# Large Scale Transductive SVMs

**Ronan Collobert**                                    COLLOBER@NEC-LABS.COM
*NEC Laboratories America,*
*4 Independence Way, Princeton, NJ08540, USA.*

**Fabian Sinz**                                        FABEE@TUEBINGEN.MPG.DE
*NEC Laboratories America,*
*4 Independence Way, Princeton, NJ08540, USA,*
*and Max Planck Insitute for Biological Cybernetics,*
*Spemannstrasse 38, 72076 Tuebingen, Germany.*

**Jason Weston**                                       JASONW@NEC-LABS.COM
*NEC Laboratories America,*
*4 Independence Way, Princeton, NJ08540, USA.*

**Léon Bottou**                                        LEONB@NEC-LABS.COM
*NEC Laboratories America,*
*4 Independence Way, Princeton, NJ08540, USA.*

## Abstract

We show how the Concave-Convex Procedure can be applied to Transductive SVMs, which traditionally requires solving a combinatorial search problem. This provides for the first time a highly scalable algorithm in the nonlinear case. Detailed experiments verify the utility of our approach. Software is available at `http://www.kyb.tuebingen.mpg.de/bs/people/fabee/transduction.html`.

**Keywords:**   transduction, transductive SVMs, semi-supervised learning, CCCP

## 1. Introduction

Transductive support vector machines (TSVMs) (Vapnik, 1995) are a method of improving the generalization accuracy of SVMs (Boser et al., 1992) by using unlabeled data. TSVMs, like SVMs, learn a large margin hyperplane classifier using labeled training data, but simultaneously force this hyperplane to be far away from the unlabeled data.

One way of justifying this algorithm, in the context of *semi-supervised learning* is that one is finding a decision boundary that lies in a region of low density, which implements the so-called cluster assumption (see e.g. Chapelle and Zien, 2005). In this framework, if you believe the underlying distribution of the two classes is such that there is a "gap" or low density region between them, then TSVMs can help because it selects a rule with exactly those properties. Vapnik (1995) has a different interpretation for the success of TSVMs, rooted in the idea that transduction (labeling a test set) is inherently easier than induction (learning a general rule). In either case, experimentally it seems clear that algorithms such as TSVM can give considerable improvement in generalization over SVMs, if the number of labeled points is small and the number of unlabeled points is large.

Unfortunately, one of the problems of research into TSVMs (and other semi-supervised approaches) is that the implementations often suffer from an inability to deal with a large number of unlabeled examples. The first implementation of TSVM appeared in Bennett and Demiriz (1998), using an integer programming method, intractable for large problems. Joachims (1999b) then proposed a combinatorial approach, known as SVMLight TSVM, that is practical for a few thousand examples. Fung and Mangasarian (2001) introduced a sequential optimization procedure that could potentially scale well, although their largest experiment used only 1000 examples. However, their method was for the linear case only, and used a special kind of SVM with a 1-norm regularizer, to retain linearity. Finally, Chapelle and Zien (2005) proposed a primal method, which turned out to show improved generalization performance over the previous approaches, but still scales as $(L+U)^3$, where $L$ and $U$ are the numbers of labeled and unlabeled examples. This method also stores the entire $(L+U) \times (L+U)$ kernel matrix in memory. Other methods such as in Bie and Cristianini (2004); Xu et al. (2005) transform the non-convex transductive problem into a convex semi-definite programming problem that scales as $(L+U)^4$ or worse.

In this article we introduce a large scale training method for TSVMs using the Concave-Concave Procedure (CCCP) (Yuille and Rangarajan, 2002; Thi, 1994). The current work is an expanded version of a technical report (Collobert et al., 2005). CCCP iteratively optimizes non-convex cost functions that can be expressed as the sum of a convex function and a concave function. The optimization is carried out iteratively by solving a sequence of convex problems obtained by linearly approximating the concave function in the vicinity of the solution of the previous convex problem. This method is guaranteed to find a local minimum and has no difficult parameters to tune. This provides what we believe is the best known method for implementing Transductive SVM with an empirical scaling of $(L+U)^2$, which involves training a sequence of typically 1-10 conventional convex SVM optimization problems. As each of these problems is trained in the dual we retain the SVMs linear scaling with problem dimensionality, in contrast to the techniques of Fung et al.

## 2. The Concave-Convex Procedure for TSVMs

**Notation** We consider a set of $L$ training pairs $\mathcal{L} = \{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_L, y_L)\}$, $\boldsymbol{x} \in \mathbb{R}^n$, $y \in \{1, -1\}$ and a (unlabeled) set of $U$ test vectors $\mathcal{U} = \{x_{L+1}, \ldots, x_{L+U}\}$. SVMs have a decision function $f_\theta(.)$ of the form

$$f_\theta(x) = w \cdot \Phi(x) + b\,,$$

where $\boldsymbol{\theta} = (\boldsymbol{w}, b)$ are the parameters of the model, and $\Phi(\cdot)$ is the chosen feature map, often implemented implicitly using the kernel trick (Vapnik, 1995).

**TSVM Formulation** The original TSVM optimization problem is the following (Vapnik, 1995; Joachims, 1999b; Bennett and Demiriz, 1998). Given a training set $\mathcal{L}$ and a test set $\mathcal{U}$, find among the possible binary vectors

$$\{\mathcal{Y} = (y_{L+1}, \ldots, y_{L+U})\}$$

the one such that an SVM trained on $\mathcal{L} \cup (\mathcal{U} \times \mathcal{Y})$ yields the largest margin.

This is a combinatorial problem, but one can approximate it (see Vapnik, 1995) as finding a SVM separating the training set under constraints which force the unlabeled examples to
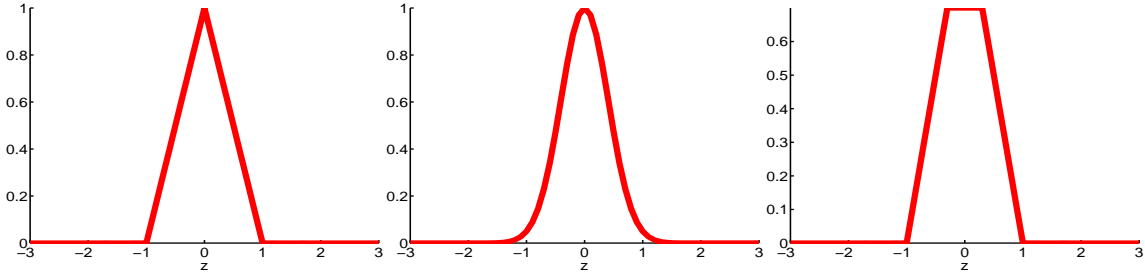
Figure 1: Three loss functions for unlabeled examples, from left to right (i) the Symmetric Hinge, (ii) Symmetric Sigmoid; and (iii) Symmetric Ramp loss.

be as far as possible from the margin. This can be written as minimizing

$$\frac{1}{2}\|\boldsymbol{w}\|^2 + C\sum_{i=1}^{L}\xi_i + C^*\sum_{i=L+1}^{L+U}\xi_i$$

subject to

$$y_i\,f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \geq 1 - \xi_i, \quad i = 1, \ldots, L$$

$$|f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)| \geq 1 - \xi_i, \quad i = L+1, \ldots, L+U$$

This minimization problem is equivalent to minimizing

$$J(\boldsymbol{\theta}) = \frac{1}{2}\|\boldsymbol{w}\|^2 + C\sum_{i=1}^{L}H_1(y_i\,f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + C^*\sum_{i=L+1}^{L+U}H_1(|f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)|), \tag{1}$$

where the function $H_1(\cdot) = \max(0, 1 - \cdot)$ is the classical Hinge Loss (Figure 2, center). The loss function $H_1(|\cdot|)$ for the unlabeled examples can be seen in Figure 1, left. For $C^* = 0$ in (1) we obtain the standard SVM optimization problem. For $C^* > 0$ we penalize unlabeled that is inside the margin. This is equivalent to using the hinge loss on the labeled data, but where we assume the label for the unlabeled example is $y_i = \text{sign}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i))$.

