

Linear subclass support vector machines

Nikolaos Gkalelis, Vasileios Mezaris, *Member, IEEE*, Ioannis Kompatsiaris, *Senior Member, IEEE*,
and Tania Stathaki

Abstract—In this letter, linear subclass support vector machines (LSSVMs) are proposed that can efficiently learn a piecewise linear decision function for binary classification problems. This is achieved using a nongaussianity criterion to derive the subclass structure of the data, and a new formulation of the optimization problem that exploits the subclass information. LSSVMs provide low computation cost during training and evaluation, and offer competitive recognition performance in comparison to other popular SVM-based algorithms. Experimental results on various datasets confirm the advantages of LSSVMs.

Index Terms—Support vector machines, subclasses, mixture of Gaussians, pattern recognition, classification, machine learning.

I. INTRODUCTION

Binary pattern classifiers, based on ensembles of linear support vector machines (LSVMs) [1], are recently receiving increasing attention due to some important advantages they offer over kernel SVMs (KSVMs) [2], [3]: a) faster training and testing times, and, b) competitive or superior classification performance. These methods employ a clustering algorithm to divide the feature space into several partitions and train a number of LSVMs in order to derive a piecewise linear decision function. The existing work in this area can be roughly divided into two major categories: a) Mixture of LSVM experts (MixLSVM) [2] utilize a gating network to provide a soft partition of the feature space, and for each partition an LSVM is used to learn the hyperplane that separates the positive from the negative samples of the partition. The gating network is then used to implicitly select the appropriate LSVM expert for classifying an unlabelled sample. b) Subclass-based approaches [3] divide each class to a number of subclasses (i.e., in contrary to MixLSVMs, each partition contains samples of only one class) and train one LSVM for each subclass, to derive a hyperplane separating the samples of a specific subclass from the samples in one or more other subclasses. During evaluation an appropriate framework, e.g. error correcting output codes (ECOC) [3], [4], is usually applied to combine the results of the binary classifiers for classifying a test sample.

Copyright (c) 2012 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

N. Gkalelis is with the Information Technologies Institute/Centre for Research and Technology Hellas (CERTH), Themi 57001, Greece, and also with the Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, U.K. (email: gkalelis@iti.gr).

V. Mezaris and I. Kompatsiaris are with the Information Technologies Institute/Centre for Research and Technology Hellas (CERTH), Themi 57001, Greece (email: bmezaris@iti.gr; ikom@iti.gr).

T. Stathaki is with the Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, U.K (email: t.stathaki@imperial.ac.uk).

To the best of our knowledge, the approaches presented until now in the literature break the multi-subclass problem into multiple independent binary subclass problems and train one LSVM for each problem. However, this yields binary classifiers that may not generalize well, since correlations between different subclasses are not adequately captured [5]. Moreover, most of the approaches combining binary LSVMs to solve multi-subclass binary problems do not provide an adequate theoretical analysis on their generalization properties. These issues can be addressed using multiclass SVMs (MSVMs) that naturally extend the concept of margin considering all data in a single constrained optimization problem (see for instance [5], [6] and the many references therein). However, a naive application of MSVMs directly in the derived subclasses will compute the parameters that enforce separability between every pair of subclasses, including subclasses of the same class. This will cause an unnecessary increase in the computational cost of the algorithm, and may additionally degrade the classification accuracy. To this end, inspired from similar approaches in discriminant analysis [7], we propose a new formulation of the optimization problem, that extends the MSVM formulation, so that only those hyperplanes that separate subclasses of different classes are computed. We refer to this new method as linear subclass SVMs (LSSVMs). For the efficient implementation of LSSVMs we exploit the sequential dual method (SDM) described in [8]. Moreover, we introduce a new nongaussianity measure for subclass partitioning, and exploit an ECOC framework for combining the binary subclass classifiers.

The rest of the paper is organized as follows: in section II the proposed approach is described, while in section III the performance of LSSVMs is evaluated using one artificial dataset and eleven publicly available datasets. Finally, conclusions are drawn in section IV.

