
The Margin Perceptron with Unlearning

Constantinos Panagiotakopoulos
Petroula Tsampouka

COSTAPAN@ENG.AUTH.GR
PETROULA@GEN.AUTH.GR

Physics Division, School of Technology, Aristotle University of Thessaloniki, 54124 Thessaloniki, Greece

Abstract

We introduce into the classical Perceptron algorithm with margin a mechanism of unlearning which in the course of the regular update allows for a reduction of possible contributions from “very well classified” patterns to the weight vector. The resulting incremental classification algorithm, called Margin Perceptron with Unlearning (MPU), provably converges in a finite number of updates to any desirable chosen before running approximation of either the maximal margin or the optimal 1-norm soft margin solution. Moreover, an experimental comparative evaluation involving representative linear Support Vector Machines reveals that the MPU algorithm is very competitive.

1. Introduction

Support Vector Machines (SVMs) (Boser et al., 1992; Vapnik, 1998; Cristianini & Shawe-Taylor, 2000) have been extensively used as linear classifiers either in the space where the patterns originally reside or in high dimensional feature spaces induced by kernels. SVMs appear to be very successful at addressing the problem of minimising an objective function involving the empirical risk while at the same time keeping low the complexity of the classifier. As measures of the empirical risk various quantities have been proposed with the 1- and 2-norm loss functions being the most widely accepted ones giving rise to the optimisation problems known as L1- and L2-SVMs (Cortes & Vapnik, 1995). In the case that the 2-norm loss function takes the place of the empirical risk an equivalent formulation exists which renders the dataset linearly separable in a high dimensional feature space. Therefore, one can either solve the general optimisation problem consisting of two terms, namely the capacity term and the 2-norm loss or just attempt to minimise the capacity

term in the high dimensional space as long as it ensures zero loss for the patterns mapped in that space. In the second treatment minimisation of the capacity term coincides with maximisation of the margin. SVMs traditionally follow the second approach.

On the other hand solving the maximum margin problem has also been the goal of various Perceptron-like algorithms (Li & Long, 2002; Gentile, 2001; Tsampouka & Shawe-Taylor, 2007). The key characteristic of such algorithms is that they work in the primal space in an online manner, i.e. processing one example at a time. Cycling repeatedly through the patterns they update their internal state stored in the weight vector each time a misclassification condition is satisfied by a pattern. More specifically, the Perceptron-like algorithms mentioned above are able to attain any approximation of the maximum margin, thereby producing infinitely close approximations of the L2-SVM solution. On the contrary, the classical Perceptron algorithm with margin (Duda & Hart, 1973) provably achieves only up to 1/2 of the maximum margin (Krauth & Mézard, 1987). The weight vector of the Perceptron algorithm (Rosenblatt, 1958) is built progressively by adding to it the patterns found misclassified multiplied by the corresponding labels. This procedure is characterised as learning. However, by introducing in the present work what we call unlearning, a mechanism allowing for a reduction or even elimination of the contribution that some patterns have to the solution weight vector when they are “very well classified”, we find that the Perceptron’s inability to provably achieve maximum margin is miraculously remedied. The possibility that the contribution of a given pattern to the current solution may increase or decrease with running time is very common to SVM algorithms although the procedure by which this is achieved does not tempt one to pay particular attention to a clear-cut distinction between learning and unlearning.

Since the L1-SVM problem is not known to admit an equivalent maximum margin interpretation via a mapping to an appropriate space Perceptron-like algorithms appeared so far unable to deal with such a task. In early attempts (Guyon & Stork, 1999) at ad-

Appearing in *Proceedings of the 27th International Conference on Machine Learning*, Haifa, Israel, 2010. Copyright 2010 by the author(s)/owner(s).

addressing the L1-SVM problem with a Perceptron an upper bound was set on the number of updates associated with any pattern before that pattern ended up being ignored in the process, thereby establishing artificially linear separability. The defect of the method was the impossibility for the patterns ignored to reenter the process rendering the effect of any early updates performed towards wrong directions totally irreversible. As a consequence, the outcome lacked the uniqueness dictated by optimisation theory since it depended heavily on the sequence according to which the patterns were presented to the algorithm. Here, we follow a procedure which does not suffer from the previous restrictions and through which the Perceptron is able in a finite number of steps to achieve any desirable fixed before running approximation of the optimal solution. Again, the mechanism of unlearning may be identified as the decisive factor in phasing out the misleading effect that some patterns have, mostly during early stages of the algorithm.

The decision of whether the contribution of a pattern to the weight vector should be reduced concerns only the pattern that is currently presented to the algorithm. Consequently, the mechanism of unlearning becomes vacuous if our algorithm is implemented in a strictly online setting in which all examples are considered to be distinct. This reveals an important difference from the Budget Perceptron (Crammer et al., 2003), an online algorithm featuring a removal of past examples from the expression of the weight vector which, however, only takes place whenever the pattern currently presented to the algorithm is inserted. Furthermore, in a batch setting we aim at convergence to the known SVM solution.

The paper is organised as follows. Section 2 describes the algorithm. Section 3 is devoted to a theoretical analysis. In Section 4 we give details concerning the implementation of our algorithm while in Section 5 we deliver our experimental results. Finally, Section 6 contains our conclusions.

2. The Margin Perceptron Algorithm with Unlearning

Given a training set which may be already mapped into a feature space of higher dimensionality we place all patterns in the same position at a distance ρ in an additional dimension, thereby constructing an embedding of our data into the so-called augmented space (Duda & Hart, 1973). The advantage of this embedding is that the linear hypothesis in the augmented space becomes homogeneous. Throughout our discussion we assume that the augmented patterns are mul-

tiplied by the corresponding labels in order to allow for a uniform treatment of both categories of patterns. Also, $R \equiv \max_k \|\mathbf{y}_k\|$ with \mathbf{y}_k the k -th augmented pattern multiplied by its label.

