
Learning Data Manifolds with a Cutting Plane Method

SueYeon Chung
Harvard University
schung@fas.harvard.edu

Uri Cohen
Hebrew University of Jerusalem
uri.cohen@alice.nc.huji.ac.il

Haim Sompolinsky
Harvard University
Hebrew University of Jerusalem
haim@fiz.huji.ac.il

Daniel D. Lee
University of Pennsylvania
Facebook AI Research
ddlee@seas.upenn.edu

Abstract

We consider the problem of classifying data manifolds where each manifold represents invariances that are parameterized by continuous degrees of freedom. Conventional data augmentation methods rely upon sampling large numbers of training examples from these manifolds; instead, we propose an iterative algorithm called M_{CP} based upon a cutting-plane approach that efficiently solves a quadratic semi-infinite programming problem to find the maximum margin solution. We provide a proof of convergence as well as a polynomial bound on the number of iterations required for a desired tolerance in the objective function. The efficiency and performance of M_{CP} are demonstrated in high-dimensional simulations and on image manifolds generated from the ImageNet dataset. Our results indicate that M_{CP} is able to rapidly learn good classifiers and shows superior generalization performance compared with conventional maximum margin methods using data augmentation methods.

1 Introduction

Handling object variability is a major challenge for machine learning systems. For example, in visual recognition tasks, changes in pose, lighting, identity or background can result in large variability in the appearance of objects [1]. Techniques to deal with this variability has been the focus of much recent work, especially with convolutional neural networks consisting of many layers. The manifold hypothesis states that natural data variability can be modeled as lower-dimensional manifolds embedded in higher dimensional feature representations [2]. A deep neural network can then be understood as disentangling or flattening the data manifolds so that they can be more easily read out in the final layer [3]. Manifold representations of stimuli have also been utilized in neuroscience, where different brain areas are believed to untangle and reformat their representations [4, 5, 6, 7, 8].

This paper addresses the problem of classifying data manifolds that contain invariances with a number of continuous degrees of freedom. These invariances may be modeled using prior knowledge, manifold learning algorithms [9, 10, 11, 12, 13, 14] or as generative neural networks via adversarial training [15]. Based upon knowledge of these structures, other work has considered building group-theoretic invariant representations [16] or constructing invariant metrics [17]. On the other hand, most approaches today rely upon data augmentation by explicitly generating “virtual” examples from these manifolds [18][19]. Unfortunately, the number of samples needed to successfully learn the underlying manifolds may increase the original training set by more than a thousand-fold [20].

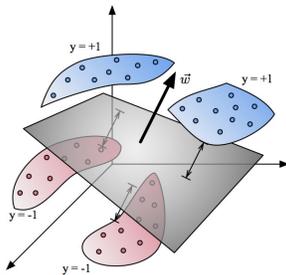


Figure 1: The maximum margin binary classification problem for a set of manifolds. The optimal linear hyperplane is parameterized by the weight vector \vec{w} which separates positively labeled manifolds from negatively labeled manifolds. Conventional data augmentation techniques resort to sampling a large number of points from each manifold to train a classifier.

We propose a new method, called the Manifold Cutting Plane algorithm or M_{CP} , that uses knowledge of the manifolds to efficiently learn a maximum margin classifier. Figure 1 illustrates the problem in its simplest form, binary classification of manifolds with a linear hyperplane with extensions to this basic model discussed later. Given a number of manifolds embedded in a feature space, the M_{CP} algorithm learns a weight vector \vec{w} that separates positively labeled manifolds from negatively labeled manifolds with the maximum margin. Although the manifolds consist of uncountable sets of points, the M_{CP} algorithm is able to find a good solution in a provably finite number of iterations and training examples.

Support vector machines (SVM) can learn a maximum margin classifier given a finite set of training examples [21]; however, with conventional data augmentation methods, the number of training examples increase exponentially rendering the standard SVM algorithm intractable. Methods such as shrinkage and chunking to reduce the complexity of SVM have been studied before in the context of dealing with large-scale datasets [22], but the resultant kernel matrix may still be very large. Other methods which subsample the kernel matrix [23] or reduce the number of training samples [24],[22] may result in suboptimal solutions that do not generalize well.

Our M_{CP} algorithm directly handles the uncountable set of points in the manifolds by solving a quadratic semi-infinite programming problem (QSIP). M_{CP} is based upon a cutting-plane method which iteratively refines a finite set of training examples to solve the underlying QSIP [25, 26, 27]. The cutting-plane method was also previously shown to efficiently handle learning problems with a finite number of examples but an exponentially large number of constraints [28]. We provide a novel analysis of the convergence of M_{CP} with both hard and soft margins. When the problem is realizable, the convergence bound explicitly depends upon the margin value whereas with a soft margin and slack variables, the bound depends linearly on the number of manifolds.

The paper is organized as follows. We first consider the hard margin problem and analyze the simplest form of the M_{CP} algorithm. Next, we introduce slack variables in M_{CP} , one for each manifold, and analyze its convergence with additional auxiliary variables. We then demonstrate the application of M_{CP} to both high-dimensional synthetic data manifolds and to feature representations of images undergoing a variety of warpings. We compare its performance, both in efficiency and generalization error, with conventional SVMs using data augmentation techniques. Finally, we discuss some natural extensions and potential future work on M_{CP} and its applications.

2 Manifolds Cutting Plane Algorithm with Hard Margin

In this section, we first consider the problem of classifying a set of manifolds when they are linearly separable. This allows us to introduce the simplest version of the M_{CP} algorithm along with the appropriate definitions and QSIP formulation. We analyze the convergence of the simple algorithm and prove an upper bound on the number of errors the algorithm can make in this setting.

