



# A Hybrid Random Forest based Support Vector Machine Classification supplemented by boosting

By Tarun Rao & T.V. Rajinikanth

*Acharya Nagarjuna University, India*

*Abstract* - This paper presents an approach to classify remote sensed data using a hybrid classifier. Random forest, Support Vector machines and boosting methods are used to build the said hybrid classifier. The central idea is to subdivide the input data set into smaller subsets and classify individual subsets. The individual subset classification is done using support vector machines classifier. Boosting is used at each subset to evaluate the learning by using a weight factor for every data item in the data set. The weight factor is updated based on classification accuracy. Later the final outcome for the complete data set is computed by implementing a majority voting mechanism to the individual subset classification outcomes.

*Keywords:* *boosting, classification, data mining, random forest, remote sensed data, support vector machine.*

*GJCST-C Classification:* G.4



A HYBRID RANDOM FOREST BASED SUPPORT VECTOR MACHINE CLASSIFICATION SUPPLEMENTED BY BOOSTING

*Strictly as per the compliance and regulations of:*

# A Hybrid Random Forest based Support Vector Machine Classification Supplemented by Boosting

Tarun Rao<sup>α</sup> & T.V. Rajinikanth<sup>σ</sup>

**Abstract-** This paper presents an approach to classify remote sensed data using a hybrid classifier. Random forest, Support Vector machines and boosting methods are used to build the said hybrid classifier. The central idea is to subdivide the input data set into smaller subsets and classify individual subsets. The individual subset classification is done using support vector machines classifier. Boosting is used at each subset to evaluate the learning by using a weight factor for every data item in the data set. The weight factor is updated based on classification accuracy. Later the final outcome for the complete data set is computed by implementing a majority voting mechanism to the individual subset classification outcomes. This method is evaluated and the results compared with that of Support Vector machine classifiers and Random forest classification methods. The proposed hybrid method when applied to classify the plant seed data sets gives better results when compared to traditional random forest and support vector machine classification methods used individually without any compromise in classification accuracy.

**Keywords:** *boosting, classification, data mining, random forest, remote sensed data, support vector machine.*

## 1. INTRODUCTION

Many organizations maintain huge data repositories which store data collected from various sources in different formats. The said data repositories are also known as data warehouses. One of the prominent sources of data is remote sensed data collected via satellites or geographical information systems software's[1].

The data thus collected can be of use in various applications including and not restricted to land use[2] [3], species distribution modeling [4] [5] [6] [7], mineral resource identification[8], traffic analysis[10], network analysis [9] and environmental monitoring systems [11] [12]. Data mining is used to extract information from the said data repositories. The information thus mined can help various stakeholders in an organization in taking strategic decisions. Data can be mined from the data repositories using various methodologies

like anomaly detection, supervised classification, clustering, association rule learning, regression, characterization and summarization and sequential pattern mining. In this paper we shall be applying a hybrid classification technique to classify plant seed remote sensed data.

A lot of research has been undertaken to classify plant functional groups, fish species, bird species etc... [7][13][14].The classification of various species shall help in conserving the ecosystem by facilitating ins predicting of endangered species distribution[15]. It can also help in identifying various resources like minerals, water resources and economically useful trees. Various technologies in this regard have been developed. Machine learning methods, image processing algorithms, geographical information systems tools etc..have added to the development of numerous systems that can contribute to the study of spatial data and can mine relevant information which can be of use in various applications. The systems developed can help constructing classification models that in turn facilitate in weather forecasting, crop yield classification, mineral resource identification, soil composition analysis and also locating water bodies near to the agricultural land.

Classification is the process wherein a class label is assigned to unlabeled data vectors. It can be categorized into supervised and un-supervised classification which is also known as clustering. In supervised classification learning is done with the help of supervisor ie. learning through example. In this method the set of possible class labels is known apriori to the end user. Supervised classification can be subdivided into non-parametric and parametric classification. Parametric classifier method is dependent on the probability distribution of each class. Non parametric classifiers are used when the density function is unknown. Examples of parametric supervised classification methods are Minimal Distance Classifier, Bayesian, Multivariate Gaussian, Support Vector machines, Decision Tree and Classification Tree. Examples of non-parametric supervised classification methods are K- nearest Neighbor, Euclidean Distance, Logistic Regression, Neural Network Kernel Density Estimation, Artificial Neural Network and Multilayer Perceptron. Unsupervised classification is just opposite to the supervised classification i.e. learning is done

**Author α:** Acharya Nagarjuna University, India. His current research interests include Data Mining. e-mail: tarun636@gmail.com

**Author σ:** Sreenidhi Institute of Science and Technology, Hyderabad, India. His current research interests include spatial data mining. Web mining, image Processing and Text mining. e-mail: rajinitv@gmail.com

without supervisor ie. learning from observations. In this method set of possible classes is not known to the end user. After classification one can try to assign a name to that class. Examples of un-supervised classification methods are Adaptive resonance theory(ART) 1, ART 2,ART 3, Iterative Self-Organizing Data Analysis Method, K-Means, Bootstrapping Local, Fuzzy C-Means, and Genetic Algorithm[17]. In this paper we shall discuss about a hybrid classification method. The said hybrid method will make use of support vector machine(SVM) classification, random forest and boosting methods. Later its performance is evaluated against traditional individual random forest classifiers and support vector machines.

A powerful statistical tool used to perform supervised classification is Support Vector machines. Herein the data vectors are represented in a feature space. Later a geometric hyperplane is constructed in the feature space which divides the space comprising of data vectors into two regions such that the data items get classified under two different class labels corresponding to the two different regions. It helps in solving equally two class and multi class classification problem. The aim of the said hyper plane is to maximize its distance from the adjoining data points in the two regions. Moreover, SVM's do not have an additional overhead of feature extraction since it is part of its own architecture. Latest research have proved that SVM classifiers provide better classification results when one uses spatial data sets as compared to other classification algorithms like Bayesian method, neural networks and k-nearest neighbors classification methods[18][19].

In Random forest(RF) classification method many classifiers are generated from smaller subsets of the input data and later their individual results are aggregated based on a voting mechanism to generate the desired output of the input data set. This ensemble learning strategy has recently become very popular. Before RF, Boosting and Bagging were the only two ensemble learning methods used. RF can be applied for supervised classification, unsupervised learning and regression. RF has been extensively applied in various areas including modern drug discovery, network intrusion detection, land cover analysis, credit rating analysis, remote sensing and gene microarrays data analysis etc...[20][21].