II. LINEAR SUBCLASS SUPPORT VECTOR MACHINES

A. Problem formulation

Let $\mathcal{U} = \{(\bar{\mathbf{x}}_\kappa, (y_\kappa, u_\kappa)), \kappa = 1, \dots, N\}$ be a subclass partition of an annotated dataset consisting of N training samples, where $\bar{\mathbf{x}}_\kappa \in \mathbb{R}^F$ is the feature vector representation of the κ -th sample (F is the feature vector length), $y_\kappa \in [1, 2]$ and $u_\kappa \in [1, H_{y_\kappa}]$ are the class and subclass labels of $\bar{\mathbf{x}}_\kappa$, H_i is the number of subclasses of the i -th class, and $H = \sum_{i=1}^2 H_i$ is the total number of subclasses. MSVMs can effectively utilize the subclass information considering all data in one optimization problem [1], [5], [6]. For instance, the method of [5] will yield the following formulation

$$J_P(\mathbf{w}, \boldsymbol{\xi}) = \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^{H_i} \|\mathbf{w}_{i,j}\|^2 + C \sum_{\kappa=1}^N \xi_\kappa, \quad (1)$$

subject to the constraints

$$(\mathbf{w}_{y_\kappa, u_\kappa} - \mathbf{w}_{i,j})^T \mathbf{x}_\kappa \geq e_{\kappa, i, j} - \xi_\kappa, \quad \forall \kappa, i, j \quad (2)$$

where $\mathbf{w}_{i,j} = [\bar{\mathbf{w}}_{i,j}^T, b_{i,j}]^T$ contains the weight vector $\bar{\mathbf{w}}_{i,j} \in \mathbb{R}^F$ and bias $b_{i,j} \in \mathbb{R}$ referring to (i, j) subclass; $\mathbf{x}_\kappa = [\bar{\mathbf{x}}_\kappa^T, 1]^T$, $\xi_\kappa \geq 0$ is the slack variable corresponding to the κ -th sample, $C > 0$ is the penalty term, and $e_{\kappa, i, j} = 1 - \delta_{\kappa, i, j}$. In the latter, $\delta_{\kappa, i, j}$ is the subclass indicator function, i.e., $\delta_{\kappa, i, j} = 1$ if $(y_\kappa, u_\kappa) = (i, j)$, $\delta_{\kappa, i, j} = 0$ otherwise. The number of constraints in this formulation is $\sum_{\kappa=1}^N \sum_{i=1}^2 H_i = NH$. However, the use of constraints that enforce separability between subclasses of the same class increases the computational complexity without necessarily improving the classification performance. In LSSVM, we extend the MSVM formulation so that only the constraints that involve subclasses of different classes are accounted in the optimization problem. Thus, we seek to optimize (1) subject to the following set of constraints:

$$c_{\kappa, i, j} = (\mathbf{w}_{y_\kappa, u_\kappa} - \mathbf{w}_{i,j})^T \mathbf{x}_\kappa - e_{\kappa, i, j} + \xi_\kappa \geq 0, \quad (3)$$

$$\forall \kappa, (i \neq y_\kappa, j) \vee (i, j) = (y_\kappa, u_\kappa).$$

It can be seen that LSSVMs require $\sum_{i=1}^2 H_i - H_{y_\kappa} + 1$ constraints for the κ -th training sample, and the total number of constraints is $J = \sum_{\kappa=1}^N (\sum_{i=1}^2 H_i - H_{y_\kappa} + 1) = N_1 H_2 + N_2 H_1 + N$. Therefore, LSSVMs have $N_1 H_1 + N_2 H_2 - N$ fewer constraints compared to MSVMs. The number of constraints controls the number of Lagrange multipliers in the Lagrangian formulation, which in turn determines the complexity of the quadratic problem for identifying the parameters of the decision functions. This reveals a significant advantage of the LSSVM over naive application of MSVMs using subclasses: LSSVMs have much lower computational cost during training.

