

Engineering fast multilevel support vector machines

Ehsan Sadrfaridpour¹, Talayeh Razzaghi², Ilya Safro^{1*}

1- School of Computing, Clemson University, Clemson SC, USA

2- Department of Industrial Engineering, New Mexico State University, Las Cruces NM, USA

February 1, 2018

Abstract

The computational complexity of solving nonlinear support vector machine (SVM) is prohibitive on large-scale data. In particular, this issue becomes very sensitive when the data represents additional difficulties such as highly imbalanced class sizes. Typically, nonlinear kernels produce significantly higher classification quality to linear kernels but introduce extra kernel and model parameters which computationally expensive fitting. This increases the quality but also reduces the performance dramatically. We introduce a generalized fast multilevel framework for regular and weighted SVM and discuss several versions of its algorithmic components that lead to a good trade-off between quality and time. Our framework is implemented using PETSc which allows an easy integration with scientific computing tasks. The experimental results demonstrate significant speed up compared to the state-of-the-art nonlinear SVM libraries.

1 Introduction

Support vector machine (SVM) is one of the most well-known supervised classification methods that has been extensively used in such fields as disease diagnosis, text categorization, and fraud detection. Training nonlinear SVM classifier (such as Gaussian kernel based) requires solving convex quadratic programming (QP) model whose running time can be prohibitive for large-scale instances without using specialized acceleration techniques such as sampling, boosting, and hierarchical training. Another typical reason of increased running time is complex data sets (e.g., when the data is noisy, imbalanced, or incomplete) that require using model selection techniques for finding the best model parameters.

The motivation behind this work was extensive applied experience with hard, large-scale, industrial (not necessarily highly heterogeneous) data sets for which fast *linear* SVMs produced extremely low quality results (as well as many other fast methods), and various nonlinear SVMs exhibited a strong trade off between running time and quality. It has been noticed in multiple works that many different real-world data sets have a strong underlying multiscale (in some works called hierarchical) structure [35, 31, 37, 66] that can be discovered through careful definitions of coarse-grained resolutions. Not surprisingly, we found that among fast methods the hierarchical nonlinear SVM was the best candidate for producing most satisfying results in a reasonable time

*Corresponding author isafro@clemson.edu

[2]. Although, several successful hierarchical SVM techniques [75, 26] have been developed since massive popularization of SVM, we found that most existing algorithms do not sustainably produce high-quality results in a short running time, and the behavior of hierarchical training is still not well studied. This is in contrast to a variety of well studied unsupervised clustering approaches [6, 52, 60].

In this paper, we discuss several multilevel techniques for engineering multilevel SVMs demonstrating their (dis)advantages and generalizing them in a framework inspired by the algebraic multigrid and multiscale optimization strategies [7]. We deliberately omit the issues related to parallelization of multilevel frameworks as it has been discussed in a variety of works related to multilevel clustering, partitioning, and SVM QP solvers. Our goal is to demonstrate fast and scalable sequential techniques focusing on different aspects of building and using multilevel learning with regular and weighted SVM. Also, we focus only on nonlinear SVMs because (a) not much improvement can be introduced or required in practice to accelerate linear SVMs, and (b) in many hard practical cases, the quality of linear SVMs is incomparable to that of nonlinear SVMs. The most promising and stable version of our multilevel SVMs are implemented in PETSc [3] which is a well known scientific computing library. PETSc was selected because of its scalability of linear algebra computations on large sparse matrices and available software infrastructure for future parallelization. Our implementation also addresses a critical need [5] of adding data analysis functionality to broadly used scientific computing software.

1.1 Computational challenges

There is a number of basic challenges one has to address when applying SVM which we successfully tackle with the multilevel framework, namely, QP solver complexity for large-scale data, imbalanced data, and SVM model parameter fitting.

Large-scale data The baseline SVM classifier is typically formulated as a convex QP problem whose solvers scale between $\mathcal{O}(n^2)$ to $\mathcal{O}(n^3)$ [25]. For example, the solver we compare our algorithm with, namely, LibSVM [11], which is one of the most popular and fast QP solvers, scales between $\mathcal{O}(n_f n_s^2)$ to $\mathcal{O}(n_f n_s^3)$ subject to how effectively the cache is exploited in practice, where n_f and n_s are the number of features and samples, respectively. Clearly, this complexity is prohibitive for nonlinear SVM models applied on practical big data without using parallelization, high-performance computing systems or another special treatment.

Imbalanced data The imbalanced data is one of the issues in which SVM often outperforms many fast machine learning methods. This problem occurs when the number of instances of one class (negative or majority class) is substantially larger than the number of instances that belong to the other class (positive or minority class). In multi-class classification, the problem of imbalanced data is even bolder and use of the standard classification methods become problematic in the presence of big and imbalanced data [45]. This might dramatically deteriorate the performance of the algorithm. It is worth noticing that there are cases in which correct classification of the smaller class is more important than misclassification of the larger class [67]. Fault diagnosis [72, 80], anomaly detection [34, 68], medical diagnosis [48] are some of applications which are known to suffer of this problem. Imbalanced data was one of our motivating factors because we have noticed that most standard SVM solvers do not behave well on it. In order to reduce the effect of minority class misclassification in highly imbalanced data, an extension of SVM, namely, the cost-sensitive SVM

(whose extensions are also known as weighted or fuzzy SVM) [43], was developed for imbalanced classification problems. In cost-sensitive SVM, a special control of misclassification penalization is introduced as a part of SVM model.

Parameter tuning The quality of SVM models is very sensitive to the parameters (such as penalty factors of misclassified data) especially in case of using kernels that typically introduce extra parameters. There are many different parameter tuning approaches such as [4, 77, 44, 12, 10, 1, 47, 38]. However, in any case, tuning parameters requires multiple executions of the training process for different parameters and due to the k-fold cross-validation which significantly increases the running time of the entire framework. In our experiments with industrial and healthcare data, not surprisingly, we were unable to find an acceptable quality SVM models without parameter fitting (also known as model selection [17, 76, 65]) which also motivated our work.

1.2 Related work

Multiple approaches have been proposed to improve the performance of SVM solvers. Examples include efficient serial algorithms that use a cohort of decomposition techniques [54], shrinking and caching [30], and fast second order working set selection [22]. A popular LibSVM solver [11] implements the sequential minimal optimization (SMO) algorithm. In the cases of simple data for which nonlinear SVM is not required such approaches as LibLINEAR [21] demonstrate excellent performance for linear SVM using a coordinate descent algorithm which is very fast but, typically, not suitable for complex or imbalanced data. Another approach to accelerate the QP solvers is a chunking [30], in which the models are solved iteratively on the subsets of training data until the global optimum is achieved.

A typical acceleration of support vector machines is done through parallelization and training on high-performance computing systems using interior-point methods (IPM) [49] applied on the dual problem which is a convex QP. The key idea of the primal-dual IPM is to remove inequality constraints using a barrier function and then resort to the iterative Newton’s method to solve the KKT system of the dual problem. For example, in PSVM [78], the algorithm reduces memory use, and parallelizes data loading and computation in IPM. It improves the decomposition-based LibSVM from $O(n^2)$ to $O(np^2/m)$, where m is a number of processors (or heterogeneous machines used), and p is a column dimension of a factorized matrix that is required for effective distribution of the data. The HPSVM solver [40] is also based on solving the primal-dual IPM and uses effective parallelism of factorization. The approach is specifically designed to take maximal advantage of the CPU-GPU collaborative computation with the dual buffers 3-stage pipeline mechanism, and efficiently handles large-scale training datasets. In HPSVM, the heterogeneous hierarchical memory is explored to optimize the bottleneck of data transfer. The P-packSVM [79] parallelizes the stochastic gradient descent solver of SVM that directly optimizes the primal objective with the help of a distributed hash table and sophisticated data packing strategy. Other works utilize many-core GPUs to accelerate SMO [55], and other architectures [74].

One of the most well known works in which hierarchical SVM technique was introduced to improve the performance and quality of a classifier is [75]. The coarsening consists of creating a hierarchical clustered representation of the data points that are merged pairwise using Euclidean distance criterion. In this work, only linear classifiers are discussed and no inheritance and refinement of model parameters was introduced. A similar hierarchical clustering framework was proposed for non-linear SVM kernels in combination with feature selection techniques to develop

an advanced intrusion detection system [27]. Another coarsening approach that uses k-means clustering was introduced in [28]. *In all these works, the quality of classifiers strictly depends on how well the data is clustered using a particular clustering method applied on it.* Our coarsening scheme is more gradual and flexible than the clustering methods in these papers. Most of them, however, can be generalized as algebraic multigrid restriction operators (will be discussed further) in special forms. Also, in our frameworks, we emphasize several important aspects of training such as coarse level models, imbalanced coarsening, and parameter learning that are typically not considered in hierarchical SVM frameworks.

Multilevel Divide-and-Conquer SVM (DC-SVM) was developed using adaptive clustering and early prediction strategy [28]. It outperforms previously mentioned methods, so we compare the computational performance and quality of classification for both DC-SVM and our proposed framework. The training time of DC-SVM for a *fixed set of parameters* is fast. However, in order to achieve high quality classifiers a parameter fitting is typically required. While DC-SVM with parameter fitting is faster than state-of-the-art *exact* SVMs, it is significantly slower than our proposed framework. Our experimental results (that include the parameter fitting component) show significant performance improvement on benchmark data sets in comparison to DC-SVM.

In several works, a scalable parallelization of hierarchical SVM frameworks is developed to minimize the communication [73, 25, 18]. Such techniques can be used on top of our framework. Successful results obtained using hierarchical structures have been shown specifically for multi-class classification [14, 26, 33, 56]. Another relevant line of research is related to multilevel clustering and segmentation methods [36, 23, 66]. They produce solutions at different levels of granularity which makes them suitable for visualization, aggregation of data, and building a hierarchical solution.

1.3 Multilevel algorithmic frameworks

In this paper, we discuss a practical construction of multilevel algorithmic frameworks (MAF) for SVM. These frameworks are inspired by the multiscale optimization strategies [7]. (We note that there exist several frameworks termed multilevel SVMs. These, however, correspond to completely different ideas. We preserve the terminology of multilevel, and multiscale optimization algorithms.) The main objective of multilevel algorithms is to construct a hierarchy of problems (coarsening), each approximating the original problem but with fewer degrees of freedom. This is achieved by introducing a chain of successive restrictions of the problem domain into low-dimensional or smaller-size domains and solving the coarse problems in them using local processing (uncoarsening) [46, 19]. The MAF combines solutions obtained by the local processing at different levels of coarseness into one global solution. Such frameworks have several key advantages that make them attractive for applying on large-scale data: they typically exhibit linear complexity (see Sec. 3.3), and are relatively easily parallelized. Another advantage of the MAF is its heterogeneity, expressed in the ability to incorporate external appropriate optimization algorithms (as a refinement) in the framework at different levels. These frameworks are extremely successful in various practical machine learning tasks such as clustering [53], segmentation [66], and dimensionality reduction [46].

The major difference between typical computational optimization MAF, and those that we introduce for SVM is the output of the model. In SVM, the main output is the set of the support vectors which is usually much smaller at all levels of the multilevel hierarchy than the total number of data points at the corresponding levels. We use this observation in our methods by redefining the training set during the uncoarsening and making MAF scalable. In particular, we inherit the support vectors from the coarse scales, add their neighborhoods, and refine the support vectors

at all scales. In other words, we improve the separating hyperplane throughout the hierarchy by gradual refinement of the support vectors until a global solution at the finest level is reached. In addition, we inherit the parameters of model selection and kernel from the coarse levels, and refine them throughout the uncoarsening.

1.4 Our contribution

We introduce novel methods of engineering fast and high quality multilevel frameworks for efficient and effective training of nonlinear SVM classifiers. We also summarize and generalize existing [58, 57] approaches. We discuss various coarsening strategies, and introduce a weighted aggregation framework inspired by the algebraic multigrid [7] which significantly improves and generalizes all of them. Without any loss in the quality of classifiers, multilevel SVM frameworks exhibit substantially faster running times and are able to generate *several* classifiers at different coarse-grained resolutions in one complete training iteration which also helps to interpret these classifiers qualitatively (see Section 3.4.8).

The proposed multilevel frameworks are particularly effective on imbalanced data sets where fitting model parameters is the most computationally expensive component. The proposed multilevel frameworks can be parallelized as any algebraic multigrid algorithm and their superiority is demonstrated on several publicly available and industrial data sets. The performance improvement over the best sequential state-of-the-art nonlinear SVM libraries with high classification quality is significant. For example, on the average, for large data sets we boost the performance 491 times over LibSVM and 45 times over the DC-SVM (which was chosen because of its superiority over other hierarchical methods mentioned above). On some large datasets, a full comparison was impossible because of non-realistic running time of the competitive approaches which demonstrates superiority of the proposed method.

