



OLLAWV: OnLine Learning Algorithm using Worst-Violators

Gabriella Melki^a, Vojislav Kecman^a, Sebastián Ventura^{b,c}, Alberto Cano^{a,*}

^a Department of Computer Science, Virginia Commonwealth University, USA

^b Department of Computer Science and Numerical Analysis, University of Cordoba, Spain

^c Knowledge Discovery and Intelligent Systems in Biomedicine Laboratory, Maimonides Biomedical Research Institute of Cordoba, Spain

ARTICLE INFO

Article history:

Received 1 December 2017

Received in revised form 6 February 2018

Accepted 20 February 2018

Available online 25 February 2018

Keywords:

Machine learning

Online learning

Support vector machines

Worst-violators

Stochastic gradient descent

ABSTRACT

Due to the ever-growing nature of dataset sizes, the need for scalable and accurate machine learning algorithms has become evident. Stochastic gradient descent methods are popular tools used to optimize large-scale learning problems because of their generalization performance, simplicity, and scalability. This paper proposes a novel stochastic, also known as online, learning algorithm for solving the L1 support vector machine (SVM) problem, named *OnLine Learning Algorithm using Worst-Violators* (OLLAWV). Unlike other stochastic methods, OLLAWV eliminates the need for specifying the maximum number of iterations and the use of a regularization term. OLLAWV uses early stopping for controlling the size of the margin instead of the regularization term. The experimental study, performed under very strict nested cross-validation (a.k.a., double resampling), evaluates and compares the performance of this proposal with state-of-the-art SVM kernel methods that have been shown to outperform traditional and widely used approaches for solving L1-SVMs such as *Sequential Minimal Optimization*. OLLAWV is also compared to 5 other traditional non-SVM algorithms. The results over 23 datasets show OLLAWV's superior performance in terms of accuracy, scalability, and model sparseness, making it suitable for large-scale learning.

© 2018 Elsevier B.V. All rights reserved.

1. Introduction

Over the past decade, dataset sizes have grown disproportionately to the speed of processors and memory capacity, limiting statistical learning methods to computational time. The exponential growth of data, in size and complexity, has produced a pressing need to develop scalable machine learning algorithms to learn from data. Many real-world applications, such as human activity recognition, operations research, and video and signal processing [1,2], require algorithms that are scalable and accurate, while being able to process data and provide insightful information in a timely fashion. The scope of this paper is concerned with developing a unique learning algorithm for such large data problems without the use of parallelization techniques and distributed systems.

Support vector machines, proposed by Cortes and Vapnik [3], represent popular linear and non-linear (kernelized) learning algorithms based on the idea of a large-margin classifier. They have been shown to improve generalization performance for binary classification problems. SVMs are similar to other machine learning techniques, but literature shows that they usually outperform them

in terms of scalability, computational efficiency, and robustness against outliers [4,5]. They are known for creating sparse and non-linear classifiers which makes them suitable for handling large datasets. A traditional approach for training SVMs is the *Sequential Minimal Optimization* (SMO) algorithm [6], a method for solving the L1-SVM's Quadratic Programming (QP) task. Although SMO provides an exact solution to the SVM QP problem, its performance is highly dependent on the SVM hyperparameters. More recent approaches, which have been shown to surpass SMO in terms of scalability while remaining competitive in accuracy, include the *Minimal Norm SVM* (MNSVM) [7] and the *Non-Negative Iterative Single Data Algorithm* (NNISDA) [8,9].

To deal with issues of scalability, this paper introduces a different approach which focuses on the minimization of the regularized L1-SVM through Stochastic Gradient Descent (SGD), a well-known simple, yet efficient technique for learning classifiers under convex loss functions. Recently, SGD algorithms have been shown to have considerable performance and generalization capabilities in the context of large-scale learning [10], and have been used to solve the SVM problem, such as *NORMA* [11] and *PEGASOS* [12,13]. Although stochastic and iterative algorithms are very simple to implement and efficient, they also have their limitations. One of these limitations is the lack of meaningful stopping criteria for the algorithm; without a pre-specified number of iterations to train, the

* Corresponding author.

E-mail address: acano@vcu.edu (A. Cano).

algorithms continue running [14]. Another limitation stems from the superlinear increase in training time as the number of samples increases. Incremental algorithms attempt to alleviate this issue, but they cannot guarantee a bound on the number of updates per iteration [15,16].

To address the limitations presented by current popular SVM solvers, this paper proposes a novel *OnLine Learning Algorithm using Worst-Violators* (OLLAWV). This unique method iterates over samples, updates the model, and utilizes a novel stopping criterion. The model is updated by iteratively selecting (without replacement) the worst violating sample, i.e. the sample with the largest error according to the current decision function, and stops training when there are no more violating samples left to update. In other words, the algorithm is implicitly identifying support vectors and stopping when it has found them all. Because samples are selected and updated without replacement, coupled with the fact that the maximum number of iterations never exceeds the size of the dataset, OLLAWV does not use the regularization updating term. Instead, the regularization is achieved by early stopping. In [17], it has been shown that the smaller the number of updates (determined by here presented stopping criterion) is, the larger will be the margin. On the other hand, the larger the margin, the better generalization of the model is. The experimental results presented here confirm both the theoretical statements in [17] and the validity of the approach proposed and taken in OLLAWV algorithm. Combining this method of updating the model and stopping criteria with the fact that SVMs are known for creating sparse kernel classifiers, the contribution aims to speed up the model training time without sacrificing the model's accuracy. The key contributions of this work include:

- Devising a unique iterative procedure for solving the L1-SVM problem, as well as a novel method for identifying support vectors, or *worst-violators*. Rather than randomly iterating over the data samples, OLLAWV aims to reduce training time by selecting and updating the samples that are most incorrectly classified with respect to the current decision hyper-plane.
- Designing a novel stopping criteria by utilizing the worst-violator identification method. This aims to eliminate the added parameterization that is included with most online methods, where the number of iterations of the algorithm needs to be set in advance. Once there are no incorrectly classified samples left, the algorithm terminates.

This work is organized as follows. Section 2 presents the notation used throughout the paper, reviews the support vector machine problem, and describes some popular SVM solvers. Section 3 presents the contribution: the algorithm design, method for selecting and updating worst-violators, and the stopping criteria. Section 4 presents the experimental environment, including results based on accuracy, runtime, and percentage of support vectors, over 23 datasets compared with current state-of-the-art algorithms. Finally, Section 5 presents the conclusions of this contribution.

2. Background

This section defines the notation that will be used throughout the paper, describes the support vector machine problem, and reviews related works on current popular solvers and online learning algorithms aimed at training support vector machines.

2.1. Notation

Let \mathcal{D} be the full training dataset of n d -dimensional samples. Let $\mathbf{Y} \in \mathcal{D}$ be a vector of n labels corresponding to each sample, such that $\mathbf{Y} \in \{-1, 1\}^n$. In the non-binary (more than two classes)

Table 1

Summary of notation used throughout the paper.

Definition	Notation
Number of samples	n
Number of input attributes	d
Input space	$\mathbf{X} \in \mathbb{R}^{n \times d}$
Labels	$\mathbf{Y} \in \{-1, 1\}^n$
Sample i	$\mathbf{x}_i = (x_{i1}, \dots, x_{id}), \forall i \in \{1, \dots, n\}$
Sample label i	$y_i \in \{-1, 1\}, \forall i \in \{1, \dots, n\}$
Full training dataset	$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_i, y_i), \dots, (\mathbf{x}_n, y_n)\}$

classification cases, $\mathbf{Y} \in \mathbb{Z}^n$. Let $\mathbf{X} \in \mathcal{D}$ be a matrix consisting of n samples that are d -dimensional, $\mathbf{X} \in \mathbb{R}^{n \times d}$. Table 1 provides a summary of the notation used in this paper.