B. Computation of the LSSVM parameters

The primal Lagrangian of the LSSVMs can be formed using the equality constraints

$$\mathcal{L}_P(\mathbf{w}, \xi) = \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^{H_i} \|\mathbf{w}_{i,j}\|^2 + C \sum_{\kappa=1}^N \xi_\kappa$$

$$+ \sum_{\kappa=1}^N \sum_{i=1}^2 \sum_{j=1}^{H_i} \varrho_{\kappa, i, j} [(\mathbf{w}_{i,j} - \mathbf{w}_{y_\kappa, u_\kappa})^T \mathbf{x}_\kappa + e_{\kappa, i, j} - \xi_\kappa], \quad (4)$$

where, $\varrho_{\kappa, i, j}$, $\forall \kappa, i, j$, is the set of dual variables for the set of constraints in (3), and $e_{\kappa, y_\kappa, j} = \varrho_{\kappa, y_\kappa, j} = 0, \forall \kappa, j \neq u_\kappa$ are dummy variables introduced for notational convenience.

For the primal problem above the Karush-Kuhn-Tucker (KKT) conditions may be stated: $\frac{\partial \mathcal{L}_P}{\partial \xi_\kappa} = 0$, $\frac{\partial \mathcal{L}_P}{\partial \mathbf{w}_{m,n}} = 0$, $c_{\kappa, i, j} \geq 0$, $\varrho_{\kappa, i, j} \geq 0$, $\varrho_{\kappa, i, j} c_{\kappa, i, j} = 0$. Evaluating the equality constraints above we get

$$C = \sum_{i=1}^2 \sum_{j=1}^{H_i} \varrho_{\kappa, m, n}. \quad (5)$$

$$\mathbf{w}_{m,n} = \sum_{\kappa=1}^N (C \delta_{\kappa, m, n} - \varrho_{\kappa, m, n}) \mathbf{x}_\kappa, \quad (6)$$

Substituting (5), (6), back to (4), and defining $\alpha_{\kappa, i, j} = (C \delta_{\kappa, i, j} - \varrho_{\kappa, i, j})$ as the new dual variable we get the dual

formulation of the Lagrangian

$$\mathcal{L}_D(\alpha) = -\frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^{H_i} \|\mathbf{w}_{i,j}\|^2 - C \mathbf{e}^T \alpha, \quad (7)$$

where, $\mathbf{w}_{i,j} = \sum_{\kappa=1}^N \alpha_{\kappa, i, j} \mathbf{x}_\kappa$ is the weight vector expressed via the dual variables, $\mathbf{e} = [\mathbf{e}_1^T, \dots, \mathbf{e}_N^T]^T$, $\alpha = [\alpha_1^T, \dots, \alpha_N^T]^T$, $\mathbf{e}, \alpha \in \mathbb{R}^J$ are block vectors whose κ -th blocks are $\mathbf{e}_\kappa = [e_{\kappa, y_\kappa, u_\kappa}, e_{\kappa, \tilde{y}_\kappa, 1}, \dots, e_{\kappa, \tilde{y}_\kappa, H_{\tilde{y}_\kappa}}]^T$, $\alpha_\kappa = [\alpha_{\kappa, y_\kappa, u_\kappa}, \alpha_{\kappa, \tilde{y}_\kappa, 1}, \dots, \alpha_{\kappa, \tilde{y}_\kappa, H_{\tilde{y}_\kappa}}]^T$, $\mathbf{e}_\kappa, \alpha_\kappa \in \mathbb{R}^{J_{y_\kappa}}$, $J_{y_\kappa} = 1 + H_{\tilde{y}_\kappa}$, and, \tilde{y}_κ is 1 for $y_\kappa = 2$ and 2 for $y_\kappa = 1$. That is, the κ -th block contains the variables associated with the κ -th sample excluding the dummy variables $e_{\kappa, y_\kappa, j} = \alpha_{\kappa, y_\kappa, j} = 0, \forall j \neq u_\kappa$. Consequently, exploiting the KKT conditions the dual optimization problem can now be stated as

$$\min_{\alpha} -\mathcal{L}_D(\alpha) \quad \text{subject to} \quad \begin{cases} \alpha \leq \mathbf{C}, \\ \mathbf{1}^T \alpha = 0, \end{cases} \quad (8)$$

