

Sparse kernel SVMs via cutting-plane training

Thorsten Joachims · Chun-Nam John Yu

Received: 12 June 2009 / Revised: 12 June 2009 / Accepted: 16 June 2009 / Published online: 23 July 2009
Springer Science+Business Media, LLC 2009

Abstract We explore an algorithm for training SVMs with Kernels that can represent the learned rule using arbitrary basis vectors, not just the support vectors (SVs) from the training set. This results in two benefits. First, the added flexibility makes it possible to find sparser solutions of good quality, substantially speeding-up prediction. Second, the improved sparsity can also make training of Kernel SVMs more efficient, especially for high-dimensional and sparse data (e.g. text classification). This has the potential to make training of Kernel SVMs tractable for large training sets, where conventional methods scale quadratically due to the linear growth of the number of SVs. In addition to a theoretical analysis of the algorithm, we also present an empirical evaluation.

Keywords Support vector machines · Kernel methods · Sparse kernel methods · Cutting plane algorithm · Basis pursuit

1 Introduction

While Support Vector Machines (SVMs) with kernels offer great flexibility and prediction performance on many application problems, their practical use is often hindered by the following two problems. Both problems can be traced back to the number of Support Vectors (SVs), which is known to generally grow linearly with the data set size (Steinwart 2003). First, training is slower than other methods and linear SVMs, where recent advances in training algorithms vastly improved training time. Second, since the prediction rule takes the form $h(\mathbf{x}) = \text{sign}[\sum_{i=1}^{\#\text{SV}} \alpha_i K(\mathbf{x}_i, \mathbf{x})]$, it is too expensive to evaluate in many applications when the number of SVs is large.

Editors: Aleksander Kołcz, Dunja Mladenić, Wray Buntine, Marko Grobelnik, and John Shawe-Taylor.

T. Joachims · C.-N.J. Yu (✉)
Dept. of Computer Science, Cornell University, Ithaca, NY 14853, USA
e-mail: cnyu@cs.cornell.edu

T. Joachims
e-mail: tj@cs.cornell.edu

This paper tackles these two problems by generalizing the notion of Support Vector to arbitrary points in input space, not just training vectors. Unlike Wu et al. (2006), who explore making the location of the points part of a large non-convex optimization problem, we propose an algorithm that iteratively constructs the set of basis vectors from a cutting-plane model. This makes our algorithm, called *Cutting-Plane Subspace Pursuit* (CPSP), efficient and modular. We analyze the training efficiency and the solution quality of the CPSP algorithm both theoretically and empirically. We find that its classification rules can be orders of magnitude sparser than the conventional support-vector representation while providing comparable prediction accuracy. The sparsity of the CPSP representation not only makes predictions substantially more efficient, it also allows the user to control training time. Especially for large datasets with sparse feature vectors (e.g. text classification), the CPSP method is substantially faster than methods that only consider basis vectors from the training set.

2 Related work

Most existing algorithms for training kernel SVMs follow the Representer Theorem and search for the optimal weight vector in the span of the training vectors $\mathbf{w} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$. This includes decomposition methods (Platt 1999; Joachims 1999) and all other dual approaches. To overcome the problems resulting from the growing number of support vectors, Burges and Schölkopf (1997) propose to post-process their solution and replace the support vector expansion with an approximation that is more sparse. Clearly, this can improve only the prediction efficiency, while it is still necessary to compute a full solution during training. For large datasets, this is intractable.

An alternative to post-processing are methods for selecting a set of basis vectors a priori. This includes sampling randomly from the training set in the NYSTROM method (Williams and Seeger 2001), greedily minimizing reconstruction error (Smola and Schölkopf 2000), and variants of the Incomplete Cholesky factorization (Fine and Scheinberg 2001; Bach and Jordan 2005). However, these selection methods are not part of the optimization process, which makes a goal-directed choice of basis vectors difficult. In fact, all but Bach and Jordan (2005) ignore label information, and all methods are limited to selecting basis vectors from the training set.

Methods like the Core Vector Machine (CVM) (Tsang et al. 2005), the Ball Vector Machine (BVM) (Tsang et al. 2007), and the active selection strategy of the LASVM method (Bordes et al. 2005) greedily select which basis vectors to include in the classification rule. While they allow the user to sacrifice solution quality to gain sparsity and training efficiency, they are also limited to selecting basis vectors from the training set.

Another set of methods are basis pursuit approaches (Keerthi et al. 2006; Vincent and Bengio 2002). They repeatedly solve the optimization problem for a given set of basis vector, and then greedily search for vectors to add or remove. The Cutting-Plane Subspace Pursuit method we propose is similar in the respect that it iteratively constructs the basis set. However, the construction of the basis set is part of the optimization algorithm itself, and the cutting-plane model makes it straightforward to add basis vectors that are not in the training set. It is not clear how to efficiently add such general basis vectors in other basis pursuit approaches.

The method most closely related to ours was proposed in Wu et al. (2006). They treat the basis vectors as variables in the SVM optimization problem, and solve the resulting non-convex program via gradient descent to a local optimum. However, training efficiency is a

bottleneck in this approach and they focus only on small datasets in their evaluation. We will consider datasets that are several orders of magnitude larger. Furthermore, we will provide theoretical results giving insight into the quality of the CPSP solution.

3 Cutting-plane algorithm for SVMs

We first introduce the Cutting-Plane Algorithm for training SVMs (Joachims 2006; Joachims et al. 2009), since it is the basis for the CPSP algorithm proposed in this paper. For a training sample, $(x_1, y_1), \dots, (x_n, y_n)$, the following is a general formulation of the large-margin training problem for learning a rule $h : \mathcal{X} \rightarrow \mathcal{Y}$ mapping from some input space \mathcal{X} to some output space \mathcal{Y} (Tsochantaridis et al. 2005). For different choices of the joint feature map $\Psi(x, y)$ and the loss function $\Delta(y, \hat{y})$, it can be specialized to classification SVMs, to Maximum-Margin Markov Networks, or various structured prediction problems.

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + \frac{C}{n} \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall i, \forall \hat{y} \in \mathcal{Y} : \langle \mathbf{w}, \Psi(x_i, y_i) - \Psi(x_i, \hat{y}) \rangle \geq \Delta(y_i, \hat{y}) - \xi_i \end{aligned} \tag{1}$$

$\langle \cdot, \cdot \rangle$ denotes an inner product. For the sake of simplicity, this paper only deals with the special case of binary classification with $\mathcal{X} = \mathfrak{R}^N$ and $\mathcal{Y} = \{-1, +1\}$, where the joint feature map is $\Psi(\mathbf{x}, y) = \frac{1}{2}y\mathbf{x}$ (or $\Psi(\mathbf{x}, y) = \frac{1}{2}y\phi(\mathbf{x})$ for the kernel $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$) and where the loss function $\Delta(y, \hat{y})$ is the zero/one-loss. In this case, it is easy to verify that (1) is equivalent to the following program, which corresponds to a binary classification SVM without explicit offset.

$$\min_{\mathbf{w}, \xi} \quad \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + \frac{C}{n} \sum_{i=1}^n \xi_i \quad \text{s.t.} \quad \forall i : y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i$$

3.1 Linear SVMs

Instead of solving (1) directly, Joachims (2006) proposes to solve the following equivalent program.

