

# On Unsupervised Training of Multi-Class Regularized Least-Squares Classifiers

Tapio Pahikkala<sup>1</sup>, *Member, ACM, IEEE*, Antti Airola<sup>1</sup>, Fabian Gieseke<sup>2</sup>, and Oliver Kramer<sup>3</sup>

<sup>1</sup>*Department of Information Technology, University of Turku, Turku 20520, Finland*

<sup>2</sup>*Department of Computer Science, University of Copenhagen, Copenhagen K 1017, Denmark*

<sup>3</sup>*Computer Science Department, Carl von Ossietzky University of Oldenburg, Oldenburg 26111, Germany*

E-mail: {tapio.pahikkala, antti.airola}@utu.fi; fabian.gieseke@diku.dk; oliver.kramer@uni-oldenburg.de

Received September 1, 2013; revised October 28, 2013.

**Abstract** In this work we present the first efficient algorithm for unsupervised training of multi-class regularized least-squares classifiers. The approach is closely related to the unsupervised extension of the support vector machine classifier known as maximum margin clustering, which recently has received considerable attention, though mostly considering the binary classification case. We present a combinatorial search scheme that combines steepest descent strategies with powerful meta-heuristics for avoiding bad local optima. The regularized least-squares based formulation of the problem allows us to use matrix algebraic optimization enabling constant time checks for the intermediate candidate solutions during the search. Our experimental evaluation indicates the potential of the novel method and demonstrates its superior clustering performance over a variety of competing methods on real world datasets. Both time complexity analysis and experimental comparisons show that the method can scale well to practical sized problems.

**Keywords** unsupervised learning, multi-class regularized least-squares classification, maximum margin clustering, combinatorial optimization

## 1 Introduction

Unsupervised learning belongs to the most important tasks at the beginning of each data mining process: in an early phase, no labeled data at all are given, and the task consists in extracting reasonable information based on the patterns only. Various unsupervised learning tasks like clustering or dimensionality reduction have been proposed in the literature over the years<sup>[1]</sup>. In this work we concentrate on clustering, which plays a central role in a variety of real world applications in computer vision, information retrieval, marketing, context modeling<sup>[2]</sup> and many other fields<sup>[3]</sup>. Roughly speaking, clustering techniques aim at grouping objects into clusters, so that objects with similar characteristics belong to the same cluster, and those with different properties to different ones.

In recent years, the supervised learning method known as the support vector machine (SVM)<sup>[4–5]</sup> as well as other regularized learning schemes have been extended to unsupervised learning settings, in most cases under the name of maximum margin clustering

(MMC)<sup>[6]</sup>. These extensions aim at finding a partition of the unlabeled patterns into classes, so that a subsequent application of the underlying supervised model yields the overall best result. In general, these unsupervised extensions induce combinatorial or non-convex optimization tasks that are difficult to address. However, since the obtained models have been proven to outperform standard clustering techniques in many experimental analyses, they have received considerable attention over the last years.

A preliminary version of this paper was published in the Proceedings of ICDM 2012<sup>[7]</sup>. The theory behind the algorithm framework as well as the discussion and motivation behind the proposed meta-heuristic have been considerably extended in the current article, and the performance evaluation now includes a running time comparison with the baseline methods.

### 1.1 Related Work

Xu *et al.*<sup>[6]</sup> were among the first ones who formalized the extension of support vector machines to unsupervised learning scenarios. Their optimization approach

is based on reformulating the original combinatorial task as semidefinite programming problem<sup>[8]</sup>, which can then be addressed via standard solvers. An extension of this framework based on semidefinite programming as well was provided by Valizadegan and Jin<sup>[9]</sup>. In contrast to the work of Xu *et al.*<sup>[6]</sup>, however, they showed how to reduce the number of involved optimization variables, which yields a more efficient optimization framework. A recent local search approach for the linear case was given by Zhao *et al.*<sup>[10]</sup>. Their optimization framework is based on a combination of recently proposed cutting plane schemes and concave-convex procedures. Similar ideas have also been presented by Li *et al.*<sup>[11]</sup>.