Other popular ensemble classification methods are bagging and boosting. Herein the complex data set is divided into smaller feature subsets. An ensemble of classifiers is formed with the classifiers being used to classify data items in each feature subset. The said feature subsets are regrouped together iteratively depending on penalty factor also known as the weight factor applied based on the degree of misclassification in the feature subsets. The class label of data items in the complete data set is computed by aggregating the

individual classification outcomes at each feature subset[22][23].

A hybrid method is being proposed in this paper which makes use of ensemble learning from RF classification and boosting algorithm and SVM classification method. The processed seed plant data is divided randomly into feature subsets. SVM classification method is used to derive the output at each feature subset. Boosting learning method is applied so as to boost the classification adeptness at every feature subset. Later majority voting mechanism is applied to arrive at the final classification result of the original complete data set.

Our next section describes Background Knowledge about Random Forest classifier, SVM and Boosting. In section 3 proposed methodology has been discussed. Performance analysis is discussed in Section 4. Section 5 concludes this work and later acknowledgement is given to the data source followed by references.

## II. BACKGROUND KNOWLEDGE

### a) Overview of SVM Classifier

Support vector machine (SVM) is a statistical tool used in various data mining methodologies like classification and regression analysis. The data can be present either in the form of a multi class or two class problem. In this paper we shall be dealing with a two class problem wherein the seed plant data sets need to be categorized under two class labels one having data sets belonging to North America and the other having data sets belonging to South America. It has been applied in various areas like species distribution, locating mineral prospective areas etc..It has become popular for solving problems in regression and classification, consists of statistical learning theory based heuristic algorithms. The advantage with SVM is that the classification model can be built using minimal number of attributes which is not the case with most other classification methods[24]. In this paper we shall be proposing a hybrid classification methodology to classify seed plant data which would lead to improving the efficiency and accuracy of the traditional classification approach.

The seed plant data sets used in the paper have data sets with known class labels. A classification model is constructed using the data sets which can be authenticated against a test data set and can later be used to predict class labels of unlabeled data sets. Since class labels of data sets are known apriori this approach is categorized as supervised classification. In unsupervised classification method also known as clustering the class label details is not known in advance. Each data vector in the data set used for classification comprises of unique attributes which is used to build the classification model[25][19]. The SVM model can be demonstrated geometrically. As demonstrated in fig 1

SVM is represented by a separating hyper plane  $f(x)$  that geometrically bisects the data space thus dividing it into two diverse regions thus resulting in classification of the input data space into two categories.

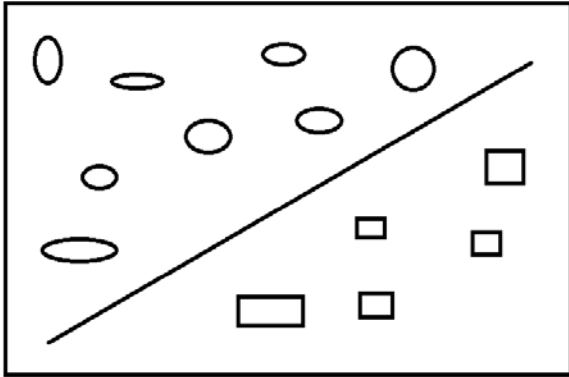


Figure 1 : The Hyperplane

The function  $f(x)$  denotes the hyperplane that separates the two regions and facilitates in classification of the data set. The two regions geometrically created by the hyperplane correspond to the two categories of data under two class labels. A data point  $x_n$  belongs to either of the region depending on the value of  $f(x_n)$ . If  $f(x_n) > 0$  it belongs to one region and if  $f(x_n) < 0$  it belongs to another region. There are many such hyperplanes which can split the data into two regions. But SVM ensures that it selects the hyperplane that is at a maximum distance from the nearest data points in the two regions. There are only few hyperplanes that shall satisfy this criterion. By ensuring this condition SVM provides accurate classification results[27].

SVM's can be represented mathematically as well. Assume that the input data consists of  $n$  data vectors where each data vector is represented by  $x_i \in R_n$ , where  $i (=1, 2, \dots, n)$ . Let the class label that needs to be assigned to the data vectors to implement supervised classification be denoted by  $y_i$ , which is  $+1$  for one category of data vectors and  $-1$  for the other category of data vectors. The data set can be geometrically separated by a hyperplane. Since the hyperplane is represented by a line it can also be mathematically represented by[8][3][28]:

$$\begin{aligned} mx_i + b &\geq +1 \\ mx_i + b &\leq -1 \end{aligned} \tag{1}$$

The hyperplane can also be represented mathematically by [31][32][33]:

$$\begin{aligned} f(x) &= \text{sgn}(mx + b) \\ &= \text{sgn}\left(\sum_{i=1}^n \alpha_i y_i x_i \cdot x + b\right) \end{aligned} \tag{2}$$

where  $\text{sgn}()$  is known as a sign function, which is mathematically represented by the following equation:

$$\text{sgn}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases} \tag{3}$$

The data vectors are said to be optimally divided by the hyperplane if the distance amid the adjoining data vectors in the two different regions from the given hyperplane is maximum.

This concept can be illustrated geometrically as in Figure 2, where the distance between the adjoining data points close to the hyperplane and the hyperplane is displayed[29][30][28].

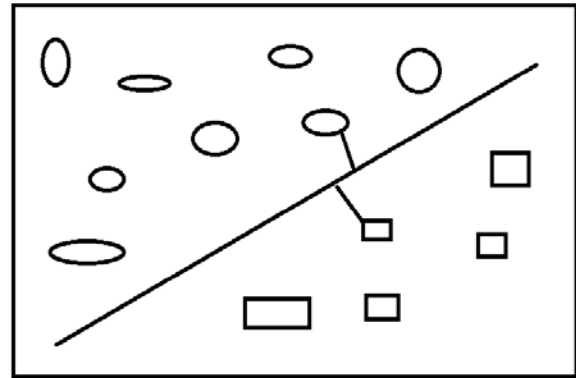


Figure 2 : Distance of the nearest data vectors from the Hyperplane

The distance  $d$  of a data point  $x$  from the hyperplane is represented mathematically by the equation:

$$d = \frac{|(m \cdot x) + b|}{|m|} \tag{4}$$

This hyperplane which has maximum distance  $d$  from adjoining points is computed to implement the said classification. This SVM can be represented as a primal formulation given by the equation [8][5][31]:

$$\begin{aligned} h(m) &= \frac{1}{2} \|m\|^2 + \text{Training error} \\ &\text{subject to } y_i(m \cdot x_i + b) > 1, \forall i \end{aligned} \tag{5}$$

The idea is to increase the margin and reduce the training error. The data sample records in the training data set belong to input set. Each of the data vectors have precise attributes based on which the classification model is built. These set of attributes are said to form a feature space. The kernel function bridges the gap between the feature space and the input space and enables to carry out classification on input space rather than complicated feature space. [29].