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + C\xi \\ \text{s.t.} \quad & \forall \hat{y}_1 \dots \hat{y}_n \in \mathcal{Y}^n : \\ & \left\langle \mathbf{w}, \frac{1}{n} \sum_{i=1}^n (\Psi(x_i, y_i) - \Psi(x_i, \hat{y}_i)) \right\rangle \geq \frac{1}{n} \sum_{i=1}^n \Delta(y_i, \hat{y}_i) - \xi \end{aligned} \tag{2}$$

This program has only a single slack variable ξ , and is therefore called the 1-slack formulation. It is shown in Joachims (2006) that any \mathbf{w} solving (2) is also a solution of (1), and that $\xi = \frac{1}{n} \sum_{i=1}^n \xi_i$. While (2) has a huge number of constraints, Algorithm 1 is a cutting-plane procedure that always constructs a solution of precision ϵ with at most $O(\frac{C}{\epsilon})$ active constraints (Joachims 2006; Teo et al. 2007; Joachims et al. 2009). In the experiments from Sect. 6, the number of active constraints was typically around 30—independent of the size of the training set.

Algorithm 1 maintains a working set of m constraints $\langle \mathbf{w}, \bar{\Psi}_i \rangle \geq \bar{\Delta}_i - \xi$ over which it solves the QP in Line 4. In each iteration I , the algorithm finds the most violated constraint

Algorithm 1 Cutting-plane algorithm for structural SVM (primal)

```

1: Input:  $S = ((x_1, y_1), \dots, (x_n, y_n))$ ,  $C$ ,  $\epsilon$ 
2:  $\bar{\Delta} \leftarrow \mathbf{0}$ ,  $\bar{\Psi} \leftarrow \mathbf{0}$ ,  $m \leftarrow 0$ 
3: repeat
4:    $(\mathbf{w}, \xi) \leftarrow \operatorname{argmin}_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + C\xi$ 
      s.t.  $\forall i: \langle \mathbf{w}, \bar{\Psi}_i \rangle \geq \bar{\Delta}_i - \xi$ 
5:   for  $i = 1, \dots, n$  do
6:      $\hat{y}_i \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}} \{ \Delta(y_i, \hat{y}) + \mathbf{w}^T \Psi(x_i, \hat{y}) \}$ 
7:   end for
8:    $(\bar{\Psi}, \bar{\Delta}, m) = \operatorname{remove\_inactive}(\bar{\Psi}, \bar{\Delta}, \mathbf{w}, \xi)$ 
9:    $m \leftarrow m + 1$ 
10:   $\bar{\Psi}_m \leftarrow \frac{1}{n} \sum_{i=1}^n [\Psi(x_i, y_i) - \Psi(x_i, \hat{y}_i)]$ 
11:   $\bar{\Delta}_m \leftarrow \frac{1}{n} \sum_{i=1}^n \Delta(y_i, \hat{y}_i)$ 
12: until  $\langle \mathbf{w}, \bar{\Psi}_m \rangle \geq \bar{\Delta}_m - \xi - \epsilon$ 
13: return  $(\mathbf{w}, \xi)$ 

```

from (2) (Lines 5–7) and adds it to the working set, so $m \leq I$. Typically, however, $m \ll I$, since constraints become inactive in later iterations and can be removed from the working set (Line 8). Therefore, the size m of the working set is roughly equal to the number of active constraints (i.e. $m \approx 30$). The algorithm is known to need at most $I \in O(\frac{C}{\epsilon})$ iterations to converge to an ϵ -accurate solution (Joachims 2006; Teo et al. 2007; Joachims et al. 2009). This means that the number of iterations is independent of the number of training examples n and the number of features N .

3.2 SVMs with kernels

While originally proposed for linear SVMs, the cutting-plane method can be extended to the non-linear case with kernels. Since \mathbf{w} now lies in the Reproducing Kernel Hilbert Space (RKHS) of the kernel $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$, we need to move to the dual representation. Algorithm 2 is this dual variant of Algorithm 1 and specialized to the case of binary classification as implemented in the SVM^{perf} software. It replaces the primal QP with its Wolfe dual in Line 5, and the solution vector in the RKHS is $\mathbf{w} \equiv \sum_i \alpha_i \bar{\Psi}_i$. Note that sums cannot be computed efficiently in the RKHS. Therefore, the assignment operator \leftarrow is replaced with a rewrite operator \equiv where appropriate. However, it is easy to verify that all inner products (Lines 4, 9, 15) can be computed as sums of kernel evaluations. The $O(\frac{C}{\epsilon})$ bound on the number of iterations (Joachims 2006; Joachims et al. 2009) holds independent of whether a kernel is used or not, but how does the time complexity per iteration change when moving from a linear to a kernelized SVM?

Without kernels, any iteration in Algorithm 2 takes at most $O(m^3)$ for solving the QP, $O(m^2)$ for ξ , $O(mN)$ for \mathbf{w} , $O(nN)$ for computing the most violated constraint, $O(n)$ for $\bar{\Delta}$, $O(nN)$ for $\bar{\Psi}$, and $O(mN)$ for adding a row/column to \mathbf{H} . So the overall time complexity is $O(m^3 + mN + nN)$.

When using a kernel, however, computing $\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle$ and \mathbf{H} becomes more expensive than in the linear case. Denote with \mathbf{Y} the matrix with $Y_{ij} = (y_i - \hat{y}_i)$ for the j -the constraint in $\bar{\Psi}$. To find the most violated constraint, for each example one now needs to evaluate

$$\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle = \sum_{k=1}^n \left(\sum_{j=1}^m \alpha_j Y_{kj} \right) K(\mathbf{x}_i, \mathbf{x}_k). \tag{3}$$

Algorithm 2 Cutting-plane algorithm for classification SVM (dual)

```

1: Input:  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ ,  $C, \epsilon, K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ 
2:  $\bar{\Delta} \leftarrow \mathbf{0}, \bar{\Psi} \leftarrow \mathbf{0}, \mathbf{H} \leftarrow \mathbf{0}, m \leftarrow 0$ 
3: repeat
4:    $\mathbf{H} \leftarrow (\mathbf{H}_{ij})_{1 \leq i, j \leq m}$ , where  $\mathbf{H}_{ij} = \langle \bar{\Psi}_i, \bar{\Psi}_j \rangle$ 
5:    $\alpha \leftarrow \operatorname{argmax}_{\alpha \geq 0} \alpha^T \bar{\Delta} - \frac{1}{2} \alpha^T \mathbf{H} \alpha$  s.t.  $\alpha^T \mathbf{1} \leq C$ 
6:    $\xi \leftarrow \frac{1}{C} (\alpha^T \bar{\Delta} - \alpha^T \mathbf{H} \alpha)$ 
7:    $\mathbf{w} \equiv \sum_i \alpha_i \bar{\Psi}_i$ 
8:   for  $i = 1, \dots, n$  do
9:      $\hat{y}_i \leftarrow \operatorname{sign}(\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle - y_i)$ 
10:  end for
11:   $(\bar{\Psi}, \bar{\Delta}, m) = \operatorname{remove\_inactive}(\bar{\Psi}, \bar{\Delta}, \alpha)$ 
12:   $m \leftarrow m + 1$ 
13:   $\bar{\Psi}_m \equiv \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i) \phi(\mathbf{x}_i)$ 
14:   $\bar{\Delta}_m \leftarrow \frac{1}{2n} \sum_{i=1}^n |y_i - \hat{y}_i|$ 
15: until  $\langle \mathbf{w}, \bar{\Psi}_m \rangle \geq \bar{\Delta}_m - \xi - \epsilon$ 
16: return  $(\mathbf{w}, \xi)$ 

