

©Gabriella Angela Melki, September 2018

All Rights Reserved.

# NOVEL SUPPORT VECTOR MACHINES FOR DIVERSE LEARNING PARADIGMS

A Dissertation submitted in partial fulfillment of the requirements for the degree of Doctor  
of Philosophy at Virginia Commonwealth University.

by

GABRIELLA ANGELA MELKI

Ph.D. Candidate

Director: Alberto Cano,

Assistant Professor, Department of Computer Science,

Virginia Commonwealth University

Director: Sebastian Ventura,

Professor, Department of Computer Science & Numerical Analysis,

University of Córdoba

Virginia Commonwealth University

Richmond, Virginia

September 2018

## CHAPTER 1

### NOVEL ONLINE SVM USING WORST-VIOLATORS

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 chapter 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.

#### 1.1 Online Learning Background

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

Table 1.1.: 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	$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_i, y_i), \dots, (\mathbf{x}_n, y_n)\}$

### 1.1.1 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) 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 ?? provides a summary of the notation used in this paper.

### 1.1.2 Online Learning Methods

**THIS MAYBE SHOULD GO IN THE SVM BACKGROUND taken from Pegasos paper.** Before indulging into the detailed description and analysis of Pegasos, we would like to draw connections to and put our work in context of some of the more recent work on SVM. For a more comprehensive and up-to-date overview of relevant work see the references in the papers cited below as well as the web site dedicated to kernel methods at <http://www.kernel-machines.org>. Due to the centrality of the SVM optimization problem, quite a few methods were devised and analyzed. The different approaches can be roughly divided into the following categories. **taken from Pegasos paper.**

- **Interior Point (IP) methods:** IP methods (see for instance [7] and the references therein) cast the SVM learning task as a quadratic optimization problem subject to linear constraints. The constraints are replaced with a barrier function. The result is

a sequence of unconstrained problems which can be optimized very efficiently using Newton or Quasi-Newton methods. The advantage of IP methods is that the dependence on the accuracy is double logarithmic, namely,  $\log(\log(1/\epsilon))$ . Alas, IP methods typically require run time which is cubic in the number of examples  $m$ . Moreover, the memory requirements of IP methods are  $O(m^2)$  which renders a direct use of IP methods very difficult when the training set consists of many examples. It should be noted that there have been several attempts to reduce the complexity based on additional assumptions (see e.g. [15]). However, the dependence on  $m$  remains super linear. In addition, while the focus of the paper is the optimization problem cast by SVM, one needs to bear in mind that the optimization problem is a proxy method for obtaining good classification error on unseen examples. Achieving a very high accuracy in the optimization process is usually unnecessary and does not translate to a significant increase in the generalization accuracy. The time spent by IP methods for finding a single accurate solution may, for instance, be better utilized for trying different regularization values.

- **Decomposition methods:** To overcome the quadratic memory requirement of IP methods, decomposition methods such as SMO [29] and SVM-Light [20] tackle the dual representation of the SVM optimization problem, and employ an active set of constraints thus working on a subset of dual variables. In the extreme case, called row-action methods [8], the active set consists of a single constraint. While algorithms in this family are fairly simple to implement and entertain general asymptotic convergence properties [8], the time complexity of most of the algorithms in this family is typically super linear in the training set size  $m$ . Moreover, since decomposition methods find a feasible dual solution and their goal is to maximize the dual objective function, they often result in a rather slow convergence rate to the optimum of the primal objective function. (See also the discussion in [19].)

- **Primal optimization:** Most existing approaches, including the methods discussed above, focus on the dual of Eq. (1), especially when used in conjunction with non-linear kernels. However, even when non-linear kernels are used, the Representer theorem [23] allows us to re-parametrize  $w$  as  $w = \sum_i y_i x_i$  and cast the primal objective Eq. (1) as an unconstrained optimization problem with the variables  $1, \dots, m$  (see Sec. 4). Tackling the primal objective directly was studied, for example, by Chapelle [10], who considered using smooth loss functions instead of the hinge loss, in which case the optimization problem becomes a smooth unconstrained optimization problem. Chapelle then suggested using various optimization approaches such as conjugate gradient descent and Newton's method. We take a similar approach here, however we cope with the non-differentiability of the hinge loss directly by using sub-gradients instead of gradients. Another important distinction is that Chapelle views the optimization problem as a function of the variables  $i$ . In contrast, though Pegasos maintains the same set of variables, the optimization process is performed with respect to  $w$ , see Sec. 4 for details.
- **Stochastic gradient descent:** The Pegasos algorithm is an application of a stochastic sub-gradient method (see for example [25,34]). In the context of machine learning problems, the efficiency of the stochastic gradient approach has been studied in [26,1,3,27,6,5]. In particular, it has been claimed and experimentally observed that, Stochastic algorithms yield the best generalization performance despite being the worst optimization algorithms. This claim has recently received formal treatment in [4,32]. Two concrete algorithms that are closely related to the Pegasos algorithm and are also variants of stochastic sub-gradient methods are the NORMA algorithm [24] and a stochastic gradient algorithm due to Zhang [37]. The main difference between Pegasos and these variants is in the procedure for setting the step size. We elaborate on this issue in Sec. 7. The convergence rate given in [24] implies that the number of iterations required to achieve  $\epsilon$ -accurate solution is  $O(1/\epsilon^2)$ . This bound is inferior to

