

New Clustering Algorithms for the Support Vector Machine Based Hierarchical Classification

Hakan Cevikalp

*Electrical and Electronics Engineering Department, Eskisehir Osmangazi University,
Meselik 26480 Eskisehir, Turkey.
Email:hakan.cevikalp@gmail.com*

Abstract

This study presents two new clustering algorithms for partition of data samples for the Support Vector Machine (SVM) based hierarchical classification. A divisive (top-down) approach is considered in which a set of classes is automatically separated into two smaller groups at each node of the hierarchy. The first algorithm splits the data samples based on a variation of the Normalized Cuts (NCuts) clustering algorithm wherein the weights of adjacency matrix are modified to utilize class membership in the process. The second algorithm also uses the NCuts clustering; however, it considers the involved classes rather than the individual data samples. It uses the minimum distances between the convex hulls of classes as a distance measure for determining the weights of the graph. Splits are determined for both algorithms based on the eigenvector corresponding to the second smallest eigenvalue of a Laplacian matrix, and it is observed that the proposed algorithms generate well-separated and well-balanced clusters. Unlike other clustering methods used for this purpose, the methods in the present study are found to be more suitable when SVMs are used as base classifiers. As demonstrated in

24 the experiments, the proposed clustering algorithms are integrated into the
25 hierarchical SVM classifiers, which results in significantly improved testing
26 times with a negligible decrease in classification accuracies as compared to
27 the traditional multi-class SVMs.

28 *Key words:* hierarchical classification, support vector machines, multi-class
29 classification, clustering, normalized cuts

30 **1. Introduction**

31 Automatic classification of data samples is very important for several
32 applications including visual object classification, text classification, speech
33 recognition, etc. This leads to a requirement for efficient and accurate clas-
34 sifiers. For example, in visual object localization problems, hundreds of sub-
35 windows need to be classified. This task will become tedious if the chosen
36 classifier is not fast. It is known that there is a direct relationship between
37 the real-time performance of a classifier and the number of classes. Human
38 beings can classify between 10^4 and 10^5 object categories, and therefore,
39 this seems to be a practical goal for machines as well (Griffin and Perona,
40 2008). It is thus crucial that classification algorithms must scale well with
41 the number of classes.

42 The Support Vector Machine (SVM) classifier is a successful method that
43 simultaneously minimizes the empirical classification error and maximizes the
44 geometric margin (Cortes and Vapnik, 1995; Burges, 1998; Schölkopf and
45 Smola, 2002). Essentially, it finds a separating hyperplane that yields the
46 largest margin (separation gap) between two class samples. The SVM for-
47 mulation was originally designed for binary classification; however, extend-

48 ing this formulation to more than two classes makes it very complex and is
49 therefore generally avoided. Yet, many classification applications have more
50 than two classes. The multi-class SVM problems are dealt with construct-
51 ing several binary classifiers. There are various strategies to achieve this
52 goal. Among these, the earliest and the most popular are the one-against-
53 rest (OAR) strategy and the one-against-one (OAO) strategy (Hsu and Lin,
54 2002). For a C -class classification problem, the former strategy trains C bi-
55 nary classifiers, in which each classifier separates one class from the remaining
56 $C - 1$ classes. All classifiers are trained on the entire training set, and the
57 class label of a test sample is determined based on the highest output value
58 of the classifier in the ensemble. The latter strategy constructs all possible
59 $C(C - 1)/2$ binary classifiers out of C classes. The decision of the ensemble
60 is typically made using the max wins algorithm: Each OAO classifier casts
61 one vote for its preferred class, and the final decision is made for the class
62 with the most votes. The OAO strategy builds more classifiers than the OAR
63 strategy, and in general, it is considerably faster than OAR during training
64 since it operates on a smaller number of training samples (Hsu and Lin, 2002).
65 However, the OAO classifiers grow in size quadratically with the number of
66 classes, which makes the OAO strategy very expensive for applications with
67 large number of classes.

68 Recently, SVM based hierarchical classifiers have gained significant at-
69 tention for large class problems owing to their capability to scale well with
70 the number of classes (Platt et al., 2000; Vural and Dy, 2004; Casasent and
71 Wang, 2005; Marszalek and Schmid, 2008; Chen et al., 2004; Griffin and
72 Perona, 2008; Zhigang et al., 2005; Liu et al., 2005). Two popular methods

73 fall into this category, namely, Decision Directed Acyclic Graphs (DDAGs)
74 (Platt et al., 2000) and Binary Hierarchical Decision Trees (BHDTs) (Vural
75 and Dy, 2004; Casasent and Wang, 2005; Chen et al., 2004; Griffin and Per-
76 ona, 2008; Zhigang et al., 2005). The DDAG method first trains $C(C - 1)/2$
77 binary classifiers, and subsequently uses a Directed Acyclic Graph (DAG)
78 during the testing phase. This is equivalent to operating on a list where
79 each node of the DAG eliminates one class from the list. Thus, the method
80 requires only $C - 1$ decision nodes to be evaluated for labeling a test sample
81 rather than $C(C - 1)/2$ classifier evaluations, which results in a significant
82 speeding up of the testing phase. However, this algorithm makes some un-
83 necessary comparisons which are considered as irrelevant for the classification
84 of a particular test sample. As an example, consider a visual object classi-
85 fication problem: When a test sample belonging to the ‘dogs’ class arrives,
86 it is rather unnecessary to make comparisons between unrelated classes such
87 as buildings, airplanes, and cars. In this way, the real-time performance
88 could be further improved. BHDT algorithms have been introduced in or-
89 der to improve the efficiency of SVM classifiers by reducing the unnecessary
90 comparisons while maintaining the high classification accuracy. To reduce
91 unnecessary class comparisons, a BHDT algorithm uses a decision tree that
92 divides the data hierarchically into two subsets until each subset consists of
93 only one class. The SVM classifier is then used for separating those sub-
94 sets at each node of the binary tree. The data partition is often achieved
95 using a clustering algorithm, and the accuracy of the SVM classifier at each
96 internal node depends on the generated clusters. Different hierarchy strate-
97 gies (top-down (Vural and Dy, 2004; Casasent and Wang, 2005; Chen et

