

## Statistical Measures to Determine Optimal Structure of Decision Tree: One versus One Support Vector Machine

Manju Bala and R.K. Agrawal\*

*Jawaharlal Nehru University, New Delhi-110 067*  
*E-mail: manjugpm@gmail.com, \*rka@mail.jnu.ac.in*

### ABSTRACT

In this paper, one versus one optimal decision tree support vector machine (OvO-ODT SVM) framework is proposed to solve multi-class problems where the optimal structure of decision tree is determined using statistical measures, i.e., information gain, gini index, and chi-square. The performance of proposed OvO-ODT SVM is evaluated in terms of classification accuracy and computation time. It is also shown that proposed OvO-ODT SVM using all the three measures is more efficient in terms of time complexity for both training and testing phases in comparison to conventional OvO and support vector machine binary decision tree (SVM-BDT). Experiments on University of California, Irvine (UCI) repository dataset illustrates that ten cross-validation accuracy of our proposed framework is comparable or better in comparison to conventional OvO and SVM-BDT for most of the datasets. However, the proposed framework outperforms the conventional OvO and SVM-BDT for all the datasets in terms of both training and testing time.

**Keywords:** Support vector machine, decision tree, information gain, gini index, chi-square, SVM-BDT, binary decision tree

### 1. INTRODUCTION

Support vector machine (SVM) is most commonly used classifier for its better performance in comparison to other classifiers C4.5, artificial neural networks, bayesian classification in data mining and pattern recognition community. SVM based on statistical learning theory developed by Vapnik<sup>1</sup> was originally designed for binary classification. Most of the real-world applications, i.e., text categorisation<sup>2</sup>, optical character recognition<sup>3</sup>, speech recognition<sup>4</sup> are essentially multi-class classification problems. How to effectively extend SVM for multi-class classification is still an ongoing research issue<sup>5</sup>. The most common way to build a multi-class SVM is by constructing and combining several binary classifiers<sup>6</sup>. To solve multi-class classification problems, the whole problem is divided into a number of binary classification problems. The two representative ensemble schemes are one against all (OAA) and one versus one<sup>5</sup> (OvO). In literature, it has been stated that conventional OvO SVM approach has the problem of unclassifiable region. To resolve unclassifiable region for OvO SVM, Platt<sup>7</sup>, *et al.* proposed decision directed acyclic graph (DDAG) scheme. Madzarov<sup>8</sup>, *et al.* proposed support vector machine binary decision tree (SVM-BDT) architecture that uses binary SVMs for making decisions at each node of decision tree which takes advantage of both the efficient computation of the tree architecture and high accuracy of SVMs. The hierarchy of binary decision subtasks using SVMs is designed with clustering algorithms which employs Euclidean distance as similarity measure for creating

two disjoint groups of patterns. Euclidean similarity measure does not take into account within class variability of patterns. Hence, it may not be suitable for measuring class separability between two different classes of patterns. In literature, other than Euclidean distance, many similarity measures were used to determine the class separability, each of them associates with some advantages and disadvantages. Among them, few more realistic and effective statistical measures used in literature are information gain, gini index and chi-square for measuring class separability.

In this paper, construction of OvO-ODT SVM has been proposed, where class separability is determined using statistical measures. The performance of the proposed method is evaluated in terms of classification accuracy and computation time required for training and testing. It is also shown theoretically that OvO-ODT SVM is more efficient in terms of computation time for both training and testing phases in comparison to conventional OvO and SVM-BDT.

### 2. DECISION TREE SUPPORT VECTOR MACHINE

The SVM is based on the idea of structural risk minimisation which minimises the generalisation error<sup>1</sup>. It classifies data by determining a set of support vectors which are members of the set of training inputs that outline a hyperplane in feature space<sup>9</sup>. Consider  $n$  training samples  $\{(\bar{x}_i, y_i) \mid 1 \leq i \leq n\}$ , where  $\bar{x}_i$  represents  $i^{\text{th}}$  input feature vector and  $y_i$  is the corresponding target class. Given a set of  $n$  training instances, each represented as  $(\bar{x}_i, y_i)$ ,

the basic problem for training an SVM can be reformulated as:

$$\text{Maximise } J = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (x_i^T x_j) \quad (1)$$

subject to  $\sum_{i=1}^n y_i \alpha_i = 0$  and  $\alpha_i \geq 0$ , for  $1 \leq i \leq n$