```

Over all n examples, this has a cost of $O(n^2 + mn)$. Similarly, adding a row/column for the new $\bar{\Psi}_m$ to the Gram matrix \mathbf{H} now requires computing

$$\forall i : \mathbf{H}_{mi} = \mathbf{H}_{im} = \langle \bar{\Psi}_i, \bar{\Psi}_m \rangle = \sum_{k=1}^n \sum_{l=1}^n Y_{ki} Y_{lm} K(\mathbf{x}_k, \mathbf{x}_l) \tag{4}$$

This takes time $O(mn^2)$, counting a single kernel evaluation as $O(1)$. So, the overall time complexity of an iteration when kernels are used is $O(m^3 + mn^2)$. This $O(n^2)$ scaling is not practical for any reasonably-sized dataset, and the algorithm has worse constants than decomposition methods like SVM^{light} that also typically scale $O(n^2)$. However, Algorithm 2 does provide a path to a substantially more efficient algorithm that is explored in the next section.

4 Cutting-plane subspace pursuit

Is it possible to remove the $O(n^2)$ scaling behavior? Here is the intuition for the approach we take. A property of the cutting-plane algorithm is that it iteratively constructs a low-dimensional subspace $\mathcal{W} = \operatorname{span}(\bar{\Psi}_1, \dots, \bar{\Psi}_m) = \{\sum_{i=1}^m \beta_i \bar{\Psi}_i : \beta \in \mathfrak{R}^m\}$ in which the final solution

$$\mathbf{w} = \sum_{i=1}^m \alpha_i \bar{\Psi}_i \tag{5}$$

is guaranteed to lie. Instead of using the Representer Theorem and considering the larger subspace $\mathcal{F} = \operatorname{span}(\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))$ to express the optimal weight vector as $\mathbf{w} = \sum_{i=1}^n \alpha'_i \phi(\mathbf{x}_i)$, the cutting-plane method tells us that we only need to consider the subspace $\mathcal{W} \subset \mathcal{F}$ in each iteration, where $m \ll n$ and m does not grow with n . Our core idea is to find a small set of basis vectors $\mathbf{b}_1, \dots, \mathbf{b}_k$ so that

$$\mathcal{W}' = \operatorname{span}(\phi(\mathbf{b}_1), \dots, \phi(\mathbf{b}_k)) \approx \mathcal{W} \tag{6}$$

Algorithm 3 Cutting-plane subspace pursuit (CPSP) algorithm

```

1: Input:  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ ,  $C, \epsilon, k_{\max}, K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ 
2:  $\bar{\Delta} \leftarrow \mathbf{0}, \bar{\Psi} \leftarrow \mathbf{0}, \hat{\Psi} \leftarrow \mathbf{0}, \mathbf{H} \leftarrow \mathbf{0}, \mathbf{B} \leftarrow \emptyset, m \leftarrow 0$ 
3: repeat
4:    $\mathbf{H} \leftarrow (\mathbf{H}_{ij})_{1 \leq i, j \leq m}$ , where  $\mathbf{H}_{ij} = \langle \hat{\Psi}_i, \hat{\Psi}_j \rangle$ 
5:    $\boldsymbol{\alpha} \leftarrow \operatorname{argmax}_{\boldsymbol{\alpha} \geq \mathbf{0}} \boldsymbol{\alpha}^T \bar{\Delta} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{H} \boldsymbol{\alpha}$  s.t.  $\boldsymbol{\alpha}^T \mathbf{1} \leq C$ 
6:    $\xi \leftarrow \frac{1}{C} (\boldsymbol{\alpha}^T \bar{\Delta} - \boldsymbol{\alpha}^T \mathbf{H} \boldsymbol{\alpha})$ 
7:    $\mathbf{w} \equiv \sum_i \alpha_i \hat{\Psi}_i$ 
8:   for  $i = 1, \dots, n$  do
9:      $\hat{y}_i \leftarrow \operatorname{sign}(\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle - y_i)$ 
10:  end for
11:   $(\bar{\Psi}, \bar{\Delta}, m) = \operatorname{remove\_inactive}(\bar{\Psi}, \bar{\Delta}, \boldsymbol{\alpha})$ 
12:   $m \leftarrow m + 1$ 
13:   $\bar{\Psi}_m \equiv \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i) \phi(\mathbf{x}_i)$ 
14:   $\bar{\Delta}_m \leftarrow \frac{1}{2n} \sum_{i=1}^n |y_i - \hat{y}_i|$ 
15:  if  $|\mathbf{B}| < k_{\max}$  then  $\mathbf{B} \leftarrow \operatorname{extend\_basis}(\mathbf{B}, \bar{\Psi}_m)$ 
16:  for  $i = 1, \dots, k$  do
17:     $\hat{\Psi}_i \equiv \operatorname{project}(\bar{\Psi}_i, \mathbf{B})$ 
18:  end for
19: until  $\langle \mathbf{w}, \bar{\Psi}_m \rangle \geq \bar{\Delta}_m - \xi - \epsilon$ 
20: return( $\mathbf{w}, \xi$ )

```

which means that we can express the final solution from (5) as

$$\mathbf{w} \approx \sum_{i=1}^k \alpha_i'' \phi(\mathbf{b}_i) \tag{7}$$

This enables efficient prediction using the rule $h(\mathbf{x}) = \operatorname{sign}[\sum_{i=1}^k \alpha_i'' K(\mathbf{b}_i, \mathbf{x})]$, given that k is small. Furthermore, we will elaborate in the following how projecting into the subspace \mathcal{W}' allows computing \mathbf{H} and $\langle \mathbf{w}, \phi(\mathbf{x}) \rangle$ in time independent of n .