the corresponding bound of Pegasos. The analysis in [37] for the case of regularized loss shows that the squared Euclidean distance to the optimal solution converges to zero but the rate of convergence depends on the step size parameter. As we show in Sec. 7, tuning this parameter is crucial to the success of the method. In contrast, Pegasos is virtually parameter free. Another related recent work is Nesterovs general primal-dual subgradient method for the minimization of non-smooth functions [28]. Intuitively, the ideas presented in [28] can be combined with the stochastic regime of Pegasos. We leave this direction and other potential extensions of Pegasos for future research.

- **Cutting Planes Approach:** Recently, Joachims [21] proposed SVM-Perf, which uses a cutting planes method to find a solution with accuracy in time  $O(md/(2))$ . This bound was later improved by Smola et al [33] to  $O(md/())$ . The complexity guarantee for Pegasos avoids the dependence on the data set size  $m$ . In addition, while SVM-Perf yields very significant improvements over decomposition methods for large data sets, our experiments (see Sec. 7) indicate that Pegasos is substantially faster than SVM-Perf.

## 1.2 OnLine Learning Algorithm using Worst-Violator

OLLAWV 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 OLLAWV came from [1] 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 OLLAWV resides in its iterative method, where the weight  $\alpha_i$  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, OLLAWV

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. Equation 1.1 shows the method for selecting the worst-violator, where  $y_o \in \mathbb{R}$  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 OLLAWV is equal to the number of support vectors that the resulting model has. This is an interesting property of OLLAWV; if the number of iterations is set beforehand, one is implicitly setting a bound on the number of support vectors.

$$[y_o, wv] = \min \{y_{wv} \cdot o_{wv}\}, \forall wv \in \{\neg \mathbf{S}\} \quad (1.1)$$

Algorithm 1.1 lists OLLAWV's pseudocode and Figure 1.2 illustrates the steps taken by OLLAWV. First, the model parameters ( $\alpha$ ,  $b$ ,  $\mathbf{S}$ ) and the algorithm variables ( $\mathbf{o}$ , iteration counter ( $t$ ), initial worst-violator index  $wv$  and it's error  $y_o$ ) 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 func-



tion given below:

$$\min_{\mathbf{w} \in \mathcal{H}_o \times \mathbb{R}} R = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max \{0, 1 - y_i o_{(w)}(\mathbf{x}_i)\}, \quad (1.2)$$

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 Equation 1.2 is non-differentiable, but it has a subgradient due to a knick-point at  $yo = 1$ . The loss function's gradient after the knick-point equals zero, which leads to a sparse model. Hence, when the value of  $yo \geq 1$  the loss is zero, and for  $yo < 1$  the loss increases linearly. The subgradient of the above cost function is given by:

$$\frac{\partial R}{\partial \mathbf{w}} = \begin{cases} \mathbf{w} - C \sum_{i=1}^n y_i \mathbf{x}_i & y_i o_i < 1 \\ \mathbf{w} & \text{otherwise.} \end{cases} \quad (1.3)$$

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 [2] 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 Equation 1.3 becomes:

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

where  $\eta > 0 \in \mathbb{R}$  is the learning rate. According to the Representer Theorem, a vector  $\boldsymbol{\alpha} \in \mathbb{R}^n$  exists such that  $\mathbf{w} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i)$  is an optimal solution to Equation 1.2, where  $\phi(\cdot)$  is a mapping from feature space to Hilbert space [3]. On the basis of the Representer Theorem, Equation 1.2 can be optimized with respect to  $\boldsymbol{\alpha}$  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:

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

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

$$\forall i : \alpha_i \leftarrow \alpha_i + \eta(Cy_i - \alpha_i) \quad (1.4)$$

If the bias term  $b$  is included in Equation 1.2, it's stochastic update is as follows:

$$\forall i : b \leftarrow b + \eta \frac{Cy_i}{n} \quad (1.5)$$

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  $\alpha$  update can be expressed as:  $\alpha_i \leftarrow \alpha_i + (\Lambda - P)$ . Note that  $\Lambda$  is the update resulting from the loss function and  $P$  is derived from the regularizer term in Equation 1.2. In the case of OLLAWV,  $P = 0$  because samples are never updated more than once and their initial  $\alpha$  value is always 0. It is important to note that in OLLAWV's case,  $\Lambda$  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