98 al., 2004; Griffin and Perona, 2008) and bottom-up (Zhigang et al., 2005))
99 and different clustering algorithms, such as k-means (Vural and Dy, 2004;
100 Zhigang et al., 2005), kernel k-means (Casasent and Wang, 2005), spherical
101 shells (Vural and Dy, 2004) and balanced subset clustering (Vural and Dy,
102 2004), have been used in the literature. It is known that BHDT algorithms
103 employ various distance measures for partition such as the Euclidean dis-
104 tance between class means (Vural and Dy, 2004; Zhigang et al., 2005), the
105 Kullback-Leibler distance between class densities (Chen et al., 2004), or the
106 number of misclassifications between classes (Griffin and Perona, 2008). In
107 addition to these algorithms, some BHDT methods determine the partitions
108 based on the clustering of data samples rather than the class sets (Marsza-
109 lek and Schmid, 2008). A well-balanced tree requires approximately $\log_2 C$
110 classifier evaluations for traversing a path from the top to a bottom decision
111 node. This results in a more efficient structure as compared to the DDAG
112 algorithm in terms of testing time.

113 The present study is focused on SVM based BHDTs wherein two clus-
114 tering algorithms are proposed for the partition of classes. Similar to other
115 BHDT algorithms, the main objective is to improve the real-time efficiency
116 (testing time) while maintaining the high classification accuracy. The first
117 clustering algorithm operates on data samples whereas the second algorithm
118 considers the class sets. It is found that both methods yield well-separated
119 and well-balanced partitions, which is compatible with the goal of SVM clas-
120 sification. The remaining sections of this article are organized as follows:
121 In section 2, the proposed clustering methods are introduced. In section
122 3, the data sets and the experimental procedure are described. Lastly, the

123 conclusion is provided in Section 4.

124 **2. Method**

125 *2.1. Design Issues*

126 BHDT methods use clustering algorithms for the partition of data; thus,
127 the classification accuracy and computational efficiency of the hierarchical
128 classification system depend heavily on the generated clusters. More pre-
129 cisely, well-balanced separable clusters at each node of the tree would signifi-
130 cantly improve the performance of the overall system. This requires that the
131 employed clustering algorithm must be compatible with the base classifier,
132 that is, the SVM classifier in our study.

133 To design optimal clustering algorithm, we should first examine the base
134 classifier SVM since the clustering algorithm as well as the base classifier
135 must aim achieving the same goal for a satisfactory performance. The SVM
136 classifier finds a separating hyperplane that maximizes the margin, which is
137 defined as the distance between the hyperplane and the closest samples from
138 the classes. To do so, SVM first approximates each class with a convex hull
139 (Bennett and Bredensteiner, 2000). A convex hull consists of all points that
140 can be written as a convex combination of the points in the original set, and
141 a convex combination of points is a linear combination of data points where
142 all coefficients are nonnegative and sum up to 1. More formally, the convex
143 hull of samples $\{\mathbf{x}_i\}_{i=1,\dots,n}$ can be written as

$$H_{convex} = \left\{ \mathbf{x} = \sum_{i=1}^n \alpha_i \mathbf{x}_i \mid \sum_{i=1}^n \alpha_i = 1, \alpha_i \geq 0 \right\}. \quad (1)$$

144 Convex hulls of two classes are illustrated in Fig. 1. Following this approx-
 145 imation, SVM finds the closest points in these convex hulls (Bennett and
 146 Bredensteiner, 2000). Then, these two points are connected with a line seg-
 147 ment. The plane, orthogonal to the line segment that bisects the line, is
 148 selected as the separating hyperplane as shown in Fig. 1. From this geomet-
 149 rical point of view, in a separable case, the two closest points on the convex
 150 hulls determine the separating hyperplane and the SVM margin is merely
 151 equivalent to the minimum distance between the convex hulls that represent
 classes.

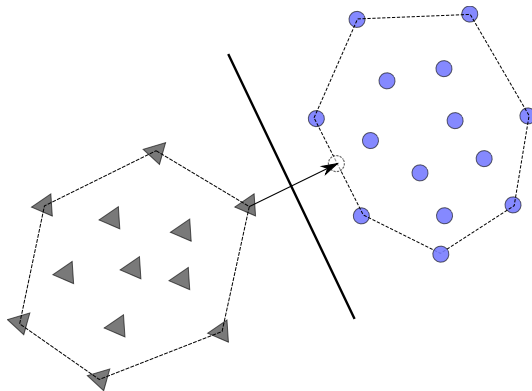


Figure 1: Two closest points on the convex hulls determine the separating hyperplane of the hard-margin SVM classifier.

