



Repositorio Institucional de la Universidad Autónoma de Madrid

<https://repositorio.uam.es>

Esta es la **versión de autor** del artículo publicado en:

This is an **author produced version** of a paper published in:

Pattern Recognition 45.12 (2012): 4414 – 4427

DOI: <http://dx.doi.org/10.1016/j.patcog.2012.06.002>

Copyright: © 2012 Elsevier B.V.

El acceso a la versión del editor puede requerir la suscripción del recurso
Access to the published version may require subscription

Hierarchical Linear Support Vector Machine

I. Rodriguez-Lujan^{a,*}, C. Santa Cruz^a, R. Huerta^b

^a*Departamento de Ingeniería Informática and Instituto de Ingeniería del Conocimiento, Universidad Autónoma de Madrid, 28049 Madrid, Spain.*

^b*BioCircuits Institute, University of California San Diego, La Jolla CA 92093-0404, USA.*

Abstract

The increasing size and dimensionality of real-world datasets make it necessary to design efficient algorithms not only in the training process but also in the prediction phase. In applications such as credit card fraud detection, the classifier needs to predict an event in 10 milliseconds at most. In these environments the speed of the prediction constraints heavily outweighs the training costs. We propose a new classification method, called a Hierarchical Linear Support Vector Machine (H-LSVM), based on the construction of an oblique decision tree in which the node split is obtained as a Linear Support Vector Machine. Although other methods have been proposed to break the data space down in subregions to speed up Support Vector Machines, the H-LSVM algorithm represents a very simple and efficient model in training but mainly in prediction for large-scale datasets. Only a few hyperplanes need to be evaluated in the prediction step, no kernel computation is required and the tree structure makes parallelization possible. In experiments with medium and large datasets, the H-LSVM reduces the prediction cost considerably while achieving classification results closer to the non-linear SVM than that of the linear case.

Keywords: Large-Scale Learning, Real-Time Prediction, Support Vector Machine, Decision Tree, Pegasos Algorithm

1. Introduction

Support Vector Machines (SVMs) have been widely used in classification problems as a result of their effectiveness. However, the increasing size of real-world datasets in domains such as bioinformatics, document categorization or credit card fraud detection compromises their application. The computational complexity of the SVM decision function scales with respect to the number of support vectors n_{SV} and Steinwart [1] showed that the number of support vectors scales linearly with respect to the number of training patterns. Consequently, other machine learning techniques are preferred in those large-scale domains in which an

*Corresponding author. Instituto de Ingeniería del Conocimiento; C/Francisco Tomás y Valiente, 11; E.P.S., Edificio B; UAM-Cantoblanco, 28049 Madrid, Spain. Tel: +34 91 497 2339; fax: +34 91 497 2334

Email addresses: irene.rodriguez@iic.uam.es (I. Rodriguez-Lujan), carlos.santacruz@iic.uam.es (C. Santa Cruz), rhuerta@ucsd.edu (R. Huerta)

8 efficient prediction step is needed, especially in real-time applications such as credit card fraud detection
9 which requires a response time of less than 10 milliseconds.

10 Although the machine learning community has been mainly focused on speeding up the training of
11 the SVM, the emergence of applications requiring fast classification makes the design of new algorithms
12 necessary whilst maintaining as much as possible the effectiveness of non-linear SVMs and improving its
13 classification complexity at the same time. Linear SVMs are the best alternative for fast execution because
14 their decision boundary is made up of a single hyperplane. However, their performance for non-linear
15 problems is uncompetitive and a compromise between performance and classification speed is needed.

16 The computational complexity of testing a pattern using a non-linear SVM is $O(n_{SV} \times d \times n_K)$ where
17 n_{SV} are the number of support vectors, d is the dimension of the samples and n_K is the cost of evaluating
18 the kernel function. In large-scale problems, the number of support vectors is usually much higher than the
19 dimension of the problem ($n_{SV} \gg d$) which is why almost all methods proposed in the literature aim at
20 reducing n_{SV} . The methods for reducing the number of support vectors can be divided into two groups [2]:

21 • **Numerical techniques** find a reduced set of basis functions necessary to classify a pattern. These
22 algorithms usually consider all of the training patterns and find a sparse representation of the support
23 vectors. A more detailed overview of these methods is given by Keerthi et al. [2]. According to
24 Keerthi's categorization, the support vector reduction can be carried out as a post-processing phase
25 after training the SVM model or during the training phase thus imposing a certain sparsity in the
26 basis functions. The post-processing techniques ([3, 4]) reduce the number of support vectors once
27 the SVM model has been trained. Therefore, they still depend on the standard SVM training which
28 can be extremely costly in large-scale problems. Among the direct simplification approaches based
29 on imposing sparseness on the basis functions in the primal space, several methods can be found in
30 the literature [5, 6, 2]. These approaches considerably reduce the SVM prediction cost while having
31 a competitive classification accuracy, but in some datasets the number of basis functions needed to
32 maintain a competitive classification accuracy is still high for efficient training and prediction phases
33 [2].

34 • **Data-reduction methods** reduce the number of SVM training patterns dividing the original training
35 set into one or several smaller datasets to train an SVM in each partition. Boosting [7], bagging [8],
36 parallel mixture of SVMs [9] and SVM-cascade [10] algorithms can be categorized into this group. A
37 recent work [11] proposes the DTSVM (Decision Tree Support Vector Machine) algorithm to build a
38 decision tree with axis-parallel splits via the CART method [12] and to train an SVM with an RBF
39 kernel in each leaf of the tree. This method reports a significant reduction of the number of support
40 vector evaluations in the test or prediction phase. However, the number is still too high for large-scale
41 datasets at the level used in credit card fraud detection.

42 Our approach does not consider the use of non-linear SVMs because of their high classification and
43 training cost. The aim of our work is to provide a model which generates non-linear decision boundaries via
44 piecewise linear decision functions. This approach is motivated by the low classification cost of linear SVMs.
45 Moreover, recent algorithms [13, 14] have shown the efficiency of stochastic gradient descent approaches for
46 training linear SVMs and their usual fast convergence for large-scale datasets. Our work is not the first
47 attempt at approximating non-linear SVMs through the linear case. Recent contributions have proposed the
48 use of linear SVMs in the manifold coordinates such as sparse coding or local coordinate coding [15, 16, 17].
49 The MLSVM method [18] is based on a mixture of linear SVMs defining an underlying probabilistic model
50 which implicitly selects the linear SVMs to be used to classify each pattern. A test sample is classified by
51 the weighted average over the mixture of classifiers.

52 Our work approaches the task as the construction of a binary decision tree whose nodes are linear SVMs.
53 The combination of linear SVMs and decision trees is motivated by the results of Bennett et al. [19] and
54 some research combining decision trees and SVMs. An interesting comparison of the classification cost of
55 decision trees and SVMs is given by Kumar and Gopal [20]. Basically, decision trees are much faster than
56 SVMs in classifying new instances whereas the classification accuracy of SVMs is superior. Pursuing the
57 objective of speeding up the prediction phase of a classifier, Zapién et al. [21] proposes a tree structure where
58 the split of each node is a linear SVM. The tree presents a particular structure, which could be considered as
59 a cascade of linear SVMs as the tree only expands the right branches. Then, it is assumed that each split in
60 the tree is able to classify correctly all of the patterns belonging to the left child. The main difference with
61 our method is that our tree is a *complete* binary decision tree in the sense that both children of each node
62 can be expanded in the following steps. Although Zapién et al. provide the most straightforward approach,
63 a balanced tree search is on average faster at classifying a datapoint since the cascade structure needs to
64 run through all of the decision nodes to evaluate the worst datapoints. In addition, the hypothesis class
65 (disjunctions of conjunctions) of H-LSVM is more general than that of the Zapién’s model (conjunctions)
66 because (i) the cascade structure (also known as decision list) can be viewed as a special type of decision
67 trees [22] and (ii) the number of decision tree skeletons with k decision nodes is given by the k -th Catalan
68 Number [19, 23] while the Zapién’s cascade structure has only one possibility. The algorithms proposed by
69 Fehr et al. [24] and Sun et al. [25] represent an extension of the Zapién model in which the linear SVM
70 is the split in each node and nonlinear SVMs make up the leaves of the tree. These models still depend
71 on a non-linear SVM which means a large number of support vector evaluations to classify a test sample.
72 The DTO-SVM algorithm [26] builds an oblique decision tree whose node split is selected between the C4.5
73 [27] parallel-split calculated from the categorical variables and the SVM-SMO [28] classifier obtained from
74 continuous attributes. The method still depends on the large number of support vectors given by the SMO
75 which makes large-scale predictions costly.

76 Another interesting approach to combine decision trees and SVMs is the one proposed by Bennett and

77 Blue [29] in which the decision tree structure is set beforehand so that the model is formulated in the primal
78 space as the minimization of a nonconvex objective function with respect to polyhedral constraints. This
79 alternative is substantially different from the aforementioned ones and that adopted in this paper since they
80 obtain the structure of the tree in an on-line manner. In addition, the margin is locally maximized in each
81 node of the H-LSVM tree whereas the Bennett and Blue model looks for a global maximization.

82 The main advantage of our H-LSVM is its ability to classify a pattern in a few milliseconds even for
83 large-scale datasets thus speeding up the prediction phase of the SVMs by several orders of magnitude while
84 maintaining a classification accuracy close of that of the non-linear SVMs. Moreover, the decision tree
85 structure is easily parallelizable which favors training acceleration [30]. As H-LSVM is a piecewise linear
86 classifier, its classification accuracy is not as good as those models based on non-linear SVMs. However, in
87 those systems which do require real-time predictions, the reduction in accuracy is bearable when compared to
88 the runtime savings. As regards other combinations of decision trees and linear SVMs models, the proposed
89 method represents an improvement in the state-of-the-art not only in terms of classification accuracy but
90 also in terms of prediction cost.

91 The paper is organized as follows: Section 2 presents the H-LSVM algorithm including the explanation
92 of several design aspects. Section 3 analyzes and compares the training and prediction complexities of linear
93 SVMs, non-linear SVMs and H-LSVM. Section 4 provides a generalization error bound for the H-LSVM
94 method and Section 5 presents the empirical results in terms of classification accuracy and prediction cost of
95 the proposed method compared to linear and non-linear SVMs and other algorithms based on the speeding
96 up of SVMs via linear SVMs. A numerical analysis of the H-LSVM scalability and generalization error
97 bound are also given in this section.

98 **2. The H-LSVM Algorithm**

99 The proposed algorithm called a Hierarchical Linear Support Vector Machine (H-LSVM) is based on the
100 construction of a decision tree. As described by Breiman et al. [12], four elements must be considered in
101 the construction:

- 102 1. The goodness of the node split which needs to be evaluated in each node of the tree.
- 103 2. The type of test carried out in each node of the tree to decide whether a pattern belongs to the left
104 or to the right child of the current node.
- 105 3. The stop-splitting rule.
- 106 4. The criteria for assigning the class label to a pattern when it reaches a leaf of the tree.

107 The H-LSVM algorithm node split is a linear SVM. The linear SVM is trained using the Pegasos algorithm
108 [13] with weighted patterns. The weight of each pattern is not fixed and it depends on which node of the
109 tree we are working on. For the rest of the elements, well-known techniques and criteria have been used.

110 Once the complete tree is trained, a pruning step can improve the generalization capability of the H-LSVM
 111 model. The four key elements for the construction of a decision tree with the pruning algorithm will be
 112 described in this section.