where, $\mathbf{1} \in \mathbb{R}^J$ is a vector of ones and $\mathbf{C} = [\mathbf{C}_1^T, \dots, \mathbf{C}_N^T]^T$ is a block vector whose κ -th block $\mathbf{C}_\kappa = [C_{\kappa, y_\kappa, u_\kappa}, C_{\kappa, \tilde{y}_\kappa, 1}, \dots, C_{\kappa, \tilde{y}_\kappa, H_{\tilde{y}_\kappa}}]^T$ contains the penalty terms associated with the κ -th sample given by $C_{\kappa, i, j} = \delta_{\kappa, i, j} C$. This optimization problem is quadratic in terms of α with linear constraints and therefore it can be solved using an appropriate technique. Here we extend the sequential dual method (SDM) presented in [8] to derive an efficient sequential derivation of the LSSVM dual variables. This algorithm uses the gradient information to optimize the dual variables and the weight vector in a sequential manner. At each iteration an additive update $\delta \alpha_\kappa$ is computed for updating the dual variables of the block vector α_κ and the weight vectors $\mathbf{w}_{i,j}$, solving the following reduced optimization problem

$$\min_{\delta \alpha_\kappa} \frac{1}{2} A_\kappa \|\delta \alpha_\kappa\|^2 + \mathbf{g}_\kappa^T \delta \alpha_\kappa \quad \text{subject to} \quad \begin{cases} \delta \alpha_\kappa \leq \mathbf{C}_\kappa \\ \mathbf{1}_\kappa^T \delta \alpha_\kappa = 0, \end{cases} \quad (9)$$

where, $A_\kappa = \|\mathbf{x}_\kappa\|^2$, $\mathbf{1}_\kappa$ is a vector of the same length with α_κ and with all elements equal to 1, $\mathbf{g}_\kappa = [g_{\kappa, y_\kappa, u_\kappa}, g_{\kappa, \tilde{y}_\kappa, 1}, \dots, g_{\kappa, \tilde{y}_\kappa, H_{\tilde{y}_\kappa}}]^T$, and $g_{\kappa, i, j}$ is the gradient of (7) with respect to the dual variable $\alpha_{\kappa, i, j}$ given by

$$g_{\kappa, i, j} = \mathbf{w}_{i,j}^T \mathbf{x}_\kappa + e_{\kappa, i, j}. \quad (10)$$

Optimality of α_κ is checked using the quantity $v_\kappa = \hat{g}_\kappa - \check{g}_\kappa$, where $\hat{g}_\kappa = \operatorname{argmax}_{i,j} g_{\kappa, i, j}$ and $\check{g}_\kappa = \operatorname{argmin}_{i,j: \alpha_{\kappa, i, j} < C_{\kappa, i, j}} g_{\kappa, i, j}$. That is, we consider that the dual variables have converged to their optimal values when $v_\kappa < \epsilon, \forall \kappa$, where ϵ is a positive tolerance parameter. The overall procedure is described in Algorithm 1.

C. Classification

In MSVM formulations a test sample \mathbf{x}_t is classified to one of the subclasses usually according to the rule $\operatorname{argmax}_{i,j} \mathbf{w}_{i,j}^T \mathbf{x}_t$. This procedure is similar to the one-versus-all classifier in multiclass problems. In contrast to this, we use the results of the LSSVM to construct the separating hyperplanes between the different subclasses, similarly to [9]. The advantage of the latter approach is that an appropriate framework can be exploited to combine the derived binary

Algorithm 1 LSSVM training

Input: Annotated data set \mathbf{X} , penalty C
Output: α , $\mathbf{w}_{i,j}$, $\forall i, j$

