

Simple Incremental One-Class Support Vector Classification

Kai Labusch, Fabian Timm, and Thomas Martinetz

Institute for Neuro- and Bioinformatics, University of Lübeck,
Ratzeburger Allee 160, D-23538 Lübeck, Germany

Abstract. We introduce the OneClassMaxMinOver (OMMO) algorithm for the problem of one-class support vector classification. The algorithm is extremely simple and therefore a convenient choice for practitioners. We prove that in the hard-margin case the algorithm converges with $\mathcal{O}(1/\sqrt{t})$ to the maximum margin solution of the support vector approach for one-class classification introduced by Schölkopf et al. Furthermore, we propose a 2-norm soft margin generalisation of the algorithm and apply the algorithm to artificial datasets and to the real world problem of face detection in images. We obtain the same performance as sophisticated SVM software such as libSVM.

1 Introduction

Over the last years, the support vector machine [1] has become a standard approach in solving pattern recognition tasks. There are several training techniques available, e.g. SMO [2]. Although these methods seem to be simple, they are hard to understand without a background in optimisation theory. Hence, they are difficult to motivate when explained to practitioners. In many cases, in particular in industrial contexts, where external libraries or other third party software cannot be used due to various reasons, these techniques are not applied, even though they might be beneficial for solving the problem.

In many applications one has to cope with the problem that only samples of one class are given. The task is to separate this class from the *other* class that consists of all outliers. Either only few samples of the outlier class are given or the outlier class is missing completely. In these cases two-class classifiers often show bad generalisation performance and it is advantageous to employ one-class classification.

Approaches to one-class classification can be divided into three groups: density estimators, reconstruction methods, and boundary methods. The first and the second group are the most powerful because they derive a model of the data that is defined everywhere in the input space. An advantage of the boundary methods is that they consider an easier problem, that is, describing only the class boundaries, instead of describing the whole distribution of the data.

In the present work, we describe a very simple and incremental boundary method based on the support vector approach. It provides the same solution

as comparable techniques such as SMO, despite being extremely simple and therefore applicable for practitioners who are not within the field of machine learning.

2 Previous Work

In the context of one-class classification several boundary methods have been developed. We only want to give a brief description of two approaches that have been introduced almost simultaneously. Tax et al [3] consider the problem of finding the smallest enclosing ball of given data samples $\mathbf{x}_i \in \mathcal{X}, i = 1, \dots, L$ that is described by the radius R and centre \mathbf{w} :

$$\min_{\mathbf{w}, R} \left(R + \frac{1}{\nu l} \sum_i \xi_i \right) \quad \text{s.t. } \forall i : \|\phi(\mathbf{x}_i) - \mathbf{w}\| \leq R - \xi_i \wedge \xi_i \geq 0 . \quad (1)$$

This is the soft version of the problem. It deals with outliers by using slack variables ξ_i in order to allow for samples that are not located inside the ball defined by \mathbf{w} and R . For $\nu \rightarrow 0$ one obtains the hard-margin solution, where all samples are located inside the ball. Here ϕ denotes a mapping of the data samples to some feature space.

Schölkopf et. al [4] show that one-class classification can be taken as two-class classification, where the *other* class is represented by the origin. They consider the problem of finding the hyperplane \mathbf{w} that separates the data samples from the origin with maximum distance ρ :

$$\min_{\mathbf{w}, \xi, \rho} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu l} \sum_i \xi_i - \rho \right) \quad \text{s.t. } \forall i : \mathbf{w}^T \phi(\mathbf{x}_i) \geq \rho - \xi_i \wedge \xi_i \geq 0 . \quad (2)$$

Again, the soft-margin problem is shown. It allows for samples that are misclassified by using slack variables ξ_i . For $\nu \rightarrow 0$ one obtains the hard-margin solution that enforces correct classification of all given samples.