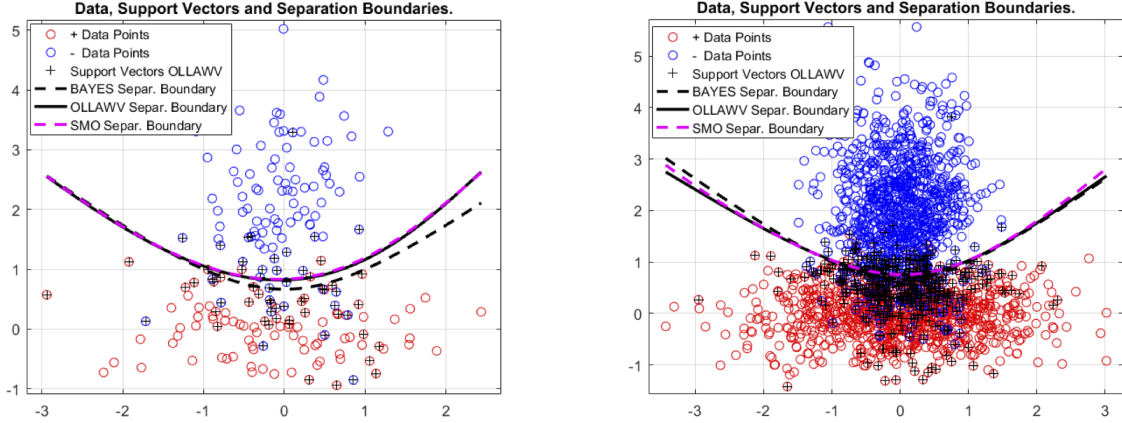


Fig. 1.1.: A case of classifying 2-dimensional normally distributed data with different covariance matrices, (left) for 200 and (right) 2000 data points. The theoretical separation boundary (denoted as the Bayes Separation Boundary) is quadratic and is shown as the dashed black curve. The other two separation 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.

vector  $\mathbf{o} \in \mathbb{R}^n$  in Algorithm 1.1), from which a worst-violator is found, changes per iteration based on the influence of the support-vectors that have been previously updated. From Equation ??, the output vector update becomes the following:

$$\mathbf{o} \leftarrow \mathbf{o} + \Lambda * \mathcal{K}(\mathbf{x}_{-\mathcal{S}}, \mathbf{x}_{wv}, \gamma) + B \quad (1.6)$$

where  $\mathcal{K}(\cdot)$  is the Gaussian radial basis function (RBF) kernel,  $\gamma \in \mathbb{R}$  is its parameter, and  $B = (\Lambda * \beta) / n$  denotes the bias update. Only non-support vector output values are calculated per iteration, as denoted by  $\mathbf{x}_{-\mathcal{S}}$  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 Equation 1.6,  $\beta \in \{0, 1\}$  indicates whether the bias term is to be used. If  $b$  is not a part of the model, it should be omitted from Equations ?? and 1.6 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. Figure 1.1 shows the decision boundaries achieved by OLLAWV versus SMO (implemented within LIBSVM) and Bayes for toy datasets.

### 1.3 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 1.3.1, compares OLLAWV to two other SVM kernel methods, and the second compares OLLAWV to 5 non-SVM methods, shown in Section 1.3.2. In each section, the experimental setups are first described and the state-of-the-art methods are listed. The results and statistical analysis [Derrac2011] 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 1.3.2 was conducted to emphasize the better performance of OLLAWV against non-SVM algorithms.

Table 1.2 presents a summary of the 23 datasets used throughout the experiments, where