113 Let us establish some notation. Given a training set $S = \{(\vec{x}_i, y_i)\}_{i=1}^N$, where $\vec{x}_i \in \mathbb{R}^d$ and $y_i \in \{+1, -1\}$,
 114 we define:

- 115 • H_k : a node in the tree.
- 116 • S_{H_k} : subset of samples in the node H_k .
- 117 • $S_{H_k}^+$: subset of positive samples in the node H_k .
- 118 • $S_{H_k}^-$: subset of negative samples in the node H_k .
- 119 • $N_{H_k} = |S_{H_k}|$: number of samples in the node H_k .
- 120 • $N_{H_k}^+ = |S_{H_k}^+|$: number of positive samples in the node H_k .
- 121 • $N_{H_k}^- = |S_{H_k}^-|$: number of negative samples in the node H_k .
- 122 • $\vec{x}_i^{H_k}$: i -th sample in the subset S_{H_k} .
- 123 • \vec{w}_{H_k} : normal vector to the hyperplane associated to the node H_k .
- 124 • b_{H_k} : bias term of the hyperplane associated to the node H_k .
- 125 • $h_{H_k}(\vec{x}_i)$: evaluation of the i -th pattern in the node H_k , that is
 126 $h_{H_k}(\vec{x}_i) = \vec{w}_{H_k} \cdot \vec{x}_i + b_{H_k}$.
- 127 • $S_{H_k}^l$: left child of the node H_k : $\{\vec{x} \in S_{H_k} \mid h_{H_k}(\vec{x}) \leq 0\}$.
- 128 • $S_{H_k}^r$: right child of the node H_k : $\{\vec{x} \in S_{H_k} \mid h_{H_k}(\vec{x}) > 0\}$.
- 129 • $p_i^{H_k}$: weight of the i -th pattern in the node H_k verifying $\sum_{i=1}^{N_{H_k}} p_i^{H_k} = 1 \forall k$.

130 *Splitting Goodness.* The definition of the splitting goodness is based on the impurity function concept [12].
 131 Two different concepts need to be defined: the impurity of a node and the impurity of a split. The impurity
 132 of a node H_k , $I(H_k)$, does not depend on the splits and it is a function of the number of patterns of each
 133 class in the node, $I(H_k) = I(N_{H_k}^+, N_{H_k}^-)$. The impurity of a split is the impurity induced by the node
 134 split which divides the samples into the subsets $S_{H_k}^l$ and $S_{H_k}^r$. The impurity of a split given by \vec{w}_{H_k} and
 135 b_{H_k} , $I(\vec{w}_{H_k}, b_{H_k})$, can be defined straightforwardly from the impurity of the children, $I(H_k^l)$ and $I(H_k^r)$, as
 136 follows,

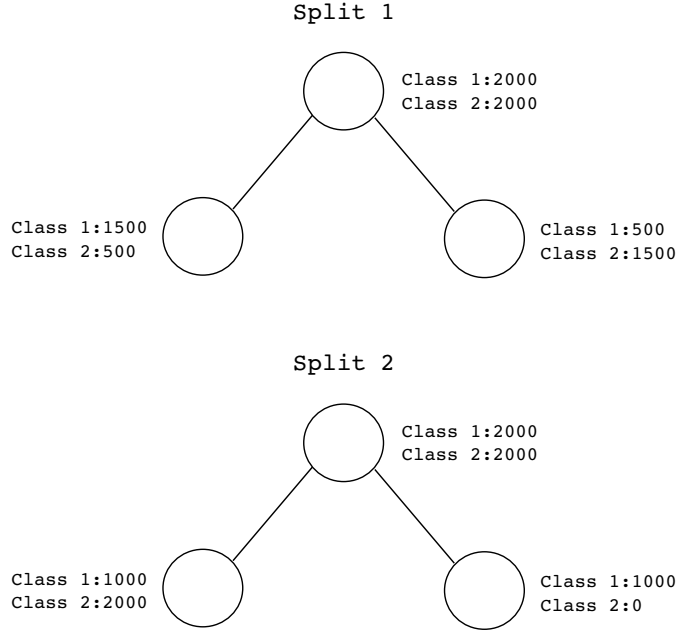


Figure 1: An example of two different splits in a decision tree. If the classification error is used as an impurity measure, both splits misclassified 1,000 samples. Nevertheless, the second split seems more desirable for the future expansion of the tree.

$$I(\vec{w}_{H_k}, b_{H_k}) = \frac{|S_{H_k}^l|}{|S_{H_k}|} I(H_k^l) + \frac{|S_{H_k}^r|}{|S_{H_k}|} I(H_k^r). \quad (1)$$

137 As the aim of the decision tree is to minimize the overall misclassification rate of the tree, it would be
 138 natural to choose the classification error as an impurity measure. However, as pointed out by Breiman et
 139 al. [12, Chapter 4], this measure has two significant limitations: i) The improvement in the impurity can be
 140 zero for all the splits in S_{H_k} , and ii) the inadequacy for an iterative-split decision tree method (see Figure 1
 141 extracted from [12, Chapter 4]).

142 As an alternative, entropy was chosen as impurity function because it is one of the most common impurity
 143 functions in recent methods. The entropy of a node H_k in a binary decision tree is formulated as follows,

$$I(H_k) = -\frac{|(S_{H_k})^+|}{|S_{H_k}|} \times \log \left(\frac{|(S_{H_k})^+|}{|S_{H_k}|} \right) - \frac{|(S_{H_k})^-|}{|S_{H_k}|} \times \log \left(\frac{|(S_{H_k})^-|}{|S_{H_k}|} \right) \quad (2)$$

144 where the superscripts + and - represents the category of the samples.

145 *Splitting Criteria.* The H-LSVM algorithm uses a linear SVM as splitting criteria because a single hyperplane
 146 vector \vec{w} is obtained as a result of the training process which makes prediction much more efficient. More
 147 precisely, the Pegasos algorithm [13] was used because it is an efficient method for training linear SVMs in

148 large-scale datasets. The Pegasos algorithm minimizes the objective function of a linear SVM in the primal
 149 space,

$$\min_{\vec{w}} \frac{\lambda}{2} \|\vec{w}\|^2 + \frac{1}{N} \sum_{(x,y) \in S} \mathcal{L}(\vec{w}; (\vec{x}, y)) \quad (3)$$

150 where $\mathcal{L}(\vec{w}; (\vec{x}, y))$ represents the loss function,

$$\mathcal{L}(\vec{w}; (\vec{x}, y)) = \max\{0, 1 - y(\vec{w} \cdot \vec{x})\}. \quad (4)$$

151 To solve the problem in Equations 3 and 4, the Pegasos algorithm alternates between stochastic gradient
 152 descent steps and projection steps:

- 153 • **Stochastic gradient descent.** On iteration t of the algorithm, a set $A_t \subset S$ of size k is chosen.
 154 Then, the objective function given in Equation 3 is approximated by,

$$\min_{\vec{w}} f(\vec{w}; A_t) = \min_{\vec{w}} \frac{\lambda}{2} \|\vec{w}\|^2 + \frac{1}{k} \sum_{(\vec{x}, y) \in A_t} \mathcal{L}(\vec{w}; (\vec{x}, y)). \quad (5)$$

155 The update of the \vec{w} based on the gradient descent method is given by $\vec{w}_{t+\frac{1}{2}} = \vec{w}_t - \eta_t \vec{\nabla}_t^w$, where
 156 $\eta_t = \frac{1}{\lambda t}$ is the learning-rate and $\vec{\nabla}_t^w$ is the subgradient of $f(\vec{w}; A_t)$ with respect to \vec{w} on the iteration t ,

$$\vec{\nabla}_t^w = \lambda \vec{w}_t - \frac{1}{k} \sum_{(\vec{x}, y) \in A_t^+} y \vec{x}, \quad (6)$$

157 A_t^+ being the set of samples in A_t with non-zero loss that is, $A_t^+ = \{(\vec{x}, y) \in A_t \mid y(\vec{w}_t \cdot \vec{x}) < 1\}$.

- 158 • **Projection step.** Projection of $\vec{w}_{t+\frac{1}{2}}$ onto the set $B = \{\vec{w} \mid \|\vec{w}\| \leq \frac{1}{\sqrt{\lambda}}\}$ since it can be shown that
 159 the optimal solution of SVM is in the set B [13].

160 The Pegasos algorithm has been used in the H-LSVM to obtain the oblique splitting hyperplane in each
 161 node of the tree but some changes have been applied:

- 162 • **Weighted-patterns.** The H-LSVM algorithm generates a piecewise linear model using a decision
 163 tree to divide the input space into disjoint regions. In each region, the proportion of patterns of
 164 each class might be unbalanced and might not necessarily be the same as in the original problem.
 165 In addition, some classification problems, such as fraud detection [31] or medical diagnosis [32], are
 166 unbalanced by nature. If the original formulation of the primal SVM objective function is used, the
 167 misclassification cost for each pattern is the same and independent of the class. However, this scheme

168 can give undesirable classifiers which assign the majority class label to all patterns [33] while we are
 169 interested in separating the classes with successive decision tree splits. To overcome the imbalance,
 170 the H-LSVM method computes the weight $p_i^{H_k}$ of the sample \vec{x}_i in the node H_k according to,

$$p_i^{H_k} = \begin{cases} \frac{1}{2N_{H_k}^+} & \text{if } \vec{x}_i \in S_{H_k}^+ \\ \frac{1}{2N_{H_k}^-} & \text{if } \vec{x}_i \in S_{H_k}^- . \end{cases} \quad (7)$$

171 Now, the objective function of the Pegasos algorithm incorporates the sample weight in the loss term,

$$\begin{aligned} \min_{\vec{w}} f(\vec{w}; A_t) = \\ \min_{\vec{w}} \frac{\lambda}{2} \|\vec{w}\|^2 + \frac{1}{k} \sum_{(p, \vec{x}, y) \in A_t} p \max\{0, 1 - y(\vec{w} \cdot \vec{x})\} \end{aligned} \quad (8)$$

172 and the subgradient of Equation 8 respect to \vec{w} on the iteration t is given by,

$$\vec{\nabla}_t^w = \lambda \vec{w}_t - \frac{1}{k} \sum_{(p, \vec{x}, y) \in A_t^+} py \vec{x} . \quad (9)$$

173 It can be easily shown that the Weighted-Pegasos algorithm still verifies that the norm of the optimum
 174 in Equation 8 is upper bounded by $\frac{1}{\sqrt{\lambda}}$ and the number of iterations required for achieving a solution
 175 of accuracy ϵ is $O(\frac{1}{\lambda \epsilon})$.

176 • **The bias term.** The presence of a bias term in the hyperplane is essential for the H-LSVM as a
 177 result of the multiple separation of the feature space. There are different approaches to estimate the
 178 bias term of the hyperplane [13]. Following the heuristics implemented in standard SGD packages¹,
 179 the bias is updated via a subgradient descent and by using a smaller learning rate (scaled by the
 180 heuristically chosen parameter τ) because the bias term is updated more frequently than the weights.
 181 At each epoch t , not only is the stochastic gradient descent applied to the \vec{w} vector but also to the
 182 bias term b : $b_{t+1} = b_t - \tau \eta_t \nabla_t^b$. The subgradient of the bias is given by $\nabla_t^b = -\frac{1}{k} \sum_{(\vec{x}, y) \in A_t^+} py$.

- 183 • **Pegasos Parameters.** Some meta parameters have to be set in the Weighted Pegasos Algorithm,
- 184 – λ **Regularization Parameter**: Obtained via a validation subset or cross validation (Section 5).
 - 185 – T **Maximum number of iterations** in Pegasos Training.
 - 186 – k **size of the subset of samples** A_t used to update the subgradient.

¹<http://leon.bottou.org/projects/sgd>

187 – ϵ_P **Tolerance** in Pegasos Training. Allowable tolerance for the norm of the difference between
188 \vec{w} vectors in consecutive iterations.

189 – τ **bias scale**. In our experiments we set $\tau = 0.01$.

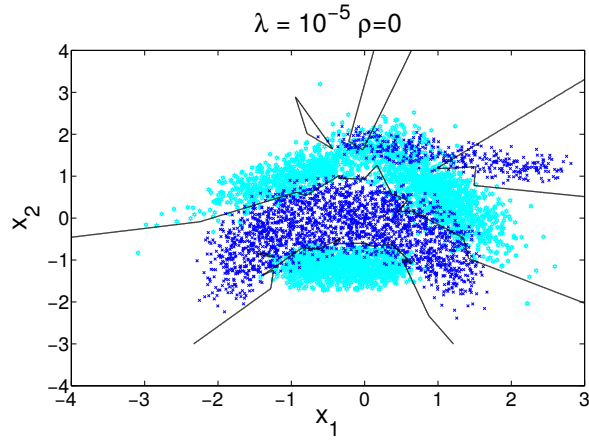
190 *Splitting-Stop Criteria.* A node split is stopped when it does not represent an improvement in the impurity
191 measure or when the rate of training samples associated to this node is lower than a parameter δ . If δ is too
192 high, the tree might be not expanded enough. Small values of δ which yield an overfitted tree are preferable
193 because this tree will be pruned later. That is why, δ was set to 10^{-i} , $i = \lfloor \log_{10} N \rfloor$ in our experiments.