In this paper we have used Gaussian radial basis functions (RBF). SVM's make use of the radial basis kernel function to be able to work at the simpler input space level. The RBF kernel used is represented mathematically by[3][29]:

$$K(x_1, x_2) = \exp\left(\frac{|x_1 - x_2|^2}{2\sigma^2}\right) \tag{6}$$

SVM selects a local Gaussian function and later the global Gaussian function is computed by aggregating all the local Gaussian function.

SVM can be used to solve either two class or multi-class problems. Multiclass classification problems



can be solved using various methods. One method is to move the data vectors to a different space thereby making the problem linear. The other method is to split the multi class problem into numerous two class problems and later with a voting mechanism combine the solutions of individual two class problems to get the solution of the original multi class problem. [8].

The steps followed while using SVM in classifying data are mentioned in the below algorithm [16]:

-----  
Algorithm 2 Classification using linear kernel based SVM Classifier  
-----

Input: I: Input data  
Output: V: Support vectors set  
begin  
Step 1: Divide the given data set into two set of data items having different class labels assigned to them  
Step 2: Add them to support vector set V  
Step 3: Loop the divided n data items  
Step 4: If a data item is not assigned any of the class labels then add it to set V  
Step 5: Break if insufficient data items are found  
Step 6: end loop  
Step7: Train using the derived SVM classifier model and test so as to validate over the unlabelled data items.  
end  
-----

#### b) Overview of Random Forest Classifier

Ensemble learning algorithms use an ensemble or a group of classifiers to classify data. Hence they give more accurate results as compared to individual classifiers. Random forest classifier is an example for ensemble classifiers. Random forests make use of an ensemble of classification trees [34][35][36][37][38][41].

In RF classification method the input data set is first subdivided into two subsets, one containing two thirds of the data points and the other containing the remaining one third. Classification tree models are constructed using the subset comprising of two thirds of data points The subset which contains one third data of data points which are not used at any given point of time to construct classification trees and are used for validation are called out of bag(OOB) data samples of the trees. There is no truncation applied at every classification tree. Hence every classification tree used in RF classification method is maximal in nature. Later RF classification method follows a majority voting process wherein classification output of every classification tree casts a vote to decide the final outcome of the ensemble classifier ie.. assigning a class label to a data item x[21]. The set of features are used to create a classification tree model at every randomly chosen subset[37]. This set of features shall remain constant throughout the growing of random forest.

In RF, the test set is used to authenticate the classification results and also used for predicting the class labels for unlabeled data after the classification model is built. It also helps in cross validation of results among different classification results provided by various classification trees in the ensemble. To perform the said cross validation the out of bag(OOB) samples are used.. The individual classification tree outcomes are aggregated with a majority vote and the cumulative result of the whole ensemble shall be more accurate and prone to lesser classification error than individual classification tree results[26].

Every classification tree in the random forest ensemble is formed using the randomly selected two thirds of input variables, hence there is little connection between different trees in the forest. One can also restrict the number of variables that split a parent node in a classification tree resulting in the reduction of connection between classification trees. The Random forest classification method works better even for larger data sets. This is not the case with other ensemble methods[1][2]. In this paper we shall be using the both boosting and random forest ensemble classification methods along with support vector machines to give a more accurate classification output. This hybrid method shall be more robust to noise as compared to individual classification method.

RF classification method works with both discreet and continuous variables which is not the case with other statistical classification modeling methods. Furthermore, there is no limit on the total number of classification trees that are generated in the ensemble process and the total number of variable or data samples(generally two thirds are used) in every random subset used to build the classification trees[36].

RF rates variables based on the classification accuracy of the said variable relative to other variables in the data set. This rank is also known as importance index. It reflects the relative importance of every variable in the process of classification. The importance index of a variable is calculated by averaging the importance of the variable across classification trees generated in the ensemble. The more the value of this importance index, the greater is a variables importance for classification. Another parameter obtained by dividing the variable's importance index by standard error is called z-score. Both importance index as well as z-score play a significant role in ensuring the efficiency of the classification process[25][36][39][38].

The importance of a variable can also be assessed by using two parameters, Gini Index decrease and OOB error estimation. Herein relative importance of variables are calculated which is beneficial in studies wherein the numbers of attributes are very high and thus leading to relative importance gaining prominence[40]. For any training data set T Gini index value can be computed for two variables i and j as[2]:

$$\sum \sum_{j \neq i} \left( \frac{k(C_i, X)}{|X|} \right) \cdot \left( \frac{k(C_j, X)}{|X|} \right) \quad (7)$$

where  $\frac{k(C_i, X)}{|X|}$  is the probability that a selected case belongs to class  $C_i$ .

RF method provides precise results with respect to variation and bias[39]. The performance of the RF classification method is better compared to other classifiers like support vector machines, Neural Networks and discriminant analysis. In this paper a hybrid classification method coalescing the advantages of both Random forest and Support vector machines in addition to boosting is used. The RF algorithm is becoming gradually popular with applications like forest classification, credit rate analysis, remote sensing image analysis, intrusion detection etc.

Yet another parameter that can contribute in assessing the classification is proximity measure of two samples. The proximity measure is the number of classification trees in which two data samples end up in the same node. This parameter when divided by the number of classification trees generated can facilitate in detecting outliers in the data sets. This computation requires large amount of memory space, depending on the total number of sample records and classification trees in the ensemble[1]. The pseudo code for Random Forest algorithm is mentioned below[42]:

-----  
 Random Forest Algorithm:  
 -----

- Input: D: training sample  
 a: number of input instance to be used to generate classification tree  
 T: total number of classification trees in random forest  
 OT: Classification Output from each tree T  
 1) OT is empty  
 2) for  $i=1$  to T  
 3)  $D_b$  = Form random sample subsets after selecting 2/3rd instances randomly from D  
 /\* For every tree this sample would be randomly selected\*/  
 4)  $C_b$  = Build classification trees using random subsets  $D_b$   
 5) Validate the classifier  $C_b$  using remaining 1/3rd instances //Refer Step 3.  
 6) OT=store classification outputs of classification trees  
 7) next i  
 8) Apply voting mechanism to derive output ORT of the Random forest(ensemble of classification trees)  
 9) return ORT  
 -----

c) *Overview of Boosting*

Ensemble learning is a process wherein a data set is divided into subsets. Individual learners are then used to classify and build the model for each of these subsets. Later the individual learning models are combined so as to determine the final classification

model of the complete data set. As the complex large data set is divided into smaller random subsets and classification model is applied on these smaller subsets the said process of ensemble learning results in improving classification efficiency and gives more accurate results. Numerous classification methodologies like bagging, boosting etc...can also be used in learning by constructing an ensemble[43][44][45].