To understand the intuition behind our approach, consider the ideal case where for every $\bar{\Psi}_i$ there exists a vector \mathbf{b}_i in input space (not necessarily from the training set) so that $\bar{\Psi}_i = \phi(\mathbf{b}_i)$ (as it does in the linear case, where $\mathbf{b}_i = \frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j) \mathbf{x}_j$). Then we could replace each $\bar{\Psi}_i$ with $\phi(\mathbf{b}_i)$, and it is easy to verify that the time complexity of an iteration goes down to $O(m^3 + mn)$ —almost like in the linear case. Furthermore, the resulting classifier would only have $k = m \approx 30$ “Support Vectors”—or, more generally named, “Basis Vectors”—, making it much faster than conventional SVM classifiers that often have 10000’s of SVs.

Unfortunately, in most cases there will be no single pre-image \mathbf{b} so that $\bar{\Psi} = \phi(\mathbf{b})$. However, in any iteration it suffices to find a set of pre-image vectors so that $\bar{\Psi}_1, \dots, \bar{\Psi}_m$ lie (approximately) in their span. In particular, we are looking for a set of basis vectors $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_k)$, $\mathbf{b}_i \in \mathfrak{R}^N$, so that for every $\bar{\Psi}_i$ in $\bar{\Psi}$

$$\min_{\boldsymbol{\beta}} \left\| \bar{\Psi}_i - \sum_{j=1}^k \beta_j \phi(\mathbf{b}_j) \right\|^2 \leq \delta \tag{8}$$

for some (small) $\delta \geq 0$. When computing \mathbf{H} and $\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle$ in Algorithm 2, we then replace $\bar{\Psi}_i$ with its projection $\hat{\Psi}_i \equiv \sum_{j=1}^k \beta_j \phi(\mathbf{b}_j)$. This is summarized in Algorithm 3, which we

call the *Cutting-Plane Subspace Pursuit* (CPSP) algorithm. It is easy to verify that \mathbf{H} and all $(\mathbf{w}, \phi(\mathbf{x}_i))$ can now be computed in time $O(m^2k^2)$ (or $O(k^3 + m^2k + mk^2)$) and $O(mk + kn)$, respectively.

Using $\hat{\Psi}$ instead of $\bar{\Psi}$ in Algorithm 3 is straightforward. However, we still have to define how the function *extend_basis*($\mathbf{B}, \bar{\Psi}_m$) (Line 15) computes the set of basis vectors $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_k)$ and how the function *project*($\bar{\Psi}_i, \mathbf{B}$) (Line 17) computes the approximate cutting-planes $\hat{\Psi}_i$. This is addressed in the following.

4.1 Projecting cutting-planes $\bar{\Psi}$ onto B

For a given subspace $\text{span}(\phi(\mathbf{b}_1), \dots, \phi(\mathbf{b}_k))$, the function *project*($\bar{\Psi}_i, \mathbf{B}$) computes the projection $\hat{\Psi}_i$ of a cutting-plane $\bar{\Psi}_i$ via the following least-squares problem:

$$\hat{\Psi}_i = \sum_{j=1}^k \beta_j \phi(\mathbf{b}_j) \quad \text{where } \beta = \min_{\beta} \left\| \bar{\Psi}_i - \sum_{j=1}^k \beta_j \phi(\mathbf{b}_j) \right\|^2 \tag{9}$$

To accomodate kernels, we maintain the $k \times k$ -matrix \mathbf{G} with $G_{ij} = K(\mathbf{b}_i, \mathbf{b}_j)$ and the $k \times n$ -matrix \mathbf{K} with $K_{ij} = K(\mathbf{b}_i, \mathbf{x}_j)$. The solution of the least-squares problem can then be written as $\beta = \frac{1}{2n} \mathbf{G}^{-1} \mathbf{K} \mathbf{Y}_{*i}$. It is more efficient, however, to use the Cholesky decomposition L_G of \mathbf{G} (i.e. $\mathbf{G} = L_G L_G^T$). With L_G , the solution can be computed via forward and back-substitution from $L_G \gamma = \frac{1}{2n} \mathbf{K} \mathbf{Y}_{*i}$ and $L_G^T \beta = \gamma$ in time $O(k^2 + kn)$. This excludes the time for computing \mathbf{K} , \mathbf{G} , and its Cholesky decomposition L_G , since these need to be computed only once and can then be used until \mathbf{B} changes. This is further discussed in the next section.

4.2 Constructing the set of basis vectors B

The method for constructing the set of basis vectors \mathbf{B} is the final part of Algorithm 3 that still needs to be specified. The goal is to find a set of basis vectors $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_k)$ such that for some small $\delta \geq 0$, all $\hat{\Psi}_i$ that are active in the current iteration fulfill (8). Recomputing \mathbf{B} in each iteration would be costly, but fortunately it is unnecessary. Only $\bar{\Psi}_m$ is new and all other $\bar{\Psi}_i$ are already well approximated by the set of basis vectors from the previous iteration. The function *extend_basis*($\mathbf{B}, \bar{\Psi}_m$) therefore only adds some new basis vectors to \mathbf{B} that are required to fit $\bar{\Psi}_m$. Note that this can only improve the fit for the other $\bar{\Psi}_i$.

To decide which basis vectors to add, we follow (Burges and Schölkopf 1997) and take a greedy approach. We search for the basis vector \mathbf{b}_{k+1} that minimizes the residual error for $\bar{\Psi}_m$, where $\hat{\Psi}_m$ is the projection for the current \mathbf{B} .

$$(\beta', \mathbf{b}') = \underset{\beta_{k+1}, \mathbf{b}_{k+1}}{\text{argmin}} \left\| \bar{\Psi}_m - \hat{\Psi}_m - \beta_{k+1} \phi(\mathbf{b}_{k+1}) \right\|^2 \tag{10}$$

This optimization problem is commonly referred to as the “preimage” problem. While exact solutions are difficult to obtain, approximate solutions can be found with gradient-based methods (Burges 1996; Schölkopf and Smola 2002) or randomized search. In this paper, we use the fix-point iteration approach described in Schölkopf and Smola (2002, Sect. 18.2.2) for the RBF kernel to solve (10) to a local optimum. In this way we can efficiently produce arbitrary vectors as basis vectors to add to \mathbf{B} . We refer to the Cutting-Plane Subspace Pursuit algorithm with this preimage method as “CPSP” in the following.

To evaluate in how far general basis vectors improve sparsity, we also explore a second preimage method that is restricted to using basis vectors from the training set. We refer

to this method as “CPSP(tr)”. As proposed by Smola and Schölkopf (2000) (and used by most of the methods we compare against), we randomly sample 59 feature vectors from the training set and pick the one with maximum objective value in (10). Note that this alternative strategy is introduced only to evaluate the benefit of selecting “support vectors” outside the training set.

The number of new basis vectors to add for each $\bar{\Psi}_m$ is a design choice. One could either use a fixed number, or keep adding until a certain δ is achieved. In the following experiments, we use the simplest choice and add exactly one basis vector for each $\bar{\Psi}_m$ until the maximum size k_{\max} specified by the user has been reached. At that point, no further vectors are added and $extend_basis(\mathbf{B}, \bar{\Psi}_m)$ returns \mathbf{B} unchanged.