**Losses for transduction** TSVMs implementing formulation (1) were first introduced in SVMLight (Joachims, 1999b). As shown above, it assigns a Hinge Loss $H_1(\cdot)$ on the labeled examples (Figure 2, center) and a "Symmetric Hinge Loss" $H_1(|\cdot|)$ on the unlabeled examples (Figure 1, left). More recently, Chapelle and Zien (2005) proposed to handle unlabeled examples with a smooth version of this loss (Figure 1, center). While we also use the Hinge Loss for labeled examples, we use for unlabeled examples a slightly more general form of the Symmetric Hinge Loss, that we allow to be "non-peaky" (Figure 1, right). Given an unlabeled example $\boldsymbol{x}$ and using the notation $z = f_{\boldsymbol{\theta}}(\boldsymbol{x})$, this loss can be written as

$$z \mapsto R_s(z) + R_s(-z), \tag{2}$$

where $s < 1$ is a hyper-parameter to be chosen and $R_s$ is what we call the "Ramp Loss", a "cut" version of the Hinge Loss (Figure 2, left).

Losses similar to the Ramp Loss have been already used for different purposes, like in the Doom II algorithm introduced by Mason et al. (2000) or in the context of "Ψ-learning"

3

by Shen et al. (2003). The $s$ parameter controls where we "cut" the Ramp Loss, and as a consequence it also controls the wideness of the flat part of the loss (2) we use for transduction: when $s = 0$, this reverts to the Symmetric Hinge $H_1(|\cdot|)$. When $s \neq 0$, we obtain a non-peaked loss function (Figure 1, right) which can be viewed as a simplification of Chapelle's loss function. We call this loss function (2) the "Symmetric Ramp Loss".
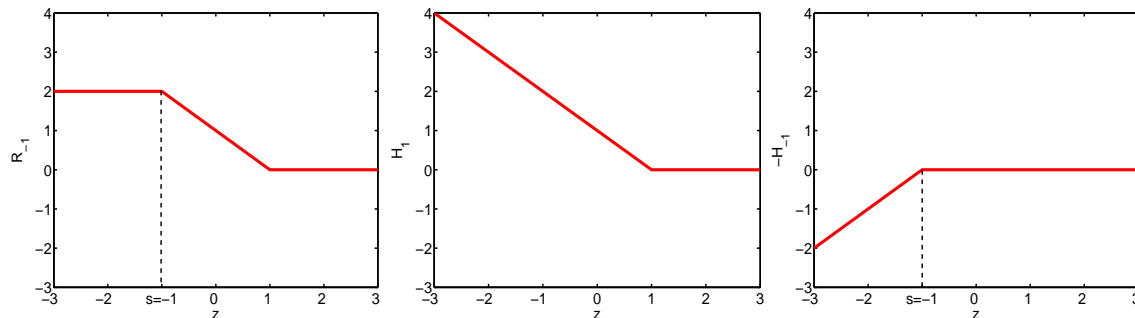


Figure 2: The Ramp Loss function (left) can be decomposed into the sum of the convex Hinge Loss (center) and a concave loss (right).

Training a TSVM using the loss function (2) is equivalent to training a SVM using the Hinge loss $H_1(\cdot)$ for labeled examples, and using the Ramp loss $R_s(\cdot)$ for unlabeled examples, where each unlabeled example appears as two examples labeled with both possible classes. More formally, after introducing

$$
\begin{aligned}
y_i &= 1 & i \in [L+1 \ldots L+U] \\
y_i &= -1 & i \in [L+U+1 \ldots L+2U] \\
\boldsymbol{x}_i &= \boldsymbol{x}_{i-U} & i \in [L+U+1 \ldots L+2U],
\end{aligned}
$$

we can rewrite (1) as

$$
J^s(\boldsymbol{\theta}) = \frac{1}{2}\|\boldsymbol{w}\|^2 + C \sum_{i=1}^{L} H_1(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + C^* \sum_{i=L+1}^{L+2U} R_s(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)). \tag{3}
$$

This is the minimization problem we now consider in the rest of the paper.

**Balancing constraint** One problem with TSVM as stated above is that in high dimensions with few training data, it is possible to classify all the unlabeled examples as belonging to only one of the classes with a very large margin, which leads to poor performance. To cure this problem, one further constrains the solution by introducing a balancing constraint that ensures the unlabeled data are assigned to both classes. Joachims (1999b) directly enforces that the fraction of positive and negatives assigned to the unlabeled data should be the same fraction as found in the labeled data. Chapelle and Zien (2005) use a similar but slightly relaxed constraint, which we also use in this work:

$$
\frac{1}{U} \sum_{i=L+1}^{L+U} f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) = \frac{1}{L} \sum_{i=1}^{L} y_i. \tag{4}
$$

4

**Concave-Convex Procedure (CCCP)** Unfortunately, the TSVM optimization problem as given above is not convex, and minimizing a non-convex cost function is usually difficult. Gradient descent techniques, such as conjugate gradient descent or stochastic gradient descent, often involve delicate hyper-parameters (LeCun et al., 1998). In contrast, convex optimization seems much more straight-forward. For instance, the SMO algorithm (Platt, 1999) locates the SVM solution efficiently and reliably.

We propose instead to optimize non convex problems using the "Concave-Convex Procedure" (CCCP) (Yuille and Rangarajan, 2002). The CCCP procedure is closely related to the "Difference of Convex" (DC) methods that have been developed by the optimization community during the last two decades (Thi, 1994). Such techniques have already been applied for dealing with missing values in SVMs (Smola et al., 2005), for improving boosting algorithms (Krause and Singer, 2004), and in the "$\Psi$-learning" framework (Shen et al., 2003).

Assume that a cost function $J(\boldsymbol{\theta})$ can be rewritten as the sum of a convex part $J_{vex}(\boldsymbol{\theta})$ and a concave part $J_{cav}(\boldsymbol{\theta})$. Each iteration of the CCCP procedure (Algorithm 1) approximates the concave part by its tangent and minimizes the resulting convex function.

---

**Algorithm 1** : The Concave-Convex Procedure (CCCP)

---

Initialize $\boldsymbol{\theta}^0$ with a best guess.

**repeat**

$$\boldsymbol{\theta}^{t+1} = \arg\min_{\boldsymbol{\theta}} \left( J_{vex}(\boldsymbol{\theta}) + J'_{cav}(\boldsymbol{\theta}^t) \cdot \boldsymbol{\theta} \right) \tag{5}$$

**until** convergence of $\boldsymbol{\theta}^t$

---

One can easily see that the cost $J(\boldsymbol{\theta}^t)$ decreases after each iteration by summing two inequalities resulting from (5) and from the concavity of $J_{cav}(\boldsymbol{\theta})$.

$$
\begin{aligned}
J_{vex}(\boldsymbol{\theta}^{t+1}) + J'_{cav}(\boldsymbol{\theta}^t) \cdot \boldsymbol{\theta}^{t+1} &\leq J_{vex}(\boldsymbol{\theta}^t) + J'_{cav}(\boldsymbol{\theta}^t) \cdot \boldsymbol{\theta}^t \quad\quad (6)\\
J_{cav}(\boldsymbol{\theta}^{t+1}) &\leq J_{cav}(\boldsymbol{\theta}^t) + J'_{cav}(\boldsymbol{\theta}^t) \cdot \left( \boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t \right) \quad\quad (7)
\end{aligned}
$$

The convergence of CCCP has been shown by Yuille and Rangarajan (2002) by refining this argument. The authors also showed that the CCCP procedure remains valid if $\boldsymbol{\theta}$ is required to satisfy some linear constraints. Note that no additional hyper-parameters are needed by CCCP. Furthermore, each update (5) is a convex minimization problem and can be solved using classical and efficient convex algorithms.