2.2. Support vector machines

Support vector machines represent a popular set of learning techniques that have been introduced under Vapnik–Chervonenkis theory and *structured risk minimization* (SRM) [3,18–21]. They minimize the expected probability of classification error, resulting in a generalized model without making assumptions about the data distribution [19]. In the context of classification, SVMs are particularly useful for finding linear predictors, or hyper-planes, in high dimensional feature spaces, which is a computationally complex learning problem [3]. This problem is approached by searching for the optimal *maximal margin* of separability between classes, and is equivalent to solving the following optimization problem:

$$\min_{(\mathbf{w}, b) \in \mathcal{H}_o \times \mathbb{R}} R = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n L(y_i, o_{(\mathbf{w}, b)}(\mathbf{x}_i)), \quad (1)$$

where \mathcal{H}_o is a general Hilbert space, $o_{(\mathbf{w}, b)}$ is the output (decision) function representing an affine mapping function $o : \mathbb{R}^d \rightarrow \mathbb{R}$ that approximately minimizes the risk and is parameterized by \mathbf{w} and b . The weights $\mathbf{w} \in \mathbb{R}^d$ and bias term $b \in \mathbb{R}$ define the linear predictor. The loss function used for the L1-SVM is the Hinge loss function: $L(y_i, o_{(\mathbf{w}, b)}(\mathbf{x}_i)) = \max\{0, 1 - y_i o_{(\mathbf{w}, b)}(\mathbf{x}_i)\}$, which penalizes errors satisfying the following: $y_i o_{(\mathbf{w}, b)}(\mathbf{x}_i) < 1$ and is a crucial element that facilitates the SVM model's sparseness. The penalty parameter $C \in \mathbb{R}$ controls the trade-off between margin maximization and classification error minimization, penalizing large norms and errors. To handle cases when the data are non-linearly separable, while enhancing the classifier's generalization capabilities, a kernel function can be used [22], as shown in Eq. (2):

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle, \quad (2)$$

where $\phi(\cdot)$ represents a mapping function from the original feature space to a higher dimensional space. The advantage of utilizing kernels is being able to calculate the inner product in the input space rather than in the very high feature dimensional space (including the infinite dimensional ones). The SVM model output, o shown in Eq. (3), for a given input vector \mathbf{x} is defined by the kernel as given below:

$$o(\mathbf{x}) = \sum_{i=1}^n \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i) + b, \quad (3)$$

where $\alpha_i \in \mathbb{R}$ are the coefficients, or weights, of the expansion in feature space, and $b \in \mathbb{R}$ is the so-called bias term. Note that if a positive definite kernel is used, there is no need for a bias term b , but b can nevertheless be used. This is why OLLAWV can be implemented with or without the bias term b , as shown in Section 3. The two terms, α and b , parameterize the SVM model. A model is called *dense* if the absolute value of all its weights are greater than

0, while a *sparse* model would be one that contains some $\alpha_i = 0$. The level of sparseness may vary, but the sparser the model, the more scalable the applications.

2.3. Methods for solving the SVM problem

Although support vector machines represent a major development in machine learning algorithms, in the case of large-scale problems (hundreds of thousands to several millions of samples), the design of SVM training algorithms still has room for improvement [23,24]. So far, there have been several different approaches for tackling large-scale SVM classification problems.

The first attempts at speeding up the training time and decreasing algorithm memory requirements were aimed at decomposing the underlying SVMs quadratic programming (QP) problem. First, Boser et al. [20] implemented Vapnik's *chunking* method. *Sequential Minimal Optimization* by Platt [6] and its improvement by Keerthi et al. [25] are an alternative approach to decomposing the QP problem, and are implemented in popular, widely used software package LIBSVM [26]. SMO is an iterative procedure that divides the SVM dual problem into a series of sub-problems, which are solved analytically by finding the optimal α values that satisfy the Karush–Kuhn–Tucker conditions [27]. Although SMO is guaranteed to converge, heuristics are used to choose α values in order to accelerate the convergence rate. This is a critical step because the convergence speed of the SMO algorithm is highly dependent on the dataset size and SVM hyperparameters [21].

Some advancements in handling large scale problems are based on a geometric interpretation of SVM problem. Some of these *geometric SVMs* include approaches that use convex hulls [28] and minimum enclosing balls such as *Core Vector Machines* (CVM) [29]. Tsang et al. [30] then improved the scalability of CVMs by introducing *Ball Vector Machines* (BVM) which do not require a QP solver. Other geometric approaches include the novel algorithms introduced by Strack [7], known as the *Sphere Support Vector Machine* (SphereSVM) and *Minimal Norm Support Vector Machine* (MNSVM), which utilize the connection between minimal enclosing balls and convex hull problems, while demonstrating a high capability for learning from large datasets. The *Non-Negative Iterative Single Data Algorithm* [9] is an efficient approach for solving the SVM problem, shown to be faster than SMO and equal in terms of accuracy [8]. NNISDA is an iterative algorithm that finds a solution to the L2-SVM using coordinate descent, inspired by ISDA [31] which was originally introduced by Kecman et al. [32].

Recently, several authors have proposed the use of a standard stochastic gradient descent (SGD) approach for SVMs to optimize large-scale learning problems [12,19,21,33,34]. Kivinen et al. [11] and Bousquet and Bottou [35] showed that stochastic algorithms can be both the fastest, and have the best generalization performances. Shalev-Shwartz and Ben-David [19] have also demonstrated that the basic SGD algorithm is very effective when data are sparse, taking less than linear $[O(d)]$ time and space per iteration to optimize a system with d parameters. It can greatly surpass the performance of more sophisticated batch methods on large data sets. The previously mentioned approaches are extended variants of a classic kernel perceptron algorithm [17].

Notable representatives of this method of learning include the *Naïve Online R Minimization Algorithm* (NORMA) by Kivinen et al. [11] and the *Primal Estimated Sub-Gradient Solver for SVM* (PEGASOS) by Shalev-Shwartz et al. [12]. NORMA is an online kernel based algorithm designed to utilize SGD for solving the SVM problem, exploiting the kernel trick in an online setting. It can be regarded as a generalization of the kernel perceptron algorithm with regularization [11]. PEGASOS solves the primal SVM problem using stochastic sub-gradient descent, implementing both linear and non-linear kernels, and showed that the algorithm does not directly

depend on the size of the data, making it suitable for large-scale learning problems.

A more recent approach that is the inspiration behind this contribution, named *OnLine Learning Algorithm* (OLLA) [36] is a unification, simplification, and expansion of the somewhat similar approaches presented in [12,19,21,33,34,37,38]. This algorithm is unique because it is not only designed to optimize the SVM cost function, but also the cost functions of several other popular non-linear (kernel) classifiers using SGD in the primal domain. Collobert and Bengio [17] provided justification for not using regularization, and thus OLLA was designed to handle cost functions with and without the regularization term. Comparisons of performances of OLLA with the popular SMO algorithm highlighted the merits of OLLA in terms of speed, as well as accuracy, when the number of samples was increased, making it suitable for large-scale learning. Comparisons using various different classifiers against SMO were also shown, but for the scope of this paper the L1-SVM was mentioned [36].

Although the SGD approaches mentioned above have many merits when it comes to solving large-scale machine learning problems, stochastic procedures also have their disadvantages. One of them stems from the lack of meaningful stopping criteria. The only specified stopping criteria is a user defined input for the number of iterations, which gives rise to the question of what it should be set to. Another disadvantage of kernelized online algorithms is that the training time for each update increases superlinearly with the number of samples. This paper aims to deal with these limitations without sacrificing accuracy for scalability with respect to large-scale problems, while maintaining applicability to small and medium sized datasets.

3. OnLine Learning Algorithm using Worst-Violator

OLLAVV is an iterative, online learning algorithm for solving the L1-SVM problem using a novel model update procedure while implementing a self-stopping condition. The inspiration behind OLLAVV came from [36] which presented a generic online learning algorithm tailored not only for SVMs, but also for various other popular classifiers that use different risk functions, with or without a regularization term. The difference, novelty, and advantage of OLLAVV resides in its iterative method, where the weight α_i of the most violating sample i.e., of the *worst-violator*, is only updated in each iteration. A worst violating sample is defined as the sample that has the largest error with respect to the current decision function. Rather than randomly selecting samples to update per iteration, OLLAVV selects (without replacement) the most incorrectly classified sample and updates the model accordingly. By iteratively updating the model using only the worst-violator, the model is essentially finding its support vectors, as well as implicitly defining a stopping criterion. If there are no more violating samples, the algorithm terminates, eliminating the need to define the number of iterations for an algorithm to perform before returning the model, as is the case with most state-of-the-art online algorithms.