---

**Algorithm 1.1** OnLine Learning Algorithm using Worst-Violators (OLLAWV)
 

---

**Input:**  $\mathcal{D}, C, \gamma, \beta, M$ 
**Output:**  $\alpha, b, \mathcal{S}$ 

```

1:  $\alpha \leftarrow \mathbf{0}, b \leftarrow 0, \mathcal{S} \leftarrow \mathbf{0}$  ▷ Initialize OLLAWV model parameters
2:  $\mathbf{o} \leftarrow \mathbf{0}, t \leftarrow 0$  ▷ Initialize the output vector and iteration counter
3:  $wv \leftarrow 0, yo \leftarrow y_{wv} * \mathbf{o}_{wv}$  ▷ Initialize hinge loss error and worst-violator index
4: while  $yo < M$  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:   $\mathbf{o} \leftarrow \mathbf{o} + \Lambda * \mathcal{K}(\mathbf{x}_{-\mathcal{S}}, \mathbf{x}_{wv}, \gamma) + B$  ▷ Update output vector
11:   $\alpha_{wv} \leftarrow \alpha_{wv} + \Lambda$  ▷ Update worst-violator's alpha value
12:   $b \leftarrow b + B$  ▷ Update bias term
13:
14:   $\mathcal{S}_t \leftarrow wv$  ▷ Save index of worst-violator
15:   $[yo, wv] \leftarrow \min_{wv \in \{-\mathcal{S}\}} \{y_{wv} \cdot \mathbf{o}_{wv}\}$  ▷ Find the worst-violator
16: end while
  
```

---

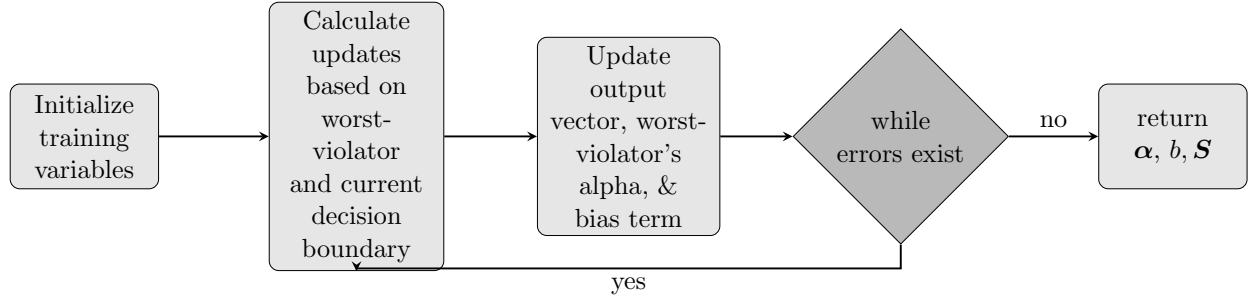


Fig. 1.2.: A summary of the steps performed by OLLAWV. The model parameters  $(\alpha, b, \mathcal{S})$  and the algorithm variables  $(\mathbf{o}, t, wv, \text{ and } yo)$  are first initialized. The worst-violator with respect to the current hyperplane is then found and the model parameters are then updated. Once no more violating samples are found, the model is returned.

Table 1.2.: Datasets

Dataset	# Samples	# Attributes	# Classes
<b><i>small datasets</i></b>			
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
<b><i>medium datasets</i></b>			
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
<b><i>large datasets</i></b>			
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

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 repository<sup>1</sup> [Lichman:2013], and the LibCVM<sup>2</sup> [tsang2005core] and LIBSVM<sup>3</sup> [CC01a] sites.

### 1.3.1 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*

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

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

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

(MNSVM) [**strack2013geometric**] and *Non-Negative Iterative Single Data Algorithm* (NNISDA) [**zigic2016**].

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 [**strack2013geometric**], it was shown that MNSVM outperforms both the L1 and L2 implementations of LIBSVM, and BVM embedded in LibCVM. NNISDA was then compared to MNSVM in [**zigic2016**], 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 Training<sup>4</sup>”. 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 [**varma2006bias**, **wu2017two**]. 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

---

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

pattern search, was performed. The best hyper-parameters were chosen from the following  $8 \times 8$  possible combinations, shown in Equations (1.7a) and (1.7b), and were also used for the competing SVM methods.

$$C \in \{4^n\}, \quad n = \{-2, \dots, 5\} \quad (1.7a)$$

$$\gamma \in \{4^n\}, \quad n = \{-5, \dots, 2\} \quad (1.7b)$$

The  $\gamma$  parameter refers to that of the Gaussian RBF kernel, given by:

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}. \quad (1.8)$$

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.

## 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 1.3 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 [Derrac2011]. 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 1.3, 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