**CCCP for TSVMs** Interestingly, the Ramp Loss can be rewritten as the difference between two Hinge losses (see Figure 2):

$$R_s(z) = H_1(z) - H_s(z). \tag{8}$$

Because of this decomposition, the TSVM minimization problem as stated in (3) is amenable to CCCP optimization. The cost $J^s(\boldsymbol{\theta})$ can indeed be decomposed into a convex $J^s_{vex}(\boldsymbol{\theta})$

and concave $J_{cav}^s(\boldsymbol{\theta})$ part as follows:

$$
\begin{aligned}
J^s(\boldsymbol{\theta}) &= \frac{1}{2}\|\boldsymbol{w}\|^2 + C\sum_{i=1}^{L} H_1(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + C^*\sum_{i=L+1}^{L+2U} R_s(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) \\
&= \underbrace{\frac{1}{2}\|\boldsymbol{w}\|^2 + C\sum_{i=1}^{L} H_1(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + C^*\sum_{i=L+1}^{L+2U} H_1(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i))}_{J_{vex}^s(\boldsymbol{\theta})} \\
&\underbrace{-C^*\sum_{i=L+1}^{L+2U} H_s(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i))}_{J_{cav}^s(\boldsymbol{\theta})}\,.
\end{aligned}
\tag{9}
$$

This decomposition allows us to apply the CCCP procedure as stated in Algorithm 1. The convex optimization problem (5) that constitutes the core of the CCCP algorithm is easily reformulated into dual variables $\boldsymbol{\alpha}$ using the standard SVM technique.

After some algebra, we show in Appendix A that enforcing the balancing constraint (4) can be achieved by introducing an extra Lagrangian variable $\alpha_0$ and an example $\boldsymbol{x}_0$ implicitely defined by

$$
\Phi(\boldsymbol{x}_0) = \frac{1}{U}\sum_{i=L+1}^{L+U} \Phi(\boldsymbol{x}_i)\,,
$$

with label $y_0 = 1$. Thus, if we note $K$ the kernel matrix such that

$$
K_{ij} = \Phi(\boldsymbol{x}_i) \cdot \Phi(\boldsymbol{x}_j)\,,
$$

the column corresponding to the example $\boldsymbol{x}_0$ is computed as follow:

$$
K_{i0} = K_{0i} = \frac{1}{U}\sum_{j=L+1}^{L+U} \Phi(\boldsymbol{x}_j) \cdot \Phi(\boldsymbol{x}_i) \quad \forall i\,.
\tag{10}
$$

The computation of this special column can be achieved very efficiently by computing it only once, or by approximating the sum (10) using an appropriate sampling method.

Given the decomposition of the cost (9) and the trick of the special extra example (10) to enforce the balancing constraint, we can easily apply Algorithm 1 to TSVMs. To simplifiy the first order approximation of the concave part in the CCCP procedure (5), we denote[1]

$$
\beta_i \;=\; y_i\, \frac{\partial J_{cav}^s(\boldsymbol{\theta})}{\partial f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)} \;=\; \begin{cases} C^* & \text{if } y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) < s \\ 0 & \text{otherwise} \end{cases}\,,
\tag{11}
$$

for unlabeled examples (that is $i \geq L+1$). The concave part $J_{cav}^s$ does not depend on labeled examples ($i \leq L$) so we obviously have $\beta_i = 0$ for all $i \leq L$. This yields Algorithm 2, after some standard derivations detailed in Appendix A.

---

1. Note that $J_{cav}^s(\cdot)$ is non-differentiable at $z = s$, because $H_s(\cdot)$ is not. It can be shown that the CCCP remains valid when using any super-derivative of the concave function. Alternatively, the function $H_s(\cdot)$ could be made smooth in a small interval $[s - \varepsilon, s + \varepsilon]$.

---

**Algorithm 2** : CCCP for TSVMs

---

Initialize $\boldsymbol{\theta}^0 = (\boldsymbol{w}^0, b^0)$ with a standard SVM solution on the labeled points.

Compute $\beta_i^0 = \begin{cases} C^* & \text{if } y_i\, f_{\boldsymbol{\theta}^0}(\boldsymbol{x}_i) < s \text{ and } i \geq L+1 \\ 0 & \text{otherwise} \end{cases}$

Set $\zeta_i = y_i$ for $1 \leq i \leq L + 2U$ and $\zeta_0 = \frac{1}{L}\sum_{i=1}^{L} y_i$

**repeat**

- **Solve** the following convex problem ( with $K_{ij} = \Phi(\boldsymbol{x}_i) \cdot \Phi(\boldsymbol{x}_j)$ )

$$\max_{\boldsymbol{\alpha}} \left( \boldsymbol{\alpha} \cdot \boldsymbol{\zeta} - \frac{1}{2}\,\boldsymbol{\alpha}^{\mathrm{T}} \boldsymbol{K}\,\boldsymbol{\alpha} \right) \text{ subject to } \begin{cases} \boldsymbol{\alpha} \cdot \mathbf{1} = 0 \\ 0 \leq y_i\, \alpha_i \leq C \quad \forall 1 \leq i \leq L \\ -\beta_i \leq y_i\, \alpha_i \leq C^* - \beta_i \quad \forall i \geq L+1 \end{cases}$$

- Define $\boldsymbol{w}^{t+1} = \displaystyle\sum_{i=0}^{L+2U} \alpha_i\, \Phi(\boldsymbol{x}_i).$

- **Compute** $b^{t+1}$ using

$$0 < y_i\, \alpha_i < C \implies y_i\,(\boldsymbol{w}^{t+1} \cdot \Phi(x_i) + b^{t+1}) = 1 \quad \forall 1 \leq i \leq L$$

$$-\beta_i < y_i\, \alpha_i < C^* - \beta_i \implies y_i\,(\boldsymbol{w}^{t+1} \cdot \Phi(x_i) + b^{t+1}) = 1 \quad \forall i \geq L+1$$

- **Compute** $\beta_i^{t+1} = \begin{cases} C^* & \text{if } y_i\, f_{\boldsymbol{\theta}^{t+1}}(\boldsymbol{x}_i) < s \text{ and } i \geq L+1 \\ 0 & \text{otherwise} \end{cases}$

**until** $\boldsymbol{\beta}^{t+1} = \boldsymbol{\beta}^t$

---

Convergence of Algorithm 2 in finite time $t^*$ is guaranteed because variable $\boldsymbol{\beta}$ can only take a finite number of distinct values, because $J(\boldsymbol{\theta}^t)$ is decreasing, and because inequality (7) is strict unless $\boldsymbol{\beta}$ remains unchanged.

**Complexity** The main point we want to emphasize in this paper is the advantage in terms of training time of our method compared to existing approaches. Training a CCCP-TSVM amounts to solve a series of SVM optimization problems with $L + 2U$ variables. Although SVM training has a worst case complexity of $\mathcal{O}((L + 2U)^3)$ it typically scales quadratically (see Joachims, 1999a; Platt, 1999), and we find this is the case for our TSVM subproblems as well. Assuming a constant number of iteration steps the whole optimization of TSVMs with CCCP should scale quadratically in most practical cases (cf. Figure 3 and Figure 6). From our experience, around five iteration steps are enough in most of the cases to reach the minimum, see e.g. Figure 4.

## 3. Previous Work

**SVMLight-TSVM** Like our work, the heuristic optimization algorithm implemented in SVMLight (Joachims, 1999b) solves successive SVM optimization problems, but on $L + U$ instead of $L + 2U$ data points. It improves the objective function by iteratively switching the labels of two unlabeled points $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$ with $\xi_i + \xi_j > 2$. It uses two nested loops to optimize a TSVM which solves a quadratic program in each step. The convergence proof of the inner loop relies on the fact that there is only a finite number $2^U$ of labelings of $U$ unlabeled points, even though it is unlikely that all of them are examined. But since the heuristic only swaps the labels of two unlabeled examples at each step in order to enforce the balancing constraint it might need many iterations to reach a minimum, which makes it intractable for big dataset sizes in practice (cf. Figure 3).

**∇TSVM** The ∇TSVM of Chapelle and Zien (2005) is optimized by performing gradient descent in the primal space: minimize

$$\frac{1}{2}\|\boldsymbol{w}\|^2 + C \sum_{i=1}^{L} H^2(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + C^* \sum_{i=L+1}^{L+U} H^*(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)),$$

where $H^2(t) = \max(1 - t)$ and $H^*(t) = \exp(-3t^2)$ (cf. Figure 1, center). This optimization problem can be considered a smoothed version of (1).

Since the gradient descent is carried out in the primal, to learn nonlinear functions it is necessary to perform kernel PCA (Schölkopf et al., 1997). The overall algorithm has a time complexity equal to the square of the number of variables times the costs of evaluating the cost function. In this case, evaluating the objective scales linear in the number of examples, so the overall worst case complexity of solving the optimization problem for ∇TSVM is $\mathcal{O}((U+L)^3)$. The KPCA calculation alone also has a time complexity of $\mathcal{O}((U+L)^3)$. This method also requires to store the entire kernel matrix in memory, which clearly becomes infeasible for large datasets.