First we assume that the training set is linearly separable in the appropriate feature space. The relation characterising optimally correct classification of the training patterns by a weight vector \mathbf{u} of unit norm in the augmented space is

$$\mathbf{u} \cdot \mathbf{y}_k \geq \gamma_d \equiv \max_{\mathbf{u}' : \|\mathbf{u}'\|=1} \min_i \{\mathbf{u}' \cdot \mathbf{y}_i\} \quad \forall k. \quad (1)$$

Here γ_d is the maximum margin in the augmented space with respect to hyperplanes passing through the origin which should be distinguished from the maximum geometric margin γ in the non-augmented feature space.

The Margin Perceptron with Unlearning (MPU) at its t -th step maintains two vectors. The first is the usual augmented weight vector \mathbf{a}_t in the primal representation having the dimensionality M of the feature space. The second is a vector of non-negative integer-valued components I_k^t and dimensionality equal to the number N of patterns which keeps track of the update history of the algorithm. Both are initially set to zero, i.e. $\mathbf{a}_0 = \mathbf{0}$ and $I_k^0 = 0 \quad \forall k$, and are updated according to the rule

$$\begin{aligned} \mathbf{a}_{t+1} &= \mathbf{a}_t + \eta_t \mathbf{y}_k \\ I_k^{t+1} &= I_k^t + \eta_t \end{aligned} \quad (2)$$

each time the training pattern \mathbf{y}_k presented to the algorithm satisfies either one of the conditions appearing in the following definition of η_t

$$\eta_t = \begin{cases} +1, & \text{if } \mathbf{a}_t \cdot \mathbf{y}_k \leq b \\ -1, & \text{if } \mathbf{a}_t \cdot \mathbf{y}_k \geq b + \delta b \text{ and } I_k^t > 0. \end{cases} \quad (3)$$

Here b and δb are positive constants. We shall refer to an update with $\eta_t = +1$ as a learning update (step) whereas to an update with $\eta_t = -1$ as an unlearning update (step). Thus, the variable I_k^t is equal to the difference between the number of learning and unlearning steps associated with the k -th pattern. The total number of learning or unlearning steps will be denoted as t^+ or t^- , respectively. Obviously, the total number of updates t is

$$t = t^+ + t^- \quad (4)$$

and

$$\sum_k I_k^t = \sum_{\tau=0}^{t-1} \eta_\tau = t^+ - t^-. \quad (5)$$

Algorithm 1 The Margin Perceptron with Unlearning

Input: A dataset $S = (\mathbf{y}_1, \dots, \mathbf{y}_k, \dots, \mathbf{y}_N)$ with \mathbf{y}_k the k -th augmented pattern multiplied by its label
Define: $R \equiv \max_k \|\mathbf{y}_k\|$
Require: $I, b > 0, \delta b > R^2$
Initialise: $t = 0, \mathbf{a}_0 = \mathbf{0}, I_k^0 = 0 \forall k$
repeat
 for $k = 1$ **to** N **do**
 $p_{tk} = \mathbf{a}_t \cdot \mathbf{y}_k$
 if ($p_{tk} \leq b$ **and** $I_k^t < I$) **then**
 $\mathbf{a}_{t+1} = \mathbf{a}_t + \mathbf{y}_k$
 $I_k^{t+1} = I_k^t + 1$
 $t \leftarrow t + 1$
 else if $I_k^t = 0$ **then**
 continue
 else if $p_{tk} \geq b + \delta b$ **then**
 $\mathbf{a}_{t+1} = \mathbf{a}_t - \mathbf{y}_k$
 $I_k^{t+1} = I_k^t - 1$
 $t \leftarrow t + 1$
 end if
 end for
until no update made within the **for** loop

It should be clear that with the variables I_k^t produced by the algorithm as described above the weight vector admits the so-called dual representation

$$\mathbf{a}_t = \sum_k I_k^t \mathbf{y}_k \quad (6)$$

as a linear combination with positive coefficients of the training patterns. Consequently, the I_k^t 's play the role of the dual variables. A learning step takes place whenever the pattern \mathbf{y}_k satisfies the misclassification condition $\mathbf{a}_t \cdot \mathbf{y}_k \leq b$ meaning that \mathbf{y}_k is either a classification or a margin error. We shall refer to such a pattern as misclassified. Instead, an unlearning step takes place whenever the pattern \mathbf{y}_k is ‘‘very well classified’’, i.e. $\mathbf{a}_t \cdot \mathbf{y}_k \geq b + \delta b$, and provided this pattern does have a non-vanishing contribution to the dual representation (6) of the current weight vector. Thus, a prerequisite for unlearning is that learning has previously taken place involving the pattern in question. Moreover, unlearning may be regarded as a mechanism preventing the dual representation (6) from expanding unnecessarily. Notice, that due to the ‘‘discrete’’ nature of the update rule the dual variables I_k^t remain naturally non-negative without having to be clipped as in the case of SVMs.

An important feature of the algorithm is the existence of the gap δb between the threshold b of the misclassification condition $\mathbf{a}_t \cdot \mathbf{y}_k \leq b$ and the threshold $b + \delta b$ of the condition $\mathbf{a}_t \cdot \mathbf{y}_k \geq b + \delta b$ for exceedingly good clas-

sification. Actually, it is not very difficult to see that such a gap should necessarily satisfy the constraint

