

Breast Cancer Survivability Predictor Using Adaboost and CART Algorithm

R.K.Kavitha¹, Dr.D.DoraiRangasamy²Ph.D Research Scholar, Vinayagar Mission University, Tamilnadu, India¹Professor, Dept. of CSE, Vinayagar Mission University, Tamilnadu, India²

ABSTRACT—Breast cancer is the second leading cancer for women in developed countries including India. Many new cancer detection and treatment approaches were developed, the cancer incidences and death of breast cancer decreased constantly. The patients are concerned about survival time after diagnosis in order to plan regarding their treatments. It is difficult for a physician to have accurate answers about prognosis. Data mining techniques are used to obtain useful information from the large amounts of data which can help the physician for decision making regarding the prognosis. This paper studies the performance comparison of Adaboost algorithm which classifies data as linear combination and CART (Classification and regression trees) which classifies data by constructing decision tree in predicting the survivability of breast cancer patients.

Keywords— prognosis, Adaboost, survival

I INTRODUCTION

Breast cancer is the second most common cause of cancer death in women in developed countries. The most effective way to reduce breast cancer deaths is detect it earlier. Many treatments have developed to reduce the number of mortalities and increase the survival time for patients. In order to predict the survivability of cancer patients, data mining algorithms obtain useful information from the large amounts of data which helps the physician for decision making regarding the prognosis. This paper compares the performance of Adaboost algorithm and CART (Classification and regression trees) algorithm.

II ADABOOST

As a successor of the boosting algorithm, it is used to combine a set of weak classifiers to form a model with higher prediction outcomes. AdaBoost is the most popular ensemble method and has been

shown to significantly enhance the prediction accuracy of the base learner. With this method, medical practitioners are able to focus on finding weak learning algorithms that only should be better than the original algorithm (weak learner). It is a learning algorithm used to generate multiple classifiers and to utilize them to build the best classifier. AdaBoost technique has become an attractive ensemble method in machine learning since it is low in error rate, performing well in the low noise data set. The advantage of this algorithm is that it requires less input parameters and needs little prior knowledge about the weak learner. As a result, several research studies have successfully applied the AdaBoost algorithm to solve classification problems in object detection, including face recognition, video sequences and signal processing systems. AdaBoost algorithm is not only used for predicting in Classification tasks, but also for presenting self-rated confidence scores which estimate the reliability of their predictions. This algorithm requires user less knowledge of computing in order to improve accuracy of models over data sets.

III CLASIFICACION AND REGRESSION TREE

CART stands for Classification And Regression Trees, a decision-tree procedure representing a classification system or predictive model introduced in 1984 by statisticians, Leo Breiman, Jerome Friedman, Richard Olshen, and Charles Stone. . CART builds classification and regression trees for predicting continuous dependent variables and categorical or predictor variables, and by predicting the most likely value of the dependent variable. The decision tree produced by CART is strictly binary, contain two branches of each decision node. CART recursively partitions the record into subsets of record with similar values of target attributes. The CART algorithm grows the tree by conducting for each decision node, an exhaustive search of all variables and all possible

splitting values, selecting the optimal split. It gives an estimate of the misclassification rate.

IV COMPARISON OF ADABOOST AND CART

Gentle AdaBoost, which is originated from the setting of weights over the training set. The training set $(x_1, y_1), \dots, (x_n, y_n)$ where each x_i belongs to instance space X , and each label y_i is in the label set Y , which is equal to the set of $\{-1, +1\}$. It assigns the weight on the training example i on round k as $D_k(i)$. The same weight will be set at the starting point $(D_k(i)=1/N, i=1, \dots, N)$. Then the weight of the misclassified example from base learning algorithm (called weak hypothesis) is increased to concentrate the hard examples in the training set in each round.

The AdaBoost algorithm is presented in seven steps below:

- 1) Assign N example $(x_1, y_1), \dots, (x_n, y_n); x_i \in X, y_i \in \{-1, +1\}$
- 2) Initialise the weights of $D_1(i)=1/N, i=1, \dots, N$
- 3) for $k=1, \dots, K$
- 4) Train weak learner using distribution D_k
- 5) Get weak hypothesis $h_k: X \rightarrow R$ with its error :
 $\epsilon_k = \sum D_k(i)$
- 6) choose $\alpha_k = R$
- 7) Output the final hypothesis:

$$H(x) = \text{sign} \left(\sum_{k=1}^K \alpha_k h_k(x) \right)$$

Classification and regression tree (CART) is a rule-based method that generates a binary tree. through a binary recursive partitioning process that splits a node based on the yes and no answer of the predictors. Although some variables may be used many times, others may not be used at all. A single variable is used to split the tree by using purity criterion. The rule generated at each step is to maximize the class purity within the two resulting subsets. Each subset is split further based on the independent rules to find the threshold among the descriptive variables at the node of all dimensions and they separate the training sample with least error. The steps of constructing the tree included:

1. Create root node;
2. Select leaf with Largest error;

3. Create node, using only those training samples, that are associated with the chosen leaf;
4. Replace selected leaf with created node;
5. Repeat 2-4 until leaves become zero