**CS$^3$VM** The work of Fung and Mangasarian (2001) is algorithmically the closest TSVM approach to our proposal. Following the formulation of transductive SVMs found in Bennett and Demiriz (1998), the authors consider transductive linear SVMs with a 1-norm regularizer, which allow them to decompose the corresponding loss function as a sum of a

| data set | classes | dims | points | labeled |
|----------|---------|------|--------|---------|
| g50c | 2 | 50 | 500 | 50 |
| Coil20 | 20 | 1024 | 1440 | 40 |
| Text | 2 | 7511 | 1946 | 50 |
| Uspst | 10 | 256 | 2007 | 50 |

Table 1: Small-Scale Datasets. We used the same datasets and experimental setup in these experiments as found in Chapelle and Zien (2005).

linear function and a concave function. Bennett proposed the following formulation which is similar to (1): minimize

$$||\boldsymbol{w}||_1 + C \sum_{i=1}^{L} \xi_i + C^* \sum_{i=L+1}^{U} \min(\xi_i, \xi_i^*)$$

subject to

$$y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \geq 1 - \xi_i, \quad i = 1, \ldots, L$$

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \geq 1 - \xi_i, \quad i = L+1, \ldots, L+U$$

$$-(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geq 1 - \xi_i^*, \quad i = L+1, \ldots, L+U$$

$$\xi_i \geq 0, \xi_i^* \geq 0.$$

The last term of the objective function is nonlinear and corresponds to the loss function given in Figure 1, left. To deal with it, the author's suggest to iteratively linearly approximate the concave part as a linear function, leading to a series of linear programming problems. This can be viewed as a simplified subcase of CCCP (a linear function being convex) applied to a special kind of SVM. Note also that the algorithm presented in their paper did not implement a balancing constraint for the labeling of the unlabeled examples as in (4). Our transduction algorithm is nonlinear and the use of kernels, solving the optimization in the dual, allows for large scale training with high dimensionality and number of examples.

## 4. Experiments

Our experimental investigation consisted of (i) small scale experiments where we could compare with existing TSVM approaches; and (ii) large scale experiments where we compare only our CCCP-TSVM approach against SVMs.

### 4.1 Small Datasets

We first performed experiments using the setup of Chapelle and Zien (2005) on the datasets given in Table 1. All datasets except g50c are real world datasets. The labels in g50c correspond to two Gaussians in a 50-dimensional space. The means of those Gaussians are placed in such a way, that the Bayes error is 5%. The coil20 data is a set of gray-scale images of 20 different objects taken from different angles, in steps of 5 degrees (S.A.Nene and H.Murase, 1996). The text dataset consists of the classes mswindows and mac of the

9

|  | Coil20 | g50c | Text | Uspst |
|---|---|---|---|---|
| SVM | 24.64 | 8.32 | 18.86 | 23.18 |
| SVMLight-TSVM | 26.26 | 6.87 | 7.44 | 26.46 |
| $\nabla$TSVM | 17.56 | 5.80 | 5.71 | 17.61 |
| CCCP-TSVM$|_{s=0}$ | 17.86 | 5.04 | 7.94 | 18.38 |
| CCCP-TSVM | 17.28 | 4.88 | 5.71 | 18.08 |

Table 2: Results on Small-Scale Datasets. We report the best test error over the hyperparameters of the algorithms, as in the methodology of Chapelle and Zien (2005). SVMLight-TSVM is the implementation in SVMLight, $\nabla$TSVM is the primal gradient descent method of Chapelle and Zien (2005) and CCCP-TSVM is our method. CCCP-TSVM with $s = 0$ is our method using the Symmetric Hinge Loss (Figure 1, left). The latter does not perform as well as the same loss function with a plateau, the Symmetric Ramp Loss (Figure 1, right), termed CCCP-TSVM.

`Newsgroup20` dataset preprocessed as in Szummer and Jaakkola (2001a). The `usps` dataset is the test part of the `USPS` hand written digit data. All datasets are split into ten parts with each part having a small amount of labeled examples and using the remaining part as unlabeled data.

To be able to compare our results to those of Chapelle and Zien (2005) we chose the same experimental setup. All methods use an RBF kernel and $\gamma$ and $C$ are tuned on the test set. For CCCP-TSVMs we also tune $C^*$ and $s$. Following Chapelle et al., we report the test error of the best fixed parameters optimized on the mean test error over all splits. Chapelle and Zien (2005) calculated this measure to calculate the best possible performance of the algorithms they compared their proposed method with, which was trained using a cross validation scheme.

The results are reported in Table 1. CCCP-TSVM achieves approximately the same error rates on all datasets as the $\nabla$TSVM in Chapelle and Zien (2005), and appears to be superior to SVMLight-TSVM. Training our algorithm with $s = 0$ (using the Symmetric Hinge loss - Figure 1, left) rather than the non-peaked loss of the Symmetric Ramp function - (Figure 1, right) also yields results similar to, but sometimes slightly worse than, $\nabla$TSVM. It appears that the Symmetric Ramp function is a better choice of loss function. Figure 5 highlights the importance of the parameter $s$ of the loss function (2), at least for the text dataset. Our conjecture is that the peaked loss of the Symmetric Hinge function forces early decisions for the $\boldsymbol{\beta}$ variables and might lead to a poor local optimum. This effect then disappears as soon as we clip the loss.

The CCCP algorithm in these experiments was implemented in C++ and the source code is available at `http://www.kyb.tuebingen.mpg.de/bs/people/fabee/transduction.html`. Without any particular optimization, CCCP TSVMs run orders of magnitude faster than SVMLight TSVMs, see Figure 3. Finally, Figure 4 shows the value of the objective function and test error during the CCCP iterations of training TSVM on two datasets. CCCP-TSVM tends to converge after only a few (around 5-10) iterations.
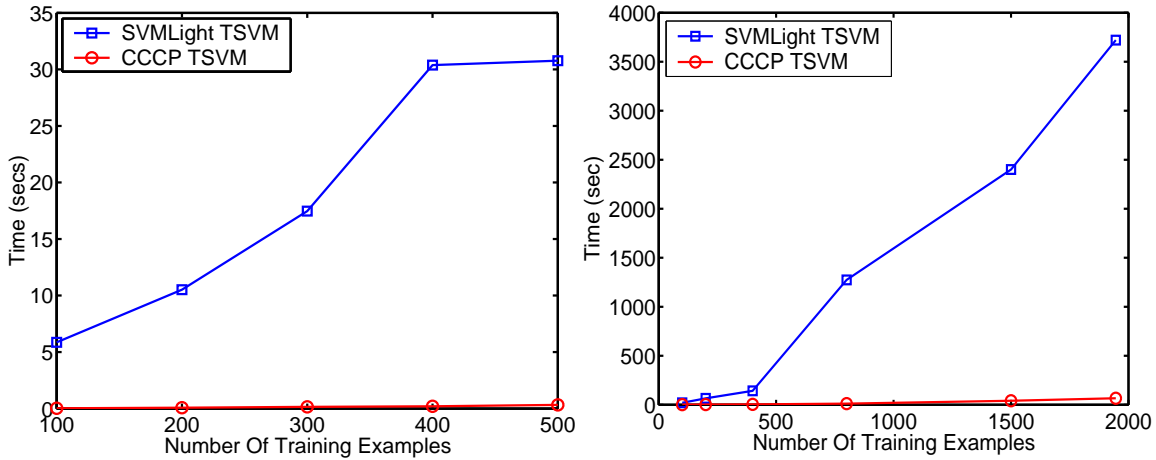
10