At every iteration, the algorithm selects a worst violating sample that has not been previously chosen, stores its index in vector \mathbf{S} , and then updates the model. Eq. (4) shows the method for selecting the worst-violator, where $y_0 \in \mathbb{R}$ is the error value, $wv \in \{1, \dots, n\}$ is the error value's index, $\mathbf{o} \in \mathbb{R}^n$ is the decision function output, and \neg is the 'not' symbol. For the L1-SVM, an error value will always be negative which is why the minimum function is used (i.e. the most negative output value or incorrectly classified sample). The worst violating sample becomes the model's support vector because its weight is updated and non-zero. Therefore, the number of iterations of OLLAVV is equal to the number of support vectors that the resulting model has. This is an interesting property of OLLAVV; if

the number of iterations is set beforehand, one is implicitly setting a bound on the number of support vectors.

$$[y_0, wv] = \min \{y_{wv} \cdot o_{wv}\}, \forall wv \in \{-S\} \tag{4}$$

Algorithm 1 lists OLLAWV’s pseudocode and Fig. 2 illustrates the steps taken by OLLAWV. First, the model parameters (α, b, S) and the algorithm variables (\mathbf{o} , iteration counter (t), initial worst-violator index wv and it’s error y_0) are first initialized. The worst-violator with respect to the current hyperplane is then found and the model parameters are updated. Once no more violating samples are found or the maximum number of iterations is reached, the model is returned.

OLLAWV performs stochastic gradient descent on the primal L1-SVM objective function given below:

$$\min_{\mathbf{w} \in \mathcal{H}_o \times \mathbb{R}} = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max \{0, 1 - y_i o_{(w)}(\mathbf{x}_i)\}, \tag{5}$$

where $o_{(w)}(\mathbf{x}_i) = \langle \mathbf{w}, \mathbf{x}_i \rangle$ is a linear predictor, with the bias term excluded for simplicity. Note that the loss function used in Eq. (5) is non-differentiable, but it has a subgradient due to a knick-point at $y_0 = 1$. The loss function’s gradient after the knick-point equals zero, which leads to a sparse model. Hence, when the value of $y_0 \geq 1$ the loss is zero, and for $y_0 < 1$ the loss increases linearly. The subgradient of the above cost function is given by:

$$\frac{\partial R}{\partial \mathbf{w}} = \begin{cases} \mathbf{w} - C \sum_{i=1}^n y_i \mathbf{x}_i & y_i o_i < 1 \\ \mathbf{w} & \text{otherwise.} \end{cases} \tag{6}$$

In the stochastic case, the calculation of the gradient needed for the weight update, is *pattern based*, not *epoch based* as in batch gradient descent. It has been shown [18] that the ideal gradient is equal to the sum of the gradients calculated after each sample is presented for fixed weights during the whole epoch. Thus, the stochastic update of \mathbf{w} from the subgradient shown in Eq. (6) becomes:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial R}{\partial \mathbf{w}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \begin{cases} Cy_i \mathbf{x}_i - \mathbf{w} & y_i o_i < 1 \\ -\mathbf{w} & \text{otherwise,} \end{cases}$$

where $\eta > 0 \in \mathbb{R}$ is the learning rate. According to the Representer Theorem, a vector $\alpha \in \mathbb{R}^n$ exists such that $\mathbf{w} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$ is an optimal solution to Eq. (5), where $\phi(\cdot)$ is a mapping from feature space to Hilbert space [19]. On the basis of the Representer Theorem, Eq. (5) can be optimized with respect to α instead of \mathbf{w} . By expressing \mathbf{w} this way and mapping input sample \mathbf{x}_i to $\phi(\mathbf{x}_i)$, the kernelized SGD update becomes:

$$\sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) \leftarrow \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) + \eta \begin{cases} Cy_i \phi(\mathbf{x}_i) - \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) & y_i o_i < 1 \\ -\sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) & \text{otherwise.} \end{cases}$$

However, OLLAWV optimizes in a stochastic manner, resulting in the following update:

$$\forall i : \alpha_i \phi(\mathbf{x}_i) \leftarrow \alpha_i \phi(\mathbf{x}_i) + \eta \begin{cases} (Cy_i \phi(\mathbf{x}_i) - \alpha_i \phi(\mathbf{x}_i)) & y_i o_i < 1 \\ (-\alpha_i \phi(\mathbf{x}_i)) & \text{otherwise} \end{cases}$$

$$\forall i : \alpha_i \leftarrow \alpha_i + \eta \begin{cases} (Cy_i - \alpha_i) & y_i o_i < 1 \\ (-\alpha_i) & \text{otherwise.} \end{cases}$$

The case when the worst violating sample is correctly classified $y_0 \geq 1$ is OLLAWV’s termination condition, i.e. is used as the stopping criterion in the algorithm. Hence the update for α is reduced to the following:

$$\forall i : \alpha_i \leftarrow \alpha_i + \eta (Cy_i - \alpha_i) \tag{7}$$

If the bias term b is included in Eq. (5), it’s stochastic update is as follows:

$$\forall i : b \leftarrow b + \eta \frac{Cy_i}{n} \tag{8}$$

In this experimental study, $\eta = 2/\sqrt{t}$ is used, where t is the current iteration; however, other learning rates such as $\eta = 1/t$ can also be used. Let $\Lambda = \eta Cy_i$ and $P = \eta \alpha_i$ be the update parameters for OLLAWV, and the α update can be expressed as: $\alpha_i \leftarrow \alpha_i + (\Lambda - P)$. Note that Λ is the update resulting from the loss function and P is derived from the regularizer term in Eq. (5). In the case of OLLAWV, $P = 0$ because samples are never updated more than once and their initial α value is always 0. It is important to note that in OLLAWV’s case, Λ never equals 0 because the samples being updated are worst-violators, meaning they are misclassified or incur some loss. The values of the decision function (output vector $\mathbf{o} \in \mathbb{R}^n$ in Algorithm 1), from which a worst-violator is found, changes per iteration based on the influence of the support-vectors that have been previously updated. From Eq. (3), the output vector update becomes the following:

$$\mathbf{o} \leftarrow \mathbf{o} + \Lambda * \mathcal{K}(\mathbf{x}_{-S}, \mathbf{x}_{wv}, \gamma) + B \tag{9}$$

where $\mathcal{K}(\cdot)$ is the Gaussian radial basis function (RBF) kernel, $\gamma \in \mathbb{R}$ is it’s parameter, and $B = (\Lambda * \beta)/n$ denotes the bias update. Only non-support vector output values are calculated per iteration, as denoted by \mathbf{x}_{-S} in the kernel function, because samples are never selected to be updated more than once. Because the output values scale with the C value, the stopping criteria for OLLAWV is also set to scale with C , rather than the classic formulation $y_i o_i \geq 1$. If the value of C is very large, $y_i o_i$ will never be greater than 1 and the algorithm will never terminate. Therefore, the stopping criteria is set to be $y_i o_i \geq M$, where $M \in \mathbb{R}$ is a scaled value of C . For the B calculation in Eq. (9), $\beta \in \{0, 1\}$ indicates whether the bias term is to be used. If b is not a part of the model, it should be omitted from Eqs. (3) and (9) by setting $\beta = 0$, otherwise $\beta = 1$.

OLLAWV is a stochastic gradient method (SGM) that has a convex cost function. Its learning rate coefficient can decrease linearly or semi-linearly during the learning stage. Hence, OLLAWV shares the complexity characteristics of SGM methods. Primarily, it can achieve linear convergence, making it a particularly convenient and practical method for solving very large machine learning problems. OLLAWV also works over a cost function without local minima, always leading towards the global minimum, even though it stops the learning process as soon as all samples are outside the prescribed margin. Fig. 1 shows the decision boundaries achieved by OLLAWV versus SMO (implemented within LIBSVM) and Bayes for toy datasets.

