

Fast and Scalable Local Kernel Machines

Nicola Segata

Enrico Blanzieri

Department of Information Engineering and Computer Science

University of Trento

Trento, Italy

SEGATA@DISI.UNITN.IT

BLANZIER@DISI.UNITN.IT

Editor: Léon Bottou

Abstract

A computationally efficient approach to local learning with kernel methods is presented. The **Fast Local Kernel Support Vector Machine (FaLK-SVM)** trains a set of local SVMs on redundant neighbourhoods in the training set and an appropriate model for each query point is selected at testing time according to a proximity strategy. Supported by a recent result by Zakai and Ritov (2009) relating consistency and localizability, our approach achieves high classification accuracies by dividing the separation function in local optimisation problems that can be handled very efficiently from the computational viewpoint. The introduction of a fast local model selection further speeds-up the learning process. Learning and complexity bounds are derived for FaLK-SVM, and the empirical evaluation of the approach (with data sets up to 3 million points) showed that it is much faster and more accurate and scalable than state-of-the-art accurate and approximated SVM solvers at least for non high-dimensional data sets. More generally, we show that locality can be an important factor to sensibly speed-up learning approaches and kernel methods, differently from other recent techniques that tend to dismiss local information in order to improve scalability.

Keywords: locality, kernel methods, local learning algorithms, support vector machines, instance-based learning

1. Introduction

Efficiently processing large amount of data is one of the challenges of current research in kernel methods. Although most of the recently proposed techniques are based on different approaches, their common assumption is that scalability can be obtained by limiting or reducing the complexity of the decision function. In fact, very fast training algorithms have been developed for linear SVM (Keerthi and DeCoste, 2005; Collins et al., 2008; Chang et al., 2008; Bordes et al., 2009; Fan et al., 2008), and indeed they are effective when the linear separation is a good choice such as in high-dimensionality problems. Other approaches permit the non-linear feature space setting, but they limit the complexity by working with a reduced number of examples or a small set of support vectors (Lee and Mangasarian, 2001), using active and online example selection (Bordes et al., 2005; Bordes and Bottou, 2005) or bounding the number of basis functions (Keerthi et al., 2006; Joachims and Yu, 2009).

In the works referenced above, computational efficiency is sought bounding some aspects of the optimisation problem. The result is an *approximation* of the optimal separation and a *smoothing* of the decision function which is more influenced by the global distribution of the examples than by the local behaviour of the unknown target function in each specific sub-region. The emerging approach

is thus to trade locality for scalability permitting, with a potentially high level of under-fitting, to achieve a fast convergence to an approximated solution of the optimisation problem.

We show here that locality is not necessary related to computational inefficiency, but, instead, it can be the key factor to obtain very fast kernel methods without the need to smooth locally the global decision function. In our proposed approach, the model is formed by a set of accurate local models trained on fixed-cardinality sub-regions of the training set and the prediction module uses for each query point the more appropriate local model. In this setting, we are not approximating with some level of inaccuracy the original SVM optimisation problem, but we are separately considering different parts of the decision function with the potential advantage of better capturing the local separation. So, instead of locally under-fit the decision function by globally smoothing it like approximated SVM solvers do, we search for decision functions that are locally-calculated and they are very similar (or even better) in terms of accuracy to the global decision function in the proximity of each testing point. This approach is theoretically supported also by the recent result obtained by Zakai and Ritov (2009) that showed how, roughly speaking, “consistency implies local behaviour”.

In this work we present **Fast Local Kernel Support Vector Machine (FaLK-SVM)**, that pre-computes a set of local SVMs covering with adjustable redundancy the whole training set and uses for prediction a model which is the nearest (in terms of neighbourhood rank in feature space) to each testing point. FaLK-SVM is obtained introducing various strategies, detailed below, to speed-up the Local SVM approach (see Blanzieri and Melgani, 2006 and Section 3.3). Scalability is obtained approximating the Local SVM approach softening the assumption that the query point must be the central example of the neighbourhood on which the local SVM is trained; in this way we use the same local SVM model for more than one testing point and we can also pre-compute the local models during training. The locality of the approach is regulated by the neighbourhood size k and the method uses all the training points. Starting from the theory of local learning algorithms (Bottou and Vapnik, 1992; Vapnik and Bottou, 1993) we derive generalisation bounds for FaLK-SVM, and we analyse the computational complexity stating that, under reasonable assumptions, the training of our technique scales as $N \log N$ and the testing as $\log N$ where N is the training set size. We also introduce a procedure for local model selection in order to speed-up the selection of the parameters and better capturing local properties of the data. The empirical evaluation (with data sets with up to 3 million examples) shows that FaLK-SVM outperforms accurate and approximated SVM solvers both in term of generalisation accuracy and computational performances.

The effectiveness and efficiency of our approach is directly related to the role that locality plays in the learning problem. It is well known, for example, that for very high-dimensional problems such as text and document classification, the linear kernel performs better than non-linear kernels which are hard to tune and can be subject to the “curse of dimensionality” (Bengio et al., 2005). On the other hand, there are problems (Blackard and Dean, 1999; Uzilov et al., 2006) which inherently require non-linear approaches to be tackled. This is due to the combination of an intrinsic dimensionality which is low with respect to the training set size and of a decision function which is not simple to learn. In general, locality plays a more important role as the number of training examples increases because the ratio between training set cardinality and the dimensionality is more favourable and the local characteristics are more evident. Other signals for the need of a non-linear kernel are the detection of uneven distributions in the data sets (typical of real-world problems), the monotonic increasing of accuracy with respect to training size also for already large amount of data and the inclusion of a high fraction of training examples in the support vector set. A representative

of this class of problems is the Forest CoverType data set (Blackard and Dean, 1999) which is a large real data set (more than half a million examples) with bounded dimensionality (54 features) that needs as many examples as possible to increase accuracy. We already showed in a very preliminary study (Segata and Blanzieri, 2009c) that our approach on this data set is more accurate than SVM and much faster than both accurate and approximated SVM solvers.

The present contribution can be seen from multiple viewpoints. (i) FaLK-SVM modifies the Local SVM approach (Blanzieri and Melgani, 2006; Zhang et al., 2006) that showed excellent classification performances but had dramatic computational problems, leading to a scalable Local SVM classifier asymptotically much faster than SVM. (ii) The approach is also an enhancement of the local learning algorithms because the learning process is not delayed until the prediction phase (*lazy learning*) but the construction of the local models occurs during training (*eager learning*). (iii) From a practical viewpoint, FaLK-SVM is a novel kernel method which outperforms accurate and approximated SVM solvers for non high-dimensional data sets. (iv) For complex classification problems that require a high fraction of support vectors (SVs), we exploit locality to avoid the need of bounding the number of total SVs as existing approximated SVM solvers do for computational reasons. (v) More generally, our approach can also be seen as a framework for localising and make scalable any kernel method, classifier and regressor and in general every data analysis that can be applied on sub-regions of the entire data set. The proposed FaLK-SVM classifier and related tools are freely available with source code as part of the Fast Local Kernel Machine Library (Segata, 2009, FaLKM-lib).

In the next Section we analyse the work on local learning algorithms, Local SVM and fast large margin classifiers that are all related with our work. Section 3 formally introduces some machine learning tools that we need in order to introduce FaLK-SVM in Section 4 and analyse its learning bounds, complexity bounds, implementation, local model selection procedure and intuitive interpretation. Section 5 details the empirical evaluation with respect to accurate and approximated approaches.

2. Related Work

Locality is often a crucial component of machine learning systems, although we are not aware of approaches exploiting locality for improving the computational performances. We review in this section those areas that are more related with our approach: local learning algorithms, local support vector machines, approximated and scalable SVM solvers.

2.1 Local Learning Algorithms

Local learning algorithms (LLAs) are a class of learning approaches introduced by Bottou and Vapnik (1992). Instead of estimating a decision function which is optimal (with respect to some criteria) for all possible unseen testing examples, the idea underlying LLAs consists in estimating the optimal decision function for each single testing point. The value of the function is estimated in a small sub-region of the input space around the query point. For a local learning algorithm, the points in the proximity of the query point have an higher influence in the training of the local model. The approach is particularly effective for uneven distributed data sets, that is, data sets presenting regions in which the examples have different spatial resolutions. In fact, with LLAs, the characteristics of the learning process can be locally adjusted. A proper choice of the locality parameter can reduce the generalisation error with respect to a global classifier as formalised by the

Local Risk Minimization principle (Vapnik and Bottou, 1993; Vapnik, 2000). Notice that there are various ways of specifying the degree of locality for LLAs as discussed for instance by Atkeson et al. (1997). Examples of LLAs are the well-known k -Nearest Neighbours (kNN) classifier, the Radial Basis Function networks (Broomhead and Lowe, 1988), and the Local SVM classifier (Blanzieri and Melgani, 2006; Zhang et al., 2006) described in Section 2.2.

Despite their theoretical and practical appeal, LLAs seem not to have been studied in depth in the last few years. This is probably due to the fact that LLAs, as formulated by Bottou and Vapnik (1992), fall in the class of *lazy learning* (or *memory-based learning*) that have great overhead on the testing phase, as opposed to *eager learning* in which the function estimation is performed during training increasing the computational performances of the testing phase.

2.2 Local Support Vector Machines

The main idea of Local SVM, described in details in Section 3.3, is to build at prediction time an example-specific maximal marginal hyperplane based on the set of k -neighbours.

Local SVM is a LLA and was independently proposed by Blanzieri and Melgani (2006, 2008) and by Zhang et al. (2006) and applied respectively to remote sensing and visual recognition tasks. Other successful applications of the approach are detailed by Segata and Blanzieri (2009a) for general real data sets, by Blanzieri and Bryl (2007) for spam filtering and by Segata, Blanzieri, Delany, and Cunningham (2009b) for noise reduction. Similar approaches have been presented by Yang and Kecman (2008) and applied in the medical domain (Yang and Kecman, 2009) and for face recognition problems (Yang and Kecman, 2010).

However, Local SVM suffers from the high computational cost of the testing phase that comprises, for each example, (i) the selection of the k nearest neighbours and (ii) the computation of the maximal separating hyperplane on the k examples. An attempt to computationally improve the Local SVM approach of Zhang et al. (2006) has been proposed by Cheng et al. (2007) where the idea is to train multiple SVMs on clusters found by a variant of k -means, called MagKmeans, that introduces in the clustering criterion the requirement that the clusters cannot have unbalanced class cardinalities. However the method does not follow directly the idea of Local SVM, the main difference being that it can build only local linear models and the size of the clusters is not fixed (MagKmeans does not have constraints on the cardinalities and the balancing requirement can cause the detection of clusters with high cardinalities). The achieved computational performances are better than their formulation of Local SVM, but worse than global SVM.

2.3 Fast Large Margin Classifiers

The need for fast and scalable kernel-based classifiers led to the development of several methods in the last few years, although considerable attention seems to have been focused especially on linear SVM classifiers. Below, we initially consider the works applicable also to non-linear kernels, successively we review the works on the linear case.

One of the first large-scale maximal margin learning that can use non-linear kernel functions is represented by Core Vector Machines (Tsang et al., 2005, CVM); reformulating the SVM approach as a minimum enclosing ball problem, the authors proved that it is possible to obtain approximated optimal solution in competitive training times by using the core sets. Good results have been achieved using non-linear kernels although it has been pointed out that the choice of the stopping criteria is crucial for the trade-off between computational efficiency and generalisation accuracy.

Ball Vector Machines (Tsang et al., 2007, BVM) are a modification of CVM in which the minimality of the enclosing balls is not required, because the radius of the ball is fixed. The resulting classifier improves the computational performances. Another approach based on an online setting of the SVM optimisation problem has been proposed by Bordes et al. (2005, LASVM) and by Bordes and Bottou (2005) and it is an algorithm that converges to the SVM solution. It has been shown that competitive accuracies can be achieved also after a single pass over the training set. The approach can be seen as a SVM solver that includes a support vector removal step. In addition, several strategies for active training-points selection can further improve computational and generalisation performances. Formulating the optimisation problem in the primal, Keerthi et al. (2006, SpSVM) proposed a method that bounds the number of basis functions considered and thus the computational complexity. Increasing the cardinality of the basis function set allows the method to converge to the SVM solution. A greedy strategy guides the choice of the basis functions to be included in the working set. Collobert et al. (2006, USVM) showed that softening the convex setting of maximal margin classifiers using a non-convex loss function can bring computational advantage over the corresponding standard convex problem. The non-convex problem is solved using the *concave-convex procedure* (Yuille and Rangarajan, 2003). Recently, the Cutting-Plane Subspace Pursuit (Joachims and Yu, 2009, CPSP) based on cutting-plane training (Joachims et al., 2009) has been proposed; it permits to learn maximal-margin decision functions in the feature space using arbitrary basis vectors instead of the support vectors only. This can result in sparser solutions increasing the testing and training computational performances especially for high-dimensional data sets. Although not always considered a method for large-scale learning, LibSVM (Chang and Lin, 2001) demonstrated to be competitive with approximated approaches from the computational viewpoint. LibSVM is a SVM solver implementing a SMO-type decomposition method proposed by Fan et al. (2005) integrating it with caching and shrinking (Joachims, 1999).

Large margin classifiers can also achieve scalability using subsampling-based approaches that train the model on a relatively small subset of the whole training set. However, the accuracy of SVM with subsampling can decrease due to the loss of information contained in the discarded training points. The decreasing of accuracy with respect to SVM without subsampling is more dramatic when a complex decision function is needed. In these cases the accuracy problems can be mitigated or reduced by developing an ensemble of classifiers. Bootstrap aggregating (bagging) by Breiman (1996) is an effective strategy to perform accurate classification using an ensemble of classifiers trained on subsets of the training set (using uniform sampling with replacement) that can also overcome the accuracies of SVM. Bagging with SVM can thus be used for obtaining scalability as long as the advantage of training smaller SVM models on subsets of the training set (that can scale cubically) overcome the disadvantage of training multiple SVMs.

Recently a lot of work has been performed in order to develop very fast and scalable solvers applicable to *linear* SVM only. Keerthi and DeCoste (2005) modified the Finite Newton method of Mangasarian (2002) introducing robust conjugate gradient techniques and other heuristics. Joachims (2006) developed an alternative formulation of the SVM optimisation problem exploiting a different form of sparsity. Lin et al. (2007) used logistic regression with Trust Region Newton Methods. Variants of coordinate descent methods for linear SVM are developed by Chang et al. (2008) in the primal and by Hsieh et al. (2008) in the dual. A different gradient approach was developed by Smola et al. (2008). Other approaches are based on Stochastic Gradient Descent (SGD) like those developed by Shalev-Shwartz et al. (2007) and by Bordes et al. (2009) which work in the primal, whereas Collins et al. (2008) apply SGD in the dual. Although SGD methods can be theoretically

used for non-linear SVM the performances are analysed for the linear case only. LIBLINEAR (Fan et al., 2008) is a fast software package implementing some of the cited works. The common idea of all the proposed methods is that the advantage of having a method that uses a huge number of training points overcomes the disadvantage of approximating the decision function with a linear model. This is effective, as explicitly noticed in almost all the cited works, when the dimensionality is very large and thus the problem is very sparse. This is, for example, the typical situation of text document classification. However, when the needed decision function is highly non-linear and the intrinsic dimensionality of the space is relatively small, the linear SVM approach cannot compete with SVM using non-linear kernels in terms of generalisation accuracy. Apart from the generalisation ability also the computational performances can be compromised in these cases, because the algorithm cannot find a good decision function and so convergence problems can occur.

3. Preliminaries

In order to introduce our approach, we need to analyse the formulation of kNN, SVM, kNNSVM and cover trees.

Here and in the following of the paper, we consider a binary class classification with examples $(\mathbf{x}_i, y_i) \in \mathcal{H} \times \{-1, +1\}$ for $i = 1, \dots, N$ and $\mathcal{X} = \{\mathbf{x}_i \mid i = 1, \dots, N\}$, where \mathcal{H} is an Hilbert space with inner product $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\|$. Extensions to multi-class problems will be explicitly discussed.

3.1 The k Nearest Neighbour Algorithm

Given an example $\mathbf{x}' \in \mathcal{H}$, it is possible to order an entire set of points \mathcal{X} with respect to \mathbf{x}' . This corresponds to define a function $r_{\mathbf{x}'} : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ that recursively reorders the indexes of the N points in \mathcal{X} :

$$\begin{cases} r_{\mathbf{x}'}(1) = \operatorname{argmin}_{i=1, \dots, N} \|\mathbf{x}_i - \mathbf{x}'\| \\ r_{\mathbf{x}'}(j) = \operatorname{argmin}_{i=1, \dots, N} \|\mathbf{x}_i - \mathbf{x}'\| \quad i \neq r_{\mathbf{x}'}(1), \dots, r_{\mathbf{x}'}(j-1) \quad \text{for } j = 2, \dots, N. \end{cases}$$

In this way, $\mathbf{x}_{r_{\mathbf{x}'}(j)}$ is the example in the j -th position in terms of distance from \mathbf{x}' , namely the j -th nearest neighbour, $\|\mathbf{x}_{r_{\mathbf{x}'}(j)} - \mathbf{x}'\|$ is its distance from \mathbf{x}' and $y_{r_{\mathbf{x}'}(j)}$ is its class. In other terms:

$$j < k \Rightarrow \|\mathbf{x}_{r_{\mathbf{x}'}(j)} - \mathbf{x}'\| \leq \|\mathbf{x}_{r_{\mathbf{x}'}(k)} - \mathbf{x}'\|.$$

Given the above definition, the majority decision rule of kNN for binary classification problems is defined by

$$\text{kNN}(\mathbf{x}) = \operatorname{sign} \left(\sum_{i=1}^k y_{r_{\mathbf{x}}(i)} \right).$$

For problems with more than two classes, the decision rule of kNN is the usual majority rule, namely the method selects the class with the highest number of representatives in the k -neighbourhood instead of taking the sign of the summation.