An alternative approach for the binary case was proposed by Zhang *et al.*<sup>[12]</sup>. Basically, their simple but surprisingly effective approach is based on iteratively applying a support vector machine model to improve an “initial guess” that is obtained via an auxiliary clustering framework. One of the key ingredients of their framework is the replacement of the original hinge loss by the  $\varepsilon$ -insensitive or the square loss. As pointed out by Zhang *et al.*<sup>[12]</sup>, the resulting models “*can more easily get out of a poor solution*”. These ideas were extended by Gieseke *et al.*<sup>[13]</sup>, who proposed matrix-based update strategies that can be used to significantly speed up stochastic search frameworks. In line with the approach of Zhang *et al.*<sup>[12]</sup>, they resorted to the square loss instead of the hinge loss.

## 1.2 Contribution

While there exists a significant body of research on extending supervised regularized classifiers to clustering under the framework of maximum margin clustering, almost all of the work has concentrated on the binary case. The exception is the cutting plane multi-class method of Zhao *et al.*<sup>[14]</sup>. The method can be efficiently trained, but is, as demonstrated in our experiments, prone to getting stuck at bad local minima leading to inferior clustering performance. While Xu and Schuurmans<sup>[15]</sup> also formalized a method for the multi-class setting, their method has a runtime complexity of  $O(n^7)$  for  $n$  patterns, which becomes impractical for real world problems.

In this work we extend the concept of supervised *one-vs-all multi-class regularized least-squares classification*<sup>[16]</sup> to unsupervised learning settings. As reported by Zhang *et al.*<sup>[12]</sup> and Gieseke *et al.*<sup>[13]</sup>, the square loss depicts a very reasonable choice in the context of such clustering settings and offers desirable computational shortcuts for corresponding optimization strategies. The contribution provided in this work is twofold:

1) Firstly, we show how to enhance simple steepest descent strategies by means of a powerful meta-heuristic that effectively avoids bad local optima. While being a seemingly simple modification, we demonstrate that this minor adaptation provides major improvements to the clustering accuracy compared with straightforward stochastic search and steepest descent implementations.

2) Secondly, in line with our previous work<sup>[13]</sup>, we provide computational shortcuts for assessing the quality of the intermediate clustering candidate solutions. As we show, these shortcuts render function calls possible to be conducted in  $O(1)$  time, which paves the way for an exhaustive search in the large combinatorial search space.

The meta-heuristic (which we call a *shaking* strategy) is an important algorithmic ingredient for the unsupervised one-vs-all extension which we address. We experimentally analyze our approach on various datasets; the results demonstrate that our approach is capable of yielding better clustering accuracies than conventional techniques in most cases.

So far, no publicly available implementation can be found in the literature that takes care of the interesting multi-class maximum margin principle, and we consider the approach presented in this work to be a valuable candidate for such difficult multi-class clustering settings. Our implementation is made available as part of the RLScore software library at <http://users.utu.fi/aatapa/RLScore/>.

## 2 Mathematical Background

In this section we provide the mathematical notations and the mathematical background related to the general concept of regularized kernel methods<sup>[4–5]</sup>, which encompasses the regularized least-squares classification framework and support vector machines as a special case. We start from the standard supervised setting and proceed to re-formalize the central concepts for the unsupervised learning setting. To simplify the notation, we denote the set of all  $n \times m$  matrices with real coefficients by  $\mathbb{R}^{n \times m}$ . Given a particular matrix  $\mathbf{M} \in \mathbb{R}^{n \times m}$ , we denote its element in the  $i$ -th row and  $j$ -th column by  $\mathbf{M}_{i,j}$ , and we use  $y_i$  to denote the  $i$ -th coordinate of a vector  $\mathbf{y} \in \mathbb{R}^n$ .

### 2.1 Binary Classification Scenarios

We start by depicting the binary cases, for both supervised and unsupervised learning settings. The multi-class scenarios that are central for the work at hand are described afterwards.

### 2.1.1 Supervised Regularized Kernel Methods

Regularized least-squares and support vector machines can be seen as special cases of so-called regularized kernel methods<sup>[4-5]</sup>. We briefly define these settings and then show how to extend the corresponding supervised models to unsupervised (multi-class) learning settings.