152

153 The clustering algorithms based on k-means clustering are not good choices
 154 for SVM based BHDTs since they do not directly target at margin maxi-
 155 mization in the sense described above. In a similar manner, the clustering
 156 algorithm maximizing the Kullback-Leibler distance between class densities
 157 is also not compatible with SVMs owing to the same reason. When the clus-
 158 tering algorithms using k-means or Kullback-Leibler distance are employed,

159 it may be more difficult to separate the resulting clusters for SVM classifiers,
160 which in turn would degrade both the classification accuracy and the real-
161 time efficiency of the overall classification system. On the other hand, the
162 Normalized Cuts (NCuts) clustering of (Shi and Malik, 2000) is more suitable
163 for partition of data since its objective is similar to that of SVMs. The NCuts
164 clustering algorithm maps the data samples into an infinite-dimensional fea-
165 ture space and cuts through the data by passing a hyperplane through the
166 maximum gap in the mapped data (Rahimi and Recht, 2004). It then labels
167 points that fall on the same side of the hyperplane as belonging to the same
168 cluster. However, NCuts is an unsupervised approach and its split does not
169 guarantee that samples belonging to the same classes are always grouped
170 together in the same node unless class-specific samples are very close to each
171 other and they are far from the other class samples. As a result, NCuts clus-
172 tering may create overlapping classes that have samples in different clusters.
173 If the created clusters are not compact, the real-time performance of the
174 BHDTs degrades and they behave like K-D trees (Friedman et al., 1977). In
175 the following section, the NCuts algorithm is modified and two variations are
176 proposed that generate well-balanced compact clusters with the maximum
177 margin.

178 *2.2. Sample Based Large Margin Clustering*

179 Let the training samples be $\mathbf{x}_{ci} \in \mathbb{R}^d$, where $c = 1, \dots, C$ indexes the
180 C classes and $i = 1, \dots, n_c$ indexes the n_c samples of class c . Prior to
181 introduction of the first proposed method, a summary of the NCuts clustering
182 algorithm is provided since the proposed method is built on this approach.

183 *2.2.1. Normalized Cuts Clustering*

184 Given a dataset of m samples, the NCuts algorithm constructs a weighted
 185 graph with m vertices $\{v_1, \dots, v_m\}$ (one for each sample) and a set of edges
 186 containing these vertices. Each edge between vertices v_i and v_j carries a non-
 187 negative weight $w_{ij} = w_{ji} \geq 0$ based on the similarity between the samples
 188 associated to the vertices. In the present study, a fully connected graph is
 189 considered in which all edges are connected. A common choice for weighting
 190 the edges, is the heat kernel (Gaussian kernel) $w_{ij} = \exp(-d(\mathbf{x}_i, \mathbf{x}_j)^2/t)$,
 191 where $d(\mathbf{x}_i, \mathbf{x}_j)$ represents the distance between samples \mathbf{x}_i and \mathbf{x}_j computed
 192 using the preferred distance function, and t is the width of the kernel.

193 The matrix $\mathbf{W} = (w_{ij})_{i,j=1,\dots,m}$ is the weighted adjacency matrix of the
 194 graph. In the case of binary clustering, assigning a label $y_i \in \{-1, +1\}$ to
 195 each sample \mathbf{x}_i cuts the graph into set A of the vertices with label $+1$ and
 196 set B of vertices with label -1 . The cost function of the method is defined
 197 as

$$NCut(A, B) = \left(\frac{1}{Vol(A)} + \frac{1}{Vol(B)} \right) \sum_{i \in A, j \in B} w_{ij} \quad (2)$$

198 where Vol is the sum of the weights in a set and $\sum_{i \in A, j \in B} w_{ij}$ is the total
 199 weight of the edges that must be removed to make A and B disjoint. This
 200 cost function penalizes the cuts that are not well-balanced and ensures that
 201 the sets A and B have approximately the same number of elements (Shi and
 202 Malik, 2000). However, optimizing the above criterion is NP hard. Thus,
 203 by resorting a relaxation, the problem is reduced to minimization of the
 204 Laplacian of the graph. If the Laplacian matrix is denoted with $\mathbf{L} = \mathbf{D} -$
 205 \mathbf{W} , where \mathbf{D} is the diagonal matrix whose entries are the column (or row)
 206 sums of \mathbf{W} , the algorithm would consist of solving the following generalized

207 eigenvalue problem

$$\mathbf{L} = \lambda \mathbf{D} \mathbf{a}. \quad (3)$$

208 Subsequently, the components of the eigenvector \mathbf{a}^* corresponding to the
209 second smallest eigenvalue of (3) are thresholded to split data into two sets,
210 i.e.,

$$\begin{cases} y_i = -1 & \text{if } a_i^* > \Delta, \\ y_i = +1 & \text{if } a_i^* \leq \Delta, \end{cases} \quad (4)$$

211 where Δ is the chosen threshold which is typically equal to zero.

212 *2.2.2. Modified Normalized Cuts Clustering*

213 As mentioned earlier, NCuts clustering is an unsupervised method which
214 does not take class membership information into consideration. Thus, the
215 resulting clusters may have overlapping classes. To overcome this pitfall
216 and reduce the overlapping regions among classes, the following similarity
217 function is adopted to weight the edges

$$w_{ij} = \begin{cases} \exp(-d(\mathbf{x}_{ci}, \mathbf{x}_{cj})/\alpha t) & \text{if } c = \acute{c}, \\ \exp(-d(\mathbf{x}_{ci}, \mathbf{x}_{cj})/t) & \text{if } c \neq \acute{c}, \end{cases} \quad (5)$$

218 where $\alpha \geq 1$ is a tuning parameter that is used to change the similarities
219 between class-specific samples. As opposed to the original heat kernel func-
220 tion, this new similarity function has a supervised nature. In particular, if
221 α parameter is equal to 1, the algorithm is equivalent to the original NCuts.
222 However, if α is set to a value higher than 1, the similarity between any two
223 patterns in the same class is artificially increased. As a result, similarities
224 between samples in a same class usually become larger than the similarities
225 between any two patterns belonging to different classes. Thus, samples be-
226 longing to the same classes tend to group in the same clusters as illustrated

227 in Fig. 2. This results in improvement of the real-time performance of the
 228 BHDT classification system.