Algorithm 1. OnLine Learning Algorithm using Worst-Violators (OLLAWV)

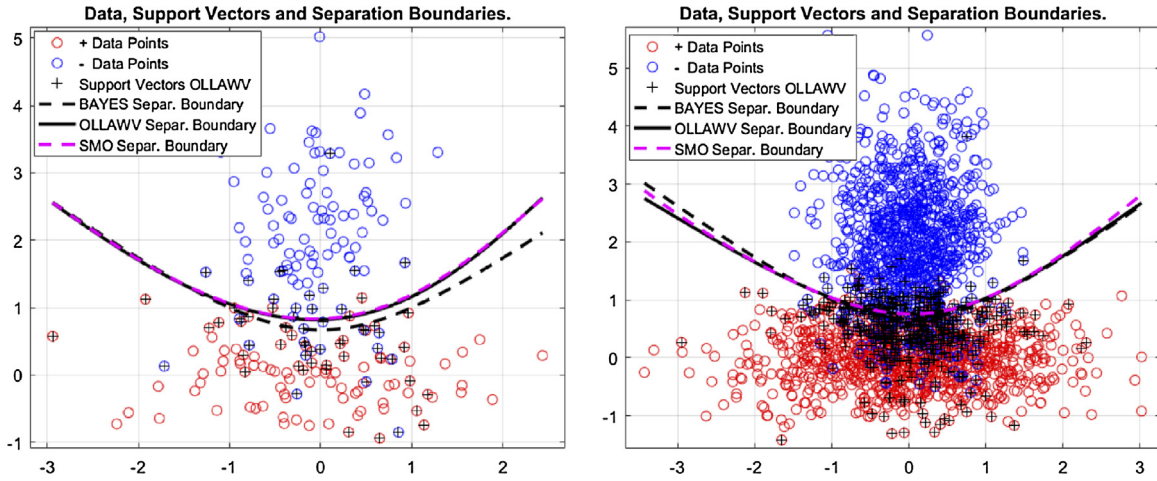


Fig. 1. A case of classifying 2-dimensional normally distributed data with different covariance matrices (left) for 200 and (right) 2000 data points. The theoretical separation boundary (denoted as the Bayes Separation Boundary) is quadratic and is shown as the dashed black curve. The other two separation boundaries shown are the ones obtained by OLLAWV and SMO (implemented within LIBSVM), respectively. In this particular case (left), the difference between the OLLAWV boundary and the SMO boundary is hardly visible. The case presented on the right shows that, with an increase of training samples, the OLLAWV and SMO boundaries converge to the theoretical Bayesian solution.

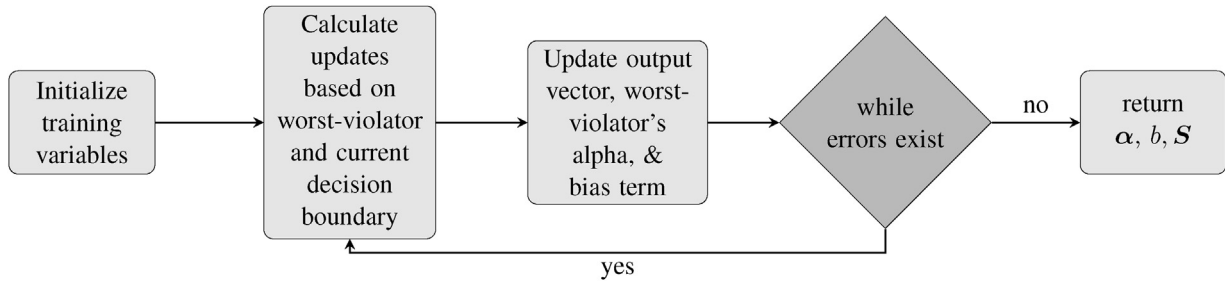


Fig. 2. A summary of the steps performed by OLLAWV. The model parameters (α, b, S) and the algorithm variables $(\mathbf{o}, t, wv, \text{ and } y_o)$ are first initialized. The worst-violator with respect to the current hyperplane is then found and the model parameters are then updated. Once no more violating samples are found, the model is returned.

Input: D, C, γ, β, M

Output: α, b, S

```

1:  $\alpha \leftarrow \mathbf{0}, b \leftarrow 0, S \leftarrow \mathbf{0}$            ▷ Initialize OLLAWV model parameters
2:  $\mathbf{o} \leftarrow \mathbf{0}, t \leftarrow 0$              ▷ Initialize the output vector and iteration counter
3:  $wv \leftarrow 0, y_o \leftarrow y_{wv} * \mathbf{o}_{wv}$    ▷ Initialize hinge loss error and worst-violator index
4: While  $y_o < M$  do
5:    $t \leftarrow t + 1$ 
6:    $\eta \leftarrow 2/\sqrt{t}$                        ▷ Learning rate
7:
8:    $\Lambda \leftarrow \eta * C * y_{wv}$              ▷ Calculate hinge loss update
9:    $B \leftarrow (\Lambda * \beta) / n$              ▷ Calculate bias update
10:   $\mathbf{o} \leftarrow \mathbf{o} + \Lambda *$              ▷ Update output vector
11:   $\mathcal{K}(\mathbf{x}_{-S}, \mathbf{x}_{wv}, \gamma) + B$ 
12:   $\alpha_{wv} \leftarrow \alpha_{wv} + \Lambda$          ▷ Update worst-violator's alpha value
13:   $b \leftarrow b + B$                          ▷ Update bias term
14:   $S_t \leftarrow wv$                            ▷ Save index of worst-violator
15:   $[y_o, wv] \leftarrow \min_{wv \in \{-S\}} [y_{wv} * \mathbf{o}_{wv}]$    ▷ Find the worst-violator
16: end while
  
```

4. Experimental environment, results, and analysis

This section presents two experimental setups of our contribution against other state-of-the-art algorithms on 23 different benchmark datasets. The first study, presented in Section 4.1, compares OLLAWV to two other SVM kernel methods, and the second compares OLLAWV to 5 non-SVM methods, shown in Section 4.2. In each section, the experimental setups are first described and

the state-of-the-art methods are listed. The results and statistical analysis [39] are then presented and analyzed. The main aim of the experiments is to compare our contribution to other support vector machine solvers that have been shown to surpass popular and widely used SVM kernel methods in terms of memory consumption, runtime, and accuracy. The supplemental experimental study in 4.2 was conducted to emphasize the better performance of OLLAWV against non-SVM algorithms.

Table 2 presents a summary of the 23 datasets used throughout the experiments, where the number of attributes (dimensionality), classes, and samples are shown. The datasets used and the results obtained are divided into three groups: *small*, *medium* and *large*. The datasets were obtained from the UCI Machine Learning repository¹ [40], and the LibCVM² [29] and LIBSVM³ [26] sites.

4.1. SVM experimental setup

The experimental setup was designed to evaluate differences in performance of the proposed OLLAWV method against the state-of-the-art algorithms: *Minimal Norm SVM* (MNSVM) [7] and *Non-Negative Iterative Single Data Algorithm* (NNISDA) [9]. These algorithms were chosen because they have shown considerable performance in runtime, memory consumption, and accuracy against the popular and widely used LIBSVM and LibCVM packages. In [7], it was shown that MNSVM outperforms both the L1 and

¹ <http://archive.ics.uci.edu/ml/index.php>.

² <http://c2inet.sce.ntu.edu.sg/ivor/cvm.html>.

³ <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>.

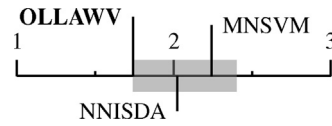
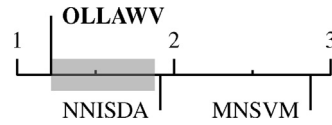
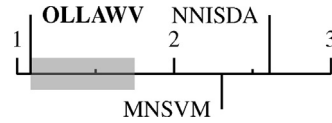
Table 2
Datasets.