The computation of dot products between vectors without explicitly mapping to another space is performed by a kernel function  $K(x_i, x_j)$ . Use of a kernel function<sup>10</sup> enables the curse of dimensionality to be addressed and the solution implicitly contains support vectors that provide a description of the significant data for classification. Substituting  $K(x_i, x_j)$  for  $(x_i^T x_j)$  in Eqn. (1) produces a new optimisation problem:

$$\text{Maximise } L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j K(x_i, x_j) \quad (2)$$

Subject to  $0 \leq \alpha_i \leq C$ ,  $i = 1, \dots, n$  and  $\sum_{i=1}^n \alpha_i y_i = 0$ , where  $C$  is trade-off variable.

Solving it for lagrangian multipliers  $\alpha_i$  gives a decision function of the form

$$f(X) = \text{sgn} \left( \sum_{i=1}^n y_i \alpha_i K(x_i, x_j) + b \right) \quad (3)$$

Conventional OvO SVM has the problem of unclassifiable region. To resolve unclassifiable region for OvO, Platt<sup>7</sup> et al. proposed DDAG approach. They have shown the existence of unclassifiable regions for  $k$ -class problem ( $k > 2$ ) which can lead to degradation of generalisation ability of classifier. In DDAG<sup>7</sup> scheme Vapnik Chernoff (VC) dimension, leave one out (LOO) error estimator, and Joachim's  $\hat{\alpha}$  LOO measures were used for estimating the generalisation ability of pairwise classifier at each level of decision tree. The measures used are computationally extensive and does not take into consideration any statistical information about the data. Madzarov<sup>8</sup>, et al. proposed SVM-BDT architecture. The hierarchy of binary decision subtasks using SVMs was designed with clustering algorithms. In SVM-BDT scheme, the classes are divided in two disjoint groups  $g_1$  and  $g_2$  using Euclidian distance as similarity measure. The two disjoint groups so obtained were then used to train a SVM classifier at the root node of the decision tree. The classes from the first and second clustering group were being assigned to left and right subtree, respectively. This process was continued recursively until there was only one class left in a group which defines a leaf in the decision tree.

### 3. PROPOSED OvO-ODT SVM FRAMEWORK USING STATISTICAL MEASURES

Euclidean similarity measure used in SVM-BDT does not take into account within class variability of patterns. Hence, it may not be suitable for measuring class separability between two different classes of patterns. To understand

better picture of the overlap of the subspaces occupied by individual classes, statistical measure is used in pattern recognition community which measures the overlap or probabilistic distance between the two class probability distribution functions (PDFs). Hence, it constitutes a natural concept of measuring class separability. Among statistical measures, information gain is one of the most commonly used measures to measure class separability of patterns in literature. Information Gain<sup>11</sup> (IG) is a measure based on entropy<sup>12</sup> which indicates degree of disorder of a system. It measures reduction in weighted average impurity of the partitions compared with the impurity of the complete set of samples when one knows the value of a specific attribute. Thus, the value of IG signifies how the whole system is related to an attribute. IG is calculated as

$$IG(C|E) = H(C) - H(C|E) \quad (4)$$

where  $IG(C|E)$  is the information gain of the label  $C$  for a given attribute  $E$ ,  $H(C)$  represents the system's entropy and  $H(C|E)$  is the system's relative entropy when the value of the attribute  $E$  is known. The system's entropy indicates its degree of disorder and is given by the following formula:

$$H(C) = - \sum_{i=1}^m p(C_i) \log p(C_i) \quad (5)$$

where  $p(C_i)$  is the probability of class  $i$ . The relative entropy is calculated as follows:

$$H(C|E) = \sum_{j=1}^{|E|} p(e_j) \left( - \sum_{i=1}^m p(C_i|e_j) \log p(C_i|e_j) \right) \quad (6)$$

where  $p(e_j)$  is the probability of value  $j$  for attribute  $e$ , and  $p(C_i|e_j)$  is the probability of class label  $C_i$  with regard to  $e_j$ .