194 *Class Assignment Criteria.* Once a pattern reaches a leaf of the decision tree, it is assigned to the majority
195 class in this leaf. It can be shown [12, Chapter 2] that this rule minimizes the expected misclassification
196 probability of the leaf assuming that the cost of misclassifying a pattern is independent of its class.

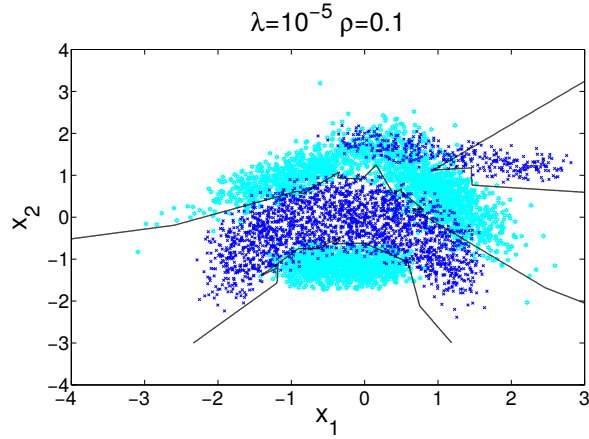
197 *Pruning.* Incorporating a pruning process into a decision tree algorithm reduces the risk of having an
198 overfitted model [34, 35]. Although the SVM formulation already incorporates a regularization term which
199 favors the generalization capability of the optimal hyperplane, a small value for δ in the splitting-stop criteria
200 might imply an overfitted model. This point can be solved by setting different δ values and evaluating the
201 performance of the model in a validation step. However, this approach is computationally costlier than
202 using a small value for δ –that is, making the tree grow as much as possible– and then applying a pruning
203 algorithm. The latter approach is used by H-LSVM and it uses the Cost Complexity (CC) pruning algorithm
204 proposed by Breiman et al. [12]. The CC method requires a pruning set not used to train the tree. This
205 set can be selected randomly or via cross-validation. The rate ρ of those training patterns kept away for the
206 pruning phase is a parameter of the H-LSVM algorithm. The main idea of the CC pruning algorithm is to
207 construct a set of decreasing-sized subtrees of the original tree and evaluate the goodness of each subtree
208 as its classification accuracy on the pruning set. In the original CC method, the smallest subtree with a
209 classification accuracy in k standard deviations of the original tree is selected. In our experiments, we set
210 $k = 0$ and, therefore, the subtree selected is that which has the highest classification accuracy and, in the
211 case of several subtrees with the highest accuracy, the smallest one is chosen. For more details, see [12,
212 Chapters 10,11].

213 Figure 2 shows the decision boundary of the H-LSVM model on the synthetic banana dataset ² for
214 different pruning rates (ρ). The H-LSVM parameters were $\lambda = 10^{-5}$ and $\delta = 10^{-3}$. Blue and light blue
215 points correspond to positive and negative samples. The problem is not linearly separable. The classification
216 accuracy of the linear SVM is 54.44% while the Gaussian Kernel SVM achieves a classification rate of 90.60%
217 and 1,152 support vectors. Figure 2(a) shows the H-LSVM decision boundary when no pruning is applied

²Dataset available at <http://www.fml.tuebingen.mpg.de/Members/raetsch/benchmark>



(a)



(b)

Figure 2: Best viewed in color. The decision boundary of the H-LSVM model in the banana dataset. Figure 2(a) shows the decision boundary when no pruning is applied. Figure 2(b) shows the decision boundary after a pruning process with $\rho = 0.1$.

218 ($\rho = 0.0$). The model is clearly overfitted. Figure 2(b) shows the H-LSVM decision boundary when pruning
 219 is applied ($\rho = 0.1$). This model only needs to evaluate at most 10 hyperplanes to classify a new pattern
 220 thus achieving a classification rate of 90.00%. In this case, H-LSVM obtains the same classification accuracy
 221 but with a classification time two orders of magnitude lower than the non-linear case.

222 2.1. Pseudocode

223 Once the model parameters $T, k, \epsilon_P, \delta, \tau$ have been fixed and the parameters λ and ρ have been estimated
 224 in the validation phase, the H-LSVM training procedure can be summarized as follows,

- 225 1. Select randomly $(1 - \rho)N$ samples from the initial training set S to form the subset S_0 . The remaining
 226 ρN samples, subset P , is used by the pruning algorithm.

- 227 2. Initialize the weight of each pattern in S_0 as described in Equation 7.
- 228 3. Train recursively the **H-LSVM Tree** following the steps given in Figure 3.
- 229 4. As a result of the H-LSVM tree construction, a set of N_H hyperplanes $\{\vec{w}_n, b_n\}_{n=1}^{N_H}$ is obtained.
- 230 5. **Pruning step:** If $\rho > 0$ apply the pruning algorithm on the set P to get $\left\{\left(\vec{w}_n, \tilde{b}_n\right)\right\}_{n=1}^{\tilde{N}_H}$ where
- 231 $\tilde{N}_H \leq N_H$; otherwise, set $\left\{\left(\vec{w}_n, b_n\right)\right\}_{n=1}^{N_H} = \left\{\left(\vec{w}_n, \tilde{b}_n\right)\right\}_{n=1}^{\tilde{N}_H}$.
- 232 6. **Prediction step:** Let \vec{x} be a new sample and the H-LSVM tree defined by $\left\{\left(\vec{w}_n, \tilde{b}_n\right)\right\}_{n=1}^{\tilde{N}_H}$. The
- 233 target \tilde{y} of the pattern \vec{x} is calculated as the majority class in the leaf node of the tree associated to \vec{x} .

234 3. Training and Prediction Complexity

235 In this section we analyze the training and classification cost as a function of the number of operations
 236 needed by the linear SVM, the non-linear SVM and the proposed H-LSVM method. As already mentioned,
 237 the main advantage of the H-LSVM method is the speeding up of the prediction phase of non-linear SVMs.
 238 SVMs have very good results in performance in off-line problems, but when they are placed in a real time
 239 operation, such as the credit card fraud detection, they are not viable. Thus, focal attention has to be placed
 240 on prediction complexity. Training complexity of H-LSVM is also provided for completeness.

241 3.1. Training Complexity

242 The linear SVMs were trained using the popular LIBLINEAR classification package [36]. The algorithm
 243 behind LIBLINEAR is coordinate descent in the dual SVM formulation [37]. As pointed out by Menon [38],
 244 this algorithm is very attractive because it converges in only $O\left(\log \frac{1}{\epsilon}\right)$ iterations, ϵ being the optimization
 245 tolerance. Menon's experiments show that this algorithm achieves a lower generalization error solution faster
 246 than Pegasos. However, for large-scale datasets Pegasos' training time *decreases* to get a fixed generalization
 247 error [39] while this is not clearly true for LIBLINEAR. The use of the LIBLINEAR package does not affect
 248 to our analysis focused on the classification complexity. The non-linear SVMs have been trained using the
 249 SMO algorithm whose training cost is $O(N^2 d)$ [28] using N d -dimensional patterns. A detailed analysis of
 250 these costs can be found in [38].

251 The H-LSVM cost is that of training as many linear SVMs as nodes in the H-LSVM tree via the
 252 Weighted-Pegasos algorithm. More precisely, if the H-LSVM decision tree has N_H internal nodes, the
 253 training complexity is given by the cost of training N_H linear SVMs with the Weighted-Pegasos algorithm.
 254 Then, considering that the number of iterations needed by the Weighted-Pegasos algorithm to achieve
 255 a solution with tolerance ϵ is $O\left(\frac{1}{\lambda\epsilon}\right)$ and the cost per iteration is $O(kd)$, the total cost of H-LSVM is
 256 $O\left(\frac{N_H kd}{\lambda\epsilon}\right)$. For simplicity, the tolerance ϵ is fixed for every node in the tree, but as suggested by Shalev-
 257 Shwartz and Srebro [39], it could be adapted as a function of the number of training samples n_i reaching
 258 the i -th node to get some fixed generalization error in each node.

```

INPUT:  $S_0, \lambda, T, k, \epsilon_P, \delta, \tau$ 
 $I_0 = I(H_0)$ 
if  $I_0 = 0$  then
    FINISH {Homogeneous node}
end if
if  $(\frac{|S_0|}{N} > \delta)$  then
     $\{\vec{w}, b\} = \text{Weighted-Pegasos}(S_0, \lambda, T, k, \epsilon_P, \tau)$ 
else
    FINISH {There are not enough patterns.}
end if
if  $I(\vec{w}, b) \geq I_0$  then
    FINISH {Cannot find any split}
end if
 $S_{H^l} = \{\vec{x} \in S_0 \mid \vec{w} \cdot \vec{x} + b \leq 0\}$ 
 $S_{H^r} = \{\vec{x} \in S_0 \mid \vec{w} \cdot \vec{x} + b > 0\}$ 
if  $|S_{H^l}| > 0$  then
    Compute the weight of each pattern in  $S_{H^l}$  using Equation 7 where  $H_k = H^l$ 
     $\text{H-LSVM.Tree}(S_{H^l}, \lambda, T, k, \epsilon_P, \delta, \tau)$ 
end if
if  $|S_{H^r}| > 0$  then
    Compute the weight of each pattern in  $S_{H^r}$  using Equation 7 where  $H_k = H^r$ 
     $\text{H-LSVM.Tree}(S_{H^r}, \lambda, T, k, \epsilon_P, \delta, \tau)$ 
end if
OUTPUT:  $\{(\vec{w}_n, b_n)\}_{n=1}^{N_H}$ 

```

Figure 3: H-LSVM Tree Construction.

	Training	Classification
<i>Linear SVM</i>	$Nd \log\left(\frac{1}{\epsilon}\right)$	d
<i>SMO-SVM</i>	$N^2 \times d$	$n_{SV} \times n_K(d)$
<i>H-LSVM</i>	$\frac{N_H k d}{\lambda \epsilon}$	$N_H^P(\vec{x}) \times d$

Table 1: Number of operations needed to train a set S of N patterns in a d -dimensional space (*Training* column) and to classify a new pattern (*Classification* column) by Linear SVM, SVM-SMO and the H-LSVM algorithm. λ : regularization parameter in Pegasos formulation. ϵ : optimization tolerance. n_{SV} : number of support vectors of the non-linear SVM model. $n_K(d)$: operations are needed to compute the kernel between each support vector and the test pattern. N_H : total number of internal nodes in the H-LSVM tree. n_i : number of training samples which reach the node i in the H-LSVM tree. $N_H^P(\vec{x})$: number of nodes encountered by pattern \vec{x} in the H-LSVM tree.

259 Table 1 (column *Training*) shows the training time complexities of the three algorithms: linear SVM,
260 SVM-SMO and H-LSVM. The H-LSVM cost is highly dependent on each dataset as it is determined by
261 the structure of the tree (N_H). As expected, the lowest training cost corresponds to the linear SVM. The
262 comparison between the training times of non-linear SVM and H-LSVM is not straightforward as it depends
263 on the H-LSVM tree structure and the λ and ϵ parameters. H-LSVM would be faster than SMO-SVM in
264 the training phase if $\frac{N_H k}{\lambda \epsilon} \ll N^2$.

