

Bounded Kernel-Based Online Learning*

Francesco Orabona

Joseph Keshet[†]

Barbara Caputo

Idiap Research Institute

CH-1920 Martigny, Switzerland

FORABONA@IDIAP.CH

JKESHET@IDIAP.CH

BCAPUTO@IDIAP.CH

Editor: Yoav Freund

Abstract

A common problem of kernel-based online algorithms, such as the kernel-based Perceptron algorithm, is the amount of memory required to store the online hypothesis, which may increase without bound as the algorithm progresses. Furthermore, the computational load of such algorithms grows linearly with the amount of memory used to store the hypothesis. To attack these problems, most previous work has focused on discarding some of the instances, in order to keep the memory bounded. In this paper we present a new algorithm, in which the instances are not discarded, but are instead projected onto the space spanned by the previous online hypothesis. We call this algorithm Projectron. While the memory size of the Projectron solution cannot be predicted before training, we prove that its solution is guaranteed to be bounded. We derive a relative mistake bound for the proposed algorithm, and deduce from it a slightly different algorithm which outperforms the Perceptron. We call this second algorithm Projectron++. We show that this algorithm can be extended to handle the multiclass and the structured output settings, resulting, as far as we know, in the first online bounded algorithm that can learn complex classification tasks. The method of bounding the hypothesis representation can be applied to any conservative online algorithm and to other online algorithms, as it is demonstrated for ALMA₂. Experimental results on various data sets show the empirical advantage of our technique compared to various bounded online algorithms, both in terms of memory and accuracy.

Keywords: online learning, kernel methods, support vector machines, bounded support set

1. Introduction

Kernel-based discriminative online algorithms have been shown to perform very well on binary and multiclass classification problems (see, for example, Freund and Schapire, 1999; Crammer and Singer, 2003; Kivinen et al., 2004; Crammer et al., 2006). Each of these algorithms works in rounds, where at each round a new instance is provided. On rounds where the online algorithm makes a prediction mistake or when the confidence in the prediction is not sufficient, the algorithm adds the instance to a set of stored instances, called the *support set*. The online classification function is defined as a weighted sum of kernel combination of the instances in the support set. It is clear that if the problem is not linearly separable or the target hypothesis is changing over time, the classification function will never stop being updated, and consequently, the support set will grow unboundedly.

*. A preliminary version of this paper appeared in the Proceedings of the 25th International Conference on Machine Learning under the title "The Projectron: a Bounded Kernel-Based Perceptron".

†. Current affiliation: Toyota Technological Institute at Chicago, 6045 S Kenwood Ave., Chicago, IL 60637.

This leads, eventually, to a memory explosion, which limits the applicability of these algorithms for those tasks, such as autonomous agents, for example, where data must be acquired continuously over time.

Several authors have tried to address this problem, mainly by bounding a priori the size of the support set with a fixed value, called a *budget*. The first algorithm to overcome the unlimited growth of the support set was proposed by Crammer et al. (2003), and refined by Weston et al. (2005). In these algorithms, once the size of the support set reaches the budget, an instance from the support set that meets some criterion is removed, and replaced by the new instance. The strategy is purely heuristic and no mistake bound is given. A similar strategy is also used in NORMA (Kivinen et al., 2004) and SILK (Cheng et al., 2007). The very first online algorithm to have a fixed memory budget and a relative mistake bound is the Forgetron algorithm (Dekel et al., 2007). A stochastic algorithm that on average achieves similar performance to Forgetron, and with a similar mistake bound was proposed by Cesa-Bianchi et al. (2006). Unlike all previous work, the analysis presented in the last paper is within a probabilistic context, and all the bounds derived there are in expectation. A different approach to address this problem for online Gaussian processes is proposed in Csató and Opper (2002), where, in common with our approach, the instances are not discarded, but rather projected onto the space spanned by the instances in the support set. However, in that paper no mistake bound is derived and there is no use of the hinge loss, which often produces sparser solutions. Recent work by Langford et al. (2008) proposed a parameter that trades accuracy for sparseness in the weights of online learning algorithms. Nevertheless, this approach cannot induce sparsity in online algorithms with kernels.

In this paper we take a different route. While previous work focused on discarding some of the instances in order to keep the support set bounded, in this work the instances are not discarded. Either they are projected onto the space spanned by the support set, or they are added to the support set. By using this method, we show that the support set and, hence, the online hypothesis, is guaranteed to be bounded, although we cannot predict its size before training. Instead of using a budget parameter, representing the maximum size of the support set, we introduce a parameter trading accuracy for sparseness, depending on the needs of the task at hand. The main advantage of this setup is that by using all training samples, we are able to provide an online hypothesis with high online accuracy. Empirically, as suggested by the experiments, the output hypotheses are represented with relatively small number of instances, and have high accuracy.

We start with the most simple and intuitive kernel-based algorithm, namely the kernel-based Perceptron. We modify the Perceptron algorithm so that the number of stored samples needed to represent the online hypothesis is always bounded. We call this new algorithm *Projectron*. The empirical performance of the Projectron algorithm is on a par with the original Perceptron algorithm. We present a relative mistake bound for the Projectron algorithm, and deduce from it a new online bounded algorithm which outperforms the Perceptron algorithm, but still retains all of its advantages. We call this second algorithm *Projectron++*. We then extend Projectron++ to the more general cases of multiclass and structured output. As far as we know, this is the first bounded multiclass and structured output online algorithm, with a relative mistake bound.¹ Our technique for bounding the size of the support set can be applied to any conservative kernel-based online algorithm and to other online algorithms, as we demonstrate for ALMA₂ (Gentile, 2001). Finally, we present some experiments with common data sets, which suggest that Projectron is comparable to

1. Note that converting the budget algorithms presented by other authors, such as the Forgetron, to the multiclass or the structured output setting is not trivial, since these algorithms are inherently binary in nature.

Perceptron in performance, but it uses a much smaller support set. Moreover, experiments with Projectron++ shows that it outperforms all other bounded algorithms, while using the smallest support set. We also present experiments on the task of phoneme classification, which is considered to be difficult and naturally with a relatively very high number of support vectors. When comparing the Projectron++ algorithm to the Passive-Aggressive multiclass algorithm (Crammer et al., 2006), it turns out that the cumulative online error and the test error, after online-to-batch conversion, of both algorithms are comparable, although Projectron++ uses a smaller support set.

In summary, the contributions of this paper are (1) a new algorithm, called Projectron, which is derived from the kernel-based Perceptron algorithm, which empirically performs equally well, but has a bounded support set; (2) a relative mistake bound for this algorithm; (3) another algorithm, called Projectron++, based on the notion of large margin, which outperforms the Perceptron algorithm and the proposed Projectron algorithm; (4) the multiclass and structured output Projectron++ online algorithm with a bounded support set; and (5) an extension of our technique to other online algorithms, exemplified in this paper for ALMA₂.

The rest of the paper is organized as follows: in Section 2 we state the problem definition and the kernel-based Perceptron algorithm. Section 3 introduces Projectron, along with its theoretical analysis. Next, in Section 4 we derive Projectron++. Section 5 presents the multiclass and structured learning variant of Projectron++. In Section 6 we apply our technique for another kernel-based online algorithm, ALMA₂. Section 7 describes experimental results of the algorithms presented, on different data sets. Section 8 concludes the paper with a short discussion.

2. Problem Setting and the Kernel-Based Perceptron Algorithm

The basis of our study is the well known *Perceptron* algorithm (Rosenblatt, 1958; Freund and Schapire, 1999). The Perceptron algorithm learns the mapping $f : \mathcal{X} \rightarrow \mathbb{R}$ based on a set of examples $\mathcal{T} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots\}$, where $\mathbf{x}_t \in \mathcal{X}$ is called an *instance* and $y_t \in \{-1, +1\}$ is called a *label*. We denote the prediction of the Perceptron algorithm as $\text{sign}(f(\mathbf{x}))$ and we interpret $|f(\mathbf{x})|$ as the confidence in the prediction. We call the output f of the Perceptron algorithm a *hypothesis*, and we denote the set of all attainable hypotheses by \mathcal{H} . In this paper we assume that \mathcal{H} is a Reproducing Kernel Hilbert Space (RKHS) with a positive definite kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ implementing the inner product $\langle \cdot, \cdot \rangle$. The inner product is defined so that it satisfies the reproducing property, $\langle k(\mathbf{x}, \cdot), f(\cdot) \rangle = f(\mathbf{x})$.

The Perceptron algorithm is an online algorithm, where learning takes place in rounds. At each round a new hypothesis function is estimated, based on the previous one. We denote the hypothesis estimated after the t -th round by f_t . The algorithm starts with the zero hypothesis, $f_0 = \mathbf{0}$. At each round t , an instance $\mathbf{x}_t \in \mathcal{X}$ is presented to the algorithm, which predicts a label $\hat{y}_t \in \{-1, +1\}$ using the current function, $\hat{y}_t = \text{sign}(f_t(\mathbf{x}_t))$. Then, the correct label y_t is revealed. If the prediction \hat{y}_t differs from the correct label y_t , the hypothesis is updated as $f_t = f_{t-1} + y_t k(\mathbf{x}_t, \cdot)$, otherwise the hypothesis is left intact, $f_t = f_{t-1}$. The hypothesis f_t can be written as a kernel expansion according to the representer theorem (Schölkopf et al., 2001),

$$f_t(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathcal{S}_t} \alpha_i k(\mathbf{x}_i, \mathbf{x}), \quad (1)$$

where $\alpha_i = y_i$ and \mathcal{S}_t is defined to be the set of instances for which an update of the hypothesis occurred, that is, $\mathcal{S}_t = \{\mathbf{x}_i, 0 \leq i \leq t \mid \hat{y}_i \neq y_i\}$. The set \mathcal{S}_t is called the *support set*. The Perceptron algorithm is summarized in Figure 1.

```

Initialize:  $\mathcal{S}_0 = \emptyset, f_0 = \mathbf{0}$ 
For  $t = 1, 2, \dots$ 
  Receive new instance  $\mathbf{x}_t$ 
  Predict  $\hat{y}_t = \text{sign}(f_{t-1}(\mathbf{x}_t))$ 
  Receive label  $y_t$ 
  If  $y_t \neq \hat{y}_t$ 
     $f_t = f_{t-1} + y_t k(\mathbf{x}_t, \cdot)$  (update the hypothesis)
     $\mathcal{S}_t = \mathcal{S}_{t-1} \cup \mathbf{x}_t$  (add instance  $\mathbf{x}_t$  to the support set)
  Else
     $f_t = f_{t-1}$ 
     $\mathcal{S}_t = \mathcal{S}_{t-1}$ 

```

Figure 1: The kernel-based Perceptron Algorithm.