Table 1.3.: 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	<b>97.33</b>	94.00	96.67	<b>0.05</b>	0.27	3.57	<b>13.50</b>	40.20	29.80
teach	52.32	52.31	<b>52.95</b>	<b>0.12</b>	0.44	8.85	<b>69.19</b>	99.80	87.40
wine	<b>98.87</b>	96.60	96.60	<b>0.28</b>	0.43	4.84	<b>15.02</b>	44.40	48.60
cancer	80.36	<b>81.86</b>	81.38	<b>0.49</b>	0.85	4.46	<b>42.79</b>	83.80	89.60
sonar	<b>92.32</b>	89.48	87.57	<b>0.59</b>	0.98	3.03	<b>31.26</b>	73.00	66.00
glass	<b>72.41</b>	67.81	69.30	<b>0.46</b>	1.01	11.94	<b>62.84</b>	90.80	87.60
vote	<b>96.54</b>	96.11	93.99	<b>0.26</b>	0.46	1.49	<b>13.36</b>	33.20	34.00
heart	82.22	<b>83.33</b>	<b>83.33</b>	<b>0.50</b>	0.91	6.45	<b>37.69</b>	73.00	82.00
dermatology	97.82	<b>98.36</b>	<b>98.36</b>	<b>1.62</b>	2.47	11.68	<b>36.94</b>	59.00	59.80
prokaryotic	88.96	88.86	<b>88.97</b>	<b>6.09</b>	10.64	50.86	<b>29.01</b>	51.20	49.00
eukaryotic	77.38	79.56	<b>81.21</b>	61.95	<b>49.16</b>	342.76	<b>54.11</b>	76.40	72.60
<i>medium datasets</i>									
optdigits	99.11	99.29	<b>99.31</b>	<b>411</b>	528	787	<b>28.64</b>	31.60	30.60
satimage	91.66	<b>92.39</b>	92.35	1,334	<b>687</b>	1,094	<b>20.72</b>	45.00	44.80
usps	97.49	98.05	<b>98.24</b>	10,214	<b>5,245</b>	7,777	<b>11.22</b>	29.40	28.00
pendigits	99.56	<b>99.62</b>	99.61	<b>723</b>	909	1,500	<b>10.27</b>	17.60	16.60
reuters	98.03	<b>98.08</b>	97.99	<b>954</b>	1,368	1,657	<b>8.770</b>	18.20	18.60
letter	96.99	99.11	<b>99.13</b>	<b>5,259</b>	12,009	26,551	<b>43.56</b>	57.60	56.60
<i>large datasets</i>									
adult	84.75	85.07	<b>85.13</b>	<b>21,025</b>	72,552	123,067	<b>34.66</b>	56.00	56.60
w3a	<b>98.86</b>	98.82	98.82	<b>6,532</b>	15,951	24,562	<b>3.270</b>	14.60	12.40
shuttle	99.77	99.83	<b>99.87</b>	<b>2,833</b>	7,420	45,062	<b>2.010</b>	6.00	16.40
web	98.94	<b>99.00</b>	99.00	<b>12,067</b>	30,583	38,040	<b>4.320</b>	13.20	10.80
ijcnn1	98.31	99.34	<b>99.41</b>	<b>162,587</b>	296,917	370,144	16.36	11.00	<b>7.600</b>
intrusion	<b>99.77</b>	99.67	99.66	<b>2,402,804</b>	4,646,810	3,772,113	<b>0.780</b>	2.000	1.700
Average	<b>91.29</b>	91.15	91.25	<b>114,209</b>	221,350	191,861	<b>25.66</b>	44.65	43.79
Ranks	<b>1.739</b>	2.022	2.239	<b>1.217</b>	1.913	2.869	<b>1.087</b>	2.609	2.304

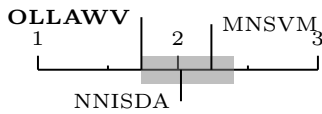


Fig. 1.3.: Bonferroni-Dunn test for Accuracy

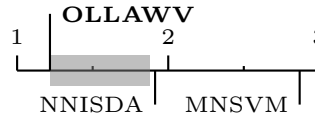


Fig. 1.4.: Bonferroni-Dunn test for Runtime

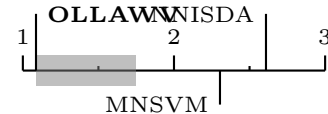


Fig. 1.5.: Bonferroni-Dunn test for % Support Vectors

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 1.3, Figures 1.3, 1.4, and 1.5 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