$$\delta b > R^2 \quad (7)$$

in order for the algorithm to have a chance of converging in a finite number of steps. Indeed, let us assume that the pattern \mathbf{y}_k presented to the algorithm satisfies $\mathbf{a}_t \cdot \mathbf{y}_k = b$. Then, a learning update takes place as a result of which the weight vector becomes $\mathbf{a}_{t+1} = \mathbf{a}_t + \mathbf{y}_k$. Now, consider a situation in which the same pattern is presented once again to the algorithm without another update having intervened. We have $\mathbf{a}_{t+1} \cdot \mathbf{y}_k = b + \|\mathbf{y}_k\|^2$ and $I_k^{t+1} > 0$. Therefore, if $\delta b \leq \|\mathbf{y}_k\|^2$ an unlearning step will take place which will exactly cancel the effect of the learning step. This procedure may be repeated again and again leading to a closed loop. If, instead, $\delta b > \|\mathbf{y}_k\|^2 \forall k$, i.e. $\delta b > R^2$ such unwanted situations are naturally avoided. The same condition also ensures that

$$t^+ > t^-, \text{ if } t > 0.$$

Obviously, one step before we arrive at the situation where the number of learning steps equals the number of unlearning steps we should have $t^+ = t^- + 1$ meaning that all dual variables vanish except for one which takes the value 1. So, let us assume that this happens at the t -th step of the algorithm and that the only non-vanishing dual variable is associated with the k -th pattern, i.e. $I_k^t = 1$ holds. Then, $\mathbf{a}_t = \mathbf{y}_k$ from where $\mathbf{a}_t \cdot \mathbf{y}_k = \|\mathbf{y}_k\|^2 < \delta b < b + \delta b$ forbidding the unlearning step that would lead to equality between learnings and unlearnings.

In the case that the assumption of linear separability in the feature space is relaxed we impose the further condition that a learning update is allowed provided that the dual variable I_k^t associated with the misclassified pattern \mathbf{y}_k satisfies the constraint $I_k^t < I$, where I is a positive integer. In such a case the update rule (2) remains the same but η_t is now defined by

$$\eta_t = \begin{cases} +1, & \text{if } \mathbf{a}_t \cdot \mathbf{y}_k \leq b \quad \text{and } I_k^t < I \\ -1, & \text{if } \mathbf{a}_t \cdot \mathbf{y}_k \geq b + \delta b \quad \text{and } I_k^t > 0. \end{cases} \quad (8)$$

Notice that again as a result of the ‘‘discrete’’ nature of the update rule the dual variables satisfy naturally the condition

$$0 \leq I_k^t \leq I \quad \forall k \quad (9)$$

without any need for clipping. Setting $I = \infty$ in (8) we recover (3).

Finally, in the case that the mapping into the feature space is induced by a kernel and a primal representation of \mathbf{a}_t is unavailable it suffices to update the dual variables I_k^t according to (2) and perform the inner products using the dual representation (6) of \mathbf{a}_t .

3. Theoretical Analysis

We first consider a linearly separable training set possessing maximum margin γ_d in the augmented space with respect to hyperplanes passing through the origin.

Theorem 1 *The Margin Perceptron with Unlearning and $I = \infty$ converges in a finite number t_c of updates satisfying the bound*

$$t_c \leq \frac{R^2 + 2b}{\gamma_d^2} \left(1 + \frac{R^2 + 2b}{4(\delta b - R^2)} \right). \quad (10)$$

Moreover, the zero-threshold solution hyperplane possesses margin γ'_d which is a fraction f of the maximum margin γ_d obeying the inequality

$$f \equiv \gamma'_d / \gamma_d > (1 + \min\{\delta b, b + R^2\}/b)^{-1}. \quad (11)$$

Finally, an after-running lower bound on f involving the margin γ'_d achieved, the length $\|\mathbf{a}_{t_c}\|$ of the solution weight vector \mathbf{a}_{t_c} and the final difference between the number of learning and unlearning steps ($t_c^+ - t_c^-$) is given by

$$f = \gamma'_d / \gamma_d \geq \gamma'_d (t_c^+ - t_c^-) / \|\mathbf{a}_{t_c}\|. \quad (12)$$

Proof Taking the inner product with the optimal direction \mathbf{u} in (6) and using (1) and (5) we obtain

$$\|\mathbf{a}_t\| \geq \mathbf{a}_t \cdot \mathbf{u} = \sum_k I_k^t (\mathbf{y}_k \cdot \mathbf{u}) \geq \gamma_d \sum_k I_k^t = \gamma_d (t^+ - t^-). \quad (13)$$

From the update rule (2) taking into account the conditions under which the update takes place and the definition of η_t we get

$$\begin{aligned} \|\mathbf{a}_{t+1}\|^2 - \|\mathbf{a}_t\|^2 &= \|\mathbf{y}_k\|^2 + 2\eta_t \mathbf{a}_t \cdot \mathbf{y}_k \\ &\leq R^2 + 2\eta_t b - \delta b(1 - \eta_t) \end{aligned}$$

a repeated application of which, using also (4) and (5) and the initial condition $\mathbf{a}_0 = \mathbf{0}$, gives

$$\begin{aligned} \|\mathbf{a}_t\|^2 &\leq R^2 t + 2b \sum_{\tau=0}^{t-1} \eta_\tau - \delta b \left(t - \sum_{\tau=0}^{t-1} \eta_\tau \right) \\ &= R^2 t + 2b(t^+ - t^-) - \delta b(t - (t^+ - t^-)) \\ &= (R^2 + 2b)(t^+ - t^-) - 2(\delta b - R^2)t^-. \end{aligned} \quad (14)$$

Combining (13) and (14) we obtain a Novikoff-like (Novikoff, 1962) squeezing relation

$$\gamma_d^2 (t^+ - t^-)^2 \leq \|\mathbf{a}_t\|^2 \leq (R^2 + 2b)(t^+ - t^-) - 2(\delta b - R^2)t^-. \quad (15)$$

From (15) we get

$$\begin{aligned} 2(\delta b - R^2)t^- &\leq (R^2 + 2b)(t^+ - t^-) - \gamma_d^2 (t^+ - t^-)^2 \\ &\leq (R^2 + 2b)^2 / 4\gamma_d^2 \end{aligned}$$

by making use of the fact that the quantity $(R^2 + 2b)(t^+ - t^-) - \gamma_d^2 (t^+ - t^-)^2$ acquires a maximum with respect to $(t^+ - t^-)$ for $(t^+ - t^-) = (R^2 + 2b)/2\gamma_d^2$. Then, taking into account (7) we immediately obtain

$$t^- \leq \frac{(R^2 + 2b)^2}{8\gamma_d^2(\delta b - R^2)}. \quad (16)$$

Again from (15) by dropping the last term in the upper bound on $\|\mathbf{a}_t\|^2$ we have $\gamma_d^2 (t^+ - t^-)^2 \leq (R^2 + 2b)(t^+ - t^-)$ or

$$t^+ - t^- \leq (R^2 + 2b)/\gamma_d^2. \quad (17)$$

Combining (16) and (17) we obtain the upper bound on the number t of updates given by the r.h.s. of (10).