In [4] it is shown that (1) and (2) turn out to be equivalent, if the $\phi(\mathbf{x}_i)$ lie on the surface of a sphere. Then, the radius R of problem (1) and the margin ρ in (2) can easily be computed by choosing a support vector on the boundary. If the Gaussian kernel is used

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \exp(-\|\mathbf{x} - \mathbf{y}\|/2\sigma^2) \quad (3)$$

in order to implicitly map the given samples to some feature space, the $\phi(\mathbf{x}_i)$ have unit norm and this condition is satisfied. To make the problem solvable, the origin has to be linearly separable from the target class. This precondition is also given if a Gaussian kernel is used. In the following we require that the data has been mapped to some feature space where these conditions hold, i.e. linear separability of the origin and unit norm of all samples.

3 OneClassMaxMinOver

In this section we describe a very simple incremental algorithm for one-class-classification called OneClassMaxMinOver(OMMO). This algorithm is closely connected to problem (2). It is inspired by the MaxMinOver algorithm for two-class classification proposed in [5]. We consider the problem of finding the hyperplane \mathbf{w}_* passing the origin and having maximum margin ρ_* with respect to the given data samples. Finding this hyperplane is equivalent to solving the optimisation problem (2), that is finding the hyperplane \mathbf{w}_* that separates the given data samples with maximum margin ρ_* from the origin (see Fig. 1).

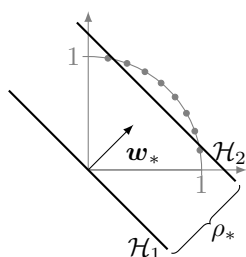


Fig. 1. Some data samples having unit norm are shown as well as the solutions of the optimisation problem (4), i.e. \mathcal{H}_1 and the solution of (2), i.e. \mathcal{H}_2 .

Mathematically, we are looking for the solution of the following optimisation problem

$$\mathbf{w}_* = \arg \max_{\mathbf{w}} \left(\min_{\mathbf{x}_i} (\mathbf{w}^T \mathbf{x}_i) \right) \quad \text{s.t. } \|\mathbf{w}\| = 1 . \quad (4)$$

The margin ρ_* is obtained by

$$\rho_* = \min_{\mathbf{x}_i} (\mathbf{w}_*^T \mathbf{x}_i) . \quad (5)$$

In the following \mathbf{w}_t denotes the approximation of \mathbf{w}_* at time t . During the learning process the constraint $\|\mathbf{w}\| = 1$ is dropped. The algorithm starts with $\mathbf{w}_0 = \mathbf{0}$ and after t_{\max} learning iterations the norm of the final approximation $\mathbf{w}_{t_{\max}}$ is set to one. In each learning iteration, the algorithm selects the sample that is closest to the current hyperplane defined by \mathbf{w}_t

$$\mathbf{x}_{\min}(t) = \arg \min_{\mathbf{x}_i} (\mathbf{w}_t^T \mathbf{x}_i) . \quad (6)$$

For each given training sample \mathbf{x}_i there is a counter variable α_i that is increased by 2 whenever the sample is selected as $\mathbf{x}_{\min}(t)$:

$$\alpha_i = \alpha_i + 2 \quad \text{for } \mathbf{x}_{\min}(t) = \mathbf{x}_i . \quad (7)$$

$\mathcal{X}'(t)$ denotes the set of samples \mathbf{x}_j for which $\alpha_j > 0$ holds at time t . Out of this set, the algorithm selects the sample being most distant with respect to the current hyperplane defined by \mathbf{w}_t :

$$\mathbf{x}_{\max}(t) = \arg \max_{\mathbf{x}_j \in \mathcal{X}'(t)} (\mathbf{w}_t^T \mathbf{x}_j) . \quad (8)$$

Whenever a sample is selected as $\mathbf{x}_{\max}(t)$, its associated counter variable is decreased by 1:

$$\alpha_i = \alpha_i - 1 \quad \text{for } \mathbf{x}_{\max}(t) = \mathbf{x}_i . \quad (9)$$

The approximation of \mathbf{w}_* in learning iteration $t + 1$ is given by

$$\mathbf{w}_{t+1} = \sum_{i=1}^L \alpha_i \mathbf{x}_i . \quad (10)$$

Note that (7) and (9) can be combined to the learning rule

$$\mathbf{w}_{t+1} = \mathbf{w}_t + 2\mathbf{x}_{\min}(t) - \mathbf{x}_{\max}(t) . \quad (11)$$

Altogether, we obtain algorithm 1 for incremental one-class classification.