2 Preliminaries

We define the optimization problems underlying SVM models for binary classification. Given a set \mathcal{J} that contains n data points $x_i \in \mathbb{R}^d$, $1 \leq i \leq n$, we define the corresponding labeled pairs (x_i, y_i) , where each x_i belongs to the class determined by a given label $y_i \in \{-1, 1\}$. Data points with positive labels are called the *minority* class which is denoted by \mathbf{C}^+ with $|\mathbf{C}^+| = n^+$. The rest of the points belongs to the *majority* class which is denoted by \mathbf{C}^- , where $|\mathbf{C}^-| = n^-$. Solving the following convex optimization problem by finding w , and b produces a hyperplane with maximum margin between \mathbf{C}^+ , and \mathbf{C}^-

$$\begin{aligned}
 & \text{minimize} && \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i && (1) \\
 & \text{subject to} && y_i(w^T \phi(x_i) + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\
 & && \xi_i \geq 0, \quad i = 1, \dots, n.
 \end{aligned}$$

The mapping of data points to higher dimensional space is done by $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ ($d \leq p$) to make two classes separable by a hyperplane. The term slack variables $\{\xi_i\}_{i=1}^n$ are used to penalize misclassified points. The parameter $C > 0$ controls the magnitude of the penalization. The primal formulation is shown at (1) which is known as the *soft margin* SVM [71].

The weighted SVM (WSVM) addresses imbalanced problems with assigning different weights to classes with parameters C^+ and C^- . The set of slack variables is split into two disjoint sets $\{\xi_i^+\}_{i=1}^{n^+}$, and $\{\xi_i^-\}_{i=1}^{n^-}$, respectively. In WSVM, the objective of (1) is changed into

$$\text{minimize} \quad \frac{1}{2}\|w\|^2 + C^+ \sum_{i=1}^{n^+} \xi_i^+ + C^- \sum_{j=1}^{n^-} \xi_j^-. \quad (2)$$

Solving the Lagrangian dual problem using kernel functions $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ produces a reliable convergence which is faster than methods for primal formulations (1) and (2). In our framework, we use the sequential minimal optimization solver implemented in LibSVM library [11]. The role of kernel functions is to measure the similarity for pairs of points x_i and x_j . We present computational results with the Gaussian kernel (RBF), $\exp(-\gamma\|x_i - x_j\|^2)$, which is known to be generally reliable when no additional assumptions about the data are known. Experiments with other kernels exhibit improvements that are similar to those with RBF.

In order to achieve an acceptable quality of the classifier, many difficult data sets require reinforcement of (W)SVM with tuning methods for such model parameters as C , C^+ , C^- , and kernel function parameters (e.g., the bandwidth parameter γ for RBF kernel function). This is one of the major sources of running time complexity of (W)SVM models which we are aiming to improve.

In our framework we use the adapted nested uniform design (NUD) model selection algorithm to fit the parameters [29] which is a popular model selection technique for (W)SVM. The main intuition behind NUD is that it finds the close-to-optimal parameter set in an iterative nested manner. The optimal solution is calculated in terms of maximizing the required performance measure (such as accuracy and G-mean). Although, we study binary classification problems, it can easily be extended to the multi-class classification using either directed multi-class classification or transforming the problem into multiple independent binary (W)SVMs that can be processed independently in parallel.

Two-level problem In order to describe the (un)coarsening algorithms, we introduce the two-level problem notation that can be extended into full multilevel hierarchy (see Figure 1). We will use subscript $(\cdot)_f$ and $(\cdot)_c$ to represent fine and coarse variables, respectively. For example, the data points of two consecutive levels, namely, fine and coarse, will be denoted by \mathcal{J}_f , and \mathcal{J}_c , respectively. The sets of fine and coarse support vectors are denoted by \mathbf{sv}_f , and \mathbf{sv}_c , respectively. We will also use a subscript in the parentheses to denote the number of level in the hierarchy where appropriate. For example, $\mathcal{J}_{(i)}$ will denote the set of data points at level i .

Proximity graphs All multilevel (W)SVM frameworks discussed in subsequent sections are based on different coarsening schemes for creating a hierarchy of data proximity graphs. Initially, at the finest level, \mathcal{J} is represented as two k -nearest neighbor (k NN) graphs $G_0^+ = (\mathbf{C}^+, E^+)$, and $G_0^- = (\mathbf{C}^-, E^-)$ for minority and majority classes, respectively, where each $x_i \in \mathbf{C}^{+(-)}$ corresponds to a node in $G_0^{+(-)}$. A pair of nodes is connected with an edge if one of them belongs to a set of k -nearest neighbors of another. In practice, we are using *approximate* k -nearest neighbors graphs (Ak NN) as our experiments with the *exact* nearest neighbor graphs do not demonstrate any improvement in the quality of classifiers whereas computing them is computationally expensive. In the computational experiments, we used FLANN library [50, 51]. Results obtained with

other approximate nearest neighbor search algorithms are found to be not significantly different. Throughout the multilevel hierarchies, in two-level representation, the fine and coarse level graphs will be denoted by $G_f^{+(-)} = (\mathcal{J}_f^{+(-)}, E_f^{+(-)})$, and $G_c^{+(-)} = (\mathcal{J}_c^{+(-)}, E_c^{+(-)})$, respectively. All coarse graphs, except $G_0^{+(-)}$ are obtained using respective coarsening algorithm.

3 Multilevel support vector machines

The multilevel frameworks discussed in this paper include three phases (see Figure 1), namely, gradual training set coarsening, coarsest support vector learning, and gradual support vector refinement (uncoarsening). In the training set coarsening phase we create a hierarchy of coarse training set representations, in which each next-coarser level contains a fewer number of points than in the previous level such that the coarse level learning problem approximates the fine level problem. The coarse level training points are not necessarily the same fine level points (such as in [58]) or their strict small clusters (such as in [75]). When the size of training set is sufficiently small to apply a high quality training algorithm for given computational resources, the set of coarsest support vectors and model parameters are trained. We denote by $M^+(-)$ the upper limit for the sizes of coarsest training sets. In the uncoarsening, both the support vectors and model parameters are inherited from the coarse level and improved using local refinement at the fine level. The uncoarsening is continued from the coarsest to the finest levels as is shown in Figure 1. Separate coarsening hierarchies are created for classes \mathbf{C}^+ , and \mathbf{C}^- , independently. The main driving routine, `mlsvm`, of a multilevel WSVM framework is presented in Algorithm 1. The SVM cost-sensitive framework is designed similarly with a parameter C , see Eq. (1).

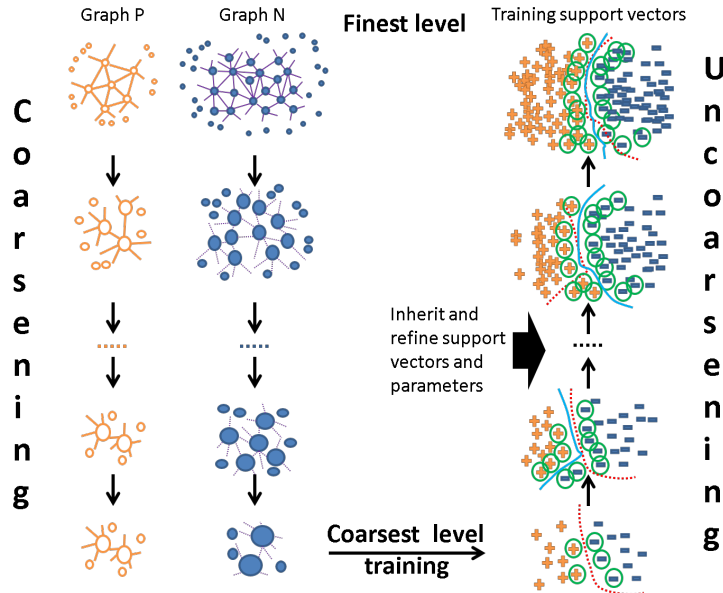


Figure 1: Multilevel SVM coarsening-uncoarsening framework scheme.

Algorithm 1 mlsvm: multilevel (W)SVM main driving routine

input: $\mathcal{J}_f = \mathbf{C}_f^+ \cup \mathbf{C}_f^-, G_f^+, G_f^-, M^+, M^-$

- 1: **if** $|\mathcal{J}_f| \leq M^+ + M^-$ **then**
- 2: $(\mathbf{sv}_f, C^+, C^-, \gamma) \leftarrow$ train (W)SVM model on \mathcal{J}_f (including NUD)
- 3: **else**
- 4: **if** $|\mathbf{C}_f^+| \leq M^+$ **then** $\mathbf{C}_c^+ \leftarrow \mathbf{C}_f^+; G_c^+ \leftarrow G_f^+$
- 5: **else** $(\mathbf{C}_c^+, G_c^+) \leftarrow$ coarsen(\mathbf{C}_f^+, G_f^+)
- 6: **if** $|\mathbf{C}_f^-| \leq M^-$ **then** $\mathbf{C}_c^- \leftarrow \mathbf{C}_f^-; G_c^- \leftarrow G_f^-$
- 7: **else** $(\mathbf{C}_c^-, G_c^-) \leftarrow$ coarsen(\mathbf{C}_f^-, G_f^-)
- 8: $(\mathbf{sv}_c, C^+, C^-, \gamma) \leftarrow$ mlsvm($\mathbf{C}_c^+ \cup \mathbf{C}_c^-, G_c^+, G_c^-, M^+, M^-$)
- 9: $\mathbf{sv}_f \leftarrow$ uncoarsen(\mathbf{sv}_c)
- 10: $(\mathbf{sv}_f, C^+, C^-, \gamma) \leftarrow$ refine($\mathbf{sv}_f, C^+, C^-, \gamma$)
- 11: **return** $\mathbf{sv}_f, C^+, C^-, \gamma$

3.1 Iterative Independent Set Coarsening

The iterative independent set (mlsvm-IIS) coarsening consists of several iterative passes in each of which a set of fine points is selected and added to the set of coarse points \mathbf{C}_c^+ or \mathbf{C}_c^- . In order to cover the space of points uniformly, this is done by selecting independent sets of nodes in G_f^+ , and G_f^- . The independent set is a set of vertices in a graph whose node-induced subgraph has no edges. We describe this coarsening in details in [58].

Coarsening We start with selecting a random independent set of nodes (or points), I_0 , using one pass over all nodes (i.e., choose a random node to I_0 , eliminate it with its neighbors from the graph, and choose the next node). The obtained independent set I_0 is added to the set of coarse points. Then, we remove I_0 from the graph and repeat the same process of finding another independent set I_1 and adding it to the coarse set on the rest of the graph. The iterations are repeated until $\sum_k |I_k| \leq Q|\mathcal{J}_f^{+(-)}|$, where Q is a parameter controlling the size of coarse level space.

Coarsest level At the coarsest level r , when $|\mathcal{J}_{(r)}| \leq M^+ + M^- \ll |\mathcal{J}_{(0)}|$, we can apply an exact (or computationally expensive) algorithm for training the coarsest classifier. Typically, $|\mathcal{J}_{(r)}|$ depends on the available computational resources. However, one can also consider some criteria of separability between $\mathbf{C}_{(r)}^+$, and $\mathbf{C}_{(r)}^-$ [70], i.e., if a fast test exists or some helpful data properties are known. In all our experiments, we used a simple criterion limiting $|\mathcal{J}_{(r)}|$ to 500. Processing the coarsest level includes an application of NUD [29] model selection to get high-quality classifiers on the difficult data sets. To this end, we obtained a solution of the coarsest level, namely, $\mathbf{sv}_{(r)}, C_{(r)}^+, C_{(r)}^-$, and $\gamma_{(r)}$.

Uncoarsening Given the solution of coarse level c , the primary goal of the uncoarsening is to interpolate and refine this solution for the current fine level f . Unlike many other multilevel algorithms, in which the inherited coarse solution contains projected variables only, in our case, we inherit not only \mathbf{sv}_c but also parameters for model selection. This is important because the model selection is an extremely time-consuming component of (W)SVM, and can be prohibitive at fine levels of the hierarchy. However, at the coarse levels, when the problem is much smaller than the original, we can apply much heavier methods for the model selection almost without any loss in the total complexity of the framework.