- 1: Initialize: $\alpha = [0, \dots, 0]^T$, $\mathbf{w}_{i,j} = [0, \dots, 0]^T$, $\forall i, j$
- 2: **repeat**
- 3: **for** $\kappa = 1$ to N **do**
- 4: Compute \mathbf{g}_κ , v_κ (10)
- 5: **if** $v_\kappa \geq \epsilon$ **then**
- 6: Compute $\delta\alpha_\kappa$ (9); $\alpha_\kappa \leftarrow \alpha_\kappa + \delta\alpha_\kappa$
- 7: $\mathbf{w}_{i,j} \leftarrow \mathbf{w}_{i,j} + \delta\alpha_{\kappa,i,j}\mathbf{x}_\kappa, \forall (i \neq y_\kappa, j) \vee (i, j) = (y_\kappa, u_\kappa)$
- 8: **end if**
- 9: **end for**
- 10: **until** $v_\kappa < \epsilon$, $\forall \kappa$

classifiers, as explained in the following. The hyperplane $\eta_{p,q}$ that separates subclass $\mathcal{X}_{1,p}$ and $\mathcal{X}_{2,q}$ is defined as $\eta_{p,q} = \mathbf{w}_{1,p} - \mathbf{w}_{2,q}$, ($p \in [1, H_1], q \in [1, H_2]$). Hence, the respective binary subclass classifier

$$g_{p,q}(\mathbf{x}_t) = \text{sign}[(\mathbf{w}_{1,p} - \mathbf{w}_{2,q})^T \mathbf{x}_t] = \text{sign}(\eta_{p,q}^T \mathbf{x}_t) \quad (11)$$

assigns the test sample to $\mathcal{X}_{1,p}$ if $g_{p,q}(\mathbf{x}_t) > 0$ and to $\mathcal{X}_{2,q}$ otherwise. The binary subclass classifiers can then be combined to a binary classifier using a subclass ternary error correcting output code (ECOC) framework [3]. In our case, the ternary coding matrix $\mathbf{M} \in \{1, 0, -1\}^{H \times H_1 H_2}$ is defined so that the rows represent the subclasses, $(1, 1), \dots, (1, H_1), (2, 1), \dots, (2, H_2)$, and the columns represent the LSSVM decision functions, $g_{1,1}, \dots, g_{H_1, H_2}$. That is, each column corresponds to a distinct subclass pair $\{(1, p), (2, q)\}$, and for such a column, \mathbf{M} has +1 in row corresponding to subclass $(1, p)$, -1 to the row corresponding to subclass $(2, q)$, and zeros in all other rows. This simulates a one-versus-one subclass coding design, where the subclasses belong to different classes. For classifying an unlabelled sample \mathbf{x}_t , we use the linear loss-weighted (LLW) decoding strategy on the output of the binary subclass classifiers [4]

$$LLW((i, j), \mathbf{x}_t) = \sum_{p=1}^{H_1} \sum_{q=1}^{H_2} -M_{i,j}^{p,q} g_{p,q}(\mathbf{x}_t) \tilde{M}_{i,j}^{p,q}, \quad (12)$$

where $M_{i,j}^{p,q}, \tilde{M}_{i,j}^{p,q}$ are the elements of the coding and weighting matrix [4], respectively, that correspond to the (i, j) subclass and the hyperplane that separates the subclasses $(1, p)$ and $(2, q)$. Then, the test sample is classified according to the rule $\text{argmin}_{i,j} LLW((i, j), \mathbf{x}_t)$.

D. Subclass partitioning

In the above, we exploit a partitioning of the data to subclasses. Assuming that the data have a Gaussian subclass structure, we use an iterative algorithm to derive this subclass division. Starting from the initial class partition (i.e. $H = 2$), at each iteration we increase the number of subclasses of the k -th class according to the following rule $k = \text{argmax}_{i=1,2}(\Phi_i)$, where Φ_i measures the nongaussianity of the i -th class. Exploiting the work in [10] we define this measure as

$$\Phi_i = \sum_{j=1}^{H_i} \left(\frac{N_{i,j} \hat{\beta}_{i,j}}{6} + \left| \frac{\hat{\gamma}_{i,j} - F(F+2)}{\sqrt{8F(F+2)/N_{i,j}}} \right| \right), \quad (13)$$

where F is the dimensionality of the feature vectors. The quantities $\hat{\beta}_{i,j}, \hat{\gamma}_{i,j}$ above are estimates of the multivariate