265 3.2. Prediction Complexity

266 The cost of classifying a new pattern $\vec{x} \in \mathbb{R}^d$ by a linear SVM is the cost of computing the dot product
267 between the model hyperplane and the pattern to be classified: $O(d)$. In the case of non-linear SVMs, the
268 classification of the pattern \vec{x} is carried out according to: $\sum_{i=1}^{n_{SV}} \alpha_i \times K(\vec{x}_i, \vec{x})$, n_{SV} being the number of
269 support vectors. If $n_K(d)$ is the number of operations needed to compute $K(\vec{x}_i, \vec{x})$, the SVM prediction
270 complexity is $n_{SV} \times n_K(d)$. The proposed H-LSVM algorithm needs to find the leaf of the tree for the pattern
271 \vec{x} which leads to $N_H^P(\vec{x}) \times d$ operations, $N_H^P(\vec{x})$ being the number of internal nodes –oblique hyperplanes–
272 evaluated by the algorithm until the pattern \vec{x} reaches a leaf in the tree.

273 The summary of the number of operations needed by each algorithm to classify a new pattern \vec{x} is given
274 in Table 1 (column *Classification*).

275 Obviously the lowest classification cost corresponds to the linear SVM but it will be shown in Section 5
276 that the linear model is not usually competitive enough for real-world datasets. As regards the non-linear
277 models, it is reasonable to assume that the number of kernel operations $n_K(d)$ is at least d . In that case,
278 H-LSVM has the lowest cost if the number of node evaluations needed to classify the pattern \vec{x} , $N_H^P(\vec{x})$, is
279 lower than the number of support vector encountered by SVM, n_{SV} . In Section 5, the values of n_{SV} and
280 $N_H^P(\vec{x})$ for real-world datasets are given, and it is shown that, in practice, the number of evaluations needed
281 by H-LSVM is indeed several orders of magnitude lower.

282 **4. Generalization Error Bound for the H-LSVM Algorithm**

283 In this section we provide a generalization error bound for the H-LSVM algorithm. First, we show
 284 that the H-LSVM learning algorithm always converges and produces a decision tree as a final model. The
 285 number of nodes to be generated is finite and upper bounded by the number of training samples because of
 286 the stopping criteria commonly used in learning decision tree schemes: the tree expansion is finished when
 287 there is no improvement in the impurity measure or when there are not enough patterns in a node. The
 288 convergence properties of the model can be obtained by considering each node separately and applying the
 289 Pegasos convergence bounds [13] which hold in the weighted version.

290 The generalization error bound is obtained based on the results given by Golea et al. in [40]. Although
 291 other bounds for decision trees have been proposed in the literature [41, 42], some of them tighter than
 292 those of Golea et al., the latter has been considered in this paper because of its simplicity and its explicit
 293 dependence on the decision tree parameters, favoring the understandability of the empirical results obtained
 294 in Section 5.6. Among the alternative bounds, the work of Shah [42] based on the Sample Compression (SC)
 295 paradigm deserves a special mention because it generally yields tighter bounds and sparse models. These
 296 bounds assume axis-parallel decision trees and their application to H-LSVM trees is not straightforward.
 297 The formulation of the SC bounds for oblique decision trees is a direction for future work which might
 298 also help to alleviate (or even eliminate) the cost of the pruning phase by using these bounds to guide the
 299 learning process in a similar way as in [42, 43, 44, 45].

300 Suppose a two-class decision tree T whose internal decision nodes are labeled with boolean functions
 301 from some class \mathcal{U} and whose leaves are labeled as -1 or 1 . The bounds obtained depend on the effective
 302 number of leaves L_{eff} , a data-dependent quantity which reflects how uniformly the training data covers the
 303 tree's leaves and which can be considerably smaller than the total number of leaves in the tree, L [46]. This
 304 bound is different from the Vapnik–Chervonenkis one, which depends on the total number of leaves in the
 305 tree [47, 48].

306 Formally, let $P = (P_1, \dots, P_L)$ the probability vector which represents the probability of a pattern \vec{x}
 307 reaching leaf i for $i = 1 \dots L$. Then, the quadratic distance between the probability vector P and the uniform
 308 probability vector $U = (1/L, \dots, 1/L)$ is given by $\rho(P, U) = \sum_{i=1}^L (P_i - 1/L)^2$ and the effective number of
 309 leaves in the tree is defined by $L_{\text{eff}} \equiv L(1 - \rho(P, U))$.

310 A bound of misclassification probability under the distribution \mathcal{D} , $P_{\mathcal{D}} [T(\vec{x}) \neq y]$, can be estimated using
 311 the following theorem [40]:

312 **Theorem 1.** *For a fixed $\xi > 0$, there is a constant c that satisfies the following. Let \mathcal{D} be a distribution*
 313 *on $X \times \{-1, +1\}$. Consider the class of decision trees of a depth of up to D , with decision functions in \mathcal{U} .*
 314 *With a probability of at least $1 - \xi$ on the training set S (of size N), every decision tree T that is consistent*
 315 *with S has*

$$P_{\mathcal{D}} [T(\vec{x}) \neq y] \leq c \left(\frac{L_{\text{eff}} \text{VCdim}(\mathcal{U}) \log^2 N \log D}{N} \right)^{\frac{1}{2}}$$

316 where L_{eff} is the effective number of leaves of T and VCdim is the Vapnik Dimension.

317 The H-LSVM algorithm is in line with this framework identifying the class \mathcal{U} with the Linear SVM. It
 318 is known that the Vapnik Dimension of a hyperplane in a d -dimensional space is $(d + 1)$ [49] therefore, the
 319 error bound for the H-LSVM method is reformulated as,

320 **Lemma 2.** For a fixed $\xi > 0$, there is a constant c that satisfies the following. Let \mathcal{D} be a distribution on
 321 $X \times \{-1, +1\}$. Consider the class of decision trees of a depth of up to D , with H-LSVM decision functions.
 322 With a probability of at least $1 - \xi$ on the training set S (of size N), every decision tree T that is consistent
 323 with S has

$$P_{\mathcal{D}} [T(\vec{x}) \neq y] \leq c \left(\frac{L_{\text{eff}} (d + 1) \log^2 N \log D}{N} \right)^{\frac{1}{2}} \quad (10)$$

324 where L_{eff} is the effective number of leaves of T .

325 In practice it is quite difficult to have a consistent tree with the training data S . In that case, a bound
 326 of the misclassification probability can be obtained as a function of the misclassification probability in S ,
 327 $P_S [T(\vec{x}) \neq y]$. Now, the probability vector P is reformulated according to the training set as:

$$P'_i = \frac{P_i P_S [T(\vec{x}) = y \mid \vec{x} \text{ reaches leaf } i]}{P_S [T(\vec{x}) = y]}$$

328 By applying the theorem given in [40] for the particular case of the H-LSVM tree, we obtain the following
 329 result,

330 **Lemma 3.** For a fixed $\xi > 0$, there is a constant c that satisfies the following. Let \mathcal{D} be a distribution
 331 on $X \times \{-1, +1\}$. Consider the class of decision trees of a depth of up to D with H-LSVM internal node
 332 decision functions. With a probability of at least $1 - \xi$ on the training set S (of size N), every decision tree
 333 T has

$$P_{\mathcal{D}} [T(\vec{x}) \neq y] \leq P_S [T(\vec{x}) \neq y] + c \left(\frac{L'_{\text{eff}} (d + 1) \log^2 N \log D}{N} \right)^{\frac{1}{3}} \quad (11)$$

334 where c is a universal constant, and $L'_{\text{eff}} = L(1 - \rho(P', U))$ is the empirical effective number of leaves of
 335 T .

336 Therefore, for a given training set S of N patterns, the parameters of the tree which determine the error
337 bound for the H-LSVM algorithm are the depth D of the tree and the effective number of leaves L_{eff} : the
338 lower these parameters are, the better generalization error. The parameter values and an estimation of the
339 model complexity according to Equation 11 are given in Section 5.

340 5. Experiments

341 The aim of the experiments described in the following subsections is fourfold:

- 342 • Compare H-LSVM with linear SVMs and non-linear SVMs in terms of classification accuracy and
343 prediction complexity (Section 5.2).
- 344 • Compare H-LSVM with Zapién’s [21, 24] and Adaboost [50] algorithms in terms of classification
345 accuracy and prediction complexity (Sections 5.3 and 5.4).
- 346 • Analyze numerically the H-LSVM scalability (Section 5.5).
- 347 • Analyze numerically the H-LSVM error bound studied in Section 4 (Section 5.6).

348 The H-LSVM has been implemented in C language and the code can be found at:

349 <https://sites.google.com/site/irenerodriguezlujan/HLSVM-1.1.zip>.

350 As the H-LSVM algorithm has been designed for binary classification domains, the experiments have
351 been conducted in large-scale binary classification problems. We have considered the large-scale datasets
352 used by Keerthi et al. [2]: *IJCNN*, *Shuttle*, *M3VO* and *Vehicle*. The *Shuttle* dataset has been converted
353 to a binary classification problem by differentiating class 1 from the rest. In the same way, the *Vehicle*
354 dataset has been reformulated as a binary classification task consisting of differentiating class 3 from the
355 rest. The *M3VO* dataset corresponds to differentiate digit 3 from all the other digits in the *MNIST* problem.
356 An extension of the *MNIST* dataset with 8,100,000 patterns (*MNIST8m*) has also been included since it
357 represents a very large-scale classification problem. Again, the classification of digit 3 from all of the others
358 has been considered (*M3VOM8*). In order to compare the performance of our method with the Zapién
359 algorithm [21, 24], initially we chose the binary datasets used in this work: *Heart* and *Faces*. However
360 in the case of the *Heart* dataset, the classification accuracy of the linear SVM is the same as that of the
361 non-linear SVM and thus, we decided not to include this dataset in our experiments. Finally, we added the
362 binary version of the covtype dataset (class 2 versus others) because of its large number of patterns. The
363 characteristics of the datasets are shown in Table 2 as well as the repositories where they are available.

364 In most of the datasets (*IJCNN*, *Shuttle*, *M3VO* and *Vehicle*), the training and test subsets are given
365 beforehand. In the *Faces* dataset, we followed the experimental setup described in [52] which uses two

	# Train	# Test	# Feat.	Repository
<i>IJCNN</i>	49,990	91,701	22	LIBSVM [51]
<i>Shuttle</i>	43,500	14,500	9	LIBSVM [51]
<i>M3VO</i>	60,000	10,000	780	LIBSVM [51]
<i>M3VOM8</i>	810,000	7,290,000	784	LIBSVM [51]
<i>Vehicle</i>	78,823	19,705	100	LIBSVM [51]
<i>Faces</i>	8,525	4,263	576	http://www.cs.ubc.ca/~pcarbo/#data
<i>Binary Covtype</i>	522,910	58,102	54	LIBSVM Repository [51]

Table 2: Binary datasets used to compare H-LSVM with linear SVMs and non-linear SVMs.

366 thirds of the observations for the training and the rest as a testing set. Moreover, data was normalized to
367 minimum and maximum feature values. We run the experiments on 10 different randomly chosen training-
368 test partitions of the dataset. In the case of the *M3VOM8* and *Covtype* datasets, we have tried to use as
369 many training patterns as possible in order to simulate a large-scale system with a high number of support
370 vectors ³. Then, the first 810,000 patterns in the *M3VOM8* dataset were used for training and the remaining
371 samples for test. In the *Covtype* dataset, according to the experiments carried out in [9, 53], 9/10 of the
372 samples for training and the remaining patterns for test. In both cases, the experiments were run on 10
373 different randomly chosen training-test partitions of the dataset.

374 In all of the experiments, linear SVMs and non-linear SVMs implemented in LIBLINEAR [36] and
375 LIBSVM [51] packages were used. The Gaussian kernel, $k(x_i, x_j) = \exp(-\gamma\|x_i - x_j\|^2)$, was used for
376 non-linear SVMs.