Algorithm 1: OneClassMaxMinOver. With $h(\mathbf{x}_i) = \sum_{j=1}^L \alpha_j \mathbf{x}_j^T \mathbf{x}_i$

```

 $\alpha_i \leftarrow 0 \quad \forall i = 1, \dots, N$ 
for  $t = 0$  to  $t_{\max}$  do
   $\mathbf{x}_{\min}(t) \leftarrow \arg \min_{\mathbf{x}_i \in \mathcal{X}} h(\mathbf{x}_i)$ 
   $\mathbf{x}_{\max}(t) \leftarrow \arg \max_{\mathbf{x}_i \in \mathcal{X}'(t)} h(\mathbf{x}_i)$ 
   $\alpha_{\min} \leftarrow \alpha_{\min} + 2$ 
   $\alpha_{\max} \leftarrow \alpha_{\max} - 1$ 
end
 $\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} / \sqrt{\sum_i \alpha_i h(\mathbf{x}_i)}$ 
 $\rho \leftarrow \min_{\mathbf{x}_i \in \mathcal{X}} h(\mathbf{x}_i)$ 

```

In section 3.1 we are going to prove that for $t \rightarrow \infty$ the following propositions hold:

- $\mathbf{w}_t / \|\mathbf{w}_t\|$ converges at least as $\mathcal{O}(1/\sqrt{t})$ to \mathbf{w}_* .
- $\alpha_i > 0$ only holds for samples \mathbf{x}_i having distance ρ_* with respect to the hyperplane \mathbf{w}_* , i.e. support vectors.

As mentioned before, we require that the data set has been mapped into some feature space where all samples have unit norm and can be linearly separated from the origin. However, this has not to be done explicitly. It can also be achieved by replacing the standard scalar product with a kernel that implements an implicit mapping to a feature space having the required properties. In case of the OMMO algorithm this corresponds to replacing the function $h(\mathbf{x}_i)$ with

$$h(\mathbf{x}_i) = \sum_{j=1}^L \alpha_j K(\mathbf{x}_j, \mathbf{x}_i) , \quad (12)$$

where $K(\mathbf{x}_j, \mathbf{x}_i)$ is an appropriate kernel function. For a discussion of kernel functions see, for example, [6].

3.1 Proof of Convergence

Our proof of convergence for the OneClassMaxMinOver algorithm is based on the proof of convergence for MaxMinOver by Martinetz [5], who showed a convergence speed of $\mathcal{O}(1/\sqrt{t})$ in the case of two-class classification.

Proposition 1. *The length of \mathbf{w}_t is bounded such that $\|\mathbf{w}_t\| \leq \rho_* t + 3\sqrt{t}$.*

Proof. This is done by induction and using the properties that

$$\rho_* \geq \rho_t = \frac{\mathbf{w}_t^T \mathbf{x}_{\min}(t)}{\|\mathbf{w}_t\|} \Rightarrow \rho_* \|\mathbf{w}_t\| \geq \mathbf{w}_t^T \mathbf{x}_{\min}(t) , \quad (13)$$

$$\forall i : \|\mathbf{x}_i\| = 1 , \text{ and} \quad (14)$$

$$\forall t : \mathbf{x}_{\min}(t)^T \mathbf{x}_{\max}(t) = \cos \beta \underbrace{\|\mathbf{x}_{\min}(t)\|}_{=1} \underbrace{\|\mathbf{x}_{\max}(t)\|}_{=1} \geq -1 . \quad (15)$$