3.2 Support Vector Machines

SVMs (Cortes and Vapnik, 1995) are classifiers with sound foundations in statistical learning theory (Vapnik, 2000). The decision rule is

$$\text{SVM}(\mathbf{x}) = \text{sign}(\langle w, \Phi(\mathbf{x}) \rangle_{\mathcal{F}} + b)$$

where $\Phi(\mathbf{x}) : \mathcal{H} \rightarrow \mathcal{F}$ is a mapping in a transformed Hilbert feature space, called \mathcal{F} , with inner product $\langle \cdot, \cdot \rangle_{\mathcal{F}}$. The parameters $w \in \mathcal{F}$ and $b \in \mathbb{R}$ are such that they minimise an upper bound on the expected risk while minimising the empirical risk. The minimisation of the complexity term is achieved by the minimisation of the quantity $\frac{1}{2} \cdot \|w\|^2$, which is equivalent to the maximisation of the margin between the classes. In the optimisation problem, the violation of the margin is prevented by the following set of constraints:

$$y_i (\langle w, \Phi(\mathbf{x}_i) \rangle_{\mathcal{F}} + b) \geq 1. \tag{1}$$

If a linear separation cannot be found in the input or feature space, the soft-margin variant of SVM permits the violation of the margin and the presence of misclassified training examples. This is possible introducing slack variables ξ_i (the empirical risk):

$$y_i (\langle w, \Phi(\mathbf{x}_i) \rangle_{\mathcal{F}} + b) \geq 1 - \xi_i \quad \xi_i \geq 0, i = 1, \dots, N. \tag{2}$$

For soft-margin SVM the optimisation problem with linear penalisation of ξ_i (L1-norm), becomes the minimisation of $\frac{1}{2} \cdot \|w\|^2 + C \sum_i \xi_i$ subject to (2). Reformulating such an optimisation problem with Lagrange multipliers α_i ($i = 1, \dots, N$), and introducing a positive definite kernel (PD) function¹ $K(\cdot, \cdot)$ that substitutes the scalar product in the feature space $\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}) \rangle_{\mathcal{F}}$ the decision rule can be expressed as:

$$\text{SVM}(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right).$$

Throughout this work, SVM denotes the soft-margin SVM.

The kernel trick avoids the explicit definition of the feature space \mathcal{F} and of the mapping function Φ (Schölkopf and Smola, 2002). Popular kernels are the linear kernel, the radial basis function kernel, and the homogeneous and inhomogeneous polynomial kernels. Their definitions are:

$$\begin{aligned} K^{lin}(\mathbf{x}, \mathbf{x}') &= \langle \mathbf{x}, \mathbf{x}' \rangle & K^{rbf}(\mathbf{x}, \mathbf{x}') &= \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\sigma} \right), \\ K^{hpol}(\mathbf{x}, \mathbf{x}') &= \langle \mathbf{x}, \mathbf{x}' \rangle^d & K^{ipol}(\mathbf{x}, \mathbf{x}') &= (\langle \mathbf{x}, \mathbf{x}' \rangle + 1)^d. \end{aligned}$$

The maximal separating hyperplane defined by SVM has been shown to have important generalisation properties and nice bounds on the VC dimension (Vapnik, 2000).

Multiple methods has been proposed in order to apply the maximal margin principle of SVM on multiple class problems. The more popular are the one-against-all method (Bottou et al., 1994) which builds a number of binary decision functions equal to the number of classes N_{cl} , the one-against-one method (Knerr et al., 1990; Kressel, 1999) which builds $N_{cl} \cdot (N_{cl} - 1) / 2$ binary decision functions using voting in the prediction phase, and the Directed Acyclic Graph SVM (Platt et al.,

1. For convention we refer to kernel functions with the capital letter K and to the number of nearest neighbours with the lower-case letter k .

2000, DAGSVM) which is a modification of the one-against-all method. Other general strategies for reducing the multi-class classification setting to a binary classification problem have been analysed and developed by Allwein et al. (2000). The study carried on by Hsu and Lin (2002) shows that, for SVM, the more effective strategies are the one-against-one and DAGSVM approaches.

3.3 Local SVM: The k NNSVM Classifier

We already introduced the idea of Local SVM in Section 2.2, here we detail k NNSVM which is the formulation of Local SVM proposed by Blanzieri and Melgani (2006, 2008). k NNSVM can be seen as a modification of the SVM approach in order to obtain a LLA able to locally adjust the capacity of the training systems.

In order to classify a given example $\mathbf{x}' \in \mathcal{H}$, we need first to retrieve its k -neighbourhood in the transformed feature space \mathcal{F} and, then, to search for an optimal separating hyperplane only over this k -neighbourhood. In practice, this means that an SVM is built over the neighbourhood in \mathcal{F} of each test example \mathbf{x}' . Accordingly, the constraints in (1) become:

$$y_{r_{\mathbf{x}'}(i)} (w \cdot \Phi(\mathbf{x}_{r_{\mathbf{x}'}(i)}) + b) \geq 1 - \xi_{r_{\mathbf{x}'}(i)}, \text{ with } i = 1, \dots, k$$

where $r_{\mathbf{x}'} : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ is a function that reorders the indexes of the training examples defined as:

$$\begin{cases} r_{\mathbf{x}'}(1) = \operatorname{argmin}_{i=1, \dots, N} \|\Phi(\mathbf{x}_i) - \Phi(\mathbf{x}')\|_{\mathcal{F}}^2 \\ r_{\mathbf{x}'}(j) = \operatorname{argmin}_{i=1, \dots, N} \|\Phi(\mathbf{x}_i) - \Phi(\mathbf{x}')\|_{\mathcal{F}}^2 \quad i \neq r_{\mathbf{x}'}(1), \dots, r_{\mathbf{x}'}(j-1) \text{ for } j = 2, \dots, N. \end{cases} \quad (3)$$

In this way, $\mathbf{x}_{r_{\mathbf{x}'}(j)}$ is the example in the j -th position in terms of distance from \mathbf{x}' and thus $j < k \Rightarrow \|\Phi(\mathbf{x}_{r_{\mathbf{x}'}(j)}) - \Phi(\mathbf{x}')\|_{\mathcal{F}} \leq \|\Phi(\mathbf{x}_{r_{\mathbf{x}'}(k)}) - \Phi(\mathbf{x}')\|_{\mathcal{F}}$ because of the monotonicity of the quadratic operator. The computation is expressed in terms of kernels as:

$$\begin{aligned} \|\Phi(\mathbf{x}) - \Phi(\mathbf{x}')\|_{\mathcal{F}}^2 &= \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}) \rangle_{\mathcal{F}} + \langle \Phi(\mathbf{x}'), \Phi(\mathbf{x}') \rangle_{\mathcal{F}} - 2 \cdot \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathcal{F}} = \\ &= K(\mathbf{x}, \mathbf{x}) + K(\mathbf{x}', \mathbf{x}') - 2 \cdot K(\mathbf{x}, \mathbf{x}'). \end{aligned}$$

If the kernel is the RBF kernel or any polynomial kernels with degree 1, the ordering function is equivalent to the one defined by the Euclidean metric. In general, for some non-linear kernels (other than the RBF kernel) the ordering function can be quite different to that produced using the Euclidean metric.

The decision rule associated with the method for an example \mathbf{x} is:

$$k\text{NNSVM}(\mathbf{x}) = \operatorname{sign} \left(\sum_{i=1}^k \alpha_{r_{\mathbf{x}}(i)} y_{r_{\mathbf{x}}(i)} K(\mathbf{x}_{r_{\mathbf{x}}(i)}, \mathbf{x}) + b \right).$$

For $k = N$, the k NNSVM method is the usual SVM whereas, for $k = 2$, the method implemented with the linear or Gaussian radial basis function kernel corresponds to the standard 1-NN classifier. Notice that in situations where the neighbourhood contains only one class the local SVM does not find any separation and so considers all the neighbourhood to belong to the predominant class similarly to the behaviour of the majority rule. Considering k NNSVM as a local SVM classifier

built in the feature space, the method has been shown to have a potentially favourable bound on the expectation of the probability of test error with respect to SVM (Blanzieri and Melgani, 2008).

The generalisation of k NNSVM for multi-class classification can occur locally, that is solving the local multi-class SVM problem, or globally, that is applying the binary k NNSVM classifier on multiple global binary problems. In Segata and Blanzieri (2009a) the adopted strategy for multi-class classification with k NNSVM is the one-against-one strategy applied on the local problems. The choice of the one-against-one approach gave good results in comparison with the same strategy on SVM, but no specific empirical studies have been performed yet to identify the most appropriate strategy for multi-class classification with Local SVM.

3.4 Cover Trees

A cover tree is a data structure introduced by Beygelzimer et al. (2006) for performing exact nearest-neighbour operations in a fast and efficient way. Cover trees can be applied in general metric spaces without any other assumption on their structure and thus also in Hilbert spaces calculating the distances by means of kernel functions using the kernel trick.

In more detail, a cover tree can be viewed as a sub-graph of a navigating net (Krauthgamer and Lee, 2004) and it is a levelled tree in which each level (indexed by a decreasing integer i) is a cover (i.e., is representative) for the level beneath it. Every node of a cover tree T is associated with a point of a data set S . Denoting with C_i the set of points associated with nodes in T at level i , with $b > 1$ a constant, and with $dist(\cdot, \cdot)$ the distance function defining the metric of the space, the invariants of a cover tree are:

Nesting $C_i \subset C_{i-1}$

Covering tree For every $\mathbf{p} \in C_{i-1}$ there exists a $\mathbf{q} \in C_i$ such that $dist(\mathbf{p}, \mathbf{q}) < b^i$ and the node in level i associated with \mathbf{q} is a parent of the node in level $i - 1$ associated with \mathbf{p} .

Separation For all distinct $\mathbf{p}, \mathbf{q} \in C_i$, $dist(\mathbf{p}, \mathbf{q}) > b^i$.

Intuitively, the nesting invariant means that once a point appears in a level, it is present for every lower level. The covering tree invariant implies that every node has a parent in a higher level such that the distance between the respective points is less than b^i , while separation invariant assures that the distance between every pair of points associated to the nodes of a level i is higher than b^i . In addition, the root of the tree (called C_∞ and containing only one example) is a randomly chosen example.

Cover trees have state-of-the-art performance for exact nearest neighbour operations for general metrics in low-dimensional spaces both in terms of computational complexity and space requirements. As theoretically proved by Beygelzimer et al. (2006), the space required by the cover tree data-structure is linear in the data set size ($O(n)$), the computational time of single point insertions, deletions and exact nearest neighbour queries is logarithmic ($O(\log n)$) while the cover tree can be built in $O(n \log n)$.

4. FaLK-SVM: A Fast and Scalable Local Kernel Machine

In this section we introduce our novel technique. Initially we detail the way to pre-compute the local models during training (Section 4.1) and the strategies to reduce the number of local models

(Section 4.2). We then describe the prediction mechanism in Section 4.2.2 and our approach for fast local model selection in Section 4.3. Successively, we derive learning bounds for the approach in Section 4.4 before discussing the computational complexity in Section 4.5 and some details about the implementation (Section 4.6).

4.1 Pre-computing the Local Models during Training Phase

For the local approach we are proposing here, we need to generalise the decision rule of k NNSVM to the case in which the local model is trained on the k -neighbourhood of a point distinct, in the general case, from the query point. A modified decision function for a query point $\mathbf{q} \in \mathcal{H}$ and another (possibly different) point $\mathbf{t} \in \mathcal{H}$ is:

$$k\text{NNSVM}_{\mathbf{t}}(\mathbf{q}) = \text{sign} \left(\sum_{i=1}^k \alpha_{r_{\mathbf{t}}(i)} y_{r_{\mathbf{t}}(i)} K(\mathbf{x}_{r_{\mathbf{t}}(i)}, \mathbf{q}) + b \right) \quad (4)$$

where $r_{\mathbf{t}}(i)$ is the k NNSVM ordering function (see above Section 3.3) and $\alpha_{r_{\mathbf{t}}(i)}$ and b come from the training of an SVM on the k -neighbourhood of \mathbf{t} in the feature space. In the following we will refer to $k\text{NNSVM}_{\mathbf{t}}(\mathbf{q})$ as being centred on \mathbf{t} , to \mathbf{t} as the centre of the model, and, if $\mathbf{t} \in \mathcal{X}$, to $V_{\mathbf{t}}$ as the Voronoi cell induced by \mathbf{t} in \mathcal{X} , formally:

$$V_{\mathbf{t}} = \{\mathbf{p} \in \mathcal{H} \text{ s.t. } \|\mathbf{p} - \mathbf{t}\| \leq \|\mathbf{p} - \mathbf{x}\|, \forall \mathbf{x} \in \mathcal{X} \text{ with } \mathbf{x} \neq \mathbf{t}\}.$$

The original decision function of k NNSVM corresponds to the case in which $\mathbf{t} = \mathbf{q}$, and thus $k\text{NNSVM}_{\mathbf{q}}(\mathbf{q}) = k\text{NNSVM}(\mathbf{q})$.

k NNSVM requires that the training of an SVM on the k -neighbourhood of the query point must be performed in the prediction step. This approach is computationally feasible only for problems with few points to test which is a condition that rarely holds in real-world classification problems. In general, we need to speed-up the prediction phase. The first modification of k NNSVM consists in predicting the label of a test point \mathbf{q} using the local SVM model built on the k -neighbourhood of its nearest neighbour in \mathcal{X} . Formally, this can be written as:

$$k\text{NNSVM}_{\mathbf{t}}(\mathbf{q}) \text{ with } \mathbf{t} = \mathbf{x}_{r_{\mathbf{q}}(1)}. \quad (5)$$

Notice that in situations where the k -neighbourhood contains only one class the local model does not find any separation and so it can adopt the majority rule for improving the computational performances.

With this formulation the local learning can switch from the *lazy learning* (Aha, 1997) setting of the original formulation of k NNSVM to the *eager learning* setting with clear advantages in terms of prediction step complexity. This is possible computing a local SVM model for each $\mathbf{x} \in \mathcal{X}$ during the training phase obtaining the sets $\{(\mathbf{t}, k\text{NNSVM}_{\mathbf{t}}) \mid \mathbf{t} \in \mathcal{X}\}$ and applying the precomputed $k\text{NNSVM}_{\mathbf{t}}$ model such that $\mathbf{t} = \mathbf{x}_{r_{\mathbf{q}}(1)}$ for each query point \mathbf{q} during the testing phase.

This approximation slightly modifies the approach of k NNSVM as a local learning algorithm. Instead of estimating the decision function for a *given* test example \mathbf{q} and thus for a specific point in the input metric space, we estimate a decision function for *each* Voronoi cell $V_{\mathbf{x}}$ induced by the training set in the input metric space. In this way, the construction of the models in the training phase requires the estimation of N local decision functions. The prediction of a test point \mathbf{q} is done using the model built for the Voronoi region in which \mathbf{q} lies ($V_{\mathbf{h}}$ with $\mathbf{h} = \mathbf{x}_{r_{\mathbf{q}}(1)}$) that can be retrieved by searching for the nearest neighbour of \mathbf{q} in \mathcal{X} .

4.2 Reducing the Number of Local Models that Need to Be Trained

The pre-computation of the local models during the training phase introduced above, increases the computational efficiency of the prediction step. However, a considerable overhead is added to the training phase. In fact, the training of an SVM for each training point can be slower than the training of a unique global SVM (especially for non small k values), so we introduce another modification of the method which aims to dramatically reduce the number of SVMs that need to be pre-computed. The idea is that we can relax the constraint that a query point \mathbf{x}' is always evaluated using the model trained around its nearest training point. The decision function of this approach is

$$\text{FastLSVM}(\mathbf{x}) = k\text{NNSVM}_{f(\mathbf{x})}(\mathbf{x}) \quad (6)$$

where $f : \mathcal{H} \mapsto \mathcal{C} \subseteq \mathcal{X}$ is a function mapping each unseen example \mathbf{x} to a unique training example $f(\mathbf{x})$ which is, accordingly to Equation 4, the centre of the local model that is used to evaluate \mathbf{x} . The set \mathcal{C} is the image of $f(\cdot)$, so $\mathcal{C} = f(\mathcal{H})$.

Notice that if $f(\cdot) = \mathbf{x}_{r_x(1)}$, we have that $\mathcal{C} = \mathcal{X}$ and that $\text{FastLSVM}(\mathbf{x})$ is equivalent to the $k\text{NNSVM}$ formulation of Equation 5, and this can happen if we use *all* the examples in the training set as centres for local SVM models. In the general case, however, we select only a proper subset $\mathcal{C} \subset \mathcal{X}$ of points to be used as centres of $k\text{NNSVM}$ models. In this case, if $\mathbf{x}_{r_x(1)} \in \mathcal{C}$ then $f(\mathbf{x})$ can be defined as $f(\mathbf{x}) = \mathbf{x}_{r_x(1)}$, but if $\mathbf{x}_{r_x(1)} \notin \mathcal{C}$ then $f(x)$ must be defined in a way such that the principle of locality is preserved and the retrieval of the model is fast at prediction time.

Two aspects need to be addressed now: the strategy to select the subset \mathcal{C} of \mathcal{X} , and the formulation of the function f associating each query example with an example in \mathcal{C} .

4.2.1 SELECTING THE CENTRES OF THE LOCAL MODELS

The approach we developed for selecting the set \mathcal{C} of the centres of the local models is based on the idea that each training point must be in the k' -neighbourhood of at least one centre with k' being a fixed parameter and $k' \leq k$. From a slightly different viewpoint, we need to cover the entire training set with a set of hyper-spheres whose centres will be the examples in \mathcal{C} and each hyper-sphere contains exactly k' points. We can formalise this idea with the concept of k' -neighbourhood covering set:

Definition 1 Given $k' \in \mathbb{N}$, a k' -neighbourhood covering set of centres $\mathcal{C} \subseteq \mathcal{X}$ is a subset of the training set such that the following holds:

$$\bigcup_{c \in \mathcal{C}} \{\mathbf{x}_{r_c(i)} \mid i = 1, \dots, k'\} = \mathcal{X}.$$

Definition 1 means that the union of the sets of the k' -nearest neighbours of \mathcal{C} corresponds to the whole training set. Theoretically, for a fixed k' , the minimisation of the number of local SVMs that we need to train can be obtained computing the SVMs centred on the points contained in the *minimal* k' -neighbourhood covering set of centres.

Definition 2 The Minimal k' -neighbourhood covering set of centres is a k' -neighbourhood covering set $\mathcal{C} \subseteq \mathcal{X}$ which have the minimal cardinality.

This problem is related to the *Set Cover Problem* (SC) (Garey and Johnson, 1979; Kearns and Vazirani, 1994; Marchand and Shawe-Taylor, 2003) and to the *Minimum Sphere Set Covering Problem* (MSSC) (Chen, 2005). However, in the SC and MSSC problems one specifies the radius of the spheres rather than their cardinality in terms of points they contain and it is not required that the centres of the hyperspheres correspond to points in the set. It is easy to show that MSSC is NP-hard but some efficient approximated results are available based on greedy approaches (Chvatal, 1979; Wang et al., 2006), integer and linear programming (Wei and Li, 2008).