Although the Perceptron algorithm is very simple, it produces an online hypothesis with good performance. Our goal is to derive and analyze a new algorithm, which outputs a hypothesis that attains almost the same performance as the Perceptron hypothesis, but can be represented using many fewer instances, that is, an online hypothesis that is “close” to the Perceptron hypothesis but represented by a smaller support set. Recall that the hypothesis f_t is represented as a weighted sum over all the instances in the support set. The size of this representation is the cardinality of the support set, $|\mathcal{S}_t|$.

3. The Projectron Algorithm

This section starts by deriving the Projectron algorithm, motivated by an example of a finite dimensional kernel space. It continues with a description of how to calculate the projected hypothesis and describes some other computational aspects of the algorithm. The section concludes with a theoretical analysis of the algorithm.

3.1 Definition and Derivation

Let us first consider a finite dimensional RKHS \mathcal{H} induced by a kernel such as the polynomial kernel. Since \mathcal{H} is finite dimensional, there are a finite number of linearly independent hypotheses in this space. Hence, any hypothesis in this space can be expressed using a finite number of examples. We can modify the Perceptron algorithm to use only one set of independent instances as follows. On each round the algorithm receives an instance and predicts its label. On a prediction mistake, we check if the instance \mathbf{x}_t can be spanned by the support set, namely, for scalars $d_i \in \mathbb{R}, 1 \leq i \leq |\mathcal{S}_{t-1}|$, not all zeros, such that

$$k(\mathbf{x}_t, \cdot) = \sum_{\mathbf{x}_i \in \mathcal{S}_{t-1}} d_i k(\mathbf{x}_i, \cdot).$$

If we can find such scalars, the instance is not added to the support set, but instead, the coefficients $\{\alpha_i\}$ in the expansion Equation (1) are changed to reflect the addition of this instance to the support

set. Namely, for every i

$$\alpha_i = y_i + y_t d_i .$$

On the other hand, if the instance and the support set are linearly independent, the instance is added to the set with $\alpha_t = y_t$ as before. This technique reduces the size of the support set without changing the hypothesis. A similar approach was used by Downs et al. (2001) to simplify SVM solutions.

Let us consider now the more elaborate case of an infinite dimensional RKHS \mathcal{H} induced by a kernel such as the Gaussian kernel. In this case, it is not possible to find a finite number of linearly independent vectors which span the whole space, and hence there is no guarantee that the hypothesis can be expressed by a finite number of instances. However, we can approximate the concept of linear independence with a finite number of vectors (Csató and Opper, 2002; Engel et al., 2004; Orabona et al., 2007).

In particular, let us assume that at round t of the algorithm there is a prediction mistake, and that the mistaken instance \mathbf{x}_t should be added to the support set, \mathcal{S}_{t-1} . Let \mathcal{H}_{t-1} be an RKHS which is the span of the kernel images of the instances in the set \mathcal{S}_{t-1} . Formally,

$$\mathcal{H}_{t-1} = \text{span}(\{k(\mathbf{x}, \cdot) | \mathbf{x} \in \mathcal{S}_{t-1}\}) . \quad (2)$$

Before adding the instance to the support set, we construct two hypotheses: a *temporary hypothesis*, f'_t , using the function $k(\mathbf{x}_t, \cdot)$, that is, $f'_t = f_{t-1} + y_t k(\mathbf{x}_t, \cdot)$, and a *projected hypothesis*, f''_t , that is the projection of f'_t onto the space \mathcal{H}_{t-1} . That is, the projected hypothesis is that hypothesis from the space \mathcal{H}_{t-1} which is the *closest* to the temporary hypothesis. In a later section we will describe an efficient way to calculate the projected hypothesis. Denote by δ_t the distance between the hypotheses, $\delta_t = f''_t - f'_t$. If the norm of distance $\|\delta_t\|$ is below some threshold η , we use the projected hypothesis as our next hypothesis, that is, $f_t = f''_t$, otherwise we use the temporary hypothesis as our next hypothesis, that is, $f_t = f'_t$. As we show in the following theorem, this strategy assures that the maximum size of the support set is always finite, regardless of the dimension of the RKHS \mathcal{H} . Guided by these considerations we can design a new Perceptron-like algorithm that projects the solution onto the space spanned by the previous support vectors whenever possible. We call this algorithm *Projectron*. The algorithm is given in Figure 2.

The parameter η plays an important role in our algorithm. If η is equal to zero, we obtain exactly the same solution as the Perceptron algorithm. In this case, however, the Projectron solution can still be sparser when some of the instances are linearly dependent or when the kernel induces a finite dimensional RKHS \mathcal{H} . If η is greater than zero we trade precision for sparseness. Moreover, as shown in the next section, this implies a bounded algorithmic complexity, namely, the memory and time requirements for each step are bounded. We analyze the effect of η on the classification accuracy in Subsection 3.3.

3.2 Practical Considerations

We now consider the problem of deriving the projected hypothesis f''_t in a Hilbert space \mathcal{H} , induced by a kernel function $k(\cdot, \cdot)$. Recall that f'_t is defined as $f'_t = f_t + y_t k(\mathbf{x}_t, \cdot)$. Denote by $P_{t-1} f'_t$ the projection of $f'_t \in \mathcal{H}$ onto the subspace $\mathcal{H}_{t-1} \subseteq \mathcal{H}$. The projected hypothesis f''_t is defined as $f''_t = P_{t-1} f'_t$. Schematically, this is depicted in Figure 3.

Expanding f'_t we have

$$f''_t = P_{t-1} f'_t = P_{t-1} (f_{t-1} + y_t k(\mathbf{x}_t, \cdot)) .$$

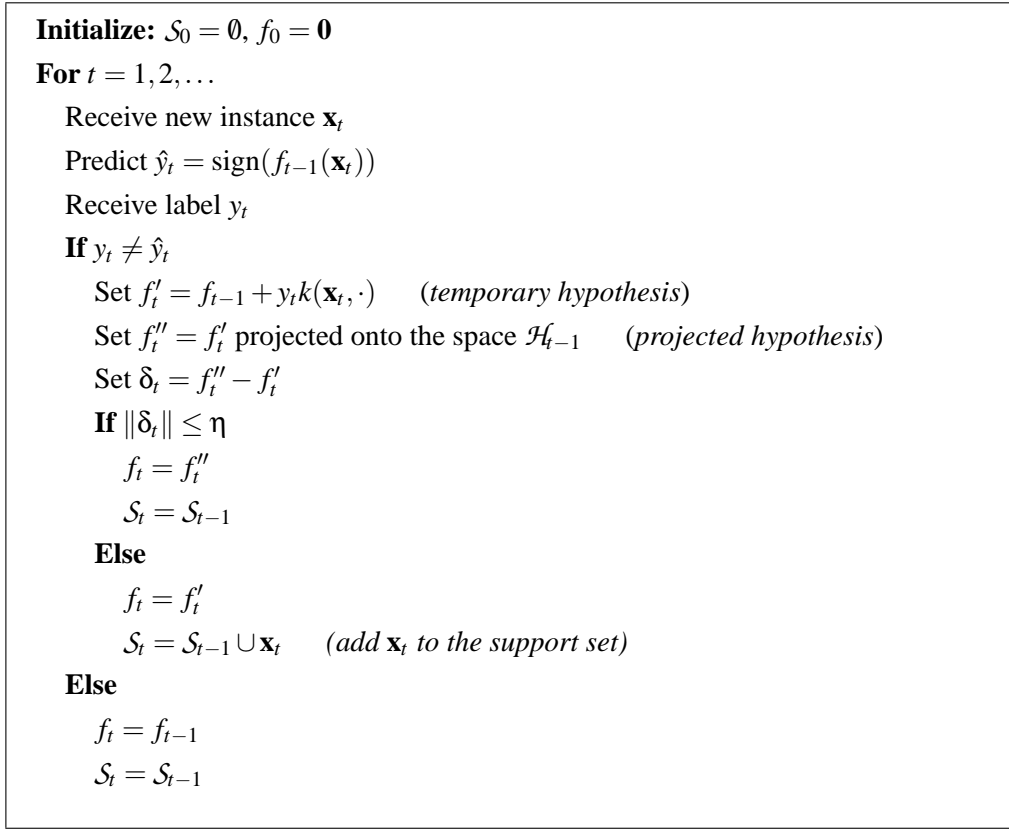


Figure 2: The Projectron Algorithm.

The projection is a linear operator, hence

$$f''_t = f_{t-1} + y_t P_{t-1} k(\mathbf{x}_t, \cdot). \quad (3)$$

Recall that $\delta_t = f''_t - f'_t$. By substituting f''_t from Equation (3) and f'_t we have

$$\delta_t = f''_t - f'_t = y_t P_{t-1} k(\mathbf{x}_t, \cdot) - y_t k(\mathbf{x}_t, \cdot). \quad (4)$$

The projection of $f'_t \in \mathcal{H}$ onto a subspace $\mathcal{H}_{t-1} \subset \mathcal{H}$ is defined as the hypothesis in \mathcal{H}_{t-1} closest to f'_t . Hence, let $\sum_{\mathbf{x}_j \in \mathcal{S}_{t-1}} d_j k(\mathbf{x}_j, \cdot)$ be an hypothesis in \mathcal{H}_{t-1} , where $\mathbf{d} = (d_1, \dots, d_{|\mathcal{S}_{t-1}|})$ is a set of coefficients, with $d_i \in \mathbb{R}$. The closest hypothesis is the one for which it holds that

$$\|\delta_t\|^2 = \min_{\mathbf{d}} \left\| \sum_{\mathbf{x}_j \in \mathcal{S}_{t-1}} d_j k(\mathbf{x}_j, \cdot) - k(\mathbf{x}_t, \cdot) \right\|^2. \quad (5)$$

Expanding Equation (5) we get

$$\|\delta_t\|^2 = \min_{\mathbf{d}} \left(\sum_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{S}_{t-1}} d_j d_i k(\mathbf{x}_j, \mathbf{x}_i) - 2 \sum_{\mathbf{x}_j \in \mathcal{S}_{t-1}} d_j k(\mathbf{x}_j, \mathbf{x}_t) + k(\mathbf{x}_t, \mathbf{x}_t) \right).$$

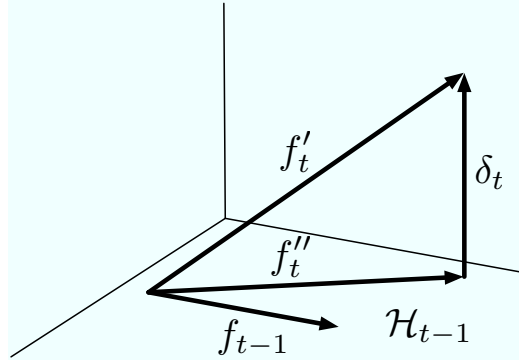


Figure 3: Geometrical interpretation of the projection of the hypothesis f_t'' onto the subspace \mathcal{H}_{t-1} .