Let  $X$  be an arbitrary set and let  $k: X \times X \rightarrow \mathbb{R}$  be a kernel function that can be seen as a similarity measure for the elements in this space. For a given labeled training set

$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \subset X \times Y$$

with  $Y = \{-1, +1\}$  the regularized risk minimization problem is defined as

$$\arg \min_{f \in \mathcal{H}_k} \left\{ \sum_{i=1}^n l(y_i, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}_k}^2 \right\}, \quad (1)$$

where  $f$  is the prediction function (also called *model*) that maps a given data pattern to a real-valued prediction and  $\|\cdot\|_{\mathcal{H}_k}$  is a norm in a reproducing kernel Hilbert space  $\mathcal{H}_k$  induced by the kernel function  $k$ . The disagreement between the predictions and the true labels is measured via a loss function  $l: Y \times \mathbb{R} \rightarrow [0, \infty)$  that gives rise to the *empirical risk*, which, in turn, measures how well the prediction function fits to all training patterns. The regularization parameter  $\lambda \in \mathbb{R}_+$  determines the trade-off between the first term of the task (1) and the complexity of the prediction function  $f$ <sup>[4-5]</sup>.

Two prominent representatives of this family of regularization methods are support vector machines and the concept of regularized least-squares classification<sup>[17]</sup>. The first one stems from the use of the hinge loss  $l(y, f(\mathbf{x})) = \max(0, 1 - yf(\mathbf{x}))$ , whereas the latter one is based on the square loss  $l(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$ . By the representer theorem<sup>[18]</sup>, any solution  $f^* \in \mathcal{H}_k$  of the task (1) has the form

$$f^*(\cdot) = \sum_{i=1}^n a_i k(\mathbf{x}_i, \cdot) \quad (2)$$

with coefficients  $\mathbf{a} = (a_1, \dots, a_n)^T \in \mathbb{R}^n$ . Hence, by plugging in the square loss into the objective and by using  $\|f^*\|_{\mathcal{H}_k}^2 = \mathbf{a}^T \mathbf{K} \mathbf{a}$ <sup>[5]</sup> with kernel matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$  consisting of entries  $\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ , we can rewrite the task at hand as

$$\arg \min_{\mathbf{a} \in \mathbb{R}^n} J(\mathbf{a}) \quad (3)$$

with

$$J(\mathbf{a}) = (\mathbf{y} - \mathbf{K} \mathbf{a})^T (\mathbf{y} - \mathbf{K} \mathbf{a}) + \lambda \mathbf{a}^T \mathbf{K} \mathbf{a}, \quad (4)$$

where  $\mathbf{y} = (y_1, \dots, y_n)^T$ . The objective  $J(\mathbf{a})$  of the above optimization task is convex and differentiable with respect to  $\mathbf{a} \in \mathbb{R}^n$ . Thus a global minimizer can analytically be obtained by enforcing  $\frac{\partial}{\partial \mathbf{a}} J(\mathbf{a}) = 0$ , that is, a minimizer of the objective is given by

$$\mathbf{a}^* = \mathbf{G} \mathbf{y} \quad (5)$$

with

$$\mathbf{G} = (\mathbf{K} + \lambda \mathbf{I})^{-1},$$

where  $\mathbf{I} \in \mathbb{R}^{n \times n}$  is the identity matrix<sup>[17]</sup>. Although the resulting models are, in general, not sparse (as it is often the case for standard support vector machines), the above closed-form solution is a desirable property<sup>[5,17]</sup>. As we will see below, this is especially the case in the context of the considered clustering scenarios.

By substituting the minimizer (5) into (2), it becomes rewritten as the following linear combination of the training labels

$$\begin{aligned} f(\cdot) &= \sum_{i=1}^n (\mathbf{G} \mathbf{y})_i k(\mathbf{x}_i, \cdot) \\ &= \sum_{i,j=1}^n \mathbf{G}_{i,j} y_j k(\mathbf{x}_i, \cdot) \\ &= \sum_{j=1}^n y_j b_j(\cdot), \end{aligned} \quad (6)$$

where the  $n$  basis functions  $b_j(\cdot)$  are defined as

$$b_j(\cdot) = \sum_{i=1}^n \mathbf{G}_{i,j} k(\mathbf{x}_i, \cdot).$$

Note that the basis functions are completely independent on the labels as they only depend on the input data points and on the regularization parameter. This observation is central with regard to this work, since several of the computational shortcuts we take advantage of in the algorithms arise from this property.