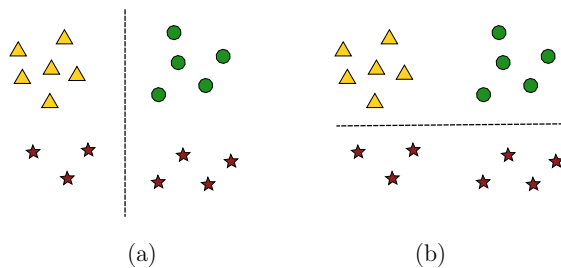


Figure 2: NCuts clustering algorithm splits data samples in order to create maximum gap between the separated samples and it creates overlapping classes as in (a). Modified NCuts algorithm on the other hand takes class membership information into consideration, and samples in the same classes tend to group in the same clusters. So, this process creates more compact clusters as in (b).

228

229 Given a set of samples at each internal node of the hierarchical tree,
 230 the proposed method applies the modified NCuts clustering algorithm to
 231 split data. In some cases, there might still be classes that lie on both sides
 232 of the cluster decision boundaries. In such cases, the approach introduced
 233 in (Marszalek and Schmid, 2008) is followed. In particular, the overlapping
 234 classes are introduced into both clusters if the ratio of the overlapping samples
 235 is greater than a selected threshold, and the resulting uncertain classification
 236 decisions are postponed until the number of classes is reduced and it becomes
 237 tractable for learning good decisions.

238 2.3. Class Based Large Margin Clustering

239 In this method, we focus on the separability of classes rather than samples
 240 and apply original NCuts clustering for partitioning class sets directly. Given

241 a set of m ($m \leq C$) classes at each internal node of the hierarchical tree, the
 242 proposed clustering algorithm would now construct a weighted graph with
 243 m vertices (one for each class at a node) and a set of edges containing these
 244 vertices. It should be noted that the size of the graph is greatly reduced since
 245 the number of classes in a node is considerably smaller than the number of
 246 samples in those classes. Since the SVM margin is equivalent to minimum
 247 distance between the convex hulls of classes, this distance measure is used
 248 during computation of the weights of the edges. As in SVM, the shortest
 249 distance between convex hulls of two classes ω_i and ω_j is determined by the
 250 two closest points in these convex hulls. The problem of finding these two
 251 points can be represented as a quadratic optimization problem

$$\begin{aligned} \min_{\mathbf{u}, \mathbf{v}} \quad & \frac{1}{2} \|\mathbf{X}_i \mathbf{u} - \mathbf{X}_j \mathbf{v}\|^2 \\ \text{s.t.} \quad & \sum_{k=1}^{n_i} u_k = 1, \quad \sum_{k=1}^{n_j} v_k = 1, \quad \mathbf{u} \geq 0, \quad \mathbf{v} \geq 0, \end{aligned} \tag{6}$$

252 where \mathbf{X}_i represents the matrix whose columns are sample vectors belonging
 253 to the class ω_i and $\mathbf{u} \geq 0$ implies that its all elements are greater than or
 254 equal to zero. It must be noted that the objective function of this quadratic
 255 optimization problem is convex and a global minimum exists. This for-
 256 mulation is also equivalent to the hard-margin SVM formulation (Bennett
 257 and Bredensteiner, 2000). Let \mathbf{u}^* and \mathbf{v}^* be an optimal solution of (6).
 258 The minimum distance between the convex hulls of classes is then given by
 259 $d(\omega_i, \omega_j) = \|\mathbf{X}_i \mathbf{u}^* - \mathbf{X}_j \mathbf{v}^*\|$. This distance is equivalent to $2/\|\mathbf{w}\|$ in SVM
 260 formulation where \mathbf{w} represents the normal of the optimal separating hyper-
 261 plane returned by the SVM algorithm. However, a problem may arise if the
 262 convex hulls of classes overlap, i.e., classes are not linearly separable. In this

263 case, the distances between those classes become zero and it may not reflect
 264 the actual similarity between classes. If the classes are close to being linearly
 265 separable and they overlap because of a few outliers, the influence of those
 266 outliers can be reduced by contracting or reducing the convex hulls during
 267 distance computation by introducing an upper bound on the coefficients in
 268 (6)¹. In this case, the new optimization problem becomes

$$\begin{aligned} \min_{\mathbf{u}, \mathbf{v}} \quad & \frac{1}{2} \|\mathbf{X}_i \mathbf{u} - \mathbf{X}_j \mathbf{v}\|^2 \\ \text{s.t.} \quad & \sum_{k=1}^{n_i} u_k = 1, \quad \sum_{k=1}^{n_j} v_k = 1, \quad 0 \leq \mathbf{u} \leq \tau, \quad 0 \leq \mathbf{v} \leq \tau, \end{aligned} \quad (7)$$

269 where $\tau \leq 1$ is the user-chosen positive bound. However, if the classes are
 270 not linearly separable, the data can be mapped into a higher-dimensional
 271 space where the classes become linearly separable. It should be noted that
 272 the objective function of (7) can be written in terms of the dot products of
 273 the samples, which allows the use of the kernel trick. As a result, the data
 274 can be implicitly mapped into a higher-dimensional space where the convex
 275 hulls do not overlap and the distances between classes in the mapped space
 276 can be computed.

277 There is only one design parameter t , the width of the heat kernel, to be
 278 fixed. Well chosen values of t generate well-balanced clusters at each node
 279 of the decision tree, which is critical for an efficient and reliable classifica-
 280 tion. On the other hand, if the width is too small, the algorithm will favor
 281 separating a single isolated class from the remaining classes. For unbalanced

¹In fact this is equivalent to soft margin formulation of SVMs. The reader is referred to (Bennett and Bredensteiner, 2000) for more information.