Upon convergence of the algorithm in t_c updates it is obvious that all the patterns \mathbf{y}_k which have a non-vanishing contribution to the expansion (6) of the solution weight vector \mathbf{a}_{t_c} necessarily have their inner product with \mathbf{a}_{t_c} in the interval $(b, b + \delta b)$, i.e.

$$b < \mathbf{a}_{t_c} \cdot \mathbf{y}_k < b + \delta b, \quad \text{if } I_k^{t_c} > 0. \quad (18)$$

From (6) with $t = t_c$ we get

$$\|\mathbf{a}_{t_c}\|^2 = \mathbf{a}_{t_c} \cdot \left(\sum_k I_k^{t_c} \mathbf{y}_k \right) = \sum_k I_k^{t_c} \mathbf{a}_{t_c} \cdot \mathbf{y}_k$$

or, taking into account (5) and (18),

$$b(t_c^+ - t_c^-) < \|\mathbf{a}_{t_c}\|^2 < (b + \delta b)(t_c^+ - t_c^-). \quad (19)$$

Since upon convergence all patterns violate the misclassification condition we have $\mathbf{a}_{t_c} \cdot \mathbf{y}_k > b \quad \forall k$ or equivalently $\mathbf{u}_{t_c} \cdot \mathbf{y}_k > b / \|\mathbf{a}_{t_c}\| \quad \forall k$ with $\mathbf{u}_{t_c} \equiv \mathbf{a}_{t_c} / \|\mathbf{a}_{t_c}\|$. Stated differently, the margin achieved satisfies $\gamma'_d = \min_k \{\mathbf{u}_{t_c} \cdot \mathbf{y}_k\} > b / \|\mathbf{a}_{t_c}\|$. Additionally, from (13)

$$1/\gamma_d \geq (t_c^+ - t_c^-) / \|\mathbf{a}_{t_c}\|. \quad (20)$$

Thus, we obtain

$$f = \gamma'_d / \gamma_d > (b / \|\mathbf{a}_{t_c}\|) / \gamma_d \geq b(t_c^+ - t_c^-) / \|\mathbf{a}_{t_c}\|^2. \quad (21)$$

If we now make use in (21) of the upper bound on $\|\mathbf{a}_{t_c}\|^2$ from (19) we get $f > (1 + \delta b/b)^{-1}$. If, instead, we employ (15) to obtain the bound $\|\mathbf{a}_{t_c}\|^2 \leq (R^2 + 2b)(t_c^+ - t_c^-)$ we get $f > (2 + R^2/b)^{-1}$. The above lower bounds on f written compactly lead to (11).

Finally, (12) follows readily from (20). \square

Remark 1 From (11) it becomes apparent that as $b \rightarrow \infty$ with δb kept fixed the MPU algorithm with $I = \infty$ is guaranteed before running to converge to the maximum margin solution, i.e. $\gamma'_d \rightarrow \gamma_d$. If, instead, we first take the limit $\delta b \rightarrow \infty$ (thereby recovering the classical Perceptron with margin) and then the limit $b \rightarrow \infty$ we get $\gamma'_d > \gamma_d/2$.

Remark 2 The lower bound on $\|\mathbf{a}_{t_c}\|^2$ from (15) with the upper bound from (19) lead to $t_c^+ - t_c^- < (b + \delta b)/\gamma_d^2$. Also, from $\gamma_d \geq \gamma'_d > b/\|\mathbf{a}_{t_c}\|$ we get $\|\mathbf{a}_{t_c}\|^2 > b^2/\gamma_d^2$ which combined with the upper bound on $\|\mathbf{a}_{t_c}\|^2$ from (19) leads to $t_c^+ - t_c^- > b^2(b + \delta b)^{-1}\gamma_d^{-2}$. Thus, upon convergence of the MPU algorithm it holds that

$$(1 + \delta b/b)^{-1}b\gamma_d^{-2} < t_c^+ - t_c^- < (1 + \delta b/b)b\gamma_d^{-2}$$

meaning that $t_c^+ - t_c^-$ is strongly constrained. Moreover, $(t_c^+ - t_c^-)/b \rightarrow \gamma_d^{-2}$ as $b \rightarrow \infty$. This is the analog of the well-known result $\sum_k \alpha_k^* = \gamma^{-2}$ for the hard margin SVM with α_k^* being the optimal dual variables.

Remark 3 As we already pointed out the MPU algorithm in a strictly online setting reduces to the classical Perceptron with margin. Nevertheless, our derivation of the upper bound on the number of updates given by the r.h.s. of (10) does hold in the context of any related strictly online scenario, e.g. a variation of the variable-size cache Budget Perceptron, in the t -th step of which the insertion ($\eta_t = +1$) or the removal ($\eta_t = -1$) of the pattern \mathbf{y}_k is performed according to the rule (2) under the conditions stated in (3) (with $I_k^t \in \{0, 1\}$). Crucial in this respect is the validity of $\delta b > R^2$. The r.h.s. of (17) was shown in (Crammer et al., 2003) to be an upper bound on the variable cache size of the Budget Perceptron but no mistake bound was obtained.