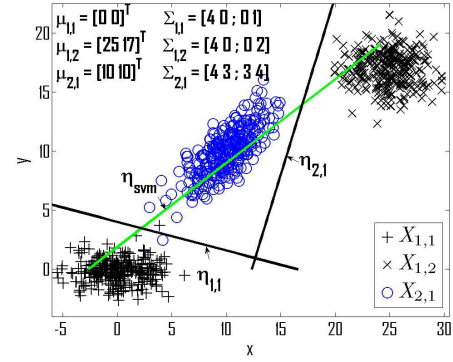


Fig. 1. Artificial dataset with a linearly separable subclass structure.

skewness and kurtosis of the (i, j) subclass respectively, computed as follows

$$\hat{\beta}_{i,j} = \frac{1}{N^2} \sum_{\kappa=1}^{N_{i,j}} \sum_{\nu=1}^{N_{i,j}} [(\bar{\mathbf{x}}_\kappa - \hat{\boldsymbol{\mu}}_{i,j})^T \hat{\Sigma}_{i,j}^{-1} (\bar{\mathbf{x}}_\nu - \hat{\boldsymbol{\mu}}_{i,j})]^3 \quad (14)$$

$$\hat{\gamma}_{i,j} = \frac{1}{N} \sum_{\kappa=1}^{N_{i,j}} [(\bar{\mathbf{x}}_\kappa - \hat{\boldsymbol{\mu}}_{i,j})^T \hat{\Sigma}_{i,j}^{-1} (\bar{\mathbf{x}}_\kappa - \hat{\boldsymbol{\mu}}_{i,j})]^2 \quad (15)$$

where $N_{i,j}, \hat{\boldsymbol{\mu}}_{i,j} = \frac{1}{N_{i,j}} \sum_{\kappa=1}^{N_{i,j}} \bar{\mathbf{x}}_\kappa$, $\hat{\Sigma}_{i,j} = \frac{1}{N_{i,j}} \sum_{\kappa=1}^{N_{i,j}} (\bar{\mathbf{x}}_\kappa - \hat{\boldsymbol{\mu}}_{i,j})(\bar{\mathbf{x}}_\kappa - \hat{\boldsymbol{\mu}}_{i,j})^T$ are the number of samples, the sample mean and the sample covariance matrix of (i, j) subclass, and the summands above are over the samples $\bar{\mathbf{x}}_\kappa, \bar{\mathbf{x}}_\nu$ that belong to the (i, j) subclass, i.e. the labels of the samples are $(y_\nu, u_\nu) = (y_\kappa, u_\kappa) = (i, j)$. It has been shown in [10] that for large $N_{i,j}$, the quantities $N_{i,j} \hat{\beta}_{i,j} / 6$ and $(\hat{\gamma}_{i,j} - F(F+2)) / \sqrt{8F(F+2)/N_{i,j}}$ follow a chi-square distribution with $F(F+1)(F+2)/6$ degrees of freedom and a standard normal distribution $\mathcal{N}(0, 1)$ respectively. Consequently, Φ_i defined in (13) expresses the departure of the i -th class distribution from a multivariate Gaussian mixture.

Using the iterative procedure described above, the best partition is selected using either cross-validation, i.e., we repeat the splitting procedure until a maximum number of subclasses is reached and then select the partition that provides the best empirical recognition rate, or by inspecting the convergence of a total nongaussianity measure $\Phi = \Phi_1 + \Phi_2$ to a steady-state solution. We used the former approach in our experiments.