377 5.1. Hyperparameter Tuning

378 Linear SVMs, non-linear SVMs and H-LSVM need to determine the values of a few parameters. In all
379 datasets, except *Covtype*, the hyperparameter selection has been made using a 5-fold cross validation on
380 the training set. The cost parameters in linear SVMs and non-linear SVMs were selected from the grid
381 $10^i, i = -6, \dots, 6$. The γ parameter of the Gaussian kernel was taken from the range $10^i, i = -3, \dots, 3$.
382 Finally, for the H-LSVM model, we fixed the maximum number of Pegasos iterations $T = 10^7$ with a
383 tolerance of $\epsilon_P = 10^{-4}$ and the minimum proportion of patterns needed to split a node δ was chosen as 10^{-i}
384 with $i = \lfloor \log_{10} N \rfloor$ to guarantee that the H-LSVM grows to a sufficient size (pruning is applied if necessary).
385 The regularization parameter λ was chosen from the grid $\frac{10^i}{N}, i = -6, \dots, 6$, N being the number of training

³LIBSVM for the *M3VOM8* dataset did not finish in reasonable time when training with all the available patterns.

	LIBLINEAR	LIBSVM		H-LSVM	
	c	c	γ	λ	ρ
<i>IJCNN</i>	10^1	10^1	10^0	10^{-5}	0
<i>Shuttle</i>	10^2	10^6	10^0	10^{-7}	0.2
<i>M3VO</i>	10^0	10^2	10^{-2}	10^{-4}	0.2
<i>M3Vom8</i>	10^0	10^2	10^{-2}	10^{-5}	0.2
<i>Vehicle</i>	10^{-1}	10^1	10^{-1}	10^{-6}	0.1
<i>Faces</i>	10^{-1}	10^1	10^{-2}	10^{-5}	0.1
<i>Covtype</i>	10^0	10^0	0.346	10^{-7}	0.0

Table 3: Parameters used in the linear SVM, non-linear SVM and H-LSVM models for each binary dataset.

386 samples. The grid was obtained from the equivalence $\lambda = \frac{1}{cN}$ between the LIBLINEAR and LIBSVM cost
387 parameter c and the λ regularizer in H-LSVM. The prune rate ρ took values in $[0.0, 0.1, 0.2]$. For the *Covtype*
388 dataset, we used the non-linear SVM hyperparameters provided in [9]. The resulting parameters for each
389 dataset and each model are given in Table 3.

390 5.2. Results

391 The results in terms of classification error (*Error (%)*) and classification cost are shown in Table 4. In
392 the case of the linear SVM, the number of hyperplane evaluations is shown whereas the number of support
393 vectors is indicated for the non-linear SVM (n_{SV} or *Hyp*). While the classification cost of linear/non-linear
394 SVMs is independent of the test sample, the H-LSVM prediction cost depends on the path of the pattern in
395 the H-LSVM tree. Thus, the mean number of H-LSVM hyperplanes encountered per test sample together
396 with the maximum number of H-LSVM hyperplane evaluations written in parentheses are shown. In those
397 cases in which there were several training/test partitions, the average and standard deviation on the 10 runs
398 of the experiment are shown.

399 The quantification of the performance of the algorithms considering the linear and non-linear SVMs as
400 the points of reference is given by the quantities *Relative Error (RE)* and *Relative Complexity (RC)*,

$$RE = \frac{e_{\text{LSVM}} - e}{e_{\text{LSVM}} - e_{\text{SVM}}} \quad (12)$$

$$RC = \frac{Hyp - 1}{n_{SV} - 1}, \quad (13)$$

401 where e represents the classification error rate. A value equals 0 at these magnitudes *RE/RC* indicate
402 that the classification accuracy/complexity is the same as that of the linear SVM while a value of 1 represents

403 the equivalence with the non-linear case. Therefore, it would be desirable to have a *Relative Error* close to
404 1 and a *Relative Complexity* close to 0.

405 As expected, the classification results of the non-linear SVMs are greater than those of the linear SVM
406 and H-LSVM. However, the classification accuracy of the H-LSVM is significantly better than that of the
407 linear model in all cases. These results are not surprising because the proposed H-LSVM method is simpler
408 than the non-linear SVM but more sophisticated than linear SVMs. While the classification error of H-
409 LSVM is closer to that of the non-linear SVM in most cases, the H-LSVM classification accuracy is closer
410 to the linear model for the *Faces* and *M3VO* datasets. Nevertheless, in the case of the *Faces* dataset the
411 H-LSVM model represents an improvement of 41% in respect to the linear SVM and it will be shown later
412 that it yields significantly better results than the Zapién et al. [21] algorithm. These results show that the
413 non-linear SVMs cannot be approximated by the proposed method in certain domains. It is worth pointing
414 out that H-LSVM outperforms the non-linear SVM in the *Covtype* dataset. Although, the classification
415 error obtained for the non-linear SVM is comparable to the results reported in [9], a thorough search of
416 the non-linear SVM parameters might provide better results. Unfortunately, to apply the hyperparameter
417 procedure described in Section 5.1 is unfeasible because of the size of the dataset and the number of support
418 vectors.

419 Our main interest is not having the best classification error rates but providing a method capable of
420 classifying a pattern in few milliseconds while obtaining a competitive performance. In this respect, the
421 non-linear SVM needs the largest number of operations in prediction while the lowest cost is that of the
422 linear SVM. However, the performance of the linear SVM can be extremely poor as in the *IJCNN* or *Covtype*
423 datasets. The classification complexity of H-LSVM is between these two models: it is higher than that of the
424 linear SVM –in the worst case it increases the cost of the linear model in one order of magnitude– but much
425 lower than the cost of the non-linear SVM –H-LSVM can accelerate the prediction cost of the non-linear
426 SVM even by a factor of 10^4 as in the case of the *M3Vom8* and *Covtype* datasets. In fact, the *Relative*
427 *Complexity* is lower than 10^{-1} in all cases.

428 5.3. Results: Comparison with SVM Trees Algorithm

429 Having compared H-LSVM to baseline models, we can contrast the results with the Zapién decision
430 tree [21]. As mentioned above, only one of the two binary classification problems used in this work cannot
431 be classified accurately by a linear SVM (*Faces*). Despite the fact that our method has been designed for
432 binary classification problems, the performance of our model in the multiclass USPS dataset was measured.
433 The USPS dataset for handwritten text recognition is available in the LIBSVM Repository [51]. It consists
434 of 7291 training samples and 2007 test samples. Each example is described by 256 features. Following
435 the methodology described in [52], we normalized the data to minimum and maximum feature values and
436 we applied *one against one* approach (1A1) for the multiclass problem. The 1A1 strategy consists of

IJCNN			
	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	7.82	1.01	2.36
<i>n_{SV} or Hyp</i>	1	3,154	7.28 (16)
<i>RE / RC</i>	0 / 0	1 / 1	0.80 / $2.0 \cdot 10^{-3}$

Shuttle			
	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	2.21	0.062	0.10
<i>n_{SV} or Hyp</i>	1	66	5.18 (12)
<i>RE / RC</i>	0 / 0	1 / 1	0.98 / $6.43 \cdot 10^{-2}$

M3VO			
	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	2.09	0.33	1.79
<i>n_{SV} or Hyp</i>	1	2,873	3.15 (8)
<i>RE / RC</i>	0 / 0	1 / 1	0.17 / $7.49 \cdot 10^{-4}$

M3VOM8			
	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	3.90	0.03	1.43 ± 0.02
<i>n_{SV} or Hyp</i>	1	13,471	4.73 ± 0.003 (11.00 ± 0.00)
<i>RE / RC</i>	0 / 0	1 / 1	$0.64 / 2.77 \cdot 10^{-4}$

Vehicle			
	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	14.18	11.88	12.61
<i>n_{SV} or Hyp</i>	1	23,642	2.84 (10)
<i>RE / RC</i>	0 / 0	1 / 1	$0.68 / 7.78 \cdot 10^{-5}$

Faces			
	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	8.81 ± 0.35	2.97 ± 0.24	6.39 ± 0.43
<i>n_{SV} or Hyp</i>	1	$1,260.3 \pm 14.81$	2.59 ± 0.06 (5.30 ± 0.15)
<i>RE / RC</i>	0 / 0	1 / 1	$0.41 / 1.26 \cdot 10^{-3}$

Covtype			
	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	23.66 ± 0.21	18.57 ± 0.20	11.39 ± 0.08
<i>n_{SV} or Hyp</i>	1	$245,687.2 \pm 167.8$	12.93 ± 0.087 (44.00 ± 1.08)
<i>RE / RC</i>	0 / 0	1 / 1	$2.41 / 4.86 \cdot 10^{-5}$

Table 4: Test error rate (*Class. Err (%)*) and classification complexity (*n_{SV} or Hyp*) of Linear SVMs, non-linear SVMs and H-LSVM. The mean number of hyperplane evaluations per test sample is indicated for linear SVMs and H-LSVM. The maximum number of H-LSVM hyperplane evaluations is shown in parentheses. In the case of non-linear SVMs, the number of support vectors (*n_{SV}*) is shown. The reference measures RE and RC (Equations 12 and 13) are also provided.

Faces

	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>SVM Trees</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	8.81	2.97	8.99	6.39
<i>n_{SV} or Hyp</i>	1	1260.3	4	2.59 (5.30)
<i>RE / RC</i>	0 / 0	1 / 1	-0.03 / 0.002	0.41 / 0.001

USPS

	<i>Linear SVM</i>	<i>Non-linear SVM</i>	<i>SVM Trees</i>	<i>H-LSVM</i>
<i>Class. Err (%)</i>	8.67	4.53	6.24	5.38
<i>n_{SV} or Hyp</i>	1	1521	49	64.77 (117)
<i>RE / RC</i>	0 / 0	1 / 1	0.59 / 0.03	0.79 / 0.04

Table 5: Comparison of the SVM Trees method by Zapién et al. [21, 52] and H-LSVM. Misclassification error (*Class. Err (%)*) and the mean number of hyperplane evaluations per test sample (*Hyp*) are shown for both methods and for the linear and non-linear SVMs (*n_{SV}*). The maximum number of H-LSVM hyperplane evaluations is indicated in parentheses. The number of hyperparameter evaluations was computed as the sum of the hyperplanes evaluated in every binary classifier. The *Relative Error (RE)* and the *Relative Complexity (RC)* of the SVM Trees method and H-LSVM are also given.

437 training a classifier for every pair of classes and classifying a new pattern based on majority voting. The
 438 hyperparameters were chosen by using 5-CV as described above. For the linear SVM, the cost parameter
 439 was set to $c = 1$, for the non-linear SVM $c = 10^1$ and $\gamma = 10^{-2}$ and for the H-LSVM algorithm $\lambda = 10^{-5}$
 440 and $\rho = 0$. The results in terms of the misclassification error and classification cost for both methods are
 441 given in Table 5. The performance of the Zapién method has been extracted from [21, 52].