The case $t = 0$ is trivial and for $t \rightarrow t + 1$ it follows that

$$\begin{aligned} \|\mathbf{w}_{t+1}\|^2 &\stackrel{(11)}{=} \mathbf{w}_t^T \mathbf{w}_t + 2\mathbf{w}_t^T (2\mathbf{x}_{\min}(t) - \mathbf{x}_{\max}(t)) + (2\mathbf{x}_{\min}(t) - \mathbf{x}_{\max}(t))^2 \\ &= \mathbf{w}_t^T \mathbf{w}_t + 2\mathbf{w}_t^T \mathbf{x}_{\min}(t) + 2 \underbrace{(\mathbf{w}_t^T \mathbf{x}_{\min}(t) - \mathbf{w}_t^T \mathbf{x}_{\max}(t))}_{\leq 0} + \\ &\quad 4\mathbf{x}_{\min}(t)^T \mathbf{x}_{\min}(t) + \mathbf{x}_{\max}(t)^T \mathbf{x}_{\max}(t) - 4\mathbf{x}_{\min}(t)^T \mathbf{x}_{\max}(t) \\ &\stackrel{(14),(15)}{\leq} \mathbf{w}_t^T \mathbf{w}_t + 2\mathbf{w}_t^T \mathbf{x}_{\min}(t) + 9 \\ &\stackrel{(13)}{\leq} \mathbf{w}_t^T \mathbf{w}_t + 2\rho_* \|\mathbf{w}_t\| + 9 \\ &\leq (\rho_* t + 3\sqrt{t})^2 + 2\rho_* (\rho_* t + 3\sqrt{t}) + 9 \\ &= \rho_*^2 t^2 + 2\rho_*^2 t + 6\rho_* t\sqrt{t} + 6\rho_* \sqrt{t} + 9t + 9 \\ &\leq \rho_*^2 (t^2 + 2t) + 6\rho_* (t+1)\sqrt{t} + 9(t+1) + \rho_*^2 \\ &\leq \rho_*^2 (t+1)^2 + 6\rho_* (t+1)\sqrt{t+1} + 9(t+1) \\ &= (\rho_* (t+1) + 3\sqrt{t+1})^2 . \end{aligned}$$

□

Theorem 1. *For $t \rightarrow \infty$ the angle γ_t between the optimal direction \mathbf{w}_* and the direction \mathbf{w}_t found by OMMO converges to zero, i.e. $\lim_{t \rightarrow \infty} \gamma_t = 0$.*

Proof.

$$\begin{aligned} \cos \gamma_t &= \frac{\mathbf{w}_*^T \mathbf{w}_t}{\|\mathbf{w}_t\|} \\ &= \frac{1}{\|\mathbf{w}_t\|} \sum_{i=0}^{t-1} \mathbf{w}_*^T (2\mathbf{x}_{\min}(i) - \mathbf{x}_{\max}(i)) \end{aligned} \quad (16)$$

$$\begin{aligned}
&= \frac{1}{\|\mathbf{w}_t\|} \sum_{i=0}^{t-1} \underbrace{\mathbf{w}_*^T \mathbf{x}_{\min}(i)}_{\geq \rho_*} \geq \frac{1}{\|\mathbf{w}_t\|} \rho_* t \quad (17) \\
&\stackrel{\text{Prop.1}}{\geq} \frac{\rho_* t}{\rho_* t + 3\sqrt{t}} = \frac{1}{1 + \frac{3\sqrt{t}}{\rho_* t}} \\
&\geq 1 - \frac{3}{\rho_* \sqrt{t}} \xrightarrow{t \rightarrow \infty} 1
\end{aligned}$$

From (16) to (17) we have used that a sample can only be *forgotten*, if it was *learnt* before, such that $\forall \mathbf{x}_{\max}(t) \exists \mathbf{x}_{\min}(t'), t' < t : \mathbf{x}_{\max}(t) = \mathbf{x}_{\min}(t')$. \square