2.1 Hard Margin QSIP

Formally, we are given a set of P manifolds $M_p \subset \mathbb{R}^N$, $p = 1, \dots, P$ with binary labels $y_p = \pm 1$ (all points in the same manifold share the same label). Each manifold M_p is defined by $\vec{x} = M_p(\vec{s})$ where $\vec{s} \in S_p$, S_p is a compact, convex subset of \mathbb{R}^D representing the parameterization of the invariances and $M_p(\vec{s}) : \mathbb{R}^D \rightarrow \mathbb{R}^N$ is a continuous function of $\vec{s} \in S_p$ so that the manifolds are bounded: $\forall \vec{s}, \|M_p(\vec{s})\| < L$ by some L . We would like to solve the following semi-infinite quadratic programming problem for the weight vector $\vec{w} \in \mathbb{R}^N$:

$$SVM^{simple} : \underset{\vec{w}}{\operatorname{argmin}} \frac{1}{2} \|\vec{w}\|^2 \text{ s.t. } \forall p, \forall \vec{x} \in M_p : y_p \langle \vec{w}, \vec{x} \rangle \geq 1 \quad (1)$$

This is the primal formulation of the problem, where maximizing the margin $\kappa = \frac{1}{\|\vec{w}\|}$ is equivalent to minimizing the squared norm $\frac{1}{2} \|\vec{w}\|^2$. We denote the maximum margin attainable by κ^* , and the optimal solution as \vec{w}^* . For simplicity, we do not explicitly include the bias term here. A non-zero bias can be modeled by adding an additional feature of constant value as a component to all the \vec{x} . Note that the dual formulation of this QSIP is more complicated, involving optimization of non-negative measures over the manifolds. In order to solve the hard margin QSIP, we propose the following simple M_{CP} algorithm.

2.2 M_{CP}^{simple} Algorithm

The M_{CP}^{simple} algorithm is an iterative algorithm to find the optimal \vec{w} in (1), based upon a cutting plane method for solving the QSIP. The general idea behind M_{CP}^{simple} is to start with a finite number of training examples, find the maximum margin solution for that training set, augment the training set by finding a point on the manifolds that violates the constraints, and iterating this process until a tolerance criterion is reached.

At each stage k of the algorithm there is a finite set of training points and associated labels. The training set at the k -th iteration is denoted by the set: $T_k = \{(\vec{x}^i \in M_{p_i}, y_{p_i})\}$ with $i = 1, \dots, |T_k|$ examples. For the i -th pattern in T_k , p_i is the index of the manifold, and y_{p_i} is its associated label.

On this set of examples, we solve the following finite quadratic programming problem:

$$SVM_{T_k} : \underset{\vec{w}}{\operatorname{argmin}} \frac{1}{2} \|\vec{w}\|^2 \text{ s.t. } \forall \vec{x}^i \in T_k : y_{p_i} \langle \vec{w}, \vec{x}^i \rangle \geq 1 \quad (2)$$

to obtain the optimal weights $\vec{w}^{(k)}$ on the training set T_k . We then find a constraint-violating point $\vec{x}_{k+1} \in M_{p_{k+1}}$ from one of the manifolds such that

$$\{p_{k+1}, \vec{x}_{k+1} \in M_{p_{k+1}}\} : y_{p_{k+1}} \langle \vec{w}^{(k)}, \vec{x}_{k+1} \rangle < 1 - \delta \quad (3)$$

with a required tolerance $\delta > 0$. If there is no such point, the M_{CP}^{simple} algorithm terminates. If such a point exists, it is added to the training set, defining the new set $T_{k+1} = T_k \cup \{(\vec{x}_{k+1}, y_{p_{k+1}})\}$. The algorithm then proceeds at the next iteration to solve $SVM_{T_{k+1}}$ to obtain $\vec{w}^{(k+1)}$. For $k = 1$, the set T_1 is initialized with at least one point from each manifold. The pseudocode for M_{CP}^{simple} is shown in Alg. 1.

In step 3 of the M_{CP}^{simple} algorithm, a point among the manifolds that violates the margin constraint needs to be found. The use of a ‘‘separation oracle’’ is common in other cutting plane algorithms such as those used for structural SVM’s [28] or linear mixed-integer programming [29]. In our case, this requires determining the feasibility of $y_p \langle \vec{w}^{(k)}, M_p(\vec{s}) \rangle < 1 - \delta$ over the D -dimensional convex parameter set $\vec{s} \in S_p$. When the manifold mapping is convex, feasibility can be determined sometimes analytically or more generally by a variety of modern convex optimization techniques. For non-convex mappings, a feasible separating point can be found with search techniques in the

Algorithm 1 Pseudocode for the M_{CP}^{simple} algorithm.

Input: δ (tolerance), P manifolds and labels $\{M_p, y_p = \pm 1\}, p = 1, \dots, P$.

1. Initialize $k = 1$, and the set $T_1 = \{(\vec{x}^i \in M_{p_i}, y_{p_i})\}$ with at least one sample from each manifold M_p .
 2. Solve for $\vec{w}^{(k)}$ in $SVMT_{T_k}$: $\min \frac{1}{2} \|\vec{w}\|^2$ s.t. $y^{p_i} \langle \vec{w}, \vec{x}^i \rangle \geq 1$ for all $(\vec{x}^i, y^{p_i}) \in T_k$.
 3. Find a point $\vec{x}_{k+1} \in M_{p_{k+1}}$ among the manifolds $\{p_{k+1} = 1, \dots, P\}$ with a margin s.t. $y_{p_{k+1}} \langle \vec{w}^{(k)}, \vec{x}_{k+1} \rangle < 1 - \delta$.
 4. If there is no such point, then stop. Else, augment the point set: $T_{k+1} = T_k \cup \{(\vec{x}_{k+1}, y_{p_{k+1}})\}$.
 5. $k \leftarrow k + 1$ and go to 2.
-

D -dimensional parameter set using gradient information or finite differences or approximately via convex relaxation techniques. In our experiments below, we provide some specific examples of how separating points can be efficiently found.

2.3 Convergence of M_{CP}^{simple}

The M_{CP}^{simple} algorithm will converge asymptotically to an optimal solution when it exists. Here we show that the $\vec{w}^{(k)}$ asymptotically converges to an optimal \vec{w}^* . Denote the change in the weight vector in the k -th iteration as $\Delta \vec{w}^{(k)} = \vec{w}^{(k+1)} - \vec{w}^{(k)}$. We present a set of lemmas and theorems leading up to the bounds on the number of iterations for convergence, and the estimation of the objective function. More detailed proofs can be found in Supplementary Materials (SM).