III. EXPERIMENTAL RESULTS

A. Artificial dataset

A classification example with two non-linearly separable classes $\mathcal{X}_1, \mathcal{X}_2$ is depicted in Figure 1. \mathcal{X}_1 consists of two Gaussian subclasses $\mathcal{X}_{1,1}, \mathcal{X}_{1,2}$, whereas \mathcal{X}_2 is a single Gaussian $\mathcal{X}_{2,1}$. In contrary to the main classes, this dataset has a linearly separable subclass structure. From each subclass distribution 300 samples were randomly drawn to form the artificial dataset. The subclass partitioning (II-D) and the LSSVM algorithm were then applied, correctly identifying the subclass structure of the data and the separating hyperplanes $\eta_{1,1}, \eta_{2,1}$. For comparison, the hyperplane η_{svm} computed by LSVM does not provide a valid solution for separating the two classes. At the same time, LSSVM avoids the unnecessary

computations for estimating the hyperplane that separates subclasses $\mathcal{X}_{1,1}$ and $\mathcal{X}_{1,2}$, which MSVM would perform.

B. UCI and Delve datasets

We compare the performance of LSSVM with LSVM, KSVM with radial basis function (RBF) kernel [1], MLSVM [5], MixLSVM [2] and subclass ECOC (SECOC) [3], using 11 datasets from the UCI [11] and Delve [12] repositories: WDBC, Ionosphere, Breast Cancer, Pima, German, Heart, Splice, Titanic, and Monk-1,-2,-3. For the evaluation of the algorithms, a division of the datasets to training and test sets is necessary. This already exists for the Monk sets. For Ionosphere, similarly to other works, we use the first 200 samples for training and the rest 151 as testing instances, while, for WDBC we created a random 50–50 train–test split. For the Breast Cancer, Pima, German, Heart, Splice and Titanic datasets we used 10 random partitions for each dataset from the Gunnar Rätsch’s benchmark collection [13]. The recognition performance of a method regarding a dataset is measured using the average correct recognition rate (CCR) along all the partitions. Moreover, for assessing the statistical significance of the difference in the performance between any two algorithms, we used the McNemar’s hypothesis test [14].

For the experiments, LSSVM and MixLSVM were implemented in Matlab, while for LSVM, KSVM, MSVM and the base classifiers of SECOC, the libsvm [15] or the liblinear [16] packages were used. In SECOC, for fair comparison with LSSVM, we used LSVMs as base classifiers and the LLW strategy for classifying a test sample. In the stage of model selection, for LSVM we need to estimate the penalty term C , while for KSVM and LSSVM we additionally optimize over the scale parameter σ of the RBF function [1] and the number of subclasses H , respectively. Three parameters are optimized for MixLSVM, i.e. C , H and the scale parameter τ of the gating function [2]. Finally, in SECOC we search for the optimal values of the three splitting criteria (θ_{size} , θ_{perf} and θ_{impr}) [3] and the penalty term C of the base classifiers. The optimal parameters for each method are selected through a grid search strategy where the primary metric is the CCR.

The evaluation results are presented in Table I. We can see that LSSVM provides top recognition rate in 8 out of the 11 datasets. Using the McNemar’s test with significance level 0.025 we verified that the differences in performance between LSSVM and MSVM are statistically significant in 8 of the datasets; and similarly in 7 of the datasets for the differences between LSSVM and MixLSVM, and between LSSVM and SECOC. We experimentally verified that the training of LSSVM is much faster than MixLSVM (on average, 25 times faster) mainly due to the slow convergence of the expectation maximization-based MixLSVM, and because LSSVM requires the optimization of fewer training parameters; similarly, we expect the training stage of LSSVM to be faster than SECOC, due to the lower number of parameters that need to be optimized. The LSSVM training and testing are also less computationally demanding in comparison to MSVM and KSVM, due to the lower number of constraints and the absence of a kernel evaluation step, respectively.

	LSVM	LSSVM	MSVM	MixLSVM	SECOC	KSVM
Monk-1	70.6%	98.8%	94.7%	86.3%	81.7%	86.1%
Monk-2	67.1%	75.5%	74.5%	76.9%	67.3%	81.9%
Monk-3	81%	96.5%	86.6%	96.5%	93%	97%
WDBC	95.8%	96.1%	95.4%	96.1%	96.1%	94.4%
Ionosph.	97.4%	98.7%	95.4%	98%	97.4%	98%
B. Cancer	69.7%	77.9%	72.6%	74.6%	74%	74.7%
Pima	76.9%	79%	75.6%	77.7%	77.9%	77%
German	76.9%	79.7%	74%	76.9%	80.3%	77.9%
Heart	81.7%	85.6%	82.7%	83.7%	83.2%	84.4%
Splice	84.1%	89.6%	88.6%	87.5%	88%	88.4%
Titanic	76.6%	78.6%	76.8%	77.8%	77.6%	78%

TABLE I
CORRECT RECOGNITION RATES.