In this research paper boosting method has been used to create the said ensemble. It works by rewarding successful classifiers and by applying penalties to unsuccessful classifiers. In the past it has been used in various applications like machine translation [46], intrusion detection [47], forest tree regression, natural language processing, unknown word recognition [48] etc.

Boosting is applied to varied types of classification problems. It is an iterative process wherein the training data set is regrouped together into subsets and various classifiers are used to classify data samples in the subsets. The data samples which were difficult to classify by a classifier also known as a weak learner at one stage are classified using new classifiers that get added to the ensemble at a later stage[49][50][51]. In this way at each stage a new classifier gets augmented to the ensemble. The difficulty in classifying a data item  $X_i$  at stage  $k$  is represented by a weight factor  $W_k(i)$ . The regrouping of training sets at each step of learning is done depending on the weight factor  $W_k(i)$ [22]. The value of the weight factor is proportional to the misclassification of the data. This way of forming regrouped data samples at every stage depending on the weight factor is called re-sampling version of boosting. Yet another way of implementing boosting is by reweighting wherein weight factor is assigned iteratively to every data item in the data set and the complete data set is used at every subsequent iteration by modifying the weights at every stage[48][52].

The most popular boosting algorithm called Adaboost[23]. Adaboost stands for Adaptive Boosting. It adapts or updates weights of the data items based on misclassification of training samples due to weak learners and regroups the data subsets depending on the new weights. The steps of Adaboost algorithm is mentioned below:

-----  
 Adaboost Algorithm  
 -----

- Divide the data set into random subsets  
 Set uniform example weights to all data items in the subsets  
 For each base learner  
 do  
 Train base learner with weighted sample of a sample data set  
 Test base learner on complete data set.  
 Update learner weight based on



misclassifications by base learners also called weak learners if misclassification occurs Set example weights based on predictions by ensemble of classifiers.  
end for

In the next section the proposed hybrid methodology is discussed in detail.

### III. PROPOSED METHODOLOGY

In this paper we shall construct a hybrid classification model which shall facilitate in predicting the class label of seed plant data from test data sets. The methodology recommended has been denoted as a schematic diagram as mentioned in Fig 3 and the detailed explanation of the steps followed has been given in the following subsections.

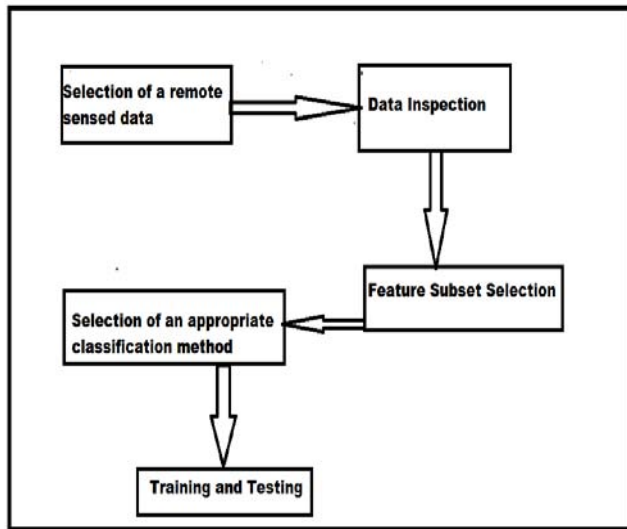


Figure 3 : Proposed Model

#### a) Selection of Remote sensed data

The data sets collected in this paper belong to various types of seed plant family viz Pinopsida, Dicotyledoneae, Monocotyledoneae, Cycadopsida, Pinopsida, Gnetopsida, Lycopodiopsida, Agaricomycetes and Marchantiopsida.

#### b) Data Inspection

The seed plant data sets are pre-processed and any missing values or duplicate values are eliminated by either ignoring tuples comprising of duplicate values or by manually entering values or replacing with a global constant or a mean value into tuples with missing values[53].

#### c) Feature subset Selection

The data sets are randomly divided into n different random subsets each subset comprising of two third of the whole data set. Classification methods are applied to each of these random subsets. The remaining

one third data sets at each subsets is used as a test set. At each random subset the following attributes were used so as to implement the classification method discussed in the next subsection: id, continent, specificEpithet and churn. Now churn is a variable that is set to yes if the seed plant data belongs to North America or if it belongs to South America it is set to no.

#### d) Selection of an appropriate classification method

In this paper seed plant data sets are classified using a hybrid classification method which makes use of Random forest, SVM classifier and boosting ensemble learning method. In the hybrid methodology the input data set is randomly subdivided into subsets. Each data item in each of the subset has a weight factor associated with it. The data items in the subsets are classified by SVM classifier. If a misclassification has occurred then the weight factor of the data items is increased otherwise it is reduced. The data subsets are rearranged and again SVM classifier is used to perform classification at each subset. The weights are again updated depending on whether it is a proper classification or a misclassification. These steps are iteratively repeated till all the weights get updated to a very low value. The output of the input data set is computed by applying voting mechanism to all the random subsets classification outputs[34]. The algorithm for the proposed hybrid methodology is given in the sample code herein:

Algorithm 1 Hybrid classification using RF and SVM supplemented by boosting

```

Input: D Training Instances
Intermediate Output: Osvm, Classification output at each feature subset
Output: O, Classification Output for the hybrid method
Step 1)Begin
Step 2)Initialize the weight  $w_i$  for each data vector  $i \in D$ .
Step 3) Generate a new data feature subset  $D_i$  from D using random replacement method.
Step 4)begin
Step 5)For each random feature subset  $D_i$  do
Step 6)begin
Step 7)Apply SVM to each feature subset
Step 8) Generate  $O_{svm}$ , the classification output from
Step 9) end
Step 10) Update the weights of all the data vectors in the training set depending on the classification outcome. If an example was misclassified then its weight is increased, or else the weight is decreased.
Step 11) Repeat steps 2 to 10 by regenerating random subsets till all the input data vectors are appropriately classified or apply iteration limit.
Step 12) Compute output O of the complete data set by applying majority voting mechanism among the final outputs of each of the Random feature subsets  $D_i$  of The original set D obtained after Step 11.
Step 13) Return O
Step 14)End
    
```

e) *Training and Testing*

The obtained classification output at each random subset is validated by using the hybrid classifier model to test against the complete data set.