In the literature, the representation (6) is sometimes called the dual form of (2) and the basis functions are called the equivalent kernels, because (6) resembles the kernel smoothing technique. For more in depth discussion about these issues, we refer to Girosi *et al.*<sup>[19]</sup> and the references therein.

Accordingly, we obtain the minimum of the objective function (3) as the following closed form expression:

$$\min_{\mathbf{a} \in \mathbb{R}^n} J(\mathbf{a}) = (\mathbf{y} - \mathbf{K} \mathbf{G} \mathbf{y})^T (\mathbf{y} - \mathbf{K} \mathbf{G} \mathbf{y}) + \lambda \mathbf{y}^T \mathbf{G} \mathbf{K} \mathbf{G} \mathbf{y}.$$

Inspired by this, we define a new objective function in which the above closed form expression is written as a function of the label vector  $\mathbf{y}$ :

$$F(\mathbf{y}) = (\mathbf{y} - \mathbf{K} \mathbf{G} \mathbf{y})^T (\mathbf{y} - \mathbf{K} \mathbf{G} \mathbf{y}) + \lambda \mathbf{y}^T \mathbf{G} \mathbf{K} \mathbf{G} \mathbf{y}. \quad (7)$$

This provides a powerful and easy to use approach for testing different labelings for the given data. Below, we will use this function as a basic tool when considering the least-squares based clustering algorithms.

### 2.1.2 Unsupervised Least-Squares Extension

As pointed out by Zhang *et al.*<sup>[12]</sup>, the square loss depicts an ideal candidate for the maximum margin principle from a practical point of view. Further, the above closed-form solution can also be used to greatly speed up the computations induced by a variety of search strategies<sup>[13]</sup>. The direct extension of the supervised regularized kernel methods (1) to the unsupervised case for a given unlabeled training set  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset X$  has the form<sup>[6]</sup>:

$$\arg \min_{\mathbf{y} \in \{-1, +1\}^n, f \in \mathcal{H}_k} \left\{ \sum_{i=1}^n l(y_i, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}_k}^2 \right\}.$$

Hence, the difficult part is the additional integer optimization variable  $\mathbf{y} \in \{-1, +1\}^n$  that encodes the partition of the given unlabeled patterns. To avoid trivial solutions, some form of a *balancing constraint* is usually added for such clustering settings, which is of the form

$$\left| \frac{1}{n} \sum_{i=1}^n y_i \right| < b_c$$

with a user-defined parameter  $b_c \in [0, 1]$ . By again considering the square loss, we can take advantage of the tools established in the previous subsection and can rewrite the above optimization problem in terms of the label vector only:

$$\arg \min_{\mathbf{y} \in \{-1, +1\}^n} F(\mathbf{y}) \quad (8)$$

with  $F(\mathbf{y})$  being the closed form given in (7).

## 2.2 Multi-Class Classification Scenarios

We can now address the multi-class learning settings that are the basis of the optimization schemes derived in this work. Like above, we start by outlining the supervised models followed by their unsupervised extensions.

### 2.2.1 Supervised Multi-Class Extension

In the literature, several ways to extend the concept of support vector machines and their variants to multi-class settings have been proposed. As reported by Rifkin and Klautau<sup>[16]</sup> the so-called one-versus-all multi-class classification setting depicts a valuable candidate for such learning scenarios, and we follow this line of research for the unsupervised case.

In such multi-class supervised settings, we are given a training set

$$\{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subset X \times \mathcal{C}$$

with  $\mathcal{C} = \{1, \dots, |\mathcal{C}|\}$  as set of all possible class labels. In a nutshell, one aims at deriving models  $f_1, \dots, f_{|\mathcal{C}|}$  such that a new pattern  $\mathbf{x} \in X$  is assigned to the class whose associated model is the most confident. Formally, given the  $|\mathcal{C}|$  models, the class label for a data point  $\mathbf{x}$  is obtained as:

$$f(\mathbf{x}) = \arg \max_{c \in \mathcal{C}} f_c(\mathbf{x}), \quad (9)$$

that is, the model with the largest real-valued prediction determines the class label.