IV. CONCLUSIONS

In this letter, LSSVMs were proposed that exploit the subclass structure of the data to efficiently compute a separating piecewise linear hyperplane. Experiments on various datasets demonstrated the effectiveness of LSSVMs.

ACKNOWLEDGMENT

This work was supported by the EC under contracts FP7-248984 GLOCAL and FP7-287911 LinkedTV.

REFERENCES

- [1] V. Vapnik, *Statistical learning theory*. New York: Wiley, 1998.
- [2] Z. Fu, A. Robles-Kelly, and J. Zhou, “Mixing linear SVMs for nonlinear classification,” *IEEE Trans. Neural Netw.*, vol. 21, no. 12, pp. 1963–1975, Dec. 2010.
- [3] S. Escalera, D. M. Tax, O. Pujol, P. Radeva, and R. P. Duin, “Subclass problem-dependent design for error-correcting output codes,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 6, pp. 1041–1054, Jun. 2008.
- [4] S. Escalera, O. Pujol, and P. Radeva, “On the decoding process in ternary error-correcting output codes,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 1, pp. 120–134, Jan. 2010.
- [5] K. Crammer and Y. Singer, “On the algorithmic implementation of multiclass kernel-based vector machines,” *J. Mach. Learn. Res.*, vol. 2, pp. 265–292, Mar. 2002.
- [6] Y. Guermeur, A. Elisseeff, and D. Zelus, “A comparative study of multiclass support vector machines in the unifying framework of large margin classifiers,” *Appl. Stoch. Model. Bus. Ind.*, vol. 21, no. 2, pp. 199–214, Mar. 2005.
- [7] N. Gkalelis, V. Mezaris, and I. Kompatsiaris, “Mixture subclass discriminant analysis,” *IEEE Signal Process. Lett.*, vol. 18, no. 5, pp. 319–322, May 2011.
- [8] S. S. Keerthi, S. Sundararajan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin, “A sequential dual method for large scale multi-class linear SVMs,” in *Proc. 14th ACM SIGKDD*, Las Vegas, USA, Aug. 2008, pp. 408–416.
- [9] E. J. Breidensteiner and K. P. Bennett, “Multicategory classification by support vector machines,” *Comput. Optim. Appl.*, vol. 12, no. 1–3, pp. 53–79, Jan. 1999.
- [10] K. V. Mardia, “Measures of multivariate skewness and kurtosis with applications,” *Biometrika*, vol. 57, no. 3, pp. 519–530, Dec. 1970.
- [11] A. Frank and A. Asuncion, “UCI machine learning repository,” University of California, Irvine, School of Information and Computer Sciences, 2010. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [12] “Data for evaluating learning in valid experiments,” <http://www.cs.toronto.edu/~delve/index.html>, accessed 2012-06-10.
- [13] G. Rätsch, T. Onoda, and K.-R. Müller, “Soft margins for adaboost,” *Mach. Learn.*, vol. 42, no. 3, pp. 287–320, Mar. 2001.
- [14] Q. McNemar, “Note on the sampling error of the difference between correlated proportions or percentages,” *Psychometrika*, vol. 12, no. 2, pp. 153–157, Jun. 1947.
- [15] C.-C. Chang and C.-J. Lin, “LIBSVM: A library for support vector machines,” *ACM Trans. Intell. Syst. Technol.*, vol. 2, pp. 27:1–27:27, 2011.
- [16] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, “LIBLINEAR: A library for large linear classification,” *Journal of Machine Learning Research*, vol. 9, pp. 1871–1874, 2008.