In our case, however, we do not need the minimality of the constraints of the k' -neighbourhood covering set of centres to be strictly satisfied, because training some more local SVMs is acceptable instead of solving an NP-hard problem.

The heuristic procedure we developed can be seen as a modification of the greedy approach for the MSSC problem (Chvatal, 1979; Wang et al., 2006). The first k' -neighbourhood is selected randomly choosing its centre in \mathcal{X} , the following k' -neighbourhoods are retrieved selecting the centres that are still not members of other k' -neighbourhoods and are as far as possible from the already selected centres. The selection of the farthest example, still not included in the k' -neighbourhoods, as the centre of the next k' -neighbourhood, is the counterpart of the selection of the set of points having the minimum overlapping with the already covered set of points used by the greedy approach to the MSSC and SC problems.

For detailing the greedy approach we adopt, we need the concepts of minimum and maximum distance between the elements of a set of points A defined respectively as:

$$d(A) = \min \|\mathbf{x} - \mathbf{x}'\| \text{ with } \mathbf{x}, \mathbf{x}' \in A \text{ and } \mathbf{x} \neq \mathbf{x}'$$

and

$$D(A) = \max \|\mathbf{x} - \mathbf{x}'\| \text{ with } \mathbf{x}, \mathbf{x}' \in A.$$

In particular, the minimum distance between points in \mathcal{X} is $m = d(\mathcal{X})$ and the maximum is $M = D(\mathcal{X})$. Our intention is to identify a system of subsets $S_i \subseteq \mathcal{X}$ with decreasing minimum distances $d(S_i)$; we can in this way define an ordering on the sets $\dots \subset S_{i+1} \subset S_i \subset S_{i-1} \subset \dots$ such that $\dots > d(S_{i+1}) > d(S_i) > d(S_{i-1}) > \dots$. With this strategy we can choose the centres of the local models first in the set S_{i+1} , then in the set S_i and so on, thus selecting first the centres that are assured to be distant at least $d(S_{i+1})$, then at least $d(S_i) < d(S_{i+1})$ and so on. More in detail, we require that in the i th set $S_i \subseteq \mathcal{X}$ the two nearest points are farther than b^i with $b > 1$, that is, they are subject to the constraint $d(S_i) > b^i$ with $b > 1$. The bound on the minimum distance $d(S_i)$ thus varies as powers of b depending on the set S_i .

Let us define precisely the system of sets $\{S_i\}$. The maximum i index of S_i is named *top* and the minimum is named *bot*, and they are univocally defined as those indexes satisfying $b^{\text{top}-1} \leq M < b^{\text{top}}$ and $b^{\text{bot}} < m \leq b^{\text{bot}+1}$. The S_i are recursively defined as:

$$\begin{cases} S_{\text{top}} &= \{\text{choose}(\mathcal{X})\} \\ S_i &= S_{i+1} \cup \underset{S \in \mathcal{X} \setminus S_{i+1}}{\text{argmax}}(|S| \text{ s.t. } d(S_{i+1} \cup S) > b^i) \quad \text{for } i = \text{top} - 1, \dots, \text{bot} \end{cases}, \quad (7)$$

where $\text{choose}(A)$ is a function that selects only one element of the non-empty set A . An example of $\text{choose}()$ for our case can be the following definition that selects the example with the minimum index:

$$\text{choose}(A) = \mathbf{x}_i \text{ with } i = \min(z \in \mathbb{N} | \mathbf{x}_z \in A).$$

Notice that, since S_i contains S_{i+1} we have that

$$S_{top} = \{\text{choose}(\mathcal{X})\} \subseteq S_{top-1} \subseteq \dots \subseteq S_{bot+1} \subseteq S_{bot} = \mathcal{X} \quad (8)$$

and, forcing for definition that $d(A) = \infty$ if $|A| = 1$,

$$d(S_{top}) = \infty > d(S_{top-1}) = M > d(S_{top-2}) > \dots > d(S_{bot+1}) > d(S_{bot}) = m.$$

We can now formalise the selection of the centres from \mathcal{X} using the S_i sets. The first centre \mathbf{c}_1 is simply the (only) example in S_{top} . The next centre \mathbf{c}_2 is chosen among the non-empty S_l sets obtained removing from S_i the first centre \mathbf{c}_1 and the points in its k' -neighbourhood; in particular \mathbf{c}_2 is chosen from the non-empty S_l with highest l . The general case for the \mathbf{c}_j centre is similar, with the only difference being that we remove from the S_i sets all the centres \mathbf{c}_t with $t < j$ and their k' -neighbourhood. More formally:

$$\begin{cases} \mathbf{c}_1 &= \text{choose}(S_{top}) \\ \mathbf{c}_j &= \text{choose}(S_l) \text{ with } l = \max(m \in \mathbb{N} | S_m \setminus \mathcal{X}_{\mathbf{c}_{j-1}} \neq \emptyset) \end{cases}, \quad (9)$$

where

$$\mathcal{X}_{\mathbf{c}_{j-1}} = \bigcup_{l=1}^j \left\{ \mathbf{x}_{r_{\mathbf{c}_l}}(h) \mid h = 1, \dots, k' \right\}.$$

is the union of all the k' -neighbourhoods of the centres already included in \mathcal{C} .

We can briefly show that the \mathcal{C} set found with Equation 9 is a k' -neighbourhood covering set of centres. In fact, the iterative procedure for selecting the centres in \mathcal{C} terminates when the $\text{choose}()$ function cannot select a point from S_l because all S_j with $j = bot, \dots, top$ are empty. Since for the set S_{bot} we always have that $S_{bot} = \mathcal{X}$, this happens only when $\mathcal{X}_{\mathbf{c}_{i-1}} = \mathcal{X}$. Noticing that $\mathcal{X}_{\mathbf{c}_i}$ in this situation is equivalent to the constraint of Definition 1, we can conclude that \mathcal{C} is a k' -neighbourhood covering set of centres.

Computationally, the selection of the centres from the S_j sets with Equation 9 can be performed efficiently once the S_j are identified. More problematic is the construction of the nested set of S_j sets. We can however notice that the S_j sets share some characteristics with the levels of cover trees. First, from Equation 7 we can easily see that for each S_j set with $j < top$ all the points in it are at least distant as b^j because $d(S_j) > b^j$; this is equivalent to the separation invariant of cover trees reported in Section 3.4. Second, always from Equation 7 we can conclude that each S_j is contained in every S_t set with $t < j$ as also explicated in Equation 8; this is equivalent to the nesting invariant of cover trees. The only constraint of our strategy to identify the S_j sets that is not respected by cover trees is the maximality of the set added to each S_j set to obtain S_{j+1} . However, the procedure to insert a new point in a cover tree is based on adding it to the highest possible level, and this is an efficient approximation of the maximality constraint we have in Equation 7. Taking all these facts into consideration, we chose to use the levels of cover tree as the S_j sets from which we select the centres as reported in Equation 9.

Consequently with the goal of reducing the number of local models, this approach no longer requires that a local SVM is trained for each training example, but we need to train only $|\mathcal{C}|$ SVMs centred on each $c \in \mathcal{C}$ obtaining the following models:

$$k\text{NNSVM}_{\mathbf{c}}(\mathbf{x}), \quad \forall \mathbf{c} \in \mathcal{C}.$$

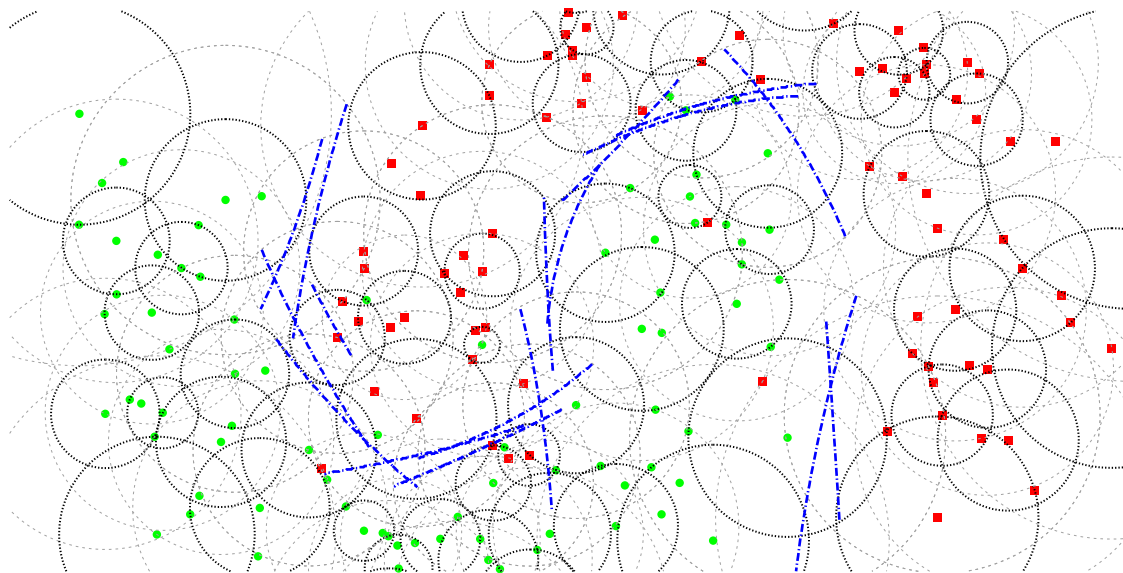


Figure 1: Graphical representation of the proposed approach using local models with $k' = 4$, $k = 15$, and local SVM with RBF kernel. The bold dotted circles highlights the k' -neighbourhoods covering all the training set (with some unavoidable redundancy), the thin dotted circles denotes the k -neighbourhoods on which the local models are trained. Some k -neighbourhoods do not produce an explicit decision function because entirely composed by points of the same class. The local SVM (with RBF kernel) decision functions are drawn in blue. Notice that, due both to the adoption of the k' -neighbourhood cover set and to the fact that only a fraction of the neighbourhoods need to be trained, we have only 17 local decision functions for 185 points.

Moreover if a neighbourhood contains only points belonging to one class the local model is the majority rule (specifically, unanimity) and the training of the SVM is avoided.

Figure 1 graphically shows the result of adopting the approach described above on a simple artificial data set with k and k' chosen for illustrative purposes. In fact, the example just aims to show the intuition behind the approach that is instead developed for large data sets and for non-extreme values of the neighbourhood parameters.

From Figure 1 we can also notice that the level of overlapping between k' -neighbourhoods and thus between k -neighbourhoods depends on the value of k' . If k' is low, a large number of k' -neighbourhoods are required to cover the entire training set, whereas if k' is large fewer k' -neighbourhoods are needed. The k' parameter thus tune the level of redundancy of the local models.

4.2.2 SELECTING THE LOCAL MODELS FOR TESTING POINTS

Once the set of centres \mathcal{C} is defined and the corresponding local models are trained, we need to select the proper model to use for predicting the label of a test point. A simple strategy we can adopt consists in selecting the model whose centre $\mathbf{c} \in \mathcal{C}$ is the nearest centre with respect to the

testing example. Using the general definition of FastLSVM of Equation 6 with $f(x) = r_{\mathbf{x}}^C(1)$ where r^C corresponds to the reordering function defined in Equation 3 performed on the C set instead of \mathcal{X} , the method, called FaLK-SVMc, is defined as:

$$\text{FaLK-SVMc}(\mathbf{x}) = k\text{NNSVM}_{\mathbf{c}}(\mathbf{x}) \text{ where } \mathbf{c} = \mathbf{x}_{r_{\mathbf{x}}^C(1)}. \quad (10)$$

FaLK-SVMc is satisfactory from the computational viewpoint, for it performs the nearest neighbour search on C only. However, it does not assure that the testing point is evaluated with the model centred on the point for which the testing point itself is the nearest in terms of neighbour ranking. For example, a testing point \mathbf{q} can be closer to \mathbf{c}_1 than \mathbf{c}_2 using the Euclidean distance, but at the same time we can have that \mathbf{q} is the i -th nearest neighbour of \mathbf{c}_1 in \mathcal{X} and the j -th nearest neighbour of \mathbf{c}_2 with $i > j$. This is a problem because using the model centred on \mathbf{c}_2 is better in terms of proximity. In order to overcome this issue of FaLK-SVMc we propose to use, for a testing point \mathbf{q} , the model centred on the training point which is the nearest in terms of the neighbourhood ranking to its training nearest neighbour. We can do this defining a function $\text{cnt} : \mathcal{X} \mapsto C$ in the following way:

$$\begin{aligned} \text{cnt}(\mathbf{x}_i) &= \text{choose}(\{\mathbf{c}_z \in C \mid \mathbf{x}_i = \mathbf{x}_{r_{\mathbf{c}_z}(h)}\}) \\ \text{where } h &= \min(t \in \{1, \dots, k'\} \mid \mathbf{x}_{r_{\mathbf{c}_j}(t)} = \mathbf{x}_i \text{ and } \mathbf{c}_j \in C). \end{aligned} \quad (11)$$

The cnt function finds, for each example \mathbf{x} , the minimum value h such that \mathbf{x} is in the h -neighbourhood of at least one centre $\mathbf{c} \in C$; then, among the centres having \mathbf{x} in their h -neighbourhoods, it selects the centre with the minimum index. The existence of h is guaranteed by the k' -neighbourhood covering strategy. In this way each training point is univocally assigned to a centre and so the decision function of this approximation of Local SVM derivable from FastLSVM of Equation 6 with $f(\mathbf{x}) = \text{cnt}(\mathbf{x})$, and called FaLK-SVM, is simply:

$$\text{FaLK-SVM}(\mathbf{x}) = k\text{NNSVM}_{\text{cnt}(\mathbf{t})}(\mathbf{x}) \text{ where } \mathbf{t} = \mathbf{x}_{r_{\mathbf{x}}(1)}. \quad (12)$$

The association between training points and centres defined by Equation 11 can be efficiently precomputed during the training phase, delaying to the testing phase only the retrieval of the nearest neighbour of the testing point and the evaluation of the local SVM model.

Figure 2 graphically shows the application of the FaLK-SVM(\mathbf{x}) prediction strategy on a toy data set; the training phase for the same data set is illustrated in Figure 1.

FaLK-SVM can be generalised for multi-class problems in the same way of $k\text{NNSVM}$, but in this paper we focus on binary problems in order to better evaluate the approach.

4.3 FaLK-SVM with Internal Model Selection: FaLK-SVMl

For training a kernel machine, once a proper kernel is chosen, it is crucial to carefully tune the kernel parameters and, for SVM, to set the soft margin regularisation constant C . Model selection is very often performed estimating the empirical error with different parameter values and a popular method is the κ -fold cross-validation² with a grid search on parameter space. Given the following loss function for the two-class classification case

$$L(y, \text{SVM}(\mathbf{x})) = \begin{cases} 0 & \text{if } y = \text{SVM}(\mathbf{x}) \\ 1, & \text{if } y \neq \text{SVM}(\mathbf{x}) \end{cases},$$

2. Although κ can be confused with the neighbourhood size k or with the kernel function K , κ is always used for denoting κ -fold CV, so the context should be sufficient to avoid ambiguity.

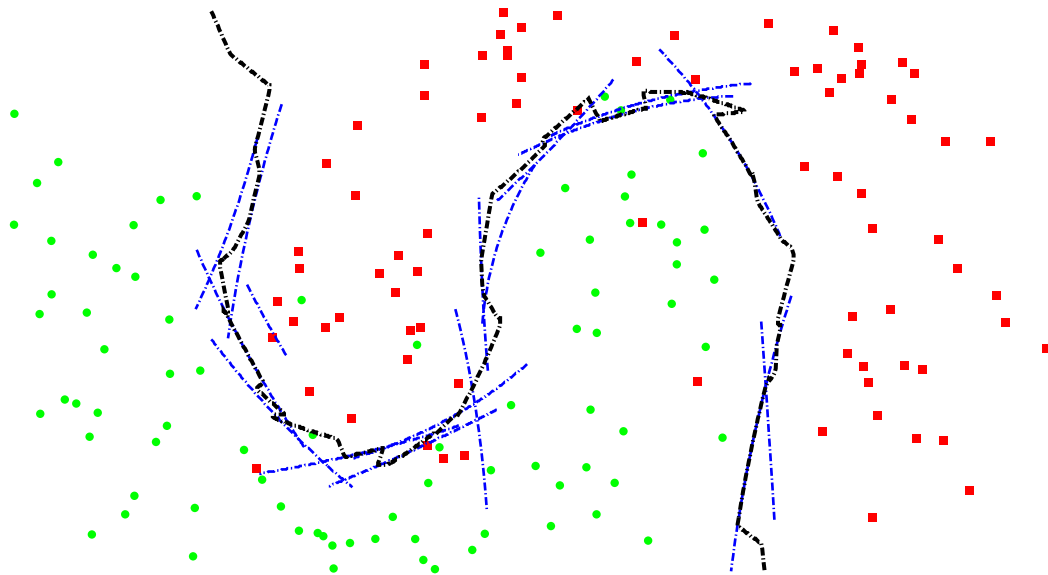


Figure 2: Graphical representation of the global decision function (black dotted line) obtained with the local decision functions (the same of Figure 1) using the described approach that uses for each query point the local decision function of the Voronoi region in which it lies.

and partitioning the training set \mathcal{X} in κ subsets each with the same cardinality³ (called folds), the κ -fold cross validation (CV) procedure consists in searching for the parameters that minimise the average of the losses on \mathcal{X}_f of the classifier trained on $\mathcal{X} \setminus \mathcal{X}_f$ for $f = 1, \dots, \kappa$. The effectiveness in terms of testing accuracies of κ -fold CV is high, but it adds a computational overhead to the training phase. In fact, the computational complexity of a κ -fold CV run on a single parameter choice is in the order of κ times the training time; if we have p parameters to set and c possible choices for each parameter, the κ -fold cross-validation with grid selection is $\kappa \cdot c^p$ times slower than a single training of the classifier.

The model selection for FaLK-SVM and FaLK-SVMc can be performed using κ -fold CV. The only difference with SVM is that our local kernel machines need to estimate an additional parameter which is the neighbourhood size k (which is however usually chosen in a small set of possible values). However, with the local setting of the classification problem we are discussing in this paper, it is also possible to efficiently tackle the complexity of the model selection phase. Basically, since FaLK-SVM trains a set of local models, we can perform the model selection in a grid-search setting on a subset of the neighbourhoods. In this way we can efficiently estimate the global parameters of FaLK-SVM without considering all the training points during model selection. The classifier implementing this approach to model selection is called FaLK-SVMl.