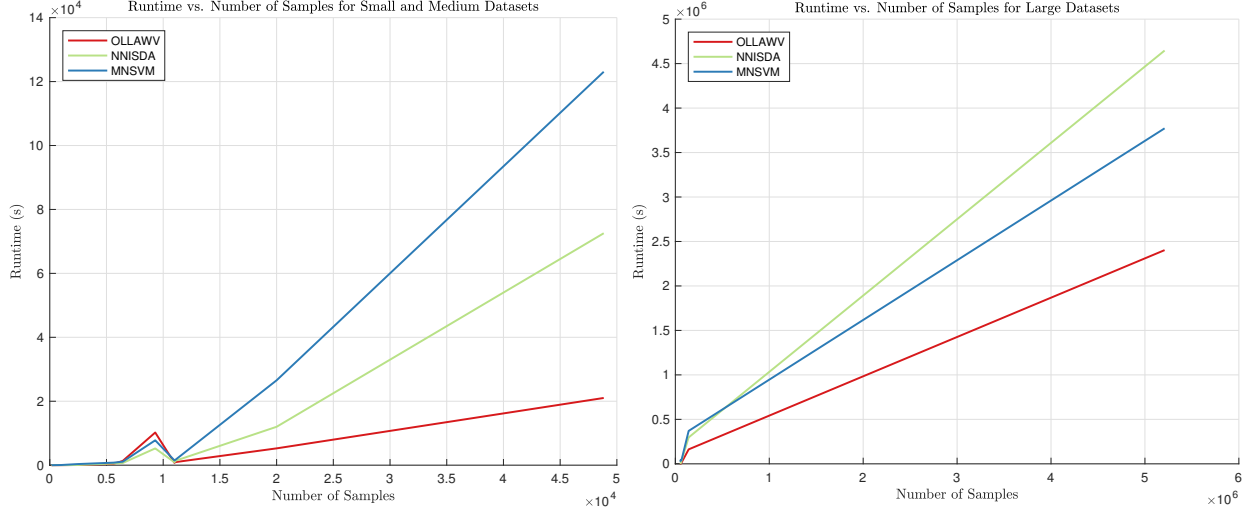


Fig. 1.6.: Runtime in seconds versus the number of samples, divided into two groups: small & medium (left) versus large (right). Note OLLAWV’s gradual increase in runtime as the number of samples increases compared to NNISDA and MNSVM’s steeper change. In almost all cases, OLLAWV displays superior runtime 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.

worse than the control algorithm, OLLAWV.

The results in Table 1.3 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 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.

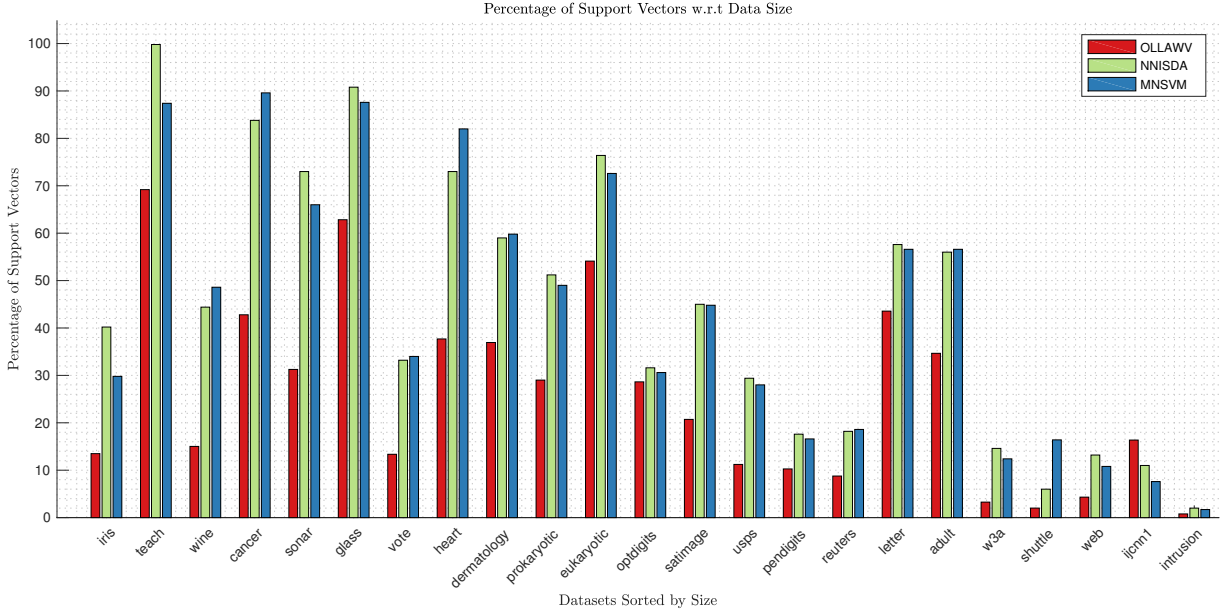


Fig. 1.7.: Size of the model given as percentage of support vectors with respect to the number of samples versus the number of samples. Note that OLLAWV's percentage of support vectors is always smaller (except in one case) than NNISDA's and MNSVM's ones.

Figures 1.3, 1.4, and 1.5 show the results of the statistical analysis for accuracy, runtime, and percentage of support vectors. Figures 1.4 and 1.5 show that OLLAWV is statistically significantly better than MNSVM and NNISDA for runtime and percentage of support vectors (model size). At the same time, Figure 1.3 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.

Figure 1.6 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. Figure 1.7 shows the correlation between OLLAWV, 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.