In this paper 10 random feature subsets were used and at every subset SVM classifier was used to perform the said classification. Voting mechanism was then applied to derive the final classification output. In this paper a total of 180 support vectors were used.

IV. PERFORMANCE ANALYSIS

a) *Environment Setting*

The study area included is from North and South America. It includes data pertaining to localities wherein seed plant species are present.

A total of 599 data set records from North American region and a total of 401 data set records from South American region are analyzed in order to execute the proposed method. Sample records used in this paper are shown in Table I shown below:

Table 1 : Sample records

id	higherGeography	continent	family	scientificName	decimalLatitude	decimalLongitude	genus	specificEpi	chum
2759	North America,	North	Lycoperdaceae	Calvatiaarctica	72	-40	Calvati	arctica	yes
86	GREENLAND	America	aeae	Empetrum					
3333	North America,	North	Ericaceae	&Wiegand	52	-56	Empetrum	eamesii	yes
01	North America,	America	Ranunculaceae	Thalictrum			Thalictrum	terrae-novae	yes
2717	North America,	North	Ranunculaceae	Greene	52	-56	Thalictrum	terrae-novae	yes
58	North America,	America							

A total of 3 features are first extracted as stated in section III. Then obtained feature vectors are fed into the hybrid classifier whose results are compared against SVM and Random forest Classifier results. A total of 80 data set records act as test data set and are used to authenticate the classification results obtained. The proposed method has been implemented under the environment setting as shown in Table II[54].

Table 2 : Environment Setting

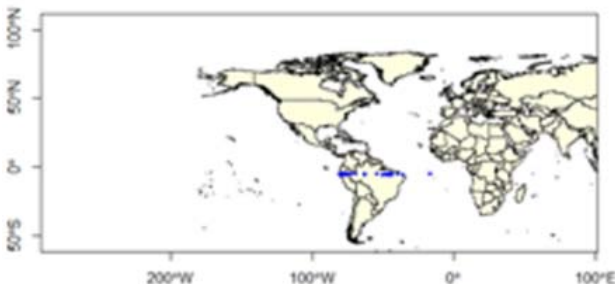
Item	Capacity
CPU	Intel CPU G645 @2.9 GHz processor
Memory	8GB RAM
OS	Windows 7 64-bit
Tools	R, R Studio

b) *Result Analysis*

Classification of the spatial data sets can be represented as a confusion or error matrix view as shown in Table III. And the classified seed plant data is demonstrated in Figure 4[54].

Table 3 : Confusion / Error Matrix View

Real group	Classification result	
	North America	South America
North America	True Negative(TN)	False Positive(FP)
South America	False Negative(FN)	True Positive(TP)



(a)

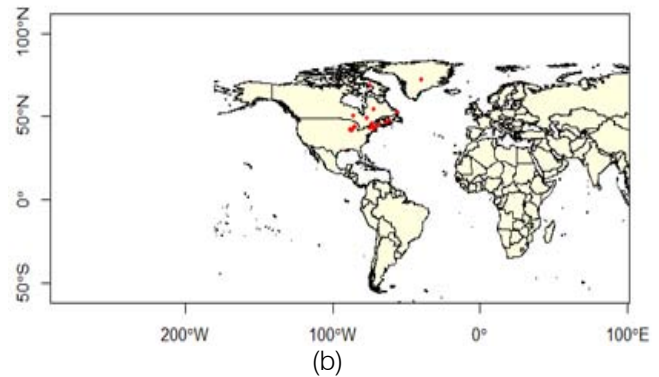


Figure 4 : Plot of seed plant data (a) Seed plant belonging to South America marked in blue. (b) Seed plant belonging to North America marked in red.

It is observed that the most conventionally utilized evaluation metrics in classification are accuracy, specificity, positive predictive value and negative predictive value. The formulae for accuracy, specificity, prevalence and negative predictive value are provided by equations (8), (9), (10) and (11)[54]:

$$\text{Accuracy} = \frac{TP + TN}{(TP + FN + FP + TN)} \times 100 \quad (8)$$

$$\text{Specificity} = \frac{TN}{(TN + FP)} \times 100 \quad (9)$$

$$\text{Prevalence} = \frac{(TN + FN)}{(TP + FP + TN + FN)} \times 100 \quad (10)$$

$$\text{Neg. Predictive Value} = \frac{TN}{TN + FN} \times 100 \quad (11)$$

The efficiency of the proposed hybrid classification is evaluated and compared with traditional RF ensemble and SVM classification methods. The confusion or error matrix view for hybrid classifier is given in Table IV.



Table 4 : Confusion Matrix for Hybrid Classifier

Prediction	Reference	
	South America	North America
South America	1	4
North America	5	70

The confusion matrix or error matrix view for SVM Classifier is given in Table V and for RF Classifier in Table VI.

Table 5 : Confusion Matrix for SVM

Prediction	Reference	
	South America	North America
South America	8	7
North America	49	16

Table 6 : Confusion Matrix for RF

Prediction	Reference	
	South America	North America
South America	36	11
North America	21	12

Performance Measures using evaluation metrics are specified in Fig 5 which are calculated using equations (8), (9), (10)and (11).

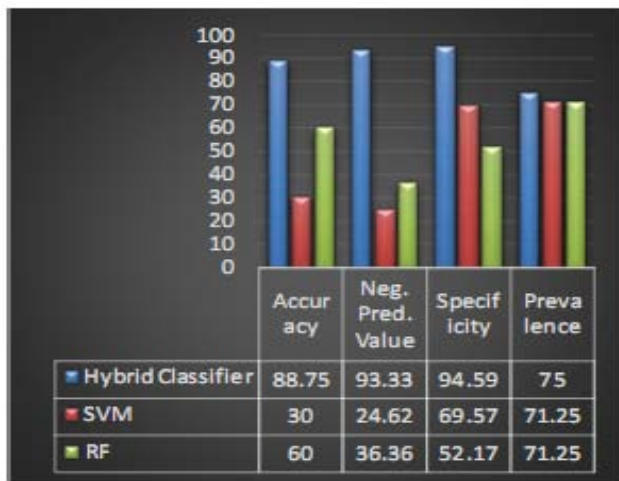


Figure 5 : Performance measures for proposed hybrid, SVM,RF classifiers

Hence, from the results obtained it is inferred that the proposed hybrid classifier is better than traditional SVM and random forest classifiers.

### V. CONCLUSION

In this paper hybrid classifier based on random forest, SVM and boosting methods is used to classify seed plant data. The hybrid classification results are compared with the results attained by implementing classification using traditional SVM and RF classifiers. The research has established that the hybrid approach of classification is more efficient as compared to traditional SVM and RF classifiers since it gives higher values of

accuracy, specificity, positive predictive value and negative predictive value.