As a first step for defining the model selection approach of FaLK-SVMl, we define a different setting of model selection for k NNSVM.

3. Without loss of generality, we assume $|\mathcal{X}| \bmod \kappa = 0$.

Definition 3 (Localised κ -fold CV model selection for k NNSVM) *The procedure applies the κ -fold CV model selection on the k -neighbourhood of the query point.*

However, since the local model is used by k NNSVM only for the central point, the model selection should be performed in order to make the local models predictive especially for the very internal points. The idea thus consists in selecting the κ validation sets exclusively from the k' most internal points, taking as each corresponding training fold the union of the remaining k' -neighbourhood points and of the $k - k'$ most external points of the k -neighbourhood.

Definition 4 (k' -internal κ -fold CV model selection for k NNSVM)

The procedure applies the localised κ -fold CV model selection on the k' -neighbourhood of the query point in the training set adding to each training fold the points in the k -neighbourhood that are not in the k' -neighbourhood with $k > k'$.

For FaLK-SVM we can apply the k' -internal κ -fold CV for k NNSVM model selection on a randomly chosen training example and use the resulting parameters for all the local models. In order to be robust the procedure is repeated on more than one k -neighbourhood choosing the parameters that minimise the average k' -internal κ -fold CV error among the k -neighbourhoods.

Definition 5 (k' -internal κ -fold CV model selection for FaLK-SVM)

The procedure applies the k' -internal κ -fold CV for k NNSVM model selection on the k -neighbourhoods of $1 \leq m \leq |C|$ randomly chosen centres selecting the parameters that minimise the average error rate among the m applications.

The variant of FaLK-SVM that adopts the k' -internal κ -fold CV described in Definition 5 is named FaLK-SVMl. Since FaLK-SVMl selects the local model parameters using a small subset of the training set, the variance of the error may be higher than the standard cross-validation strategies. However, for huge data sets the standard model selection can be too slow to be applied and, in any case, one may use large values of m to decrease the risk of selecting non-optimal parameters.

4.3.1 A SPECIFIC STRATEGY FOR SETTING THE RBF KERNEL WIDTH

As already proposed by Tsang et al. (2005) and by Segata and Blanzieri (2009b), good choices for the RBF kernel width σ of SVM are based on the median (or other percentiles) of the distribution of distances. In FaLK-SVMl we can thus efficiently estimate σ for each local model simply calculating the median of the distances in the neighbourhood. This approach has some analogies with standard SVM using a variable RBF kernel width that have good potentialities for classification (Chang et al., 2005). Since other percentiles different from the median can give better accuracy performances, in FaLK-SVMl the percentile can be a value to set using the k' -internal κ -fold CV approach.

4.4 Generalisation Bounds for k NNSVM and FaLK-SVM

The class of LLAs introduced by Bottou and Vapnik (1992) includes k NNSVM, and can be theoretically analysed using the framework based on the local risk minimisation (Vapnik and Bottou, 1993; Vapnik, 2000). On the other hand, FaLK-SVM is not a LLA as intended by Bottou and Vapnik (1992). In fact, LLAs compute the local function for each specific testing point thus delaying the neighbourhood retrieval and model training until the testing point is available. However, we show here that generalisation bounds for FaLK-SVM can be derived starting from the LLA ones.

We need to recall the bound for the local risk minimisation, which is a generalisation of the global risk minimisation theory.

Theorem 6 (Vapnik (2000)) *For a testing point \mathbf{x}' and with probability $1 - \eta$ simultaneously for all bounded functions $A \leq L(y, f(\mathbf{x}, \alpha)) \leq B$, $\alpha \in \Lambda$ (where Λ is a set of parameters), and all locality functions $0 \leq T(\mathbf{x}, \mathbf{x}_0, \beta) \leq 1$, $\beta \in (0, \infty)$, the following inequality holds true:*

$$R^{LLA}(\alpha, \beta, \mathbf{x}') \leq \frac{\frac{1}{N} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i, \alpha)) T(\mathbf{x}_i, \mathbf{x}', \beta) + (B - A) \gamma(N, h^\Sigma)}{|\frac{1}{N} \sum_{i=1}^N T(\mathbf{x}_i, \mathbf{x}', \beta) - \gamma(N, h^\beta)|},$$

where

$$\gamma(N, h) = \sqrt{\frac{h \ln(2N/h + 1) - \ln \eta / 2}{N}},$$

and h^Σ is the VC dimension of the set of functions $L(y_i, f(\mathbf{x}_i, \alpha)) T(\mathbf{x}_i, \mathbf{x}', \beta)$, $\alpha \in \Lambda$, $\beta \in (0, \infty)$ and h^β is the VC dimension of $T(\mathbf{x}_i, \mathbf{x}', \beta)$

For k NNSVM, the loss function is simply

$$L(y_i, f(\mathbf{x}_i, \alpha)) = \begin{cases} 0 & \text{if } y_i = f(\mathbf{x}_i, \alpha) \\ 1 & \text{if } y_i \neq f(\mathbf{x}_i, \alpha) \end{cases}$$

and the locality function is

$$T(\mathbf{x}_i, \mathbf{x}', k) = \begin{cases} 1 & \text{if } \exists j \leq k \text{ s.t. } i = r_{\mathbf{x}'}(j) \\ 0 & \text{otherwise} \end{cases}.$$

It is straightforward to show that $\sum_{i=1}^N T(\mathbf{x}_i, \mathbf{x}', k) = k$. Moreover $T(\mathbf{x}_i, \mathbf{x}', k)$ has VC dimension equal to 2; it is, in fact, the class of functions corresponding to hyperspheres centred on \mathbf{x}' with diameters equal to the distances of the points from \mathbf{x}' and can thus shatter any set of two points with different classes, but cannot shatter three points with the nearest and furthest points having a class different from the third point.

We observe that, in our case,

$$\sum_{i=1}^N L(y_i, f(\mathbf{x}_i, \alpha)) T(\mathbf{x}_i, \mathbf{x}', \beta) = \sum_{i=1}^k L(y_i, f(\mathbf{x}_i, \alpha))$$

and so we can obtain:

$$R^{k\text{NNSVM}}(\alpha, k, \mathbf{x}') \leq \frac{\frac{1}{N} k \cdot v_{\mathbf{x}'} + \gamma(N, h^\Sigma)}{|\frac{1}{N} k - \gamma(N, 2)|} \quad (13)$$

where $v_{\mathbf{x}'}$ is the ratio of misclassified training points in the k -neighbourhood of \mathbf{x}' .

The possibility of local approaches to obtain a lower bound on test misclassification probability acting with the locality parameter, as stated in Vapnik and Bottou (1993); Vapnik (2000) for LLA, it is even more evident for k NNSVM considering Equation 13. In fact, although choosing a $k < N$ is not sufficient to lower the bound, as the model training becomes more and more local k decreases and (very likely) the misclassification training rate $v_{\mathbf{x}'}$ decreases as well. Moreover, also the complexity of the classifier (and thus h^Σ) can decrease when the neighbourhood decreases, because simpler decision functions can be used when fewer points are considered. Taking this into

consideration, it is necessary to consider the trade-off between the degree of locality k , the function of the empirical error with respect to k and the complexity of the local classifier needed with respect to k , in order to find a minimum of the expected risk which is lower than the $k = N$ case. Multiple strategies can be used to tune this trade-off, especially if prior or high-level information are available for a specific problem; since in this work we aim to be as general as possible, the expected risk is estimated for the computational experiments using cross-validation based approaches.

FaLK-SVM pre-computes local models to be used for testing points lying in sub-regions (k -NN Voronoi cells) of the training set. The risk associated to FaLK-SVM considering a specific query point \mathbf{x}' can be defined using the risk of k NNSVM, supposing that $\mathbf{x}' \in V_{\mathbf{x}_i}$ and so $\mathbf{x}_{r_{\mathbf{x}'}(1)} = \mathbf{x}_i$:

$$R^{\text{FaLK-SVM}}(\alpha, k, \mathbf{x}') = R^{k\text{NNSVM}}(\alpha, k, \mathbf{x}') + \lambda(\mathbf{x}', \mathbf{x}_{r_{\mathbf{x}'}(1)}) \leq R^{k\text{NNSVM}}(\alpha, k, \mathbf{x}') + \lambda_{r_{\mathbf{x}'}(1)} \quad (14)$$

where $\lambda(\mathbf{x}', \mathbf{x}_{r_{\mathbf{x}'}(1)})$ is due to the approximation introduced, for the prediction of the label of the query point \mathbf{x}' , by the use of the k -neighbourhood of $r_{\mathbf{x}'}(1)$ instead of the k -neighbourhood of \mathbf{x}' itself and

$$\lambda_{r_{\mathbf{x}'}(1)} = \max_{\mathbf{x}'' \in V_{\mathbf{x}_i}} \lambda(\mathbf{x}'', \mathbf{x}_{r_{\mathbf{x}'}(1)}).$$

If we consider $k' = 1$, the approximation is due to the fact that $\{r_{\mathbf{x}}(i) | i = 1, \dots, k\}$ and $\{r_{\mathbf{x}'}(i) | i = 1, \dots, k\}$ can be slightly different; however, considering a non very low value for k , the differences between the two sets are possible only for the very peripheral points of the neighbourhoods which are those that influence less the shape of the decision function in the central region. We will empirically show that $\lambda_{r_{\mathbf{x}'}(1)}$ is, on average, a small penalising term that still permits to achieve lower risks than SVM using k' values higher than 1.

The risk of FaLK-SVM in its eager learning setting (i.e., without the explicit dependency on the query point) can thus be defined as:

$$\begin{aligned} R^{\text{FaLK-SVM}}(\alpha, k) &= \int_{\mathbf{x}'} R^{\text{FaLK-SVM}}(\alpha, k, \mathbf{x}') g(\mathbf{x}') d\mathbf{x}' & (15) \\ &\leq \int_{\mathbf{x}'} \left(R^{k\text{NNSVM}}(\alpha, k, \mathbf{x}_i) + \lambda_{r_{\mathbf{x}'}(1)} \right) g(\mathbf{x}') d\mathbf{x}' \\ &= \int_{\mathbf{x}'} R^{k\text{NNSVM}}(\alpha, k, \mathbf{x}_i) g(\mathbf{x}') d\mathbf{x}' + \int_{\mathbf{x}'} \lambda_{r_{\mathbf{x}'}(1)} g(\mathbf{x}') d\mathbf{x}' \\ &= \int_{\mathbf{x}'} R^{k\text{NNSVM}}(\alpha, k, \mathbf{x}_i) g(\mathbf{x}') d\mathbf{x}' + E[\lambda]. \end{aligned}$$

where $E[\lambda]$ is the expectation of the term due to the use of the k NNSVM risk for FaLK-SVM as discussed above.

4.5 Computational Complexity Analysis

We analyse here the computational performances of FaLK-SVM from the theoretical complexity viewpoint. The training phase of FaLK-SVM can be subdivided in four steps:

- the building of the cover tree that scales as $O(N \log N)$;
- the retrieval of the local models that scales as $O(|C| \cdot k \log N)$;
- the assignment of each point to a k' -neighbourhood that scales as $O(N)$;

- the training of the local SVM models that scales as $O(|C| \cdot k^3)$.

The overall training time, considering the worst case in which $k' = 1$ so $|C| = N$, scales as:

$$O(N \log N + C \cdot k \log N + N + C \cdot k^3) = O(kN \cdot \max(\log N, k^2))$$

which is, considering a reasonably low and fixed value for k as happens in practice for large data sets, sub-quadratic, and in particular $O(N \log N)$, in the number of training points.

For the testing phase of FaLK-SVM we can distinguish two steps (for each testing point):

- the retrieval of the nearest training point that scales as $O(\log N)$;
- the prediction of the testing label using the selected local model that scales as $O(k)$.

The testing can thus be performed in $O(\max(\log N, k))$, so it is logarithmic in N . FaLK-SVMc is even faster because it scales as $O(\max(\log |C|, k)) \leq O(\max(\log N, k))$.

FaLK-SVM is thus asymptotically faster than SVM (also considering the worst case in which SVM scales quadratically and $k' = 1$) and all the classifiers taking more than $O(N \log N)$ for training and $O(\log N)$ for testing. Moreover, FaLK-SVM can be very easily parallelised differently from SVM whose parallelisation, although possible (Zanni et al., 2006; Dong, 2005), is rather critical; for FaLK-SVM is sufficient that, every time the points for a model are retrieved, the training of the local SVM is performed on a different processor. In this way the time complexity of FaLK-SVM can be further lowered to $O(N \cdot \max(k \log N, k^3 / N_{proc}))$ where N_{proc} is the number of processors.

Another advantage of FaLK-SVM over SVM is space complexity. Since FaLK-SVM performs SVM training on small subregions (assuming a reasonable low k), there are no problems of fitting the kernel matrix into main memory. The overall required space is, in fact, $O(N + k^2)$, that is, linear in N , which is much lower than SVM space complexity of $O(N^2)$. For large data sets, FaLK-SVM can still maintain in memory the entire local kernel matrix (if k is not too large), whereas SVM must discard some kernel values thus increasing SVM time complexity due to the need of recomputing them. Analysing the space required to store the trained model in secondary storage devices (e.g., hard disks), we can notice that FaLK-SVM needs to save in the model file the entire set of local models; although we store the models with pointers to the training set points, we need to maintain the whole training set in the model file (or give as input for the testing module both the model file and the original training set). FaLK-SVM, in other words, needs to store the training set also in the model file, differently from SVM that needs to store only the support vectors (whose number however grows linearly with N).

4.5.1 CURSE OF DIMENSIONALITY

Although not explicitly considered here, cover trees have a constant in the complexity bounds depending on the so-called doubling constant (Clarkson, 1997; Krauthgamer and Lee, 2004) which is a robust estimation of the intrinsic dimensionality of the data. Notice that the intrinsic dimensionality of a data set can be much lower than the dimensionality intended simply as the number of features. Regardless of the doubling constant, FaLK-SVM maintains the derived complexity bounds⁴ with respect to N , but the overhead introduced for building the cover tree and retrieving the k -neighbourhoods can be very high. This drawback, due to the well-known problem of the *curse of*

4. The high intrinsic dimensionality can cause the need for an high value of $|C|$, but in the bound we already considered the worst case in which $k' = 1$ and thus $|C| = N$.

Algorithm 1 FaLK-SVM-train (training set $\mathbf{x}[]$, training size \mathbf{n} , neighbourhood size \mathbf{k} , assignment neighbourhood size \mathbf{k}')

```

1:  $models[] \leftarrow \mathbf{null}$  // the set of models
2:  $modelPtrs[] \leftarrow \mathbf{null}$  // the set of pointers to the models
3:  $c \leftarrow 0$  // the counter for the centres of the models
4:  $indexes[] \leftarrow \{1, \dots, N\}$  // the indexes for centres selection
5: Randomise  $indexes$  // randomise the indexes
6: for  $i \leftarrow 1$  to  $N$  do
7:    $index \leftarrow indexes[i]$  // get the i-th index
8:   if  $modelPtrs[index] = \mathbf{null}$  then // if the point has not been assigned to a model...
9:      $localPoints[] \leftarrow$  get ordered  $k$ NN of  $x[i]$  // ... retrieve its  $k$ -neighbourhood ...
10:     $models[c] \leftarrow$  SVMtrain on  $localPoints[]$  // ... train a local SVM. ...
11:     $modelPtrs[index] \leftarrow models[c]$  // ... assign the centre to the trained model.
12:    for  $j = 1$  to  $k'$  do // Assign the model to the  $k' < k$  nearest neighbours of the centre
13:       $ind \leftarrow$  get index of  $localPoints[j]$ 
14:      if  $modelPtrs[ind] = \mathbf{null}$  then // assign the points in the  $k'$ -neighbourhood ...
15:         $modelPtrs[ind] \leftarrow models[c]$  // ... to the  $c$ -th model
16:      end if
17:    end for
18:     $c \leftarrow c+1$ 
19:  end if
20: end for
21: return  $models, modelPtrs$ 

```

Algorithm 2 FaLK-SVM-predict (training set $\mathbf{x}[]$, points-to-model pointers $\mathbf{modelPtrs}$, Local SVM models \mathbf{models} , query point \mathbf{q})

```

1: Set  $p =$  get NN of  $q$  in  $x$  // retrieve the nearest training point with respect to  $q$ . ...
2: Set  $nnIndex =$  get index of  $p$  // ... retrieve its index ...
3: return  $label =$  SVMpredict  $q$  with  $modelPtrs[nnIndex]$  // ... and use the corresponding model
   for predict the label of the query point.

```

dimensionality that affects also SVM with local kernels (Bengio et al., 2005), is not however crucial here, as we are considering non-linear classification problems that are not high-dimensional. In fact, apart from computational problems, high-dimensional problems are typically tackled by approaches not related with the concept of locality (e.g., linear SVM instead of SVM with a RBF kernel).

4.6 Implementation and Availability

FaLK-SVM (and also FkNN and FkNNSVM that are the implementations of kNN and k NNSVM using cover trees) is available as part of the Fast Local Kernel Machine Library (Segata, 2009, FaLKM-lib). FaLK-SVM is written in C/C++ and it uses LibSVM v. 2.88 (Chang and Lin, 2001) for local SVM training and testing whereas we use our own implementation of the cover trees data-structure. The pseudo-code for the training phase is reported in Algorithm 1 and for the testing phase

method	brief description
FkNN	implementation of kNN (Section 3.1) with cover trees
FkNNSVM	implementation of k NNSVM (Section 3.3) with cover trees
FaLK-SVM	implementation of fast and scalable local kernel machines (see Equation 12)
FaLK-SVM-train	module for the training of FaLK-SVM (see Algorithm 1)
FaLK-SVM-predict	module for the testing of FaLK-SVM (see Algorithm 2)
FaLK-SVMc	faster prediction variant of FaLK-SVM (see Equation 10)
FaLK-SVMl	implementation of FaLK-SVM with local model selection (Section 4.3)
FkNNSVM-nr	implementation of k NNSVM for noise reduction (Segata et al., 2009b)
FaLKNR	impl. of noise reduction with FaLK-SVM (Segata et al., 2009a)

Table 1: Summary for the classifiers developed in the local kernel machine framework and implemented in FaLKM-lib.

in Algorithm 2 (use of cover trees and minimisation of t in Equation 11 are omitted for clearness). Table 1 summarizes the classifiers discussed in this paper and the modules of FaLKM-lib.