Lemma 1. *The change in the weights satisfies $\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle \geq 0$.*

Proof. Define $\vec{w}(\lambda) = \vec{w}^{(k)} + \lambda \Delta \vec{w}^{(k)}$. Then for all $0 \leq \lambda \leq 1$, $\vec{w}(\lambda)$ satisfies the constraints on the point set T_k : $y_{p_i} \langle \vec{w}(\lambda), \vec{x}_i \rangle \geq 1$ for all $(\vec{x}_i, y_{p_i}) \in T_k$. However, if $\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle < 0$, there exists a $0 < \lambda' < 1$ such that $\|\vec{w}(\lambda')\|^2 < \|\vec{w}^{(k)}\|^2$, contradicting the fact that $\vec{w}^{(k)}$ is a solution to $SVMT_{T_k}$. \square

Next, we show that the norm $\|\vec{w}^{(k)}\|^2$ must monotonically increase by a finite amount at each iteration.

Theorem 2. *In the k th iteration of M_{CP}^{simple} algorithm, the increase in the norm of $\vec{w}^{(k)}$ is lower bounded by $\|\vec{w}^{(k+1)}\|^2 \geq \|\vec{w}^{(k)}\|^2 + \frac{\delta_k^2}{L^2}$, where $\delta_k = 1 - y_{p_{k+1}} \langle \vec{w}^{(k)}, \vec{x}_{k+1} \rangle$ and $\|\vec{x}_{k+1}\| \leq L$.*

Proof. First, note that $\delta_k > \delta \geq 0$, otherwise the algorithm stops. We have: $\|\vec{w}^{(k+1)}\|^2 = \|\vec{w}^{(k)}\|^2 + \|\Delta \vec{w}^{(k)}\|^2 + 2 \langle \vec{w}^{(k)}, \Delta \vec{w}^{(k)} \rangle \geq \|\vec{w}^{(k)}\|^2 + \|\Delta \vec{w}^{(k)}\|^2$ (Lemma 1). Consider the point added to set $T_{k+1} = T_k \cup (\vec{x}_{k+1}, y_{p_{k+1}})$. At this point, $y_{p_{k+1}} \langle \vec{w}^{(k+1)}, \vec{x}_{k+1} \rangle \geq 1$, $y_{p_{k+1}} \langle \vec{w}^{(k)}, \vec{x}_{k+1} \rangle = 1 - \delta_k$, hence $y_{p_{k+1}} \langle \Delta \vec{w}^{(k)}, \vec{x}_{k+1} \rangle \geq \delta_k$. Then, from the Cauchy-Schwartz inequality,

$$\|\Delta \vec{w}^{(k)}\|^2 \geq \frac{\delta_k^2}{\|\vec{x}_{k+1}\|^2} > \frac{\delta_k^2}{L^2} > \frac{\delta^2}{L^2} \quad (4)$$

Since the solution \vec{w}^* satisfies the constraints for T_k , $\|\vec{w}^{(k)}\| \leq \frac{1}{\kappa^*}$. Thus, the sequence of iterations monotonically increase norms and are upper bounded by $\frac{1}{\kappa^*}$. Due to convexity, there is a single global optimum and the M_{simple}^4 algorithm is guaranteed to converge. \square

As a corollary, we see that this procedure is guaranteed to find a realizable solution if it exists in a finite number of steps.

Corollary 3. *The M_{CP}^{simple} algorithm converges to a zero error classifier in less than $\frac{L^2}{(\kappa^*)^2}$ iterations, where κ^* is the optimal margin and L bounds the norm of the points on the manifolds.*

Proof. When there is an error, we have $\delta_k > 1$, and $\|\vec{w}^{(k+1)}\|^2 \geq \|\vec{w}^{(k)}\|^2 + \frac{1}{L^2}$ (See (4)). This implies the total number of possible errors is upper bounded by $\frac{L^2}{(\kappa^*)^2}$. \square

With a finite tolerance $\delta > 0$, we obtain a bound on the number of iterations required for convergence:

Corollary 4. *The M_{CP}^{simple} algorithm for a given tolerance $\delta > 0$ terminates in less than $\frac{L^2}{(\kappa^*\delta)^2}$ iterations where κ^* is the optimal margin and L bounds the norm of the points on the manifolds.*

Proof. Again, $\|\vec{w}^k\|^2 \leq \|\vec{w}^*\|^2 = \frac{1}{(\kappa^*)^2}$ and each iteration increases the squared norm by at least $\frac{\delta^2}{L^2}$. \square

We can also bracket the error in the objective function after M_{CP}^{simple} terminates:

Corollary 5. *With tolerance $\delta > 0$, after M_{CP}^{simple} terminates with solution $\vec{w}_{M_{CP}}$, the optimal value $\|\vec{w}^*\|$ of $SV M^{simple}$ is bracketed by: $\|\vec{w}_{M_{CP}}\|^2 \leq \|\vec{w}^*\|^2 \leq \frac{1}{(1-\delta)^2} \|\vec{w}_{M_{CP}}\|^2$.*

Proof. The lower bound on $\|\vec{w}^*\|^2$ is as before. Since M_{CP}^{simple} has terminated, setting $\vec{w}' = \frac{1}{(1-\delta)}\vec{w}_{M_{CP}}$ would make \vec{w}' feasible for $SV M^{simple}$, resulting in the upper bound on $\|\vec{w}^*\|^2$. \square

3 M_{CP} with Slack Variables

In many classification problems, the manifolds may not be linearly separable due to their dimensionality, size, or correlation structure. In these situations, M_{CP}^{simple} will not be able to find a feasible solution. To handle these problems, the classic approach is to introduce slack variables on each point ($SV M^{slack}$). Unfortunately, this approach requires integrations over entire manifolds with an appropriate measure defined by the infinite set of slack variables. Thus, we formulate an alternative version of the QSIP with slack variables below.

3.1 QSIP with Manifold Slacks

In this work, we propose using only one slack variable per manifold for classification problems with non-separable manifolds. This formulation demands that all the points on each manifold $\vec{x} \in M_p$ obey an inequality constraint with one manifold slack variable, $y_p \langle \vec{w}, \vec{x} \rangle + \xi_p \geq 1$. As we see below, solving for this constraint is tractable and the algorithm has good convergence guarantees.