Algorithm 2 mlsvm-IIS: uncoarsening and refinement at level f

```
1: Find nearest neighbors  $N_f^+$  and  $N_f^-$  of support vectors  $\text{sv}_c$  in  $G_f^+$  and  $G_f^-$ 
2:  $T \leftarrow \text{sv}_c \cup N_f^+ \cup N_f^-$  ▷  $T$  is a training set for refinement
3: if  $|T| < Q_t$  then
4:    $C^O \leftarrow (C_c^+, C_c^-)$ ;  $\gamma^O \leftarrow \gamma_c$ 
5:   Run NUD using the initial center  $(C^O, \gamma^O)$  and train (W)SVM on  $T$ 
6: else
7:    $C_f^+ \leftarrow C_c^+$ ;  $C_f^- \leftarrow C_c^-$ ;  $\gamma_f \leftarrow \gamma_c$  ▷ Inherit the coarse parameters
8: end if
9: if  $|T| \geq Q_c$  then
10:  Partition  $T$  into  $K$  (almost) equal size clusters
11:   $\forall k \in K$  find  $P$  nearest opposite-class clusters
12:  Train (W)SVMs on pairs of nearest clusters only
13:   $\text{sv}_f \leftarrow$  support vectors from all trained models
14: else
15:   $\text{sv}_f \leftarrow$  train (W)SVM on  $T$ 
16: end if
17: Return  $\text{sv}_f, C_f^+, C_f^-, \gamma_f$ 
```

The uncoarsening and refinement are presented in Alg. 2. After the coarsest level is solved exactly and reinforced by the model selection (line 2 in Alg. 1), the coarse support vectors sv_c and their nearest neighbors (in our experiments no more than 5) in both classes (i.e., N_f^+ and N_f^-) initialize the fine level training set T (lines 1-2 in Alg. 2). If $|T|$ is still small (relatively to the existing computational resources, and the initial size of the data) for applying model selection, i.e., it less than a parameter Q_t , then we use coarse parameters $C_c^{+(-)}$, and γ_c as initializers for the current level NUD grid search, and retrain (lines 4-5 in Alg. 2). Otherwise, the coarse $C_c^{+(-)}$, and γ_c are inherited in $C_f^{+(-)}$, and γ_f (line 7 in Alg. 2). Then, being large for a direct application of the (W)SVM, T is partitioned into several equal size clusters, and pairs of nearest opposite clusters are retrained. The obtained solutions are contributed to sv_f (lines 10-13 in Alg. 2). *We note that partition-based retraining can be done in parallel, as different pairs of clusters are independent. Moreover, the total complexity of the algorithm does not suffer from reinforcing the partition-based retraining with model selection.*

This coarsening scheme is one of the fastest and easily implementable. While the entire framework (including uncoarsening) is definitely much faster than a regular (W)SVM solver such as libSVM (which is used in our implementation as a refinement), it is not the fastest among the multilevel SVM frameworks. This is a typical trade-off in discrete multilevel frameworks [15, 62], namely, when the quality of coarsening suffers, the most work is done at the refinement. A similar independent set coarsening approach was used in multilevel dimensionality reduction [23]. However, in contrast to that coarsening scheme, we found that using only one independent set (including possible maximization of it) does not lead to the best quality of classifiers. Instead, a more gradual coarsening makes the framework much more robust to the changes in the parameters and the shape of data manifold.

3.2 Algebraic multigrid coarsening

The algebraic multigrid (AMG) (W)SVM coarsening scheme (mlsvm-AMG) is inspired by the AMG aggregation solvers for computational optimization problems such as [63, 39, 35, 64]. Its first version was briefly presented in [61]. The AMG coarsening generalizes the independent set and clustering [28] based approaches leveraging a high quality coarsening and flexibility of AMG which belongs to the same family of multiscale learning strategies with the same main phases, namely, coarsening, coarsest scale learning, and uncoarsening. Instead of eliminating a subset of the data points, in AMG coarsening, the original problem is gradually restricted to smaller spaces by creating *aggregates of fine data points and their fractions* (which is an important feature of AMG), and turning them into the data points at coarse levels. The main mechanism underlying the coarsening phase is the AMG [69, 7] which successfully helps to identify the interpolation operator for obtaining a fine level solution from the coarse aggregates. In the uncoarsening phase, the solution obtained at the coarsest level (i.e., the support vectors and parameters) is gradually projected back to the finest level by interpolation and further local refinement of support vectors and parameters. A critical difference between AMG approach and the earlier work of Razzaghi et al. [58] is that in AMG approach the coarse level support vectors are not the original data points prolonged from the finest level. Instead, they are centroids of aggregates that contain both full fine-level data points and their fractions.

Framework initialization The AMG framework is initialized with $G_0^{+(-)}$ with the edge weights that represent the strength of connectivity between nodes in order to “simulate” the following interpolation scheme applied at the uncoarsening, in which strongly coupled nodes can interpolate solution to each other. In the classifier learning problems, this is expressed as a similarity measure between points. We define a distance function between nodes (or corresponding data points) as an inverse of the Euclidean distance. More advanced distance measure approaches such as [8, 13] are often essential in similar multilevel frameworks.

Coarsening Phase We describe the two-level process of obtaining the coarse level training set \mathcal{J}_c^+ with corresponding G_c^+ from the current fine level G_f^+ and its training set (e.g., the transition from level f to c). The majority class is coarsened similarly.

The process is started with selecting seed nodes that will serve as centers of coarse level nodes, i.e., the aggregates at level f . Coarse nodes will correspond to the coarse data points at level c . Structurally, each aggregate can include one full seed f -level point, and possibly several other f -level points and their fractions. Intuitively, it is equivalent to grouping points in \mathcal{J}_f^+ into many small subsets allowing intersections, where each subset of nodes will correspond to a coarse point at level c . During the aggregation process, most coarse points will correspond to aggregates of size greater than 1 (because, throughout the hierarchy, they accumulate many fine points and their fractions), so we introduce the notion of a volume $v_i \in \mathbb{R}_+$ for all $i \in \mathcal{J}_f^+$ to reflect the importance of a point or its capacity that includes finest-level aggregated points and their fractions. We also introduce the edge weighting function $w : E_f^+ \rightarrow \mathbb{R}_{\geq 0}$ to reflect the strength of connectivity and similarity between nodes.

In Algorithm 3, we show the details of AMG coarsening. In the first step, we compute the future-volumes ϑ_i for all $i \in \mathcal{J}_f^+$ to determine the order in which f -level points will be tested for

declaring them as seeds (line 2), namely,

$$\vartheta_i = v_i + \sum_{j \in \mathcal{J}_f^+} v_j \cdot \frac{w_{ji}}{\sum_{k \in \mathcal{J}_f^+} w_{jk}}. \quad (3)$$

The future-volume ϑ_i is defined as a measure (that is often used in multilevel frameworks [62]) of how much an aggregate seeded by a point i may potentially grow at the next level c .

We assume that in the finest level, all volumes are ones. We start with selecting a dominating set of seed nodes $C \subset \mathcal{J}_f^+$ to initialize aggregates. Nodes that are not selected to C will belong to F such that $\mathcal{J}_f^+ = F \cup C$. Initially, the set F is set to be \mathcal{J}_f^+ , and $C = \emptyset$ since no seeds have been selected (line 1). After that, points with ϑ_i that is exceptionally larger than the average $\bar{\vartheta}$ are transferred to C as the most “representative” points (line 3). Then, all points in F are accessed in the decreasing order of ϑ_i updating C iteratively (lines 6-11), namely, if with the current C , and F , for point $i \in F$, $\sum_{j \in C} w_{ij} / \sum_{j \in \mathcal{J}_f^+} w_{ij}$ is less than or equal to some threshold Q , i.e., the point is not strongly coupled to already selected points in C , then i is moved from F to C . The points with larger future-volumes usually have a better chance to be selected to C to serve as centers of future coarse points. Selecting too few seeds (and then coarse level points) causes “overcompressed” coarser level which typically leads to the classification quality drop. Therefore, in order to keep sufficiently many points at the coarse level, the parameter Q is set to 0.4-0.6. It has been observed that in most AMG algorithms, $Q > 0.6$ is not required (however, this depends on the type and goals of aggregation). In our experiments $Q = 0.5$, and $\eta = 2$. Other similar values do not significantly change the results.

Algorithm 3 mlsvm-AMG: selecting seeds for coarse nodes

```

1:  $C \leftarrow \emptyset, F \leftarrow \mathcal{J}_f^+$ 
2: Calculate  $\forall i \in F$   $\vartheta_i$ , and the average  $\bar{\vartheta}$ 
3:  $C \leftarrow$  nodes with  $\vartheta_i > \eta \cdot \bar{\vartheta}$ 
4:  $F \leftarrow V_f \setminus C$ 
5: Recompute  $\vartheta_i \forall i \in F$ 
6: Sort  $F$  in descending order of  $\vartheta$ 
7: for  $i \in F$  do
8:   if  $\left( \frac{\sum_{j \in C} w_{ij}}{\sum_{j \in \mathcal{J}_f^+} w_{ij}} \right) \leq Q$  then
9:     move  $i$  from  $F$  to  $C$ 
10:  end if
11: end for
12: return  $C$ 

```

When the set C is selected, we compute the AMG interpolation matrix $P \in \mathbb{R}^{|\mathcal{J}_f^+| \times |C|}$ that is defined as

$$P_{ij} = \begin{cases} w_{ij} / \sum_{k \in \Gamma_i} w_{ik} & \text{if } i \in F, j \in \Gamma_i \\ 1 & \text{if } i \in C, j = I(i) \\ 0 & \text{otherwise} \end{cases}, \quad (4)$$

where $\Gamma_i = \{j \in C \mid ij \in E_f^+\}$ is the set of i th seed neighbors, and $I(i)$ denotes the index of a coarse point at level c that corresponds to the fine level aggregate around seed $i \in C$. Typically, in AMG methods, the number of non-zeros in each row is limited by the parameter called the interpolation order or caliber [7] (see discussion about R and Table 7). This parameter controls the complexity of a coarse-scale system (the number of non-zero elements in the matrix of coarse k -NN graph). It limits the number of fractions a fine point can be divided into (and thus attached to the coarse points). If a row in P contains too many non-zero elements then it is likely to increase the number of non-zeros in the coarse graph matrix. In multigrid methods, this number is usually controlled by different approaches that measure the strength of connectivity (or importance) between fine and coarse variables (see discussion and implementation in [59]).

Using the matrix P , the aggregated data points and volumes for the coarse level are calculated. The edge between points $p = I(i)$ and $q = I(j)$ is assigned with weight $w_{pq} = \sum_{k \neq l} P_{ki} \cdot w_{kl} \cdot P_{lj}$. The volume for the aggregate $I(i)$ in the coarse graph is computed by $\sum_j v_j P_{ji}$, i.e., the total volume of all points is preserved at all levels during the coarsening. The coarse point $q \in \mathcal{J}_c^+$ seeded by $i = I^{-1}(q) \in \mathcal{J}_f^+$ is represented by

$$\sum_{j \in \mathcal{A}_i} P_{j,q} \cdot j, \quad (5)$$

where \mathcal{A}_i is a set of fine points in aggregate i .

The stopping criteria for the coarsening depends on the available computational resources that can be used in order to train the classifier at the coarsest level. In our experiments, the coarsening stops when the size is less than a threshold (typically, 500 points) that ensures a fast performance of the LibSVM dual solver.

Uncoarsening The uncoarsening of AMG multilevel framework is similar to that of the `mlsvm-IIS`. The main difference is in lines 1-2 in Alg. 2. Instead of defining the training set for the refinement at level f as

$$T \leftarrow \text{sv}_c \cup N_f^+ \cup N_f^-,$$

all coarse support vectors are uncoarsened by adding to T all elements of the corresponding aggregates, namely,

$$T \leftarrow \emptyset; \quad \forall p \in \text{sv}_c \quad \forall j \in \mathcal{A}_p \quad T \leftarrow T \cup j. \quad (6)$$

3.3 Complexity of multilevel framework

The complexity of MAF for (W)SVM consists of three parts, namely, generating approximated k -NN graphs of both classes, coarsening and uncoarsening. The complexity of generating approximated k -NN graphs is beyond the scope of this work. We refer the reader to the analysis and comparison of FLANN library that was used in our implementation [50, 51]. When we compare the performance of 1 V-cycle of our solver and FLANN, we do not observe a significant difference between them when running FLANN to find 10 nearest neighbors.