5. Empirical Analysis

The empirical analysis is organised into three experiments performed with different objectives and using different data sets. Experiment 1 (Section 5.1) has the objective of assessing the generalisation performances of FaLK-SVM with respect to SVM (using LibSVM) and to k NNSVM (using FkNNSVM) and thus assessing if FaLK-SVM is more accurate than SVM and if it is a good approximation of k NNSVM. For this experiment we use 25 non-large data sets. Experiment 2 (Section 5.2) focuses on comparing the classification accuracies and the computational performances of FaLK-SVM (and its variants FaLK-SVMc and FaLK-SVMl) with respect to SVM (using LibSVM) on large data sets. For this experiment we use 8 data sets with training set cardinalities ranging from about 50k examples to more than 1 million. Experiment 3 (Section 5.3) aims to understand (i) whether FaLK-SVM has better scalability and accuracy performances than LibSVM, a number of approximated SVM solvers (CVM, BVM, LASVM, CPSP and USVM) and SVM-bagging and (ii) which are the computational and accuracy differences between FaLK-SVM, FaLK-SVMc and FaLK-SVMl. For this last experiment we use 4 data sets with increasing training set size up to 3 million examples. The experiments, unless otherwise specified, are carried out on an AMD Athlon 64 X2 Dual Core Processor 5000+, 2600MHz, with 3.56Gb of RAM with Linux operating system.

5.1 Experiment 1: Comparison of FaLK-SVM with LibSVM and FkNNSVM

In this evaluation we compare SVM (using LibSVM), k NNSVM (using FkNNSVM) and FaLK-SVM on 25 non-large data sets, with the objective of studying the generalisation performances of k NNSVM with respect to SVM and the level of approximation introduced by FaLK-SVM to the FkNNSVM algorithm.

data set name	# of features	# of points	class balancing	data set name	# of features	# of points	class balancing
sonar	60	208	53%/47%	fourclass	2	862	64%/36%
heart	13	270	56%/44%	tic-tac-toe	9	958	65%/35%
mushrooms	112	300	53%/47%	mam	5	961	54%/46%
haberman	3	306	74%/26%	numer	24	1000	70%/30%
liver	6	345	58%/42%	splice	60	1000	52%/48%
ionosphere	34	351	64%/36%	spambase	57	1000	57%/43%
vote	15	435	61%/39%	vehicle	21	1243	76%/24%
musk1	166	476	57%/43%	cmc	7	1473	57%/43%
hill-valley	100	606	51%/49%	ijcnn1	22	1500	68%/32%
breast	10	683	65%/35%	a1a	123	1605	76%/24%
australian	14	690	56%/44%	chess	35	2130	52%/48%
transfusion	4	748	76%/24%	astro	4	3089	65%/35%
diabetes	8	768	65%/35%				

Table 2: The 25 binary-class data sets of Experiment 1.

5.1.1 EXPERIMENTAL PROTOCOL

The data sets are listed in Table 2; they are retrieved from the UCI (Asuncion and Newman, 2007) and STATLOG (Michie et al., 1994) repositories, with cardinality between 200 and 3100 points (some data sets have been randomly sub-sampled), dimensionality lower than 200, not very unbalanced, and they are all scaled in the $[0, 1]$ interval. The comparison is carried out using three different kernel functions (linear, RBF and homogeneous polynomial), in a 10-fold CV experimental setting. Internal to each training fold the model selection is performed with a nested 10-fold CV choosing the parameters in the following ranges. The regularisation parameter C is chosen for all methods in the set $\{2^{-2}, 2^{-1}, \dots, 2^9, 2^{10}\}$, the width parameter σ of the RBF kernel in $\{2^{-5}, 2^{-4}, \dots, 2^2, 2^3\}$, the degree of the polynomial kernels in $\{1, 2, 3\}$. The neighbourhood parameter k for FkNNSVM and FaLK-SVM is selected by the cross-validation procedure in the set $\{2^1, 2^2, \dots, 2^9, 2^{10}, |\mathcal{X}|\}$ where $|\mathcal{X}|$ is the cardinality of the training set,⁵ while the k' parameter of FaLK-SVM is fixed to $k/2$ which is a value that privileges scalability over accuracy because we want to test a value that can permit good computational results for large and very large data sets.

5.1.2 RESULTS AND DISCUSSION

Table 3 reports the accuracy results of all tested methods and kernels. In addition to the mean ranks reported in the figure, we assessed the statistical significance of the differences between pairs of methods using the Wilcoxon Signed Rank Test (Wilcoxon, 1945; Demšar, 2006) with $\alpha = 0.05$. The test highlights that FkNNSVM is significantly better than LibSVM for the linear and polynomial kernels, whereas for the RBF kernel no significant differences are detected, although the mean rank of FkNNSVM with RBF kernel is lower than LibSVM with RBF kernel. Applied to FaLK-SVM, the Wilcoxon Signed Rank Test detects a significant difference with respect to LibSVM only for the linear kernel. If we perform the Friedman test (Friedman, 1940) ($\alpha = 0.05$), the null hypothesis is

5. For data set with less than 1024 points some k values are of course not tested.

data set	LibSVM			FkNNSVM			FaLK-SVM		
	K^{lin}	K^{rbf}	K^{hpol}	K^{lin}	K^{rbf}	K^{hpol}	K^{lin}	K^{rbf}	K^{hpol}
sonar	74.52	87.83	83.16	89.36	86.90	87.40	84.55	87.88	84.05
heart	84.81	82.22	84.81	84.81	81.11	84.81	83.70	81.85	83.70
mushrooms	97.99	98.33	98.32	98.67	98.33	98.6	99.00	99.00	99.00
haberman	73.20	73.20	72.89	75.82	75.16	74.18	73.25	73.20	73.87
liver	68.71	74.24	71.90	73.64	73.96	73.94	70.73	71.92	71.92
ionosphere	88.04	93.72	88.88	93.75	94.59	93.75	86.91	94.01	89.18
vote	94.95	96.32	94.95	96.32	96.33	96.32	94.94	96.32	94.94
musk1	86.55	94.54	93.07	89.44	94.96	91.17	87.18	93.90	92.43
hill-valley	63.70	66.00	63.70	64.86	65.18	64.86	65.17	64.03	65.00
breast	96.78	96.78	96.78	96.49	96.49	96.35	96.19	96.49	96.19
australian	85.50	84.78	84.20	84.78	85.50	84.92	85.07	85.07	84.78
transfusion	76.21	77.40	76.47	79.81	78.74	79.81	79.67	78.87	79.94
diabetes	76.54	76.54	76.68	76.81	78.24	77.07	75.90	76.68	75.12
fourclass	77.39	100.00	78.66	100.00	100.00	100.00	100.00	100.00	100.00
tic-tac-toe	98.33	99.68	100.00	100.00	100.00	100.00	100.00	100.00	100.00
mam	82.10	82.63	81.27	82.95	82.73	82.85	81.80	82.63	80.97
numer	77.00	75.90	76.50	76.30	75.70	76.00	76.70	74.70	75.90
splICE	80.41	86.70	86.60	80.41	86.30	86.60	78.30	86.20	86.60
spambase	89.80	90.60	89.80	90.60	90.50	90.60	90.70	90.60	90.70
vehicle	82.71	84.16	84.80	82.78	84.64	84.71	83.27	84.72	85.04
cmc	59.26	65.45	64.16	62.46	67.72	63.61	63.61	65.31	64.36
ijcnn1	85.53	93.94	92.73	93.93	93.47	93.60	92.80	94.47	93.20
a1a	83.43	81.94	83.43	82.87	82.06	82.87	82.87	82.06	82.87
chess	96.57	98.45	98.03	97.84	98.50	98.08	97.32	98.45	98.08
astro	95.34	96.73	96.89	96.96	96.92	97.05	96.96	96.67	96.86
mean rank	7.04	4.60	5.80	4.38	3.86	4.02	5.72	4.56	5.02

Table 3: 10-fold CV accuracy results for the 25 data set of Experiment 1. The best results for each data set are highlighted in bold (taking into account all decimal values).

violated, but, according to the Nemenyi post-hoc test (Nemenyi, 1963) ($\alpha = 0.05$) the only method that is statistically significantly different from the others is SVM with linear kernel.

The observation that FkNNSVM is significantly better than SVM if a non-local kernel is used, is a confirmation of what we already noticed (Segata and Blanzieri, 2009a). Using the RBF kernel, instead, no significant differences are detected, although the mean rank of FkNNSVM with RBF kernel is lower than LibSVM with RBF kernel. This is mainly due to the fact that SVM with RBF kernel is already very accurate and significant improvements over it are very difficult. We may also say that locality is already included in the RBF kernel and thus, at least for non-large data sets, the adoption of a local method is somehow equivalent. Regarding FaLK-SVM, significant differences with respect to LibSVM are detected only for the linear kernel. Although FaLK-SVM does not achieve the accuracy results of FkNNSVM, if we look to the mean ranks, we can conclude that the approximation on the k NNSVM approach introduced in FaLK-SVM still permits to achieve slightly

data set name	# of feat.	train. points	testing points	class balancing	original source
ijcnn1	22	49990	91701	90%/10%	LibSVM rep. (Chang and Lin, 2001)
cov-type *	54	100000	481010	51%/49%	LibSVM rep. (Chang and Lin, 2001)
census-inc	41	199523	99762	94%/6%	UCI rep. (Asuncion and Newman, 2007)
cod-rna	8	364651	121549	67%/33%	(Uzilov et al., 2006)
intr-det	40	1026588	311029	79%/21%	UCI KDD rep. (Hettich and Bay, 1999)
2-spirals *	2	100000	100000	50%/50%	Synthetic (Segata and Blanzieri, 2009c)
ndcc *	5	100000	100000	61%/39%	Synthetic (Thompson, 2006)
checker-b *	2	300000	100000	50%/50%	Synthetic (e.g., see Tsang et al., 2005)

Table 4: The 8 large data sets of the second empirical experiment. The data sets whose extensions are used also in Experiment 3 are denoted with *.

better results than SVM also on non-large data sets, confirming our preliminary analysis (Segata and Blanzieri, 2009c). These results also indicates that the $E[\lambda]$ term introduced in the risk of FaLK-SVM (Section 4.4), due to the approximations introduced to the k NN SVM approach, is small enough to assure higher generalisation accuracies with respect to SVM.

The overall outcome of this experiment is that FaLK-SVM is a good approximation of FkNN SVM that maintains a little advantage over SVM and it is particularly effective with the RBF kernel with respect to linear and polynomial kernels. Notice that the experiment is carried out using small data sets in which locality is very likely to play a marginal role differently from large data sets in which it can be crucial.

5.2 Experiment 2: FaLK-SVM, FaLK-SVMc and FaLK-SVMI vs. LibSVM and FkNN on Large Data Sets

In this experiment we apply FaLK-SVM, FaLK-SVMc, FaLK-SVMI, LibSVM on 8 large data sets comparing the computational and generalisation performances using the RBF kernel, because preliminary experiments showed that the linear or polynomial kernels have very low accuracy results on the considered problems. We also add to the comparison the k NN classifier (implemented with cover trees and called FkNN) using the Euclidean distance.

5.2.1 EXPERIMENTAL PROTOCOL

The data sets considered in this experiment are listed in Table 4 with the corresponding sources and are all scaled in the $[0, 1]$ interval. They range from a training set cardinality of about 50k points to more than one million, whereas the dimensionality is not high (always under 60) with separated test sets. In order to select the parameters a 10-fold CV procedure is performed in the training set (apart from FaLK-SVMI) choosing the values in the following sets: $C \in \{2^{-2}, 2^{-1}, \dots, 2^9, 2^{10}\}$, $\sigma \in \{2^{-15}, 2^{-14}, \dots, 2^4, 2^5\}$, k for FaLK-SVM in $\{250, 500, 1000, 2000, 4000, 8000\}$ with $k' = k/2$, and k for FkNN SVM in $\{1, 3, 5, 9, 15, 21, 31, 51, 71, 101, 151\}$. FaLK-SVM does not necessarily test all values for k because if the maximum empirical accuracy is found for a specific value of k , for example $k = 500$, and for the following value, in this case $k = 1000$, the maximum is lower, the remaining higher values of k are not tested. Due to the computational resources necessary

data set	FkNN		LibSVM		FaLK-SVM		FaLK-SVMc		FaLK-SVMl
	10f-CV	test	10f-CV	test	10f-CV	test	10f-CV	test	test
ijcnn1	97.37	96.64	98.99	97.98	99.04	98.04	98.96	97.98	98.03
cov-type	91.73	91.99	92.60	92.83	92.68	92.89	92.44	92.60	92.84
census-inc	94.53	94.52	95.14	95.13	95.07	95.07	95.00	94.99	94.99
cod-rna	95.88	96.25	97.18	97.17	97.19	97.23	97.06	97.09	97.29
intr-det	99.74	92.04	99.89	91.77	99.74	91.97	99.69	92.01	91.91
2-spirals	88.43	88.43	85.18	85.29	88.42	88.47	88.29	88.45	88.30
ndcc	85.47	84.99	86.66	86.21	86.63	86.29	86.33	85.93	86.24
checker-b	94.31	94.08	94.46	94.21	94.46	94.21	94.45	94.19	94.23
test acc.	4.25		3.25		1.63		3.38		2.50
mean rank	4.25		3.25		1.63		3.38		2.50

Table 5: Empirical (using 10-fold CV) and generalisation accuracies of FkNN, LibSVM, FaLK-SVM, FaLK-SVMc and FaLK-SVMl on the 8 large data sets of Experiment 2. The best generalisation accuracy for each data set is highlighted in bold. The last line reports the mean rank of each method among the 8 data sets.

for performing model selection, especially for LibSVM, we performed the cross-validation runs on a Linux-based TORQUE cluster with 20 nodes. For FaLK-SVMl the local model selection is performed on 10 local models, $C \in \{2^0, 2^2, 2^4, 2^6\}$, $k \in \{500, 1000, 2000, 4000\}$, σ locally estimated with the 1st, 10th, 50th or 90th percentile of the distribution of the distances.

5.2.2 RESULTS AND DISCUSSION

Table 5 reports the generalisation accuracies of the analysed classifiers. Looking at the mean ranks, we can see that FaLK-SVM is the most accurate (it achieves the best results in half of the data sets), followed by FaLK-SVMl. LibSVM and FaLK-SVMc seem to perform very similar but little worse than FaLK-SVM and FaLK-SVMl. Not surprisingly, FkNN performs poorly in almost all the data sets, except for the intr-det data set in which it achieves the best result. According to the Wilcoxon Signed Rank Test (Wilcoxon, 1945; Demšar, 2006) FaLK-SVM is significantly more accurate than LibSVM, whereas, excluding FkNN, no other significant differences are detected. Apart for the intr-det data set that has slightly different distribution in the training and testing sets (some types of network attacks are present in the test set only), the best empirical accuracies are always very similar to the generalisation accuracies meaning that all techniques avoid over-fitting.

Table 6 reports the training times together with the speed-ups of FaLK-SVM, FaLK-SVMc and FaLK-SVMl with respect to LibSVM. We can notice that the speed-ups achieved by FaLK-SVM and FaLK-SVMc are always greater than 4.7, and in the majority of the cases they are at least one order of magnitude bigger than LibSVM. Generally, FaLK-SVMc turns out to be faster than FaLK-SVM although the two classifiers implement the same training algorithm. This happens because the model selection chooses for FaLK-SVMc a lower value of k with respect to FaLK-SVM. In fact, FaLK-SVMc is less accurate than FaLK-SVM in choosing the nearest model for a testing point, and this causes an higher value of the $E[\lambda]$ constant that increases the risk of FaLK-SVMc with respect to FaLK-SVM (see Equation 14 and Equation 15). So using a lower k (and thus a lower k') tends to have more

data set	LibSVM	FaLK-SVM		FaLK-SVMc		FaLK-SVMI	
	training time (s)	training time (s)	speed-up on LibSVM	training time (s)	speed-up on LibSVM	train. time with l.m.s.	speed-up on LibSVM
ijcnn1	102	15	6.8	15	6.8	1850	0.1
cov-type	8362	88	95.0	38	220.1	1214	6.9
census-inc	13541	6047	4.7	2391	5.7	10271	1.3
cod-rna	9777	395	24.8	225	43.5	579	16.9
intr-det	5262	286	18.4	284	18.5	450	11.7
2-spirals	4043	188	21.5	81	49.9	3442	1.2
ndcc	1487	302	4.9	92	16.2	4609	0.3
checker-b	6047	334	18.1	366	16.5	1374	4.4

Table 6: Training times for Experiment 2 of LibSVM, FaLK-SVM, FaLK-SVMc and FaLK-SVMI and the speed-ups of the three local methods with respect to LibSVM. The best training time for each data set is highlighted in bold.

data set	LibSVM	FaLK-SVM		FaLK-SVMc		FaLK-SVMI	
	testing time (s)	testing time (s)	speed-up on LibSVM	testing time (s)	speed-up on LibSVM	testing time (s)	speed-up on LibSVM
ijcnn1	43	32	1.3	5	8.6	36	1.2
cov-type	2795	202	13.8	73	38.3	191	14.6
census-inc	597	1347	0.4	58	10.3	1328	0.4
cod-rna	396	261	1.5	58	6.8	259	1.5
intr-det	192	146	1.3	76	2.5	149	1.3
2-spirals	957	10	95.7	5	191.4	18	53.2
ndcc	148	61	2.4	7	21.1	61	2.4
checker-b	167	10	16.7	7	23.9	7	23.9

Table 7: Testing times for Experiment 2 of LibSVM, FaLK-SVM, FaLK-SVMc and FaLK-SVMI and the speed-ups of the three local methods with respect to LibSVM. The best testing time for each data set is highlighted in bold.

models in the proximity of the testing point making the choice less problematic. FaLK-SVMI is sometimes slower than LibSVM, but we have to consider that FaLK-SVMI includes model selection, whereas for the other methods the time needed by model selection is not considered in the training time, so, practically speaking, FaLK-SVMI is the fastest method if the optimal parameters are not *a priori* known.

The testing times required by the analysed methods are reported in Table 7. As expected FaLK-SVMc is the fastest among all methods with speed-up over LibSVM ranging from more than 2 to almost 200. FaLK-SVM and FaLK-SVMI are also generally faster than LibSVM with only one case in which the testing time is about two times slower.