In OvO-ODT SVM,  $k(k-1)/2$  independent binary SVM's are constructed for  $k$ -class problem. The optimal SVM model is selected on the basis of maximum value of IG that signifies more separability between patterns belonging to two different classes. IG for a given independent binary OvO SVM for training data containing  $n_i$  elements of class  $i$  and  $n_j$  elements of class  $j$  can be calculated as:

$$IG(i, j) = H(i, j) - \left[ \text{prob}(i)H(t_p, f_p) + \text{prob}(j)H(f_n, t_n) \right] \quad (7)$$

$$\text{where } H(x, y) = -x \log \left( \frac{x}{x+y} \right) - y \log \left( \frac{y}{x+y} \right) \quad (8)$$

$$\text{prob}(i) = \frac{n_i}{(n_i + n_j)} \quad \text{and} \quad \text{prob}(j) = \frac{n_j}{(n_i + n_j)} \quad (9)$$

$t_p, f_p, t_n$  and  $f_n$  denote true positive, false positive, true negative and false negative, respectively.

The Gini index is another popular measure for feature selection proposed by Breiman<sup>13</sup>, et al. It measures the

impurity of given set of training data  $D$  and can be calculated as:

$$Gini(D) = 1 - \sum_{i=1}^2 (prob(i))^2 \quad (10)$$

For a binary split, a weighted sum of the impurity of each resulting partition is computed. The reduction in impurity that would be incurred by a particular binary split is calculated as:

$$\Delta Gini(i, j) = Gini(D) - Gini_{i,j}(D) \quad (11)$$

where  $Gini_{i,j}(D) = [prob(i)Gini(L) + prob(j)Gini(R)]$  (12)

$Gini(L)$  and  $Gini(R)$  are the Gini index on the left and right side of the hyperplane. OvO SVM model  $(i_k, j_k)$  that maximises the reduction in impurity *i.e.* Gini index is selected as splitting node in decision tree SVM at a particular level.

Chi-square<sup>14</sup> similarity measure is another criterion used for binary split in data mining and machine learning. It is a statistical hypothesis test in which the sampling distribution of the test statistic is a chi-square distribution when the null hypothesis is true. In this case, the null hypothesis is that a random rule would place  $t_p$  patterns from class  $i$  and  $f_p$  tuples from class  $j$  independently in the left branch of decision tree and the remainder in the right branch of decision tree. The candidate decision rule would differ significantly from the random rule if the proportions differed significantly from those given by the random rule. The chi-square statistic  $\chi^2$  will be given by

$$\begin{aligned} \chi^2 = & g\left(t_p, (t_p + f_p)P_{pos}\right) + g\left(f_n, (f_n + t_n)P_{pos}\right) \\ & + g\left(f_p, (t_p + f_p)P_{neg}\right) + g\left(t_n, (t_n + f_n)P_{neg}\right) \end{aligned} \quad (13)$$

where  $g(count, expect) = \frac{(count - expect)^2}{expect}$ ,  $P_{pos}$  is the

probability of positive samples in left subtree and  $P_{neg}$  is the probability of negative samples in right subtree.

The higher the value of  $\chi^2$ , the less likely it is that the null hypothesis is true. Thus, for a sufficiently high  $\chi^2$ , one can reject the null hypothesis and can consider candidate rule is informative. Hence OvO SVM model  $(i_k, j_k)$  that maximises  $\chi^2$  is selected as splitting node in decision tree SVM at a particular level.

The outline for decision tree SVM using IG class separability measure for  $k$ -class is given below:

Input: Training data  $X = \bigcup_i^n x_i$ . The following steps are

involved

- Step 1. Generate the initial list  $\{C_1, \dots, C_k\}$
- Step 2. Calculate  $H(i, j)$  using Eqn (8) for  $i = 1, \dots, k$  and  $j > i$
- Step 3. Calculate,  $H(t_p, f_p)$ ,  $H(f_n, t_n)$ ,  $prob(i)$  and  $prob(j)$  using Eqn (8) and Eqn (9), respectively.
- Step 4. Compute  $IG(i, j)$  using Eqn (7).
- Step 5. Determine class pair  $(C_p, C_j)$  for which  $IG(i, j)$  takes

maximum value from the list. If  $x_k, k=1, \dots, n$  belongs to class  $C_i$  then delete  $C_j$  from the list else delete class  $C_i$ .

Step 6. If the remaining classes are  $>2$ , repeat Steps 2-5 otherwise terminate the algorithm.