In the coarsening phase, we need to consider the complexity of coarsening the approximated k -NN graphs of \mathbf{C}^+ and \mathbf{C}^- including aggregation of the data points. The complexity of coarsening

is similar to that of AMG applied on graph $G = (V, E)$ which is proportional to $|V| + |E|$, where $|E| \approx k|V|$, where k is the number of nearest neighbors. In our experiments, we found that no data set requires $k > 10$ to improve the quality of classification. Because we do not anticipate to obtain exceptionally high-degree nodes during the coarsening, we also do not expect to observe very fast increasing density of nonzero features (nnz) in data points. Thus, we bound the complexity of coarsening with $\mathcal{O}(\text{nnz}(\mathcal{J}))$ without having hidden coefficients, in practice.

The complexity of the uncoarsening mostly depends on that of the underlying QP solver (such as LibSVM) applied at the refinement stage. Another factor that affects the complexity is the number of support vectors found at each scale which is typically significantly smaller than the number of data points. Typically, the complexity will be approximately $\mathcal{O}(|\mathcal{J}|) + \mathcal{O}(QP\text{Solver}(p \text{ points})) \cdot |\text{support vectors}|/p$, where p is the number of parts, the set of support vectors is split to if partitioning is applied. The overall computational time obtained in our experiments and the amount of work per unit is presented in Section 4. In particular, in Table 15 we demonstrate the computational time per data point and per feature value.

3.4 Engineering multilevel framework

The AMG framework generalizes many multilevel approaches by allowing a “soft” weighted aggregation of points (and their fractions) in contrast to the “strict” clustering [28] and subset based aggregations [58]. In this section we describe a variety of improvements we experimented with to further boost the quality of the multilevel classification framework, and improve the performance of both the training and validation processes in terms of the quality and running time. All of them are applicable in both “strict” or “soft” coarsening schemes.

3.4.1 Imbalanced classification

One of the major advantages of the proposed coarsening scheme is its natural ability to cope with the imbalanced data in addition to the cost-sensitive models solved in the refinement. When the coarsening is performed on both classes simultaneously, and in a small class the number of points reaches an allowed minimum, this level is simply copied throughout the rest of levels required to coarsen the big class. Since the number of points at the coarsest level is small, this does not affect the overall complexity of the framework. Therefore, the numbers of points in both classes are within the same range at the coarsest level regardless of how imbalanced they were at the fine levels. Such training on the balanced data mitigates the imbalance effects and improves the performance quality of trained models.

3.4.2 Coarse level density problem

Both `mlsvm-IIS` and `mlsvm-AMG` do not change the dimensionality during the coarsening which potentially may turn into a significant computational bottleneck for very large-scale data. In many applications, a high-dimensional data is sparse, and, thus, even if the number of points is large, the space requirements are still not prohibitive. Examples include text tf-idf data and categorical features that are converted into a binary representation. While the `mlsvm-IIS` coarsening selects *original* (i.e., sparse) points for the coarse levels, the `mlsvm-AMG` aggregates points using a linear combination such as in Eq. 5. Even when the original $j \in \mathcal{A}_i$ are sparse, the points at coarse levels may eventually become much denser than the original points.

The second type of coarse level density is related to the aggregation itself. When f -level data points are divided into several parts to join several aggregates, the number of edges in coarse graphs is increasing when it is generated by $L_c \leftarrow P^T L_f P$ (subject to L_c diagonal entry correction). Finest level graphs that contain high-degree nodes have a good chance to generate very dense graphs at the coarse levels if their density is not controlled. This can potentially affect the performance of the framework.

The coarse level density problem is typical to most AMG and AMG-inspired approaches. We control it by filtering weak edges and using the order of interpolation in aggregation. The weak edges not only increase the density of coarse levels but also may affect the quality of the training process during the refinement. In `mlsvm`-AMG framework, we eliminate weak edges between i and j if $w_{ij} < \theta \cdot \text{avg}_{ki}\{w_{ki}\}$ and $w_{ij} < \theta \cdot \text{avg}_{kj}\{w_{kj}\}$ where $\text{avg}\{\cdot\}$ is the average of corresponding adjacent edge weights. We experimented with different values of θ between 0.001 and 0.005 which was typically a robust parameter that does not require much attention.

The order of interpolation, r , is the number of nonzeros allowed per row in P . A single nonzero j th entry in row i , $P_{ij} = 1$, means that a fine point i fully belongs to aggregate j which leads to creation of small clusters of fine points without splitting them. Typically, in AMG methods, increasing r improves the quality of solvers making them, however, slower. We experiment with different values of r and conclude that high interpolation orders such as 2 and 4 perform better than 1. In the same time, we observed that there is no practical need to increase it more.

3.4.3 Validation for model selection

The problem of finding an optimal set of parameters (i.e., the model selection) for a better quality is important for complex data sets. Typically, this component is computationally expensive because repetitive training is required for different choices of parameters. A validation data is then required to choose the best trained model. A performance of model selection techniques is affected by the quality and size of the validation data. (We note that the test data for which the computational results are presented remains completely isolated from any training and validation.)

The problem of a validation set choice requires a special attention in multilevel frameworks because the models at the coarse levels should not necessarily be validated on the corresponding coarse data. As such, we propose different approaches to find the most suitable types of validation data. We experimented with *coarse sampling* (CS), *coarse cross k-fold* (CCkF), *finest full* (FF), and *fine sampling* (FS) methods to choose validation set for multilevel frameworks.

CS: The data in \mathcal{J}_i^+ and \mathcal{J}_i^- is sampled and one part of it (in our experiments 10% or 20%) is selected for a validation set for model selection. In other words, the validation is performed on the data at the same level. This approach is extremely fast on the data in which the coarsening is anticipated to be uniform without generating a variability in the density of aggregation in different parts of the data. Typically, its quality is acceptable on homogeneous data. However, qualitatively, this approach may suffer from a small size of the validation data in comparison to the size of test data.

CCkF: In this method we apply a complete k-fold cross validation at all levels using the coarse data from the same level. The disadvantage of this method is that it is more time consuming but the quality is typically improved. During the k-fold cross validation, all data is covered. With this method, the performance measures are improved in comparison to the CS but the quality of the finest level can degrade because of potential overfitting at the coarse levels.

FF: This method exploits a multilevel framework by combining a coarse training set $\mathcal{J}_c^{+(-)}$ with a validation set that is a whole finest level training set $\mathcal{J}_0^{+(-)}$. The idea behind this approach is to choose the best model which increases a required performance measure (such as accuracy, and G-mean) of coarse aggregates with respect to the original data rather than the aggregates. This significantly increases the quality of final models. However, this method is time consuming on very large data sets as all original points participate in validation.

FS: This method resolves the complexity of FF by sampling $\mathcal{J}_0^{+(-)}$ to serve as a validation set at the coarse levels. The size of sampling should depend on computational resources. However, we note that we have not observed any drop in quality if it is more than 10% of the $\mathcal{J}_0^{+(-)}$. Both FF and FS exhibit the best performance measures.

3.4.4 Underlying solver

At all iterations of the refinement and at the coarsest level we used LibSVM [11] as an underlying solver by applying it on the small subsets of data (see lines 5 and 12 in Alg. 2). Depending on the objective Eq. (1) or (2), SVM or WSVM solvers are applied. In this paper we report the results of WSVM in which the objective Eq. (2) is given by

$$\text{minimize} \quad \frac{1}{2}\|w\|^2 + C(W^+ \sum_{i=1}^{n^+} \xi_i^+ + W^- \sum_{j=1}^{n^-} \xi_j^-), \quad (7)$$

where the optimal C and γ are fitted using model selection, and the class importance coefficients are W^+ and W^- . While for the single-level WSVM, the typical class importance weighting scheme is

$$W^+ = \frac{1}{|\mathcal{J}^+|}, \quad W^- = \frac{1}{|\mathcal{J}^-|},$$

in MAF, the aggregated points in each class have different importance due to the different accumulated volume of finer points. The aggregated points which represent more fine points are more important than aggregated points which represent small number of fine points. Therefore, the MAF approach for calculating the class weights is based on sum of the volumes in each class, i.e.,

$$W^+ = \frac{1}{\sum_{i \in \mathcal{J}^+} v_i}, \quad W^- = \frac{1}{\sum_{i \in \mathcal{J}^-} v_i}. \quad (8)$$

This method, however, ignores the importance of each point in its class. We find that the most successful penalty scheme is the one that is personalized per point, and adapt (8) to be

$$\forall i \in \mathcal{J}^{+(-)} \quad W_i = \frac{v_i}{\left(\sum_{j \in \mathcal{J}^{+(-)}} v_j\right)^2}.$$

In other words, we consider the relative volume of the point in its class but also multiply it by an inverse of the total volume of the class which gives more weight to a small class. This helps to improve the correctness of a small class classification.

3.4.5 Expanding training set in refinement

Typically, in many applications, the number of support vectors is much smaller than $|\mathcal{J}|$. This observation allows some freedom in exploring the space around support vectors inherited from coarse levels. This can be done by adding more points to the refinement training set in attempt to improve the quality of a hyperplane. We describe several possible strategies that one can follow in designing a multilevel (W)SVM framework.

Full disaggregation: This is a basic method presented in (6) in which all aggregates of coarse support vectors contribute all their elements to the training set. It is a default method in `mlsvm-AMG`.

k -distant disaggregation: In some cases, the quality can be improved by adding to T the k -distant neighbors of aggregate elements. In other words, after (6), we apply

$$\forall p \in T \quad T \leftarrow T \cup N_f^{+(-)}(p),$$

where $N_f^{+(-)}(p)$ is a set of neighbors of p in $G_f^{+(-)}$ depending on the class of p . Similarly, one can add points within distance k from the aggregates of inherited support vectors. Clearly, this can only improve the quality. However, the refinement training is expected to be increasingly slower especially when the $G_f^{+(-)}$ contains high-degree nodes. Even if the finest level graph nodes are of a small degree, this could happen at the coarse levels if a proper edge filtering and limiting interpolation order are not applied. In very rare cases, we observed a need for adding distance 2 neighbors.

Sampling aggregates: In some cases, the coarse level aggregates may become very dense which does not improve the quality of refinement training. Instead, it may affect the running time. One way to avoid of unnecessary complexity is to sample the elements of aggregates. A better way to sample them than a random sampling (after adding the seed) is to order them by the interpolation weights P_{ij} . The ascending order which gives a preference to the fine points that are split across more than one aggregate was the most successful option in our experiments. Fine non-seed points whose $P_{ij} = 1$ are likely to have high similarity with the seeds which does not improve the quality of the support vectors.

3.4.6 Partitioning in the refinement

When the number of uncoarsened support vectors at level f is too big for the available computational resources, multiple small-size models are trained and either validated or used as a final output. Small-size models are required for applying model selection in a reasonable computational time. For this purpose we propose to partition the training set (see lines 10-13 in Alg. 2) in equal k parts of approximately equal size using fast graph partitioning solvers [9]. Note that applying similar clustering strategies may lead to highly imbalanced parts which will make the whole process of refinement acceleration useless.

In both `mlsvm-AMG` and `mlsvm-IIS`, we leverage the graphs of both classes G_f^+ and G_f^- with the inverses of Euclidean distance between nodes playing the role of edge weights. After both graphs are partitioning, two sets of approximately equal size partitions, Π_f^+ and Π_f^- are created. For each $\pi_i \in \Pi_f^{+(-)}$ we compute its centroid c_i in order to estimate the nearest parts of opposite classes and train multiple models by pairs of parts.

The training by pairs of parts works as follows. For each c_i we find the nearest c_j such that i and j are in different classes and evaluate at most $|\Pi_f^+ + \Pi_f^-|$ models for different choices of (p_i, p_j) pairs (without repetitions). The set of all models is denoted by M_f . We note that the training of such pairs is independent and can be easily parallelized. There are multiple ways one can test (or validate) a point using all models “voting”. The simplest strategy which performs well on many data sets is a majority voting. However, the most successful way to generate a prediction was a voting by the relative distance from the test point t to

$$x_{ij} = \frac{c_i \sum_{q \in p_i} v_q + c_j \sum_{q \in p_j} v_q}{\sum_{q \in p_i} v_q + \sum_{q \in p_j} v_q},$$

the weighted center of the segment connecting c_i and c_j . For all pairs of nearest parts i and j , the label is computed as

$$\text{sign}\left(\frac{\sum_{ij \in M_f} l_{ij}(t) d^{-1}(t, x_{ij})}{\sum_{ij \in M_f} d^{-1}(t, x_{ij})}\right),$$

where $l_{ij}(t)$ is a label of ij model for point t , and $d(\cdot, \cdot)$ is a distance function between two points. We experimented with several distance functions that emphasize the importance of parts’ proximity, namely, Euclidean, exponential, and Manhattan. The quality of final models obtained using Euclidean distance was the highest.