Theorem 2. *Beyond some finite number of iterations \hat{t} the set $\mathcal{X}'(t)$ will always consist only of support vectors.*

Proof. First, we show that after some finite number of iterations t' the $\mathbf{x}_{\min}(t)$ with $t > t'$ will always be a support vector. We use an orthogonal decomposition of \mathbf{w}_t as shown in Fig. 2. \mathcal{X}^{sv} will denote the set of true support vectors. For an indirect proof we assume that a finite number of iterations t' where $\mathbf{x}_{\min}(t)$ with $t > t'$ will always be a support vector does not exist, i.e. $\nexists t' < \infty \forall t, t' < t : \mathbf{x}_{\min}(t) \in \mathcal{X}^{\text{sv}}$.

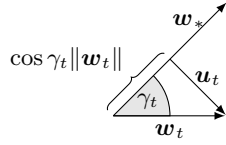
$$\begin{aligned}
\Rightarrow \rho_* \geq \rho_t &= \frac{\mathbf{w}_t^T \mathbf{x}_{\min}(t)}{\|\mathbf{w}_t\|} \stackrel{(20)}{=} \frac{(\cos \gamma_t \|\mathbf{w}_t\| \mathbf{w}_* + \mathbf{u}_t)^T \mathbf{x}_{\min}(t)}{\|\mathbf{w}_t\|} \\
&= \cos \gamma_t \mathbf{w}_*^T \mathbf{x}_{\min}(t) + \frac{\mathbf{u}_t^T \mathbf{x}_{\min}(t)}{\|\mathbf{w}_t\|} \\
&\stackrel{(21)}{=} \underbrace{\cos \gamma_t}_{\xrightarrow{t \rightarrow \infty} 1} \underbrace{\mathbf{w}_*^T \mathbf{x}_{\min}(t)}_{\geq \rho_*} + \underbrace{\frac{\mathbf{u}_t^T \mathbf{x}_{\min}(t)}{\|\mathbf{u}_t\|}}_{\leq 1} \underbrace{\sin \gamma_t}_{\xrightarrow{t \rightarrow \infty} 0} \quad (18)
\end{aligned}$$

If $\mathbf{x}_{\min}(t)$ is not a support vector, $\mathbf{w}_*^T \mathbf{x}_{\min}(t) > \rho_*$ holds. Due to (18), there is a t' where $\mathbf{x}_{\min}(t)$ being a non-support vector and $t > t'$ inevitably leads to a contradiction. Note that for $t > t'$ only support vectors are added to the set $\mathcal{X}'(t)$, i.e. there is a finite number of non-support vectors contained in the set $\mathcal{X}'(t)$. As a consequence after a finite number of iterations t'' also $\mathbf{x}_{\max}(t)$ will always be a support vector.

Now, we show that all non-support vectors in the set $\mathcal{X}'(t)$ will be removed. *Assumption:* There exists a sample \mathbf{x} that is not a support vector but it remains in the set $\mathcal{X}'(t)$, i.e. $\exists \mathbf{x} : \mathbf{x} \notin \mathcal{X}^{\text{sv}} \wedge \mathbf{x} \in \mathcal{X}'(t)$ for all t . This means that

$$\frac{\mathbf{w}_t}{\|\mathbf{w}_t\|} \mathbf{x} < \underbrace{\frac{\mathbf{w}_t}{\|\mathbf{w}_t\|} \mathbf{x}_{\max}(t)}_{\xrightarrow{t \rightarrow \infty} \rho_*} \quad (19)$$

always holds. This leads to a contradiction since after a finite number of iterations $\mathbf{x}_{\max}(t)$ will always be a support vector. \square



$$\mathbf{w}_t = \cos \gamma_t \|\mathbf{w}_t\| \mathbf{w}_* + \mathbf{u}_t \quad (20)$$

$$\|\mathbf{u}_t\| = \|\mathbf{w}_t\| \sin \gamma_t \quad (21)$$

$$\|\mathbf{w}_*\| = 1$$

Fig. 2. Orthogonal decomposition of \mathbf{w}_t and properties that hold within this decomposition.