Similar computational steps are followed for other two measures to determine the structure of OvO-ODT SVM. To compute the time complexity of training phase, assume without loss of any generality that the number of samples in each class is approximately the same, *i.e.*,  $n/k$ . To solve  $k$ -class problem using conventional OvO,  $k(k-1)/2$  binary SVM classifiers are developed. Assuming the time complexity of building a SVM with  $n$  samples and  $d$  features is  $O(n^2d)$ , it can be shown that training time of conventional OvO,  $T_{con-OvO-SVM}^{train}$ , is  $O(n^2d)$ . In worst case, the decision tree generated in SVM-BDT is skewed if classes in two groups at every level are divided into uneven size. Under the assumption that group  $g_1$  contains only one class and group  $g_2$  contains remaining classes, the decision tree so generated will be of depth  $(k-1)$  for  $k$ -class problem. Hence, the training time of SVM-BDT will be given by

$$\begin{aligned} T_{SVM-BDT-worst}^{train} &= \left( n^2d + \left( \frac{n(k-1)}{k} \right)^2 d + \dots + \left( \frac{2n}{k} \right)^2 d \right) \\ &\cong (n^2dk) \end{aligned} \quad (14)$$

In SVM-BDT approach under best case, the class in two groups at every level is divided into approximately the same size. The decision tree so generated will be almost height balanced of maximum depth  $[\log(k)]$ . The number of nodes in decision tree at depth  $i$  is  $2^{i-1}$ , each containing  $n/2^{i-1}$  samples. Hence, the training time for SVM-BDT in best case is given by

$$\begin{aligned} T_{SVM-BDT-best}^{train} &= \left( n^2d + 2 \left( \frac{n}{2} \right)^2 d + 4 \left( \frac{n}{2} \right)^2 d + \dots \right. \\ &\quad \left. + 2^{\log(k)-1} \left( \frac{n}{2^{\log(k)-1}} \right)^2 d \right) \cong (n^2d) \end{aligned} \quad (15)$$

However in general, the structure of decision tree generated using statistical measures is almost height balanced of maximum depth  $[\log(k)]$ . There are  $2^{i-1}$  nodes at  $i^{\text{th}}$  level and each node uses  $(2n/k)$  samples. Hence, the training time for OvO-ODT SVM using statistical measure is given by

$$T_{OvO-ODT-SVM}^{train} \cong \sum_{i=1}^{\log(k)} 2^{i-1} \left( \frac{2n}{k} \right)^2 d \cong n^2 \left( \frac{d}{k} \right) \quad (16)$$

During testing phase of the conventional OvO,  $k(k-1)/2$  decision functions are to be evaluated. Also, the majority voting is computed with  $k(k-1)/2$  operation. Hence, the testing time  $T_{Con-OvO}^{test}$  for each sample is given by  $k(k-1)/2$ . In worst case, the depth of SVM-BDT is  $(k-1)$ . The testing time  $T_{SVM-BDT-worst}^{test}$  for each sample in

worst case is  $k$ . However in best case the depth of SVM-BDT is  $\lceil \log(k) \rceil$  which requires  $\lceil \log(k) \rceil$  testing time  $T_{SVM-BDT-best}^{test}$  for each sample. Since, the maximum depth of OvO-ODT SVM is  $\lceil \log(k) \rceil$ , the testing time requires  $\lceil \log(k) \rceil$  operations. According to the above analysis, it is evident that the training and testing time for OvO-ODT SVM will always require less computation time in comparison to conventional OvO and SVM-BDT.

#### 4. EXPERIMENTAL RESULTS

To evaluate the performance of the proposed OvO-ODT SVM framework using information gain, gini index and chi-square. They have performed experiments on publically available UCI<sup>15</sup> benchmark datasets. Table 1 describes the datasets used in the experiments. All the experiments were performed on a computer having Pentium 4 dual-core processor with 1GB RAM. The kernel functions used in experiments are given in Table 2.