442 In both cases H-LSVM is superior in terms of classification accuracy whereas the classification cost is in
 443 the same order of magnitude. Specifically, their classification complexity is quite similar in the *Faces* dataset
 444 but the SVM Trees algorithm is slightly faster for the *USPS* database. In any case, the classification cost of
 445 both algorithms is of the same order of magnitude. In summary, the H-LSVM decision tree, expanding both
 446 children of each node as well as the weighted patterns used in the linear SVM training, provide advantages
 447 in terms of classification accuracy while maintaining the classification cost. It is also worth noting that in all
 448 cases the maximum depth of the tree is lower than the number of internal nodes (the number of linear SVMs),
 449 which means that the structure of the tree is far from being a cascade of classifiers as in [21, 24, 25, 52].
 450 The superiority of the H-LSVM tree against the Zapién’s algorithm in terms of classification accuracy is not
 451 surprising given that, as already mentioned in Section 1, the hypothesis class of H-LSVM (disjunctions of
 452 conjunctions) is more general than that of the SVM Trees algorithm (conjunctions) [22]. What is more, this
 453 difference can be quantified taking into account that the number of decision tree skeletons with k decision

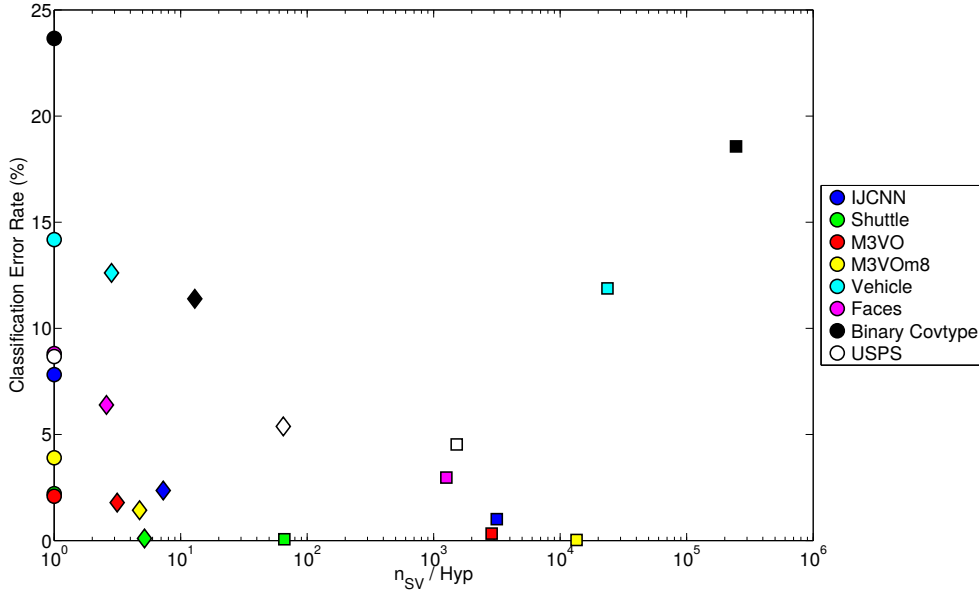


Figure 4: Best viewed in color. Classification complexity (n_{SV} / Hyp) versus classification error rate for different datasets. • Linear SVM ■ Non-linear SVM ♦ H-LSVM.

454 nodes is given by the k -th Catalan Number [19, 23] in comparison to the only one possibility for the Zapién’s
 455 method.

456 In order to visualize the trade-off between the misclassification error vs. classification cost, Figure 4 shows
 457 the dependence between these two magnitudes for the linear SVM, non-linear SVM and H-LSVM. The x-
 458 axis represents the number of support vectors or hyperplanes encountered by each method in logarithmic
 459 scale. The y-axis shows the classification error rate. Each dataset is represented by a color according to
 460 the legend. Circles, squares and diamonds represent the linear SVM, non-linear SVM and H-LSVM models,
 461 respectively. The lower left-hand area is associated to the best scenario: the lowest classification error
 462 and the lowest classification complexity. In this Figure, three clusters can be easily identified according to
 463 the classifier (circles, squares and diamonds). Clearly, the non-linear SVMs have the highest classification
 464 complexity while the H-LSVM cost is closer to the linear one. Looking at the classification error, in all cases
 465 the non-linear SVM is superior – except the *Covtype* dataset – and the H-LSVM effectiveness is greater than
 466 that of the linear model.

467 Finally, to give an idea of the quality of the H-LSVM algorithm with regard to the prediction time,
 468 Table 6 shows the time in seconds needed by a linear SVM, a non-linear SVM and H-LSVM to classify
 469 a new pattern in an Intel(R) Core(TM) i7 CPU 920 at 2.67GHz. The training time is also included for
 470 completeness. As expected, the lowest training and test times correspond to the linear SVM. As regards the

	Linear SVM		Nonlinear SVM		H-LSVM	
	<i>Training</i>	<i>Test</i>	<i>Training</i>	<i>Test</i>	<i>Training</i>	<i>Test</i>
IJCNN	$4.00 \cdot 10^{-1}$	$1.22 \cdot 10^{-7}$	$2.54 \cdot 10^1$	$2.17 \cdot 10^{-4}$	$1.44 \cdot 10^3$	$2.21 \cdot 10^{-6}$
Shuttle	$4.77 \cdot 10^{-1}$	$1.09 \cdot 10^{-7}$	$4.83 \cdot 10^0$	$3.91 \cdot 10^{-6}$	$2.37 \cdot 10^4$	$2.76 \cdot 10^{-6}$
M3VO	$4.76 \cdot 10^0$	$6.47 \cdot 10^{-7}$	$6.40 \cdot 10^3$	$3.27 \cdot 10^{-3}$	$6.25 \cdot 10^3$	$3.61 \cdot 10^{-5}$
M3Vom8	$3.38 \cdot 10^2$	$3.00 \cdot 10^{-6}$ $\pm 1.28 \cdot 10^{-8}$	$7.71 \cdot 10^4$	$2.77 \cdot 10^{-2}$ $\pm 5.88 \cdot 10^{-6}$	$7.19 \cdot 10^4$	$2.88 \cdot 10^{-5}$ $\pm 4.90 \cdot 10^{-8}$
Vehicle	$2.82 \cdot 10^0$	$4.83 \cdot 10^{-7}$	$1.85 \cdot 10^3$	$9.02 \cdot 10^{-3}$	$2.99 \cdot 10^4$	$2.67 \cdot 10^{-6}$
Faces	$1.04 \cdot 10^0$ $\pm 5.28 \cdot 10^{-2}$	$2.41 \cdot 10^{-6}$ $\pm 5.38 \cdot 10^{-9}$	$2.39 \cdot 10^1$ $\pm 9.61 \cdot 10^{-2}$	$2.52 \cdot 10^{-3}$ $\pm 1.66 \cdot 10^{-5}$	$4.70 \cdot 10^3$ $\pm 5.98 \cdot 10^1$	$1.36 \cdot 10^{-5}$ $\pm 4.71 \cdot 10^{-7}$
Covtype	$6.65 \cdot 10^1$ $\pm 1.05 \cdot 10^{-1}$	$1.30 \cdot 10^{-7}$ $\pm 1.76 \cdot 10^{-9}$	$2.10 \cdot 10^4$ $\pm 1.31 \cdot 10^2$	$1.95 \cdot 10^{-2}$ $\pm 7.43 \cdot 10^{-5}$	$3.13 \cdot 10^4$ $\pm 8.47 \cdot 10^1$	$6.83 \cdot 10^{-6}$ $\pm 3.86 \cdot 10^{-8}$
USPS	$5.46 \cdot 10^0$	$3.58 \cdot 10^{-6}$	$5.21 \cdot 10^0$	$1.17 \cdot 10^{-3}$	$5.88 \cdot 10^3$	$1.59 \cdot 10^{-4}$

Table 6: Training and testing time in seconds required by LIBLINEAR, LIBSVM and H-LSVM.

471 training cost discussed in Section 3.1, the differences between the training cost of the non-linear SVM and
472 H-LSVM are given by the structure of the H-LSVM tree. Therefore, depending on the dataset either the
473 non-linear SVM or H-LSVM is faster in the training phase. Focusing on the aim of speeding up the non-
474 linear SVM prediction cost, the H-LSVM classification time is always in the order of tenths of milliseconds
475 at most and significantly lower than those of the non-linear SVM.

476

477 Several techniques based on the use of linear SVMs on the manifold coordinates have come out recently
478 [15, 16, 17]. In particular, the Locally Linear SVM (LLSVM) model proposed by Ladicky et al. [17] reports
479 results for the USPS dataset. The classification accuracy of H-LSVM is slightly better than that of the
480 LLSVM and the LLSVM algorithm needs to compute the distance to 100 k-means centroids while H-LSVM
481 evaluates on average 64.77 hyperplanes (maximum 117). That is, both methods are comparable in terms of
482 classification accuracy and prediction complexity.

483 *5.4. Results: Comparison with Adaboost Algorithm*

484 Other natural competitors for H-LSVM are boosting algorithms [54] since they create piecewise linear
485 functions with a good generalization performance [55] and low classification cost. In particular, we have
486 considered the most known boosting algorithm: AdaBoost (Adaptive Boosting) [50]. Decision stumps were
487 used in accordance with the Adaboost algorithm originally proposed by its authors Freund and Schapire
488 [50] and motivated by its successful application in the state-of-the-art Viola-Jones face detection algorithm
489 [56]. The main advantage of Adaboost with decision stumps on competitors is the speed of learning and
490 prediction, which is particularly critical in large-scale problems.

491 Adaboost requires the establishment of the maximum number of weak classifiers H to be used. Since our
492 paper focuses on accelerating the classification times, H was fixed to make the prediction cost of Adaboost
493 comparable to that of H-LSVM. Then, by taking into account that Adaboost needs to evaluate all the weak
494 learners to classify a test pattern and considering that the prediction complexity of each decision stump is
495 $O(1)$, H is computed as the mean number of hyperplanes evaluated by H-LSVM multiplied by the dimension
496 of the patterns. The comparison of both methods in terms of misclassification rate and classification cost
497 as well as the value of the parameter H for each dataset are given in Table 7.

498 The results show that Adaboost has a better performance in the *Shuttle* and *Vehicle* datasets, the differ-
499 ence in classification accuracy being 0.15% at most. However, in some cases such as *IJCNN* and *Covtype*,
500 H-LSVM significantly outperforms Adaboost. On average, Adaboost and H-LSVM have misclassification
501 rates of 7.81% and 5.15%, respectively on all the datasets. Overall, the H-LSVM yields a better perfor-
502 mance/classification speed ratio than Adaboost with decision stumps.

503 *5.5. Numerical Analysis of H-LSVM Scalability*

504 To illustrate the applicability of the H-LSVM algorithm to real large-scale scenarios, we show the scala-
505 bility in the training time and convergence of the test error rates as the number of training samples increases.
506 In this regard, Figures 5a – 5c show the training complexity of H-LSVM in terms of the number of hyper-
507 planes in the H-LSVM tree and the computational time as a function of the number of training samples N .
508 Figure 5d shows the training and test classification accuracies as a function of N . The results represent the
509 average on the 10 training/test partitions of the *Binary Covtype* dataset. In turn, 4 subsets of size 10,000,
510 50,000, 100,000 and 200,000, respectively, have been randomly chosen from each training partition.

511 According to the training cost of H-LSVM presented in Section 3.1, $O(\frac{N_H kd}{\lambda \epsilon})$, by maintaining λ , ϵ and d
512 constant for the different training sizes, the training cost of H-LSVM depends on the subsampling rate k and
513 the number of internal nodes in the H-LSVM tree N_H as $O(N_H k)$. Unfortunately, N_H depends inextricably
514 on the problem in question and thus, a general estimation of N_H based on N cannot be provided. However,
515 it is possible to compare the number of internal nodes in the H-LSVM tree against those corresponding to
516 the best tree (balanced tree) and the worst one (a linear or cascade tree). In this regard, Figures 5a (linear

	Adaboost			H-LSVM		
	<i>Class. Err (%)</i>	<i>H</i>	<i>Cost</i>	<i>Class. Err (%)</i>	<i>Hyp</i>	<i>Cost</i>
IJCNN	6.50	170	170	2.36	7.28	160.16
Shuttle	0.08	50	50	0.10	5.18	46.62
M3VO	2.69	2460	2460	1.79	3.15	2457.00
M3VOm8	3.61 ± 0.01	3710	3710	1.43 ± 0.02	4.73	3708.32
Vehicle	12.46	290	290	12.61	2.84	284.00
Faces	6.39 ± 0.33	1500	1500	6.39 ± 0.43	2.59	1491.84
Binary Covtype	22.95 ± 0.23	700	700	11.39 ± 0.08	12.93	698.22
Average	7.81	1268.57	1268.57	5.15	5.53	1263.74

Table 7: Test error rate (*Class. Err (%)*) and classification complexity (*Cost*) of Adaboost and H-LSVM. The number of weak learners (*H*) are indicated for Adaboost and the mean number of hyperplane evaluations per test sample is indicated for H-LSVM (*Hyp*). Note that the classification cost is computed by considering that the classification complexity of each decision stump is $O(1)$ whereas it is $O(d)$ for each hyperplane in the H-LSVM tree.