However, a single slack requirement for each manifold by itself may not be sufficient for good generalization performance. Our empirical studies show that generalization performance can be improved if we additionally demand that some representative points $\vec{x}_p \in M_p$ on each manifold also obey the margin constraint: $y_p \langle \vec{w}, \vec{x}_p \rangle \geq 1$. In this work, we implement this intuition by specifying appropriate center points \vec{x}_p^c for each manifold M_p . This center point could be the center of mass of the manifold, a representative point, or an exemplar used to generate the manifolds [20]. Additional slack variables for these constraints could potentially be introduced; in the present work, we simply demand that these points strictly obey the margin inequalities corresponding to their manifold label. Formally, the QSIP optimization problem is summarized below, where the objective function is minimized over the weight vector $\vec{w} \in \mathbb{R}^N$ and slack variables $\vec{\xi} \in \mathbb{R}^P$:

$$SV M_{manifold}^{slack} : \underset{\vec{w}, \vec{\xi}}{\operatorname{argmin}} F(\vec{w}, \vec{\xi}) = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{p=1}^P \xi_p$$

$$s.t. \forall p, \forall \vec{x} \in M_p : y_p \langle \vec{w}, \vec{x} \rangle + \xi_p \geq 1 (\text{manifolds}), \forall p : y_p \langle \vec{w}, \vec{x}_p^c \rangle \geq 1 (\text{centers}), \xi_p \geq 0$$

3.2 M_{CP}^{slack} Algorithm

With these definitions, we introduce our M_{slack}^4 algorithm with slack variables below.

Algorithm 2 Pseudocode for the M_{CP}^{slack} algorithm.

Input: δ (tolerance), P manifolds and labels $\{M_p, y_p = \pm 1\}$, and centers \vec{x}_p^c

1. Initialize $k = 1$, and the set $T_1 = \{(\vec{x}^i \in M_{p_i}, y_{p_i})\}$ with at least one sample from each manifold M_p .
 2. Solve for $\vec{w}^{(k)}, \vec{\xi}^{(k)}$: $\min \frac{1}{2} \|\vec{w}\|^2 + C \sum_{p=1}^P \xi_p$ s.t. $y_{p_\mu} \langle \vec{w}, \vec{x}^\mu \rangle + \xi_{p_\mu} \geq 1$ for all $(\vec{x}^\mu, y_{p_\mu}) \in T_k$ and $y_p \langle \vec{w}, \vec{x}_p^c \rangle \geq 1$ for all p .
 3. Find a point $\vec{x}^{k+1} \in M_{p_{k+1}}$ among the manifolds $\{p = 1, \dots, P\}$ with slack violation larger than the tolerance δ : $y_{p_{k+1}} \langle \vec{w}^{(k)}, \vec{x}_{k+1} \rangle + \xi_{p_{k+1}} < 1 - \delta$
 4. If there is no such point, then stop. Else, augment the point set: $T_{k+1} = T_k \cup \{(\vec{x}_{k+1}, y_{p_{k+1}})\}$.
 5. $k \leftarrow k + 1$ and go to 2.
-

The proposed M_{CP}^{slack} algorithm modifies the cutting plane approach to solve a semi-infinite, semi-definite quadratic program. Each iteration involves a finite set: $T_k = \{(\vec{x}^i \in M_{p_i}, y_{p_i})\}$ with $i = 1, \dots, |T_k|$ examples that is used to define the following soft margin SVM:

$$\begin{aligned}
& SVM_{T_k}^{slack} : \underset{\vec{w}, \vec{\xi}}{\operatorname{argmin}} \frac{1}{2} \|\vec{w}\|^2 + C \sum_{p=1}^P \xi_p \\
& \text{s.t. } \forall (\vec{x}^i, y_{p_i}) \in T_k : y_{p_i} \langle \vec{w}, \vec{x}^i \rangle + \xi_{p_i} \geq 1; \forall p : y_p \langle \vec{w}, \vec{x}_p^c \rangle \geq 1 \text{ (centers)}, \xi_p \geq 0
\end{aligned}$$

to obtain the weights $\vec{w}^{(k)}$ and slacks $\vec{\xi}^{(k)}$ at each iteration. We then find a point $\vec{x}_{k+1} \in M_{p_{k+1}}$ from one of the manifolds so that:

$$y_{p_{k+1}} \langle \vec{w}^{(k)}, \vec{x}_{k+1} \rangle + \xi_{p_{k+1}}^{(k)} = 1 - \delta_k \quad (5)$$

where $\delta_k > \delta$. If there is no such a point, the M_{CP}^{slack} algorithm terminates. Otherwise, the point \vec{x}_{k+1} is added as a training example to the set $T_{k+1} = T_k \cup \{(\vec{x}_{k+1}, y_{p_{k+1}})\}$ and the algorithm proceeds to solve $SVM_{T_{k+1}}^{slack}$ to obtain $\vec{w}^{(k+1)}$ and $\vec{\xi}^{(k+1)}$.

3.3 Convergence of M_{CP}^{slack}

Here we show that the objective function $F(\vec{w}, \vec{\xi}) = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{p=1}^P \xi_p$ is guaranteed to increase by a finite amount with each iteration. This result is similar to [30], but here we present statements in the primal domain over an infinite number of examples. More detailed proofs can be found in SM.

Lemma 6. *The change in the weights and slacks satisfy:*

$$\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_p \Delta \xi_p^{(k)} \geq 0 \quad (6)$$

where $\Delta \vec{w}^{(k)} = \vec{w}^{(k+1)} - \vec{w}^{(k)}$ and $\Delta \xi^{(k)} = \xi^{(k+1)} - \xi^{(k)}$.

Proof. Define $\vec{w}(\lambda) = \vec{w}^{(k)} + \lambda \Delta \vec{w}^{(k)}$ and $\vec{\xi}(\lambda) = \vec{\xi}^{(k)} + \lambda \Delta \xi^{(k)}$. Then for all $0 \leq \lambda \leq 1$, $\vec{w}(\lambda)$ and $\vec{\xi}(\lambda)$ satisfy the constraints for $SVM_{T_k}^{slack}$. The resulting change in the objective function is given by:

$$\begin{aligned}
& F(\vec{w}(\lambda), \vec{\xi}(\lambda)) - F(\vec{w}^{(k)}, \vec{\xi}^{(k)}) = \\
& \lambda \left[\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_p \Delta \xi_p^{(k)} \right] + \frac{1}{2} \lambda^2 \|\Delta \vec{w}^{(k)}\|^2 \quad (7)
\end{aligned}$$