The reason for better results in the case of hybrid classification methodology used in this paper is since it makes use of the advantages of each of the individual traditional SVM, RF classifications methods. Furthermore, the classification results are supplemented using boosting ensemble classification method. In the future the proposed method can be used so as to classify vector, raster remote sensed data that can be collected via satellites and various geographical information systems.

### VI. ACKNOWLEDGMENT

We direct our frank appreciativeness to the Field Museum of Natural History(Botany)- Seed Plant Collection (accessed through GBIF data portal, <http://data.gbif.org/datasets/resource/14346,2013-06-03>) for providing us with different seed plant data sets. We also thank ANU university for providing all the support in the work conducted.

### REFERENCES RÉFÉRENCES REFERENCIAS

1. Pall Oskar Gislason, Jon Atli Benediktsson, Johannes R. Sveinsson, Random Forests for land cover classification, Pattern Recognition Letters, Volume 27, Issue 4, March 2006, Pages 294-300, ISSN 0167-8655, <http://dx.doi.org/10.1016/j.patrec.2005.08.011>.
2. V.F. Rodriguez-Galiano, B.Ghimire, J. Rogan, M. Chica-Olmo, J.P .Rigol-Sanchez, An assessment of the effectiveness of a random forest classifier for land-cover classification, ISPRS Journal of Photogrammetry and Remote Sensing, Volume 67, January 2012, Pages 93-104, ISSN0924-2716, <http://dx.doi.org/10.1016/j.isprsjprs.2011.11.002>.
3. Hassiba Nemmour, Youcef Chibani, Multiple support vector machines for land cover change detection: An application for mapping urban extensions, ISPRS Journal of Photogrammetry and Remote Sensing, Volume 61, Issue 2, November 2006, Pages 125-133,ISSN0924-2716,<http://dx.doi.org/10.1016/j.isprsjprs.2006.09.004>.
4. P.A. Aguilera, A. Fernández, F.Reche, R. Rumí, Hybrid Bayesian network classifiers: Application to species distribution models, Environmental Modelling & Software, Volume25, Issue12, December 2010, Pages 1630-1639,ISSN1364 8152,<http://dx.doi.org/10.1016/j.envsoft.2010.04.016>.
5. Till Rumpf, Christoph Römer, Martin Weis, Markus Sökefeld, Roland Gerhards, Lutz Plümer, Sequential support vector machine classification for small-grain weed species discrimination with special regard to Cirsiumarvense and Galiumaparine, Computers and Electronics in Agriculture, Volume 80, January 2012,

- Pages89-96, ISSN0168-1699,<http://dx.doi.org/10.1016/j.compag.2011.10.018>.
6. Jing Hu, Daoliang Li, Qingling Duan, Yueqi Han, Guifen Chen, XiuliSi, Fish species classification by color, texture and multi-class support vector machine using computer vision, *Computers and Electronics in Agriculture*, Volume 88, October 2012, Pages 133-140, ISSN0168-1699, <http://dx.doi.org/10.1016/j.compag.2012.07.008>.
  7. L. Naidoo, M.A. Cho, R. Mathieu, G. Asner, Classification of savanna tree species, in the Greater Kruger National Park region, by integrating hyperspectral and LiDAR data in a Random Forest data mining environment, *ISPRS Journal of Photogrammetry and Remote Sensing*, Volume 69, April 2012, Pages 167-179, ISSN0924-2716, <http://dx.doi.org/10.1016/j.isprsjprs.2012.03.005>.
  8. Maysam Abedi, Gholam-Hossain Norouzi, Abbas Bahroudi, Support vector machine for multi-classification of mineral prospectivity areas, *Computers & Geosciences*, Volume 46, September2012, Pages 272-283, ISSN 0098-004, <http://dx.doi.org/10.1016/j.cageo.2011.12.014>.
  9. D.Ruano- Ordás, J. Fdez-Glez, F. Fdez-Riverola, J.R. Méndez, Effective scheduling strategies for boosting performance on rule-based spam filtering frameworks, *Journal of Systems and Software*, Volume86, Issue12, December2013, Pages 3151-3161, ISSN 0164-1212, <http://dx.doi.org/10.1016/j.jss.2013.07.036>.
  10. N. Abdul Rahim, Paulraj M.P., A.H .Adom, Adaptive Boosting with SVM Classifier for Moving Vehicle Classification, *Procedia Engineering*, Volume 53, 2013, Pages 411-419, ISSN1877-7058, <http://dx.doi.org/10.1016/j.proeng.2013.02.054>.
  11. Konstantinos Topouzelis, Apostolos Pysillos, Oil spill feature selection and classification using decision tree forest on SAR image data, *ISPRS Journal of Photogrammetry and Remote Sensing*, Volume 68, March 2012, Pages 135-143, ISSN 0924-2716, <http://dx.doi.org/10.1016/j.isprsjprs.2012.01.005>.
  12. Yang Shao, Ross S. Lunetta, Comparison of support vector machine, neural network, and CART algorithms for the land-cover classification using limited training data points, *ISPRS Journal of Photogrammetry and Remote Sensing*, Volume70, June2012, Pages78-87, ISSN0924-2716, <http://dx.doi.org/10.1016/j.isprsjprs.2012.04.001>.
  13. Paul Bosch, Julio López, Héctor Ramírez, Hugo Robotham, Support vector machine under uncertainty: An application for hydroacoustic classification of fish-schools in Chile, *Expert Systems with Applications*, Volume 40, Issue 10, August 2013 ,Pages 4029-4034, ISSN0957-4174, <http://dx.doi.org/10.1016/j.eswa.2013.01.006>.
  14. Hongji Lin, Han Lin, Weibin Chen, Study on Recognition of Bird Species in Minjiang River Estuary Wetland, *Procedia Environmental Sciences*, Volume10, Part C, 2011, Pages 2478-2483, ISSN1878-0296, <http://dx.doi.org/10.1016/j.proenv.2011.09.386>.
  15. Rafael Pino-Mejías, María Dolores Cubiles-de-la-Vega, María Anaya-Romero, Antonio Pascual-Acosta, Antonio Jordán-López, Nicolás Bellinfante-Crocci, Predicting the potential habitat of oaks with data mining models and the R system, *Environmental Modelling & Software*, Volume 25 , Issue 7, July 2010, Pages 826-836, ISSN 1364-8152, <http://dx.doi.org/10.1016/j.envsoft.2010.01.004>.
  16. S.N.Jeyanthi, Efficient Classification Algorithms using SVMs for Large Datasets, A Project Report Submitted in partial fulfillment of the requirements for the Degree of Master of Technology in Computational Science, Supercomputer Education and Research Center, IISC, BANGALORE, INDIA, June 2007.
  17. Yugal kumar, G. Sahoo, Analysis of Parametric & Non Parametric Classifiers for Classification Technique using WEKA, *I.J. Information Technology and Computer Science*, 2012, 7,43-49 Published Online July 2012 in MECS DOI:10.5815/ijitcs.2012.07.06
  18. M.Arun Kumar, M. Gopal, A hybrid SVM based decision tree, *Pattern Recognition*, Volume 43, Issue 12, December 2010, Pages 3977-3987, ISSN0031-3203, <http://dx.doi.org/10.1016/j.patcog.2010.06.010>.
  19. Rajasekhar. N.; Babu, S.J.; Rajinikanth, T.V., "Magnetic resonance brain images classification using linear kernel based Support Vector Machine," *Engineering (NUICONe)*, 2012 Nirma University International Conference on , vol., no pp.1,5,68 Dec.2012 doi10.1109/NUICONe.2012.6493213
  20. V.F. Rodríguez-Galiano, F.Abarca-Hernández, B. Ghimire, M. Chica-Olmo, P.M .Akinson, C. Jeganathan, Incorporating Spatial Variability Measures in Land-cover Classification using Random Forest, *Procedia Environmental Sciences*, Volume3, 2011, Pages44-49, ISSN1878-0296, <http://dx.doi.org/10.1016/j.proenv.2011.02.009>.
  21. Reda M. Elbasiony, Elsayed A.Sallam, Tarek E. Eltobely, Mahmoud M. Fahmy, A hybrid network intrusion detection framework based on random forests and weighted k-means, *A in Shams Engineering Journal*, Available online 7 March 2013, ISSN 2090-4479, <http://dx.doi.org/10.1016/j.asej.2013.01003>.
  22. Gonzalo Martínez-Muñoz, Alberto Suárez, Using boosting to prune bagging ensembles, *Pattern Recognition Letters*, Volume28, Issue 1, 1 January 2007, Pages 156-165, ISSN0167-8655, <http://dx.doi.org/10.1016/j.patrec.2006.06.018>.
  23. Chun-Xia Zhang, Jiang-She Zhang, Gai-Ying Zhang, An efficient modified boosting method for solving classification problems, *Journal of Computational*