This experiment shows that for 8 non high-dimensional data sets, our approach outperforms a state-of-the-art accurate SVM solver both in terms of generalisation accuracies and computational

performances. Although we have an additional parameter to tune (k), FaLK-SVM and FaLK-SVMc are faster enough to maintain the performance advantages over LibSVM also for model selection (we choose k in a small set of values). Moreover, with FaLK-SVMl we addressed the problem of model selection with a specific approach to set the parameters; FaLK-SVMl outperforms LibSVM in generalisation accuracy, and the time it needs for both internal model selection and training is at least comparable (faster in 7 cases on a total of 8) with the time LibSVM needs for the training only.

5.3 Experiment 3: Comparison of Scalability Performances of FaLK-SVM, FaLK-SVMc, FaLK-SVMl, LibSVM and Approximated SVM Solvers

In this experiment we test the scalability performances of our techniques (FaLK-SVM, FaLK-SVMc, FaLK-SVMl) on training sets with increasing sizes using the RBF kernel against several other techniques. The techniques taken into account are LibSVM, the approximated SVM solvers called CVM, LASVM, USVM, BVM, CPSP (Section 2.3), SVM-bagging with fixed dimension of the sub-sampled training sets (SVM-B) and SVM-bagging with fixed proportion of the sub-sampled training sets with respect to the whole training set (SVM-Bs). Although we apply all the classifiers with the same protocol on the same data sets, we report, for clearness, the results in two parts: the comparison of FaLK-SVM with LibSVM and the approximated SVM solvers in Section 5.3.2, the comparison of FaLK-SVM with its variants FaLK-SVMc and FaLK-SVMl in Section 5.3.3.

5.3.1 EXPERIMENTAL PROTOCOL

We consider here the data sets of Table 4 for which we can further enlarge the training set size. The data sets for which we can add sets of new training examples are the cov-type data set (full training set of 500k points) and the three artificial data sets named 2-spirals, ndcc and checker-b (up to 3 million points). For cov-type the testing set is reduced to 50k examples (the other examples are added to the training set) so the accuracy results are not directly comparable to the previous experiment.

The model selection for all the classifiers (with the exception of FaLK-SVMl that performs internally a local model selection) is performed on the smallest training set only, using the chosen parameter for all the higher training set sizes. This is necessary, especially for LibSVM and approximated SVM solvers, for computational reasons. For LibSVM, BVM, CVM, USVM (with the convex concave procedure) and CPSP, we performed cross validation for C and σ using the same setting of the previous experiment. The default threshold value ϵ for the stopping criteria are maintained: 10^{-3} for LibSVM, FaLK-SVM, LASVM and 10^{-1} for CPSP while CVM and BVM automatically choose the value of ϵ based on the data at each application. We set the same size of the kernel cache (100M) for all the methods. The maximum number of core vectors for CVM and BVM is 50000 (the default value), the maximum number of basis vectors for CPSP is set to 1000. For SVM-B and SVM-Bs we need to set respectively the size and the proportion of the sub-sampled training sets and the parameters of the SVMs. For SVM-B the size of the sub-sampled training sets is equal to the 5% of the original training sets of each data set (namely 100000 for cov-type, 2-spirals and ndcc and 300000 for checker-b), whereas SVM-Bs maintains the same sampling rate (5%) for all the applications and so the cardinality of the sub-sampled training sets increases with the cardinality of the training set. Both SVM-B and SVM-Bs train 101 SVMs and the prediction is performed using the majority voting. The value of 101 is chosen because it is a sufficient high number for allowing good accuracies and it is an odd number preventing possible ties in the majority voting. The parameters of the

SVM for SVM-B and SVM-Bs (using LibSVM) are chosen using model selection with the same grids of parameters used for LibSVM. We also tested FaLK-SVM using the same setting of the previous experiment. Each algorithm is tested for training set sizes requiring no more than 100000 seconds (more than 27 hours) for training.

Since the authors of BVM (Tsang et al., 2005) and CVM (Tsang et al., 2007) declared the Linux implementation of their techniques deprecated (see the authors reply to Loosli and Canu (2007) available on BVM webpage), we use the Windows executables on a Intel Pentium D Dual Core CPU 3.40GHz with 2Gb of RAM running Windows XP instead of the AMD Athlon 64 X2 Dual Core Processor 5000+, 2600MHz, with 3.56Gb of RAM with Linux operating system used for all the other classifiers. Because of the use of different operating systems and hardware for BVM and CVM, their running times should not be directly compared to the others. However, the comparison is justified by preliminary tests that showed that the Linux version of BVM on the AMD Athlon machine and the Windows version of BVM on the Intel Pentium machine have similar running times.

5.3.2 RESULTS AND DISCUSSION: FALK-SVM VS LIBSVM AND APPROXIMATED SVM SOLVERS

Figure 3 shows the generalisation accuracies of the methods at increasing training set sizes. Some methods do not appear in the figures due to low generalisation results or computational difficulties that cause abnormal terminations of the algorithms, and some accuracy results for large training set sizes are not present due to the excessive computational time required for training (more than 100000 seconds). We can observe that it is very important to use as many points as possible in order to increase the accuracies for the cov-type and ndcc data sets. The same consideration can be done for the 2-spirals data, although FaLK-SVM already starts from very high accuracies and the increment is limited, while for the checker-b data set the increment of the accuracies is negligible for almost all the methods. For the checker-b data set, the enlarging of the training set is not motivated from the accuracy viewpoint, but we still use it as a benchmark for the computational performances.

Comparing the generalisation accuracies of Figure 3 among the tested methods, we can see that FaLK-SVM is almost always on top for each of the four data sets. In this experiment as well as in the previous ones, we set $k' = k/2$; lower values for k' would probably allow FaLK-SVM to achieve higher accuracy results (although with worse computational performances). However, even if the choice of the k' parameter can be non-optimal, we decided to avoid the model selection for k' since the results are already satisfactory. The methods that seem to give results comparable with FaLK-SVM (apart from the 2-spirals data set) are LibSVM and USVM and they are able, in few cases, to slightly improve the FaLK-SVM results (LibSVM for 2 training set sizes for cov-type and checker-b, USVM for 2 training set sizes for cov-type and checker-b and 1 for ndcc). The bagging techniques give high accuracies only for the checker-b data set; this is not surprising because we already noticed that for the checker-b problem the use of large data sets is not required and thus subsampling-based methods like SVM-B and SVM-Bs are competitive. The results of the online and active learning approach of LASVM are slightly lower than FaLK-SVM, LibSVM and USVM. CPSP gives acceptable results in only one case, and for the 2-spirals and checker-b data sets it suffers from numerical problems possibly due to the scaling of the features in the $[0, 1]$ interval. Enlarging the maximum number of basis functions for CPSP gives higher accuracies but the computational time needed to build the models is too high. The results we achieve here for LibSVM and LASVM on the

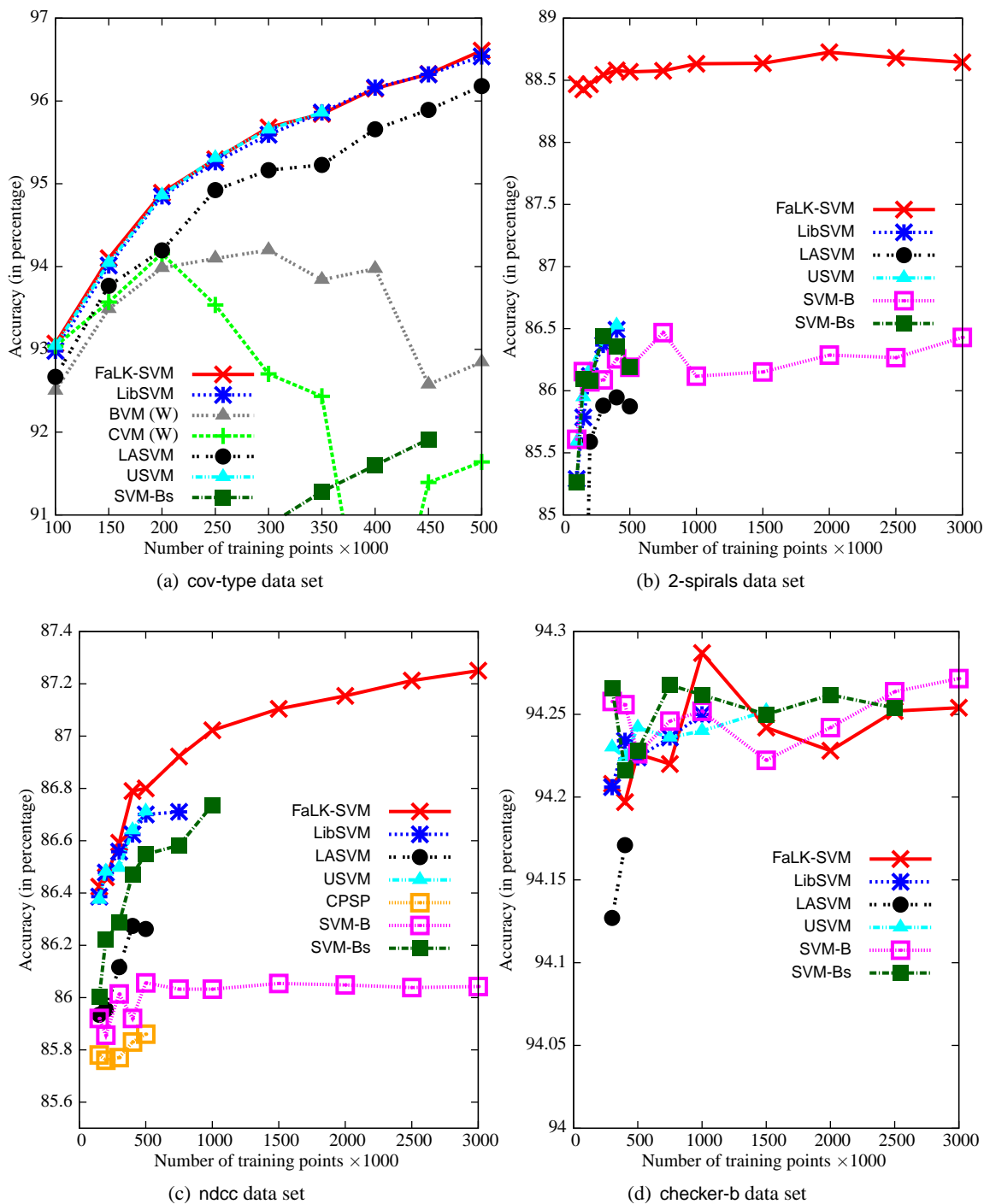


Figure 3: Generalisation accuracies of FaLK-SVM, LibSVM, BVM, CVM, LASVM, USVM, CPSP, SVM-B and SVM-Bs on the cov-type, 2-spirals, ndcc and checker-b data sets with increasing training set sizes (Experiment 3). Some accuracies are missing due to the excessive computational requirements (more than 100000 seconds for training) of the corresponding method for large training set sizes.

cov-type data sets are a little higher than the results in Bordes et al. (2005) (about 1% better both for 100k and 500k training set sizes), and we believe that this is due to the model selection approach we used here that is performed with an exhaustive cross-validation grid search for C and σ . As we can notice in Figure 3, we observed stability problems for CVM and BVM, even if we used the Windows binaries as suggested by the authors.

The training computational performances shown in Figure 4 highlight that FaLK-SVM is always much faster than the alternative techniques that are competitive from the accuracy viewpoint. In fact, although CVM, BVM and SVM-B show good scalability performances and in few cases they overcome the performances of FaLK-SVM, we noticed from Figure 3 that their generalisation abilities are poor. The scaling behaviours of LibSVM, LASVM and USVM are very similar (among the three methods LibSVM is the fastest for ndcc, LASVM is the fastest for 2-spirals and USVM is the fastest for checker-b) but substantially worse than FaLK-SVM one (FaLK-SVM is always at least one order of magnitude faster with speed-ups increasing with the training set sizes). SVM-Bs is slightly faster than LibSVM but the scalability behaviour is very similar. The methods that achieve acceptable accuracy results on the smallest training set size (i.e., LibSVM, LASVM, USVM) are not applicable when the number of training examples increases sensibly because of poor computational scalability performance; this is evident for the 2-spirals, ndcc and checker-b data sets in which the training times of LibSVM, LASVM, USVM exceed 100000 seconds as soon as the training set cardinality approaches one million (the only exception is USVM that is applicable on 1.5 training examples of the checker-b data set). On the contrary, FaLK-SVM processes data sets of 3 millions examples in the order of minutes or few hours. An experiment comparing LibSVM and LASVM on the cov-type data set with conclusions similar to ours is reported by Bordes et al. (2005) in which however LASVM is about a third faster than LibSVM whereas here LibSVM slightly overcomes LASVM; this is probably due to the fact that for LASVM the only available implementation is the original one by Bordes et al. (2005) whereas LibSVM is frequently updated and improved. Finally, CPSP performs slightly better than LibSVM, LASVM and USVM.

The computational performances of the prediction phase are reported in Figure 5. Also in this case the performance of FaLK-SVM is excellent: only CPSP and CVM are faster in 2 data sets than FaLK-SVM, but their corresponding generalisation accuracies are low. As expected, CPSP achieves very fast predictions because it limits the number of basis function to 1000 and thus for each testing points no more than 1000 kernel functions are computed. LibSVM, LASVM and USVM achieve similar results also in testing performances and, apart from small training sets for the ndcc data set, they are at least one order of magnitude slower than FaLK-SVM and the difference grows for large training set sizes. The slowest approach is SVM-Bs and this is due to the high number of models that need to be evaluated for each testing point and to the fact that the size of the models increases with the training set size. Also SVM-B has rather high prediction times but, since the size of the models is almost constant, also the performances at increasing training sizes are constant. It can be argued that the number of models used for bagging can be lowered to obtain faster prediction times; however, if we want to achieve the computational performances of FaLK-SVM we need to use no more than 20 models (in the worst case) and this seriously affects the prediction accuracies that are already much lower than FaLK-SVM ones.

The overall conclusion we can draw about the scalability of the proposed techniques is that, at least for these 4 non high-dimensional data sets, FaLK-SVM is substantially better than the state-of-the-art kernel methods for classification, and this is achieved without affecting the accuracy performances that showed to be always at least as good as the best alternative technique. Apart

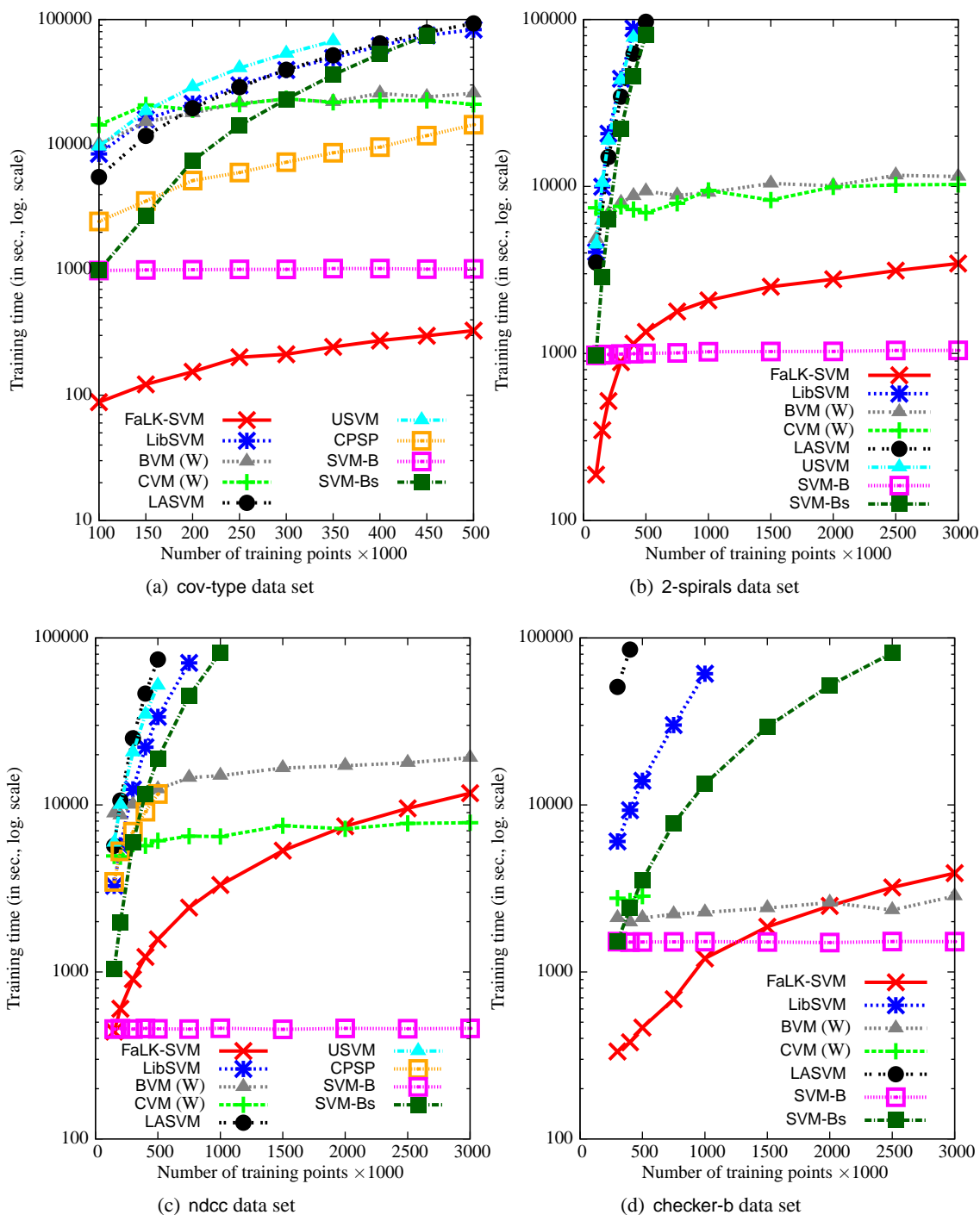


Figure 4: Training times of FaLK-SVM, LibSVM, BVM, CVM, LASVM, USVM, CPSP, SVM-B and SVM-Bs on the cov-type, 2-spirals, ndcc and checker-b data sets with increasing training set sizes (Experiment 3). The times (in seconds) are reported in logarithmic scale.