It is easy to see that the above multi-class classification rule can also be written as

$$f(\mathbf{x}) = \arg \min_{c \in \mathcal{C}} L(c, f_1(\mathbf{x}), \dots, f_{|\mathcal{C}|}(\mathbf{x})), \quad (10)$$

where  $L$  is the following type of a one-versus-all multi-class loss function:

$$L(c, f_1(\mathbf{x}), \dots, f_{|\mathcal{C}|}(\mathbf{x})) = \sum_{h=1}^{|\mathcal{C}|} l(p_h(c), f_h(\mathbf{x})),$$

and

$$p_h(c) = \begin{cases} 1, & \text{if } h = c, \\ -1, & \text{otherwise.} \end{cases} \quad (11)$$

Here, the one-versus-all multi-class loss is decomposed to a sum of ordinary type of binary loss functions.

Let  $\mathbf{c} \in \mathcal{C}^n$  be the vector containing the class labels of the training examples. Further, in order to support the forthcoming considerations below, we overload the notation  $p_h$  by defining it also for vectors of class labels, which leads to the following mapping from the class label vector to  $\{-1, +1\}$ -valued vectors, one per each class  $h \in \mathcal{C}$ :

$$p_h(\mathbf{c}) = \begin{pmatrix} -1 \\ \vdots \\ -1 \end{pmatrix} + 2 \sum_{j=1}^n \delta_{c_j h} \mathbf{e}^j, \quad (12)$$

where  $\delta$  is the Kronecker delta (i.e., we have  $\delta_{c_j h} = 1$  if  $c_j = h$  and  $\delta_{c_j h} = 0$  otherwise), and  $\mathbf{e}^j$  is the  $j$ -th standard basis vector of  $\mathbb{R}^n$ . Hence, the  $i$ -th component of vector  $p_h(\mathbf{c})$  equals  $+1$  in case the  $i$ -th training pattern belongs to the class  $h$ , and equals  $-1$  otherwise (see Fig.1 for an example).

A variety of different objectives (and loss functions) have been proposed in the literature<sup>[16]</sup>. In the following, we consider extensions of the supervised models (1) to the multi-class case having the form:

$$\begin{aligned}
\mathbf{c} &= \begin{pmatrix} 1 \\ 2 \\ 1 \\ 3 \\ 2 \\ 4 \\ 2 \end{pmatrix} & p_1(\mathbf{c}) &= \begin{pmatrix} +1 \\ -1 \\ +1 \\ -1 \\ -1 \\ -1 \\ -1 \end{pmatrix} & p_2(\mathbf{c}) &= \begin{pmatrix} -1 \\ +1 \\ -1 \\ -1 \\ +1 \\ -1 \\ +1 \end{pmatrix} & p_3(\mathbf{c}) &= \begin{pmatrix} -1 \\ -1 \\ -1 \\ +1 \\ -1 \\ -1 \\ -1 \end{pmatrix} & p_4(\mathbf{c}) &= \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \\ -1 \\ +1 \\ -1 \end{pmatrix} \\
& \text{(a)} & & & & & & \text{(b)}
\end{aligned}$$

Fig.1. (a) A vector of class labels and (b) corresponding mappings.

$$\arg \min_{f_1, \dots, f_{|C|} \in \mathcal{H}_k} M(f_1, \dots, f_{|C|}, \mathbf{c}),$$

where

$$\begin{aligned}
& M(f_1, \dots, f_{|C|}, \mathbf{c}) \\
&= \sum_{h=1}^{|C|} \sum_{i=1}^n \left( l(p_h(c_i), f_h(\mathbf{x}_i)) + \lambda \|f_h\|_{\mathcal{H}_k}^2 \right)
\end{aligned}$$

for some vector of class labels  $\mathbf{c}$ . Hence, the goal of the learning process is the search of binary-valued prediction functions  $f_1, \dots, f_{|C|}$  that minimize the above risk. Given that the one-versus-all loss is decomposed into a sum of binary ones, the process can be considered as training  $|C|$  binary models *independently*. As pointed out above, such frameworks have been shown to work as well as other sophisticated multi-class schemes for such supervised settings, see the work of Rifkin and Klautau<sup>[16]</sup>.