The classification accuracy is determined using ten cross-validations. The authors have used  $\gamma = [2^{-11}, 2^{-10}, 2^{-9}, \dots, 2^0]$  and trade-off variable  $C=1000$ . For a given kernel function, the authors determined the value of  $\gamma$  for which the maximum classification accuracy was achieved. Table 3 shows the comparison of maximum classification accuracy between conventional OvO SVM and SVM-BDT with the

**Table 1. Description of datasets**

Problem	#Train_data	# Class	#Attributes
Wine	178	3	13
Vehicle	846	4	18
Glass	214	6	9
Segmentation	210	7	19
Ecoli	336	8	7

**Table 2. Kernel functions**

Kernel Function	$K(x, x_i)$ for $\gamma > 0$
Gaussian	$exp\left(-\gamma  x - x_i ^2\right)$
Laplace	$exp\left(-\gamma  x - x_i \right)$
Cauchy	$\left(1/\left(1 + \gamma  x - x_i ^2\right)\right)$
Hypersecant	$2/\left(exp\left(\gamma  x - x_i \right) + exp\left(-\gamma  x - x_i \right)\right)$
Square sync	$sin^2\left(\gamma  x - x_i \right)/\left(\gamma  x - x_i ^2\right)$

**Table 3. Comparison of classification accuracy (per cent)**

Dataset	Kernel Choice	OvO	SVM-BDT	OvO-ODT SVM		
				$\chi^2$	Gini	IG
Wine	Gaussian	82.58	81.88	82.58	82.58	<b>83.15</b>
	Laplace	<b>82.58</b>	81.76	<b>82.58</b>	<b>82.58</b>	<b>82.58</b>
	Cauchy	<b>82.02</b>	81.76	<b>82.02</b>	<b>82.02</b>	<b>82.02</b>
	Hypersecant	<b>93.26</b>	92.63	<b>93.26</b>	<b>93.26</b>	92.13
	Square sync	75.28	75.46	<b>76.97</b>	<b>76.97</b>	75.28
Vehicle	Gaussian	<b>76.83</b>	75.24	76.60	76.60	76.60
	Laplace	<b>77.42</b>	76.24	<b>77.42</b>	<b>77.42</b>	<b>77.42</b>
	Cauchy	<b>76.48</b>	74.65	<b>76.48</b>	<b>76.48</b>	<b>76.48</b>
	Hypersecant	<b>83.33</b>	82.98	83.22	<b>83.33</b>	<b>83.33</b>
	Square sync	<b>71.51</b>	71.08	70.922	70.80	70.8
Glass	Gaussian	72.43	65.42	70.56	<b>72.90</b>	71.03
	Laplace	75.70	68.69	<b>76.17</b>	<b>76.17</b>	75.24
	Cauchy	72.90	69.42	71.03	<b>73.36</b>	71.03
	Hypersecant	<b>71.96</b>	69.47	71.09	<b>71.96</b>	70.09
	Square sync	<b>66.36</b>	58.41	62.62	58.13	62.62
Segmentation	Gaussian	84.76	82.89	<b>85.24</b>	<b>85.24</b>	84.29
	Laplace	87.14	87.05	<b>87.62</b>	<b>87.62</b>	87.14
	Cauchy	<b>86.19</b>	85.05	<b>86.19</b>	<b>86.19</b>	85.71
	Hypersecant	90	90	<b>90.95</b>	<b>90.95</b>	89.05
	Square sync	<b>81.9</b>	73.80	79.52	79.05	80.95
Ecoli	Gaussian	<b>85.42</b>	81.34	<b>85.42</b>	<b>85.42</b>	<b>85.42</b>
	Laplace	<b>87.20</b>	85.78	<b>87.20</b>	<b>87.20</b>	86.91
	Cauchy	<b>85.42</b>	79.90	<b>85.42</b>	<b>85.42</b>	84.82
	Hypersecant	<b>85.42</b>	83.89	83.33	<b>85.42</b>	84.82
	Square sync	85.12	82.78	<b>87.20</b>	82.85	85.12