If (6) is not satisfied, then there is some $0 < \lambda' < 1$ such that $F(\vec{w}(\lambda'), \vec{\xi}(\lambda')) < F(\vec{w}^{(k)}, \vec{\xi}^{(k)})$, which contradicts the fact that $\vec{w}^{(k)}$ and $\vec{\xi}^{(k)}$ are a solution to SVM_{T_k} . \square

We derive that the added point at each iteration must be a support vector for the next weight:

Lemma 7. *In each iteration of M_{CP}^{slack} algorithm, the added point $(\vec{x}_{k+1}, y_{p_{k+1}})$ must be a support vector for the new weights and slacks, s.t.*

$$y_{p_{k+1}} \langle \vec{w}^{(k+1)}, \vec{x}_{k+1} \rangle + \xi_{p_{k+1}}^{(k+1)} = 1 \quad (8)$$

$$y_{p_{k+1}} \langle \Delta \vec{w}^{(k)}, \vec{x}_{k+1} \rangle + \Delta \xi_{p_{k+1}}^{(k)} = \delta_k \quad (9)$$

Proof. Suppose $y_{p_{k+1}} \langle \vec{w}^{(k+1)}, \vec{x}_{k+1} \rangle + \xi_{p_{k+1}}^{(k+1)} = 1 + \epsilon$ for some $\epsilon > 0$. Then we can choose $\lambda' = \frac{\delta_k}{\delta_k + \epsilon} < 1$ so that $\vec{w}(\lambda') = \vec{w}^{(k)} + \lambda' \Delta \vec{w}^{(k)}$ and $\vec{\xi}(\lambda') = \vec{\xi}^{(k)} + \lambda' \Delta \vec{\xi}^{(k)}$ satisfy the constraints for $SV M_{T_{k+1}}^{slack}$. But, from Lemma 6, we have $F(\vec{w}(\lambda'), \vec{\xi}(\lambda')) < F(\vec{w}^{(k+1)}, \vec{\xi}^{(k+1)})$ which contradicts the fact that $\vec{w}^{(k+1)}$ and $\vec{\xi}^{(k+1)}$ are a solution to $SV M_{T_{k+1}}$. Thus, $\epsilon = 0$ and the point $(\vec{x}_{k+1}, y_{p_{k+1}})$ must be a support vector for $SV M_{T_{k+1}}$. (9) results from subtracting (5) from (8). \square

We also derive a bound on the following quadratic function over nonnegative values:

Lemma 8. *Given $K > 0, \delta > 0$, then $\forall x \geq 0$*

$$\frac{1}{2}(x - \delta)^2 + Kx \geq \min\left(\frac{1}{2}\delta^2, \frac{1}{2}K\delta\right) \quad (10)$$

Proof. The minimum value occurs when $x^* = [\delta - K]_+$. When $K \geq \delta$, then $x^* = 0$ and the minimum is $\frac{1}{2}\delta^2$. When $K < \delta$, the minimum occurs at $K(\delta - \frac{1}{2}K) \geq \frac{1}{2}K\delta$. Thus, the lower bound is the smaller of these two values. \square

Using the lemmas above, the lower bound on the change in the objective function can be found:

Theorem 9. *In each iteration k of M_{CP}^{slack} algorithm, the increase in the objective function for $SV M_{manifold}^{slack}$, defined as $\Delta F^{(k)} = F(\vec{w}^{(k+1)}, \vec{\xi}^{(k+1)}) - F(\vec{w}^{(k)}, \vec{\xi}^{(k)})$, is lower bounded by*

$$\Delta F^{(k)} \geq \min\left(\frac{1}{8} \frac{\delta_k^2}{L^2}, \frac{1}{2} C \delta_k\right) \quad (11)$$

Proof. We first note that the change in objective function is strictly increasing:

$$F(\vec{w}^{(k+1)}, \vec{\xi}^{(k+1)}) - F(\vec{w}^{(k)}, \vec{\xi}^{(k)}) = \left[\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_p \Delta \xi_p^{(k)} \right] + \frac{1}{2} \|\Delta \vec{w}^{(k)}\|^2 > 0 \quad (12)$$

This can be seen immediately from Lemma 6 when $\Delta \vec{w}^{(k)} \neq 0$. On the other hand, if $\Delta \vec{w}^{(k)} = 0$, we know that $\Delta \xi_{p_k}^{(k)} = \delta_k$ from Lemma 7 and $\Delta \xi_{p \neq p_k}^{(k)} = 0$ since $\vec{\xi}^{(k)}$ is the solution for $SV M_{T_k}$. So for $\Delta \vec{w}^{(k)} = 0$, $F(\vec{w}^{(k+1)}, \vec{\xi}^{(k+1)}) - F(\vec{w}^{(k)}, \vec{\xi}^{(k)}) = C\delta_k$. To compute a general lower bound on the increase in the objective function, we proceed as follows.

The added point \vec{x}_k comes from a particular manifold M_{p_k} . If $\Delta \xi_{p_k}^{(k)} \leq 0$, from Lemma 7 we have $y_{p_k} \langle \Delta \vec{w}^{(k)}, \vec{x}_k \rangle \geq \delta_k$. Then by the Cauchy-Schwarz inequality, $\|\Delta \vec{w}^{(k)}\|^2 \geq \frac{\delta_k^2}{L^2}$ which yields $F(\vec{w}^{(k+1)}, \vec{\xi}^{(k+1)}) - F(\vec{w}^{(k)}, \vec{\xi}^{(k)}) \geq \frac{1}{2} \frac{\delta_k^2}{L^2}$.

We next analyze $\Delta\xi_{p_k}^{(k)} > 0$ and consider the finite set of points $(\vec{x}^\nu, p_k) \in T_k$ that come from the p_k manifold. There must be at least one such point in T_k by initialization of the algorithm. Each of these points obeys the constraints:

$$y_{p_k} \langle \vec{w}^{(k)}, \vec{x}^\nu \rangle + \xi_{p_k}^{(k)} = 1 + \epsilon_\nu^{(k)} \quad (13)$$

$$y_{p_k} \langle \vec{w}^{(k+1)}, \vec{x}^\nu \rangle + \xi_{p_k}^{(k+1)} = 1 + \epsilon_\nu^{(k+1)} \quad (14)$$

$$\epsilon_\nu^{(k)}, \epsilon_\nu^{(k+1)} \geq 0 \quad (15)$$

We consider the minimum value of the thresholds: $\eta = \min_\nu \epsilon_\nu^{(k)}$. We have two possibilities: η is positive so that none of the points are support vectors for $SV M_{T_k}^{slack}$, or $\eta = 0$ so that at least one support vector lies in M_{p_k} .