Figure 3: Training times for g50c (left) and Text (right) with SVMLight TSVMs and CCCP TSVMs. These plots use a single trial.
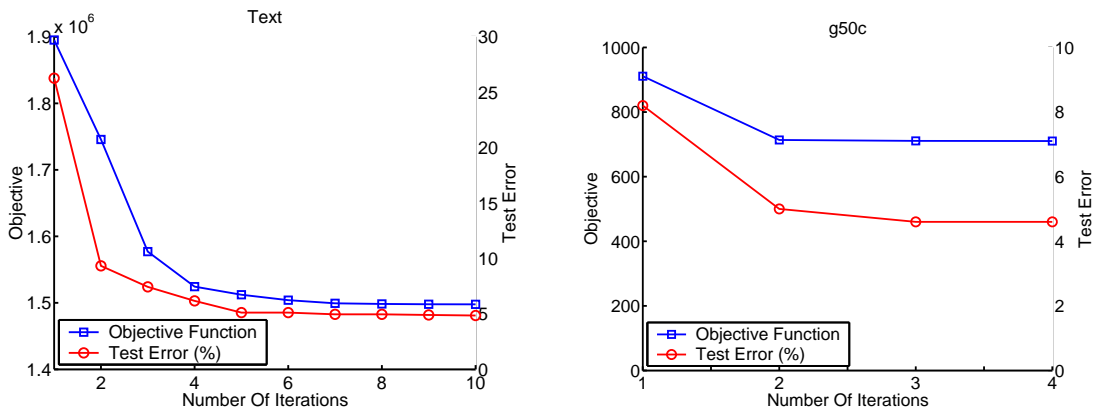


Figure 4: Value of the objective function and test error during the CCCP iterations of training TSVM on two datasets (single trial), Text (left) and g50c (right). CCCP-TSVM tends to converge after only a few iterations.

## 4.2 Large Scale Experiments

**RCV1** The first large scale experiment that we conducted was to separate the two largest top-level categories CCAT (CORPORATE/INDUSTRIAL) and GCAT (GOVERNMENT/SOCIAL) of the training part of the Reuters dataset as provided by Lewis et al. (2004). The set of these two categories consists of 17754 documents in a bag of words format, weighted with a TF.IDF scheme and normalized to length one. We performed experiments using 100 and 1000 labeled examples. We perfomed model selection on the validation set to find the parameters $s$, $C$, $C^*$ and the kernel parameter $\sigma$ using an RBF kernel. For model selection we use a validation set with 2000 and 4000 labeled examples for the two experiments. The remaining 12754 examples were used for testing. For each set of optimal hyperparameters we carried out another modeselection for the parameters $C^*$ and $s$ on subsamples of the $12.7k$ test set using the original training and validation sets but parts of the test set as unlabeled data. We tested the algorithm on the whole $12.7k$
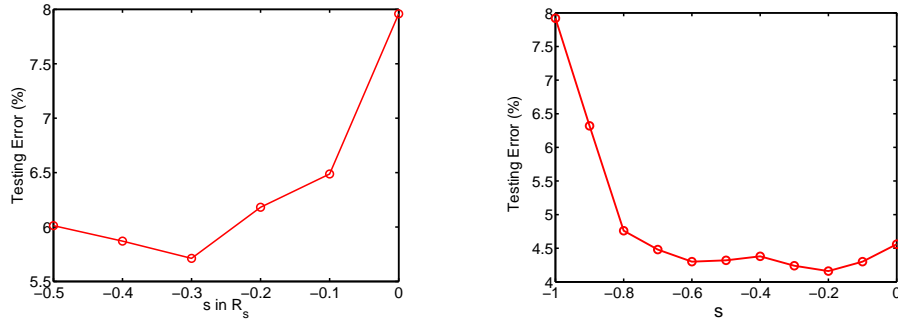
11

Figure 5: Effect of the parameter $s$ of the Symmetric Ramp loss (see Figure 1 and equation (2) ) on the Text dataset (left) and the g50c dataset (right). The peaked loss of the Symmetric Hinge function ($s = 0$) forces early decisions for the $\beta$ variables and might lead to a poor local optimum. This effect then disappears as soon as we clip the loss.

| Method | Train size | Unlabeled size | Parameters | Test Error |
|--------|------------|----------------|------------|------------|
| SVM | 100 | 0 | $C = 252.97,\ \sigma = 15.81$ | 16.61% |
| TSVM | 100 | 500 | $C = 2.597, C^* = 10,\ s = -0.2,\ \sigma = 3.95$ | 11.99% |
| TSVM | 100 | 1000 | $C = 2.597, C^* = 10,\ s = -0.2,\ \sigma = 3.95$ | 11.67% |
| TSVM | 100 | 2000 | $C = 2.597, C^* = 10,\ s = -0.2,\ \sigma = 3.95$ | 11.47% |
| TSVM | 100 | 5000 | $C = 2.597, C^* = 2.5297,\ s = -0.2,\ \sigma = 3.95$ | 10.65% |
| TSVM | 100 | 10000 | $C = 2.597, C^* = 2.5297,\ s = -0.4,\ \sigma = 3.95$ | 10.64% |
| SVM | 1000 | 0 | $C = 25.297,\ \sigma = 7.91$ | 11.04% |
| TSVM | 1000 | 500 | $C = 2.597, C^* = 10,\ s = -0.4,\ \sigma = 3.95$ | 11.09% |
| TSVM | 1000 | 1000 | $C = 2.597, C^* = 2.5297,\ s = -0.4,\ \sigma = 3.95$ | 11.06% |
| TSVM | 1000 | 2000 | $C = 2.597, C^* = 10,\ s - 0.4 =,\ \sigma = 3.95$ | 10.77% |
| TSVM | 1000 | 5000 | $C = 2.597, C^* = 2.5297,\ s = -0.2,\ \sigma = 3.95$ | 10.81% |
| TSVM | 1000 | 10000 | $C = 2.597, C^* = 25.2970,\ s = -0.4,\ \sigma = 3.95$ | 10.72% |

Table 3: Comparing CCCP-TSVMs with SVMs on the RCV1 problem for different number of labeled and unlabeled examples. See text for details.

test set using the different subsamples as unlabeled data and the optimal hyperparameters obtained by the second model selection. A selection of the results can be seen in Table 3.

The best result we obtained for $1k$ training points was 10.58% test error using $10.5k$ unlabeled points and 10.42% using $9.5k$ unlabeled points for $0.1k$ training points. Compared to the best performance of a SVM of 11.04% for the former and 16.61% for the latter this shows that unlabeled data can improve the results on this problem, especially in the case of few training data, where the improvement in test error is around 5.5%. However, when enough training data is available to the algorithm, the improvement is only in the order of one percent.

Figure 6 shows the training time of CCCP optimization as a function of the number of unlabeled examples. On a 64 bit opteron processor the optimization time for $12.5k$ unlabeled examples was approximately 18 min using the $1k$ training set and 69 min using $0.1k$ training
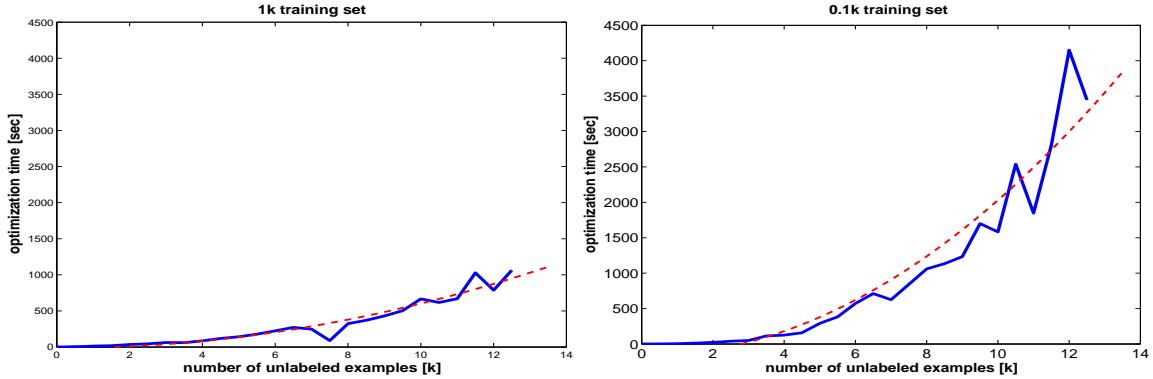
Figure 6: Optimization time for the Reuters dataset as a function of the number of unlabeled data. The algorithm was trained on 1,000 points (left) and on 100 points (right). The dashed lines represent a parabola fitted at the time measurements.

points. Although the worst case complexity of SVMs is cubic and the optimization time seems to be dependent on the ratio of the number of labeled to unlabeled examples, the training times clearly show a quadratic trend.