Let us define $\mathbf{K}_{t-1} \in \mathbb{R}^{t-1 \times t-1}$ to be the matrix generated by the instances in the support set \mathcal{S}_{t-1} , that is, $\{\mathbf{K}_{t-1}\}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ for every $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{S}_{t-1}$. Let us also define $\mathbf{k}_t \in \mathbb{R}^{t-1}$ to be the vector whose i -th element is $k_{t_i} = k(\mathbf{x}_i, \mathbf{x}_t)$. We have

$$\|\delta_t\|^2 = \min_{\mathbf{d}} (\mathbf{d}^T \mathbf{K}_{t-1} \mathbf{d} - 2\mathbf{d}^T \mathbf{k}_t + k(\mathbf{x}_t, \mathbf{x}_t)) . \quad (6)$$

Solving Equation (6), that is, applying the extremum conditions with respect to \mathbf{d} , we obtain

$$\mathbf{d}^* = \mathbf{K}_{t-1}^{-1} \mathbf{k}_t \quad (7)$$

and, by substituting Equation (7) into Equation (6),

$$\|\delta_t\|^2 = k(\mathbf{x}_t, \mathbf{x}_t) - \mathbf{k}_t^T \mathbf{d}^* . \quad (8)$$

Furthermore, by substituting Equation (7) back into Equation (3) we get

$$f_t'' = f_{t-1} + y_t \sum_{\mathbf{x}_j \in \mathcal{S}_{t-1}} \mathbf{d}_j^* k(\mathbf{x}_j, \cdot) . \quad (9)$$

We have shown how to calculate both the distance δ_t and the projected hypothesis f_t'' . In summary, one needs to compute \mathbf{d}^* according to Equation (7), and plug the result either into Equation (8) to obtain δ_t , or into Equation (9) to obtain the projected hypothesis.

In order to make the computation more tractable, we need an efficient method to calculate the matrix inversion \mathbf{K}_t^{-1} iteratively. The first method, used by Cauwenberghs and Poggio (2000) for incremental training of SVMs, directly updates the inverse matrix. An efficient way to do this, exploiting the incremental nature of the approach, is to recursively update the inverse matrix. Using the matrix inversion lemma it is possible to show (see, e.g., Csató and Opper, 2002) that after the addition of a new sample, \mathbf{K}_t^{-1} becomes

$$\mathbf{K}_t^{-1} = \begin{bmatrix} & & & 0 \\ & \mathbf{K}_{t-1}^{-1} & & \vdots \\ & & & 0 \\ 0 & \dots & 0 & 0 \end{bmatrix} + \frac{1}{\|\delta_t\|^2} \begin{bmatrix} \mathbf{d}^* \\ -1 \end{bmatrix} \begin{bmatrix} \mathbf{d}^{*T} & -1 \end{bmatrix}$$

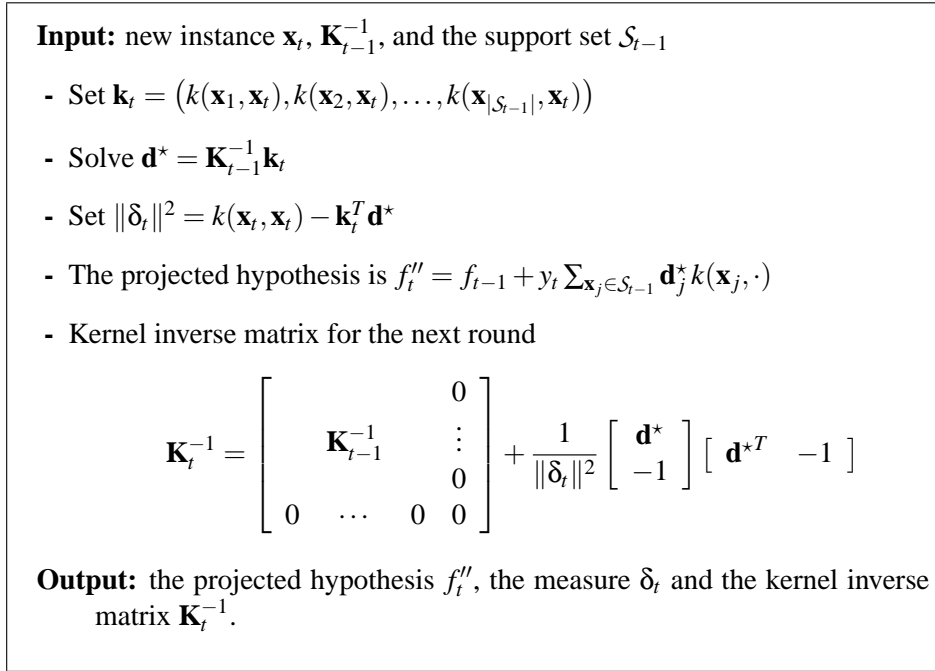


Figure 4: Calculation of the projected hypothesis f_t'' .

where \mathbf{d}^* and $\|\delta_t\|^2$ are already evaluated during the previous steps of the algorithm, as given by Equation (7) and Equation (8). Thanks to this incremental evaluation, the time complexity of the linear independence check is $O(|\mathcal{S}_{t-1}|^2)$, as one can easily see from Equation (7). Note that the matrix \mathbf{K}_{t-1} can be safely inverted since, by incremental construction, it is always full-rank.

An alternative way to derive the inverse matrix is to use the Cholesky decomposition of \mathbf{K}_{t-1} and to update it recursively. This is known to be numerically more stable than directly updating the inverse. In our experiments, however, we found out that the method presented here is as stable as the Cholesky decomposition.

Overall, the time complexity of the algorithm is $O(|\mathcal{S}_t|^2)$, as described above, and the space complexity is $O(|\mathcal{S}_t|^2)$, due to the storage of the matrix \mathbf{K}_t^{-1} , similar to the second-order Perceptron algorithm (Cesa-Bianchi et al., 2005). A summary of the derivation of f_t'' , the projection of f_t' onto the space spanned by \mathcal{S}_{t-1} , is described in Figure 4.

3.3 Analysis

We now analyze the theoretical aspects of the proposed algorithm. First, we present a theorem which states that the size of the support set of the Projectron algorithm is bounded.

Theorem 1 *Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a continuous Mercer kernel, with \mathcal{X} a compact subset of a Banach space. Then, for any training sequence $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}, i = 1, 2, \dots$ and for any $\eta > 0$, the size of the support set of the Projectron algorithm is finite.*

Proof The proof of this theorem follows the same lines as the proof of Theorem 3.1 in Engel et al. (2004). From the Mercer theorem it follows that there exists a mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}'$, where \mathcal{H}' is

an Hilbert space, $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ and ϕ is continuous. Given that ϕ is continuous and that \mathcal{X} is compact, we obtain that $\phi(\mathcal{X})$ is compact. From the definition of δ_t in Equation (5) we get that every time a new basis vector is added we have

$$\begin{aligned} \eta^2 \leq \|\delta_t\|^2 &= \min_{\mathbf{d}} \left\| \sum_{\mathbf{x}_j \in \mathcal{S}_{t-1}} d_j k(\mathbf{x}_j, \cdot) - k(\mathbf{x}_t, \cdot) \right\|^2 \leq \min_{d_j} \|d_j k(\mathbf{x}_j, \cdot) - k(\mathbf{x}_t, \cdot)\|^2 \\ &= \min_{d_j} \|d_j \phi(\mathbf{x}_j) - \phi(\mathbf{x}_t)\|^2 \leq \|\phi(\mathbf{x}_j) - \phi(\mathbf{x}_t)\|^2 \end{aligned}$$

for any $1 \leq j \leq |\mathcal{S}_{t-1}|$. Hence from the definition of packing numbers (Cucker and Zhou, 2007, Definition 5.17), we get that the maximum size of the support set in the Projectron algorithm is bounded by the packing number at scale η of $\phi(\mathcal{X})$. This number, in turn, is bounded by the covering number at scale $\eta/2$, and it is finite because the set is compact (Cucker and Zhou, 2007, Proposition 5.18). \blacksquare

Note that this theorem guarantees that the size of the support set is bounded, however it does not state that the size of the support set is fixed or that it can be estimated before training.

The next theorem provides a mistake bound. The main idea is to bound the maximum number of mistakes of the algorithm, relative to any hypothesis $g \in \mathcal{H}$, even chosen in hindsight. First, we define the loss with a margin $\gamma \in \mathbb{R}$ of the hypothesis g on the example (\mathbf{x}_t, y_t) as

$$\ell_\gamma(g(\mathbf{x}_t), y_t) = \max\{0, \gamma - y_t g(\mathbf{x}_t)\}, \quad (10)$$

and we define the cumulative loss, D_γ , of g on the first T examples as

$$D_\gamma = \sum_{t=1}^T \ell_\gamma(g(\mathbf{x}_t), y_t).$$

Before stating the bound, we present a lemma that will be used in the rest of our proofs. We will use its first statement to bound the scalar product between a projected sample and the competitor, and its second statement to derive the scalar product between the current hypothesis and the projected sample.

Lemma 2 *Let $(\hat{\mathbf{x}}, \hat{y})$ be an example, with $\hat{\mathbf{x}} \in \mathcal{X}$ and $\hat{y} \in \{+1, -1\}$. If we denote by $f(\cdot)$ an hypothesis in \mathcal{H} , and denote by $q(\cdot)$ any function in \mathcal{H} , then the following holds*

$$\hat{y} \langle f, q \rangle \geq \gamma - \ell_\gamma(f(\hat{\mathbf{x}}), \hat{y}) - \|f\| \cdot \|q - k(\hat{\mathbf{x}}, \cdot)\|.$$

Moreover, if $f(\cdot)$ can be written as $\sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \cdot)$ with $\alpha_i \in \mathbb{R}$ and $\mathbf{x}_i \in \mathcal{X}, i = 1, \dots, m$, and $q(\cdot)$ is the projection of $k(\hat{\mathbf{x}}, \cdot)$ onto the space spanned by $k(\mathbf{x}_i, \cdot), i = 1, \dots, m$, then

$$\hat{y} \langle f, q \rangle = \hat{y} f(\hat{\mathbf{x}}).$$

Proof The first inequality comes from an application of the Cauchy-Schwarz inequality and the definition of the hinge loss in Equation (10). The second equality follows from the fact that $\langle f, q - k(\hat{\mathbf{x}}, \cdot) \rangle = 0$, because $f(\cdot)$ is orthogonal to the difference between $k(\hat{\mathbf{x}}, \cdot)$ and its projection onto the space in which $f(\cdot)$ lives. \blacksquare

With these definitions at hand, we can state the following bound for Projectron.