Finally, let us consider the one-versus-all classification with the squared loss. Using analogous arguments as in the binary classification case, the minimizer of the above objective function has a closed form, which can be written simply as

$$\min_{f_1, \dots, f_{|C|} \in \mathcal{H}_k} M(f_1, \dots, f_{|C|}, \mathbf{c}) = Q(\mathbf{c}),$$

where

$$Q(\mathbf{c}) = \sum_{h=1}^{|C|} F(p_h(\mathbf{c})), \quad (13)$$

and  $F$  is defined via (7).

### 2.2.2 Unsupervised Least-Squares Extension

Exactly as for the binary case, one can extend the multi-class framework depicted above to unsupervised settings by considering the class membership vector  $\mathbf{c} \in \mathcal{C}^n$  as additional optimization variable. For the square loss, this leads to the following optimization problem:

$$\arg \min_{\substack{\mathbf{c} \in \mathcal{C}^n \\ f_1, \dots, f_{|C|} \in \mathcal{H}_k}} M(f_1, \dots, f_{|C|}, \mathbf{c}) = \arg \min_{\mathbf{c} \in \mathcal{C}^n} Q(\mathbf{c}).$$

Hence, one is again given a mixed-integer programming problem; the key problem is to find an appropriate as-

signment for the variable  $\mathbf{c}$  such that the induced  $|C|$  supervised binary classification tasks yield the overall best results. Note that, while the objectives seem to be independent from each other, they interact via the vector  $\mathbf{c}$ , since changing the class membership of a *single* training instance leads to the modification of *two* of the induced classification models  $f_1, \dots, f_{|C|}$ .

Note that additional constraints can (and should) be used to enforce appropriate ratios of the cluster sizes. In the next section, we propose an efficient algorithm for finding accurate clustering solutions that aim at minimizing the above objective subject to such cluster constraints, along with computational shortcuts for assessing the quality of intermediate candidate solutions.

## 3 Algorithmic Framework

In this section we describe the basic ideas behind the proposed algorithm without getting into the computational details. Due to the discrete nature of the clustering problem, we employ direct optimization methods for searching appropriate labels for the data points. In the literature, this type of methods are often referred to as hill climbing algorithms (see e.g., [20]).

There are different variants of hill climbing, such as stochastic and steepest hill climbing. In addition, there are so-called meta-algorithms that are built on top of the hill climbing algorithms, such as climbing with random restarts. Here, we focus mainly on the idea of the steepest descent hill climbing, in which all closest neighbors of the current solution are compared, and the current solution is replaced with the neighbor having the lowest value of the objective function. In our case, the set of closest neighbors consists of label vectors that differ from the current solution *only by one entry*. In addition, we propose a meta-algorithm we call *shaking* that uses the idea of steepest descent so that it is less likely to get stuck to local minima with inferior clustering performance than the basic steepest descent search.

For convenience, we use  $S(\mathbf{c}, j, d)$ , where  $\mathbf{c} \in \mathcal{C}^n$ ,  $j \in \{1, \dots, n\}$ , and  $d \in \mathcal{C}$ , to denote the value of the objective function  $Q$  defined in (13) for a cluster label vector, whose entries are equal to those of  $\mathbf{c}$  except that the label of the  $j$ -th data point has been switched

from  $c_j$  to  $d$ . This allows us to denote the search directions in the space of cluster label vectors so that each direction corresponds to switching the cluster label of a single data point. Armed with the above notation and a vector of initial cluster assignments  $c \in \mathcal{C}^n$  for  $n$  data points, we next consider the search algorithms for solving the clustering problems.