**MNIST** In the second large scale experiment, we conducted experiments on the MNIST handwritten digit database, as a 10-class problem. The original data has 60,000 training examples and 10,000 testing examples. We subsampled the training set for labeled points, and used the test set for unlabeled examples (or the test set + remainder of the training set when using more than 10,000 unlabeled examples). We performed experiments using 100 and 1000 labeled examples. We performed model selection for SVMs by trying a grid of values for $\sigma$ and $C$ and selecting the best ones using a separate validation set of size 1000. For efficiency reasons, for the TSVM we fixed $\sigma$ and $C$ to have the same values also for the TSVM. We then performed model selection using 2000 unlabeled examples to find the best choices of $C^*$ and $s$ using the validation set. When using more unlabeled data, we only reperformed model selection on $s$ as it appeared that this parameter was the most sensitive to changes in the unlabeled set, and kept the other parameters fixed. For the larger labeled set we took 2000, 5000 and 10000, 20000, 40000 and 60000 unlabeled examples. We always measure the error rate on the complete test set. The test error and parameter choices for each experiment are given in the Table 4.

The result show an improvement over SVM for CCCP-TSVMs which increases steadily as the number of unlabeled examples increases. Most experiments in semi-supervised learning only use a few labeled examples and do not use as many unlabeled examples as described here. It is thus reassuring to know that these methods do not apply just to toy examples with around 50 training points, but gains are still possible with more realistic dataset sizes.

## 5. Discussion

### 5.1 Cluster kernels and manifold-learning

Transductive SVM is not the only method of leveraging unlabeled data in a supervised learning task. In recent years this has become a popular research topic, and a battery of

13

| Method | Train size | Unlabeled size | Parameters | Test Error |
|--------|-----------|----------------|------------|------------|
| SVM | 100 | 0 | $C = 10$, $\gamma = 0.0128$ | 23.44% |
| TSVM | 100 | 2000 | $C^* = 1$, $s = -0.1$ | 16.81% |
| SVM | 1000 | 0 | $C = 10$, $\gamma = 0.0128$ | 7.77% |
| TSVM | 1000 | 2000 | $C^* = 5$, $s = -0.1$ | 7.13% |
| TSVM | 1000 | 5000 | $C^* = 1$, $s = -0.1$ | 6.28% |
| TSVM | 1000 | 10000 | $C^* = 0.5$, $s = -0.1$ | 5.65% |
| TSVM | 1000 | 20000 | $C^* = 0.3$, $s = -0.1$ | 5.43% |
| TSVM | 1000 | 40000 | $C^* = 0.2$, $s = -0.1$ | 5.31% |
| TSVM | 1000 | 60000 | $C^* = 0.1$, $s = -0.1$ | 5.38% |

Table 4: Comparing CCCP-TSVMs with SVMs on the MNIST problem for different number of labeled and unlabeled examples. See text for details.

techniques have been proposed. Some of the main methods, which we refer to as cluster kernels, do not change the learning algorithm at all, but merely the representation of the data as a pre-processing step. In a purely unsupervised fashion, these methods learn cluster or manifold structure from the data, and produce a new representation of it such that distances between points in the new space are small if they are in the same cluster or on the same manifold. Some of the main methods include Chapelle et al. (2002); Chapelle and Zien (2005); Sindhwani et al. (2005); Szummer and Jaakkola (2001b) and Weston et al. (2003).

Other notable methods include generalizations of nearest-neighbor or Parzen window type approaches to learning manifolds given labeled data, e.g. Zhu et al. (2003); Belkin and Niyogi (2002); Zhou et al. (2004). Finally, Bayesian approaches have also been pursued (Graepel et al., 2000; Lawrence and Jordan, 2005).

We note that some cluster kernel methods such as in Chapelle and Zien (2005) can perform significantly better than TSVM on some datasets. However, in the same paper Chapelle et al. show that in fact as these methods provide a new representation, one can just as easily run TSVM on the new representation as well. Then, the combination of TSVM and cluster kernels provides state-of-the-art results.

## 5.2 Semi-supervised versus transductive learning

From a theoretical point of view, there is much ideological debate over which underlying theory that explains TSVM is correct.

**Semi-supervised school** The majority of researchers appear to be in the *semi-supervised* school of thought, which claims that TSVMs help simply because of a regularizer that reflects prior knowledge, see e.g. Chapelle and Zien (2005). That is, one is given a set of unlabeled data, and one uses it to improve an inductive classifier to improve its generalization on an unseen test set.

**Transductive school** Transduction's aim, as argued by its inventor Vapnik, is to improve the generalization ability on a test set, given the test set as unlabeled data. Vapnik claims

|                      | Text       | MNIST-8   |
| -------------------- | ---------- | --------- |
| SVM                  | 18.86%     | 6.68%     |
| semi-supervised TSVM | 6.60%      | 5.27%     |
| transductive TSVM    | **6.12%**  | **4.87%** |

Table 5: Transductive CCCP-TSVM versus semi-supervised TSVMs compared on the Text and MNIST (8-vs-rest) datasets. Using the test set as unlabeled data appears to perform better than using a unlabeled set separate from the test set, see text for details.

that this can mean solving a less complex problem than induction, where one must first learn a general inductive rule and then apply it using deduction. One could potentially label the test data without building an inductive rule. He has also obtained a generalization bound for TSVMs that is slightly better than the one for SVMs (Vapnik, 1995).

**Experimental comparison**   If one takes transductive learning to mean just applying semi-supervised learning in the case where the unlabeled data is the test set of interest, then we can objectively compare this to semi-supervised learning (using an unlabeled set different from the test set).

We constructed some experiments to explore this difference. Taking the Text dataset, we split the test set into two pieces: (a) unlabeled data; and (b) test data. For each of the 10 splits we only report the generalization error on part (b). We then compared CCCP-TSVMs using unlabeled data (a) for training to using the test set (b) for training. If transduction is better than semi-supervised learning, using (b) should be better than (a). We performed a similar experiment on MNIST (digit 8 versus rest, 1000 unlabeled examples, but evaluating generalization on the whole test set of 10000 examples).

For Text, we fixed to a linear kernel, $C = 1000$, and $s = -0.3$. For MNIST-8 we fixed $\gamma = 0.0128$ and $C = 10$. We report the best test error over possible values of $C^*$. The results, shown in table 5.2 indicate that transductive TSVM is better than semi-supervised TSVM[2].

One intuitive argument for the improvement is simply that the test set is in a region of the space more relevant to the test error measured (the very points where the unlabeled data lie!). If the other unlabeled data in set (a) happen to be far away from these points, the changes to the hyperplane (inductive rule) may be less relevant. By this experiment we do not claim to solve the idealogical battle among researchers, but provide some evidence that perhaps transductive inference (even in its full sense, of avoiding induction) may be meaningful. In particular, if the test data is not i.i.d. with respect to the training data then we believe transduction could be particularly worthwhile.

## 6. Conclusion

In this article we have described an algorithm for TSVMs using CCCP that brings scalability improvements over existing implementation approaches. It involves the iterative solving of

---

2. Using the Wilcoxon signed rank test ($\alpha = 0.05$) one obtains p-values of 0.1641 and 0.3750 for Text and MNIST-8 respectively comparing on the 10 splits. One obtains 0.0627 and 0.0176 comparing over the every value of $C^*$ and every split (giving 60 trials in effect).

standard dual SVM QP problems, and usually requires just a few iterations. One nice thing about being an extension of standard SVM training is that any improvements in SVM scalability can immediately also be applied to TSVMs. For example in the linear case, one could easily apply fast linear SVM training such as in Keerthi and DeCoste (2005) to produce very fast linear TSVMs. For the nonlinear case, one could apply the online SVMs training scheme of Bordes et al. (2005) to give a fast online transductive learning procedure.