V STRATIFIED 10 FOLD CLASS VALIDATION

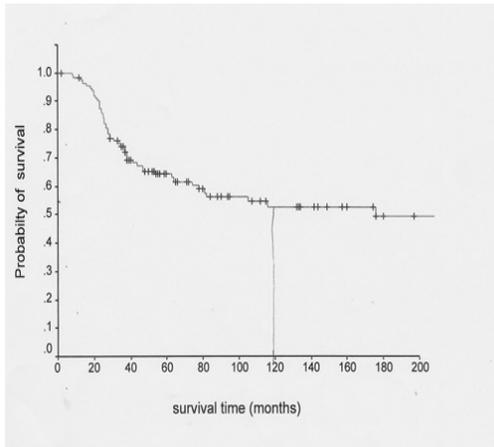
Stratified 10-fold cross-validation is a common validation method used to minimise bias and variance associated with the random sampling of the training and test sets. Moreover, it is a popular method for data selection in data mining related to medical research. In this study the process of stratified 10-fold cross-validation consists of four steps:

- 1) divide the data set into a set of subclasses;
- 2) assign a new sequence number to each set of subclasses;
- 3) randomly partition the subclass into 10 subsets or folds and;
- 4) combine each fold of each subclass into a single fold. Therefore, the size of each single fold is approximately equal to that of the original data set.

VI EXPERIMENTAL ANALYSIS

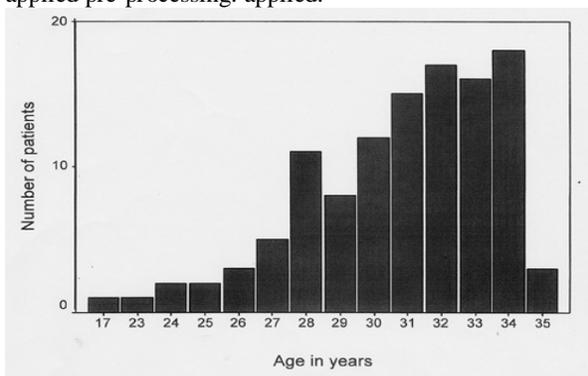
In this study, the models were evaluated based on performance measures including accuracy, sensitivity and specificity. The results were achieved by using stratified

10-fold cross-validation for each model, and were averaged from the test set (the remaining fold), for each fold. Our experiments were done in MATLAB 7 release 14 with GML AdaBoost MATLAB. The experiment results show that the accuracy of the Real and Gentle are decreasing rapidly. Comparing of accuracy, sensitivity and specificity for each classifier was measured by using our breast cancer data set. The same training and test. sets were utilized in all experiments with stratified 10-fold cross-validation and selecting 10 iterations, in order to compare the performance of classification tasks. The experiment results show that Modest AdaBoost outperforms Bagging,

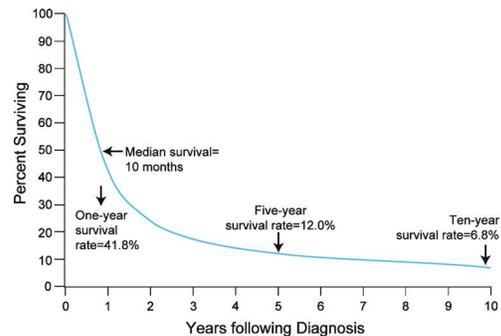


Survival of patients treated for breast cancer (Kaplan-Meier analysis).

In using AdaBoost algorithms to extract breast cancer survivability patterns in breast cancer databases at Hospital, we have successfully utilized stratified 10-fold cross-validation to divide the data set into 10 groups, with the same number in each class. Then presented the accuracy, sensitivity and specificity of classifiers in breast cancer survivability. We found that the accuracy and sensitivity of the models generated from Modest algorithm slightly improved (about 4%) after applied pre-processing. applied.



Presentation of patients with breast cancer by age.



VII CONCLUSION

Adaboost was introduced to achieve better accuracy. Experimental results conducted on the collected breast cancer data set demonstrated the effectiveness of the proposed techniques. Particularly, the proposed method mainly aims to predict the single-point in unknown data rather than estimate multiple-point survival rate in known data, which is usually done by Kaplan-Meier curve. This paper is expected to be of benefit for medical decision making systems to give an alternative choice for medical practitioners to construct more accurate predictive models and stronger classifiers.

REFERENCES

- [1].American Cancer Society, Cancer facts and figures 2006. <http://www.cancer.org/downloads/STT/CAFF2006PWSecured.pdf>. Accessed 24 Jul 2007.
- [3].Borovkova, S., Analysis of survival data. <http://www.math.leidenuniv.nl/~naw/serie5/03/dec2002/pdf/borovkova.pdf>. Accessed.
- [4].Breiman, L., Friedman, J., Olshen, R. and Stone, C. (1984): Classification and regression trees. Wadsworth, Belmont.
- [5]Cabena, P., Hadjinian, P., Stadler, R., Verhees, J. and Zanasi, A. (1998): Discovering data mining from concept to implementation. Upper Saddle River, N.J., Prentice Hall...