Dataset	# Samples	# Attributes	# Classes
<i>Small datasets</i>			
iris	150	4	3
teach	151	5	3
wine	178	13	3
cancer	198	32	2
sonar	208	60	2
glass	214	9	6
vote	232	16	2
heart	270	13	2
dermatology	366	33	6
prokaryotic	997	20	3
eukaryotic	2,427	20	2
<i>Medium datasets</i>			
optdigits	5,620	64	10
satimage	6,435	36	6
usps	9,298	256	10
pendigits	10,992	16	10
reuters	11,069	8,315	2
letter	20,000	16	26
<i>Large datasets</i>			
adult	48,842	123	2
w3a	49,749	300	2
shuttle	58,000	7	7
web (w8a)	64,700	300	2
ijcnn1	141,691	22	2
intrusion	5,209,460	127	2

L2 implementations of LIBSVM, and BVM embedded in LibCVM. NNISDA was then compared to MNSVM in [9], and showed an added improvement in runtime performance. MNSVM was implemented in an open source C++ framework called “GSVM Command Line Tool for Geometric SVM Training⁴”. Both, NNISDA and OLLAWV were implemented as additional modules within Strack–Kecman’s code, keeping the experimental environment controlled for all three algorithms. The experiments for all methods were run on the same computer containing two Intel Xeon X5680 CPUs (6-core, 3.33 GHz) and 96 GB of RAM.

Experiments were performed using double, or nested, 5-fold cross-validation in order to objectively evaluate the models’ performances and tune hyper-parameters [42,43]. In the outer loop, the data are separated into 5 equally sized folds and each part is held out in turn as the test set, and the remaining four parts are used as the training set. In the inner loop, 5-fold cross-validation is also used over the training set assigned by the outer loop, where the best hyper-parameters are chosen. The best model obtained by the inner loop is then applied on the outer loop’s test set. This procedure ensures the model’s performance is not optimistically biased as when using a single loop of k -fold cross-validation. It ensures the class labels of the test data will not be seen when tuning the hyper-parameters, which is consistent with real-world applications. Obviously, such a rigorous procedure is computationally expensive, but the goal is to fairly compare different classification models on the same data sets, with the same cross-validation procedure, and hyper-parameters. First, the datasets were normalized by linear transformation of the feature values to the range [0, 1]. Then, the training process, also involving model selection using pattern search, was performed. The best hyper-parameters were chosen from the following 8×8 possible combinations, shown in Eqs. (10a) and (10b), and were also used for the competing SVM methods.

$$C \in \{4^n\}, \quad n = \{-2, \dots, 5\} \quad (10a)$$

**Fig. 3.** Bonferroni–Dunn test for accuracy.**Fig. 4.** Bonferroni–Dunn test for runtime.**Fig. 5.** Bonferroni–Dunn test for % support vectors.

$$\gamma \in \{4^n\}, \quad n = \{-5, \dots, 2\} \quad (10b)$$

The γ parameter refers to that of the Gaussian RBF kernel, given by:

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}. \quad (11)$$

To deal with multi-class classification problems, the one-vs-one, or pairwise, approach was used. The pairwise training procedure trains $c(c-1)/2$ binary classifiers, a classifier for each possible pair of classes, where c is the number of classes. During the prediction phase, a voting scheme is used where all $c(c-1)/2$ models predict an unseen data sample and the class that received the highest number of votes is considered to be the samples true class.

4.1.1. Comparison results and analysis

The classification performance was measured using the following metrics: accuracy, runtime, and the percentage of support vectors (size of the model). Table 3 displays the results for OLLAWV and the two state-of-the-art methods. The percentage of support vectors was reported for analyzing the complexities of the resulting models over the variously sized datasets. In order to analyze the performances of the multiple models, non-parametric statistical tests are used to validate the experimental results obtained [39]. The Iman–Davenport non-parametric test is run to investigate whether significant differences exist among the performance of the algorithms by ranking them over the datasets used, using the Friedman test. The algorithm ranks for each metric are presented in the last row of Table 3, and the lowest (best) rank value is typeset in bold. After the Iman–Davenport test indicates significant differences (with p -value = 0.2397 for accuracy, and p -value = 0 for runtime and percent support vectors), the Bonferroni–Dunn post-hoc test is then used to find where they occur between algorithms by assuming the classifiers’ performances are different by at least some critical value (critical distance is 0.66 for $\alpha = 0.05$). Below Table 3, Figs. 3–5 highlight the critical distance (in gray) from the best ranking algorithm to the rest. The algorithms to the right of the critical distance bar perform statistically significantly worse than the control algorithm, OLLAWV.

The results in Table 3 indicate that OLLAWV outperforms NNISDA and MNSVM in terms of accuracy, runtime, and model complexity. Although the differences in accuracy between the methods is not very large, on average, OLLAWV is about 2 times faster than NNISDA and MNSVM. As mentioned previously, OLLAWV aims to speed up the learning process without sacrificing the model’s accuracy. This stems from OLLAWV’s ability to produce

⁴ Strack–Kecman, <https://github.com/strackr/gsvm>.

Table 3
Comparison of OLLAWV vs. NNISDA and MNSVM.

Dataset	Accuracy (%)			Runtime (s)			Support vectors (%)		
	OLLAWV	NNISDA	MNSVM	OLLAWV	NNISDA	MNSVM	OLLAWV	NNISDA	MNSVM
<i>Small datasets</i>									
iris	97.33	94.00	96.67	0.05	0.27	3.57	13.50	40.20	29.80
teach	52.32	52.31	52.95	0.12	0.44	8.85	69.19	99.80	87.40
wine	98.87	96.60	96.60	0.28	0.43	4.84	15.02	44.40	48.60
cancer	80.36	81.86	81.38	0.49	0.85	4.46	42.79	83.80	89.60
sonar	92.32	89.48	87.57	0.59	0.98	3.03	31.26	73.00	66.00
glass	72.41	67.81	69.30	0.46	1.01	11.94	62.84	90.80	87.60
vote	96.54	96.11	93.99	0.26	0.46	1.49	13.36	33.20	34.00
heart	82.22	83.33	83.33	0.50	0.91	6.45	37.69	73.00	82.00
dermatology	97.82	98.36	98.36	1.62	2.47	11.68	36.94	59.00	59.80
prokaryotic	88.96	88.86	88.97	6.09	10.64	50.86	29.01	51.20	49.00
eukaryotic	77.38	79.56	81.21	61.95	49.16	342.76	54.11	76.40	72.60
<i>Medium datasets</i>									
optdigits	99.11	99.29	99.31	411	528	787	28.64	31.60	30.60
satimage	91.66	92.39	92.35	1334	687	1094	20.72	45.00	44.80
usps	97.49	98.05	98.24	10,214	5245	7777	11.22	29.40	28.00
pendigits	99.56	99.62	99.61	723	909	1500	10.27	17.60	16.60
reuters	98.03	98.08	97.99	954	1368	1657	8.770	18.20	18.60
letter	96.99	99.11	99.13	5259	12,009	26,551	43.56	57.60	56.60
<i>Large datasets</i>									
adult	84.75	85.07	85.13	21,025	72,552	123,067	34.66	56.00	56.60
w3a	98.86	98.82	98.82	6532	15,951	24,562	3.270	14.60	12.40
shuttle	99.77	99.83	99.87	2833	7420	45,062	2.010	6.00	16.40
web	98.94	99.00	99.00	12,067	30,583	38,040	4.320	13.20	10.80
ijcnn1	98.31	99.34	99.41	162,587	296,917	370,144	16.36	11.00	7.600
intrusion	99.77	99.67	99.66	2,402,804	4,646,810	3,772,113	0.780	2.000	1.700
Average	91.29	91.15	91.25	114,209	221,350	191,861	25.66	44.65	43.79
Ranks	1.739	2.022	2.239	1.217	1.913	2.869	1.087	2.609	2.304

The best results are bolded.