### 1.3.2 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*, *Naïve Bayes*, and *Logistic*. These methods have been implemented within the Weka framework [eibe2016weka]. The experiments were performed under the same nested 5-fold cross-validation framework as the SVM algorithm experimental study which was described in Section 1.3.1. The following hyper-parameters shown in Table 1.4 were used for the non-SVM algorithms.

Table 1.4.: 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: $\{1e^{-7}, 1e^{-8}, 1e^{-9}\}$

## Results and Statistical Analysis

Table 1.5 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 accross all datasets, and the algorithm ranks. As the results indicate, OLLAWV outperforms all other state-of-the-art methods. Figure 1.9 displays the average accuracy results for OLLAWV and the non-SVM methods accross 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. Figure 1.8 shows the critical distance bar (which is 1.391), and indicates that all other algorithms perform statistically

Table 1.5.: Accuracy Comparison for Non-SVM Methods versus OLLAWV

Dataset	OLLAWV	KNN	J48	JRip	Naïve Bayes	Logistic
iris	<b>97.33</b> $\pm$ 1.49	96.00 $\pm$ 3.65	94.00 $\pm$ 2.79	90.67 $\pm$ 4.35	96.00 $\pm$ 2.79	97.33 $\pm$ 2.79
teach	52.32 $\pm$ 3.46	<b>59.64</b> $\pm$ <b>2.89</b>	49.72 $\pm$ 7.58	56.75 $\pm$ 9.60	53.75 $\pm$ 6.46	51.77 $\pm$ 6.00
wine	<b>98.87</b> $\pm$ <b>1.54</b>	97.73 $\pm$ 3.72	90.43 $\pm$ 5.83	93.24 $\pm$ 3.27	96.60 $\pm$ 3.14	96.05 $\pm$ 2.50
cancer	<b>80.36</b> $\pm$ <b>5.80</b>	77.32 $\pm$ 6.93	73.81 $\pm$ 8.57	73.78 $\pm$ 5.81	67.73 $\pm$ 5.07	77.32 $\pm$ 7.70
sonar	<b>92.32</b> $\pm$ <b>3.11</b>	88.99 $\pm$ 4.59	76.16 $\pm$ 10.6	75.18 $\pm$ 6.77	73.69 $\pm$ 7.65	75.18 $\pm$ 7.30
glass	<b>72.41</b> $\pm$ <b>2.28</b>	67.73 $\pm$ 5.91	65.06 $\pm$ 5.51	65.59 $\pm$ 9.66	49.46 $\pm$ 5.19	62.04 $\pm$ 5.70
vote	<b>96.54</b> $\pm$ <b>1.87</b>	92.26 $\pm$ 3.19	95.70 $\pm$ 2.12	96.54 $\pm$ 2.45	92.24 $\pm$ 3.24	93.54 $\pm$ 2.50
heart2	82.22 $\pm$ 2.93	79.63 $\pm$ 5.71	78.52 $\pm$ 2.81	80.74 $\pm$ 4.06	<b>84.44</b> $\pm$ <b>4.46</b>	83.33 $\pm$ 3.90
dermatology	<b>97.82</b> $\pm$ <b>0.05</b>	96.18 $\pm$ 1.78	94.52 $\pm$ 2.21	91.27 $\pm$ 5.08	97.28 $\pm$ 1.64	96.98 $\pm$ 2.20
pro	<b>88.96</b> $\pm$ <b>2.14</b>	87.96 $\pm$ 3.01	78.54 $\pm$ 1.62	79.13 $\pm$ 2.78	62.38 $\pm$ 3.54	87.57 $\pm$ 2.50
euk	77.38 $\pm$ 1.96	<b>81.42</b> $\pm$ <b>2.06</b>	65.27 $\pm$ 2.92	66.42 $\pm$ 3.47	39.27 $\pm$ 3.43	69.55 $\pm$ 1.30
optdigits	<b>99.11</b> $\pm$ <b>0.38</b>	98.74 $\pm$ 0.39	90.87 $\pm$ 1.09	91.28 $\pm$ 0.40	92.42 $\pm$ 0.75	95.05 $\pm$ 0.90
satimage	<b>91.66</b> $\pm$ <b>0.80</b>	90.38 $\pm$ 0.72	85.64 $\pm$ 1.21	85.33 $\pm$ 0.77	85.41 $\pm$ 0.92	88.14 $\pm$ 1.30
usps	<b>97.49</b> $\pm$ <b>0.22</b>	97.04 $\pm$ 0.47	88.73 $\pm$ 0.46	89.20 $\pm$ 1.00	79.45 $\pm$ 0.59	91.88 $\pm$ 0.60
pendigits	<b>99.56</b> $\pm$ <b>0.12</b>	99.33 $\pm$ 0.17	96.24 $\pm$ 0.31	96.34 $\pm$ 0.41	88.34 $\pm$ 0.65	95.59 $\pm$ 0.30
reuters	<b>98.03</b> $\pm$ <b>0.22</b>	97.15 $\pm$ 0.43	96.90 $\pm$ 0.32	97.18 $\pm$ 0.44	93.52 $\pm$ 0.02	69.54 $\pm$ 0.90
letter	<b>96.99</b> $\pm$ <b>0.21</b>	95.71 $\pm$ 0.19	87.34 $\pm$ 0.68	87.02 $\pm$ 0.66	74.12 $\pm$ 0.97	77.45 $\pm$ 0.30
adult	<b>84.75</b> $\pm$ <b>0.26</b>	83.85 $\pm$ 0.28	84.38 $\pm$ 0.28	83.73 $\pm$ 0.17	80.57 $\pm$ 0.09	82.46 $\pm$ 0.30
w3a	<b>98.86</b> $\pm$ <b>0.04</b>	98.60 $\pm$ 0.06	98.71 $\pm$ 0.05	98.41 $\pm$ 0.10	96.71 $\pm$ 0.20	98.61 $\pm$ 0.30
shuttle	99.77 $\pm$ 0.03	99.93 $\pm$ 0.03	<b>99.97</b> $\pm$ <b>0.02</b>	99.96 $\pm$ 0.02	98.57 $\pm$ 0.24	96.83 $\pm$ 0.30
web	<b>98.94</b> $\pm$ <b>0.05</b>	98.89 $\pm$ 0.06	98.79 $\pm$ 0.09	98.50 $\pm$ 0.13	96.71 $\pm$ 0.21	98.70 $\pm$ 0.60
ijcnn1	98.31 $\pm$ 0.07	<b>98.48</b> $\pm$ <b>0.04</b>	98.40 $\pm$ 0.09	98.11 $\pm$ 0.10	90.69 $\pm$ 0.26	92.29 $\pm$ 0.30
intrusion	<b>99.77</b> $\pm$ <b>0.02</b>	88.20 $\pm$ 1.06	58.01 $\pm$ 26.6	87.66 $\pm$ 3.79	49.75 $\pm$ 30.7	65.15 $\pm$ 15.00
Average	<b>91.29</b> $\pm$ <b>1.26</b>	90.05 $\pm$ 2.06	84.60 $\pm$ 3.64	86.18 $\pm$ 2.84	79.96 $\pm$ 3.58	84.45 $\pm$ 2.80
Ranks	<b>1.500</b>	2.500	4.041	3.9583	5.0625	3.9375