3.4.7 Model Selection

The MAF allows a flexible design for model selection techniques such as various types of parameter grid search [12], NUD [29] that we use in our computational experiments, and other search approaches [44, 4, 76]. A mechanism that typically works behind most of such search techniques evaluates different combinations of parameters (such as C^+ , C^- , and γ) and chooses the one that exhibits the best performance measure. Besides the general applicability of model selection because the size of inherited and disaggregated T (in the uncoarsening of `mlsvm-IIS` and `mlsvm-AMG`) is typically smaller than that of the corresponding training set, the MAF has two following advantages.

Fast parameter search: In many cases, there is no need to test all combinations of the parameters. The inherited c -level parameters can serve as a center point for their refinement only. For example, NUD suggests two-stage search strategy. In the first stage a wide range of parameters is considered. In the second stage, the best combination from the first stage is locally refined using a smaller search range. In MAF, we do not need to apply the first stage as we only refine the inherited c -level parameters. Other grid search methods can be adjusted in a similar way.

Selecting suitable performance measures for the best model: In MAF, a criterion for choosing the best model throughout the hierarchy is more influential than that at the finest level in non-MAF frameworks. Moreover, these criteria can be different at different levels. For example, when one focuses on highly imbalanced sets, a criteria such as the best G-mean could be more beneficial than the accuracy. We found that introducing 2-level criteria for imbalanced sets such as (a) choose the best G-mean, and (b) among the combinations with the best G-mean choose the best sensitivity, performs particularly good if applied at the coarse levels when the tie breaker may be often required.

3.4.8 Models at different scales

Throughout the hierarchy, we solve (W)SVM models at different scales of coarseness. Intuitively, the coarsening procedure gradually creates generalized (or summarized) descriptions of the finest level data which results in generalized coarse hyperplanes which can be used as final solutions. Indeed, the finest level, rich data can easily lead to overfitted models, a phenomenon frequently observed in practice [20]. In a multilevel framework, one can use models from multiple scales because the most correct validation is done against the fine level data in any case. Our experiments confirm that more than half of the best models are obtained from the coarse levels.

4 Computational Results

We compare our algorithms in terms of classification quality and computational performance to the state-of-the-art sequential SVM algorithms LibSVM, DC-SVM, and fast Ensemble SVM. The DC-SVM is a most recent, fast, hierarchical approach that outperforms other hierarchical methods which was the reason to choose it for comparison. The classification quality is evaluated using the following performance measures: sensitivity (SN), specificity (SP), G-mean, and accuracy (ACC), Precision (PPV), and F1, namely,

$$\text{SN} = \frac{TP}{TP + FN}, \quad \text{SP} = \frac{TN}{TN + FP}, \quad \text{G-mean} = \sqrt{\text{SP} \cdot \text{SN}},$$
$$\text{ACC} = \frac{TP + TN}{FP + TN + TP + FN}, \quad \text{Precision (PPV)} = \frac{TP}{TP + FP},$$

and

$$\text{F1-measure } (F_1) = \frac{2 \times TP}{2TP + FP + FN}$$

where TN , TP , FP , and FN correspond to the numbers of true negative, true positive, false positive, and false negative points.

In all experiments the data is normalized using z-score. Each experimental result in the following tables represents an average over 100 executions of the same type with different random seeds. The computational time reported in all experiments contains generating the k -NN graph. The computational time is reported in seconds unless it is explicitly mentioned otherwise.

In each class, a part of the data is assigned to be the test data using k -fold cross validation. We experimented with $k=5$ and 10 (no significant difference was observed). The experiments are repeated k times to cover all the data as test data. The data randomly shuffled for each k -fold cross validation. The presented results are the averages of performance measures for all k folds. Data points which are not in the test data are used as the training data in $\mathcal{J}^{+(-)}$. The test data is never used for any training or validation purposes. The Metis library [32] is used for graph partitioning during the refinement phase. We present the details about data sets in in Table 1.

The Forest data set [24] has 7 classes and different classes are reported in the literature (typically, not the difficult ones). Class 5 is used in our experiments as the most difficult and highly imbalanced. We report our results on other classes which are listed in Table 2 for convenient comparison with other methods.

Table 1: Benchmark data sets.

Dataset	r_{imb}	n_f	$ \mathcal{J} $	$ \mathbf{C}^+ $	$ \mathbf{C}^- $
Advertisement	0.86	1558	3279	459	2820
Buzz	0.80	77	140707	27775	112932
Clean (Musk)	0.85	166	6598	1017	5581
Cod-rna	0.67	8	59535	19845	39690
EEG Eye State	0.55	14	14980	6723	8257
Forest (Class 5)	0.98	54	581012	9493	571519
Hypothyroid	0.94	21	3919	240	3679
ISOLET	0.96	617	6238	240	5998
Letter	0.96	16	20000	734	19266
Nursery	0.67	8	12960	4320	8640
Protein homology	0.99	74	145751	1296	144455
Ringnorm	0.50	20	7400	3664	3736
Twonorm	0.50	20	7400	3703	3697

Table 2: The Forest data set classes with $n_f = 54$ and $|\mathcal{J}| = 581012$

Class No	r_{imb}	$ \mathbf{C}^+ $	$ \mathbf{C}^- $
Class 1	0.64	211840	369172
Class 2	0.51	283301	297711
Class 3	0.94	35754	545258
Class 4	1.00	2747	578265
Class 5	0.98	9493	571519
Class 6	0.97	17367	563645
Class 7	0.96	20510	560502

4.1 mlsvm-IIS results

The performance measures of single- (LibSVM) and multi-level (W)SVMs are computed and compared in Table 3. In our earlier work [58], it has been shown in that the multilevel (W)SVM produces similar results compared to the single-level (W)SVM, but it is much faster (see Table 4). All experiments on all data sets have been executed on a single machine Intel Core i7-4790, 3.60GHz, and 16 GB RAM. The framework ran in sequential mode with no parallelization using Ubuntu 14.04.5 LTS, Matlab 2012a, Metis 5.0.2, and FLANN 1.8.4.

4.2 mlsvm-AMG sparsity preserving coarsening

We have experimented with the light version of mlsvm-AMG in which instead of computing a linear combination of f -level points to get c -level points (see Eq. 5), we prolongate the seed to be a corresponding coarse point in attempt to preserve the sparsity of data points. In terms of quality of classifiers, the performance measures of this method are similar to that of mlsvm-IIS and in most cases (see Tables 5-6) are faster. However, for Buzz and Cod-rna datasets, although mlsvm-AMG performs faster, it results in a lower sensitivity and specificity (see Table 5) for SVM, and higher sensitivity and specificity for WSVM (see Table 5) compared to mlsvm-IIS. For Protein dataset, the sensitivity and specificity are improved compared to mlsvm-IIS (see Table 5).

Table 3: Quality comparison using performance measures for multi- and single-level of (W)SVM. Each cell contains an average over 100 executions including model selection for each of them. Column “Depth” shows the number of levels. The best results are highlighted in bold font.

	Dataset	Multilevel				Depth	Single-level			
		ACC	SN	SP	G-mean		ACC	SN	SP	G-mean
SVM	Advertisement	0.94	0.97	0.79	0.87	7	0.92	0.99	0.45	0.67
	Buzz	0.94	0.96	0.85	0.90	14	0.97	0.99	0.81	0.89
	Clean (Musk)	1.00	1.00	0.99	0.99	5	1.00	1.00	0.98	0.99
	Cod-rna	0.95	0.93	0.97	0.95	9	0.96	0.96	0.95	0.96
	EEG Eye State	0.83	0.82	0.88	0.85	6	0.88	0.90	0.86	0.88
	Forest (Class 5)	0.93	0.93	0.90	0.91	33	1.00	1.00	0.86	0.92
	Hypothyroid	0.98	0.98	0.74	0.85	4	0.99	1.00	0.71	0.83
	ISOLET	0.99	1.00	0.83	0.92	11	0.99	1.00	0.85	0.92
	Letter	0.98	0.99	0.95	0.97	8	1.00	1.00	0.97	0.98
	Nursery	1.00	0.99	0.98	0.99	10	1.00	1.00	1.00	1.00
	Protein homology	1.00	1.00	0.72	0.85	18	1.00	1.00	0.80	0.89
	Ringnorm	0.98	0.98	0.99	0.98	6	0.98	0.99	0.98	0.98
	Twonorm	0.97	0.98	0.97	0.97	6	0.98	0.98	0.99	0.98
	WSVM	Advertisement	0.94	0.96	0.80	0.88	7	0.92	0.99	0.45
Buzz		0.94	0.96	0.87	0.91	14	0.96	0.99	0.81	0.89
Clean (Musk)		1.00	1.00	0.99	0.99	5	1.00	1.00	0.98	0.99
Cod-rna		0.94	0.97	0.95	0.96	9	0.96	0.96	0.96	0.96
EEG Eye State		0.87	0.89	0.86	0.88	6	0.88	0.90	0.86	0.88
Forest (Class 5)		0.92	0.92	0.90	0.91	33	1.00	1.00	0.86	0.93
Hypothyroid		0.98	0.98	0.75	0.86	4	0.99	1.00	0.75	0.86
ISOLET		0.99	1.00	0.85	0.92	11	0.99	1.00	0.85	0.92
Letter		0.99	0.99	0.96	0.99	8	1.00	1.00	0.97	0.99
Nursery		1.00	0.99	0.98	0.99	10	1.00	1.00	1.00	1.00
Protein homology		1.00	1.00	0.87	0.92	18	1.00	1.00	0.80	0.89
Ringnorm		0.98	0.97	0.99	0.98	6	0.98	0.99	0.98	0.98
Twonorm		0.97	0.98	0.97	0.97	6	0.98	0.98	0.99	0.98

Table 4: Comparison of computational time for single- (LibSVM) and multilevel (mlsvm-IIS) solvers in seconds. Presented values include running time in seconds for both WSVM and SVM with model selection.

Dataset	Multilevel	Single-level
Advertisement	196	412
Buzz	2329	70452
Clean (Musk)	30	167
Cod-rna	172	1611
EEG Eye State	51	447
Forest (Class 5)	13785	352500
Hypothyroid	3	5
ISOLET	69	1367
Letter	45	333
Nursery	63	519
Protein homology	1564	73311
Ringnorm	4	42
Twonorm	4	45

Table 5: Performance measures of regular and weighted `mlsvm-AMG`. Column 'Depth' shows the number of levels in the multilevel hierarchy which is independent of SVM type.

Dataset	regular <code>mlsvm-AMG</code>				weighted <code>mlsvm-AMG</code>				Depth
	ACC	SN	SP	G-mean	ACC	SN	SP	G-mean	
Advertisement	0.95	0.99	0.64	0.86	0.95	0.99	0.64	0.86	2
Buzz	0.87	0.89	0.79	0.83	0.93	0.95	0.85	0.90	8
Clean	0.99	1.00	0.98	0.99	0.99	1.00	0.98	0.99	4
Cod-rna	0.86	0.85	0.88	0.87	0.89	0.89	0.90	0.90	6
EEG Eye State	0.87	0.88	0.85	0.86	0.87	0.88	0.85	0.86	4
Forest (Class 5)	0.97	0.98	0.79	0.88	0.96	0.97	0.82	0.89	9
ISOLET	0.99	1.00	0.83	0.91	0.99	1.00	0.83	0.91	3
Letter	0.99	0.99	0.95	0.97	0.99	0.99	0.93	0.96	5
Nursery	0.99	0.99	1.00	0.99	1.00	1.00	1.00	1.00	4
Protein homology	0.97	0.97	0.86	0.91	0.97	0.97	0.85	0.91	5
Ringnorm	0.98	0.98	0.98	0.98	0.98	0.99	0.98	0.98	3
Twonorm	0.98	0.97	0.98	0.98	0.98	0.97	0.98	0.98	3

Table 6: Comparison of computational time for single- (LibSVM) and multilevel (sparse `mlsvm-AMG`) solvers in seconds. Presented values include running time in seconds for both WSVM and SVM with model selection.

Dataset	Sparse <code>mlsvm-AMG</code>	Single-level
Advertisement	91	412
Buzz	957	70452
Clean	6	167
Cod-rna	92	1611
EEG Eye State	45	447
Forest (Class 5)	13328	352500
ISOLET	64	1367
Letter	18	333
Nursery	33	519
Protein	1597	73311
Ringnorm	5	42
Twonorm	4	45

We perform the sensitivity analysis of the order of interpolation denoted by r (see Eq. 4), the maximum number of fractions a point in F can be divided into, and compare the performance measures and computational time in Table 7. As r increases, the performance measures such as G-mean are improving until they do not stop changing for larger r . For example, for Buzz dataset, the G-mean is not changing for larger $r = 6$. The presented results are computed without advancements FF and FS (see Section 3.3). Using these techniques, we obtain G-mean 0.95 with $r = 1$ for Buzz data set. Higher interpolation orders increase the time but produce the same quality on that data set.