Now we relax the assumption of linear separability of the training set in the feature space and we move on to the analysis of the MPU algorithm with $I < \infty$.

Theorem 2 *The Margin Perceptron with Unlearning converges in a finite number t_c of updates satisfying the bound*

$$t_c \leq NI \frac{2b + \delta b}{\delta b - R^2}. \quad (22)$$

Let us consider the ‘‘objective’’ function

$$\mathcal{J}(\mathbf{w}, C_p) \equiv \frac{1}{2} \|\mathbf{w}\|^2 + C_p \sum_k \max\{1 - \mathbf{w} \cdot \mathbf{y}_k, 0\},$$

where $C_p > 0$ is a ‘‘penalty’’ parameter and

$$\mathcal{J}_{\text{opt}}(C_p) \equiv \min_{\mathbf{w}} \mathcal{J}(\mathbf{w}, C_p)$$

its minimal value with C_p fixed. Then, with \mathbf{a}_{t_c} being the solution weight vector and provided $b > \delta b$

$$\frac{\delta \mathcal{J}}{\mathcal{J}_{\text{opt}}} \equiv \frac{\mathcal{J}(\mathbf{a}_{t_c}/b, I/b) - \mathcal{J}_{\text{opt}}(I/b)}{\mathcal{J}_{\text{opt}}(I/b)} < \delta \equiv \frac{2\delta b/b}{1 - \delta b/b}. \quad (23)$$

Moreover, $\frac{\delta \mathcal{J}}{\mathcal{J}_{\text{opt}}}$ satisfies the after-running bound

$$\frac{\delta \mathcal{J}}{\mathcal{J}_{\text{opt}}} \leq \delta_{\text{aft}} \equiv \frac{\frac{1}{2} \|\mathbf{a}_{t_c}\|^2 + I \sum_k \max\{b - \mathbf{a}_{t_c} \cdot \mathbf{y}_k, 0\}}{b(t_c^+ - t_c^-) - \frac{1}{2} \|\mathbf{a}_{t_c}\|^2} - 1. \quad (24)$$

Proof From (5) and the restriction $I_k^t \leq I$ we have

$$t^+ - t^- = \sum_k I_k^t \leq NI. \quad (25)$$

Also, from (14) and taking into account (25) we get

$$2(\delta b - R^2)t^- \leq (R^2 + 2b)(t^+ - t^-) \leq NI(R^2 + 2b)$$

or

$$t^- \leq NI \frac{R^2 + 2b}{2(\delta b - R^2)}$$

which combined with (25) leads to the upper bound on the number t of updates given by the r.h.s. of (22).

Let us assume that the algorithm has converged in t_c updates to the solution weight vector \mathbf{a}_{t_c} . We define

$$\mathbf{w}_c = \frac{\mathbf{a}_{t_c}}{b}, \quad \alpha_k^c = \frac{I_k^{t_c}}{b} \quad (0 \leq \alpha_k^c \leq \frac{I}{b}), \quad \epsilon = \frac{\delta b}{b} < 1$$

and notice that only the α_k^c 's with k from the sets

$$\begin{aligned} \mathcal{K}_1 &= \{k : \alpha_k^c > 0, \quad 1 < \mathbf{w}_c \cdot \mathbf{y}_k < 1 + \epsilon\} \\ \mathcal{K}_2 &= \{k : \alpha_k^c = I/b, \quad \mathbf{w}_c \cdot \mathbf{y}_k \leq 1\} \end{aligned} \quad (26)$$

are non-zero and consequently

$$\mathbf{w}_c = \sum_k \alpha_k^c \mathbf{y}_k = \sum_{k \in \mathcal{K}_1} \alpha_k^c \mathbf{y}_k + \sum_{k \in \mathcal{K}_2} \alpha_k^c \mathbf{y}_k.$$

Then,

$$\begin{aligned} \|\mathbf{w}_c\|^2 &= \mathbf{w}_c \cdot \left(\sum_k \alpha_k^c \mathbf{y}_k \right) = \sum_k \alpha_k^c \mathbf{w}_c \cdot \mathbf{y}_k \\ &= \sum_{k \in \mathcal{K}_1} \alpha_k^c \mathbf{w}_c \cdot \mathbf{y}_k + \sum_{k \in \mathcal{K}_2} \alpha_k^c \mathbf{w}_c \cdot \mathbf{y}_k \\ &< (1 + \epsilon) \sum_{k \in \mathcal{K}_1} \alpha_k^c + \sum_{k \in \mathcal{K}_2} \alpha_k^c - \sum_{k \in \mathcal{K}_2} \alpha_k^c + \sum_{k \in \mathcal{K}_2} \alpha_k^c \mathbf{w}_c \cdot \mathbf{y}_k \\ &< (1 + \epsilon) \left(\sum_{k \in \mathcal{K}_1} \alpha_k^c + \sum_{k \in \mathcal{K}_2} \alpha_k^c \right) - \sum_{k \in \mathcal{K}_2} \alpha_k^c (1 - \mathbf{w}_c \cdot \mathbf{y}_k) \\ &= (1 + \epsilon) \sum_k \alpha_k^c - \frac{I}{b} \sum_k \max\{1 - \mathbf{w}_c \cdot \mathbf{y}_k, 0\}. \end{aligned} \quad (27)$$

Now, taking into account (27) we have

$$\begin{aligned} \mathcal{J}(\mathbf{w}_c, I/b) &= \|\mathbf{w}_c\|^2 - \frac{1}{2} \|\mathbf{w}_c\|^2 \\ &\quad + \frac{I}{b} \sum_k \max\{1 - \mathbf{w}_c \cdot \mathbf{y}_k, 0\} \\ &= \|\mathbf{w}_c\|^2 - \frac{1}{2} \sum_{i,j} \alpha_i^c \alpha_j^c \mathbf{y}_i \cdot \mathbf{y}_j \\ &\quad + \frac{I}{b} \sum_k \max\{1 - \mathbf{w}_c \cdot \mathbf{y}_k, 0\} \\ &< (1 + \epsilon) \sum_k \alpha_k^c - \frac{1}{2} \sum_{i,j} \alpha_i^c \alpha_j^c \mathbf{y}_i \cdot \mathbf{y}_j. \end{aligned} \quad (28)$$