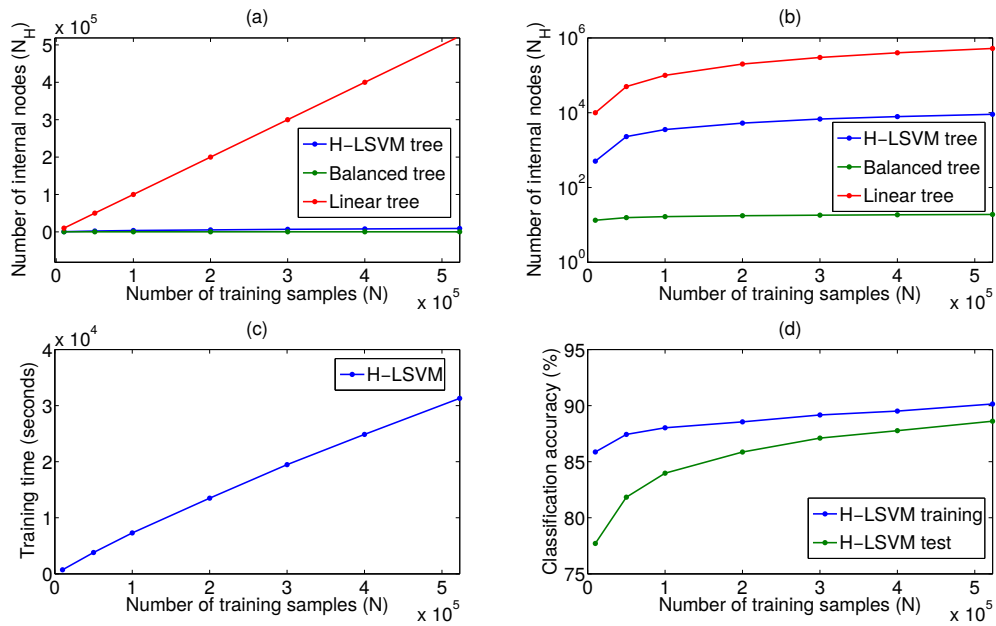


Figure 5: H-LSVM training complexity and classification rate convergence as a function of the number of training samples (N) for the *Covtype* dataset. Figure 5a: number of internal nodes N_H in the H-LSVM tree. Reference values corresponding to a balanced and cascade tree are also included. Figure 5b: Figure 5a using logarithmic axis in the y-axis. Figure 5c: training time of the H-LSVM algorithm. Figure 5d: training and test classification accuracies.

517 y-axis) and 5b (logarithmic y-axis) include the number of internal nodes associated with the balanced
 518 decision tree ($\log_2(N)$) and those encountered in the cascade structure (N). In this case, the complexity of
 519 the H-LSVM tree is closer to that of the balanced tree. Similar results are expected for the other datasets
 520 since in all cases the maximum depth of the tree is much lower than the number of internal nodes.

521 Furthermore, the variability in the distribution of the training samples throughout the decision tree also
 522 affects the exact computation of a general training cost of the Pegasos algorithm in each node. Although
 523 the subsample size k is fixed at the beginning of the algorithm – in this experiment, it was set 50,000 –,
 524 the effective subsampling size in each node is determined online as the minimum between k and the number
 525 of samples reaching the current node, which is totally linked to each particular dataset. Nevertheless, the
 526 empirical measure of the training time of H-LSVM as a function of the number of training samples N shown
 527 in Figure 5c seems to have a linear growth and, in fact, it has been proven that the polynomial curve fitting
 528 the points with the lowest error is that of degree 1. Again, a comparison with the training complexity of the
 529 balanced and the linear decision trees can be valuable. Under the assumption that the cost of computing
 530 the split of the i -th node is proportional to the number of samples n_i reaching the node, the balanced tree
 531 has a linear cost with respect to N :

$$\sum_{i=0}^{\log_2(N)} i n_i = \sum_{i=0}^{\log_2(N)} i \frac{i}{2^i} = O(N) ,$$

532 while the linear tree has a quadratic dependence:

$$\sum_{i=0}^{N-1} n_i = \sum_{i=0}^{N-1} (N - i) = O(N^2) .$$

533 Therefore, the complexity of H-LSVM training is closer to the best scenario. Finally, Figure 5d reveals
 534 that the gap between training and test errors converges with approximately 300,000 patterns. Although
 535 the classification accuracy in the test set increases with the number of training samples, the improvement
 536 becomes smaller as N grows, especially when N is larger than 300,000 in which case the difference with
 537 respect to the model trained with all the training samples is 0.52%.

538 The preceding results corroborate the applicability of H-LSVM to large-scale scenarios.

539 5.6. Numerical Analysis of H-LSVM Generalization Error Bound

540 Lemma 3 provides a generalization error bound for the H-LSVM method as a function of some data-
 541 dependent parameters according to the equation,

$$P_{\mathcal{D}} [T(\vec{x}) \neq y] \leq P_S [T(\vec{x}) \neq y] + c \left(\frac{L'_{\text{eff}} (d+1) \log^2 N \log D}{N} \right)^{\frac{1}{3}}$$

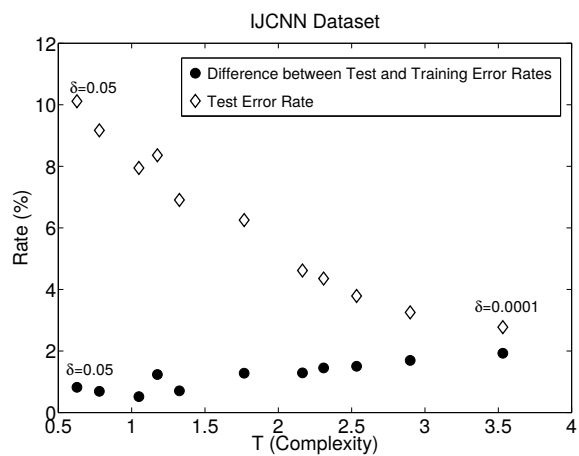
542 which establishes a linear dependence between the difference $P_{\mathcal{D}} [T(\vec{x}) \neq y] - P_S [T(\vec{x}) \neq y]$ and the
 543 complexity of the tree given by $T = \left(\frac{L'_{\text{eff}} (d+1) \log^2 N \log D}{N} \right)^{\frac{1}{3}}$. In this section, an empirical analysis of the
 544 above equation is provided.

545 The misclassification probability under a distribution \mathcal{D} , $P_{\mathcal{D}} [T(\vec{x}) \neq y]$ has been approximated with the
 546 error rate in the test set: $\hat{P}_{\mathcal{D}} [T(\vec{x}) \neq y]$. The range of values of the complexity measure T depends on the
 547 characteristic of each dataset making the comparison between the different datasets impossible. However,
 548 an interesting point of analysis is to determine whether in practice a linear correlation exists between
 549 the difference of the test and training error rates and the complexity of the model T . This relationship
 550 is analyzed for the *IJCNN* and *Faces* datasets by varying the values of the δ parameter to obtain the
 551 values for T , $P_S [T(\vec{x}) \neq y]$ and $\hat{P}_{\mathcal{D}} [T(\vec{x}) \neq y]$. The δ parameter allows the complexity of the model to be
 552 measured and controlled. If δ takes values in the grid $\{\delta_1 > \delta_2 > \dots > \delta_M\}$, the obtained trees \mathcal{T}_{δ_m} verify
 553 $\mathcal{T}_{\delta_1} \subseteq \mathcal{T}_{\delta_2} \subseteq \dots \subseteq \mathcal{T}_{\delta_M}$. In our experiment, the δ grid was: $\{5 \cdot 10^{-2}, 2.5 \cdot 10^{-2}, 1 \cdot 10^{-2}, 7.5 \cdot 10^{-3}, 5 \cdot 10^{-3},$
 554 $2.5 \cdot 10^{-3}, 1 \cdot 10^{-3}, 7.5 \cdot 10^{-4}, 5 \cdot 10^{-4}, 2.5 \cdot 10^{-4}, 1 \cdot 10^{-4}\}$, the prune rate was fixed to $\rho = 0.0$ in both cases
 555 and the λ parameter was selected as in Table 3. The obtained results are shown in Figure 6 in which the
 556 high correlation between the complexity term T and the gap between the training and test errors is shown
 557 by the points representing the difference between the training and test errors. More precisely, the linear
 558 correlation between the H-LSVM tree complexity T and the difference between the test and training error
 559 rates is 0.91 for the *IJCNN* dataset and it is 0.97 for the *Faces* dataset. These high correlations show that
 560 the generalization error bound given in Lemma 3 holds in practice.

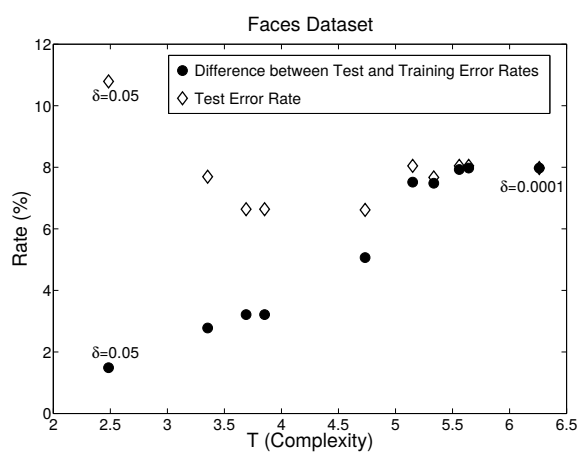
561 Finally, it is interesting to see how the underfitting and overfitting effects are reflected in Figure 6. In
 562 the case of the *IJCNN* dataset, the differences between the test and training error rates are small for the
 563 largest values of δ while the test error rate is the worst. It is a case of underfitting. On the other hand,
 564 the lowest values for δ have the largest differences between the test and training error rates but the test
 565 error rate is the lowest. This scenario is preferable to that with large values of δ . In the *Faces* dataset
 566 the underfitting/overfitting are clearly reflected for high/small δ values, respectively. Regarding how the δ
 567 parameter was chosen in the experiments (see Table 3), it makes sense that the optimal pruning rate for the
 568 *Faces* dataset was $\rho = 0.1$ in order to avoid overfitting.

569 6. Conclusions

570 This paper has presented and analyzed a new classification method for medium and large-scale datasets.
 571 As the application of non-linear SVMs in these problems is prohibitive because it generates a large number
 572 of support vectors, the proposed method takes advantage of the efficiency of linear SVMs to construct a
 573 piecewise linear model. The new algorithm, called a Hierarchical Linear Support Vector Machine (H-LSVM),
 574 is based on the construction of a decision tree whose node splits are Linear Support Vector Machine trained



(a)



(b)

Figure 6: Difference between the test and training error rates and the test error rate as a function of the complexity T of the H-LSVM model. Figure 6(a) IJCNN dataset; Figure 6(b) Faces dataset.

575 with a modified version of the Pegasos algorithm with weighted patterns.

576 Here, we provide a description of the H-LSVM algorithm, an upper bound of the H-LSVM generalization
577 error and an analysis of the H-LSVM prediction cost compared with those of linear SVM and non-linear
578 SVM. The experiments carried out in medium and large datasets show that the H-LSVM algorithm improves
579 the classification accuracy of linear SVMs. Compared with the existing methods based on the construction of
580 a decision tree with linear SVMs as splitting criteria, the H-LSVM model is superior in terms of classification
581 accuracy while maintaining a classification complexity of the same order of magnitude.

582 In summary, the H-LSVM method is an attempt at a solution to the problem of applying SVM technology
583 to industrial settings with high loads in real-time classification. In online industrial environments when
584 decisions have to be taken in a hundredth of a second, non-linear SVMs are just impossible to apply.
585 H-LSVMs may bridge this gap because it is simple and efficient.

586 Acknowledgements

587 The authors would like to thank the anonymous reviewers for their comments that help improve the
588 manuscript. I.R.-L. is supported by an FPU grant from Universidad Autónoma de Madrid, and par-
589 tially supported by the Universidad Autónoma de Madrid-IIC Chair and TIN 2010-21575-C02-01. R.H.
590 acknowledges partial support by ONR N00014-07-1-0741, USARIEM-W81XWH-10-C-0040 (ELINTRIX)
591 and JPL-1396686.