### 3.1 Basic Descent Strategies

One of the most straightforward approaches is the so-called stochastic hill climbing, in which the algorithm simply goes through the data points one at a time and switches its current class label to another one if it decreases the objective value, and stops when a local optimum is found. This type of algorithms were proposed for binary clustering in our previous work<sup>[13]</sup>. In our experiments, we use a similar algorithm (Algorithm 1) modified for multi-class clustering as a baseline method. The modification can take advantage of the computational shortcuts presented in the appendix, and it is therefore computationally as efficient as the other methods proposed in this paper.

Another framework is the basic steepest descent search (see Algorithm 2) for the multi-class clustering

#### Algorithm 1. STOCHASTIC DESCENT

```

1: Initialize  $c \in \mathcal{C}^n$  randomly
2: loop
3:    $b \leftarrow \text{True}$ 
4:   for  $j = 1, \dots, n$  do
5:      $d \leftarrow \arg \min_{d \in \mathcal{C}} S(c, j, d)$ 
6:     if  $c_j \neq d$  then
7:        $c_j \leftarrow d$ 
8:        $b \leftarrow \text{False}$ 
9:     end if
10:  end for
11:  if  $b$  then ▷ Stop if local optimum found
12:    break
13:  end if
14: end loop
    
```

#### Algorithm 2. STEEPEST DESCENT

```

1: Initialize  $c \in \mathcal{C}^n$  randomly
2: loop
3:    $j, d \leftarrow \arg \min_{j \in \{1, \dots, n\}, d \in \mathcal{C}} S(c, j, d)$ 
4:   if  $c_j = d$  then
5:     break
6:   else
7:      $c_j \leftarrow d$ 
8:   end if
9: end loop
    
```

problem. The idea is that during each iteration the algorithm finds a pair  $(j, d)$ , where  $j$  is the index of the data point, for which switching the cluster label would decrease the value of the objective function the most and  $d$  is the corresponding new cluster label. The algorithm stops when a local minimum is found, that is, switching the cluster assignment of a single data point will not decrease the objective value.

### 3.2 Avoiding Local Minima via Shaking

Both the stochastic and steepest descent methods can easily get stuck in local minima corresponding to inferior clusterings of the data; the existence and severity of this problem is confirmed in our experiments. For this reason, we propose to improve the steepest descent algorithm with a simple, but surprisingly effective trick (see Algorithm 3).

#### Algorithm 3. STEEPEST DESCENT WITH SHAKING

```

1: Initialize  $c \in \mathcal{C}^n$  randomly
2: for  $i = 0, \dots, s$  do
3:   for  $d \in \mathcal{C}$  do
4:      $\alpha \leftarrow \frac{n}{2^i |\mathcal{C}|} + \frac{n}{|\mathcal{C}|} - |\{h \mid c_h = d\}|$ 
5:     for  $j = 1, \dots, \alpha$  do
6:        $j \leftarrow \arg \min_{j \in \{1, \dots, n\}, d \neq c_j} S(c, j, d)$ 
7:        $c_j \leftarrow d$ 
8:     end for
9:   end for
10: end for
    
```

Instead of traversing the search space exactly towards the steepest descent direction, the algorithm iterates through the clusters and each cluster at a time claims a number of points from the other clusters. The point the cluster  $d$  claims next is determined by the steepest descent direction. That is, the point for which switching the cluster label to  $d$  would decrease the objective value the most (or increase the least), is assigned to the cluster  $d$ .

The number of points claimed by the clusters is determined by the following formula:

$$\alpha = \frac{n}{2^i |\mathcal{C}|} + \frac{n}{|\mathcal{C}|} - |\{h \mid c_h = d\}|. \quad (14)$$

Firstly, it depends on the round  $i$  of the outer loop, which is encoded in the first term of (14) via the exponential factor in the nominator. During the first iterations, all clusters claim a large number of points, but the value of the first term decreases exponentially with respect to the iteration index of the loop.

To prevent the smaller clusters from disappearing completely, the number of points claimed by a cluster

also depends on the number of points already assigned to the cluster in question, that is, the small clusters claim more points than the large ones. This behavior is ensured by the sum of the second and third terms of (14), which is positive for clusters smaller than average sized and negative for the larger ones. Note that this definition of  $\alpha$  is just an ad hoc heuristic and better ones may be designed, for example, if we have a prior knowledge about the cluster sizes. Nevertheless, as we show in the experiments, even this simple approach performs considerably better than the stochastic hill climbing and the basic steepest descent, since it avoids many of the inferior local minima into which the basic approaches get stuck.