From the weak duality theorem it follows that

$$\mathcal{J}(\mathbf{w}, C_p) \geq \mathcal{L}(\boldsymbol{\alpha}) = \sum_k \alpha_k - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j \mathbf{y}_i \cdot \mathbf{y}_j,$$

where $\mathcal{L}(\boldsymbol{\alpha})$ is the dual Lagrangian and $0 \leq \alpha_k \leq C_p$. Thus, setting $C_p = I/b$ and $\alpha_k = \alpha_k^c$ we get

$$\mathcal{J}(\mathbf{w}_c, I/b) \geq \mathcal{J}_{\text{opt}}(I/b) \geq \sum_k \alpha_k^c - \frac{1}{2} \sum_{i,j} \alpha_i^c \alpha_j^c \mathbf{y}_i \cdot \mathbf{y}_j. \quad (29)$$

Combining (28) with (29) we obtain

$$\mathcal{J}(\mathbf{w}_c, I/b) - \mathcal{J}_{\text{opt}}(I/b) < \epsilon \sum_k \alpha_k^c. \quad (30)$$

Moreover, taking into account (27) we get

$$\begin{aligned} \sum_k \alpha_k^c - \frac{1}{2} \sum_{i,j} \alpha_i^c \alpha_j^c \mathbf{y}_i \cdot \mathbf{y}_j &= \sum_k \alpha_k^c - \frac{1}{2} \|\mathbf{w}_c\|^2 \\ &> \sum_k \alpha_k^c - \frac{1}{2}(1 + \epsilon) \sum_k \alpha_k^c = \frac{1}{2}(1 - \epsilon) \sum_k \alpha_k^c \end{aligned}$$

which combined with (29) leads to

$$\mathcal{J}_{\text{opt}}(I/b) > \frac{1}{2}(1 - \epsilon) \sum_k \alpha_k^c. \quad (31)$$

From (30) and (31) we immediately obtain (23).

Finally, using $\sum_k \alpha_k^c = \sum_k I_k^{t_c}/b = (t_c^+ - t_c^-)/b$, the definition of \mathcal{J} and the lower bound on \mathcal{J}_{opt} from (29) we easily derive (24). \square

Remark 4 From the before-running ‘‘accuracy’’ δ of (23) it follows that as b grows with δb kept fixed the MPU algorithm is guaranteed to converge to solutions for which the objective function $\mathcal{J}(\mathbf{a}_{t_c}/b, I/b)$ approaches as close as one wishes the optimal one $\mathcal{J}_{\text{opt}}(C_p)$ with $C_p = I/b$. Thus, if we allow only b to grow while keeping δb and I fixed we will asymptotically approach an optimal objective with vanishing penalty parameter C_p . If, instead, our goal is to approach the optimal objective $\mathcal{J}_{\text{opt}}(C_p)$ with a given value of C_p we should follow a limiting process in which both I and b increase with their ratio kept all the way equal to the desirable value of C_p . A simple way of doing so is to choose first a sufficiently large value for the integer I and then determine b from the relation $b = I/C_p$ with the value of δb kept fixed. To select the appropriate value of I which ensures before running the desired accuracy δ we should, of course, express δ in terms of I and C_p as $\delta = 2C_p \delta b I^{-1} (1 - C_p \delta b I^{-1})^{-1}$ from where $I = [C_p \delta b (2 + \delta) \delta^{-1}] + 1$ (with $[x]$ the integer part of $x \geq 0$). We conclude that the size of the upper bound I on the dual variables in the MPU algorithm does not depend on the penalty parameter C_p alone but on the required accuracy δ as well.

4. Implementation

To reduce the computational cost we construct a two-member sequence of reduced ‘‘active sets’’ of data points with the second being a subset of the first. As we cycle once through the full dataset the largest active set (first-level active set) is formed from the points of the full dataset either satisfying the condition $\mathbf{a}_t \cdot \mathbf{y}_k \leq \bar{c}b$ with $\bar{c} \geq 1$ or having $I_k^t > 0$. The second-level active set is built as we cycle once through the first-level active set and comprises the points \mathbf{y}_k for which I_k^t undergoes a change and remains positive. Of course, as we cycle through the full dataset or the first-level active set along with the selection of points which form the first-level or the second-level active set we do perform updates whenever a point satisfies the conditions either for learning or unlearning. The second-level active set is presented repetitively to the algorithm for N_{ep_2} mini-epochs. Then, the first-level active set becomes the set under consideration which is presented once again to the algorithm. During each round involving the first-level set, a new second-level set is constructed and a new cycle of N_{ep_2} passes begins. By invoking the first-level set for the $(N_{\text{ep}_1} + 1)^{\text{th}}$ time we trigger the loading of the full dataset and the procedure starts all over again until no update occurs as we run over the full dataset. It is understood, of course, that the maximum number of rounds N_{ep_1} and N_{ep_2} for each active set is not exhausted in the case that no update occurs during a round. Finally, following a common practice every time we make use of the full dataset we actually employ a permuted instance of it since presenting the points to the algorithm in a different order helps avoiding the disorientation caused by spurious patterns in the sequence of examples. Evidently, the whole procedure amounts to a different way of sequentially presenting the patterns to the algorithm and does not affect the validity of our analysis.