Theorem 3 *Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ be a sequence of instance-label pairs where $\mathbf{x}_t \in \mathcal{X}$, $y_t \in \{-1, +1\}$, and $k(\mathbf{x}_t, \mathbf{x}_t) \leq 1$ for all t . Assume that the Projectron algorithm is run with $\eta \geq 0$. Then the number of prediction mistakes it makes on the sequence is bounded by*

$$\frac{\|g\|^2}{(1-\eta\|g\|)^2} + \frac{D_1}{1-\eta\|g\|} + \frac{\|g\|}{1-\eta\|g\|} \sqrt{\frac{D_1}{1-\eta\|g\|}}$$

where g is an arbitrary function in \mathcal{H} , such that $\|g\| < \frac{1}{\eta}$.

Proof Define the relative progress in each round as $\Delta_t = \|f_{t-1} - \lambda g\|^2 - \|f_t - \lambda g\|^2$, where λ is a positive scalar to be optimized. We bound the progress from above and below, as in Gentile (2003). On rounds where there is no mistake, Δ_t equals 0. On rounds where there is a mistake there are two possible updates: either $f_t = f_{t-1} + y_t P_{t-1} k(\mathbf{x}_t, \cdot)$ or $f_t = f_{t-1} + y_t k(\mathbf{x}_t, \cdot)$. In the following we start bounding the progress from below, when the update is of the former type. In particular we set $q(\cdot) = P_{t-1} k(\mathbf{x}_t, \cdot)$ in Lemma 2 and use $\delta_t = y_t P_{t-1} k(\mathbf{x}_t, \cdot) - y_t k(\mathbf{x}_t, \cdot)$ from Equation (4). Let τ_t be an indicator function for a mistake on the t -th round, that is, τ_t is 1 if there is a mistake on round t and 0 otherwise. We have

$$\begin{aligned} \Delta_t &= \|f_{t-1} - \lambda g\|^2 - \|f_t - \lambda g\|^2 = 2\tau_t y_t \langle \lambda g - f_{t-1}, P_{t-1} k(\mathbf{x}_t, \cdot) \rangle - \tau_t^2 \|P_{t-1} k(\mathbf{x}_t, \cdot)\|^2 \\ &\geq \tau_t \left(2\lambda - 2\lambda \ell_1(g(\mathbf{x}_t), y_t) - \tau_t \|P_{t-1} k(\mathbf{x}_t, \cdot)\|^2 - 2\lambda \|g\| \cdot \|\delta_t\| - 2y_t f_{t-1}(\mathbf{x}_t) \right). \end{aligned} \quad (11)$$

Moreover, on every projection update $\|\delta_t\| \leq \eta$, and $\|P_{t-1} k(\mathbf{x}_t, \cdot)\| \leq 1$ by the theorem's assumption, so we have

$$\Delta_t \geq \tau_t \left(2\lambda - 2\lambda \ell_1(g(\mathbf{x}_t), y_t) - \tau_t - 2\eta \lambda \|g\| - 2y_t f_{t-1}(\mathbf{x}_t) \right).$$

We can further bound Δ_t by noting that on every prediction mistake $y_t f_{t-1}(\mathbf{x}_t) \leq 0$. Overall we have

$$\|f_{t-1} - \lambda g\|^2 - \|f_t - \lambda g\|^2 \geq \tau_t \left(2\lambda - 2\lambda \ell_1(g(\mathbf{x}_t), y_t) - \tau_t - 2\eta \lambda \|g\| \right). \quad (12)$$

When there is an update without projection, similar reasoning yields that

$$\|f_{t-1} - \lambda g\|^2 - \|f_t - \lambda g\|^2 \geq \tau_t \left(2\lambda - 2\lambda \ell_1(g(\mathbf{x}_t), y_t) - \tau_t \right),$$

hence the bound in Equation (12) holds in both cases.

We sum over t on both sides, remembering that τ_t can be upper bounded by 1. The left hand side of the equation is a telescoping sum, hence it collapses to $\|f_0 - \lambda g\|^2 - \|f_T - \lambda g\|^2$, which can be upper bounded by $\lambda^2 \|g\|^2$, using the fact that $f_0 = \mathbf{0}$ and that $\|f_T - \lambda g\|^2$ is non-negative. Finally, we have

$$\lambda^2 \|g\|^2 + 2\lambda D_1 \geq M(2\lambda - 2\eta \lambda \|g\| - 1), \quad (13)$$

where M is the number of mistakes. The last equation implies a bound on M for any choice of $\lambda > 0$, hence we can take the minimum of these bounds. From now on we can suppose that $M > \frac{D_1}{1-\eta\|g\|}$. In fact, if $M \leq \frac{D_1}{1-\eta\|g\|}$ then the theorem trivially holds. The minimum of Equation (13) as a function of λ occurs at

$$\lambda^* = \frac{M(1-\eta\|g\|) - D_1}{\|g\|^2}.$$

By our hypothesis that $M > \frac{D_1}{1-\eta\|g\|}$ we have that λ^* is positive. Substituting λ^* into Equation (13) we obtain

$$\frac{(D_1 - M(1 - \eta\|g\|))^2}{\|g\|^2} - M \leq 0.$$

Solving for M and overapproximating concludes the proof. \blacksquare

This theorem suggests that the performance of the Projectron algorithm is slightly worse than that of the Perceptron algorithm. Specifically, if we set $\eta = 0$, we recover the best known bound for the Perceptron algorithm (see for example Gentile, 2003). Hence the degradation in the performance of Projectron compared to Perceptron is related to $\frac{1}{1-\eta\|g\|}$. Empirically, the Projectron algorithm and the Perceptron algorithm perform similarly, for a wide range of settings of η .

4. The Projectron++ Algorithm

The proof of Theorem 3 suggests how to improve the Projectron algorithm to improve upon the performance of the Perceptron algorithm, while maintaining a bounded support set. We can change the Projectron algorithm so that an update takes place not only if there is a prediction mistake, but also when the confidence of the prediction is low. We refer to this latter case as a *margin error*, that is, $0 < y_t f_{t-1}(\mathbf{x}_t) < 1$. This strategy is known to improve the classification rate but also increases the size of the support set (Crammer et al., 2006). A possible solution to this obstacle is not to update on every round in which a margin error occurs, but only when there is a margin error and the new instance *can be projected* onto the support set. Hence, the update on round in which there is a margin error would in general be of the form

$$f_t = f_{t-1} + y_t \tau_t P_{t-1} k(\mathbf{x}_t, \cdot),$$

with $0 < \tau_t \leq 1$. The last constraint comes from the proof of Theorem 3, where we upper bound τ_t by 1. Note that setting τ_t to 0 is equivalent to leaving the hypothesis unchanged.

In particular, disregarding the loss term in Equation (11), the progress Δ_t can be made positive with an appropriate choice of τ_t . Whenever this progress is non-negative the worst-case number of mistakes decreases, hopefully along with the classification error rate of the algorithm. With this modification we expect better performance, that is, fewer mistakes, but without any increase of the support set size. We can even expect solutions with a smaller support set, since new instances can be added to the support set only if misclassified, hence having fewer mistakes should result in a smaller support set. We name this algorithm *Projectron++*. The following theorem states a mistake bound for Projectron++, and guides us in how to choose τ_t .

Theorem 4 *Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ be a sequence of instance-label pairs where $\mathbf{x}_t \in \mathcal{X}$, $y_t \in \{-1, +1\}$, and $k(\mathbf{x}_t, \mathbf{x}_t) \leq 1$ for all t . Assume that Projectron++ is run with $\eta > 0$. Then the number of prediction mistakes it makes on the sequence is bounded by*

$$\frac{\|g\|^2}{(1 - \eta\|g\|)^2} + \frac{D_1}{1 - \eta\|g\|} + \frac{\|g\|}{1 - \eta\|g\|} \sqrt{\max\left(0, \frac{D_1}{1 - \eta\|g\|} - B\right)}$$

where g is an arbitrary function in \mathcal{H} , such that $\|g\| < \frac{1}{\eta}$,

$$0 < \tau_t < \min \left\{ 2 \frac{\ell_1(f_{t-1}(\mathbf{x}_t), y_t) - \frac{\|\delta_t\|}{\eta}}{\|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2}, 1 \right\}$$

and

$$B = \sum_{\{t: 0 < y_t, f_{t-1}(\mathbf{x}_t) < 1\}} \tau_t \left(2\ell_1(f_{t-1}(\mathbf{x}_t), y_t) - \tau_t \|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2 - 2\frac{\|\delta_t\|}{\eta} \right) > 0.$$

Proof The proof is similar to the proof of Theorem 3, where the difference is that during rounds in which there is a margin error we update the solution whenever it is possible to project ensuring an improvement of the mistake bound. Assume that $\lambda \geq 1$. On rounds when a margin error occurs, as in Equation (11), we can write

$$\begin{aligned} \Delta_t + 2\tau_t \lambda \ell_1(g(\mathbf{x}_t), y_t) &\geq \tau_t (2\lambda - \tau_t \|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2 - 2\lambda \|\delta_t\| \cdot \|g\| - 2y_t f_{t-1}(\mathbf{x}_t)) \\ &> \tau_t \left(2 \left(1 - \frac{\|\delta_t\|}{\eta} \right) - \tau_t \|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2 - 2y_t f_{t-1}(\mathbf{x}_t) \right) \\ &= \tau_t \left(2\ell_1(f_{t-1}(\mathbf{x}_t), y_t) - \tau_t \|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2 - 2\frac{\|\delta_t\|}{\eta} \right), \end{aligned} \tag{14}$$

where we used the bounds on $\|g\|$ and λ . Let β_t be the right hand-side of Equation (14). A sufficient condition to have β_t positive is

$$\tau_t < 2 \frac{\ell_1(f_{t-1}(\mathbf{x}_t), y_t) - \frac{\|\delta_t\|}{\eta}}{\|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2}.$$

Constraining τ_t to be less than or equal to 1 yields the update rule in the theorem.