592 References

- 593 [1] I. Steinwart, Sparseness of support vector machines—some asymptotically sharp bounds, *Neural Information Processing*
594 *Systems* 16 (2004) 1069–1076.
- 595 [2] S. S. Keerthi, O. Chapelle, D. DeCoste, Building support vector machines with reduced classifier complexity., *Journal of*
596 *Machine Learning Research* 7 (2006) 1493–1515.
- 597 [3] C. J. C. Burges, B. Schölkopf, Improving the accuracy and speed of support vector learning machines, in: *Advances in*
598 *Neural Information Processing Systems* 9, MIT Press, Cambridge, MA, 1997, pp. 375–381.
- 599 [4] T. Downs, K. E. Gates, A. Masters, Exact simplification of support vector solutions, *Journal of Machine Learning Research*
600 2 (2002) 293–297.
- 601 [5] E. E. Osuna, F. Girosi, Reducing the run-time complexity in support vector machines, *Advances in kernel methods*, MIT
602 Press, Cambridge, MA, USA, 1999, pp. 271–283.
- 603 [6] M. Wu, B. Schölkopf, G. H. Bakir, Building sparse large margin classifiers., in: *ICML '05: Proceedings of the 22th*
604 *international conference on Machine learning*, Vol. 119, 2005, pp. 996–1003.
- 605 [7] R. E. Schapire, Y. Singer, Improved boosting algorithms using confidence-rated predictions, *Machine Learning* 37 (3)
606 (1999) 297–336.
- 607 [8] L. Breiman, Bagging predictors, *Machine Learning* 24 (2) (1996) 123–140.
- 608 [9] R. Collobert, S. Bengio, Y. Bengio, A parallel mixture of SVMs for very large scale problems., *Neural Computation* 14 (5)
609 (2002) 1105–1114.

- 610 [10] H. P. Graf, E. Cosatto, L. Bottou, I. Durdanovic, V. Vapnik, Parallel support vector machines: The cascade SVM, in:
611 Advances in Neural Information Processing Systems, MIT Press, 2005, pp. 521–528.
- 612 [11] C. G. X. L. F.Chang, C. Lu, Tree decomposition for large-scale SVM problems, Journal of Machine Learning Research 11
613 (2010) 2935–2972.
- 614 [12] L. Breiman, J. H. Friedman, R. A. Olshen, C. J. Stone, Classification and Regression Trees, Chapman & Hall, New York,
615 NY, 1984.
- 616 [13] Y. Singer, N. Srebro, Pegasos: Primal estimated sub-gradient solver for SVM, in: In ICML, 2007, pp. 807–814.
- 617 [14] A. Bordes, L. Bottou, P. Gallinari, Sgd-qn: Careful quasi-newton stochastic gradient descent, Journal of Machine Learning
618 Research 10 (2009) 1737–1754.
- 619 [15] X. Zhou, N. Cui, Z. Li, F. Liang, T. S. Huang, Hierarchical gaussianization for image classification., in: ICCV, IEEE,
620 2009, pp. 1971–1977.
- 621 [16] K. Yu, T. Zhang, Y. Gong, Non-linear learning using local coordinate coding, in: Advances in Neural Information Pro-
622 cessing Systems 22, 2009, pp. 2223–2231.
- 623 [17] L. Ladicky, P. Torr, Locally Linear Support Vector Machines, in: ICML '11: Proceedings of the 28th international
624 conference on Machine learning, 2011, pp. 985–992.
- 625 [18] Z. Fu, A. Robles-Kelly, J. Zhou, Mixing linear SVMs for nonlinear classification., IEEE Transactions on Neural Networks
626 21 (12) (2010) 1963–1975.
- 627 [19] K. P. Bennett, N. Cristianini, J. Shawe-Taylor, D. Wu, Enlarging the margins in perceptron decision trees, Machine
628 Learning 41 (2000) 295–313.
- 629 [20] M. Arun Kumar, M. Gopal, A hybrid SVM based decision tree, Pattern Recognition 43 (2010) 3977–3987.
- 630 [21] K. Z. Arreola, J. Fehr, H. Burkhardt, Fast support vector machine classification using linear SVMs., in: ICPR (3), IEEE
631 Computer Society, 2006, pp. 366–369.
- 632 [22] M. Anthony, Generalization error bounds for threshold decision lists, Journal of Machine Learning Research 5 (2004)
633 189–217.
- 634 [23] J. R. Quinlan, R. L. Rivest, Inferring decision trees using the minimum description length principle, Information and
635 Computation 80 (3) (1989) 227–248.
- 636 [24] J. Fehr, K. Z. Arreola, H. Burkhardt, Fast support vector machine classification of very large datasets., in: GfKI, Studies
637 in Classification, Data Analysis, and Knowledge Organization, Springer, 2007, pp. 11–18.
- 638 [25] J. Su, G. Wang, Q. Hu, S. Li, A novel SVM decision tree and its application to face detection., in: SNPD (1), IEEE
639 Computer Society, 2007, pp. 385–389.
- 640 [26] V. Menkovski, I. Christou, S. Efreimidis, Oblique decision trees using embedded Support Vector Machines in classifier
641 ensembles, in: 7th IEEE International Conference on Cybernetic Intelligent Systems, 2008. CIS 2008., 2008, pp. 1–6.
- 642 [27] J. R. Quinlan, C4.5: programs for machine learning, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993.
- 643 [28] J. C. Platt, Fast training of support vector machines using sequential minimal optimization, in: Advances in Kernel
644 Methods: Support Vector Learning, MIT Press, Cambridge, MA, 1998, pp. 185–208.
- 645 [29] K. P. Bennett, J. Blue, A support vector machine approach to decision trees, in: The 1998 IEEE International Joint
646 Conference on Neural Networks Proceedings, 1998. IEEE World Congress on Computational Intelligence., 1998.
- 647 [30] Y. Ben-Haim, E. Tom-Tov, A streaming parallel decision tree algorithm., Journal of Machine Learning Research 11 (2010)
648 849–872.
- 649 [31] S. Viaene, R. A. Derrig, G. Dedene, Cost-sensitive learning and decision making for massachusetts pip claim fraud data.,
650 International Journal of Intelligent Systems 19 (12) (2004) 1197–1215.
- 651 [32] S.-B. Park, S. Hwang, B.-T. Zhang, Mining the risk types of human papillomavirus (hpv) by adacost., in: DEXA, Vol.
652 2736 of Lecture Notes in Computer Science, Springer, 2003, pp. 403–412.

- 653 [33] Y. Huang, S. Du, Weighted Support Vector Machine for classification with uneven training class sizes, in: Proceedings of
654 2005 International Conference on Machine Learning and Cybernetics, Vol. 7, 2005, pp. 4365–4369.
- 655 [34] J. Quinlan, Simplifying decision trees., International Journal of Man-Machine Studies 27 (1987) 221–234.
- 656 [35] W. W. Cohen, Efficient pruning methods for separate-and-conquer rule learning systems., in: IJCAI, 1993, pp. 988–994.
- 657 [36] R. Fan, K. Chang, C. Hsieh, X. Wang, C.-J. Lin, LIBLINEAR: A library for large linear classification., Journal of Machine
658 Learning Research 9 (2008) 1871–1874, software available at <http://www.csie.ntu.edu.tw/~cjlin/liblinear/>.
- 659 [37] C. Hsieh, K. Chang, C. Lin, S. S. Keerthi, S. Sundararajan, A dual coordinate descent method for large-scale linear SVM.,
660 in: ICML '08: Proceedings of the 25th international conference on Machine learning, Vol. 307, 2008, pp. 408–415.
- 661 [38] A. K. Menon, Large-scale support vector machines: Algorithms and theory, Tech. rep., University of California, San Diego
662 (2009).
- 663 [39] S. Shalev-Shwartz, N. Srebro, SVM optimization: inverse dependence on training set size., in: Proceedings of the 25th
664 International Conference on Machine Learning (ICML'08), Vol. 307, 2008, pp. 928–935.
- 665 [40] M. G. Peter, I. Llew Mason, Generalization in decision trees and dnf: Does size matter?, in: Advances in Neural Information
666 Processing Systems, The MIT Press, 1997, pp. 259–265.
- 667 [41] Y. Mansour, D. A. McAllester, Generalization bounds for decision trees., in: COLT, 2000, pp. 69–74.
- 668 [42] M. Shah, Sample compression bounds for decision trees, in: ICML, Vol. 227, 2007, pp. 799–806.
- 669 [43] M. Marchand, M. Sokolova, Learning with decision lists of data-dependent features, Journal of Machine Learning Research
670 6 (2005) 427–451.
- 671 [44] F. Laviolette, M. Marchand, M. Shah, S. Shanian, Learning the set covering machine by bound minimization and margin-
672 sparsity trade-off, Machine Learning 78 (1-2) (2010) 175–201.
- 673 [45] M. Shah, M. Marchand, J. Corbeil, Feature selection with conjunctions of decision stumps and learning from microarray
674 data, IEEE Trans. Pattern Anal. Mach. Intell. 34 (1) (2012) 174–186.
- 675 [46] R. C. Holte, Very simple classification rules perform well on most commonly used datasets, Machine Learning 11 (1993)
676 63–90.
- 677 [47] U. M. Fayyad, K. B. Irani, What should be minimized in a decision tree?, in: Association for the Advancement of Artificial
678 Intelligence (AAAI), 1990, pp. 749–754.
- 679 [48] A. Ehrenfeucht, D. Haussler, M. Kearns, Learning decision trees from random examples needed for learning, Information
680 and Computation 82 (1989) 231–246.
- 681 [49] V. N. Vapnik, Statistical Learning Theory, Wiley-Interscience, 1998.
- 682 [50] Y. Freund, R. E. Schapire, Experiments with a new boosting algorithm., in: International Conference on Machine Learning
683 (ICML), 1996, pp. 148–156.
- 684 [51] C.-C. Chang, C.-J. Lin, LIBSVM: a library for support vector machines, software available at [http://www.csie.ntu.edu.
685 tw/~cjlin/libsvm](http://www.csie.ntu.edu.tw/~cjlin/libsvm) (2001).
- 686 [52] J. Fehr, K. Z. Arreola, H. Burkhardt, Fast support vector machine classification of very large datasets, Tech. rep., University
687 of Freiburg, Department of Computer (2007).
- 688 [53] I. W. Tsang, J. T. Kwok, P.-M. Cheung, Core Vector Machines: Fast SVM Training on Very Large Data Sets., Journal
689 of Machine Learning Research 6 (2005) 363–392.
- 690 [54] Y. Freund, R. E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, in:
691 Proceedings of the Second European Conference on Computational Learning Theory, 1995, pp. 23–37.
- 692 [55] R. E. Schapire, Y. Freund, P. Bartlett, W. S. Lee, Boosting the Margin: A New Explanation for the Effectiveness of
693 Voting Methods., The Annals of Statistics 26 (5) (1998) 1651–1686.
- 694 [56] P. Viola, M. Jones, Robust real-time face detection, International Journal of Computer Vision 57 (2) (2004) 137–154.

*Author Biography

Irene Rodriguez-Lujan received her degree in Computer Engineering and Mathematics from the Universidad Autonoma de Madrid (UAM) in 2007 and the Master degree in Computer Science in 2009 from the same university. Currently she is working on her Ph.D. thesis in the UAM Computer Science Department and she collaborates with the Instituto de Ingenieria del Conocimiento (IIC). Her research interests include feature selection and classification in real-time large-scale systems.

Carlos Santa Cruz received the physics degree from the Universidad Autonoma of Madrid (UAM) in 1987. He received the Ph.D. degree in physics from the UAM in 1991. Currently, he is a Professor at the Computer Science Department of the UAM and head of the Quantitative Methods Group of the Instituto de Ingenieria del Conocimiento (IIC). His research interests are in pattern recognition, model building, and time series forecasting using neural networks.

Ramon Huerta received the B.S. and M.S. degrees in physics from the Universidad Autonoma de Madrid, Spain, and the Ph.D. degree in Computer Science from the same University in 1994. He became Associate Professor at UAM in 2000 and left on leave of absence to his current Research Scientist position at the University of California San Diego. His research interests are at the intersection of Artificial Intelligence, Physics and Biology with more than 80 journal articles in these three fields.