An additional mechanism providing a very substantial improvement of the computational efficiency is the one of performing multiple updates once a data point is presented to the algorithm¹. It is understood, of course, that in order for a multiple update to be compatible with our theoretical analysis it should be equivalent to a certain number of updates occurring as a result of repeatedly presenting to the algorithm the point in question. In the case of learning we are allowed to perform $\min\{\lambda, I - I_k^t\}$ learning updates at once with $\lambda = [(b - \mathbf{a}_t \cdot \mathbf{y}_k) \|\mathbf{y}_k\|^{-2}] + 1$. Indeed, for any positive integer $n < \min\{\lambda, I - I_k^t\} \leq \lambda$ and $\mathbf{a}_{t+n} = \mathbf{a}_t + n\mathbf{y}_k$ it holds that $\mathbf{a}_{t+n} \cdot \mathbf{y}_k = \mathbf{a}_t \cdot \mathbf{y}_k + n \|\mathbf{y}_k\|^2 \leq \mathbf{a}_t \cdot \mathbf{y}_k +$

¹The MPU with multiple updates bears some resemblance to the algorithm of (Hsieh et al., 2008).

Table 1. Results of a high accuracy comparative study of the MPU, DCD, SVM^{perf} and Pegasos algorithms.

data set	MPU		DCD			SVM ^{perf}			Pegasos		
	\mathcal{J}	Secs	ϵ	\mathcal{J}	Secs	ϵ	\mathcal{J}	Secs	iter	\mathcal{J}	Secs
Adult	11434.4	0.4	0.06	11434.4	0.3	0.004	11434.4	9.3	3×10^7	11434.9	397.0
Web	6605.17	0.1	0.07	6605.15	0.1	0.001	6605.22	11.1	2×10^7	6605.31	199.7
Physics	49643.6	0.4	0.05	49642.6	2.0	0.02	49644.1	2.2	5×10^5	49644.2	13.3
Real-sim	5402.25	0.5	0.01	5402.22	0.4	0.001	5402.58	6.1	2×10^7	5402.53	395.3
News20	2562.57	4.9	0.003	2562.57	4.2	0.001	2562.71	659.6	1×10^7	2562.69	1673.0
Webspam	69592.1	3.1	0.07	69591.4	2.8	0.001	69592.9	10.3	9×10^6	69599.6	297.7
Covertypes	340490	22.2	0.08	340491	46.6	0.008	340493	505.4	1×10^8	340505	2447.9
C11	44725.4	5.7	0.02	44725.3	5.1	0.001	44730.5	29.2	1×10^8	44730.3	2743.6
CCAT	92990.2	8.1	0.02	92989.3	7.9	0.001	92992.8	54.3	1×10^8	93000.5	3103.6

$(\lambda - 1) \|\mathbf{y}_k\|^2 \leq \mathbf{a}_t \cdot \mathbf{y}_k + (b - \mathbf{a}_t \cdot \mathbf{y}_k) \|\mathbf{y}_k\|^{-2} \|\mathbf{y}_k\|^2 = b$ meaning that after n consecutive updates the pattern \mathbf{y}_k is still misclassified with $I_k^{t+n} = I_k^t + n < I$ and consequently one more learning update involving that pattern could legitimately take place. Similarly, in the case of unlearning we can show that we are allowed to perform $\min\{\lambda, I_k^t\}$ unlearning updates at once with $\lambda = [(\mathbf{a}_t \cdot \mathbf{y}_k - (b + \delta b)) \|\mathbf{y}_k\|^{-2}] + 1$. Obviously, the gain from the mechanism of multiple updates lies in the fact that at the expense of a single inner product calculation we effectively perform numerous updates.

Finally, we point out that we may terminate the MPU algorithm with $I < \infty$ before it converges when the after-running accuracy δ_{aft} entering (24) (with \mathbf{a}_{t_c} and $(t_c^+ - t_c^-)$ now being the final weight vector and the final value of $(t^+ - t^-)$, respectively) falls below a certain desirable level δ_{stop} at the end of a complete cycle through the full dataset. Notice, that the derivation of (24) relies only on the validity of (5) and (9) which are always respected by the algorithm. Moreover, if the guaranteed before-running accuracy δ of (23) is equal or smaller than δ_{stop} it is obvious that, despite the early stopping, the algorithm will necessarily achieve the required accuracy δ_{stop} in at most the number of updates given by the r.h.s. of (22). Similarly, if $I = \infty$ we may terminate the MPU algorithm when the lower bound of (12) on f falls below the r.h.s. of (11).

5. Experiments

We concentrate on the more popular L1-SVM problem and compare the MPU algorithm with SVMs on the basis of their ability to arrive fast at a certain approximation of the optimal 1-norm soft margin solution. Moreover, in order to be able to train on rather large datasets we choose as feature space the initial instance space. Recently, a lot of effort has been devoted to developing algorithms which are able to train fast

on large datasets provided linear kernels are employed. Among such algorithms we choose for our experiments the Dual Coordinate Descent (DCD) algorithm (Hsieh et al., 2008), SVM^{perf} (Joachims, 2006), the first cutting-plane algorithm for training linear SVMs and Pegasos (Shalev-Schwartz et al., 2007), a stochastic gradient descent algorithm. The source for DCD (version 1.5) is available at <http://www.csie.ntu.edu.tw/~cjlin/liblinear>, for SVM^{perf} (version 2.50) at <http://smvlight.joachims.org> and for Pegasos at <http://ttic.uchicago.edu/~shai/code>.

The datasets we used for training are the binary Adult ($N=32561$, $M=123$) and Web ($N=49749$, $M=300$) UCI datasets as compiled by Platt, the training set of the KDD04 Physics dataset ($N=50000$, $M=70$ after removing the 8 columns containing missing features) obtainable from <http://kodiak.cs.cornell.edu/kddcup/datasets.html>, the Real-sim ($N=72309$, $M=20958$), News20 ($N=19996$, $M=1355191$) and Webspam (unigram treatment with $N=350000$ and $M=254$) datasets all available at <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>, the multiclass Covertypes ($N=581012$, $M=54$) UCI dataset and the full Reuters RCV1 dataset ($N=804414$, $M=47236$) obtainable from http://www.jmlr.org/papers/volume5/lewis04a/lyrl2004_rcv1v2_README.htm.