After a new \mathbf{b}_{k+1} is added to \mathbf{B} , a column/row needs to be added to the kernel matrices \mathbf{G} and \mathbf{K} . This takes $O(n + k)$ kernel evaluations, and the Cholesky factorization of \mathbf{G} can be updated in time $O(k^2)$.

5 Theoretical analysis

Before evaluating the CPSP algorithm empirically, we first give a theoretical characterization of the quality of its solutions and the number of iterations it takes until convergence.

The following theorem gives an upper bound on the number of iterations of Algorithm 3. It extends the general results (Joachims 2006; Teo et al. 2007; Joachims et al. 2009) for cutting-plane training of SVMs to the CPSP algorithm.

Theorem 1 *For parameter C , precision ϵ , training-set size n , and basis-set size k_{\max} , Algorithm 3 terminates after at most $O(k_{\max} + \frac{C}{\epsilon})$ iterations.*

Proof After the first k_{\max} iterations, the basis \mathbf{B} becomes fixed, and from then on we are essentially solving the optimization problem:

$$\begin{aligned} & \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C\xi \\ \text{s.t. } & \forall \hat{y} \in \{-1, 1\}^n : \left\langle \mathbf{w}, \sum_{i=1}^n (y_i - \hat{y}_i) \phi(\mathbf{x}_i) \right\rangle \geq \sum_{i=1}^n \Delta(y_i, \hat{y}_i) - \xi \quad \text{and} \quad (11) \\ & \mathbf{w} \in \sum_{\mathbf{b}_i \in \mathbf{B}} \beta_i \phi(\mathbf{b}_i) \end{aligned}$$

Let $P_{\mathbf{B}}$ be the orthogonal projection operator onto the subspace spanned by \mathbf{B} . Such an orthogonal projection operator always exists in a Hilbert Space. After folding the subspace constraint into the objective by replacing \mathbf{w} with $P_{\mathbf{B}}\mathbf{w}$, the above optimization problem can be re-written as (using the self-adjointness and linearity of $P_{\mathbf{B}}$):

$$\begin{aligned} & \min_{\mathbf{w}, \xi} \frac{1}{2} \|P_{\mathbf{B}}\mathbf{w}\|^2 + C\xi \\ \text{s.t. } & \forall \hat{y} \in \{-1, 1\}^n : \left\langle \mathbf{w}, \sum_{i=1}^n (y_i - \hat{y}_i) P_{\mathbf{B}}\phi(\mathbf{x}_i) \right\rangle \geq \sum_{i=1}^n \Delta(y_i, \hat{y}_i) - \xi \end{aligned}$$

Finally the operator $P_{\mathbf{B}}$ in the objective can be dropped since if \mathbf{w} contains any component in \mathbf{B}^{\perp} , it will only increase the objective without changing value of the LHS of the constraints. This is in the form of the general Structural SVM optimization problem solved

by Algorithm 1, with the feature space changed from being spanned by $\phi(x_i)$ to being spanned by $P_B\phi(x_i)$. The $O(\frac{C}{\epsilon})$ iteration bound from (Joachims 2006; Teo et al. 2007; Joachims et al. 2009) therefore applies. \square

The time complexity of each iteration was already discussed in Sect. 4, but can be summarized as follows. In iterations where no new basis vector is added to \mathbf{B} , the time complexity is $O(m^3 + mk^2 + kn)$, since only the new $\bar{\Psi}_m$ needs to be projected and the respective column be added to \mathbf{H} . In iterations where \mathbf{B} is extended, the time complexity is $O(m^3 + k^2m + km^2 + kmn)$ plus the time it takes to solve the preimage problem (10). Note that typical values are $m \approx 30$, $k \in [10..1000]$, and $n > 10000$.

The following theorem describes the quality of the solution at termination, accounting for the error incurred by projecting on an imperfect \mathbf{B} . Most importantly, the theorem justifies our use of (10) for deciding which basis vectors to add.

Theorem 2 *When Algorithm 3 terminates with $\|\bar{\Psi}_i - \hat{\Psi}_i\| \leq \delta$ for all $\bar{\Psi}_i$ and $\hat{\Psi}_i$, then the primal objective value o of the solution found does not exceed the exact solution o^* by more than $o - o^* \leq C(\delta\sqrt{2C} + \epsilon)$.*

Proof Let \mathbf{w}^* be the optimal solution with value o^* . We know that the optimal \mathbf{w}^* satisfies $\|\mathbf{w}^*\| \leq \sqrt{2C}$. Hence for all i ,

$$|\langle \mathbf{w}, \bar{\Psi}_i \rangle - \langle \mathbf{w}, \hat{\Psi}_i \rangle| \leq \|\mathbf{w}\| \|\bar{\Psi}_i - \hat{\Psi}_i\| \leq \delta\sqrt{2C}$$

Let P_B be the orthogonal projection on the subspace spanned by $\phi(\mathbf{b}_i)$ in the final basis \mathbf{B} . Let \mathbf{v}^* be the optimal solution to the optimization problem (11) restricted to the subspace \mathbf{B} , we have:

$$\begin{aligned} o &\leq \frac{1}{2} \|\mathbf{v}^*\|^2 + C(\xi + \epsilon) \\ &= \frac{1}{2} \|\mathbf{v}^*\|^2 + C \max_{1 \leq i \leq m} [\bar{\Delta}_i - \langle \mathbf{v}^*, \hat{\Psi}_i \rangle] + C\epsilon \\ &\leq \frac{1}{2} \|P_B \mathbf{w}^*\|^2 + C \max_{1 \leq i \leq m} [\bar{\Delta}_i - \langle P_B \mathbf{w}^*, \hat{\Psi}_i \rangle] + C\epsilon \\ &\quad \text{[since } \mathbf{v}^* \text{ is the optimal solution wrt the basis } \mathbf{B}] \\ &= \frac{1}{2} \|P_B \mathbf{w}^*\|^2 + C \max_{1 \leq i \leq m} [\bar{\Delta}_i - \langle \mathbf{w}^*, P_B \hat{\Psi}_i \rangle] + C\epsilon \\ &= \frac{1}{2} \|P_B \mathbf{w}^*\|^2 + C \max_{1 \leq i \leq m} [\bar{\Delta}_i - \langle \mathbf{w}^*, \hat{\Psi}_i \rangle] + C\epsilon \\ &\leq \frac{1}{2} \|\mathbf{w}^*\|^2 + C \max_{1 \leq i \leq m} [\bar{\Delta}_i - \langle \mathbf{w}^*, \hat{\Psi}_i \rangle] + C\epsilon \\ &\leq \frac{1}{2} \|\mathbf{w}^*\|^2 + C \max_{1 \leq i \leq m} [\bar{\Delta}_i - \langle \mathbf{w}^*, \bar{\Psi}_i \rangle + \delta\sqrt{2C}] + C\epsilon \\ &\leq o^* + C(\delta\sqrt{2C} + \epsilon) \end{aligned}$$

\square

6 Experimental analysis

The following experiments are designed to evaluate how the CPSP method compares to conventional training methods in terms of sparsity (i.e. prediction efficiency) and training efficiency. In particular, they explore whether the use of general basis vectors outside the training set improves prediction accuracy and training efficiency, and how both quantities scale with basis set size k_{\max} .