Let $B = \sum_{\{t: 0 < y_t, f_{t-1}(\mathbf{x}_t) < 1\}} \beta_t$. Similarly to the proof of Theorem 3, we have

$$\lambda^2 \|g\|^2 + 2\lambda D_1 \geq M(2\lambda - 2\eta\lambda \|g\| - 1) + B. \tag{15}$$

Again, the optimal value of λ is

$$\lambda^* = \frac{M(1 - \eta \|g\|) - D_1}{\|g\|^2}.$$

We can assume that $M(1 - \eta \|g\|) - D_1 \geq \|g\|^2$. In fact, if $M < \frac{\|g\|^2 + D_1}{1 - \eta \|g\|}$, then the theorem trivially holds. With this assumption, λ^* is positive and greater than or equal to 1, satisfying our initial constraint on λ . Substituting this optimal value of λ into Equation (15), we have

$$\frac{(D_1 - M(1 - \eta \|g\|))^2}{\|g\|^2} - M + B \leq 0.$$

Solving for M concludes the proof. ■

The proof technique presented here is very general, in particular it can be applied to the Passive-Aggressive algorithm PA-I (Crammer et al., 2006). In fact, removing the projection step and updating on rounds in which there is a margin error, with $f_t = f_{t-1} + y_t \tau_t k(\mathbf{x}_t, \cdot)$, we end up with the condition $0 < \tau_t < \min \left\{ 2 \frac{\ell_1(f_{t-1}(\mathbf{x}_t), y_t)}{\|k(\mathbf{x}_t, \cdot)\|^2}, 1 \right\}$. This rule generalizes the PA-I bound whenever $R = 1$ and $C = 1$, however the obtained bound substantially improves upon the original bound in Crammer et al. (2006).

The theorem gives us some freedom for the choice of τ_t . Experimentally we have observed that we obtain the best performance if the update is done with the following rule

$$\tau_t = \min \left\{ \frac{\ell_1(f_{t-1}(\mathbf{x}_t), y_t)}{\|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2}, 2 \frac{\ell_1(f_{t-1}(\mathbf{x}_t), y_t) - \frac{\|\delta_t\|}{\eta}}{\|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2}, 1 \right\}.$$

The added term in the minimum comes from ignoring the term $-2 \frac{\|\delta_t\|}{\eta}$ and in finding the maximum of the quadratic equation. Notice that the term $\|P_{t-1}k(\mathbf{x}_t, \cdot)\|^2$ in the last equation can be practically computed as $\mathbf{k}_t^T \mathbf{d}^*$, as can be derived using the same techniques presented in Subsection 3.2.

We note in passing that the condition on whether \mathbf{x}_t can be projected onto \mathcal{H}_{t-1} on margin error may be stated as $\ell_1(f_{t-1}(\mathbf{x}_t), y_t) \geq \frac{\|\delta_t\|}{\eta}$. This means that if the loss is relatively large, the progress is also large and the algorithm can afford “wasting” a bit of it for the sake of projecting.

The algorithm is summarized in Figure 2. The performance of the Projectron++ algorithm, the Projectron algorithm and several other bounded online algorithms are compared and reported in Section 7.

5. Extension to Multiclass and Structured Output

In this section we extend Projectron++ to the multiclass and the structured output settings (note that Projectron can be generalized in a similar way). We start by presenting the more complex decision problem, namely the structured output, and then we derive the multiclass decision problem as a special case.

In structured output decision problems the set of possible labels has a unique and defined structure, such as a tree, a graph or a sequence (Collins, 2000; Taskar et al., 2003; Tsochantaridis et al., 2004). Denote the set of all labels as $\mathcal{Y} = \{1, \dots, k\}$. Each instance is associated with a label from \mathcal{Y} . Generally, in structured output problems there may be dependencies between the instance and the label, as well as between labels. Hence, to capture these dependencies, the input and the output pairs are represented in a common feature representation. The learning task is therefore defined as finding a function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ such that

$$y_t = \arg \max_{y \in \mathcal{Y}} f(\mathbf{x}_t, y). \quad (16)$$

Let us generalize the definition of the RKHS \mathcal{H} introduced in Section 2 to the case of structured learning. A kernel function in this setting should reflect the dependencies between the instances and the labels, hence we define the structured kernel function as a function on the domain of the instances and the labels, namely, $k^S : (\mathcal{X} \times \mathcal{Y})^2 \rightarrow \mathbb{R}$. This kernel function induces the RKHS \mathcal{H}^S , where the inner product in this space is defined such that it satisfies the reproducing property, $\langle k^S((\mathbf{x}, y), \cdot), f \rangle = f(\mathbf{x}, y)$.

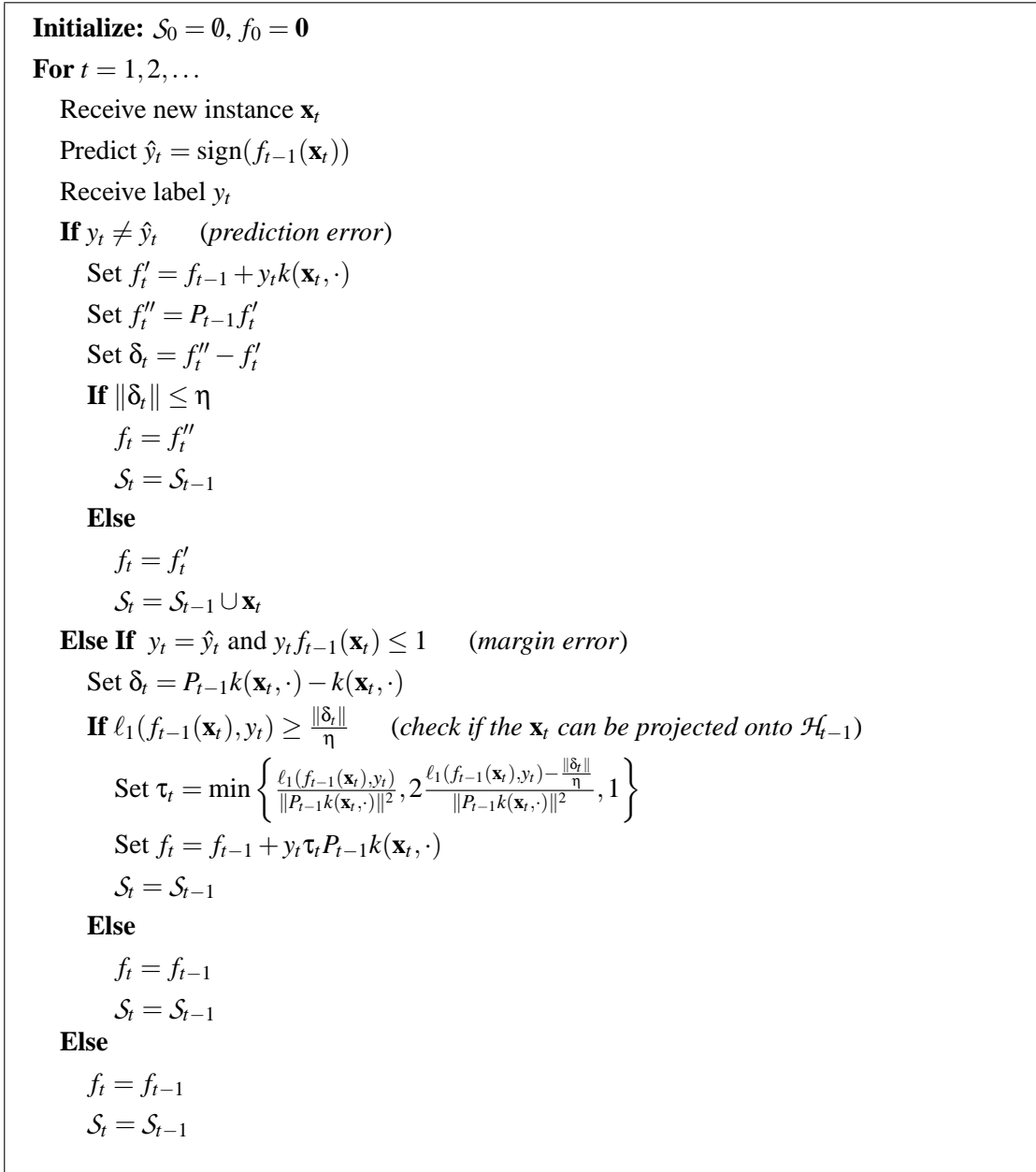


Figure 5: The Projectron++ Algorithm.

As in the binary classification algorithm presented earlier, the structured output online algorithm receives instances in a sequential order. Upon receiving an instance, $\mathbf{x}_t \in \mathcal{X}$, the algorithm predicts a label, y'_t , according to Equation (16). After making its prediction, the algorithm receives the correct label, y_t . We define the loss suffered by the algorithm on round t for the example (\mathbf{x}_t, y_t) as

$$\ell_\gamma^S(f, \mathbf{x}_t, y_t) = \max\{0, \gamma - f(\mathbf{x}_t, y_t) + \max_{y'_t \neq y_t} f(\mathbf{x}_t, y'_t)\},$$

and the cumulative loss D_γ^S as

$$D_\gamma^S = \sum_{t=1}^T \ell_\gamma^S(f, \mathbf{x}_t, y_t).$$

Note that sometimes it is useful to define γ as a function $\gamma: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ describing the discrepancy between the predicted label and the true label. Our algorithm can handle such a label cost function, but we will not discuss this issue here (see Crammer et al., 2006, for further details).

As in the binary case, on rounds in which there is a prediction mistake, $y'_t \neq y_t$, the algorithm updates the hypothesis f_{t-1} by adding $k((\mathbf{x}_t, y_t), \cdot) - k((\mathbf{x}_t, y'_t), \cdot)$ or its projection. When there is a margin mistake, $0 < \ell_\gamma^S(f_{t-1}, \mathbf{x}_t, y_t) < \gamma$, the algorithm updates the hypothesis f_{t-1} by adding $\tau_t P_{t-1}(k((\mathbf{x}_t, y_t), \cdot) - k((\mathbf{x}_t, y'_t), \cdot))$, where $0 < \tau_t < 1$ and will be defined shortly. Now, for the structured output case, δ_t is defined as

$$\delta_t = k((\mathbf{x}_t, y_t), \cdot) - k((\mathbf{x}_t, y'_t), \cdot) - P_{t-1}(k((\mathbf{x}_t, y_t), \cdot) - k((\mathbf{x}_t, y'_t), \cdot)).$$

The analysis of the structured output Projectron++ algorithm is similar to that provided for the binary case. We can easily obtain the generalization of Lemma 2 and Theorem 4 as follows

Lemma 5 *Let $(\hat{\mathbf{x}}, \hat{y})$ be an example, with $\hat{\mathbf{x}} \in \mathcal{X}$ and $\hat{y} \in \mathcal{Y}$. Denote by $f(\cdot)$ an hypothesis in \mathcal{H}^S . Let $q(\cdot) \in \mathcal{H}^S$. Then the following holds for any $y' \in \mathcal{Y}$:*

$$\langle f, q \rangle \geq \gamma - \ell_\gamma^S(f, \hat{\mathbf{x}}, \hat{y}) - \|f\| \cdot \|q - (k((\hat{\mathbf{x}}, \hat{y}), \cdot) - k((\hat{\mathbf{x}}, y'), \cdot))\|.$$