4.3 Full `mlsvm-AMG` coarsening

The best version of full `mlsvm-AMG` coarsening whose results are reported, chooses the best model from different scales (see Sec. 3.4.8). For this type of `mlsvm-AMG`, all experiments on all data sets

Table 7: Sensitivity analysis of interpolation order r in `mlsvm-AMG` for Buzz data set.

	Metric	$r = 1$	$r = 2$	$r = 3$	$r = 4$	$r = 5$	$r = 6$	$r = 10$
<code>mlsvm-AMG SVM</code>	G-mean	0.26	0.33	0.56	0.83	0.90	0.91	0.89
	SN	0.14	0.33	0.68	0.89	0.98	0.97	0.95
	SP	0.47	0.34	0.47	0.79	0.82	0.86	0.82
	ACC	0.21	0.33	0.64	0.87	0.95	0.95	0.94
<code>mlsvm-AMG WSVM</code>	G-mean	0.26	0.40	0.60	0.90	0.93	0.93	0.94
	SN	0.14	0.32	0.74	0.95	0.98	0.98	0.98
	SP	0.47	0.5	0.48	0.85	0.88	0.89	0.89
	ACC	0.21	0.35	0.69	0.93	0.96	0.97	0.97
	time(sec.)	389	541	659	957	1047	1116	1375

have been executed on a single machine with CPU Intel Xeon E5-2665 2.4 GHz and 64 GB RAM. The framework runs in sequential mode. The FLANN library is used to generate the approximated k -NN graph for $k = 10$. Once it is generated for the whole data set, its result is saved and reused. In all experiments all data points are randomly reordered as well as for each k -fold, the indices from the original test data are removed and reordered, so *no order in which points are entered into QP solver affects the solution*. Each experiment includes a full k -fold cross validation. The average performance measures over 5 experiments each of which includes 10-fold cross validation are demonstrated in Tables 9 and 10. The best model among all the levels for each fold of cross validation is selected using *validation* data (see Sec. 3.4.7). Using the best model, the performance measures over the test data are calculated and reported as the final performance for this specific fold of cross validation.

For the purpose of comparison, the results of previous work using validation techniques CS, CCKF without partitioning the training data during the refinement [61] are presented in Table 8.

Table 8: Performance measures and running time (in seconds) for weighted single level SVM (LibSVM), and weighted `mlsvm-AMG` on benchmark data sets in [41] *without partitioning*.

Dataset	Single level WSVM					<code>mlsvm-AMG</code>				
	ACC	SN	SP	G-mean	Time	ACC	SN	SP	G-mean	Time
Advertisement	0.92	0.99	0.45	0.67	231	0.83	0.92	0.81	0.86	213
Buzz	0.96	0.99	0.81	0.89	26026	0.88	0.97	0.86	0.91	233
Clean (Musk)	1.00	1.00	0.98	0.99	82	0.97	0.97	0.97	0.97	7
Cod-RNA	0.96	0.96	0.96	0.96	1857	0.94	0.97	0.92	0.95	102
Forest	1.00	1.00	0.86	0.92	353210	0.88	0.92	0.88	0.90	479
Hypothyroid	0.99	1.00	0.75	0.86	3	0.98	0.83	0.99	0.91	3
ISOLET	0.99	1.00	0.85	0.92	1367	0.99	0.89	1.00	0.94	66
Letter	1.00	1.00	0.97	0.99	139	0.98	1.00	0.97	0.99	12
Nursery	1.00	1.00	1.00	1.00	192	1.00	1.00	1.00	1.00	2
Ringnorm	0.98	0.99	0.98	0.98	26	0.98	0.98	0.98	0.98	2
Twonorm	0.98	0.98	0.99	0.98	28	0.98	0.98	0.97	0.98	1

The results using validation techniques FF, FS with partitioning the training data during the refinement phase are presented in Table 9, 10. We compare our performance and quality with those obtained by LibSVM, DC-SVM, and Ensemble SVM. All results are related to WSVM. The “Single level WSVM” column in Table 9 represents the weighted SVM results produced by LibSVM. The

LibSVM is the state-of-the-art SVM library which produces almost the best G-mean results over our experimental datasets except Advertisement, Buzz, and Forest. The DC-SVM [28] produces better G-mean on 4 datasets compare to LibSVM (see Table 9) but has lower G-mean on 4 other datasets. We choose DC-SVM not only because it has a hierarchical framework (with different principles of (un)coarsening) but also because it significantly outperforms other hierarchical techniques which are typically fast but not of high quality.

The mlsvm-AMG demonstrates significantly better computation time than DC-SVM on almost all datasets (see Table 9). Furthermore, mlsvm-AMG classification quality is significantly better on both Advertisement and Buzz datasets compared to LibSVM. In addition, the comparison between DC-SVM and mlsvm-AMG shows that the latter has higher G-mean for Advertisement, Buzz, Clean, Cod, Ringnorm, and Twonorm datasets. A better performance of DC-SVM is observed on Forest dataset if mlsvm-AMG is applying partitioning, i.e., when the number of support vectors is big. However, in another version of multilevel framework with validation techniques CS, CCKF without partitioning the training data during the refinement, the G-mean raises to 0.90 (see Table 8). In it interesting to note that the dimensionality of Advertisement dataset is the main source of complexity for the parameter fitting in both LibSVM and mlsvm-AMG. All versions of multilevel SVMs produce G-mean 0.90 for this dataset which is significantly higher than that of LibSVM which is 0.67. The results for this dataset are not significantly different for DC-SVM which is, however, 3 times slower than full mlsvm-AMG and 6 times slower than sparse mlsvm-AMG.

Table 9: Performance measures for single level WSVM (LibSVM), DC-SVM and mlsvm-AMG on benchmark data sets using partitioning and FF, FS validation techniques.

Datasets	Single level WSVM				DC-SVM				mlsvm-AMG			
	ACC	SN	SP	G-mean	ACC	SN	SP	G-mean	ACC	SN	SP	G-mean
Advertisement	0.92	0.99	0.45	0.67	0.95	0.83	0.97	0.90	0.95	0.85	0.96	0.91
Buzz	0.96	0.99	0.81	0.89	0.96	0.88	0.97	0.92	0.94	0.95	0.94	0.95
Clean (Musk)	1.00	1.00	0.98	0.99	0.96	0.91	0.97	0.94	0.99	0.99	0.99	0.99
Cod-RNA	0.96	0.96	0.96	0.96	0.93	0.93	0.94	0.93	0.93	0.97	0.91	0.94
Forest	1.00	1.00	0.86	0.92	1.00	0.88	1.00	0.94	0.77	0.96	0.80	0.88
Letter	1.00	1.00	0.97	0.99	1.00	1.00	1.00	1.00	0.98	0.99	0.98	0.99
Nursery	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Ringnorm	0.98	0.99	0.98	0.98	0.95	0.92	0.98	0.95	0.98	0.98	0.98	0.98
Twonorm	0.98	0.98	0.99	0.98	0.97	0.98	0.96	0.97	0.98	0.98	0.97	0.98

The computational time in seconds is demonstrated in Table 10. Our experiments exhibit significant performance improvement.

4.4 Large datasets

Large datasets SUSY and Higgs were downloaded from UCI repository [42]. The MNIST8M was downloaded from LibSVM dataset. Half of each class was randomly sampled to make classification more difficult. All the methods are benchmarked using Intel Xeon (E5-2680v3) with 128Gb memory.

The experiments with DC-SVM have not been finished after 3 full days of running and its performance is not presented because of unrealistic slowness of the method. Therefore, it is not comparable with mlsvm-AMG on large datasets. The LibSVM performs slower than DC-SVM on

Table 10: Computational time in seconds for single level WSVM (LibSVM), DC-SVM and mlsvm-AMG.

Dataset	Single level WSVM	DC-SVM	mlsvm-AMG
Advertisement	231	610	213
Buzz	26026	2524	31
Clean (Musk)	82	95	94
Cod-RNA	1857	420	13
Forest	353210	19970	948
Letter	139	38	30
Nursery	192	49	2
Ringnorm	26	38	2
Twonorm	28	30	1

these datasets and is also not presented. Although, fast linear SVM solvers are beyond the scope of this work, we compare the mlsvm-AMG with the LibLinear [21] that is significantly faster than both DC-SVM and LibSVM. We note that they also can be used as the refinement in multilevel frameworks. However, in practice, we do not observe a need for this because nonlinear SVM refinement is already fast enough in our multilevel framework.

The results for performance measures and computational time are presented in Tables 12, and 13. The mlsvm-AMG produces higher G-means on SUSY, HIGGS, and 8 (out of 10) of classes in the MNIST8M datasets. On classes 8, 5, and 9 of MNIST8M we have an improvement of 24%, 6% and 5%, respectively. On the average, the G-mean for all larger datasets are 5% higher for mlsvm-AMG compare to LibLinear. The mlsvm-AMG is faster than LibLinear on SUSY and HIGGS datasets and slower on MNIST8M dataset. However, this slowness is eliminated if linear SVM solver is used in the refinement. The results for seven classes of Forest dataset are presented in Table 14.

Finally, we present the computational time in terms of the amount of work per unit for all datasets in Table 15. In “ $\frac{\mu s}{point}$ ” and “ $\frac{\mu s}{value}$ ” columns, we present the computational time in microseconds per data point and one feature value in data point, respectively.

Table 11: Larger benchmark data sets.

Dataset	r_{imb}	n_f	$ \mathcal{J} $	$ \mathbf{C}^+ $	$ \mathbf{C}^- $
SUSY	0.54	18	5000000	2287827	2712173
MNIST8M (Class 0)	0.90	784	4050003	399803	3650200
MNIST8M (Class 1)	0.89	784	4050003	455085	3594918
MNIST8M (Class 2)	0.90	784	4050003	402165	3647838
MNIST8M (Class 3)	0.90	784	4050003	413843	3636160
MNIST8M (Class 4)	0.90	784	4050003	394335	3655668
MNIST8M (Class 5)	0.91	784	4050003	365918	3684085
MNIST8M (Class 7)	0.90	784	4050003	399465	3650538
MNIST8M (Class 6)	0.90	784	4050003	422888	3627115
MNIST8M (Class 8)	0.90	784	4050003	394943	3655060
MNIST8M (Class 9)	0.90	784	4050003	401558	3648445
HIGGS	0.53	28	11000000	5170877	5829123

Table 12: Performance measures for single level WSVM (LibLinear), DC-SVM/LibSVM and mlsvm-AMG on larger benchmark data sets using partitioning and FF, FS validation techniques.

Dataset	LibLinear				DC-SVM and LibSVM	mlsvm-AMG			
	ACC	SN	SP	G-mean		ACC	SN	SP	G-mean
SUSY	0.69	0.61	0.76	0.68		0.75	0.71	0.78	0.74
MNIST8M (Class 0)	0.98	0.90	0.99	0.95		0.94	0.93	0.94	0.95
MNIST8M (Class 1)	0.98	0.93	0.99	0.96		0.94	0.95	0.94	0.95
MNIST8M (Class 2)	0.97	0.77	0.99	0.87		0.91	0.87	0.91	0.89
MNIST8M (Class 3)	0.96	0.70	0.98	0.83		0.88	0.86	0.89	0.88
MNIST8M (Class 4)	0.97	0.81	0.99	0.90	Stopped or failed after 3 days without any result	0.84	0.92	0.84	0.88
MNIST8M (Class 5)	0.96	0.64	0.99	0.80		0.84	0.91	0.83	0.86
MNIST8M (Class 6)	0.98	0.86	0.99	0.92		0.95	0.98	0.95	0.96
MNIST8M (Class 7)	0.98	0.85	0.99	0.91		0.93	0.90	0.93	0.92
MNIST8M (Class 8)	0.92	0.36	0.98	0.60		0.82	0.87	0.81	0.84
MNIST8M (Class 9)	0.94	0.64	0.97	0.80		0.81	0.91	0.79	0.85
HIGGS	0.54	0.55	0.54	0.54		0.62	0.61	0.63	0.62

Table 13: Computational time in seconds for single level WSVM (LibLinear), DC-SVM/LibSVM and mlsvm-AMG on larger benchmark data sets

Dataset	LibLinear	DC-SVM and LibSVM	mlsvm-AMG
SUSY	1300		1116
MNIST8M (Class 0)	1876		21703
MNIST8M (Class 1)	859		14970
MNIST8M (Class 2)	1840		28469
MNIST8M (Class 3)	2362		22612
MNIST8M (Class 4)	1448	Stopped or failed after 3 days without any result	23490
MNIST8M (Class 5)	2360		30901
MNIST8M (Class 6)	1628		62336
MNIST8M (Class 7)	1747		17965
MNIST8M (Class 8)	2626		22439
MNIST8M (Class 9)	1650		25391
HIGGS	4406		3283

4.4.1 Results for k-distant disaggregation

The Forest, Clean, and Letter are the three data sets which demonstrate significant improvement on classification quality by adding the distant neighbors. The results for including the distant neighbors for the Letter data set experimenting with multiple coarse neighbor size reveal the largest improvement for $r = 1$ (see Figure 3).

