]. Arguably, the three most popular ensemble methods are Bagging [4], Boosting [5] and Random Forests [8]. We will portray these in some more detail, next we will quickly touch upon some different methods for developing groups.

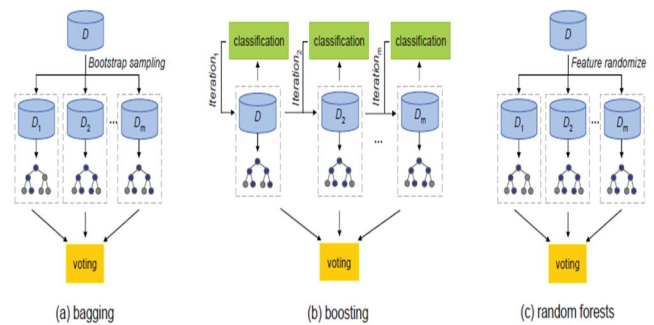


Fig.1. Diagrammatic representation of the three popular ensemble methods.

2.1 BAGGING

In the Bagging system (Bootstrap aggregating) the base learner is prepared I times on distinctive bootstrap repeats of the information. A bootstrap replicate [14] is a specimen of m preparing illustrations drawn arbitrarily with trade from the preparation set of m cases. It holds, on normal, 63.2% of the definitive preparing set, so numerous preparing cases may seem various times, others don't show up in the bootstrap.

The forecasts of the base models are joined with greater part voting to get the expectation of the Bagged classifier. Bagging works particularly well for precarious studying calculations. These are calculations where little changes in the preparation information may bring about huge changes in the yield classifiers. Plainly for stable calculations Bagging may not realize a great deal of differing qualities around the classifiers and will bring about poor groups. It may even somewhat debase the execution of stable calculations (e.g. k-

nearest neighbor), in light of the fact that viably more diminutive information sets are utilized for preparing every classifier.

Decision trees, rule learners and neural networks then again, are usually recognized temperamental calculations and are overall suited as base learners for Bagging. Empirical evaluations[12] with decision trees and neural networks have demonstrated that a Bagging gathering almost dependably outflanks a solitary classifier. [13] likewise tried different things with a few variants on the calculation and reached the accompanying conclusions (around others): bagging works preferable with unpruned trees over with pruned trees as base models the correctness of Bagging is somewhat expanded by averaging likelihood gauges rather than performing a dominant part vote. The accompanying pseudo-code outlines the fundamental thought of Bagging.

2.2. BOOSTING

Boosting [6,7] is a meta-algorithm which can be viewed as a model averaging method. It is the most widely used ensemble method and one of the most powerful learning ideas introduced in the last twenty years. Originally designed for classification, it can also be profitably extended to regression. One first creates a „weak□ classifier, that is, it suffices that its accuracy on the training set is slightly better than random guessing. A succession of models are built iteratively, each one being trained on a data set in which points misclassified (or, with regression, those poorly predicted) by the previous model are given more weight. Finally, all of the successive models are weighted according to their success and then the outputs are combined using voting (for classification) or averaging (for regression), thus creating a final model. The original boosting algorithm combined three weak learners to generate a strong learner [17].

2.2.1. ADABOOST

Adaboost [16], short for 'versatile boosting', is the most prevalent boosting calculation. It utilizes the same preparing set again and again (in this way it require not be extensive) and can likewise join together a self-assertive number of base learners. Adaboost is adjustable as in consequent classifiers

assembled are tweaked energetic about those occasions misclassified by past classifiers.

Adaboost is delicate to uproarious information and outliers. In a few issues, on the other hand, it might be less vulnerable to the overfitting issue than most studying calculations. The classifiers it uses might be frail (i.e., show a generous lapse rate), however provided that their execution is somewhat superior to irregular (i.e. their slip rate is littler than 0.5 for twofold arrangement), they will enhance the last model.

Indeed, classifiers with a failure rate higher than might be normal from an irregular classifier will be handy, since they will have negative coefficients in the last direct fusion of classifiers and subsequently carry on as their inverses. Adaboost produces and calls another powerless classifier in each of an arrangement of rounds $t=1, \dots, T$. For every call, a conveyance of weights D_t is upgraded that demonstrates the essentialness of samples in the information set for the classification.

2.3. RANDOM FOREST

Random forest is an ensemble learning method implemented by growing many classification trees and having them "vote" for a final decision according to a majority role. A random forest generates a number of M decision trees according to the following rule:

(1) Assuming that the number of cases in the training set is N , sample N cases at random with replacement from the original data (bootstrap). This sample will be the training set for growing a tree.

(2) Let the number of features be M . A small number of $m (<<M)$ features are selected at random, and the best split within these features is used to split the node. The value of m is held constant during the growth of the forest.

(3) Each tree is grown to the largest extent possible. There is no pruning. Repeating the creation of a decision tree a number of L times, we obtain L distinct decision trees, forming a randomly generated "forest" [9].

3. EXPERIMENTAL RESULTS

The Lung Cancer Data set is a gathering of lung tumor cases got from the Adyar Cancer Institute, Chennai and from the City Cancer Centre, Madurai. It holds 418 examples depicted by 11

continuous characteristics. 351 examples are tried as positive in tern portrayed as the patient's having threatening tumor and the remaining 67 specimens are tried as negative, indicate that the patient's having benevolent tumor.

For every information set, 10-fold cross acceptance is utilized for assessment. In every fold, preparing information are arbitrarily parceled into named set L and unlabeled set U for a given unlabel rate (μ), which could be figured by the extent of U over the measure of $L \cup U$. Case in point, if apreparation set holds 100 samples, part the preparation consistent with unlabel rate 80% will transform a set with 20 marked illustrations and a set with 80 unlabeled examples.