Moreover if $f(\cdot)$ can be written as $\sum_{i=1}^m \alpha_i k((\mathbf{x}_i, y_i), \cdot)$ with $\alpha_i \in \mathbb{R}$ and $\mathbf{x}_i \in \mathcal{X}$, $i = 1, \dots, m$, and q is the projection of $k((\hat{\mathbf{x}}, \hat{y}), \cdot) - k((\hat{\mathbf{x}}, y'), \cdot)$ in the space spanned by $k((\mathbf{x}_i, y_i), \cdot)$, $i = 1, \dots, m$, we have that

$$\langle f, q \rangle = f(\hat{\mathbf{x}}, \hat{y}) - f(\hat{\mathbf{x}}, y').$$

Theorem 6 *Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ be a sequence of instance-label pairs where $\mathbf{x}_t \in \mathcal{X}$, $y_t \in \mathcal{Y}$, and $\|k((\mathbf{x}_t, y), \cdot)\| \leq 1/2$ for all t and $y \in \mathcal{Y}$. Assuming that Projectron++ is run with $\eta > 0$, the number of prediction mistakes it makes on the sequence is bounded by*

$$\frac{\|g\|^2}{(1 - \eta\|g\|)^2} + \frac{D_1^S}{1 - \eta\|g\|} + \frac{\|g\|}{1 - \eta\|g\|} \sqrt{\max\left(0, \frac{D_1^S}{1 - \eta\|g\|} - B\right)}$$

where g is an arbitrary function in \mathcal{H}^S , such that $\|g\| < \frac{1}{\eta}$,

$$\begin{aligned} a &= P_{t-1}(k((\mathbf{x}_t, y_t), \cdot) - k((\mathbf{x}_t, y'_t), \cdot)) \\ 0 < \tau_t &< \min \left\{ 2 \frac{\ell_1^S(f_{t-1}, \mathbf{x}_t, y_t) - \frac{\|\delta_t\|}{\eta}}{\|a\|^2}, 1 \right\} \\ B &= \sum_{\{t: 0 < \ell_1^S(f_{t-1}, \mathbf{x}_t, y_t) < 1\}} \tau_t \left(2\ell_1^S(f_{t-1}, \mathbf{x}_t, y_t) - \tau_t \|a\|^2 - 2\frac{\|\delta_t\|}{\eta} \right) > 0. \end{aligned}$$

As in Theorem 4 there is some freedom in the choice of τ_t , and again we set it to

$$\tau_t = \min \left\{ \frac{\ell_1^S(f_{t-1}, \mathbf{x}_t, y_t)}{\|a\|^2}, 2 \frac{\ell_1^S(f_{t-1}, \mathbf{x}_t, y_t) - \frac{\|\delta_t\|}{\eta}}{\|a\|^2}, 1 \right\}.$$

In the multiclass decision problem case, the kernel $k((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2))$ is simplified to $\delta_{y_1 y_2} k(\mathbf{x}_1, \mathbf{x}_2)$, where $\delta_{y_1 y_2}$ is the Kronecker delta. This corresponds to the use of a different prototype for each class. This simplifies the projection step, in fact $k((\mathbf{x}_t, y_t), \cdot)$ can be projected only onto the functions in S_{t-1} belonging to y_t , the scalar product with the other functions being zero. So instead of storing a single matrix \mathbf{K}_{t-1}^{-1} , we need to store m matrices, where m is the number of classes, each one being the inverse matrix of the Gram matrix of the functions of one class. This results in improvements in both memory usage and computational cost of the algorithm. To see this suppose that we have m classes, each with n vectors in the support set. Storing a single matrix means having a space and time complexity of $O(m^2 n^2)$ (cf. Section 3), while in the second case the complexity is $O(mn^2)$. We use this method in the multiclass experiments presented in Section 7.

6. Bounding Other Online Algorithms

It is possible to apply the technique in the basis of the Projectron algorithm to any conservative online algorithm. A conservative online algorithm is an algorithm that updates its hypothesis only on rounds on which it makes a prediction error. By applying Lemma 2 to a conservative algorithm, we can construct a bounded version of it with worst case mistake bounds. As in the previous proofs, the idea is to use Lemma 2 to bound the scalar product of the competitor and the projected function. This yields an additional term which is subtracted from the margin γ of the competitor.

The technique presented here can be applied to other online kernel-based algorithms. As an example, we apply our technique to ALMA₂ (Gentile, 2001). Again we define two hypotheses: a temporary hypothesis f'_t , which is the hypothesis of ALMA₂ after its update rule, and a projected hypothesis, which is the hypothesis f'_t projected on the set \mathcal{H}_{t-1} as defined in Equation (2). Define the projection error δ_t as $\delta_t = f'_t - f''_t$. The modified ALMA₂ algorithm uses the projected hypothesis f''_t whenever the projection error is smaller than a parameter η , otherwise it uses the temporary hypothesis f'_t . We can state the following bound

Theorem 7 *Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ be a sequence of instance-label pairs where $\mathbf{x}_t \in \mathcal{X}$, $y_t \in \{-1, +1\}$, and $k(\mathbf{x}_t, \mathbf{x}_t) \leq 1$ for all t . Let α, B and $C \in \mathbb{R}^+$ satisfy the equation*

$$C^2 + 2(1 - \alpha)BC = 1.$$

Assume ALMA₂($\alpha; B, C$) projects every time the projection error $\|\delta_t\|$ is less than $\eta \geq 0$, then the number of prediction mistakes it makes on the sequence is bounded by

$$\frac{D_\gamma}{\gamma - \eta} + \frac{\rho^2}{2} + \sqrt{\frac{\rho^4}{4} + \frac{\rho^2}{\gamma - \eta} D_\gamma + \rho^2}$$

where $\gamma > \eta$, $\rho = \frac{1}{C^2(\gamma - \eta)^2}$, and g is an arbitrary function in \mathcal{H} , such that $\|g\| \leq 1$.

Proof The proof follows the original proof presented in Gentile (2001). Specifically, according to Lemma 2, one can replace the relation $y_t \langle g, k(\mathbf{x}_t, \cdot) \rangle \geq \gamma - \ell_\gamma(g(\mathbf{x}_t), y_t)$ with $y_t \langle g, P_{t-1} k(\mathbf{x}_t, \cdot) \rangle \geq$

Data Set	Samples	Features	Classes	Kernel	Parameters
<i>a9a</i> (Platt, 1999)	32561	123	2	Gaussian	0.04
<i>ijcnn1</i> (Prokhorov, 2001)	49990	22	2	Gaussian	8
<i>news20.binary</i> (Keerthi et al., 2005)	19996	1355191	2	Linear	-
<i>vehicle</i> (Duarte and Hu, 2004)	78823	100	2	Gaussian	0.125
<i>synthetic</i> (Dekel et al., 2007)	10000	2	2	Gaussian	1
<i>mnist</i> (Lecun et al., 1998)	60000	780	10	Polynomial	7
<i>usps</i> (Hull, 1994)	7291	256	10	Polynomial	13
<i>timit</i> (subset) (Lemel et al., 1986)	~ 150000	351	39	Gaussian	80

Table 1: Data sets used in the experiments

$\gamma - \eta - \ell_\gamma(g(\mathbf{x}_t), y_t)$, and further substitute $\gamma - \eta$ for γ . ■

7. Experimental Results

In this section we present experimental results that demonstrate the effectiveness of the Projectron and the Projectron++ algorithms. We compare both algorithms to the Perceptron algorithm, the Forgetron algorithm (Dekel et al., 2007) and the Randomized Budget Perceptron (RBP) algorithm (Cesa-Bianchi et al., 2006). For Forgetron, we choose the state-of-the-art “self-tuned” variant, which outperforms all of its other variants. We used the PA-I variant of the Passive-Aggressive algorithm (Crammer et al., 2006) as a baseline algorithm, as it gives an upper bound on the classification performance of the Projectron++ algorithm. All the algorithms were implemented in MATLAB using the DOGMA library (Orabona, 2009).

We tested the algorithms on several standard machine learning data sets:² *a9a*, *ijcnn1*, *news20.binary*, *vehicle* (combined), *usps*, *mnist*. We also used a synthetic dataset and the acoustic-phonetic dataset *timit*. The synthetic dataset was built in the same way as in Dekel et al. (2007). It is composed of 10000 samples taken from two separate bi-dimensional Gaussian distributions. The means of the positive and negative samples are $(1, 1)$ and $(-1, -1)$, respectively, while the covariance matrices for both are diagonal matrices with $(0.2, 2)$ as their diagonal. The labels are flipped with a probability of 0.1 to introduce noise. The list of the data sets, their characteristics and the kernels used, are given in Table 1. The parameters of the kernels were selected to have the best performance with the Perceptron and were used for all the other algorithms to result in a fair comparison. The C parameter of PA-I was set to 1, to give an update similar to Perceptron and Projectron. All the experiments were performed over five different permutations of the training set.

Experiments with one setting of η . In the first set of experiments we compared the online average number of mistakes and the support set size of all algorithms. Both Forgetron and RBP work by discarding vectors from the support set, if the size of the support set reaches the budget size, B . Hence for a fair comparison, we set η to some value and selected the budget sizes of Forgetron and RBP to be equal to the final size of the support set of Projectron. In particular, in Figure 6, we set $\eta = 0.1$ in Projectron and ended up with a support set of size 793, hence $B = 793$. In Figure 6(a) the average online error rate for all algorithms on the *a9a* data set is plotted. Note that Projectron closely tracks Perceptron. On the other hand Forgetron and RBP stop improving

². Downloaded from <http://www.sie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>.