3.2 Soft-OneClassMaxMinOver

So far, we only have considered the hard-margin problem. In order to realise a 2-norm soft-margin version of the OMMO algorithm, we consider the quadratic optimisation problem:

$$\min_{\mathbf{w}, \xi} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{2} \sum_i \xi_i^2 \right) \quad \text{s.t.} \quad \forall i : \mathbf{w}^T \phi(\mathbf{x}_i) \geq 1 - \xi_i \quad . \quad (22)$$

In the hard-margin case ($C \rightarrow \infty$) this is equivalent to the optimisation problem (4). Note, that compared to the optimisation problems (1) and (2) the constraint on each slack variable ($\xi_i \geq 0$) disappears. By constructing the primal Lagrangian of (22), setting the partial differentiations to zero and rearranging [6], we obtain

$$\min_{\alpha} \left(\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j \left(K(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{C} \delta_{ij} \right) \right) \quad \text{s.t.} \quad \forall i : \alpha_i \geq 0 \quad , \quad (23)$$

where δ_{ij} is the Kronecker delta which is 1 if $i = j$, and 0 otherwise. As mentioned in [6] this can be understood as solving the hard-margin problem in a modified kernel space. The modified kernel is $K(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{C} \delta_{ij}$. Hence, in order to implement a 2-norm soft-margin version of OMMO, we modify algorithm 1 such that

$$h(\mathbf{x}_i) = \sum_{j=1}^L \alpha_j \left(K(\mathbf{x}_j, \mathbf{x}_i) + \frac{1}{C} \delta_{ij} \right) \quad . \quad (24)$$

4 Experiments & Results

We applied the OMMO algorithm to artificial datasets and a real-world problem using Gaussian kernels. We created a sinusoid and an xor dataset each consisting of 250 samples (Fig. 3). The hyperparameters were set to extremal values and to more appropriate ones that can be determined, for instance, by cross-validation. The results on the artificial datasets are shown in Fig. 3. Similar to the approach (2) that implements a 1-norm slack term, different solutions ranging from hard to soft margin can be realised by controlling the parameter C , i.e. the relevance

of outliers can be controlled. Furthermore, we applied the OMMO algorithm to the problem of face detection where we used the MIT-CBCL face detection dataset¹ that contains 2901 images of faces and 28121 images of non-faces of size 19x19 pixels. The dataset is divided into a training set containing 2429 faces and 4548 non-faces and a test set containing 472 faces and 23573 non-faces. We used the raw data but performed the preprocessing steps described in [7] to reduce the within-class variance. Afterwards, we took the training set to perform a simple grid search over σ, C and chose randomly 1215 faces to train OMMO and tested the performance on a test set with 1214 faces and 4548 non-faces. To reduce the variance we performed 25 runs at all combinations of σ, C . The performance of OMMO for a fixed σ, C was evaluated by the equal-error-rate of the receiver-operator-characteristics (ROC). Having determined the optimal parameters σ, C , we trained OMMO with the whole training set of 2429 faces and computed the ROC curve of the 24045 test samples. The same steps were performed using the libSVM [8], except that here we used the parameter ν to control the softness. A comparison between both ROC curves is depicted in Fig. 4. Although these two approaches differ significantly in their implementation complexity, their performance is almost equal. The execution time of OMMO and libSVM cannot be compared directly because it depends heavily on implementation details.

5 Conclusions

Based on an existing two-class classification method, we proposed a very simple and incremental boundary approach, called **OneClassMaxMinOver**. OMMO can be realised by only a few lines of code, which makes it interesting particularly for practitioners. We proved that after a finite number of learning steps OMMO yields a maximum margin hyperplane that is described only by support vectors. Furthermore, we showed that the speed of convergence of OMMO is $\mathcal{O}(1/\sqrt{t})$ where t is number of iterations. Considering the ideas described in [9], even a $\mathcal{O}(1/t)$ convergence of OMMO can be expected.