Our implementation of the CPSP algorithm is available for download at http://svmlight.joachims.org/svm_perf.html.

We compare the CPSP algorithm with the exact solution computed by SVM^{light}, as well as approximate solutions of the NYSTROM method (NYSTROM) (Williams and Seeger 2001), the Incomplete Cholesky Factorization (INCCHOL) (Fine and Scheinberg 2001), the Core Vector Machine (CVM) (Tsang et al. 2005), the Ball Vector Machine (BVM) (Tsang et al. 2007), and LASVM with margin-based active selection and finishing (Bordes et al. 2005). Both the NYSTROM method and the Incomplete Cholesky Factorization are implemented in SVM^{perf} as described in Joachims et al. (2009). We use the RBF-Kernel $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$ in all experiments. The cache sizes of SVM^{light}, CVM, BVM, and LASVM were set to 1 GB.

We compare on the following five binary classification tasks, each split into training/validation/test set. If not mentioned otherwise, parameters (i.e. C and γ) are selected to maximize performance on the validation set for each method and k_{\max} individually. Both C and γ are explored on a log-scale. The first dataset is Adult as compiled by John Platt with 123 features and using a train/validation/test split of 20000/6281/6280. Second is the Reuters RCV1 CCAT text-classification dataset with 47236 features. We use 78127 examples from the original test set for training and split the original training set into validation and test sets of sizes 11575 and 11574 respectively. Third and fourth, we classify the digit “0” against the rest (OCR0), as well as classify the digits “01234” against the digits “56789” (OCR*) on the MNIST dataset. The MNIST datasets have 780 features and we use a training/validation/test split of 50000/5000/5000. Finally, we use the IJCNN (task 1) dataset as pre-processed by Chih-Jen Lin. It has 22 features and we use a training/validation/test split of 113533/14169/14169.

How accurate are the solutions for a given sparsity budget? We first explore a scenario where the application demands an upper bound on the number of support vectors to achieve a desired computational efficiency at prediction time. Table 1 summarizes the results. The first two lines show the performance of SVM^{light} for the linear kernel and SVM^{light} for the RBF kernel as baselines to compare against. All but the Adult dataset show substantial non-linear structure, and the RBF kernel outperforms a linear SVM. The number of SVs when using the RBF kernel is given in the third line. The remaining lines in Table 1 are for the “sparse” methods, all of which use $k_{\max} = 1000$ basis vectors. Note that this is well below the 2786 to 28748 support vectors required by the exact SVM.

The CPSP algorithm with the general preimage method matches the accuracy of SVM^{light} up to ± 0.1 . This means that the prediction accuracy is roughly the same as for the exact method, while speeding up prediction by a factor between 2.7 to 28. We will see in Sect. 6 that far fewer than $k_{\max} = 1000$ basis vectors would have sufficed on some of the tasks for the CPSP algorithm, leading to an even larger speedup.

Random sampling of the basis vectors in the NYSTROM method and the Incomplete Cholesky Factorization (INCCHOL) perform consistently worse than the CPSP method, except on the OCR0 dataset where all methods do well with $k_{\max} = 1000$ basis vectors. The

Table 1 Prediction accuracy with $k_{\max} = 1000$ basis vectors (except SVM^{light}, where the number of SVs is shown in the *third line*) using the RBF kernel (except linear)

	ADULT	CCAT	OCR0	OCR*	IJCNN
SVM-light (linear)	84.4	94.2	99.4	87.6	92.2
SVM-light (RBF)	84.4	95.1	99.8	98.6	99.4
#SV	7125	28748	2786	19309	9243
CPSP	84.5	95.0	99.8	98.5	99.3
CPSP(tr)	84.1	93.5	99.8	97.9	99.2
NYSTROM	84.3	92.5	99.7	97.0	99.1
INCCHOL	84.0	92.1	99.7	97.0	98.9
CVM	78.4	88.1	99.8	96.9	98.2
BVM	77.1	56.1	99.8	89.1	97.7
LASVM	83.8	91.7	99.8	97.2	97.5

Core Vector Machine (CVM), the Ball Vector Machine (BVM), and the LASVM algorithm with active selection are not competitive on most datasets.

How does accuracy scale with basis-set size? As mentioned above, a lower number of basis vectors $k_{\max} \ll 1000$ could have sufficed to get reasonable accuracy on some datasets. The plots in Fig. 1 investigate this question and show by how much the test accuracies for a given k_{\max} are lower than the accuracy of the exact SVM solution. In each plot, 0 corresponds to the accuracy of the exact SVM solution.

Figure 1 shows that CPSP dominates all other methods not only for $k_{\max} = 1000$, but over the whole range. For all datasets, the CPSP method using general preimages outperforms the other methods especially for small numbers of basis vectors. In particular, on three of the five datasets, CPSP already performs within 1% of the exact solution with only 50 basis vectors. Similarly, on all five datasets does CPSP perform equivalent or better than the linear SVM when using 50 basis vectors or more. Especially on ADULT, CCAT, and OCR0 far fewer than $k_{\max} = 1000$ basis vectors would have sufficed to reach an acceptable level of performance.

What is the benefit of using general basis vectors? A key premise of the paper is that using basis vectors outside the training set is beneficial. To test its validity, Fig. 1 and Table 1 include the performance of the CPSP(tr) algorithm, which is identical to CPSP except for selecting basis vectors only from the training set. Consistently over all dataset, Fig. 1 shows that the general CPSP algorithm provides improved prediction accuracy over CPSP(tr) especially for small numbers of basis vectors. The difference is largest on the CCAT dataset, where the general CPSP algorithm with 10 basis vectors already performs at an accuracy for which CPSP(tr) requires about 1000 basis vectors. This confirms our hypothesis that basis vectors outside the training set can lead to more accurate solutions at a given level of sparsity.

How accurate is the objective value? The four methods CPSP, CPSP(tr), NYSTROM, and INCCHOL all optimize the same objective function as a regular SVM. How well do they manage to minimize this objective? The plots in Fig. 2 show by what factor their primal objective value is higher than the exact SVM solution. All methods use the same parameters (i.e. C and γ), which are picked to optimize validation set accuracy of the exact SVM.