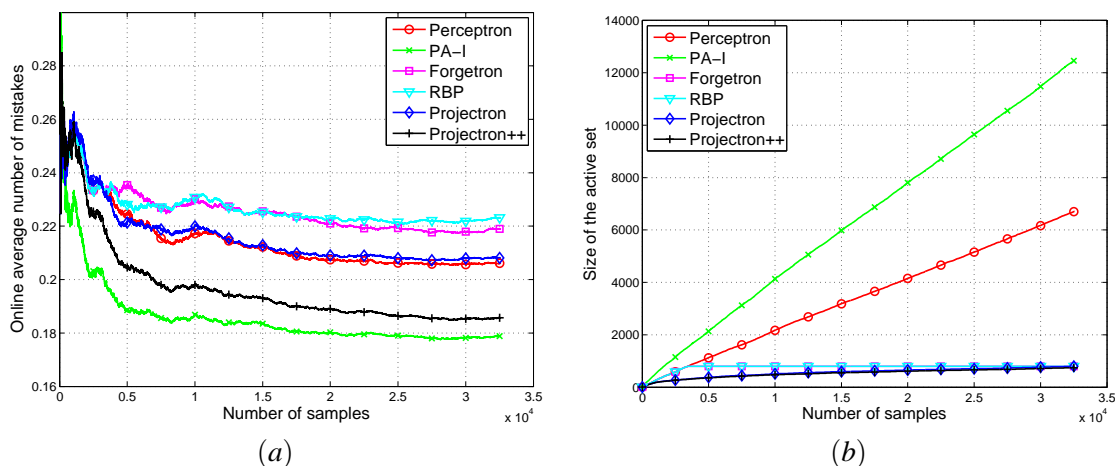


Figure 6: Average online error (left) and size of the support set (right) for the different algorithms on *a9a* data set as a function of the number of training samples (better viewed in color). B is set to 793, $\eta = 0.1$.

after reaching the support set size B , around 3400 samples. Moreover, as predicted by its theoretical analysis, Projectron++ achieves better results than Perceptron, even with fewer number of supports.

Figure 6(b) shows the growth of the support set as a function of the number of samples. While for the PA-I and the Perceptron the growth is clearly linear, it is sub-linear for Projectron and Projectron++: they will reach a maximum size and then they will stop growing, as stated in Theorem 1. Another important consideration is that Projectron++ outperforms Projectron *both* with respect to the size of the support set and number of mistakes. Using our MATLAB implementation, the running times for this experiment are ~ 35 s for RBP and Forgetron, ~ 40 s for Projectron and Projectron++, ~ 130 s for Perceptron, and ~ 375 s for PA-I. Hence Projectron and Projectron++ have a running time smaller than Perceptron and PA-I, due to their smaller support sets.

The same behavior can be seen in Figure 7, for the *synthetic* data set. Here the gain in performance of Projectron++ over Perceptron, Forgetron and RBP is even greater.

Experiments with a range of values for η - Binary. To analyze in more detail the behavior of our algorithms we decided to run other tests using a range of values of η . For each value we obtain a different size of the support set and a different number of mistakes. We used the data to plot a curve corresponding to the percentage of mistakes as a function of the support set size. The same curve was plotted for Forgetron and RBP, where the budget size was selected as described before. In this way we compared the algorithms along the continuous range of budget sizes, displaying the trade-off between sparseness and accuracy. For the remaining experiments we chose not to show the performance of Projectron, as it was always outperformed by Projectron++.

In Figure 8 we show the performance of the algorithms on different binary data sets: (a) *ijcnn1*, (b) *a9a*, (c) *news20.binary*, and (d) *vehicle (combined)*. Because Projectron++ used a different support set size for each permutation of the training samples, we plotted five curves, one for each of the five permutations. RBP and Forgetron have fixed budget sizes set in advance, hence for these algorithms we just plotted standard deviation bars, that are very small so they can be hardly seen in the figures. In all of the experiments Projectron++ outperforms Forgetron and RBP. One may

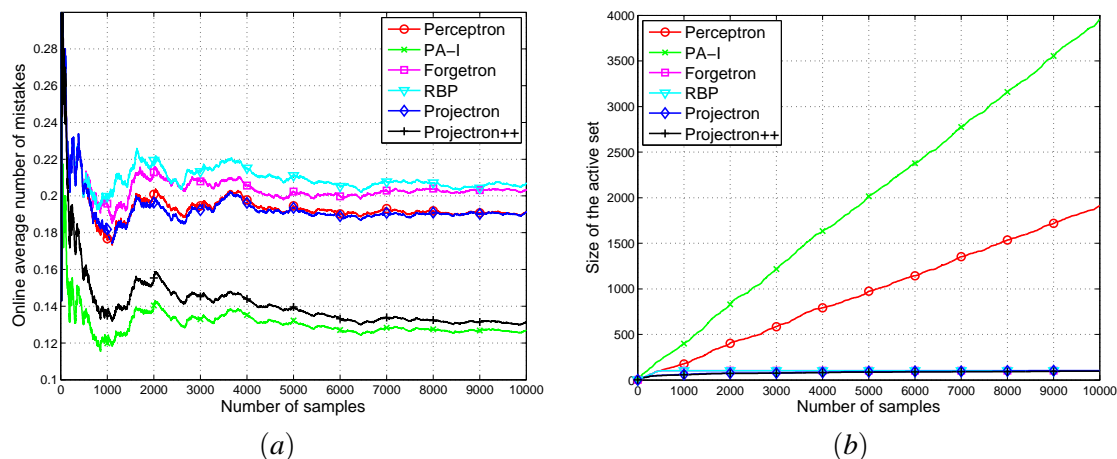


Figure 7: Average online error (left) and size of the support set (right) for the different algorithms on the *synthetic* data set as a function of the number of training samples (better viewed in color). B is set to 103, $\eta = 0.04$.

note that there is a point in all the graphs where the performance of Projectron++ is better than Perceptron, and has a smaller support set. Projectron++ gets closer to the classification rate of the PA-I, without paying the price of a larger support set. Note that the performance of Projectron++ is consistently better than RBP and Forgetron, regardless of the kernel used, particularly, on the database *news20.binary*, which is a text classification task with linear kernel. In this task the samples are almost mutually orthogonal, so finding a suitable subspace on which to project is difficult. Nevertheless Projectron++ succeeded in obtaining better performance. The reason is probably due to the margin updates, which are performed without increasing the size of the solution. Note that a similar modification would not be trivial in Forgetron and in RBP, because the proofs of their mistake bounds strongly depend on the rate of growth of the norm of the solution.

Experiments with a range of values for η - Multiclass. We have also considered multiclass data sets, using the multiclass version of Projectron++. Due to the fact that there are no other bounded online algorithms with a mistake bound for multiclass, we have extended RBP in the natural manner to multiclass. In particular we used the *max-score* update in Crammer and Singer (2003), for which a mistake bound exists, discarding a vector at random from the solution each time a new instance is added and the number of support vectors is equal to the budget size. We name it Multiclass Random Budget Perceptron (MRBP). It should be possible to prove a mistake bound for this algorithm, extending the proof in Cesa-Bianchi et al. (2006). In Figure 9 we show the results for Perceptron, Passive-Aggressive, Projectron++ and MRBP trained on (a) *usps*, and (b) *mnist* data sets. The results confirm the findings found for the binary case.

The last data set used in our experiments is a corpus of continuous natural speech for the task of phoneme classification. The data we used is a subset of the TIMIT acoustic-phonetic data set, which is a phonetically transcribed corpus of high quality continuous speech spoken by North American speakers (Lemel et al., 1986). The features were generated from nine adjacent vectors of Mel-Frequency Cepstrum Coefficients (MFCC) along with their first and second derivatives. The TIMIT corpus is divided into a training set and a test set in such a way that no speakers from the training set

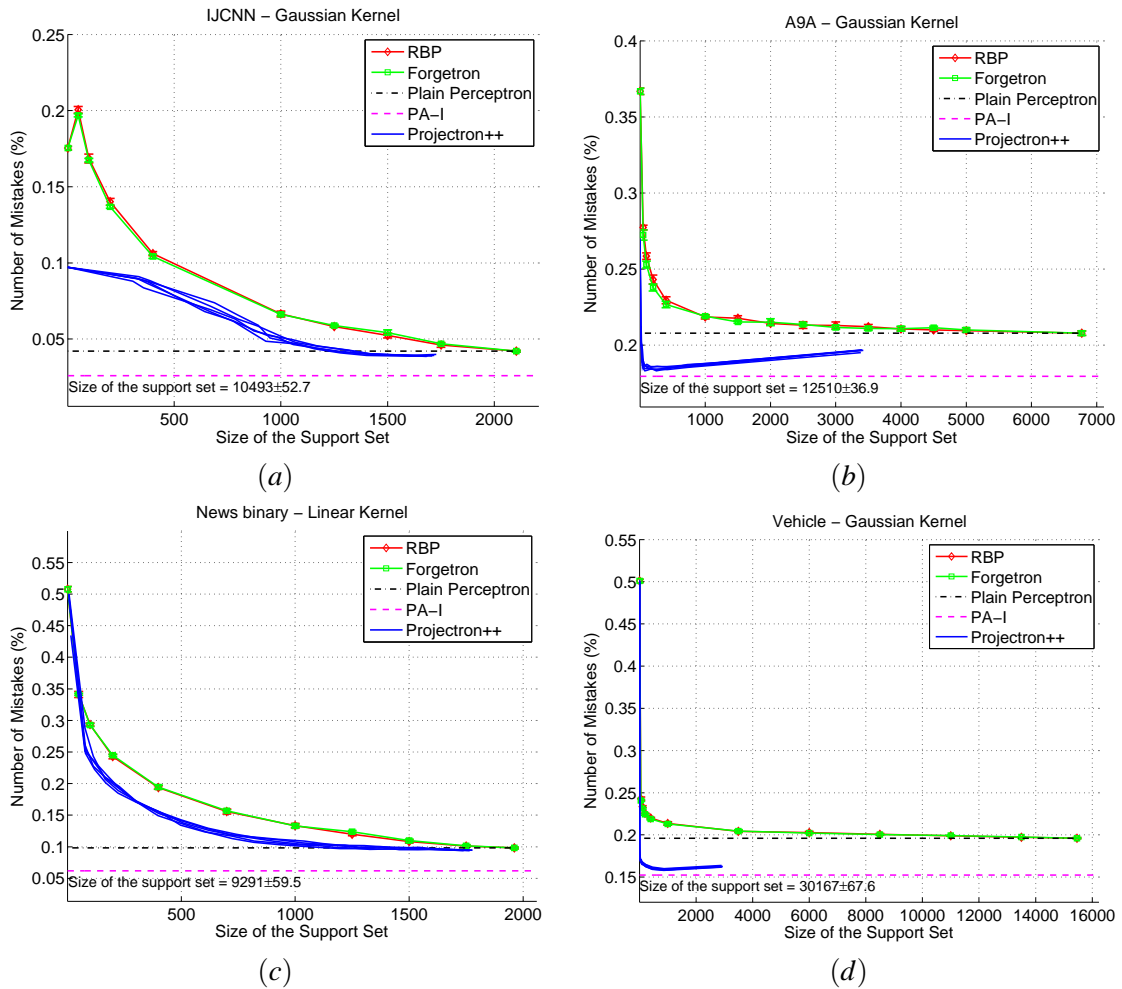


Figure 8: Average online error for the different algorithms as a function of the size of the support set on different binary data sets.

appear in the test set (speaker independent). We randomly selected 500 training utterances from the training set. The average online errors are shown in Figure 10 (a). We also tested the performance of the algorithm on the proposed TIMIT core test set composed of 192 utterances, the results of which are in Figure 10 (b). We used online-to-batch conversion (Cesa-Bianchi et al., 2004) to give a bounded batch solution. We did not test the performance of MRBP on the test set because for this algorithm the online-to-batch conversion does not produce a bounded solution. We compare the batch solution to the online-to-batch conversion of the PA-I solution. The results of Projectron++ are comparable to those of PA-I, while the former uses a smaller support set. These results also suggest that the batch solution is stable when varying the value of η , as the difference in performance on test set is less than 3%.