- and Applied Mathematics, Volume 214, Issue 2, 1 May 2008, Pages 381-392, ISSN 0377-0427, <http://dx.doi.org/10.1016/j.cam.2007.03.003>.
24. F. Löw, U. Michel, S. Dech, C. Conrad, Impact of feature selection on the accuracy and spatial uncertainty of per-field crop classification using Support Vector Machines, *ISPRS Journal of Photogrammetry and Remote Sensing*, Volume 85, November 2013, Pages 102-119, ISSN 0924-2716, <http://dx.doi.org/10.1016/j.isprsjprs.2013.08.007>.
  25. Liu Mingjun Wang, Jun Wang, Duo Li, Comparison of random forest, support vector machine and back propagation neural network for electronic tongue data classification: Application to the recognition of orange beverage and Chinese vinegar, *Sensors and Actuators B: Chemical*, Volume 177, February 2013, Pages 970-980, ISSN 0925-4005, <http://dx.doi.org/10.1016/j.snb.2012.11.071>.
  26. Ching-Chiang Yeh, Der-Jang Chi, Yi-Rong Lin, Going-concern prediction using hybrid random forests and rough set approach, *Information Sciences*, Available online 6 August 2013, ISSN 0020-0255, <http://dx.doi.org/10.1016/j.ins.2013.07.011>.
  27. Hsun-Jung Cho, Ming-Te Tseng, A support vector machine approach to CMOS-based radar signal processing for vehicle classification and speed estimation, *Mathematical and Computer Modelling*, Volume 58, Issues 1-2, July 2013, Pages 438-448, ISSN 0895-7177, <http://dx.doi.org/10.1016/j.mcm.2012.11.003>.
  28. Lam Hong Lee, Rajprasad Rajkumar, Lai Hung Lo, Chin Heng Wan, Dino Isa, Oil and gas pipeline failure prediction system using long range ultrasonic transducers and Euclidean-Support Vector Machines classification approach, *Expert Systems with Applications*, Volume 40, Issue 6, May 2013, Pages 1925-1934, ISSN 0957-4174, <http://dx.doi.org/10.1016/j.eswa.2012.10.006>.
  29. David Meyer, Support Vector Machines The Interface to lib svm in package e1071, Technische University of Wien, Austria, September, 2012.
  30. Xiaowei Yang, Qiaozhen Yu, Lifang He, Tengjiao Guo, The one-against-all partition based binary tree support vector machine algorithms for multi-class classification, *Neurocomputing*, Volume 113, 3 August 2013, Pages 1-7, ISSN 0925-2312, [10.1016/j.neucom.2012.12.048](http://dx.doi.org/10.1016/j.neucom.2012.12.048).
  31. Steve R. Gunn, Support Vector Machines for Classification and Regression, A Technical Report Faculty of Engineering, Science and Mathematics School of Electronics and Computer Science, University Of South Hampton, May 1998.
  32. Asdrúbal López Chau, Xiaou Li, Wen Yu, Convex and concave hulls for classification with support vector machine, *Neuro computing*, Volume 122, 25 December 2013, Pages 198-209, ISSN 0925-2312, <http://dx.doi.org/10.1016/j.neucom.2013.05.040>.
  33. Xinjun Peng, Yifei Wang, Dong Xu, Structural twin parametric-margin support vector machine for binary classification, *Knowledge-Based Systems*, Volume 49, September 2013, Pages 63-72, ISSN 0950-7051, <http://dx.doi.org/10.1016/j.knsys.2013.04.013>.
  34. Adnan Idris, Muhammad Rizwan, Asifullah Khan, Churn prediction in telecom using Random Forest and PSO based data balancing in combination with various feature selection strategies, *Computers & Electrical Engineering*, Volume 38, Issue 6, November 2012, Pages 1808-1819, ISSN 0045-7906, <http://dx.doi.org/10.1016/j.compeleceng.2012.09.001>.
  35. Elias Zintzaras, Axel Kowald, Forest classification trees and forest support vector machines algorithms: Demonstration using microarray data, *Computers in Biology and Medicine*, Volume 40, Issue 5, May 2010, Pages 519-524, ISSN 0010-4825, <http://dx.doi.org/10.1016/j.combiomed.2010.03.006>.
  36. Ching-hiang Yeh, Fengyi Lin, Chih-Yu Hsu, A hybrid KMV model, random forests and rough set theory approach for credit rating, *Knowledge-Based Systems*, Volume 33, September 2012, Pages 166-172, ISSN 0950-7051, <http://dx.doi.org/10.1016/j.knsys.2012.04.004>.
  37. Björn Waske, Sebastian vander Linden, Carsten Oldenburg, Benjamin Jakimow, Andreas Rabe, Patrick Hostert, imageRF – A user-oriented implementation for remote sensing image analysis with Random Forests, *Environmental Modeling & Software*, Volume 35, July 2012, Pages 192-193, ISSN 1364-8152, <http://dx.doi.org/10.1016/j.envsoft.2012.01.014>.
  38. Miao Liu, Mingjun Wang, Jun Wang, Duo Li, Comparison of random forest, support vector machine and back propagation neural network for electronic tongue data classification: Application to the recognition of orange beverage and Chinese vinegar, *Sensors and Actuators B: Chemical*, Volume 177, February 2013, Pages 970-980, ISSN 0925-4005, <http://dx.doi.org/10.1016/j.snb.2012.11.071>.
  39. Robin Genuer, Jean-Michel Poggi, Christine Tuleau-Malot, Variable selection using r random forests, *Pattern Recognition Letters*, Volume 31, Issue 14, 15 October 2010, Pages 2225-2236, ISSN 0167-8655, <http://dx.doi.org/10.1016/j.patrec.2010.03.014>.
  40. Katherine R. Gray, Paul Aljabar, Rolf A. Heckemann, Alexander Hammers, Daniel Rueckert, for the Alzheimer's Disease Neuroimaging Initiative, Random forest-based similarity measures for multimodal classification of Alzheimer's disease, *NeuroImage*, Volume 65, 15 January 2013, Pages 167-175, ISSN 10538119, <http://dx.doi.org/10.1016/j.neuroimage.2012.09.065>.
  41. Meng Xiao, Hong Yan, Jinzhong Song, Yuzhou Yang, Xianglin Yang, Sleep stages classification