Case $\eta > 0$:

In this case, we define a linear set of slack variables:

$$\xi_p(\lambda) = \begin{cases} \xi_p^{(k)} + \lambda \Delta \xi_p^{(k)} & p \neq p_k \\ \xi_{p_k}^{(k)} & p = p_k \end{cases} \quad (16)$$

and weights $\vec{w}(\lambda) = \vec{w}^{(k)} + \lambda \Delta \vec{w}^{(k)}$. Then for $0 \leq \lambda \leq \min\left(\frac{\eta_{p_k}}{\Delta \xi_{p_k}^{(k)}}, 1\right)$, $\vec{w}(\lambda)$ and $\vec{\xi}(\lambda)$ will satisfy the constraints for $SV M_{T_k}$. Following similar reasoning in Lemma 6, this implies

$$\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_{p \neq p_k} \Delta \xi_p^{(k)} \geq 0 \quad (17)$$

Then in this case, we have

$$\begin{aligned} & F(\vec{w}^{(k+1)}, \vec{\xi}^{(k+1)}) - F(\vec{w}^{(k)}, \vec{\xi}^{(k)}) \\ &= \left[\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_p \Delta \xi_p^{(k)} \right] + \frac{1}{2} \|\Delta \vec{w}^{(k)}\|^2 \\ &\geq C \Delta \xi_{p_k}^{(k)} + \frac{1}{2} \|\Delta \vec{w}^{(k)}\|^2 \end{aligned} \quad (18)$$

$$\geq \frac{1}{2} \frac{(\delta_k - \Delta \xi_{p_k}^{(k)})^2}{L^2} + C \Delta \xi_{p_k}^{(k)} \quad (19)$$

$$\geq \min\left(\frac{1}{2L^2} \delta_k^2, \frac{1}{2} C \delta_k\right) \quad (20)$$

by applying (17) in (18), Lemma 7 and Cauchy-Schwarz in (19) and Lemma 8 in (20).

Case $\eta = 0$:

In this case, we consider $\varepsilon = \min_{\epsilon_\nu^{(k)}=0} \epsilon_\nu^{(k+1)} \geq 0$, i.e. the minimal increase among the support vectors. We then define

$$\xi_p(\lambda) = \begin{cases} \xi_p^{(k)} + \lambda \Delta \xi_p^{(k)} & p \neq p_k \\ \xi_{p_k}^{(k)} + \lambda (\Delta \xi_{p_k}^{(k)} - \varepsilon) & p = p_k \end{cases} \quad (21)$$

and weights $\vec{w}(\lambda) = \vec{w}^{(k)} + \lambda \Delta \vec{w}^{(k)}$. There will then be a finite range of $0 \leq \lambda \leq \lambda_{min}$ for which $\vec{\xi}(\lambda)$ and $\vec{w}(\lambda)$ satisfy the constraints for $SV M_{T_k}$ so that

$$\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_{p \neq p_k} \Delta \xi_p^{(k)} + C (\Delta \xi_{p_k}^{(k)} - \varepsilon) \geq 0 \quad (22)$$

$$\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_p \Delta \xi_p^{(k)} \geq C \varepsilon \quad (23)$$

We also have a support vector $(\vec{x}^\nu, p_k) \in T_k$ so that

$$y_{p_k} \langle \vec{w}^{(k+1)}, \vec{x}^\nu \rangle + \xi_{p_k}^{(k+1)} = 1 + \varepsilon \quad (24)$$

$$y_{p_k} \langle \Delta \vec{w}^{(k)}, \vec{x}^\nu \rangle + \Delta \xi_{p_k}^{(k)} = \varepsilon \quad (25)$$

Using Lemma 7 and Cauchy-Schwarz, we get:

$$y_{p_k} \langle \Delta \vec{w}^{(k)}, \vec{x}_k - \vec{x}^\nu \rangle = \delta_k - \varepsilon \quad (26)$$

$$\|\Delta \vec{w}^{(k)}\|^2 \geq \frac{1}{4L^2} (\delta_k - \varepsilon)^2 \quad (27)$$

Thus, we have:

$$F(\vec{w}^{(k+1)}, \vec{\xi}^{(k+1)}) - F(\vec{w}^{(k)}, \vec{\xi}^{(k)}) \quad (28)$$

$$= \left[\langle \Delta \vec{w}^{(k)}, \vec{w}^{(k)} \rangle + C \sum_p \Delta \xi_p^{(k)} \right] + \frac{1}{2} \|\Delta \vec{w}^{(k)}\|^2 \quad (29)$$

$$\geq C\varepsilon + \frac{1}{8L^2} (\delta_k - \varepsilon)^2 \quad (30)$$

$$\geq \min \left(\frac{1}{8L^2} \delta_k^2, \frac{1}{2} C \delta_k \right) \quad (31)$$

by applying (23) and (27) to obtain the final bound. \square

Since the M_{CP}^{slack} algorithm is guaranteed to increase the objective by a finite amount, it will terminate in a finite number of iterations if we require $\delta_k > \delta$ for some positive $\delta > 0$.

Corollary 10. *The M_{CP}^{slack} algorithm for a given $\delta > 0$ will terminate after at most $P \cdot \max \left(\frac{8CL^2}{\delta^2}, \frac{2}{\delta} \right)$ iterations where P is the number of manifolds, L bounds the norm of the points on the manifolds.*

Proof. $\vec{w} = 0$ and $\xi_p = 1$ is a feasible solution for $SV M_{manifold}^{slack}$. Therefore, the optimal objective function is upper-bounded by $F(\vec{w} = 0, \vec{\xi} = 1) = PC$. The upper bound on the number of iterations is then provided by Theorem (9). \square

We can also bound the error in the objective function after M_{CP}^{slack} terminates:

Corollary 11. *With $\delta > 0$, after M_{CP}^{slack} terminates with solution \vec{w}_{MCP} , slack $\vec{\xi}_{MCP}$ and value $F_{MCP} = F(\vec{w}_{MCP}, \vec{\xi}_{MCP})$, then the optimal value F^* of $SV M_{manifold}^{slack}$ is bracketed by:*

$$F_{MCP} \leq F^* \leq F_{MCP} + PC\delta. \quad (32)$$

Proof. The lower bound on F^* is apparent since $SV M_{manifold}^{slack}$ includes $SV M_{T_k}^{slack}$ constraints for all k . Setting the slacks $\xi_p = \xi_{MCP,p} + \delta$ will make the solution feasible for $SV M_{manifold}^{slack}$ resulting in the upper bound. \square

4 Experiments

4.1 L_q Ellipsoidal Manifolds

As an illustration of our method, we have generated D -dimensional L_q -norm ellipsoids with random radii, centers, and directions. The points on each manifold M_p are parameterized as

$M_p = \left\{ \vec{x} \mid \vec{x} = \vec{x}_p^c + \sum_{i=1}^D R_i s_i \vec{u}_i^p \in \mathbb{R}^N \right\}$ where the center \vec{x}_p^c and basis vectors \vec{u}_i^p are random Gaussian and R_i , the ellipsoidal radii, are sampled from $\text{Unif}[0.5R_0, 1.5R_0]$ with mean R_0 .

We compare the performance of M_{CP} to the conventional Point SVM (SVM^{simple} , SVM^{slack}) with samples drawn from the surface of the L_q ellipsoids for training and test examples. Performance is measured by generalization error on the task of classifying points from positively labeled manifolds from negatively labeled ones, as a function of the number of training samples used during learning.

For these manifolds, the separation oracle of M_{CP} returns points that minimize $y^p \left[\vec{w} \cdot \left(\vec{x}_0^p + \sum_{i=1}^D R_i s_i \vec{u}_i^p \right) \right]$ over the set $\|\vec{s}\|_q \leq 1$. For norms with $q \geq 1$, the solution can be expressed as $s_i = -\frac{(h_i^p)^{1/(q-1)}}{\left\{ \sum (h_i^p)^{q/(q-1)} \right\}^{1/q}}$ where $h_i^p = y^p \vec{w} \cdot \vec{u}_i^p$. For M_{CP}^{slack} , we used an additional single margin constraint per manifold given by the center of the ellipsoids \vec{x}_p^c . For the parameters shown in Fig. 2, we use $N = 500$ to test M_{CP}^{simple} and $N = 490$ for the M_{CP}^{slack} simulations which are below and above the critical capacity. The results illustrate that M_{CP} achieves low generalization error very efficiently compared to a conventional maximum margin classifier using sampled training examples.

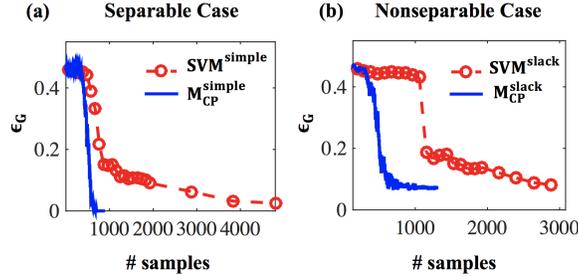


Figure 2: **Generalization error (ϵ_G) of the M_{CP} solution for L_{50} ellipsoids**, shown as a function of the total number of training samples (blue solid) compared with conventional Point SVM (red dashed). $P = 48$, $D = 10$, $R_0 = 20$ (mean elliptic radii) and (a) $P/N = 0.096$ just below the critical capacity is used for M_{CP}^{simple} and (b) $P/N = 0.098$ for M_{CP}^{slack} with slack coefficient $C_{opt} = 100$. ϵ_G is computed from 500 points per manifold.

4.2 ImageNet Object Manifolds

We also apply the M_{CP} algorithm to a more realistic class of object manifolds. Here each object manifold is constructed from a set of 6 – D affine warpings of images from the ImageNet data-set [31]. The resulting P manifolds are split into dichotomies, with $\frac{1}{2}P$ manifolds considered as positive instances and the remaining $\frac{1}{2}P$ manifolds as negative instances. M_{CP} is then used to learn a binary classifier for the manifolds. We note that M_{CP} can easily be extended to multi-class problems but was not used in this experiment.

We compared the performance of M_{CP} to conventional SVM on sampled points from two different feature representations of the images: in the original pixel space and in a V1-like representation of the same images created by applying full-wave rectification after filtering by arrays of Gabor functions [32]. We trained on $P = 4$ object manifolds, and obtained the parameter C through cross-validation. At each iteration of M_{CP} , a constraint-violating point on the manifolds is found using local search among $K = 5$ neighboring samples in the affine warping space. Fig. 3 shows how quickly M_{CP} converges to low generalization error on both the Gabor features and pixel representations. We note that the maximal number P of manifolds that are separable depends upon the feature representation. Thus, M_{CP} can also be used to investigate the benefits of alternative features representations, such as those found in different areas of brain visual processing or in the layers of a deep neural network.

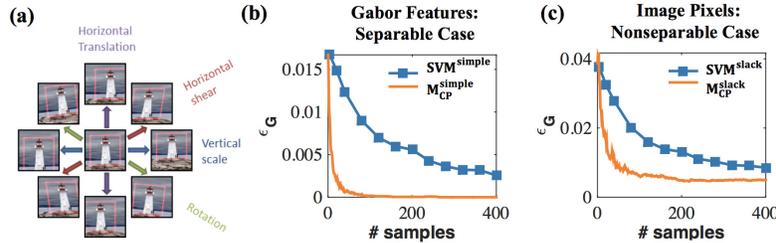


Figure 3: **Image-based object manifolds.** (a) Basic affine transformation: a template image (middle) with an object bounding box (pink) surrounded with changes along 4 axes defined by basic affine transformation. (b-c) Generalization error of the M_{CP} solution for 6-D image-based object manifolds as a function of the number of training samples (solid orange line) compared with that of conventional Point SVM (blue squares), obtained in (b) Gabor representations and (c) Pixel representations. The generalization error is averaged over all possible choices for labels.

5 Discussion

We described and analyzed a novel algorithm, M_{CP} , based upon the cutting-plane method for finding the maximum margin solution for classifying manifolds. We proved the convergence of M_{CP} , and provided bounds on the number of iterations required. Our experiments with both synthetic data manifolds and image manifolds demonstrate the efficiency of M_{CP} and superior performance in terms of generalization error, compared to conventional SVM's using many virtual examples.