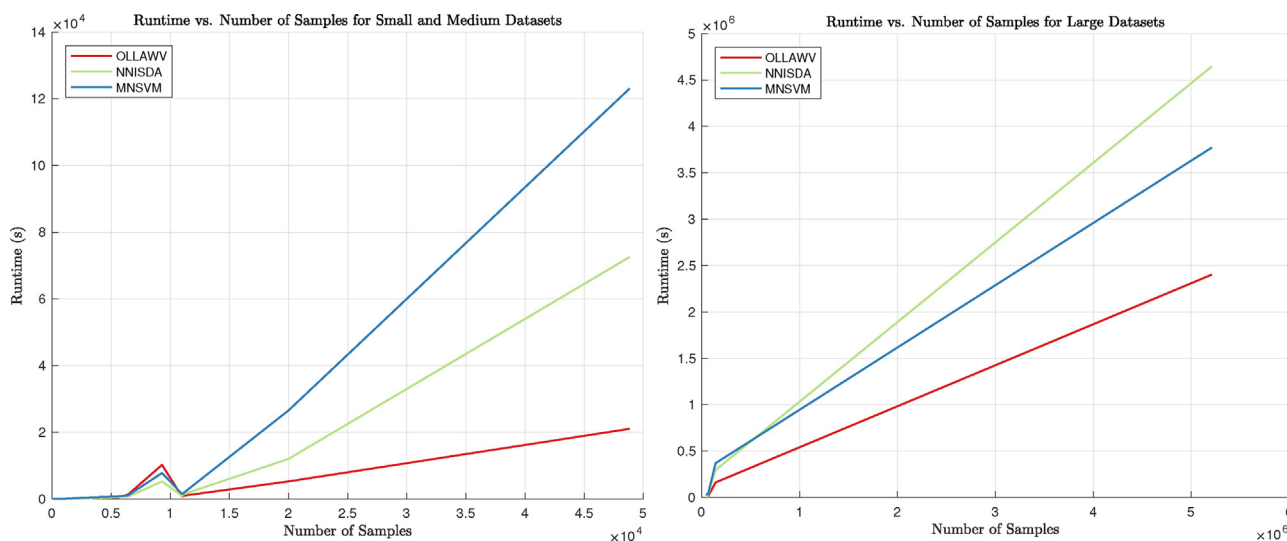


Fig. 6. Runtime in seconds versus the number of samples, divided into two groups: small & medium (left) versus large (right). Note OLLAWV's gradual increase in runtime as the number of samples increases compared to NNISDA and MNSVM's steeper change. In almost all cases, OLLAWV displays superior runtime over state-of-the-art. Runtime depends upon many characteristics: dimensionality, class-overlapping, complexity of the separation boundary, number of classes as well as upon the number of support vectors, which partly explains the tiny bump in the left figure.

sparse models, as is shown by the averaged percentage of support vectors. The speedup that OLLAWV presents is proportional to the model complexity and the experimental results show that OLLAWV produces, on average, models that are 1.7 times smaller than the two state-of-the-art methods used. This highlights the applicability and advantage that OLLAWV has for learning from large datasets.

Figs. 3–5 show the results of the statistical analysis for accuracy, runtime, and percentage of support vectors. Figs. 4 and 5 show that OLLAWV is statistically significantly better than MNSVM and NNISDA for runtime and percentage of support vectors (model

size). At the same time, Fig. 3 emphasizes what was mentioned earlier: OLLAWV is shown to speed up the learning process without sacrificing model accuracy against the state-of-the-art methods used.

Fig. 6 plots the correlation of OLLAWV, NNISDA, and MNSVM's runtime versus number of samples for the small, medium, and large datasets. The figure clearly emphasizes the benefit of using OLLAWV for large-scale learning due to its gradual increase in runtime as the number of samples increases in comparison to NNISDA and MNSVM. Fig. 7 shows the correlation between OLLAWV,

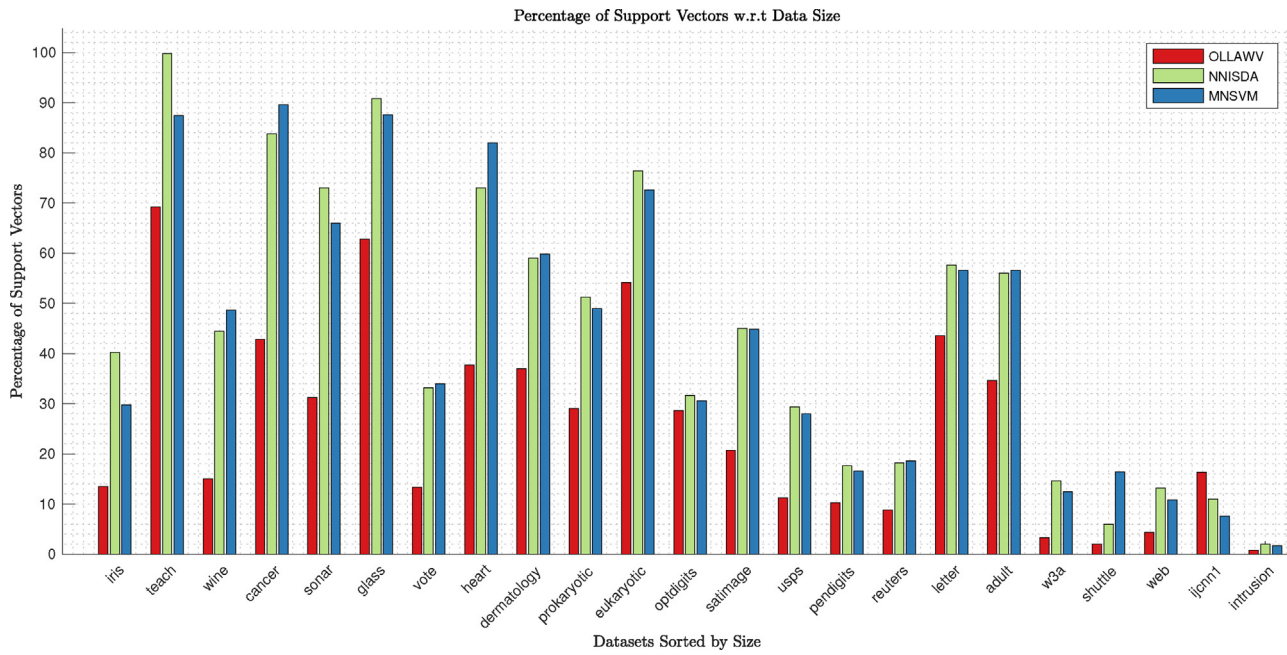


Fig. 7. Size of the model given as percentage of support vectors with respect to the number of samples versus the number of samples. Note that OLLAWV's percentage of support vectors is always smaller (except in one case) than NNISDA's and MNSVM's ones.

Table 4
Non-SVM algorithm hyper-parameters.

Algorithm	Parameters
<i>k</i> -NN	Number of neighbors: $k \in \{1, 3, 5, 7\}$
J48	Pruning: { True, False }, Pruning confidence: { 0.1, 0.25, 0.5 }
JRip	Pruning: { True, False }
Naive Bayes	Use kernel estimation: { True, False }
Logistic	Log-likelihood: { $1e^{-7}$, $1e^{-8}$, $1e^{-9}$ }

NNISDA, and MNSVM's percentage of support vectors and the number of samples for all datasets. It highlights OLLAWV's model sparseness in comparison to the competing methods, while mirroring the runtime results.

4.2. Non-SVM experimental study

The supplemental experimental setup was designed to compare the performance of the proposed OLLAWV against the following 5 popular and widely-used non-SVM algorithms: *k*-Nearest Neighbors (KNN), J48, JRip, Naive Bayes, and Logistic. These methods have been implemented within the Weka framework [41]. The experiments were performed under the same nested 5-fold cross-validation framework as the SVM algorithm experimental study which was described in Section 4.1. The following hyper-parameters shown in Table 4 were used for the non-SVM algorithms.