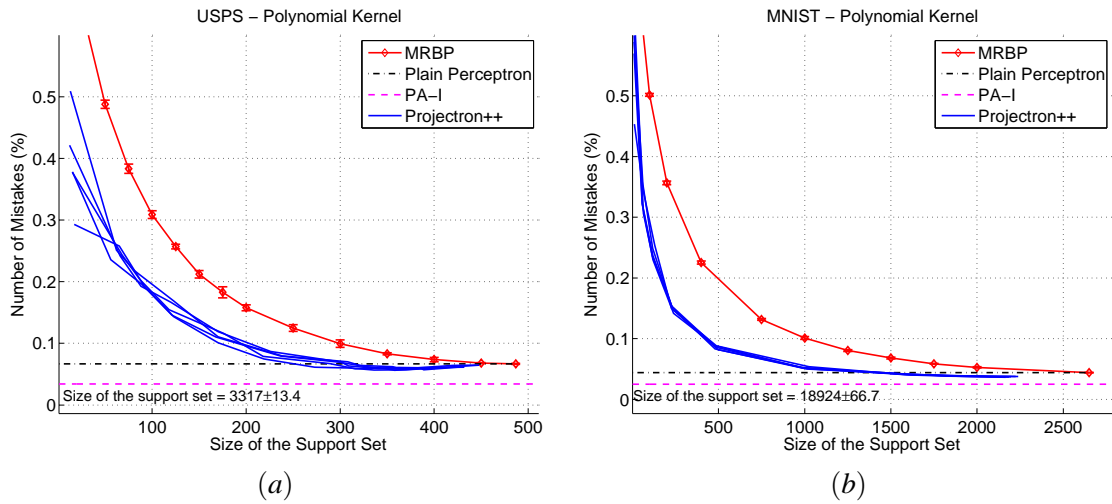


Figure 9: Average online error for the different algorithms as a function of the size of the support set on different multiclass data sets.

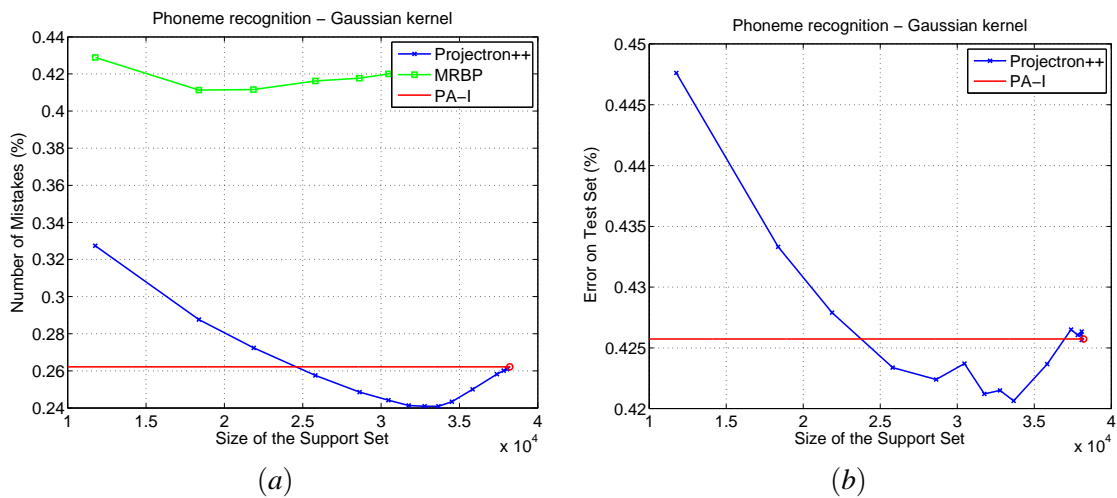


Figure 10: Average online error (a) and test error (b) for the different algorithms as a function of the size of the support set on a subset of the *timit* data set.

8. Discussion

This paper presented two different versions of a bounded online learning algorithm. The algorithms depend on a parameter that allows one to trade accuracy for sparseness of the solution. The size of the solution is always guaranteed to be bounded, although the size of this bound is unknown before the training begins. Therefore, these algorithms solve the memory explosion problem of the Perceptron and similar algorithms. Although the size of the support set cannot be determined

before training, practically, for a given target accuracy, the size of the support sets of Projectron or Projectron++ are much smaller than those of other budget algorithms such as Forgetron and RBP.

The first algorithm, Projectron, is based on the Perceptron algorithm. The empirical performance of Projectron is comparable to that of Perceptron, but with the advantage of a bounded solution. The second algorithm, Projectron++, introduces the notion of large margin and, for some values of η , outperforms the Perceptron algorithm, while assuring a bounded solution. The experimental results suggest that Projectron++ outperforms other online bounded algorithms such as Forgetron and RBP, with a similar hypothesis size.

There are two unique advantages of Projectron and Projectron++. First, these algorithms can be extended to the multiclass and the structured output settings. Second, a standard online-to-batch conversion can be applied to the online bounded solution of these algorithms, resulting in a bounded batch solution. The major drawback of these algorithms is their time and space complexity, which is quadratic in the size of the support set. Trying to overcome this acute problem is left for future work.

Acknowledgments

The authors would like to thank the anonymous reviewers for suggesting a way to improve Theorem 3. This work was supported by EU project DIRAC (FP6-0027787). The authors would like to thank Andy Cotter for proofreading the manuscript.

References

- G. Cauwenberghs and T. Poggio. Incremental and decremental support vector machine learning. In *Advances in Neural Information Processing Systems 14*, pages 409–415, 2000.
- N. Cesa-Bianchi, A. Conconi, and C. Gentile. On the generalization ability of on-line learning algorithms. *IEEE Trans. on Information Theory*, 50(9):2050–2057, 2004.
- N. Cesa-Bianchi, A. Conconi, and C. Gentile. A second-order perceptron algorithm. *SIAM Journal on Computing*, 34(3):640–668, 2005.
- N. Cesa-Bianchi, A. Conconi, and C. Gentile. Tracking the best hyperplane with a simple budget Perceptron. In *Proc. of the 19th Conference on Learning Theory*, pages 483–498, 2006.
- L. Cheng, S. V. N. Vishwanathan, D. Schuurmans, S. Wang, and T. Caelli. Implicit online learning with kernels. In *Advances in Neural Information Processing Systems 19*, pages 249–256, 2007.
- M. Collins. Discriminative reranking for natural language parsing. In *Proc. of the 17th International Conference on Machine Learning*, pages 175–182, 2000.
- K. Crammer and Y. Singer. Ultraconservative online algorithms for multiclass problems. *Journal of Machine Learning Research*, 3:951–991, 2003.
- K. Crammer, J. Kandola, and Y. Singer. Online classification on a budget. In *Advances in Neural Information Processing Systems 16*, 2003.

- K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7:551–585, 2006.
- L. Csató and M. Opper. Sparse on-line gaussian processes. *Neural Computation*, 14:641–668, 2002.
- F. Cucker and D. X. Zhou. *Learning Theory: An Approximation Theory Viewpoint*. Cambridge University Press, New York, NY, USA, 2007.
- O. Dekel, S. Shalev-Shwartz, and Y. Singer. The Forgetron: A kernel-based perceptron on a budget. *SIAM Journal on Computing*, 37(5):1342–1372, 2007.
- T. Downs, K. E. Gates, and A. Masters. Exact simplification of support vectors solutions. *Journal of Machine Learning Research*, 2:293–297, 2001.
- M. F. Duarte and Y. H. Hu. Vehicle classification in distributed sensor networks. *Journal of Parallel and Distributed Computing*, 64:826–838, 2004.
- Y. Engel, S. Mannor, and R. Meir. The kernel recursive least-squares algorithm. *IEEE Transactions on Signal Processing*, 52(8):2275–2285, 2004.
- Y. Freund and R. E. Schapire. Large margin classification using the Perceptron algorithm. *Machine Learning*, pages 277–296, 1999.
- C. Gentile. A new approximate maximal margin classification algorithm. *Journal of Machine Learning Research*, 2:213–242, 2001.
- C. Gentile. The robustness of the p-norm algorithms. *Machine Learning*, 53(3):265–299, 2003.
- J. J. Hull. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5):550–554, 1994.
- S. S. Keerthi, D. Decoste, and T. Joachims. A modified finite newton method for fast solution of large scale linear svms. *Journal of Machine Learning Research*, 6:2005, 2005.
- J. Kivinen, A. Smola, and R. Williamson. Online learning with kernels. *IEEE Trans. on Signal Processing*, 52(8):2165–2176, 2004.
- J. Langford, L. Li, and T. Zhang. Sparse online learning via truncated gradient. In *Advances in Neural Information Processing Systems 21*, pages 905–912, 2008.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998.
- L. Lemel, R. Kassel, and S. Seneff. Speech database development: Design and analysis. Report no. SAIC-86/1546, Proc. DARPA Speech Recognition Workshop, 1986.
- F. Orabona. *DOGMA: A MATLAB Toolbox for Online Learning*, 2009. Software available at <http://dogma.sourceforge.net>.
- F. Orabona, C. Castellini, B. Caputo, J. Luo, and G. Sandini. Indoor place recognition using online independent support vector machines. In *Proc. of the British Machine Vision Conference 2007*, pages 1090–1099, 2007.

- J. C. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in kernel methods – support vector learning*, pages 185–208. MIT Press, Cambridge, MA, USA, 1999.
- D. Prokhorov. IJCNN 2001 neural network competition. Technical report, 2001.
- F. Rosenblatt. The Perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–407, 1958.
- B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In *Proc. of the 14th Conference on Computational Learning Theory*, pages 416–426, London, UK, 2001. Springer-Verlag.
- B. Taskar, C. Guestrin, and D. Koller. Max-margin Markov networks. In *Advances in Neural Information Processing Systems 17*, 2003.
- I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun. Support vector machine learning for interdependent and structured output spaces. In *Proc. of the 21st International Conference on Machine Learning*, pages 823–830, 2004.
- J. Weston, A. Bordes, and L. Bottou. Online (and offline) on an even tighter budget. In Robert G. Cowell and Zoubin Ghahramani, editors, *Proc. of AISTATS 2005*, pages 413–420, 2005.