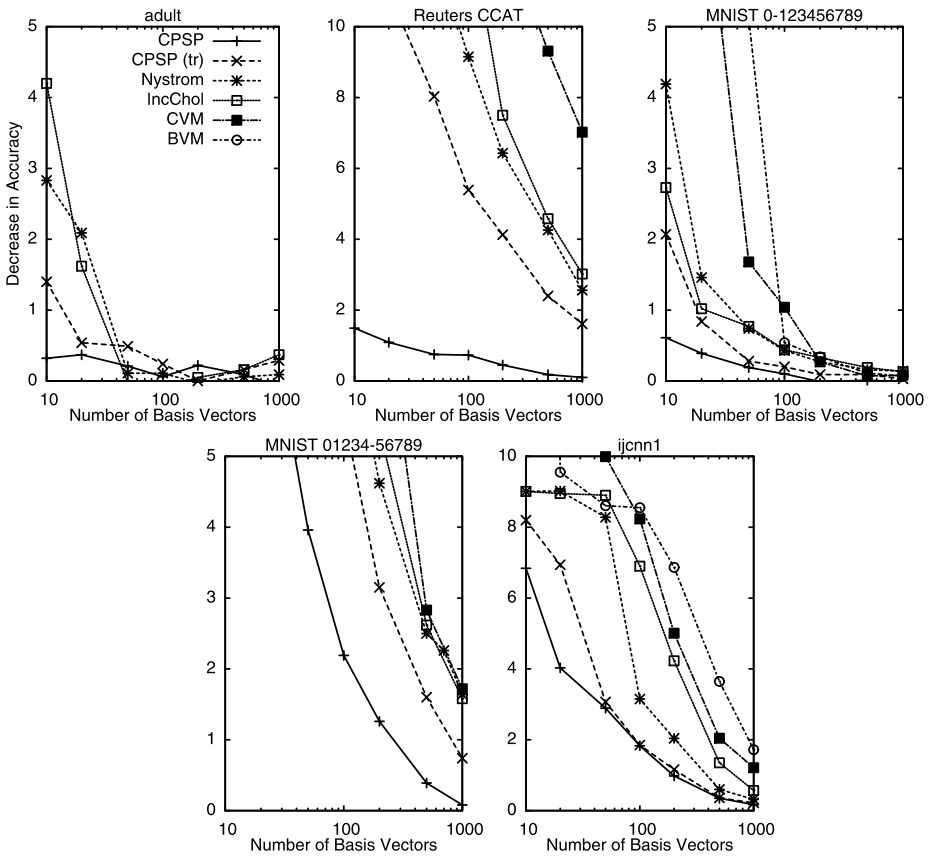


Fig. 1 Decrease in accuracy w.r.t. exact SVM for different basis-set sizes k_{max}

Again, CPSP dominates the other methods, and the curves in Fig. 2 very much resemble the curves in Fig. 1. This verifies that finding a subspace that contains a solution of low objective value is indeed crucial for good prediction accuracy, and that the subspaces found by CPSP are of superior fidelity (also compared to CPSP(tr)).

What is the training and test efficiency? While efficiency at test time may be the dominant criterion for many applications, training has to be tractable as well. Since CPSP does more work in each iteration (e.g. solve a pre-image problem), one superficial concern might be that the training process is slow. However, the following shows that the increased sparsity observed in Fig. 1 not only improves prediction efficiency, but also speeds up training. This is a key difference to the Reduced Set method (Burges and Schölkopf 1997). The Reduced Set method requires solving an exact SVM, making it intractable for large training sets.

Table 2 compares the training time and number of basis vectors that each method needs to reach a certain prediction accuracy. The experiment simulates how a user may chose to trade prediction accuracy for improved training and test efficiency. In particular, Table 2 shows the number of basis vectors (left) and the training time (right) to reach a test accuracy that is not more than 0.5% below the test accuracy of the exact SVM. Basis set sizes $k_{max} \in$

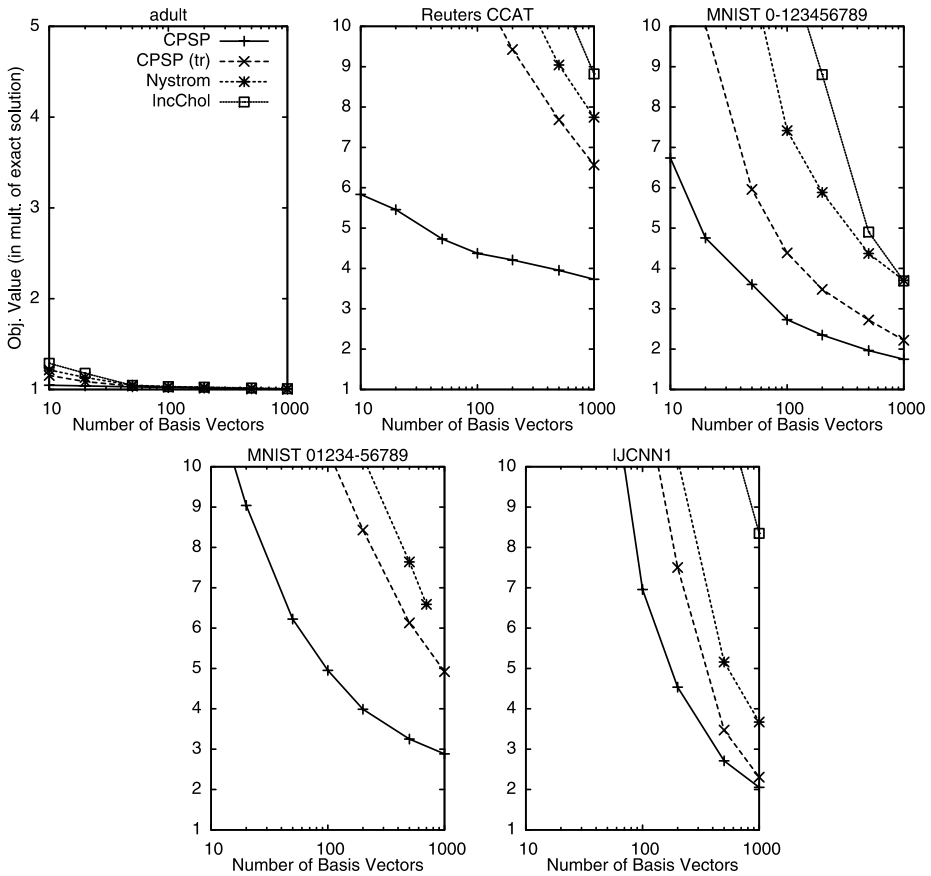
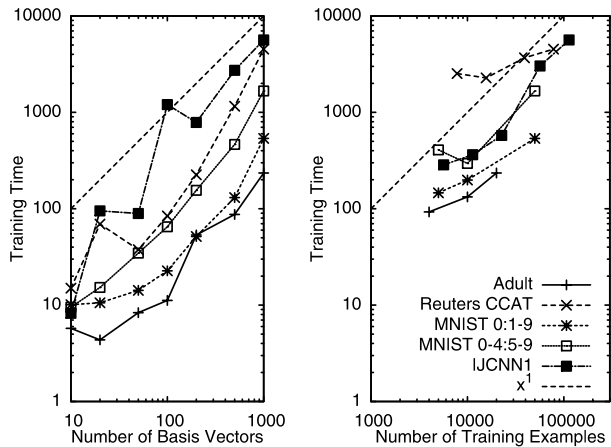


Fig. 2 Primal objective value of the approximate solutions expressed as multiples of the exact SVM solution