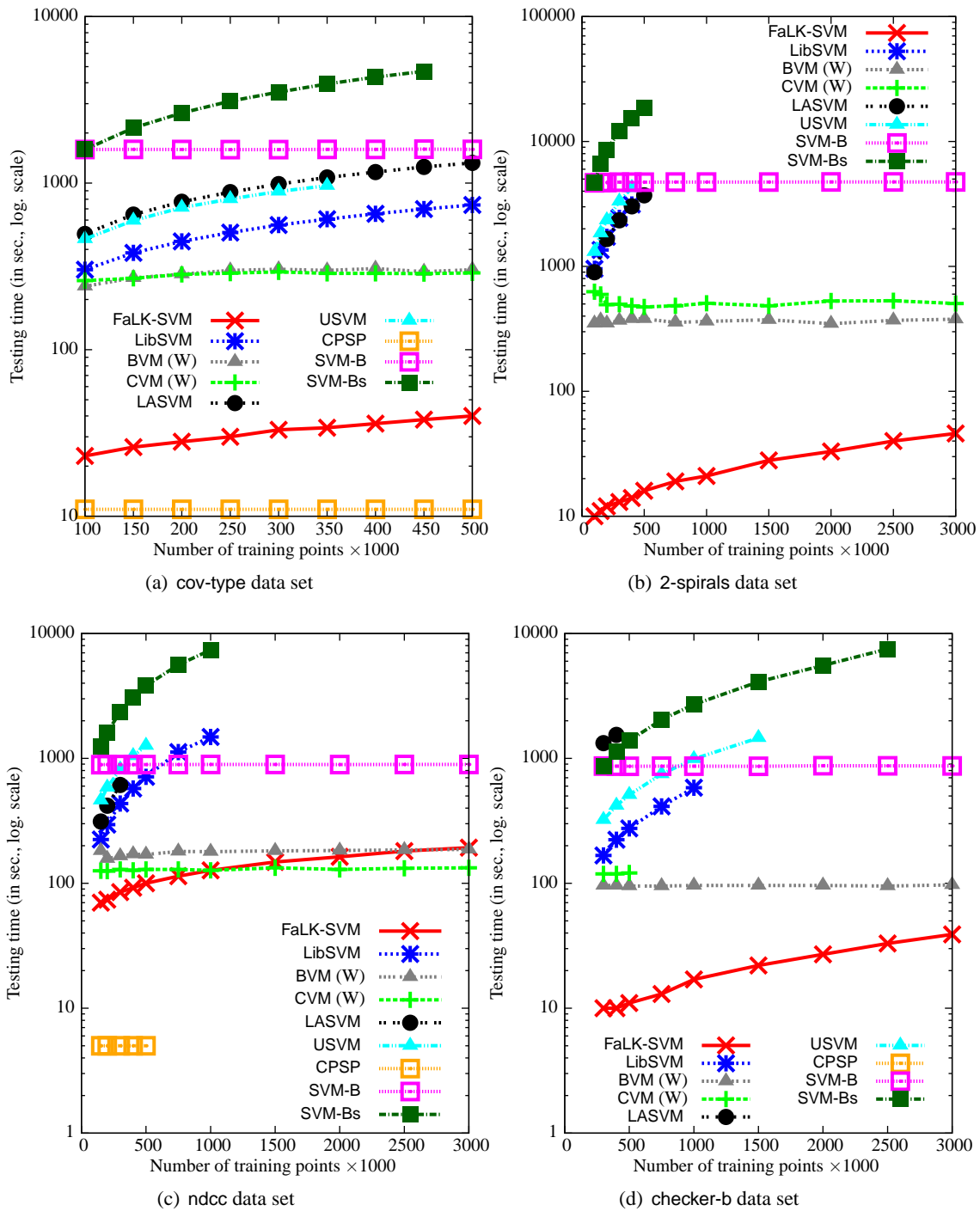


Figure 5: Testing times of FaLK-SVM, LibSVM, BVM, CVM, LASVM, USVM, CPSP, SVM-B and SVM-Bs on the data sets of Experiment 3 with increasing training set sizes. The times (in seconds) are reported in logarithmic scale. Some testing times are missing due to the excessive computational requirements (more than 100000 seconds for training) of the corresponding method for large training set sizes.

for LibSVM (and consequently for SVM-B and SVM-Bs), we have to say that the available code of the other tested techniques has not been recently updated and for this reason it is possible to argue that higher performances with more optimised implementations of the tested approaches could be reached. It is also necessary to underline that in literature LASVM, USVM, CPSP, BVM and CVM have been prevalently tested on data sets with high dimensionality or, apart for cov-type, on data sets not requiring highly non-linear decision functions. The approximated non-linear SVM solvers and bagging approaches we tested could be indicated for data in which the linear kernel is not the optimal choice, but, at the same time, the decision function can be accurately reconstructed with a reduced amount of information (number of examples, support vectors or basis functions).

5.3.3 RESULTS AND DISCUSSION: COMPARISON BETWEEN FaLK-SVM, FaLK-SVMc AND FaLK-SVMl

Figure 6 reports the comparison of the generalisation accuracies of FaLK-SVM, FaLK-SVMc and FaLK-SVMl at increasing training set size. The computational performances for the training phase are reported in Figure 7, and for the testing phase in Figure 8.

From the accuracy viewpoint, we can notice that, as expected, FaLK-SVM is almost always slightly more accurate than FaLK-SVMc. FaLK-SVMl, apart from checker-b, is less accurate than FaLK-SVM for the smaller training set sizes, and this is due to the fact that FaLK-SVM performs a full grid search for model selection whereas FaLK-SVMl adopts the very fast local model selection approach. However, FaLK-SVMl rivals FaLK-SVM as the training set sizes increases. This is reasonable because FaLK-SVM uses for all the training set sizes the parameters found for the smaller training sets, and the best cross-validated parameters can differ for sub-sampled sets with different cardinality. For example, as the number of training points increases, the radius of the local neighbourhoods decreases if we maintain the same k and k' values, and the original value for the width parameter of the RBF kernel can no longer be the optimal one. For this reason, in the case of cov-type and ndcc data sets, FaLK-SVMl achieves higher accuracies than FaLK-SVM for the largest training sets. FaLK-SVMl shows a slightly higher accuracy variability than FaLK-SVM and FaLK-SVMc at different training set sizes; this is an empirical confirmation that the parameters selected by FaLK-SVMl can be less stable than the parameters selected with standard cross validation, but the phenomenon seems to be acceptable if not negligible.

The training computational performances of Figure 7 confirm (as already discussed in Section 5.2.2) that, although FaLK-SVM and FaLK-SVMc make use of the same training algorithm, the model selection procedure selects lower values of k for FaLK-SVMc, thus assuring faster training times than FaLK-SVM. The speed-ups of FaLK-SVMc with respect to FaLK-SVM are however never higher than one order of magnitude. For FaLK-SVMl we can notice a somehow irregular behaviour for increasing dimensions of the training set and this is due to the different values of the neighbourhood, kernel and regularisation parameters it chooses during the internal fast local model selection phase. In some cases FaLK-SVMl is significantly slower than FaLK-SVM. However, the training times for FaLK-SVMl include the model selection procedure whereas for FaLK-SVM we consider only the training with the optimal parameters, so we can conclude that FaLK-SVMl is a good choice for huge training sets on which traditional model selection becomes intractable.

The testing times reported in Figure 8 confirm that FaLK-SVMc is always faster than FaLK-SVM and FaLK-SVMl. In particular, we can notice that FaLK-SVMc at least halves the testing time of FaLK-SVM. FaLK-SVMl is computationally very similar to FaLK-SVM. This is not surprising because the

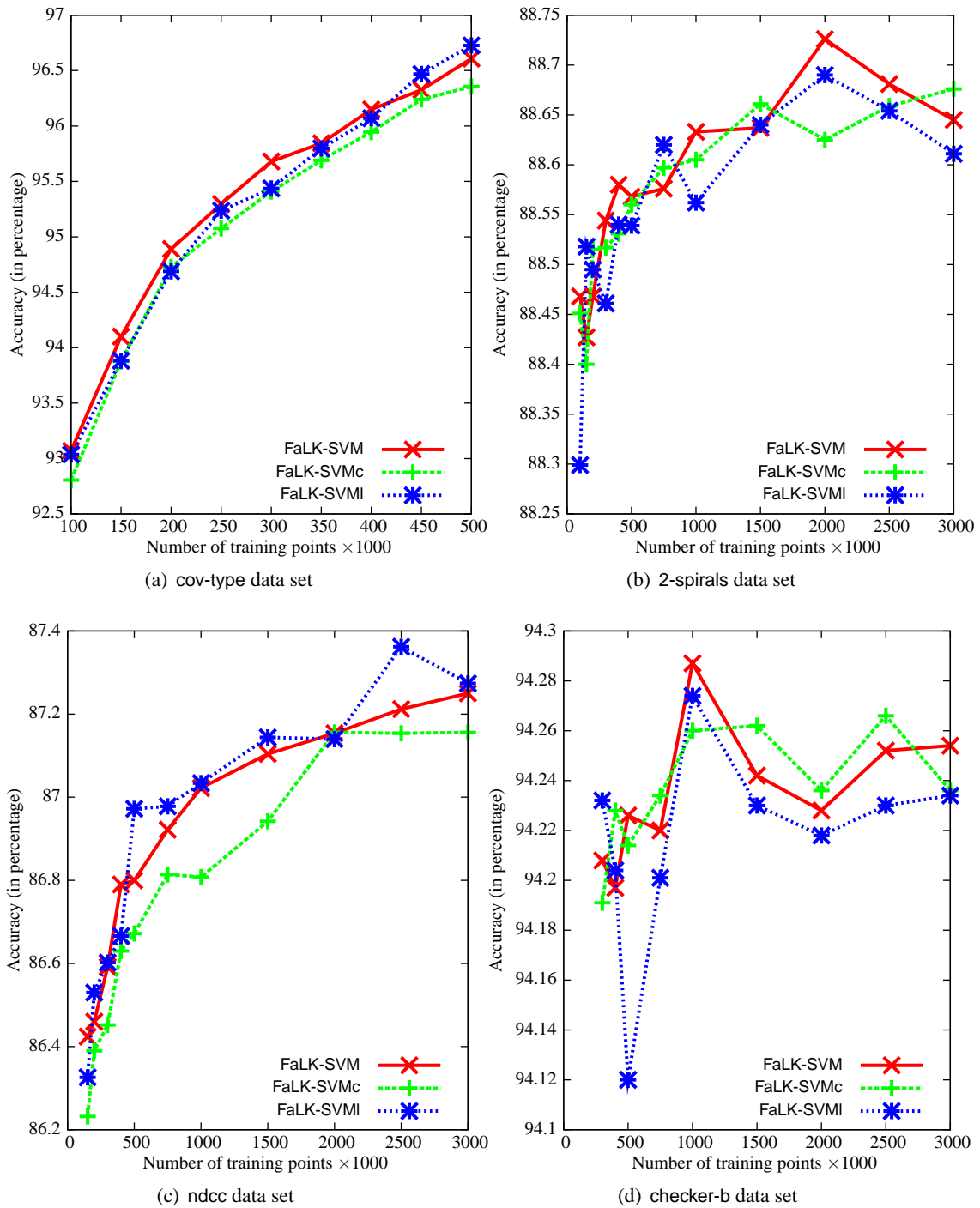


Figure 6: Generalisation accuracies of FaLK-SVM, FaLK-SVMc and FaLK-SVMl on the cov-type, 2-spirals, ndcc and checker-b data sets with increasing training set sizes (Experiment 3).

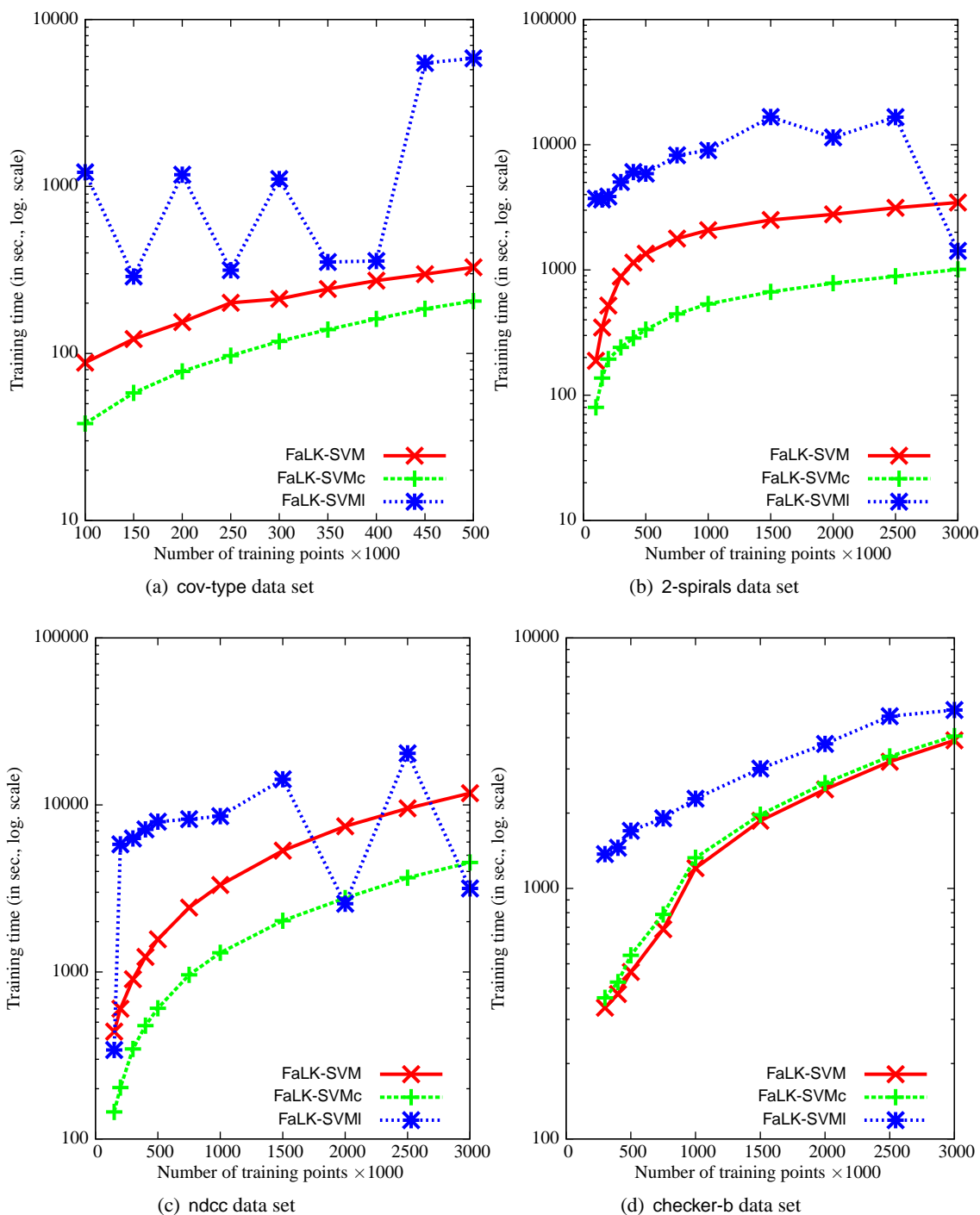


Figure 7: Training times of FaLK-SVM, FaLK-SVMc, and FaLK-SVMI on the cov-type, 2-spirals, ndcc and checker-b data sets with increasing training set sizes (Experiment 3). The times (in seconds) are reported in logarithmic scale.

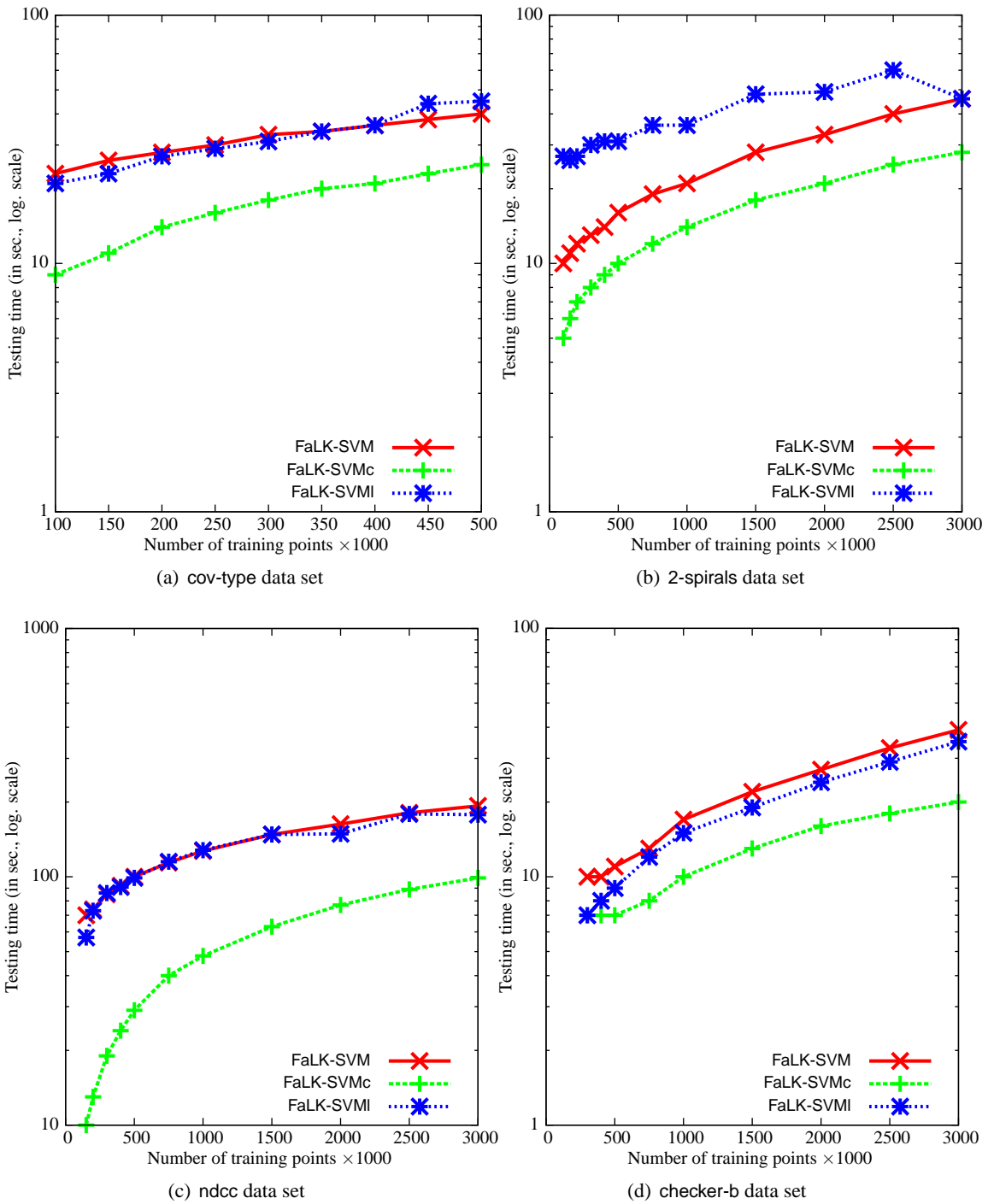


Figure 8: Testing times of FaLK-SVM, FaLK-SVMc, and FaLK-SVMI on the cov-type, 2-spirals, ndcc and checker-b data sets with increasing training set sizes (Experiment 3). The times (in seconds) are reported in logarithmic scale.

only difference between FaLK-SVM and FaLK-SVMl regards the model selection but both classifiers need, during testing, to perform a nearest neighbour search of the query points among all training examples, differently from FaLK-SVMc that performs the nearest neighbour search only among the centres of the local models.