By simply using a modified kernel function $K'(\mathbf{x}, \mathbf{y}) = K(\mathbf{x}, \mathbf{y}) + \frac{1}{C} \delta_{xy}$ OMMO can also realise a soft maximum margin solution controlled by the softness parameter C . Thus, OMMO can cope with datasets that also contain outliers.

In the future, a closer look at convergence speed and bounds on the target error will be made. Moreover, the problem of parameter validation will be examined, since in many one-class problems only target objects are available. Thus, standard validation techniques cannot be applied. If it is not possible to evaluate the target error, the complexity or volume of the boundary description have to be estimated in order to select good hyperparameters. Since, simple sampling techniques fail to measure the volume in high-dimensional input spaces, more sophisticated methods need to be derived.

¹ <http://cbcl.mit.edu/software-datasets/FaceData.html>

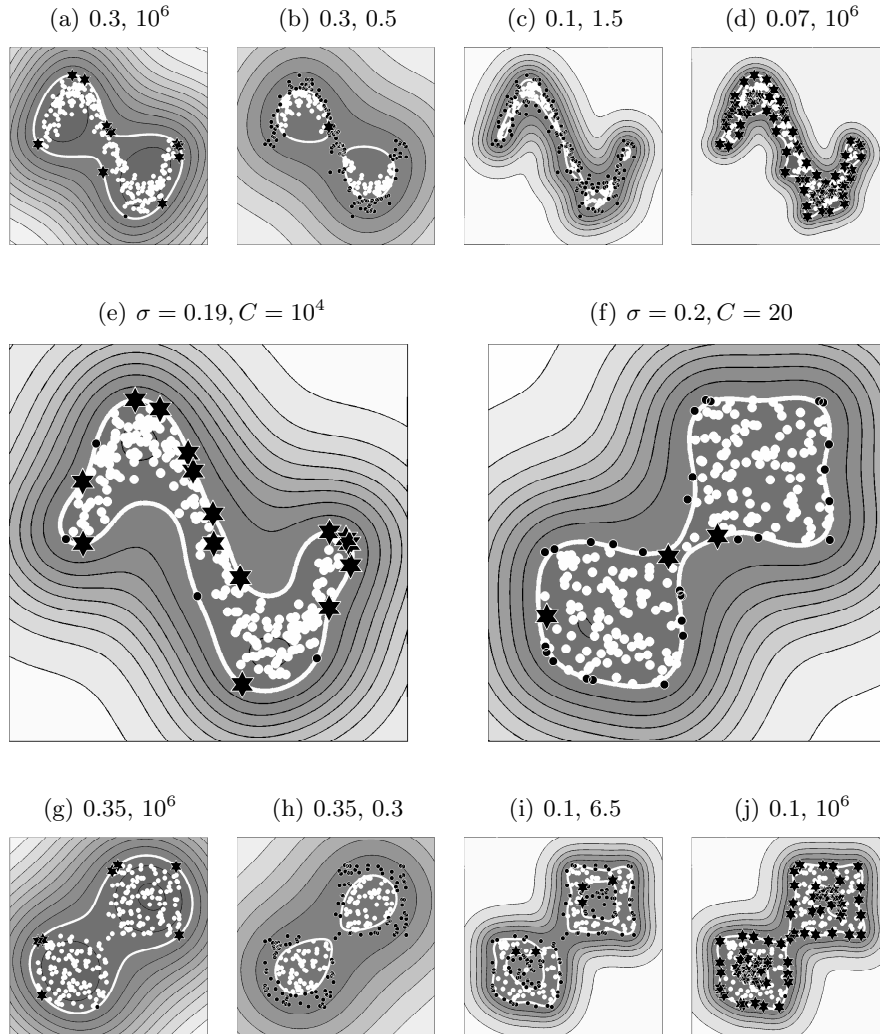


Fig. 3. Our algorithm applied to two exemplary artificial datasets – sinusoid, and xor. The parameters (σ, C) are shown above each graph. The stars depict support vectors that lie inside the hypersphere, dark circles depict support vectors outside the hypersphere. All other samples as well as the boundary are represented in white. In the first and third row extremal values for σ and C were chosen to achieve hard-margin ((a), (d), (g), (j)) as well as soft margin solutions ((b), (c), (h), (i)). In (b), (c), and (h) there is no support vector which lies inside the hypersphere and so the boundary is only influenced by support vectors from outside the hypersphere. The best solutions for the datasets are shown in the second row, where the values lie in between the extremal hard and soft margin solutions.

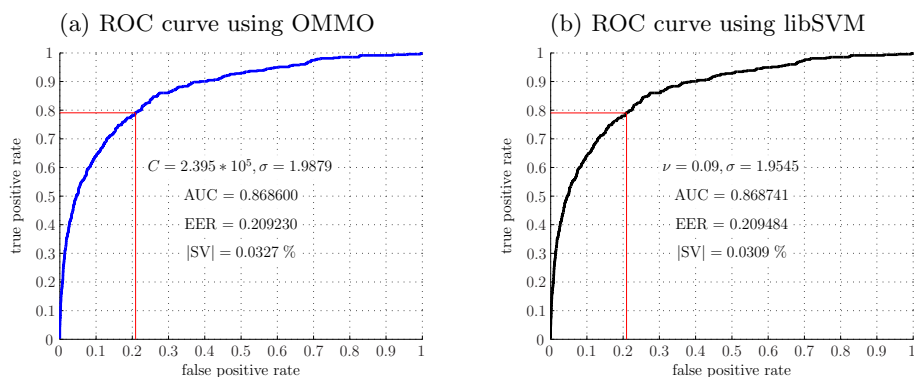