282 datasets, it is more efficient to separate the large classes (classes with many
283 samples) at the upper levels of the hierarchy, which renders the well-balanced
284 binary decision nodes less efficient. In such cases, the weights of the adja-
285 cency matrix and the width of the kernel can be adjusted to accommodate
286 this kind of supervision. More precisely, we can deliberately decrease the
287 values of the edge weights between the large classes and the others by lower-
288 ing the width to ensure that the large classes will be separated at the upper
289 levels of the hierarchy.

290 *2.4. Removing Outliers*

291 In the proposed clustering method, the presence of data outliers can sig-
292 nificantly change the true geometric structure of the convex hulls. It is un-
293 desirable to allow a few outlying points to excessively influence the distance
294 computations. Therefore, the influence of data outliers should be restricted
295 in the clustering process. As described earlier, this can be done by contract-
296 ing or reducing the convex hulls during distance computations by putting
297 an upper bound on the coefficients in (7) (Bennett and Bredensteiner, 2000).
298 Although this procedure reduces the effects of data outliers, it does not allow
299 the identification of all outlier samples. For a fair comparison with the ex-
300 isting methods, the Support Vector Data Description (SVDD) method (Tax
301 and Duin, 2005) was used for identifying and removing the outliers in the
302 experiments. Given a class, the SVDD method finds a compact bounding
303 hypersphere enclosing all samples in that class. In the case of data outliers,
304 the volume of the hypersphere is minimized for detecting those outliers. The
305 class samples that fall outside the bounding sphere are considered as outliers.
306 This method also allows the use of the kernel trick, and thus it is compatible

307 with the proposed clustering method and SVMs.

308 **3. Experiments**

309 The BHDTs using the proposed clustering methods were tested on syn-
310 thetic and real databases to assess their performance, and they were com-
311 pared to the OAO, OAR and DAG SVMs as well as the BHDTs using k-means
312 based clustering (Vural and Dy, 2004) in terms of classification accuracy and
313 testing time². We experimented with the SVM classifiers using linear ker-
314 nel $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$, polynomial kernels $k(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle)^p$ with degree
315 $p = 2, 3$, and the Gaussian kernel $\exp(-d(\mathbf{x}_i, \mathbf{x}_j)^2/t)$. For some databases, it
316 was found that the SVM classification algorithms using linear or polynomial
317 kernels either did not converge to a solution or the classification performances
318 were too low since the selected kernels failed to approximate the class de-
319 cision boundaries correctly. Thus, such results were omitted and only the
320 results for kernels yielding good classification accuracies were reported. The
321 distances between samples were computed using the Euclidean distance, and
322 the heat kernel function was used for weighting the edges of the adjacency
323 matrix of the graph during clustering for the BHDT method using the pro-
324 posed sample based clustering. However, for the BHDT method using the
325 proposed class based clustering, the selected kernel function (linear, poly-
326 nomial, or the Gaussian) was used for computing the distances between the
327 convex hulls of classes. Subsequently, the heat kernel function was again used
328 to weight the edges as usual.

²All programs can be downloaded at <http://www2.ogu.edu.tr/~mlcv/softwarelink.htm>

329 The shape of the decision boundaries and the distances between the sam-
330 ples or classes significantly affect the optimal values of the design parameters,
331 the kernel width t , and α . In other words, choosing the best design param-
332 eters is data dependent. Therefore, randomly created training and validation
333 sets were used to fix these parameters. The best values of the design param-
334 eters were determined using a global coarse-to-fine search. Specifically, the
335 minimum and maximum values of design parameters that produce acceptable
336 classification rates were first determined by coarsely searching over a wide
337 range of the parameter space. Subsequently, a coarse grid was constructed
338 over the unknown parameters using these computed values and finally a local
339 search was performed near the parameters that yield the best classification
340 rate for determining the final best values. All experiments were conducted
341 in Matlab environment using a 3-GHz machine with 3 GB of RAM.

342 *3.1. Experiments on Synthetic Data*

343 Here we illustrate some properties of the methods on two simple synthetic
344 data sets. For the first database, 3-dimensional samples drawn from normal
345 distributions were used with means $[\mu \mu \mu]^\top$ where the value of μ is changed
346 between -50 and 48 to create 100 classes. The classes were close to being
347 linearly separable with having small overlaps between them. Thus, it was
348 assumed that k-means based clustering should work well in this case. For
349 each class, we used 20 samples for training and 20 samples for testing.

350 For the second database we used 2-dimensional samples drawn from two-
351 component mixture models which are typically used in XOR problem. By
352 shifting centers of mixture components, 50 classes were created, each having
353 40 samples. The first six classes are plotted in Fig. 3. It can be observed

354 that the classes are not linearly separable in this case, and k-means based
 355 clustering would not work well since the overall means are near the origin for
 356 all classes. A total of 40 samples per class were used for both training and
 testing.

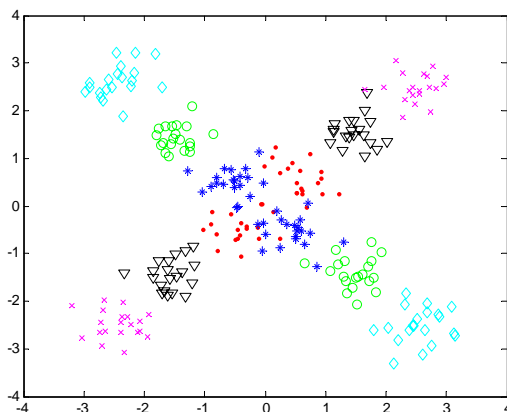


Figure 3: The first six classes having two mixture components from the second synthetic database. The mean of each class is the origin.