The FaLK-SVMl results permit us to discuss the sensitivity of the local kernel machine approaches with respect to the neighbourhood parameter k . The local model selection of FaLK-SVMl selects different values of k at each application and the k selected by FaLK-SVMl is very often different from the value selected by FaLK-SVM with standard model selection. However, the accuracy results are not very dependent on this as we notice in Figure 6 in which the accuracy variations of FaLK-SVMl as the training sets increase are rather smooth and the accuracies are very similar to FaLK-SVM. So the sensitivity of the accuracies with respect to k for local kernel machines seems to be low at least for large data sets. From the computational viewpoint, instead, we know from the complexity analysis that the advantages of local kernel machines are effective as long as k is substantially lower than N . In our experiments the value of k selected by model selection is bounded to 8000; higher values, although not tested, may decrease the computational performances.

We can conclude that FaLK-SVM, FaLK-SVMc and FaLK-SVMl achieve similar accuracy and computational results. When the model selection for FaLK-SVM and FaLK-SVMc become computationally intractable, FaLK-SVMl is an option to efficiently perform model selection and thus obtain a lower overall training time. When very low testing times are required, FaLK-SVMc is preferable to FaLK-SVM at the price of a slightly lower generalisation accuracy.

6. Conclusions

In this work, we have introduced a new local kernel-based classifier, called FaLK-SVM, that is scalable for large non high-dimensional data. The approach is developed starting from the theory of local learning algorithms and in particular from the Local SVM classifier, called k NNSVM. Various strategies are introduced to overcome the computational problems of k NNSVM and to switch from a completely lazy-learning setting to a eager learning setting with efficient predictions. Learning and complexity bounds for FaLK-SVM are favourable if compared with the SVM ones. FaLK-SVM has, in fact, a training time complexity which is sub-quadratic in the training set size, and a prediction time complexity which is logarithmic. A novel approach for model selection, again based on locality, is introduced obtaining the FaLK-SVMl classifier which substantially unburdens the model selection strategies based on cross-validation. Another variant of the algorithm, called FaLK-SVMc, permits to simplify the prediction phase. We thus showed that locality can be used to develop computationally efficient classifiers.

We carried out an extensive empirical evaluation of the introduced approaches showing that, for large classification problems requiring non linear decision functions our FaLK-SVM algorithm is much faster and accurate than traditional and approximated SVM solvers. In fact, (i) FaLK-SVM achieves very good accuracy results because it considers all the points without locally under-fitting the data and (ii) FaLK-SVM is very fast and scalable because the cardinality of the local problems can be maintained low. The variant called FaLK-SVMc further enhances testing speed at the price of a little accuracy loss, and the other variant, called FaLK-SVMl, decreases the overall training time.

In general, we have showed that locality can be the key not only for obtaining accurate classifiers, but also for effectively speeding-up kernel-based algorithms.

Further developments of the approach include a dimensionality reduction preprocessing step in order to attack also high-dimensional problems, the application of local classifiers different from SVM, and a distributed parallel version. Also the determination of the critical value of the intrinsic dimensionality (rather than the number of features) above which the local approaches are not effective is still an open question and the answer should be very data-dependent.

References

- David W. Aha. *Lazy Learning*. Kluwer Academic Publishers Norwell, MA, USA, 1997.
- Erin L. Allwein, Robert E. Schapire, and Yoram Singer. Reducing multiclass to binary : A unifying approach for margin classifiers. *Journal of Machine Learning Research*, 1:113–141, 2000.
- Arthur Asuncion and David Newman. UCI machine learning repository, 2007. URL <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- Christopher G. Atkeson, Andrew W. Moore, and Stefan Schaal. Locally weighted learning. *Artificial Intelligence Review*, 11(1-5):11–73, 1997.
- Yoshua Bengio, Olivier Delalleau, and Nicolas Le Roux. The curse of highly variable functions for local kernel machines. In *Advances in Neural Information Processing Systems*, volume 18, 2005.
- Alina Beygelzimer, Sham Kakade, and John Langford. Cover trees for nearest neighbor. In *Twenty-third International Conference on Machine Learning (ICML 06)*, pages 97–104, New York, NY, USA, 2006. ACM.
- Jock A Blackard and Denis J. Dean. Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and Electronics in Agriculture*, 24:131–151, 1999.
- Enrico Blanzieri and Anton Bryl. Evaluation of the highest probability SVM nearest neighbor classifier with variable relative error cost. In *CEAS 2007*, Mountain View, California, 2007.
- Enrico Blanzieri and Farid Melgani. An adaptive SVM nearest neighbor classifier for remotely sensed imagery. In *IEEE International Conference on Geoscience and Remote Sensing Symposium (IGARSS 06)*, pages 3931–3934, 2006.
- Enrico Blanzieri and Farid Melgani. Nearest neighbor classification of remote sensing images with the maximal margin principle. *IEEE Transactions on Geoscience and Remote Sensing*, 46(6):1804–1811, 2008.
- Antoine Bordes and Léon Bottou. The Huller: A simple and efficient online SVM. In *Machine Learning: ECML 2005*, Lecture Notes in Artificial Intelligence, LNAI 3720, pages 505–512. Springer Verlag, 2005.
- Antoine Bordes, Seyda Ertekin, Jason Weston, and Léon Bottou. Fast kernel classifiers with online and active learning. *Journal of Machine Learning Research*, 6:1579–1619, 2005.
- Antoine Bordes, Léon Bottou, and Patrick Gallinari. SGD-QN: Careful quasi-Newton stochastic gradient descent. *Journal of Machine Learning Research*, 10:1737–1754, July 2009.

- Léon Bottou and Vladimir N. Vapnik. Local learning algorithms. *Neural Computation*, 4(6):888–900, 1992.
- Léon Bottou, Corinna Cortes, John S. Denker, Harris Drucker, Isabelle Guyon, Lawrence D. Jackel, Yann Le Cun, Urs A. Muller, Eduard Säckinger, Patrice Simard, and Vladimir Vapnik. Comparison of classifier methods: A case study in handwritten digit recognition. In *Twelfth IAPR International Conference on Pattern Recognition, Conference B: Computer Vision & Image Processing*, volume 2, pages 77–82. IEEE, 1994.
- Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, August 1996.
- David Broomhead and David Lowe. Multivariable functional interpolation and adaptive networks. *Complex Systems*, 2:321–355, 1988.
- Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: A Library for Support Vector Machines*, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Kai-Wei Chang, Cho-Jui Hsieh, and Chih-Jen Lin. Coordinate descent method for large-scale l_2 -loss linear support vector machines. *Journal of Machine Learning Research*, 9:1369–1398, 2008.
- Qun Chang, Qingcai Chen, and Xiaolong Wang. Scaling Gaussian RBF kernel width to improve SVM classification. *International Conference on Neural Networks and Brain, 2005. (ICNN&B 05)*, 1:19–22, 2005.
- Long Chen. New analysis of the sphere covering problems and optimal polytope approximation of convex bodies. *Journal of Approximation Theory*, 133(1):134, 2005.
- Haibin Cheng, Pang-Ning Tan, and Rong Jin. Localized support vector machine and its efficient algorithm. *SIAM International Conference on Data Mining*, 2007.
- Vasek Chvatal. A greedy heuristic for the set-covering problem. *Mathematics of Operations Research*, pages 233–235, 1979.
- Kenneth L. Clarkson. Nearest neighbor queries in metric spaces. In *Twenty-ninth Annual ACM Symposium on Theory of computing (STOC 97)*, pages 609–617, New York, NY, USA, 1997. ACM.
- Michael Collins, Amir Globerson, Terry Koo, Xavier Carreras, and Peter L. Bartlett. Exponentiated gradient algorithms for conditional random fields and max-margin markov networks. *Journal of Machine Learning Research*, 9:1775–1822, 2008.
- Ronan Collobert, Fabian Sinz, Jason Weston, and Léon Bottou. Trading convexity for scalability. In *Twenty-third International Conference on Machine Learning (ICML 06)*, pages 201–208, New York, NY, USA, 2006. ACM.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, 2006.

- Jian-xiong Dong. Fast SVM training algorithm with decomposition on very large data sets. *IEEE Transaction Pattern Analysis and Machine Intelligence*, 27(4):603–618, 2005. Senior Member-Krzyzak, Adam and Fellow-Suen, Ching Y.
- Rong-En Fan, Pai-Hsuen Chen, and Chih-Jen Lin. Working set selection using second order information for training support vector machines. *The Journal of Machine Learning Research*, 6: 1889–1918, 2005.
- Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. LIBLINEAR: A library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874, 2008.
- Milton Friedman. A comparison of alternative tests of significance for the problem of m rankings. *The Annals of Mathematical Statistics*, 11(1):86–92, 1940.
- Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. WH Freeman & Co. New York, NY, USA, 1979.
- Seth Hettich and Stephen D. Bay. The UCI KDD archive, 1999. URL <http://kdd.ics.uci.edu>.
- Cho-Jui Hsieh, Kai-Wei Chang, Chih-Jen Lin, S. Sathiya Keerthi, and S. Sundararajan. A dual coordinate descent method for large-scale linear SVM. In *Twenty-fifth International Conference on Machine Learning (ICML 08)*, pages 408–415, New York, NY, USA, 2008. ACM.
- Chih-Wei Hsu and Chih-Jen Lin. A comparison of methods for multiclass support vector machines. *IEEE Transactions on Neural Networks*, 13(2):415–425, 2002.
- Thorsten Joachims. Making large-scale support vector machine learning practical. *Advances in kernel methods: support vector learning*, pages 169–184, 1999.
- Thorsten Joachims. Training linear SVMs in linear time. In *Twelveth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 217–226. ACM New York, NY, USA, 2006.
- Thorsten Joachims and Chun-Nam Yu. Sparse kernel SVMs via cutting-plane training. *Machine Learning*, 2009.
- Thorsten Joachims, Thomas Finley, and Chun-Nam John Yu. Cutting-plane training of structural SVMs. *Machine Learning*, 77(1):27–59, 2009.
- Michael J. Kearns and Umesh V. Vazirani. *An Introduction to Computational Learning Theory*. MIT Press Cambridge, MA, USA, 1994.
- Sathiya Keerthi and Dennis DeCoste. A modified finite Newton method for fast solution of large scale linear SVMs. *Journal of Machine Learning Research*, 6:341–361, 2005.
- Sathiya Keerthi, Olivier Chapelle, and Dennis DeCoste. Building support vector machines with reduced classifier complexity. *Journal of Machine Learning Research*, 7:1493–1515, 2006.

- Stefan Knerr, Leon Personnaz, and Gerard Dreyfus. Single-layer learning revisited: A stepwise procedure for building and training a neural network. *Optimization Methods and Software*, 1: 23–34, 1990.
- Robert Krauthgamer and James R. Lee. Navigating nets: simple algorithms for proximity search. In *Fifteenth Annual ACM-SIAM Symposium on Discrete algorithms (SODA 04)*, pages 798–807, Philadelphia, PA, USA, 2004. Society for Industrial and Applied Mathematics.
- Ulrich H.-G. Kressel. Pairwise classification and support vector machines. *Advances in Kernel Methods: Support Vector Learning*, pages 255–268, 1999.
- Yuh-jye Lee and Olvi L. Mangasarian. RSVM: Reduced support vector machines. In *First SIAM International Conference on Data Mining*, 2001.
- Chih-Jen Lin, Ruby C. Weng, and S. Sathiya Keerthi. Trust region Newton methods for large-scale logistic regression. In *Twenty-fourth International Conference on Machine learning (ICML 07)*, pages 561–568, New York, NY, USA, 2007. ACM.
- Gaëlle Loosli and Stéphane Canu. Comments on the “Core vector machines: Fast SVM training on very large data sets”. *Journal of Machine Learning Research*, 8:291–301, 2007.
- Olvi L. Mangasarian. A finite Newton method for classification. *Optimization Methods and Software*, 17(5):913–929, 2002.
- Mario Marchand and John Shawe-Taylor. The set covering machine. *The Journal of Machine Learning Research*, 3:723–746, 2003.
- Donald Michie, David J. Spiegelhalter, Charles C. Taylor, and John Campbell, editors. *Machine Learning, Neural and Statistical Classification*. Ellis Horwood, Upper Saddle River, NJ, USA, 1994. ISBN 0-13-106360-X.
- Paul Nemenyi. *Distribution-Free Multiple Comparisons*. PhD thesis, Princeton, 1963.
- John C. Platt, Nello Cristianini, and John Shawe-Taylor. Large margin DAGs for multiclass classification. *Advances in Neural Information Processing Systems*, 12(3):547–553, 2000.
- Bernard Schölkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press Cambridge, MA, USA, 2002.
- Nicola Segata. FaLKM-lib v1.0: A library for fast local kernel machines. Technical Report DISI-09-025, id 1613, DISI, University of Trento, Italy, 2009. Software available at <http://disi.unitn.it/~segata/FaLKM-lib>.
- Nicola Segata and Enrico Blanzieri. Empirical assessment of classification accuracy of Local SVM. In *Eighteenth Annual Belgian-Dutch Conference on Machine Learning (Benelearn 2009)*, pages 47–55, 2009a.
- Nicola Segata and Enrico Blanzieri. Operators for transforming kernels into quasi-local kernels that improve SVM accuracy. Technical Report DISI-09-042, id 1652, Tech. rep., DISI, University of Trento, 2009b.

- Nicola Segata and Enrico Blanzieri. Fast local support vector machines for large datasets. In *International Conference on Machine Learning and Data Mining (MLDM 2009)*, volume 5632 of *Lecture Notes in Computer Science*. Springer, 2009c.
- Nicola Segata, Enrico Blanzieri, and Pádraig Cunningham. A scalable noise reduction technique for large case-based systems. In L Ginty and D.C Wilson, editors, *Case-Based Reasoning Research and Development: 8th International Conference on Case-Based Reasoning (ICCB09)*, volume 09 of *Lecture Notes in Artificial Intelligence*, pages 755–758. Springer, 2009a.
- Nicola Segata, Enrico Blanzieri, Sarah Jane Delany, and Pádraig Cunningham. Noise reduction for instance-based learning with a local maximal margin approach. *Journal of Intelligent Information Systems*, 2009b. In Press.
- Shai Shalev-Shwartz, Yoram Singer, and Nathan Srebro. Pegasos: Primal estimated sub-gradient solver for SVM. In *Twenty-fourth International Conference on Machine Learning (ICML 07)*, pages 807–814, New York, NY, USA, 2007. ACM.
- Alexander J. Smola, SVN Vishwanathan, and Quoc V. Le. Bundle methods for machine learning. *Advances in Neural Information Processing Systems*, 20:1377–1384, 2008.
- Michael E. Thompson. NDCC: Normally distributed clustered datasets on cubes, 2006. www.cs.wisc.edu/dmi/svm/ndcc/.
- Ivor W. Tsang, James T. Kwok, and Pak-Ming Cheung. Core vector machines: Fast SVM training on very large data sets. *The Journal of Machine Learning Research*, 6:363–392, 2005.
- Ivor W. Tsang, Andras Kocsor, and James T. Kwok. Simpler core vector machines with enclosing balls. In *Twenty-fourth International Conference on Machine Learning (ICML 07)*, pages 911–918, New York, NY, USA, 2007. ACM.
- Andrew V. Uzilov, Joshua M. Keegan, and David H. Mathews. Detection of non-coding rnas on the basis of predicted secondary structure formation free energy change. *BMC Bioinformatics*, 7(1): 173, 2006.
- Vladimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 2000.
- Vladimir N. Vapnik and Léon Bottou. Local algorithms for pattern recognition and dependencies estimation. *Neural Computation*, 5(6):893–909, 1993.
- Jigang Wang, Predrag Neskovic, and N. Leon Cooper. A minimum sphere covering approach to pattern classification. *International Conference on Pattern Recognition*, 3:433–436, 2006.
- Xun-Kai Wei and Ying-Hong Li. Linear programming minimum sphere set covering for extreme learning machines. *Neurocomputing*, 71(4–6):570–575, 2008.
- Frank Wilcoxon. Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6):80–83, 1945.
- Tao Yang and Vojislav Kecman. Adaptive local hyperplane classification. *Neurocomputing*, 71 (13-15):3001–3004, 2008.

- Tao Yang and Vojislav Kecman. Adaptive local hyperplane algorithm for learning small medical data sets. *Expert Systems*, 26(4):355–359, 2009.
- Tao Yang and Vojislav Kecman. Face recognition with adaptive local hyperplane algorithm. *Pattern Analysis & Applications*, 13(1):79–83, 2010. ISSN 1433-7541.
- Alan L. Yuille and Anand Rangarajan. The concave-convex procedure. *Neural Computation*, 15(4): 915–936, 2003.
- Alon Zakai and Ya’acov Ritov. Consistency and localizability. *Journal of Machine Learning Research*, 10:827–856, 2009.
- Luca Zanni, Thomas Serafini, and Gaetano Zanghirati. Parallel software for training large scale support vector machines on multiprocessor systems. *Journal of Machine Learning Research*, 7: 1467–1492, 2006.
- Hao Zhang, Alexander C. Berg, Michael Maire, and Jitendra Malik. SVM-KNN: Discriminative nearest neighbor classification for visual category recognition. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2:2126–2136, 2006.