Fig. 4. The receiver-operator-characteristics shows that the two algorithms achieve the same performance on a test set of 472 faces and 23573 non-faces. Both models obtained by parameter validation are rather hard- than soft-margin (C large, ν small). They have a Gaussian width of $\sigma \approx 1.9$ and the fraction of support vectors is almost equal. The performance measured by the area under curve (AUC) for OMMO is the same as for libSVM. This holds also for the equal-error-rate (EER).

References

1. Vapnik, V.N.: The nature of statistical learning theory. Springer Verlag, Heidelberg, DE (1995)
2. Platt, J.: Fast training of support vector machines using sequential minimal optimization. In Schölkopf, B., Burges, C.J.C., Smola, A.J., eds.: *Advances in Kernel Methods — Support Vector Learning*, Cambridge, MA, MIT Press (1999) 185–208
3. Tax, D.M.J., Duin, R.P.W.: Data domain description using support vectors. In: *ESANN*. (1999) 251–256
4. Schölkopf, B., Platt, J.C., Shawe-Taylor, J., Smola, A.J., Williamson, R.C.: Estimating the support of a high-dimensional distribution. *Neural Computation* **13**(7) (2001) 1443–1471
5. Martinez, T.: MaxMinOver: A Simple Incremental Learning Procedure for Support Vector Classification. In: *IEEE Proceedings of the International Joint Conference on Neural Networks (IJCNN 2004)*, Budapest, Hungary (2004) 2065–2070
6. Cristianini, N., Shawe-Taylor, J.: *An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods*. Cambridge University Press, Cambridge, U.K. (2000)
7. Sung, K.K.: Learning and example selection for object and pattern detection. In: *MIT AI-TR*. (1996)
8. Chang, C.C., Lin, C.J.: LIBSVM: a library for support vector machines. (2001) Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
9. Martinez, T.: Minover revisited for incremental support-vector-classification. In Rasmussen, C.E., Bühlhoff, H.H., Schölkopf, B., Giese, M.A., eds.: *DAGM-Symposium*. Volume 3175 of *Lecture Notes in Computer Science*., Springer (2004) 187–194