Table 14: Performance measures and running time (in seconds) for all classes of Forest dataset using full mlsvm-AMG.

Dataset	ACC	SN	SP	G-mean	PPV	F_1	Time
Class 1	0.73	0.79	0.69	0.74	0.60	0.68	926
Class 2	0.70	0.78	0.62	0.70	0.67	0.72	215
Class 3	0.90	0.99	0.90	0.94	0.39	0.56	1496
Class 4	0.92	1.00	0.92	0.96	0.99	0.98	3231
Class 5	0.80	0.96	0.80	0.88	0.07	0.14	948
Class 6	0.86	0.95	0.95	0.90	0.17	0.28	2972
Class 7	0.91	0.87	0.91	0.89	0.28	0.42	2269

Table 15: Complexity Analysis

Dataset	$ \mathcal{J} $	n_f	$ \mathcal{J} \cdot n_f$	$\frac{\mu s}{point}$	$\frac{\mu s}{value}$
Nursery	13K	19	246.2K	232	12
Twonorm	7.4K	20	148K	405	20
Ringnorm	7.4K	20	148K	541	27
Letter	20K	16	320K	100	6
Cod-rna	59.5K	8	476.3K	100	13
Clean (Musk)	6.6K	166	1.1M	909	6
Advertisement	3.3K	1558	5.1M	31107	20
Buzz	140.7K	77	10.8M	1628	21
Forest	581K	54	31.4M	207	4
Susy	5M	18	90M	223	12
Higgs	11M	28	308M	298	11
mnist 4M	4.1M	784	3.2G	6673	9

4.4.2 Using partitioning in the refinement

When the training set becomes too big during the refinement (at any level), a partitioning is used to accelerate the performance. In Table 16, we compare the classification quality (G-mean), the size of training data, and the computational time. In columns “Partitioned” (“Full”), we show these three factors when (no) partitioning is applied. When no partitioning is used, we train the model with the whole training data at each level. The partitioning starts when the size of training data is 5000 points. Typically, at the very coarse levels the size of training data is small, so in the experiment demonstrated in Table 16, we show the numbers beginning level 5, the last level at which the partitioning was not applied. The results in this and many other similar experiments show significant improvement in computational time when the training data is partitioned with very minor loss in G-mean. The best level in the hierarchy is considered as the final level which is selected based on G-mean. Therefore, with no partitioning we obtain G-mean 0.79 and with partitioning it is 0.77 which are not significantly different results.

Figure 3: Effect of considering 1-distant disaggregation during the refinement phase on the G-mean for the Letter data set

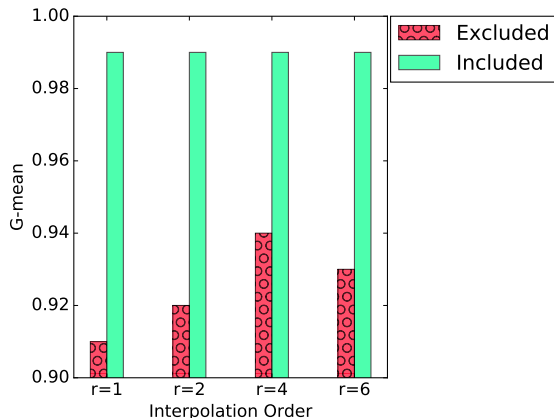


Table 16: The G-mean, training set size, and computational time are reported for levels 1-5 of Forest data set for Class 5. The partitioning is started with 5000 points.

Level	G-mean		Size of training set		Computational time	
	Full	Partitioned	Full	Partitioned	Full	Partitioned
5	0.69	0.69	4387	4387	373	373
4	0.68	0.72	18307	18673	1624	1262
3	0.79	0.77	47588	43977	6607	1528
2	0.79	0.72	95511	33763	17609	917
1	0.72	0.74	138018	24782	27033	576

4.5 Comparison with fast Ensemble SVM

A typical way to estimate the correctness of a multilevel solver is to compare its performance to those that use the local refinement techniques only. The EnsembleSVM [16] is a free software package containing efficient routines to perform ensemble learning with SVM models. The implementation exhibits very fast performance avoiding duplicate storage and evaluation of support vectors which are shared between constituent models. In fact, it is similar to our refinement and can potentially replace it in the multilevel framework. The comparison of our method with EnsembleSVM is presented in Table 17. While the running time is incomparable because of the obvious reasons (the complexity of EnsembleSVM is comparable to that of our last refinement only), the quality of our solver is significantly higher.

Table 17: Ensemble SVM on benchmark data sets

Dataset	Ensemble SVM				mlsvm-AMG			
	ACC	SN	SP	G-mean	ACC	SN	SP	G-mean
Advertisement	0.52	0.41	0.95	0.57	0.95	0.85	0.96	0.91
Buzz	0.65	0.36	0.99	0.59	0.94	0.95	0.94	0.95
Clean (Musk)	0.85	0.00	0.85	0.00	0.99	0.99	0.99	0.99
Cod-RNA	0.90	0.82	0.94	0.88	0.93	0.97	0.91	0.94
Forest	0.98	0.32	0.99	0.57	0.77	0.96	0.80	0.88
Letter	0.97	0.75	0.98	0.86	0.98	0.99	0.98	0.99
Nursery	0.68	1.00	0.68	0.82	1.00	1.00	1.00	1.00
Ringnorm	0.68	0.61	1.00	0.78	0.98	0.98	0.98	0.98
Twonorm	0.75	0.89	0.76	0.81	0.98	0.98	0.97	0.98

5 Conclusions

In this paper we introduced novel multilevel frameworks for nonlinear support vector machines. and discussed the details of several techniques for engineering multilevel frameworks that lead to a good trade-off between quality and running time. We ran a variety of experiments to compare several state-of-the-art SVM libraries and our frameworks on the classification quality and computation performance. The computation time of the proposed multilevel frameworks exhibits a significant improvement compared to the state-of-the-art SVM libraries with comparable or improved classification quality. For large data sets with more than 100,000 and up to millions of data points, we observed an improvement of computational time within an order of magnitude in comparison to DC-SVM and more two orders of magnitude in comparison to LibSVM. The improvement for larger datasets is even more significant. The code for mlsvm-AMG is available at <https://github.com/esadr/mlsvm>.

Acknowledgements

This material is based upon work supported by the National Science Foundation under Grants No. 1638321 and 1522751.

References

- [1] An, S., Liu, W., Venkatesh, S.: Fast cross-validation algorithms for least squares support vector machine and kernel ridge regression. *Pattern Recognition* **40**(8), 2154–2162 (2007)
- [2] Asharaf, S., Murty, M.N.: Scalable non-linear support vector machine using hierarchical clustering. In: *Pattern Recognition, 2006. ICPR 2006. 18th International Conference on*, vol. 1, pp. 908–911. IEEE (2006)
- [3] Balay, S., Abhyankar, S., Adams, M.F., Brown, J., Brune, P., Buschelman, K., Dalcin, L., Eijkhout, V., Gropp, W.D., Kaushik, D., Knepley, M.G., McInnes, L.C., Rupp, K., Smith, B.F., Zampini, S., Zhang, H.: PETSc users manual. Tech. Rep. ANL-95/11 - Revision 3.7, Argonne National Laboratory (2016). URL <http://www.mcs.anl.gov/petsc>

- [4] Bao, Y., Hu, Z., Xiong, T.: A pso and pattern search based memetic algorithm for svms parameters optimization. *Neurocomputing* **117**, 98–106 (2013)
- [5] Berry, M., Potok, T.E., Balaprakash, P., Hoffmann, H., Vatsavai, R., Prabhat: Machine Learning and Understanding for Intelligent Extreme Scale Scientific Computing and Discovery. Tech. Rep. 15-CS-1768, ASCR DOE Workshop Report (2015). URL <https://www.ornl.gov/machinelearning2015/>
- [6] Brandes, U., Delling, D., Gaertler, M., Gorke, R., Hoefer, M., Nikoloski, Z., Wagner, D.: On modularity clustering. *IEEE Transactions on Knowledge and Data Engineering* **20**(2), 172–188 (2008)
- [7] Brandt, A., Ron, D.: Chapter 1 : Multigrid solvers and multilevel optimization strategies. In: J. Cong, J.R. Shinnerl (eds.) *Multilevel Optimization and VLSICAD*. Kluwer (2003)
- [8] Brannick, J., Brezina, M., MacLachlan, S., Manteuffel, T., McCormick, S., Ruge, J.: An energy-based amg coarsening strategy. *Numerical linear algebra with applications* **13**(2-3), 133–148 (2006)
- [9] Buluç, A., Meyerhenke, H., Safro, I., Sanders, P., Schulz, C.: Recent advances in graph partitioning. *Algorithm Engineering: Selected Results and Surveys*. LNCS 9220, Springer-Verlag (2016)
- [10] Cawley, G.C., Talbot, N.L.: On over-fitting in model selection and subsequent selection bias in performance evaluation. *Journal of Machine Learning Research* **11**(Jul), 2079–2107 (2010)
- [11] Chang, C.C., Lin, C.J.: Libsvm: A library for support vector machines. *acm transactions on intelligent systems and technology*, 2: 27: 1–27: 27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm> (2011)
- [12] Chapelle, O., Vapnik, V., Bousquet, O., Mukherjee, S.: Choosing multiple parameters for support vector machines. *Machine learning* **46**(1), 131–159 (2002)
- [13] Chen, J., Safro, I.: Algebraic distance on graphs. *SIAM J. Scientific Computing* **33**(6), 3468–3490 (2011)
- [14] Cheong, S., Oh, S.H., Lee, S.Y.: Support vector machines with binary tree architecture for multi-class classification. *Neural Information Processing-Letters and Reviews* **2**(3), 47–51 (2004)
- [15] Chevalier, C., Safro, I.: Comparison of coarsening schemes for multilevel graph partitioning. *Learning and Intelligent Optimization* pp. 191–205 (2009)
- [16] Claesen, M., De Smet, F., Suykens, J.A., De Moor, B.: Ensemblesvm: a library for ensemble learning using support vector machines. *Journal of Machine Learning Research* **15**(1), 141–145 (2014)
- [17] Coussement, K., Van den Poel, D.: Churn prediction in subscription services: An application of support vector machines while comparing two parameter-selection techniques. *Expert systems with applications* **34**(1), 313–327 (2008)