**Table 4. Comparison of training and testing time (Seconds)**

Dataset	Kernel choice	Training time (s)					Testing time (s)				
		OvO	SVM-BDT	OvO-ODT SVM			OvO	SVM-BDT	OvO-ODT SVM		
				$\chi^2$	Gini	IG			$\chi^2$	Gini	IG
Wine	Gaussian	2.15	2.15	<b>2.05</b>	2.13	2.11	0.08	0.04	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>
	Laplace	2.98	2.83	2.67	2.18	<b>1.48</b>	0.07	0.04	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>
	Cauchy	2.89	1.98	1.69	<b>1.34</b>	1.54	0.08	0.08	0.06	0.07	<b>0.03</b>
	Hypersecant	3.55	3.22	2.15	<b>2.13</b>	2.29	0.06	0.05	0.03	0.06	<b>0.02</b>
	Square sync	3.13	2.93	2.27	<b>2.18</b>	2.60	0.07	0.03	0.03	<b>0.02</b>	0.03
Vehicle	Gaussian	189.96	185.45	<b>155.43</b>	170.92	155.79	1.85	0.56	<b>0.46</b>	0.53	<b>0.46</b>
	Laplace	188.68	187.55	<b>134.34</b>	134.54	134.54	1.45	0.57	0.56	<b>0.35</b>	0.57
	Cauchy	186.44	186.57	144.34	144.54	<b>144.23</b>	1.73	0.73	0.46	<b>0.45</b>	0.73
	Hypersecant	198.69	181.75	<b>145.45</b>	148.78	145.75	1.64	0.76	0.56	<b>0.46</b>	0.76
	Square sync	189.75	189.57	159.72	159.65	<b>149.77</b>	1.63	0.66	<b>0.53</b>	0.59	0.66
Glass	Gaussian	12.31	11.88	<b>8.99</b>	10.06	9.05	1.20	0.32	<b>0.22</b>	0.28	<b>0.22</b>
	Laplace	11.90	10.53	9.84	5.45	<b>4.80</b>	1.21	0.66	<b>0.14</b>	0.18	0.18
	Cauchy	13.62	11.78	<b>7.62</b>	9.56	10.71	1.43	0.89	<b>0.15</b>	0.17	0.17
	Hypersecant	8.36	7.35	6.36	7.46	<b>5.93</b>	1.21	0.79	<b>0.15</b>	0.18	0.19
	Square sync	8.93	8.04	5.93	<b>4.92</b>	5.42	1.25	0.99	<b>0.11</b>	0.14	0.19
Segmentation	Gaussian	3.00	2.98	<b>2.39</b>	2.45	2.49	1.40	0.55	<b>0.42</b>	0.54	0.44
	Laplace	3.15	2.90	<b>2.34</b>	2.52	2.42	1.35	0.89	0.53	0.43	<b>0.33</b>
	Cauchy	3.10	2.89	2.13	2.12	<b>2.10</b>	1.78	0.76	0.34	<b>0.32</b>	0.33
	Hypersecant	4.91	4.89	4.57	4.56	<b>4.21</b>	1.24	0.54	0.36	<b>0.31</b>	<b>0.31</b>
	Square sync	3.99	3.81	3.61	3.21	<b>3.03</b>	1.54	0.64	0.38	<b>0.32</b>	0.35
Ecoli	Gaussian	26.48	25.99	<b>17.55</b>	21.97	20.07	1.55	0.88	<b>0.54</b>	0.60	0.96
	Laplace	24.38	23.45	17.34	<b>16.87</b>	17.34	1.82	1.78	1.05	1.12	<b>1.02</b>
	Cauchy	27.85	25.83	21.85	<b>19.85</b>	24.85	2.31	1.20	0.99	0.98	<b>0.92</b>
	Hypersecant	20.98	20.23	19.78	<b>16.78</b>	19.78	1.82	1.15	<b>1.03</b>	1.10	1.04
	Square sync	23.65	21.94	20.85	<b>20.35</b>	20.85	1.88	1.33	1.04	<b>1.09</b>	1.03

proposed OvO-ODT SVM using all the three statistical measures. The best classification accuracy for each dataset is shown in bold. From Table 3 it can be observed that the classification accuracy of the proposed OvO-ODT SVM method is better or equivalent to conventional OvO SVM and SVM-BDT for different choices of kernel functions for most of datasets using all three statistical measures. Among the three measures used for determining the structure of decision tree, neither of them is clear winner over the others.

Table 4 shows the computation time for training and testing phase for all kernels with  $\gamma = 2^{-11}$  and trade-off variable  $C=1000$ . The minimum training and testing time for each dataset is shown in bold. It can be observed from Table 4 that time required for training OvO-ODT SVM using all the three statistical measures is less in comparison to conventional OvO and SVM-BDT. Among the three measures used for building OvO-ODT SVM, neither of them is a clear winner over the other measures. It can also be observed that the proposed OvO-ODT SVM is also computationally better than conventional OvO and SVM-BDT in terms of testing time for all the three measures.