3.1. CLASSIFICATION ACCURACY(CA)

The classification accuracy A_i of an individual program i depends on the number of samples correctly classified (true positive plus true negative) and is evaluated by the formula:where t is the number of samples correctly classified and n is the total number of samples.

3.2 BRIER SCORE

The Brier Score is probably the most commonly used verification measure for assessing the accuracy of probability forecasts. The score is the mean squared error of the probability forecasts over the verification sample and is expressed as:

where N is the sample size. The observations o_j are all binary, 1 if the event occurs and 0 if it doesn" t. The Brier score ranges from 0 for a perfect forecast to 1 for the worst possible forecast. Although the score can be computed on a single forecast, the result wouldn" t be very meaningful because the observation is binary and the forecast is a probability.

3.3 AREA UNDER ROC CURVE (AUC)

The accuracy of the test depends on how well the test separates the group being tested into those with and without the disease in question. Accuracy is measured by the area under the ROC curve. An area of 1 represents a perfect test; an area of 0.5 represents a worthless test. A rough guide for classifying the accuracy of a diagnostic test is the traditional academic point system:

- 0.90-1 = excellent (A)

- 0.80-0.90 = good (B)
- 0.70-0.80 = fair (C)
- 0.60-0.70 = poor (D)
- 0.50-0.60 = fail (F)

3.4. RESULTS AND DISCUSSION

This section analyses the result of the experiment. To make the study more fruitful and effective two types of group comparison have been made. In the first group, the accuracy is calculated separately for the 4 ensemble methods. The results of the accuracy and its corresponding values are given in Table I.

	CA	BRIER	AUC
TREE	0.784	0.453	0.797
BOOSTED			
BAGGED	0.790	0.407	0.823
FOREST	0.804	0.293	0.894
	0.835	0.176	0.939

CA= Classification Accuracy, BRIER = Brier score, AUC = Area Under Table I. Computing the Accuracy using different ensemble methods

In the second category, error estimation is computed for the above 4 ensemble methods. Table II summarizes the result of error estimation among the various ensemble methods.

	MSE	RMSE	R2
TREE	0.297	0.473	0.527
BOOSTED			
BAGGED	0.311	0.495	0.505
FOREST	0.270	0.430	0.570
	0.224	0.484	0.316

MSE = Mean Squared Error, RMSE = Root Mean

Squared Error, $R^2 = R$ -squared Table II. Calculating the Error using various ensemble methods.

The formula for calculating MSE, RMSE & R^2 is specified in Table III.

Measure	Formula
MSE	$\sum ()$
RMSE	$\sqrt{\sum ()}$
R^2	$\frac{\sum ()}{\sum ()}$

4.CONCLUSION

In this study Simple tree classifier with boosting, bagging and random forest systems have been recognized for the correlation of execution of exactness and blunder estimation for the arrangement of lung tumor dataset. By leading the investigations it is watched that Random Forest is the best calculation for discovering if the tumor is generous or harmful on the tumor datasets which are utilized as they are accessible.

REFERENCES

[1]. Minna. JD and Schiller JH," Harrison's Principles of Internal Medicine (17th ed.)", McGraw-Hill, pp. 551-562,2008

[2]. K.Balachandran and Dr. R.Anitha, "Supervisory expert system approach for pre-diagnosis of lung cancer IJAEA january 2010.aea

[3]. J. Han and M. Kamber, "Data Mining: Concepts and Techniques", Morgan Kaufmann Publishers, 2001

[4]. Leo Breiman, Bagging predictors, Machine Learning, v.24 n.2, p.123-140, Aug. 1996

[5]. Drucker, H. Cortes, C. 1996. Boosting decision trees In Touretsky, D., Mozer, M. Hasselmo, M., Advances in Neural Information Processing Systems, 8, 479-485 Cambridge, MA. MIT Press.

[6]. Y. Freund and R.E. Schapire" A short introduction to boosting" J. Jpn, Soc. Artificial Intelligence.14

[7]. Y. Freund and R.E. Schapire. "Experiments with a new boosting algorithm" in: Proceedings Of the 13th International Conference on Machine Learning, Pg.no: 148-156, 1996.

[8]. Leo Breiman. Random forests. In Machine Learning , pages 5–32, 2001.

[9]. R. Jiang, W. Tang, X. Wu and W. Fu, "A random forest approach to the detection of epistatic interactions in case-control studies" BMC Bioinformatics, vol. 10 (Suppl 1): S65, Jan. 2009

[10]. Thomas G. Dietterich. Ensemble methods in machine learning. Lecture Notes in Computer Science, 1857:1–15, 2000.

[11]. Dietterich, T. Ensemble methods in machine learning. In Proceedings of the 11th International Workshop on Multiple Classifier Systems, Volume 1857 of Lecture Notes in Computer Science, pp. 1–15, 2000.

[12]. Dietterich, T. An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. Machine Learning 40 (2), 139–157, 2000.

[13]. Dietterich, T. (2002). Ensemble learning. In M. Arbib (Ed.),The Handbook of Brain Theory and Neural Networks, Volume 2. The MIT Press.

[14]. Efron, B. and R. Tibshirani (1993). An Introduction to the Bootstrap . Chapman and Hall.

[15]. Opitz, D. and J. Shavlik (1996). Generating accurate and diverse members of a neural-network ensemble. Advances in Neural Information Processing Systems 8, 535–541.

[16]. Bauer, E. and R. Kohavi (1999). An empirical comparison of voting classification algorithms: Bagging, boosting, and variants.Machine Learning 36,105.

[17]. Y. Freund and R.E. Schapire, "Decision-theoretic Generalization of Online Learning and an Application to Boosting", Journal of Computer and System Sciences, vol. 55, no.1, Pg.no:119-139, 1997.