- based on heart rate variability and random forest, *Biomedical Signal Processing and Control*, Volume 8, Issue 6, November 2013, Pages 624-633, ISSN 1746-8094, <http://dx.doi.org/10.1016/j.bspc.2013.06.001>.
42. Akin Özçift, Random forests ensemble classifier trained with data resampling strategy to improve cardiac arrhythmia diagnosis, *Computers in Biology and Medicine*, Volume 41, Issue 5, May 2011, Pages 265-271, ISSN 0010-4825, <http://dx.doi.org/10.1016/j.complbiomed.2011.03.001>.
  43. Simone Borra, Agostino Di Ciaccio, Improving nonparametric regression methods by bagging and boosting, *Computational Statistics & Data Analysis*, Volume 38, Issue 4, 28 February 2002, Pages 407-420, ISSN 0167-9473, [http://dx.doi.org/10.1016/S0167-9473\(01\)00068-8](http://dx.doi.org/10.1016/S0167-9473(01)00068-8).
  44. Imed Zitouni, Hong-Kwang Jeff Kuo, Chin-Hui Lee, Boosting and combination of classifiers for natural language call routing systems, *Speech Communication*, Volume 41, Issue 4, November 2003, Pages 647-661, ISSN 0167-6393, [http://dx.doi.org/10.1016/0167-6393\(03\)00103-1](http://dx.doi.org/10.1016/0167-6393(03)00103-1).
  45. Jafar Tanha, Maarten van Someren, Hamideh Afsarmanesh, Boosting for Multiclass Semi-Supervised Learning, *Pattern Recognition Letters*, Available online 21 October 2013, ISSN 0167-8655, <http://dx.doi.org/10.1016/j.patrec.2013.10.008>.
  46. Tong Xiao, Jingbo Zhu, Tongran Liu, Bagging and Boosting statistical machine translation systems, *Artificial Intelligence*, Volume 195, February 2013, Pages 496-527, ISSN 0004-3702, <http://dx.doi.org/10.1016/j.artint.2012.11.005>.
  47. Tansel Özyer, Reda Alhajj, Ken Barker, Intrusion detection by integrating boosting genetic fuzzy classifier and data mining criteria for rule pre-screening, *Journal of Network and Computer Applications*, Volume 30, Issue 1, January 2007, Pages 99-113, ISSN 1084-8045, <http://dx.doi.org/10.1016/j.jnca.2005.06.002>.
  48. Jakkrit TeCho, Cholwich Nattee, Thanaruk Theeramunkong, Boosting-based ensemble learning with penalty profiles for automatic Thai unknown word recognition, *Computers & Mathematics with Applications*, Volume 63, Issue 6, March 2012, Pages 1117-1134, ISSN 0898-1221, <http://dx.doi.org/10.1016/j.camwa.2011.11.062>.
  49. Chun-Xia Zhang, Jiang-She Zhang, A local boosting algorithm for solving classification problems, *Computational Statistics & Data Analysis*, Volume 52, Issue 4, 10 January 2008, Pages 1928-1941, ISSN 0167-9473, <http://dx.doi.org/10.1016/j.c.sda.2007.06.015>.
  50. Shaban Shataee, Holger Weinaker, Manoucher Babanejad, Plot-level Forest Volume Estimation Using Airborne Laser Scanner and TM Data, Comparison of Boosting and Random Forest Tree Regression Algorithms, *Procedia Environmental Sciences*, Volume 7, 2011, Pages 68-73, ISSN 1878-0296, <http://dx.doi.org/10.1016/j.proenv.2011.07.013>.
  51. Song feng Zheng, QBoost: Predicting quantiles with boosting for regression and binary classification, *Expert Systems with Applications*, Volume 39, Issue 2, 1 February 2012, Pages 1687-1697, ISSN 0957-4174, <http://dx.doi.org/10.1016/j.eswa.2011.06.060>.
  52. L.I. Kuncheva, M. Skurichina, R. P. W. Duin, An experimental study on diversity for bagging and boosting with linear classifiers, *Information Fusion*, Volume 3, Issue 4, December 2002, Pages 245-258, ISSN 1566-2535, [http://dx.doi.org/10.1016/S1566-2535\(02\)00093-3](http://dx.doi.org/10.1016/S1566-2535(02)00093-3).
  53. D.Lu & Q. Weng (2007): A survey of image classification methods and techniques for improving classification performance, *International Journal of Remote Sensing*, 28:5, 823-870, <http://dx.doi.org/10.1080/01431160600746456>.
  54. {R Core Team}, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, 2013, <http://www.Rproject.org>



This page is intentionally left blank