## Acknowledgments

## Appendix A. Derivation of the optimization problem

We are interested in minimizing the TSVM cost function (3), under the constraint (4). We rewrite the problem here for convenience: minimizing

$$J^s(\boldsymbol{\theta}) = \frac{1}{2}\|\boldsymbol{w}\|^2 + C\sum_{i=1}^{L} H_1(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + C^* \sum_{i=L+1}^{L+2U} R_s(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i))\,, \tag{12}$$

under the constraint

$$\frac{1}{U}\sum_{i=L+1}^{L+U} f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \;=\; \frac{1}{L}\sum_{i=1}^{L} y_i\,. \tag{13}$$

Assume that a cost function $J(\boldsymbol{\theta})$ can be rewritten as the sum of a convex part $J_{vex}(\boldsymbol{\theta})$ and a concave part $J_{cav}(\boldsymbol{\theta})$. As mentioned above in Algorithm 1, the minimization of $J(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$ ($\boldsymbol{\theta}$ being possibly restricted to a space $\mathcal{A}$ defined by some linear constraints) can be achieved by iteratively updating the parameters $\boldsymbol{\theta}$ using the following update

$$\boldsymbol{\theta}^{t+1} = \arg\min_{\boldsymbol{\theta}\in\mathcal{A}} \left( J_{vex}(\boldsymbol{\theta}) + J'_{cav}(\boldsymbol{\theta}^t)\cdot\boldsymbol{\theta} \right)\,. \tag{14}$$

In the case of our cost (12), we showed (see (9)) that $J^s(\boldsymbol{\theta})$ can be decomposed into the sum of $J^s_{vex}(\boldsymbol{\theta})$ and $J^s_{cav}(\boldsymbol{\theta})$ where

$$J^s_{vex}(\boldsymbol{\theta}) = \frac{1}{2}\,\|\boldsymbol{w}\|^2 + C\sum_{i=1}^{L} H_1(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) + C^* \sum_{i=L+1}^{L+2U} H_1(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)) \tag{15}$$

and

$$J^s_{cav}(\boldsymbol{\theta}) = -C^* \sum_{i=L+1}^{L+2U} H_s(y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i))\,. \tag{16}$$

16

In order to apply the CCCP update (14) we first have to calculate derivative of the concave part (16) with respect to $\boldsymbol{\theta}$:

$$\frac{\partial J_{cav}^s(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -C^* \sum_{i=L+1}^{L+2U} \frac{\partial J_{cav}^s(\boldsymbol{\theta})}{\partial f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)} \frac{\partial f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)}{\partial \boldsymbol{\theta}}$$

We introduce the notation

$$\begin{aligned}
\beta_i &= y_i \frac{\partial J_{cav}^s(\boldsymbol{\theta})}{\partial f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)} \\
&= \begin{cases} C^* H_s'[y_i f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)] & \text{if } i \geq L+1 \\ 0 & \text{otherwise} \end{cases} \\
&= \begin{cases} C^* & \text{if } y_i f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) < s \text{ and } i \geq L+1 \\ 0 & \text{otherwise} \end{cases}
\end{aligned}.$$

Since $f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) = \boldsymbol{w} \cdot \Phi(\boldsymbol{x}_i) + b$ with $\theta = (w, b)$, and $\partial f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)/\partial \boldsymbol{\theta} = (\Phi(\boldsymbol{x}_i), 1)$, each update (14) of the CCCP procedure applied to the our minimization problem (12) consists in minimizing the following cost

$$\begin{aligned}
J_{vex}^s(\boldsymbol{\theta}) + \frac{\partial J_{cav}^s(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \cdot \boldsymbol{\theta} &= J_{vex}^s(\boldsymbol{\theta}) + \left( \sum_{i=L+1}^{L+2U} y_i \beta_i \frac{\partial f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)}{\partial \boldsymbol{\theta}} \right) \cdot \boldsymbol{\theta} \\
&= J_{vex}^s(\boldsymbol{\theta}) + \sum_{i=L+1}^{L+2U} \beta_i y_i [\boldsymbol{w} \cdot \Phi(\boldsymbol{x}_i) + b] ,
\end{aligned} \tag{17}$$

under the linear constraint (13).

The convex part (16) contains Hinge Losses which can be rewritten as

$$H_1(z) = \max(0, 1 - z) = \min \xi \quad \text{s.t } \xi \geq 0, \, \xi \geq 1 - z .$$

It is thus easy to see that the minimization of (17) under the constraint (13) is equivalent to the following quadratic minimization problem under constraints:

$$\underset{\boldsymbol{\theta}, \boldsymbol{\xi}}{\arg\min} \quad \frac{1}{2} ||\boldsymbol{w}||^2 + C \sum_{i=1}^{L} \xi_i + C^* \sum_{i=L+1}^{L+2U} \xi_i + \sum_{i=L+1}^{L+2U} \beta_i y_i f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

$$\text{s.t.} \quad \frac{1}{U} \sum_{i=L+1}^{L+U} f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) = \frac{1}{L} \sum_{i=1}^{L} y_i \tag{18}$$

$$y_i f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \geq 1 - \xi_i \quad \forall 1 \leq i \leq L + 2U \tag{19}$$

$$\xi_i \geq 0 \quad \forall 1 \leq i \leq L + 2U \tag{20}$$

17

Introducing lagrangian variables $\alpha_0$, $\boldsymbol{\alpha}$ and $\boldsymbol{\nu}$ corresponding respectively to constraints (18), (19) and (20), we can write the Lagrangian of this problem as

$$
\begin{aligned}
\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\nu}) =& \frac{1}{2}||\boldsymbol{w}||^2 + C \sum_{i=1}^{L} \xi_i + C^* \sum_{i=L+1}^{L+2U} \xi_i + \sum_{i=L+1}^{L+2U} \beta_i\, y_i\, f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \\
& - \alpha_0 \left( \frac{1}{U} \sum_{i=L+1}^{L+U} f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - \frac{1}{L} \sum_{i=1}^{L} y_i \right) \\
& - \sum_{i=1}^{L+2U} \alpha_i\, (y_i f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - 1 + \xi_i) \\
& - \sum_{i=1}^{L+2U} \nu_i \xi_i \,,
\end{aligned}
\tag{21}
$$

where $\alpha_0$ can be positive or negative (equality constraint) and $\alpha_i$, $i \geq 1$ are non-negative (inequality constraints).

Taking in account that $\beta_i = 0$ for $i \leq L$, calculating the derivatives with respect to the primal variables yields

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}} &= \boldsymbol{w} - \sum_{i=1}^{L+2U} y_i\, (\alpha_i - \beta_i)\, \Phi(\boldsymbol{x}_i) - \frac{\alpha_0}{U} \sum_{i=L+1}^{L+U} \Phi(\boldsymbol{x}_i) \\
\frac{\partial \mathcal{L}}{\partial b} &= -\sum_{i=1}^{L+2U} y_i\, (\alpha_i - \beta_i) - \alpha_0 \\
\frac{\partial \mathcal{L}}{\partial \xi_i} &= C - \alpha_i - \nu_i \quad \forall 1 \leq i \leq L \\
\frac{\partial \mathcal{L}}{\partial \xi_i} &= C^* - \alpha_i - \nu_i \quad \forall L+1 \leq i \leq L+2U \,.
\end{aligned}
$$

For simplifying the notation, we now define an extra special example $\boldsymbol{x}_0$ in an implicit manner:

$$
\Phi(\boldsymbol{x}_0) = \frac{1}{U} \sum_{i=L+1}^{L+U} \Phi(\boldsymbol{x}_i)\,,
$$

and we set $y_0 = 1$ and $\beta_0 = 0$. Setting the derivatives to zero gives us

$$
\boldsymbol{w} = \sum_{i=0}^{L+2U} y_i\, (\alpha_i - \beta_i)\, \Phi(x_i)
\tag{22}
$$

and

$$
\sum_{i=0}^{L+2U} y_i\, (\alpha_i - \beta_i) = 0
\tag{23}
$$

and also

$$
C - \alpha_i - \nu_i = 0 \quad \forall 1 \leq i \leq L, \quad C^* - \alpha_i - \nu_i \quad \forall L+1 \leq i \leq L+2U \,.
\tag{24}
$$