357

358 The classification rates for the different kernels on the first and the second
 359 synthetic databases are given in Tables 1- 3. The Proposed Method1 refers to
 360 the BHDT classifier using the modified NCuts clustering, and the Proposed
 361 Method2 denotes the BHDT classifier using the class based NCuts clustering.
 362 The design parameters were set by repeating the above described procedures
 363 5 times and the final reported classification accuracies are averages over 10
 364 repetitions. Asterisks in the table indicate the performance differences that
 365 are statistically significant at 5% level between the given method and the
 366 corresponding best result indicated in bold (statistical significany tests are
 367 determined based on the two-sample t test procedure (Devore, 2004) for all

Table 1: Classification Rates(%) for the Linear Kernel on the First Synthetic Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	87.27 ± 1.4	12.27
Proposed Method2	87.25 ± 1.3	1.73
BHDT of (Vural and Dy, 2004)	87.25 ± 1.3	1.76
DAG SVM	87.38 ± 1.3	28.48
OAO SVM	87.24 ± 1.3	1280.00
OAR SVM	NA	NA

368 experiments), and the testing time indicates the time consumed for classifying
 369 all test points. It should be noted that the testing time for a single sample is
 370 fixed for multi-class SVMs using the OAO and OAR approaches; however, it
 371 changes for BHDTs and DDAGs. An exact number of $C - 1$ classifiers must
 372 be evaluated in order to label a test sample for DDAGs. For BHDTs, the
 373 best case occurs if the predicted class is found at the first node, and the worst
 374 case occurs if the predicted class is found after applying all $C - 1$ decision
 375 functions.

376 For the first database, only linear and the Gaussian kernels produced sat-
 377 isfactory results where the Gaussian kernel yielded better results than the
 378 linear kernel. For both kernels, the classification rates of all tested methods
 379 are very similar except for OAR SVM, which yields a very poor classification
 380 accuracy with respect to the other methods for the Gaussian kernel and does
 381 not converge to a solution for the linear kernel (it is because all pairwise

Table 2: Classification Rates(%) for the Gaussian Kernel on the First Synthetic Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	90.23 \pm 0.4	14.84
Proposed Method2	90.21 \pm 0.4	2.08
BHDT of (Vural and Dy, 2004)	90.26 \pm 0.5	2.10
DAG SVM	90.23 \pm 0.4	33.67
OAo SVM	90.30 \pm 0.5	1489.80
OAR SVM	71.88* \pm 0.4	37.97

382 classes are close to being linearly separable, but classes are no longer linearly
383 separable when the OAR scheme is used). It can be found that BHDTs
384 using the proposed class based clustering and BHDTs using k-means clus-
385 tering offer the best performance in terms of testing time for both kernels,
386 and the BHDT classifier using the proposed class based large margin cluster-
387 ing is considerably faster than the one using the proposed modified NCuts
388 clustering.

389 For the second database it was found that only the Gaussian kernel
390 worked well since the decision boundaries between classes are nonlinear and
391 highly complex. For all other tested kernels, either the SVM classification
392 algorithm did not converge to a solution or the classification accuracies were
393 too low. For the Gaussian kernel, the BHDT classifier using the proposed
394 class based clustering achieved the best performance in terms of testing time
395 among all the tested methods. Both of our proposed methods, DAG and

Table 3: Classification Rates (%) for the Gaussian Kernel on the Second Synthetic Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	96.98 ± 0.4	17.53
Proposed Method2	96.97 ± 0.4	1.81
BHDT of (Vural and Dy, 2004)	73.73* ± 6.3	2.67
DAG SVM	96.74 ± 0.9	16.82
OAo SVM	96.95 ± 0.4	373.66
OAR SVM	94.12* ± 2.8	15.92

396 OAo SVMs yield the best classification accuracies where the BHDT classi-
 397 fier using modified NCuts clustering wins with a slight edge. As expected,
 398 the BHDT classifier using k-means clustering has the worst performance in
 399 terms of classification accuracy; its testing time is also low as compared to the
 400 second proposed method. This is because all class means are near the origin
 401 and k-means based clustering fails to measure the actual similarities among
 402 classes, and thus it becomes very difficult to separate the resulting classes
 403 using the SVM classifier. As in the first synthetic database, the BHDT clas-
 404 sifier using the class based clustering is much faster than the one using the
 405 modified NCuts clustering.

406 *3.2. Experiments on Coil-100 Object Database*

407 Here we test methods on the Coil-100 objects database³. The Coil-100
 408 database includes 72 view-images of 100 different objects taken at 5-degree-
 409 apart orientations. The size of each image is 128×128 . All images were
 410 converted to gray scale and Principal Component Analysis (PCA) was ap-
 411 plied to reduce the dimensionality to 100. A total of 36 samples were used
 412 from each class for training, and the remaining 36 samples were used for
 413 testing. The design parameters were set using 5 random training/test splits
 414 and the final reported classification accuracies were averages over 10 random
 415 training/test splits. Data outliers in the training sets were removed using
 416 the SVDD method prior to the application of the classifiers.