There is a natural extension of M_{CP} to nonlinear classifiers via the kernel trick, as all our operations involve dot products between the weight vector \vec{w} and manifold points $M_p(\vec{s})$. At each iteration, the algorithm would solve the dual version of the $SVMT_k$ problem which is readily kernelized. More theoretical work is needed with infinite-dimensional kernels when the manifold optimization problem no longer is a QSIP but becomes a fully (doubly) infinite quadratic programming problem.

Beyond binary classification, variations of M_{CP} can also be used to solve other machine learning problems including multi-class classification, ranking, ordinal regression, one-class learning, etc. We can also use M_{CP} to evaluate the computational benefits of manifold representations at successive layers of deep networks in both machine learning and in brain sensory hierarchies. We anticipate using M_{CP} to help construct novel hierarchical architectures that can incrementally reformat the manifold representations through the layers for better overall performance in machine learning tasks, improving our understanding on how neural architectures can learn to process high dimensional real-world signal ensembles and cope with large variability due to continuous modulation of the underlying physical parameters.

Acknowledgments

The work is partially supported by the Gatsby Charitable Foundation, the Swartz Foundation, the Simons Foundation (SCGB Grant No. 325207), the NIH, and the Human Frontier Science Program (Project RGP0015/2013). D. Lee also acknowledges the support of the US National Science Foundation, Army Research Laboratory, Office of Naval Research, Air Force Office of Scientific Research, and Department of Transportation.

References

- [1] Geoffrey E Hinton, Peter Dayan, and Michael Revow. Modeling the manifolds of images of handwritten digits. *Neural Networks, IEEE Transactions on*, 8(1):65–74, 1997.
- [2] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.

- [3] Pratik Prabhajan Brahma, Dapeng Wu, and Yiyuan She. Why deep learning works: A manifold disentanglement perspective. *IEEE transactions on neural networks and learning systems*, 27(10):1997–2008, 2016.
- [4] Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object recognition in cortex. *Nature neuroscience*, 2(11):1019–1025, 1999.
- [5] Thomas Serre, Lior Wolf, and Tomaso Poggio. Object recognition with features inspired by visual cortex. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 2, pages 994–1000. IEEE, 2005.
- [6] Chou P Hung, Gabriel Kreiman, Tomaso Poggio, and James J DiCarlo. Fast readout of object identity from macaque inferior temporal cortex. *Science*, 310(5749):863–866, 2005.
- [7] James J DiCarlo and David D Cox. Untangling invariant object recognition. *Trends in cognitive sciences*, 11(8):333–341, 2007.
- [8] Marino Pagan, Luke S Urban, Margot P Wohl, and Nicole C Rust. Signals in inferotemporal and perirhinal cortex suggest an untangling of visual target information. *Nature neuroscience*, 16(8):1132–1139, 2013.
- [9] Joshua B Tenenbaum et al. Mapping a manifold of perceptual observations. *Advances in neural information processing systems*, pages 682–688, 1998.
- [10] Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000.
- [11] Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- [12] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- [13] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *The Journal of Machine Learning Research*, 7:2399–2434, 2006.
- [14] Guillermo Canas, Tomaso Poggio, and Lorenzo Rosasco. Learning manifolds with k-means and k-flats. In *Advances in Neural Information Processing Systems*, pages 2465–2473, 2012.
- [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [16] Fabio Anselmi, Joel Z Leibo, Lorenzo Rosasco, Jim Mutch, Andrea Tacchetti, and Tomaso Poggio. Unsupervised learning of invariant representations in hierarchical architectures. *arXiv preprint arXiv:1311.4158*, 2013.
- [17] Patrice Y Simard, Yann Le Cun, and John S Denker. Memory-based character recognition using a transformation invariant metric. In *Pattern Recognition, 1994. Vol. 2-Conference B: Computer Vision & Image Processing., Proceedings of the 12th IAPR International Conference on*, volume 2, pages 262–267. IEEE, 1994.
- [18] Partha Niyogi, Federico Girosi, and Tomaso Poggio. Incorporating prior information in machine learning by creating virtual examples. *Proceedings of the IEEE*, 86(11):2196–2209, 1998.
- [19] Bernhard Schölkopf, Chris Burges, and Vladimir Vapnik. Incorporating invariances in support vector learning machines. In *International Conference on Artificial Neural Networks*, pages 47–52. Springer, 1996.
- [20] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [21] Vladimir Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.
- [22] Alex J Smola and Bernhard Schölkopf. *Learning with kernels*. Citeseer, 1998.
- [23] Yuh-Jye Lee and Olvi L Mangasarian. Rsvm: Reduced support vector machines. In *Proceedings of the 2001 SIAM International Conference on Data Mining*, pages 1–17. SIAM, 2001.

- [24] Jigang Wang, Predrag Neskovic, and Leon N Cooper. Training data selection for support vector machines. In *International Conference on Natural Computation*, pages 554–564. Springer, 2005.
- [25] Shu-Cherng Fang, Chih-Jen Lin, and Soon-Yi Wu. Solving quadratic semi-infinite programming problems by using relaxed cutting-plane scheme. *Journal of computational and applied mathematics*, 129(1):89–104, 2001.
- [26] Kenneth O Kortanek and Hoon No. A central cutting plane algorithm for convex semi-infinite programming problems. *SIAM Journal on Optimization*, 3(4):901–918, 1993.
- [27] Y Liu, Kok Lay Teo, and Soon-Yi Wu. A new quadratic semi-infinite programming algorithm based on dual parametrization. *Journal of Global Optimization*, 29(4):401–413, 2004.
- [28] Thorsten Joachims. Training linear svms in linear time. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 217–226. ACM, 2006.
- [29] Hugues Marchand, Alexander Martin, Robert Weismantel, and Laurence Wolsey. Cutting planes in integer and mixed integer programming. *Discrete Applied Mathematics*, 123(1):397–446, 2002.
- [30] Ioannis Tsochantaridis, Thorsten Joachims, Thomas Hofmann, and Yasemin Altun. Large margin methods for structured and interdependent output variables. *Journal of machine learning research*, 6(Sep):1453–1484, 2005.
- [31] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009.
- [32] Thomas Serre, Lior Wolf, Stanley Bileschi, Maximilian Riesenhuber, and Tomaso Poggio. Robust object recognition with cortex-like mechanisms. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 29(3):411–426, 2007.