In the case of the Covertypes dataset we study the binary classification problem of the first class versus rest while for the RCV1 we consider both the binary text classification tasks of the C11 and CCAT classes versus rest. The Physics and Covertypes datasets were rescaled by multiplying all the features with 0.001. The experiments were conducted on a 2.5 GHz Intel Core 2 Duo processor with 3 GB RAM running Windows Vista. Our codes written in C++ were compiled using the g++ compiler under Cygwin.

In our experiments we chose $C_p = 1$. Also, we did

Table 2. Results of a low (1%) accuracy comparative study of the MPU, DCD, SVM^{perf} and Pegasos algorithms. We report training times in secs keeping the dataset ordering of Table 1.

MPU	0.1	0.1	0.1	0.3	1.6	1.2	6.2	3.1	3.7
DCD	0.2	0.1	0.1	0.2	0.6	1.4	9.1	3.5	3.6
SVM ^{perf}	1.5	0.3	0.3	0.8	12.8	7.4	76.5	12.9	19.3
Pegasos	4.0	1.0	0.3	3.2	10.6	13.7	141.9	27.9	31.6

not include any bias term (i.e. we did not perform any augmentation). The implementation of the MPU algorithm followed closely the description of Section 4 with active set parameters $\bar{c} = 1.01$, $N_{ep_1} = 3$ and $N_{ep_2} = 10$. Moreover, we did make use of the mechanism of early stopping. We also assumed $\delta b = 3R^2$ throughout. For Pegasos we chose $k = 30$ and $\lambda = N^{-1}$ (since $(\lambda N)^{-1}$ plays the role of C_p). In the first experiment our procedure was to obtain a very good approximation of the optimal objective \mathcal{J}_{opt} (i.e. $\frac{\delta \mathcal{J}}{\mathcal{J}_{opt}} \leq 10^{-4}$) using MPU with $\delta = 10^{-5}$ and $\delta_{stop} = 10^{-4}$ and then try to achieve comparable values of \mathcal{J} with the other linear SVM solvers. For DCD and SVM^{perf} we employed the “accuracy” ϵ (≤ 0.1) whereas for Pegasos the number of iterations iter in order to obtain the required value of \mathcal{J} . Table 1 contains our experimental results (objective \mathcal{J} and training runtime) for this high accuracy comparative study. SVM^{perf} is clearly considerably slower than DCD and MPU but much faster than Pegasos. DCD, instead, appears to be statistically the fastest but still its runtimes are very close to the ones of MPU. We also carried out a second experiment aiming at obtaining a value of \mathcal{J} which deviates from \mathcal{J}_{opt} by 1%. The results of this low accuracy study are reported in Table 2. We see that DCD and MPU are still the fastest with close runtimes followed by SVM^{perf}. The remarkable finding is, however, the spectacular improvement that Pegasos presents. In the light of our results we may conclude that the MPU algorithm is very competitive.

6. Conclusions

By enriching the classical Margin Perceptron with a mechanism of unlearning which allows the algorithm to recover from any wrong decisions taken on the way we succeeded in tackling both the hard margin and the 1-norm soft margin problems. We derived upper bounds on the number of steps required in order for the algorithm to achieve any desirable fixed before running approximation of the optimal solution for both tasks and described strategies for its efficient implementation. Finally, the experiments provided evidence that the new algorithm is a very competitive linear SVM.

References

- Boser, B., Guyon, I., and Vapnik, V. A training algorithm for optimal margin classifiers. In *COLT*, pp. 144–152, 1992.
- Cortes, C. and Vapnik, V. Support vector networks. *Machine Learning*, 20:273–297, 1995.
- Crammer, K., Kandola, J., and Singer, Y. Online classification on a budget. In *NIPS 16*, 2003.
- Cristianini, N. and Shawe-Taylor, J. *An Introduction to Support Vector Machines*. Cambridge University Press, 2000.
- Duda, R. O. and Hart, P. E. *Pattern Classification and Scene Analysis*. Wiley, 1973.
- Gentile, C. A new approximate maximal margin classification algorithm. *JMLR*, 2:213–242, 2001.
- Guyon, I. and Stork, D. Linear discriminant and support vector classifiers. In *Advances in Large Margin Classifiers*. MIT Press, 1999.
- Hsieh, C.-J., Chang, K.-W., Lin, C.-J., Keerthi, S. S., and Sundararajan, S. A dual coordinate descent method for large-scale linear SVM. In *ICML*, pp. 408–415, 2008.
- Joachims, T. Training linear SVMs in linear time. In *KDD’06*, pp. 217–226. ACM Press, 2006.
- Krauth, W. and Mézard, M. Learning algorithms with optimal stability in neural networks. *Journal of Physics A*, 20:L745–L752, 1987.
- Li, Y. and Long, P. The relaxed online maximum margin algorithm. *Machine Learning*, 46:361–387, 2002.
- Novikoff, A. B. J. On convergence proofs on perceptrons. In *Proc. Symp. Math. Theory Automata, Vol. 12*, pp. 615–622, 1962.
- Rosenblatt, F. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65 (6):386–408, 1958.
- Shalev-Schwartz, S., Singer, Y., and Srebro, N. Pegasos: Primal estimated sub-gradient solver for SVM. In *ICML*, pp. 807–814, 2007.
- Tsampouka, P. and Shawe-Taylor, J. Approximate maximum margin algorithms with rules controlled by the number of mistakes. In *ICML*, pp. 903–910, 2007.
- Vapnik, V. *Statistical Learning Theory*. Wiley, 1998.