Table 4: Classification Rates (%) for Polynomial Kernel with $p = 2$ on Coil-100 Objects Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	93.93* \pm 0.6	17.29
Proposed Method2	94.32 \pm 0.4	15.75
BHDT of (Vural and Dy, 2004)	94.26 \pm 0.4	17.90
DAG SVM	94.14 \pm 0.5	64.24
OA0 SVM	94.15 \pm 0.5	2710.02
OAR SVM	94.50 \pm 0.5	122.72

³Available at <http://www1.cs.columbia.edu/CAVE/software/softlib/coil-100.php>

Table 5: Classification Rates (%) for Polynomial Kernel with $p = 3$ on Coil-100 Objects Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	$89.21^* \pm 0.4$	28.71
Proposed Method2	$89.70^* \pm 0.3$	25.32
BHDT of (Vural and Dy, 2004)	$89.70^* \pm 0.4$	25.67
DAG SVM	$88.82^* \pm 0.3$	64.89
OA0 SVM	$88.80^* \pm 0.3$	2698.90
OAR SVM	91.86 ± 0.2	650.96

417 The results for different kernels are given in Tables 4- 6. It can be seen
418 from the tables that the BHDT classifier using the proposed class based
419 clustering is the most efficient method in terms of testing time in all cases.
420 However, its classification rate is slightly lower than some other traditional
421 multi-class SVMs. Nevertheless, the classification performance of the pro-
422 posed method is still satisfactory for most of the cases. It is found that the
423 BHDT of (Vural and Dy, 2004) using k-means clustering has worse perfor-
424 mance than the second proposed method in terms of classification accuracy
425 for polynomial kernel with degree 2 whereas it has better performance for the
426 Gaussian kernel. Overall the best classification accuracies are obtained using
427 the Gaussian kernels, and the worst results are obtained for the polynomial
428 kernel with degree 3.

Table 6: Classification Rates (%) for the Gaussian Kernel on Coil-100 Objects Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	94.97* \pm 0.2	104.69
Proposed Method2	96.45* \pm 0.3	58.03
BHDT of (Vural and Dy, 2004)	96.97* \pm 0.3	77.78
DAG SVM	97.39 \pm 0.3	209.96
OA0 SVM	97.33 \pm 0.3	9819.70
OAR SVM	96.28* \pm 0.3	515.26

429 *3.3. Experiments on the AR Face Database*

430 The AR face database includes 26 frontal images with different facial
431 expressions, illuminations conditions, and occlusions for 126 subjects. Images
432 were recorded in two different sessions 14 days apart. Thirteen images were
433 recorded under controlled circumstances in each session. The size of the
434 images is 768×576 . A total of 50 individuals (30 males and 20 females)
435 were randomly selected for experiment. The images were aligned and scaled
436 so that the centers of the two eyes always fall on fixed coordinates. The
437 pre-processed images of a person are illustrated in Fig. 4. As in the previous
438 experiment, PCA was applied and the dimensionality was decreased to 200.
439 A total of 20 samples were randomly selected for training for each individual
440 while keeping the remaining six for testing. This process was repeated 10
441 times and the final classification rates were obtained by averaging the results
442 obtained in each run.



Figure 4: Some pre-processed images from the AR face database.

443 The results for the different kernels are given in Tables 7- 9. It can be seen
 444 that the best classification accuracies are obtained using the linear kernel
 445 for all methods except OAR SVM, which achieves the best classification
 446 accuracy for the polynomial kernel with degree 2. For the linear kernel, the
 447 best classification accuracy is obtained using DAG SVM and OAO SVM,
 448 with DAG SVM being more successful owing to a slight advantage. It can be
 449 seen that both of these methods significantly outperform the other methods.
 450 It is found that the BHDT of (Vural and Dy, 2004) is the fastest in terms
 451 of testing time followed by the Proposed Method2. For both the polynomial
 452 and the Gaussian kernels, the best classification accuracy is obtained by OAR
 453 SVM. The BHDT classifier using the proposed class based clustering and the
 454 BHDT of (Vural and Dy, 2004) are the most efficient methods in terms
 455 of testing time for the polynomial and Gaussian kernel functions whereas
 456 the Proposed Method2 outperforms the BHDT of (Vural and Dy, 2004) in
 457 terms of classification accuracy. The Proposed Method1 also outperforms
 458 the BHDT of (Vural and Dy, 2004) for the Gaussian kernel; however, its
 459 classification accuracy is slightly lower for the polynomial kernel.

Table 7: Classification Rates (%) for Linear Kernel on AR Face Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	96.87* \pm 1.0	1.58
Proposed Method2	97.47* \pm 0.7	0.65
BHDT of (Vural and Dy, 2004)	96.47* \pm 0.9	0.53
DAG SVM	99.20 \pm 0.4	2.95
OA0 SVM	99.17 \pm 0.3	64.37
OAR SVM	97.10* \pm 0.5	5.01

460 **4. Summary and Conclusion**

461 In this study, two new clustering algorithms were proposed for the par-
462 tition of data samples for SVM based BHDTs. The proposed methods have
463 two major advantages over the traditional clustering algorithms used for this
464 purpose. Firstly, the proposed methods are suitable when SVMs are used
465 as the base classifier. On the other hand, the most commonly employed k-
466 means clustering algorithm may not be compatible with the SVM classifier as
467 demonstrated in the synthetic database experiments. It must be noted that
468 the k-means clustering is based on the assumption that the class densities are
469 Gaussian and isotropic. However, in general, this assumption is not true for
470 most of the classification problems in the real world, and k-means clustering
471 algorithm might produce clusters that are difficult to separate using the SVM
472 classifier. Secondly, the proposed methods allow the use of kernel functions as
473 opposed to the other clustering algorithms such as Kullback-Leibler distance

Table 8: Classification Rates (%) for Polynomial Kernel with $p = 2$ on AR Face Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	$90.67^* \pm 1.2$	12.14
Proposed Method2	$91.12^* \pm 1.5$	2.74
BHDT of (Vural and Dy, 2004)	$90.83^* \pm 1.5$	2.74
DAG SVM	$93.00^* \pm 1.6$	3.68
OA0 SVM	$93.00^* \pm 1.6$	67.96
OAR SVM	99.14 ± 0.2	65.52

474 based clustering or the k-means algorithm. This is a significant advantage
 475 since the use of different kernel functions in SVM can significantly change
 476 the decision boundaries.