In order to find the minimum of the minimization problem (12) we want to find a saddle point of the Lagrangian (21), as classical SVMs methods. Substituting (22), (23) and (24) into the Lagrangian (21) yields the following maximization problem

$$
\arg\max_{\boldsymbol{\alpha}} \quad -\frac{1}{2} \sum_{i,j=0}^{L+2U} y_i\,y_j(\alpha_i - \beta_i)\,(\alpha_j - \beta_j)\Phi(\boldsymbol{x}_i) \cdot \Phi(\boldsymbol{x}_j)
$$
$$
+ \sum_{i=1}^{L+2U} \alpha_i + \alpha_0 \left( \frac{1}{L} \sum_{i=1}^{L} y_i \right) \tag{25}
$$

under the constraints

$$
\begin{aligned}
0 \le \alpha_i \le C \quad &\forall 1 \le i \le L \\
0 \le \alpha_i \le C^* \quad &\forall L+1 \le i \le L+2U \\
\sum_{i=0}^{L+2U} y_i\,(\alpha_i - \beta_i) = 0\,. &
\end{aligned} \tag{26}
$$

The parameter $\boldsymbol{w}$ is then given by (22) and $b$ is obtained using one of the following Karush-Kuhn-Tucker (KKT) constraints:

$$
\alpha_0 \ne 0 \implies \frac{1}{U} \sum_{i=L+1}^{L+U} [\boldsymbol{w} \cdot \Phi(\boldsymbol{x}_i) + b] = \frac{1}{L} \sum_{i=1}^{L} y_i
$$

$$
\forall 1 \le i \le L,\ 0 < \alpha_i < C \implies y_i[\boldsymbol{w} \cdot \Phi(\boldsymbol{x}_i) + b] = 1
$$
$$
\forall L+1 \le i \le L+2U,\ 0 < \alpha_i < C^* \implies y_i[\boldsymbol{w} \cdot \Phi(\boldsymbol{x}_i) + b] = 1
$$

If we define $\zeta_i = y_i$ for $1 \le i \le L+2U$ and $\zeta_0 = \frac{1}{L}\sum_{i=1}^{L} y_i$, if we consider the kernel matrix $K$ such that

$$
K_{ij} = \Phi(\boldsymbol{x}_i) \cdot \Phi(\boldsymbol{x}_j)\,,
$$

and if we perform the substitution

$$
\tilde{\alpha}_i = y_i\,(\alpha_i - \beta_i)\,,
$$

then we can rewrite the maximization problem (25) under the constraints (26) as the following

$$
\arg\max_{\tilde{\boldsymbol{\alpha}}} \quad \boldsymbol{\zeta} \cdot \tilde{\boldsymbol{\alpha}} - \frac{1}{2} \tilde{\boldsymbol{\alpha}}^{\mathrm{T}} K \tilde{\boldsymbol{\alpha}}
$$

under the constraints

$$
\begin{aligned}
0 \le y_i\,\tilde{\alpha}_i \le C \qquad\qquad &\forall 1 \le i \le L \\
-\beta_i \le y_i\,\tilde{\alpha}_i \le C^* - \beta_i \qquad\qquad &\forall L+1 \le i \le L+2U \\
\sum_{i=0}^{L+2U} \tilde{\alpha}_i = 0\,. &
\end{aligned} \tag{27}
$$

Obviously this optimization problem is very close to a SVM optimization problem. It is thus possible to optimize it with a standard optimizer for SVMs. Note that only the bounds in (27) on the $\tilde{\alpha}_i$ have to be adjusted after each update of $\boldsymbol{\beta}$.

19

# References

M. Belkin and P. Niyogi. Using manifold structure for partially labelled classification. In *Advances in Neural Information Processing Systems*. MIT Press, 2002.

K. Bennett and A. Demiriz. Semi-supervised support vector machines. In M. S. Kearns, S. A. Solla, and D. A. Cohn, editors, *Advances in Neural Information Processing Systems 12*, pages 368–374. MIT Press, Cambridge, MA, 1998.

T. De Bie and N. Cristianini. Convex methods for transduction. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems 16*. MIT Press, Cambridge, MA, 2004.

Antoine Bordes, Seyda Ertekin, Jason Weston, and Léon Bottou. Fast kernel classifiers with online and active learning. *Journal of Machine Learning Research*, May 2005. http://jmlr.csail.mit.edu/papers/v6/bordes05a.html.

B. E. Boser, I. M. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152, Pittsburgh, PA, 1992. ACM Press.

O. Chapelle, J. Weston, and B. Schölkopf. Cluster kernels for semi-supervised learning. *Neural Information Processing Systems 15*, 2002.

O. Chapelle and A. Zien. Semi-supervised classification by low density separation. In *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*, 2005.

R. Collobert, J. Weston, and L. Bottou. Working document: Trading convexity for scalability. June 2005. URL `http://www.kyb.tuebingen.mpg.de/bs/people/fabee/transduction.html`.

G. Fung and O. Mangasarian. Semi-supervised support vector machines for unlabeled data classification. In *Optimisation Methods and Software*, pages 1–14. Kluwer Academic Publishers, Boston, 2001.

T. Graepel, R. Herbrich, and K. Obermayer. Bayesian transduction. In *Advances in Neural Information Processing Systems 12, NIPS*, pages 456–462, 2000.

T. Joachims. Making large-scale support vector machine learning practical. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods*. The MIT Press, 1999a.

T. Joachims. Transductive inference for text classification using support vector machines. In *International Conference on Machine Learning, ICML*, 1999b.

S. Keerthi and D. DeCoste. A modified finite newton method for fast solution of large scale linear svms. *Journal of Machine Learning Research*, 6:341–361, 2005.

N. Krause and Y. Singer. Leveraging the margin more carefully. In *International Conference on Machine Learning, ICML*, 2004.

N. D. Lawrence and M. I. Jordan. Semi-supervised learning via gaussian processes. In *Advances in Neural Information Processing Systems, NIPS*. MIT Press, 2005.

Y. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller. Efficient backprop. In G.B. Orr and K.-R. Müller, editors, *Neural Networks: Tricks of the Trade*, pages 9–50. Springer, 1998.

D. D. Lewis, Y. Yang, T. Rose, and F. Li. Rcv1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research*, 5:361–397, 2004. URL http://www.jmlr.org/papers/volume5/lewis04a/lewis04a.pdf.

L. Mason, P. L. Bartlett, and J. Baxter. Improved generalization through explicit optimization of margins. *Machine Learning*, 38(3):243–255, 2000.

J. C. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods*. The MIT Press, 1999.

S.K.Nayar S.A.Nene and H.Murase. Columbia object image libary (coil-20). Technical Report CUS-005-96, Columbia Univ. USA, Febuary 1996.

B. Schölkopf, A. Smola, and K.-R. Müller. Kernel principal component analysis. In *Proceedings ICANN97*, Springer Lecture Notes in Computer Science, page 583, 1997.

X. Shen, G. C. Tseng, X. Zhang, and W. H. Wong. On (psi)-learning. *Journal of the American Statistical Association*, 98(463):724–734, 2003.

V. Sindhwani, P. Niyogi, and M. Belkin. Beyond the point cloud: from transductive to semi-supervised learning. In *International Conference on Machine Learning, ICML*, 2005.

A. J. Smola, S. V. N. Vishwanathan, and T. Hofmann. Kernel methods for missing variables. In *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*, 2005.

M. Szummer and T. Jaakkola. Partially labeled classification with markov random walks. *NIPS*, 14, 2001a.

M. Szummer and T. Jaakkola. Partially labeled classification with Markov random walks. *Neural Information Processing Systems 14*, 2001b.

H. A. Le Thi. *Analyse numérique des algorithmes de l'optimisation D.C. Approches locales et globale. Codes et simulations numériques en grande dimension. Applications.* PhD thesis, INSA, Rouen, 1994.

V. Vapnik. *The Nature of Statistical Learning Theory.* Springer, second edition, 1995.

J. Weston, C. Leslie, D. Zhou, A. Elisseeff, and W. S. Noble. Cluster kernels for semi-supervised protein classification. *Advances in Neural Information Processing Systems 17*, 2003.

L. Xu, J. Neufeld, B. Larson, and D. Schuurmans. Maximum margin clustering. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 1537–1544. MIT Press, Cambridge, MA, 2005.

A. L. Yuille and A. Rangarajan. The concave-convex procedure (CCCP). In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, Cambridge, MA, 2002. MIT Press.

D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf. Learning with local and global consistency. In *Advances in Neural Information Processing Systems, NIPS*, 2004.

X. Zhu, Z. Ghahramani, and J. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *The Twentieth International Conference on Machine Learning*, pages 912–919, 2003.