## 5. CONCLUSIONS

In this paper, OvO-ODT SVM framework is proposed to solve multi-class problems where the optimal structure of decision tree is determined using statistical measures, i.e., information gain, gini index, and chi-square. It has also been shown theoretically that the computation time of training and testing of OvO-ODT SVM using statistical measures is better in comparison to conventional OvO SVM and SVM-BDT. The performance of the proposed framework is evaluated in terms of classification accuracy and computation time of training and testing phases. Experimental results on UCI repository dataset demonstrate better or equivalent performance of the proposed OvO-ODT SVM scheme in comparison to conventional OvO SVM and SVM-BDT in terms of classification accuracy for most of the datasets using all the three measures employed for construction of decision tree. The experimental results obtained also demonstrate that the training and testing time of proposed OvO-ODT SVM takes less computation time in comparison to conventional OvO SVM and SVM-BDT.

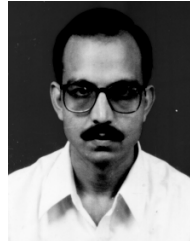
## REFERENCES

1. Vapnik, V.N. Statistical Learning Theory. New York, John Wiley & Sons, 1998.
2. Joachim's, T.; Cristianini, N.; & Taylor, J. S. Composite Kernels for hypertext categorisation. *In* The International Conference on Machine Learning, 2001.
3. Mori, S.; Suen, C.Y. & Yamamoto, K. Historical review of OCR research and development. *Proceedings IEEE*, 1992, **80**, 1029-058.
4. Schmidt, M. Identifying speaker with support vector networks. *In* Interface '96 Proceedings, 1996, Sydney.
5. Rifkin, R. & Klautau, A. In defence of one-vs-all classification. *J. Machine Learning*, 2004, **5**, 101-41.
6. Hsu, C.W. & Lin, C.J. A comparison of methods for multi-class support vector machine. *IEEE Trans. Neural Networks*, 2002, **13**(2), 415-25.
7. Platt, Cristianini N. & Shawe-Taylor, J. Large margin DAGSVM's for multi-class classification. *Adv. Neural Info. Proc. Syst.*, 2000, **12**, 547-53.
8. Gjorgji, M.; Dejan, G. & Ivan, C. A multi-class SVM classifier utilising binary decision tree. *Informatica*, 2009, **33**, 233-41.
9. Kittler, J. & Hojjatoleslami, A. A weighted combination of classifiers employing shared and distinct representations. *IEEE Comp. Vision Pattern Recog.*, 1998, 924-29.
10. Scholkopf, B. & Smola, A. Learning with kernels. Cambridge, MA, MIT Press, 2002.
11. Bala, M. & Agrawal, R. K. Evaluation of decision tree SVM framework using different statistical measures. *ARTcom*, 2009, 341-45.
12. Quinlan, J.R. Induction of decision trees. *Machine Learning*, 1986, **1**, 81-106.
13. Breiman, L.; Friedman, J.; Ohlsen, R. & Stone, C. Classification and regression trees. Belmont, CA, Wadsworth, 1984.
14. Duda, R. O.; Hart, P. E. & Stork, D. E. Pattern classification, Ed 2. John Wiley & Sons, 2000.
15. [www.ics.uci.edu/~mllearn/MLRepository.html](http://www.ics.uci.edu/~mllearn/MLRepository.html).

## Contributors



**Ms Manju Bala** doing her PhD at the School of Computer and Systems Sciences, Jawaharlal Nehru University, New Delhi. She received BE(Comp. Sci. and Engg.) from Maharishi Dayanand University and MTech (Comp. Sci.) from IASE Deemed University, Rajasthan. Presently, she is working as an Assistant Professor in Computer Science department, I.P. College for Women, University of Delhi, New Delhi. Her current area of research is pattern recognition.



**Mr Ramesh Kumar Agrawal** obtained MTech (Computer Application) from Indian Institute of Technology Delhi, New Delhi. He has done his PhD (Computational Physics) from Delhi University. Presently, he is working as an Associate Professor at the School of Computer and Systems Sciences, Jawaharlal Nehru University, New Delhi. His current areas of research are: Classification, feature extraction and selection for pattern recognition problems in domains of image processing, security, and bioinformatics.