Table 2 Number of SV (*left*) and training time (*right*) to reach an accuracy that is not more than 0.5% below the accuracy of the exact solution of SVM-light (see Table 1). The RBF kernel is used for all methods. '>' indicates that the largest tractable solution did not achieve the target accuracy

	Number of SV					Training Time (CPU-Seconds)				
	ADULT	CCAT	OCR0	OCR*	IJCNN	ADULT	CCAT	OCR0	OCR*	IJCNN
SVM-light	7125	28748	2786	19309	9243	56	9272	400	4629	1175
CPSP	10	200	20	500	500	6	225	11	465	2728
CPSP(tr)	50	5000	50	2000	500	30	88873	57	8967	2178
NYSTROM	50	> 5000	100	5000	1000	10	> 2281	37	2270	1572
INCCHOL	50	> 2000	100	> 2000	2000	14	> 21673	66	> 12330	59454
CVM	5000	20000	200	5000	2000	43	23730	2	497	29
BVM	5000	20000	200	5000	5000	67	11004	2	538	229
LASVM	2000	10000	100	2000	5000	51	3433	5	295	705

Fig. 3 Training times of CPSP for varying basis-set sizes (*left*) and training-set sizes with $k_{\max} = 1000$ (*right*)



{10, 20, 50, 100, 200, 500, 1000, 2000, 5000, 10000, 20000} were tried for each method and the results for the smallest k_{\max} that achieved the target accuracy are shown.

Table 2 shows that the solutions found by CPSP are typically substantially more sparse than those of the other methods. Compared to the exact solution, they lead to an 18 to 712 fold speed-up at prediction time. Compared to the other approximate methods, the speed-up is still typically between 5 and 10.

The increased sparsity also leads to very efficient training times for CPSP. While it is difficult to rank methods by aggregated training times, CPSP is clearly among the fastest methods in the comparison, especially on those tasks where general basis vectors provide a substantial gain in sparsity. On the CCAT text-classification dataset, it is orders of magnitude faster than any of the other methods (and CPSP(tr)). For such large and sparse data, there simply does not appear to be a small subset of training vectors that can represent an accurate classifier, and the increased sparsity from allowing general basis vectors greatly improves training efficiency. More generally, all training methods for SVMs scale super-linearly with the number of SVs, so that improving sparsity is the key to making large-scale training tractable. The scaling properties of CPSP are explored in more detail in the following section.

Comparing to the results published in Wu et al. (2006), our method is substantially faster. They focus mostly on small training sets with less than 1000 examples. The USPS OCR dataset with 7291 and 256 features is their largest dataset, and they report an average training time of ~ 2400 seconds. This dataset roughly compares to our OCR0 task. However, the OCR0 dataset is an order of magnitude larger.

How does training time scale with basis-set size? Finally, let us investigate the efficiency of CPSP in more detail. The left-hand plots in Fig. 3 show training time for different values of k_{\max} . Parameters (regularization C , RBF γ) are individually picked via CV for each method and k_{\max} . While the theoretical time complexity is $O(k_{\max}^3)$, the actual scaling shown in Fig. 3 (left) is much more benign. For $k_{\max} < 1000$, the time contribution of the cubic parts of the algorithm (e.g. repeatedly updating the Cholesky factorization L_G) is still rather small, and the scaling behavior is only modestly super-linear.

How does training time scale with training-sample size? Finally, the right-hand plot in Fig. 3 shows training time of CPSP for different training set sizes. Parameters (regularization C , RBF γ) are individually picked via CV for each method and training set size. As

expected from the theoretical analysis, the scaling behavior is roughly linear, making CPSP particularly attractive for large datasets.

7 Conclusions

We presented a training algorithm for kernel SVMs that constructs a sparse set of basis vectors as part of the cutting-plane optimization process. The algorithm's efficiency and effectiveness is characterized theoretically, and an experimental comparison shows that it produces solutions of a sparsity that is superior to NYSTROM, INCCHOL, CVM, BVM, and LASVM. We find that the ability to use basis vectors outside the training set substantially contributes to this gain in sparsity and efficiency, especially on large datasets with sparse feature vectors.

Acknowledgements This work was funded in part under NSF award IIS-0713483.

References

- Bach, F., & Jordan, M. (2005). Predictive low-rank decomposition for kernel methods. In *ICML* (pp. 33–40).
- Bordes, A., Ertekin, S., Weston, J., & Bottou, L. (2005). Fast kernel classifiers with online and active learning. *JMLR*, 6, 1579–1619.
- Burges, C. (1996). Simplified support vector decision rules. In *ICML* (pp. 71–77).
- Burges, C., & Schölkopf, B. (1997). Improving the accuracy and speed of support vector learning machines. *NIPS*, 9, 375–381.
- Fine, S., & Scheinberg, K. (2001). Efficient SVM training using low-rank kernel representations. *JMLR*, 2, 243–264.
- Joachims, T. (1999). Making large-scale SVM learning practical. In Schölkopf, B., Burges, C., Smola, A. (Eds.), *Advances in kernel methods—support vector learning* (pp. 169–184). Cambridge: MIT Press.
- Joachims, T. (2006). Training linear SVMs in linear time. In *SIGKDD* (pp. 217–226).
- Joachims, T., Finley, T., & Yu, C. N. (2009). Cutting-plane training of structural SVMs. *Machine Learning*, 76(1).
- Keerthi, S., Chapelle, O., & DeCoste, D. (2006). Building support vector machines with reduced classifier complexity. *JMLR*, 7, 1493–1515.
- Platt, J. (1999). Using analytic QP and sparseness to speed training of support vector machines. In *NIPS* (pp. 557–563).
- Schölkopf, B., & Smola, A. J. (2002). *Learning with kernels*. Cambridge: MIT Press.
- Smola, A., & Schölkopf, B. (2000). Sparse greedy matrix approximation for machine learning. In *ICML* (pp. 911–918).
- Steinwart, I. (2003). Sparseness of support vector machines. *JMLR*, 4, 1071–1105.
- Teo, C. H., Smola, A., Vishwanathan, S. V., & Le, Q. V. (2007). A scalable modular convex solver for regularized risk minimization. In *SIGKDD* (pp. 727–736).
- Tsang, I., Kwok, J., & Cheung, P. M. (2005). Core vector machines: Fast SVM training on very large data sets. *JMLR*, 6, 363–392.
- Tsang, I. W., Kocsor, A., & Kwok, J. T. (2007). Simpler core vector machines with enclosing balls. In *ICML* (pp. 911–918).
- Tsochantaridis, I., Joachims, T., Hofmann, T., & Altun, Y. (2005). Large margin methods for structured and interdependent output variables. *JMLR*, 6, 1453–1484.
- Vincent, P., & Bengio, Y. (2002). Kernel matching pursuit. *Machine Learning*, 48(1–3), 165–187.
- Williams, C., & Seeger, M. (2001). Using the Nystrom method to speed up kernel machines. In *NIPS*.
- Wu, M., Schölkopf, B., & Bakir, G. H. (2006). A direct method for building sparse kernel learning algorithms. *JMLR*, 7, 603–624.