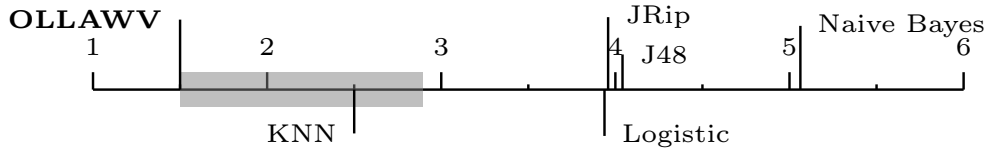


Fig. 1.8.: Bonferroni-Dunn test for Accuracy

significantly worse than OLLAWV, except for KNN.

#### 1.4 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

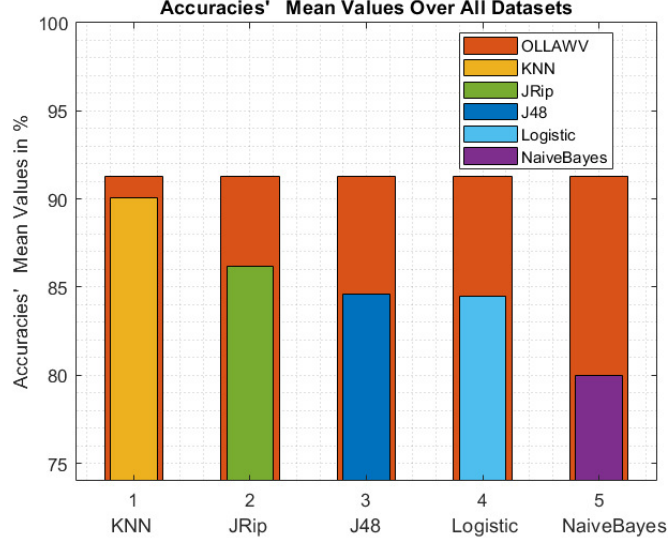


Fig. 1.9.: Mean accuracy over all datasets for OLLAWV and the 5 non-SVM state-of-the-art methods.

proposed. The model is updated by changing the weight  $\alpha_i$  of 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.