477 In the proposed methods, the most important design parameter is the
 478 width of the heat kernel function. Selection of the best value of this param-
 479 eter is data dependent since the distances between the convex class sets or
 480 data samples significantly affect the optimal values of the parameter. More-
 481 over, for unbalanced data sets it must be specially fixed so that the larger
 482 classes will be separated at the upper levels of the hierarchical tree. Thus, it
 483 is better to fix this parameter based on a cross-validation scheme as done in
 484 the present study. The use of the proposed clustering algorithms with well-
 485 chosen parameters in BHDTs will generate well-balanced separable clusters
 486 at each internal node of the decision tree, which is crucial for a reliable and
 487 efficient classification. Although both of the proposed clustering methods

Table 9: Classification Rates (%) for the Gaussian kernel on the AR Face Database.

Methods	Classification Rate	Testing Time (secs)
Proposed Method1	92.00* \pm 1.6	15.36
Proposed Method2	91.90* \pm 1.2	6.07
BHDT of (Vural and Dy, 2004)	90.87* \pm 1.5	6.45
DAG SVM	94.83* \pm 1.5	16.97
OA0 SVM	94.90* \pm 1.4	408.17
OAR SVM	97.93 \pm 0.8	102.05

488 worked well in the experiments, the proposed class based NCuts clustering
489 method seemed to be more efficient than the proposed sample based cluster-
490 ing method. The average testing time for labeling an unseen sample using
491 the proposed methods is $O(\log_2 C)$. Thus, it can be concluded that the pro-
492 posed methods can considerably increase the speed of classification with a
493 small decrease in the classification accuracies when the number of classes is
494 large.

495 References

- 496 Bennett, K. P. and Bredensteiner, E. J., 2000. Duality and geometry in SVM
497 classifiers, International Conference on Machine Learning.
- 498 Burges, C.J.C., 1998. A tutorial on support vector machines for pattern
499 recognition, Data Mining and Knowledge Discovery, pp. 121-167.

- 500 Casasent, D. and Wang, Y.-C., 2005. A hierarchical classifier using new sup-
501 port vector machines for automatic target recognition, *Neural Networks*,
502 vol. 18, pp. 541-548.
- 503 Chen, Y., Crawford, M.M., and Ghosh, J., 2004. Integrating support vector
504 machines in a hierarchical output space decomposition framework, *Proc.*
505 *2004 International Geoscience and Remote Sensing Symposium*, pp. 949-
506 952.
- 507 Cortes, C. and Vapnik, V., 1995. Support vector networks, *Machine Learning*,
508 vol. 20, pp. 273-297.
- 509 Devore, J. L., 2004. *Probability and Statistics for Engineering and the Sci-*
510 *ences*, Brooks/Cole-Thomson Learning, Belmont CA.
- 511 Friedman, J. H. and Bentley, J. I. and Finkel, R. A., 1977. An algorithm
512 for finding best matches in logarithmic expected time, *ACM Trans. Math*
513 *Software*, vol. 3, pp. 209-226.
- 514 Griffin, G. and Perona, P., 2008. Learning and using taxonomies for fast
515 visual categorization, *IEEE Computer Society Conference on Computer*
516 *Vision and Pattern Recognition*.
- 517 Hsu, C. and Lin, C., 2002. A comparison of Methods for Multi-Class Support
518 Vector Machines, *IEEE Transactions on Neural Networks*, vol. 13, pp. 415-
519 425.
- 520 Liu, S., Yi, H., Chia, L.T., and Rajan, D., 2005. Adaptive Hierarchical Multi-
521 Class SVM Classifier for Texture-Based Image Classification, *IEEE Inter-*
522 *national Conference on Multimedia and Expo*.

- 523 Marszalek, M. and Schmid, C., 2008. Constructing category hierarchies for
524 visual recognition, European Conference on Computer Vision.
- 525 Martinez, A. M. and Benavente, R., 1998. The AR Face Database, CVC
526 Tech. Report 24.
- 527 Platt, J. C., Cristianini, N., and Shawe-Taylor, J., 2000. Large margin dags
528 for multiclass classification, Advances in Neural Information Processing
529 Systems, pp. 547-553.
- 530 Rahimi, A. and Recht, B., 2004. Clustering with Normalized Cuts is Clus-
531 tering with a Hyperplane, Statistical Learning in Computer Vision.
- 532 Schölkopf, B., Smola, A., 2002. Learning with Kernels: Support Vector Ma-
533 chines, Regularization, Optimization and Beyond.
- 534 Shi J., and Malik, J., 2000. Normalized Cuts and Image Segmentation, IEEE
535 Transactions on Pattern Analysis and Machine Intelligence, vol. 22, pp.
536 888-905.
- 537 Tax, D. M. J., and Duin, R. P. W., 2005. Support Vector Data Description,
538 Machine Learning, vol. 54, pp. 45-66.
- 539 Vural V. and Dy, J.G., 2004. A hierarchical method for multi-class support
540 vector machines, Proceedings of the twenty-first International Conference
541 on Machine Learning.
- 542 Zhigang, L., Wenzhong, S., Qianqing, Q., Xiaowen, L., and Donghu, X.,
543 2005. Hierarchical support vector machines, Geoscience and Remote Sens-
544 ing Symposium.