4.2.1. Results and statistical analysis

Table 5 displays the accuracy results for OLLAWV and five state-of-the-art methods. The table also shows the standard deviation for accuracy per outer fold, the average values across all datasets, and the algorithm ranks. As the results indicate, OLLAWV outperforms all other state-of-the-art methods. Fig. 9 displays the average accuracy results for OLLAWV and the non-SVM methods across all datasets and highlights OLLAWV's better performance. The Friedman test indicates that OLLAWV performs significantly better than the competing methods for $\alpha = 0.05$ and is ranked first. Fig. 8 shows the critical distance bar (which is 1.391), and indicates that all other

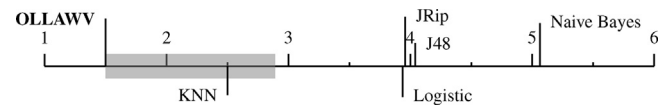


Fig. 8. Bonferroni–Dunn test for accuracy.

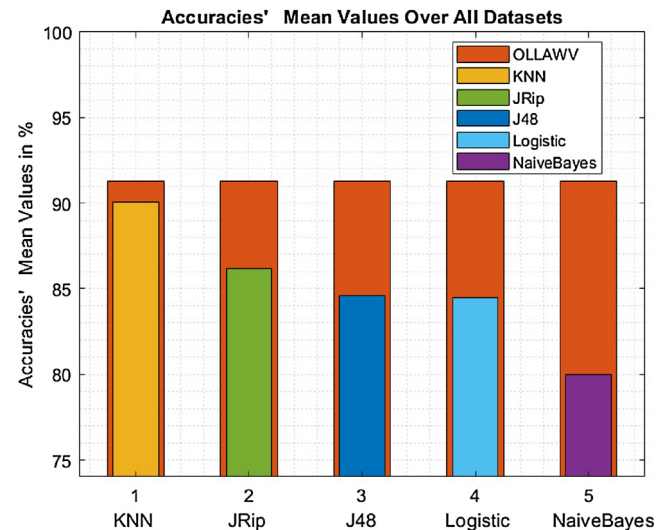


Fig. 9. Mean accuracy over all datasets for OLLAWV and the 5 non-SVM state-of-the-art methods.

algorithms perform statistically significantly worse than OLLAWV, except for KNN.

5. Conclusion

This paper proposed a novel online learning procedure and algorithm for solving the L1-SVM problem, which is a unique method in terms of both iterating over samples and updating the model. A new stopping criterion for the stochastic gradient procedure is also proposed. The model is updated by changing the weight α_i of

Table 5
Accuracy comparison for non-SVM methods versus OLLAWV.

Dataset	OLLAWV	KNN	J48	JRip	Naïve Bayes	Logistic
iris	97.33 ± 1.49	96.00 ± 3.65	94.00 ± 2.79	90.67 ± 4.35	96.00 ± 2.79	97.33 ± 2.79
teach	52.32 ± 3.46	59.64 ± 2.89	49.72 ± 7.58	56.75 ± 9.60	53.75 ± 6.46	51.77 ± 6.68
wine	98.87 ± 1.54	97.73 ± 3.72	90.43 ± 5.83	93.24 ± 3.27	96.60 ± 3.14	96.05 ± 2.58
cancer	80.36 ± 5.80	77.32 ± 6.93	73.81 ± 8.57	73.78 ± 5.81	67.73 ± 5.07	77.32 ± 7.78
sonar	92.32 ± 3.11	88.99 ± 4.59	76.16 ± 10.6	75.18 ± 6.77	73.69 ± 7.65	75.18 ± 7.31
glass	72.41 ± 2.28	67.73 ± 5.91	65.06 ± 5.51	65.59 ± 9.66	49.46 ± 5.19	62.04 ± 5.75
vote	96.54 ± 1.87	92.26 ± 3.19	95.70 ± 2.12	96.54 ± 2.45	92.24 ± 3.24	93.54 ± 2.59
heart2	82.22 ± 2.93	79.63 ± 5.71	78.52 ± 2.81	80.74 ± 4.06	84.44 ± 4.46	83.33 ± 3.93
dermatology	97.82 ± 0.05	96.18 ± 1.78	94.52 ± 2.21	91.27 ± 5.08	97.28 ± 1.64	96.98 ± 2.28
pro	88.96 ± 2.14	87.96 ± 3.01	78.54 ± 1.62	79.13 ± 2.78	62.38 ± 3.54	87.57 ± 2.56
euk	77.38 ± 1.96	81.42 ± 2.06	65.27 ± 2.92	66.42 ± 3.47	39.27 ± 3.43	69.55 ± 1.34
optdigits	99.11 ± 0.38	98.74 ± 0.39	90.87 ± 1.09	91.28 ± 0.40	92.42 ± 0.75	95.05 ± 0.91
satimage	91.66 ± 0.80	90.38 ± 0.72	85.64 ± 1.21	85.33 ± 0.77	85.41 ± 0.92	88.14 ± 1.11
usps	97.49 ± 0.22	97.04 ± 0.47	88.73 ± 0.46	89.20 ± 1.00	79.45 ± 0.59	91.88 ± 0.65
pendigits	99.56 ± 0.12	99.33 ± 0.17	96.24 ± 0.31	96.34 ± 0.41	88.34 ± 0.65	95.59 ± 0.18
reuters	98.03 ± 0.22	97.15 ± 0.43	96.90 ± 0.32	97.18 ± 0.44	93.52 ± 0.02	69.54 ± 0.28
letter	96.99 ± 0.21	95.71 ± 0.19	87.34 ± 0.68	87.02 ± 0.66	74.12 ± 0.97	77.45 ± 0.16
adult	84.75 ± 0.26	83.85 ± 0.28	84.38 ± 0.28	83.73 ± 0.17	80.57 ± 0.09	82.46 ± 0.14
w3a	98.86 ± 0.04	98.60 ± 0.06	98.71 ± 0.05	98.41 ± 0.10	96.71 ± 0.20	98.61 ± 0.12
shuttle	99.77 ± 0.03	99.93 ± 0.03	99.97 ± 0.02	99.96 ± 0.02	98.57 ± 0.24	96.83 ± 0.12
web	98.94 ± 0.05	98.89 ± 0.06	98.79 ± 0.09	98.50 ± 0.13	96.71 ± 0.21	98.70 ± 0.08
ijcnn1	98.31 ± 0.07	98.48 ± 0.04	98.40 ± 0.09	98.11 ± 0.10	90.69 ± 0.26	92.29 ± 0.16
intrusion	99.77 ± 0.02	88.20 ± 1.06	58.01 ± 26.6	87.66 ± 3.79	49.75 ± 30.7	65.15 ± 15.7
Average	91.29 ± 1.26	90.05 ± 2.06	84.60 ± 3.64	86.18 ± 2.84	79.96 ± 3.58	84.45 ± 2.83
Ranks	1.500	2.500	4.041	3.9583	5.0625	3.9375

The best results are bolded.

a single worst-violator per iteration and stops when all violating samples i.e., support vectors, are found. Finding the *worst-violators* is done without replacement. Such an approach results in a significant shortening of training time, as well as in a huge decrease in the resulting model size. The key features of the proposed algorithm, OLLAWV, stem from its implicit ability of finding support vectors and its self-stopping condition. This design was devised to address the limitations presented by current SVM solvers.

The first experimental study demonstrates the better performance of OLLAWV compared with state-of-the-art SVM solvers (MNSVM and NNISDA) which have been shown to outperform the popular SMO implementation in the LIBSVM package. The results for accuracy, runtime, and percentage of support-vectors, obtained by the strict nested cross-validation procedure, were compared and further validated using statistical analysis with non-parametric tests. They highlighted the advantages and major speedup achieved by OLLAWV against the competing MNSVM and NNISDA. The second, supplemental experimental study evaluated the performance of OLLAWV against 5 popular non-SVM methods, showing the better performance of OLLAWV against all five non-SVM algorithms (k -Nearest Neighbors (KNN), J48, JRip, Naïve Bayes, and Logistic). The proposal, OLLAWV, performs statistically better in terms of runtime and model size across all 23 evaluated benchmark datasets, without compromising accuracy.