- [18] Cui, L., Wang, C., Li, W., Tan, L., Peng, Y.: Multi-Modes Cascade SVMs: Fast Support Vector Machines in Distributed System, pp. 443–450. Springer Singapore, Singapore (2017). DOI 10.1007/978-981-10-4154-9_51. URL http://dx.doi.org/10.1007/978-981-10-4154-9_51
- [19] Dhillon, I., Guan, Y., Kulis, B.: A fast kernel-based multilevel algorithm for graph clustering. In: Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD’05), pp. 629–634. ACM Press (2005). DOI <http://doi.acm.org/10.1145/1081870.1081948>
- [20] Dietterich, T.: Overfitting and undercomputing in machine learning. *ACM computing surveys (CSUR)* **27**(3), 326–327 (1995)
- [21] Fan, R.E., Chang, K.W., Hsieh, C.J., Wang, X.R., Lin, C.J.: Liblinear: A library for large linear classification. *Journal of machine learning research* **9**(Aug), 1871–1874 (2008)
- [22] Fan, R.E., Chen, P.H., Lin, C.J.: Working set selection using second order information for training support vector machines. *The Journal of Machine Learning Research* **6**, 1889–1918 (2005)
- [23] Fang, H.r., Sakellaridi, S., Saad, Y.: Multilevel manifold learning with application to spectral clustering. In: Proceedings of the 19th ACM international conference on Information and knowledge management, pp. 419–428. ACM (2010)
- [24] Frank, A., Asuncion, A.: UCI machine learning repository. [<http://archive.ics.uci.edu/ml>]. irvine, ca: University of california, School of Information and Computer Science” **vol. 213** (2010)
- [25] Graf, H.P., Cosatto, E., Bottou, L., Dourdanovic, I., Vapnik, V.: Parallel support vector machines: The cascade SVM. In: Advances in neural information processing systems, pp. 521–528 (2004)
- [26] Hao, P.Y., Chiang, J.H., Tu, Y.K.: Hierarchically svm classification based on support vector clustering method and its application to document categorization. *Expert Systems with applications* **33**(3), 627–635 (2007)
- [27] Horng, S.J., Su, M.Y., Chen, Y.H., Kao, T.W., Chen, R.J., Lai, J.L., Perkasa, C.D.: A novel intrusion detection system based on hierarchical clustering and support vector machines. *Expert Systems with Applications* **38**(1), 306 – 313 (2011). DOI <http://doi.org/10.1016/j.eswa.2010.06.066>. URL <http://www.sciencedirect.com/science/article/pii/S0957417410005701>
- [28] Hsieh, C.J., Si, S., Dhillon, I.: A divide-and-conquer solver for kernel support vector machines. In: E.P. Xing, T. Jebara (eds.) Proceedings of the 31st International Conference on Machine Learning, *Proceedings of Machine Learning Research*, vol. 32, pp. 566–574. PMLR, Beijing, China (2014). URL <http://proceedings.mlr.press/v32/hsieha14.html>
- [29] Huang, C., Lee, Y., Lin, D., Huang, S.: Model selection for support vector machines via uniform design. *Computational Statistics & Data Analysis* **52**(1), 335–346 (2007)
- [30] Joachims, T.: Making large scale svm learning practical. Tech. rep., Universität Dortmund (1999)

- [31] Karypis, G., Han, E.H., Kumar, V.: Chameleon: Hierarchical clustering using dynamic modeling. *Computer* **32**(8), 68–75 (1999)
- [32] Karypis, G., Kumar, V.: MeTiS: A Software Package for Partitioning Unstructured Graphs, Partitioning Meshes, and Computing Fill-Reducing Orderings of Sparse Matrices, Version 4.0. University of Minnesota, Minneapolis, MN (1998)
- [33] Khan, L., Awad, M., Thuraisingham, B.: A new intrusion detection system using support vector machines and hierarchical clustering. *The VLDB Journal* **16**(4), 507–521 (2007). DOI 10.1007/s00778-006-0002-5. URL <http://dx.doi.org/10.1007/s00778-006-0002-5>
- [34] Khreich, W., Granger, E., Miri, A., Sabourin, R.: Iterative boolean combination of classifiers in the roc space: an application to anomaly detection with hmms. *Pattern Recognition* **43**(8), 2732–2752 (2010)
- [35] Kushnir, D., Galun, M., Brandt, A.: Fast multiscale clustering and manifold identification. *Pattern Recognition* **39**(10), 1876–1891 (2006). DOI 10.1016/j.patcog.2006.04.007. URL <http://dx.doi.org/10.1016/j.patcog.2006.04.007>
- [36] Kushnir, D., Galun, M., Brandt, A.: Fast multiscale clustering and manifold identification. *Pattern Recognition* **39**(10), 1876–1891 (2006)
- [37] Lee, H., Grosse, R., Ranganath, R., Ng, A.Y.: Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In: *Proceedings of the 26th annual international conference on machine learning*, pp. 609–616. ACM (2009)
- [38] Lessmann, S., Stahlbock, R., Crone, S.F.: Genetic algorithms for support vector machine model selection. In: *Neural Networks, 2006. IJCNN'06. International Joint Conference on*, pp. 3063–3069. IEEE (2006)
- [39] Leyffer, S., Safro, I.: Fast response to infection spread and cyber attacks on large-scale networks. *Journal of Complex Networks* **1**(2), 183–199 (2013)
- [40] Li, T., Liu, X., Dong, Q., Ma, W., Wang, K.: Hpsvm: Heterogeneous parallel svm with factorization based ipm algorithm on cpu-gpu cluster. In: *2016 24th Euromicro International Conference on Parallel, Distributed, and Network-Based Processing (PDP)*, pp. 74–81. IEEE (2016)
- [41] Lichman, M.: UCI machine learning repository (2013). URL <http://archive.ics.uci.edu/ml>
- [42] Lichman, M.: UCI machine learning repository (2013). URL <http://archive.ics.uci.edu/ml>
- [43] Lin, C.F., Wang, S.D.: Fuzzy support vector machines. *Neural Networks, IEEE Transactions on* **13**(2), 464–471 (2002)
- [44] Lin, S.W., Lee, Z.J., Chen, S.C., Tseng, T.Y.: Parameter determination of support vector machine and feature selection using simulated annealing approach. *Applied soft computing* **8**(4), 1505–1512 (2008)

- [45] López, V., del Río, S., Benítez, J.M., Herrera, F.: Cost-sensitive linguistic fuzzy rule based classification systems under the mapreduce framework for imbalanced big data. *Fuzzy Sets and Systems* **258**, 5–38 (2015)
- [46] Lovaglio, P., Vittadini, G.: Multilevel dimensionality-reduction methods. *Statistical Methods & Applications* **22**(2), 183–207 (2013). DOI 10.1007/s10260-012-0215-2. URL <http://dx.doi.org/10.1007/s10260-012-0215-2>
- [47] Luts, J., Ojeda, F., Van de Plas, R., De Moor, B., Van Huffel, S., Suykens, J.A.: A tutorial on support vector machine-based methods for classification problems in chemometrics. *Analytica Chimica Acta* **665**(2), 129–145 (2010)
- [48] Mazurowski, M.A., Habas, P.A., Zurada, J.M., Lo, J.Y., Baker, J.A., Tourassi, G.D.: Training neural network classifiers for medical decision making: The effects of imbalanced datasets on classification performance. *Neural networks* **21**(2), 427–436 (2008)
- [49] Mehrotra, S.: On the implementation of a primal-dual interior point method. *SIAM Journal on optimization* **2**(4), 575–601 (1992)
- [50] Muja, M., Lowe, D.G.: Fast approximate nearest neighbors with automatic algorithm configuration. In: *International Conference on Computer Vision Theory and Application VISS-APP'09*, pp. 331–340. INSTICC Press (2009)
- [51] Muja, M., Lowe, D.G.: Scalable nearest neighbor algorithms for high dimensional data. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **36**(11), 2227–2240 (2014)
- [52] Noack, A., Rotta, R.: Multi-level algorithms for modularity clustering. In: *Experimental Algorithms*, pp. 257–268. Springer (2009)
- [53] Noack, A., Rotta, R.: Multi-level algorithms for modularity clustering. In: J. Vahrenhold (ed.) *Experimental Algorithms, Lecture Notes in Computer Science*, vol. 5526, pp. 257–268. Springer Berlin Heidelberg (2009). DOI 10.1007/978-3-642-02011-7_24. URL http://dx.doi.org/10.1007/978-3-642-02011-7_24
- [54] Osuna, E., Freund, R., Girosi, F.: An improved training algorithm for support vector machines. In: *Neural Networks for Signal Processing [1997] VII. Proceedings of the 1997 IEEE Workshop*, pp. 276–285. IEEE (1997)
- [55] Platt, J.C.: Fast training of support vector machines using sequential minimal optimization. In: *Advances in kernel methods*, pp. 185–208. MIT press (1999)
- [56] Puget, R., Baskiotis, N.: Hierarchical label partitioning for large scale classification. In: *Data Science and Advanced Analytics (DSAA), 2015. 36678 2015. IEEE International Conference on*, pp. 1–10. IEEE (2015)
- [57] Razzaghi, T., Roderick, O., Safro, I., Marko, N.: Multilevel weighted support vector machine for classification on healthcare data with missing values. *PloS one* **11**(5), e0155119 (2016)
- [58] Razzaghi, T., Safro, I.: Scalable multilevel support vector machines. In: *International Conference on Computational Science (ICCS), Procedia Computer Science*, vol. 51, pp. 2683–2687. Elsevier (2015)

- [59] Ron, D., Safro, I., Brandt, A.: Relaxation-based coarsening and multiscale graph organization. *Multiscale Modeling & Simulation* **9**(1), 407–423 (2011)
- [60] Rotta, R., Noack, A.: Multilevel local search algorithms for modularity clustering. *Journal of Experimental Algorithmics (JEA)* **16**, 2–3 (2011)
- [61] Sadrfaridpour, E., Jeeredy, S., Kennedy, K., Luckow, A., Razzaghi, T., Safro, I.: Algebraic multigrid support vector machines. accepted in European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning (ESANN), arXiv preprint arXiv:1611.05487 (2017)
- [62] Safro, I., Ron, D., Brandt, A.: Multilevel algorithms for linear ordering problems. *ACM Journal of Experimental Algorithmics* **13** (2008)
- [63] Safro, I., Sanders, P., Schulz, C.: Advanced coarsening schemes for graph partitioning. *ACM Journal of Experimental Algorithmics (JEA)* **19**, 2–2 (2015)
- [64] Safro, I., Temkin, B.: Multiscale approach for the network compression-friendly ordering. *J. Discrete Algorithms* **9**(2), 190–202 (2011)
- [65] Schölkopf, B., Smola, A.J.: Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT press (2002)
- [66] Sharon, E., Galun, M., Sharon, D., Basri, R., Brandt, A.: Hierarchy and adaptivity in segmenting visual scenes. *Nature* **442**(7104), 810–813 (2006). DOI 10.1038/nature04977. URL <http://dx.doi.org/10.1038/nature04977>
- [67] Sun, Y., Kamel, M.S., Wong, A.K., Wang, Y.: Cost-sensitive boosting for classification of imbalanced data. *Pattern Recognition* **40**(12), 3358–3378 (2007)
- [68] Tavallaee, M., Stakhanova, N., Ghorbani, A.A.: Toward credible evaluation of anomaly-based intrusion-detection methods. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* **40**(5), 516–524 (2010)
- [69] Trottenberg, U., Schuller, A.: Multigrid. Academic Press, Orlando, FL (2001)
- [70] Wang, L.: Feature selection with kernel class separability. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **30**(9), 1534–1546 (2008)
- [71] Wu, Q., Zhou, D.X.: Svm soft margin classifiers: linear programming versus quadratic programming. *Neural computation* **17**(5), 1160–1187 (2005)
- [72] Yang, Z., Tang, W., Shintemirov, A., Wu, Q.: Association rule mining-based dissolved gas analysis for fault diagnosis of power transformers. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* **39**(6), 597–610 (2009)
- [73] You, Y., Demmel, J., Czechowski, K., Song, L., Vuduc, R.: Ca-svm: Communication-avoiding support vector machines on distributed systems. In: *Parallel and Distributed Processing Symposium (IPDPS), 2015 IEEE International*, pp. 847–859. IEEE (2015)

- [74] You, Y., Fu, H., Song, S.L., Randles, A., Kerbyson, D., Marquez, A., Yang, G., Hoisie, A.: Scaling support vector machines on modern hpc platforms. *Journal of Parallel and Distributed Computing* **76**, 16–31 (2015)
- [75] Yu, H., Yang, J., Han, J.: Classifying large data sets using svms with hierarchical clusters. In: *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 306–315. ACM (2003)
- [76] Zhang, X., Chen, X., He, Z.: An aco-based algorithm for parameter optimization of support vector machines. *Expert Systems with Applications* **37**(9), 6618–6628 (2010)
- [77] Zhou, L., Lai, K.K., Yu, L.: Credit scoring using support vector machines with direct search for parameters selection. *Soft Computing-A Fusion of Foundations, Methodologies and Applications* **13**(2), 149–155 (2009)
- [78] Zhu, K., Wang, H., Bai, H., Li, J., Qiu, Z., Cui, H., Chang, E.Y.: Parallelizing support vector machines on distributed computers. In: *Advances in Neural Information Processing Systems*, pp. 257–264 (2008)
- [79] Zhu, Z.A., Chen, W., Wang, G., Zhu, C., Chen, Z.: P-packsvm: Parallel primal gradient descent kernel svm. In: *Data Mining, 2009. ICDM'09. Ninth IEEE International Conference on*, pp. 677–686. IEEE (2009)
- [80] Zhu, Z.B., Song, Z.H.: Fault diagnosis based on imbalance modified kernel fisher discriminant analysis. *Chemical Engineering Research and Design* **88**(8), 936–951 (2010)