Acknowledgements

This research was supported by the VCU Research Support Funds, the Spanish Ministry of Economy and Competitiveness and the European Regional Development Fund, project TIN2017-83445-P.

References

- [1] Y. Lu, K. Boukharouba, J. Boonært, A. Fleury, S. Lecoeuche, Application of an incremental SVM algorithm for on-line human recognition from video surveillance using texture and color features, *Neurocomputing* 126 (2014) 132–140.

- [2] P. Ksieniewicz, B. Krawczyk, M. Woźniak, Ensemble of Extreme Learning Machines with trained classifier combination and statistical features for hyperspectral data, *Neurocomputing* 271 (2018) 28–37.
- [3] C. Cortes, V. Vapnik, Support-vector networks, *Mach. Learn.* 20 (3) (1995) 273–297.
- [4] G. Melki, A. Cano, V. Kecman, S. Ventura, Multi-target support vector regression via correlation regressor chains, *Inf. Sci.* 415 (2017) 53–69.
- [5] L. Gonzalez-Abriel, C. Angulo, H. Nuñez, Y. Leal, Handling binary classification problems with a priority class by using Support Vector Machines, *Appl. Soft Comput.* 61 (2017) 661–669.
- [6] J. Platt, Sequential Minimal Optimization: A Fast Algorithm for Training Support Vector Machines, Technical Report: MSR-TR-98-14, 1998.
- [7] R. Strack, Geometric Approach to Support Vector Machines Learning for Large Datasets, Virginia Commonwealth University, 2013.
- [8] V. Kecman, L. Zigic, Algorithms for direct L2 support vector machines, *Proceedings of the IEEE International Symposium on Innovations in Intelligent Systems and Applications (2014)* 419–424.
- [9] L. Zigic, Direct L2 Support Vector Machine, Ph.D. Dissertation, Virginia Commonwealth University, 2016 <http://scholarscompass.vcu.edu/etd/4274/>.
- [10] L. Bottou, Large-scale machine learning with stochastic gradient descent, *Proceedings of 19th International Conference on Computational Statistics (2010)* 177–186.
- [11] J. Kivinen, A.J. Smola, R.C. Williamson, Online learning with kernels, *IEEE Trans. Signal Process.* 52 (8) (2004) 2165–2176.
- [12] S. Shalev-Shwartz, Y. Singer, N. Srebro, A. Cotter, Pegasos: primal estimated sub-gradient solver for SVM, *Math. Program.* 127 (1) (2011) 3–30.
- [13] T. Zhang, Solving large scale linear prediction problems using stochastic gradient descent algorithms, *Proceedings of the 21st International Conference on Machine Learning, ACM (2004)* 116.
- [14] C. Panagiotakopoulos, P. Tsampouka, The stochastic gradient descent for the primal L1-SVM optimization revisited, *Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases (2013)* 65–80.
- [15] G. Cauwenberghs, T. Poggio, Incremental and decremental support vector machine learning, *Advances in Neural Information Processing Systems (2001)* 409–415.
- [16] X. Song, L. Jian, Y. Song, A chunk updating LS-SVMs based on block Gaussian elimination method, *Appl. Soft Comput.* 51 (2017) 96–104.
- [17] R. Collobert, S. Bengio, Links between perceptrons, MLPs and SVMs, *Proceedings of the 21st International Conference on Machine Learning (2004)* 23.
- [18] V. Kecman, *Learning and Soft Computing: Support Vector Machines, Neural Networks, and Fuzzy Logic Models*, MIT Press, 2001.
- [19] S. Shalev-Shwartz, S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*, Cambridge University Press, 2014.
- [20] B.E. Boser, I.M. Guyon, V.N. Vapnik, A training algorithm for optimal margin classifiers, *Proceedings of the 5th Annual Workshop on Computational Learning Theory (1992)* 144–152.
- [21] B. Schölkopf, A. Smola, *Learning with Kernels; Support Vector Machines, Regularization, Optimization, and Beyond*, MIT Press, 2002.

- [22] M.A. Aizerman, E.A. Braverman, L. Rozonoer, Theoretical foundations of the potential function method in pattern recognition learning, *Automation and Remote Control*, vol. 25 (1964) 821–837.
- [23] B. Cyganek, B. Krawczyk, M. Woźniak, Multidimensional data classification with chordal distance based kernel and support vector machines, *Eng. Appl. Artif. Intell.* 46 (2015) 10–22.
- [24] B. Krawczyk, J.A. Sáez, M. Woźniak, Tackling label noise with multi-class decomposition using fuzzy one-class support vector machines, *IEEE International Conference on Fuzzy Systems (2016)* 915–922.
- [25] S.S. Keerthi, S.K. Shevade, C. Bhattacharyya, K.R.K. Murthy, Improvements to Platt's SMO algorithm for SVM classifier design, *Neural Comput.* 13 (3) (2001) 637–649.
- [26] C.-C. Chang, C.-J. Lin, LIBSVM: a library for support vector machines, *ACM Trans. Intell. Syst. Technol.* 2 (2011), 27:1–27:27, Software available at: <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [27] S. Boyd, L. Vanderberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [28] K.P. Bennett, E.J. Bredensteiner, Duality and geometry in SVM classifiers, *International Conference on Machine Learning (2000)* 57–64.
- [29] I.W. Tsang, J.T. Kwok, P.-M. Cheung, Core vector machines: fast SVM training on very large data sets, *J. Mach. Learn. Res.* 6 (2005) 363–392.
- [30] I.W. Tsang, A. Kocsor, J.T. Kwok, Simpler core vector machines with enclosing balls, *Proceedings of the 24th International Conference on Machine Learning (2007)* 911–918.
- [31] T. Huang, V. Kecman, I. Kopriva, *Kernel Based Algorithms for Mining Huge Data Sets, Supervised, Semi-supervised, and Unsupervised Learning*, Springer-Verlag, 2006.
- [32] V. Kecman, T. Huang, M. Vogt, Iterative single data algorithm for training kernel machines from huge data sets: theory and performance, *Stud. Comput. Intell.* 177 (2005) 255–274.
- [33] J. Kivinen, A.J. Smola, R.C. Williamson, Large margin classification for moving targets, *International Conference on Algorithmic Learning Theory*, vol. 2 (2002) 113–127.
- [34] R. Herbrich, *Learning Kernel Classifiers*, MIT Press, 2016.
- [35] O. Bousquet, L. Bottou, The tradeoffs of large scale learning, *Advances in Neural Information Processing Systems (2008)* 161–168.
- [36] V. Kecman, Fast online algorithm for nonlinear support vector machines and other alike models, *Opt. Mem. Neural Netw.* 25 (4) (2016) 203–218.
- [37] V. Kecman, G. Melki, Fast online algorithms for support vector machines, *Proceedings of the IEEE Southeast Conference (2016)* 1–6.
- [38] G. Melki, V. Kecman, Speeding up online training of L1 support vector machines, *Proceedings of the IEEE Southeast Conference (2016)* 1–6.
- [39] J. Derrac, S. García, D. Molina, F. Herrera, A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms, *Swarm Evol. Comput.* 1 (1) (2011) 3–18.
- [40] M. Lichman, *UCI Machine Learning Repository*, 2013 <http://archive.ics.uci.edu/ml>.
- [41] F. Eibe, M. Hall, I. Witten, J. Pal, *The WEKA Workbench, Online Appendix for Data Mining: Practical Machine Learning Tools and Techniques*, 4, 2016.
- [42] S. Varma, R. Simon, Bias in error estimation when using cross-validation for model selection, *BMC Bioinformatics* 7 (1) (2006) 91.
- [43] S.-J. Wu, V.-H. Pham, T.-N. Nguyen, Two-phase optimization for support vectors and parameter selection of support vector machines: two-class classification, *Appl. Soft Comput.* 59 (2017) 129–142.