Intuitively, the idea of this approach can be considered as “shaking” the solution so that the local minima can be avoided. In the beginning, the solution is shaken profusely, but the intensity quickly decreases with the rounds of the outer loop. Because of the exponential decrease, the number  $s$  of shakings can be set to a small constant. In our experiments, we always use the value  $s = 20$ . This turned out to be a comparatively robust parameter choice in preliminary experiments.

The behavior of the proposed method is demonstrated in Fig.2 with a toy example consisting of three Gaussian clusters. First, the patterns are labeled randomly. Then, each cluster claims points from other clusters at a time and the number of points claimed decreases during each iteration.

The shaking heuristic is a powerful optimization approach due to its dynamic mechanism that controls the magnitude of changes according to a cooling scheme depending on the iteration number. In this regard, the shaking heuristic is similar to dynamic mutation rate control employing an external cooling scheme as simulated annealing (SA)<sup>[21]</sup> does. SA adapts the probability to accept a solution deterioration according to an exponential cooling scheme. A change of solution  $\mathbf{c}$  that increases the fitness value by  $\Delta S$  is accepted with a decreasing probability

$$p(\Delta S, T) = e^{-\Delta S/T}$$

depending on the fitness change  $\Delta S$  and temperature  $T$  that must be increased in the course of the optimization run (e.g., depending on the iteration number). Our shaking heuristic differs from SA, as each solution change is accepted, but the magnitude of change is controlled dynamically. The dynamic scheme that is employed depends on the iteration number in an exponential kind of way. It is possible to introduce a further parameter that controls the speed of the decrease of  $\alpha$ . Extensions with this regards will be subject to future work.

A further characteristic of the shaking heuristic is that not a random part of the solution is optimized, but the search concentrates on particular solution parts, i.e., one cluster in each iteration employing greedily the steepest descent approach. A greedy cluster assignment is a recommendable optimization strategy, in particular in case of large solution spaces. Another important part of the shaking heuristic is the second and the third summands of (14). They give small clusters the chance to develop, discouraging the optimization process from losing found clusters. Equipped with this knowledge, the algorithm has an advantage over pure direct search.

### 3.3 Efficient Optimization via Matrix-Based Updates

As observed above, all considered search algorithms can be formulated in terms of the operation  $S(\mathbf{c}, j, d)$  that computes the objective value for a certain neighbor of the label vector  $\mathbf{c}$ . In particular, we note that the shaking heuristic requires roughly  $O(s|\mathcal{C}|n^2)$  calls of  $S(\mathbf{c}, j, d)$ .

A naive approach for computing it would require re-training the classifiers from scratch each time the operation is used, which would clearly be computationally infeasible except for very small datasets. For example, the computational complexity of training support vector machine classifiers with non-linear kernels is  $O(n^3)$  in the worst case, and hence the overall complexity of running the shaking heuristic together with those would be of order  $O(s|\mathcal{C}|n^5)$ .

The next theorem characterizes one of the main contributions of this paper. This offers massive runtime savings compared with the naive implementation mentioned above.

**Theorem 1.** *The computational complexity of Algorithm 3 is*

$$O((s|\mathcal{C}| + r)n^2),$$

where  $r$  is the rank of the kernel matrix.

*Proof.* For the sake of exposition, we defer the lengthy proof to the appendix. The general idea is that given a certain amount of time spent in the preprocessing phase, one can test each possible flip of a coordinate in  $O(1)$  time, which is far less compared with a naive implementation that would take  $O(n^2)$  time per flip. Algorithm 3 performs  $O(s|\mathcal{C}|n^2)$  calls for the function  $S(\mathbf{c}, j, d)$ , each with  $O(1)$  cost, and the initialization of the caches requires  $O(n^2r)$  time.  $\square$

As we show in the experimental evaluation, both the shaking framework and the computational shortcuts are important algorithmic ingredients to address the challenging combinatorial task induced by the unsupervised multi-class extension of the maximum margin principle considered in this work.

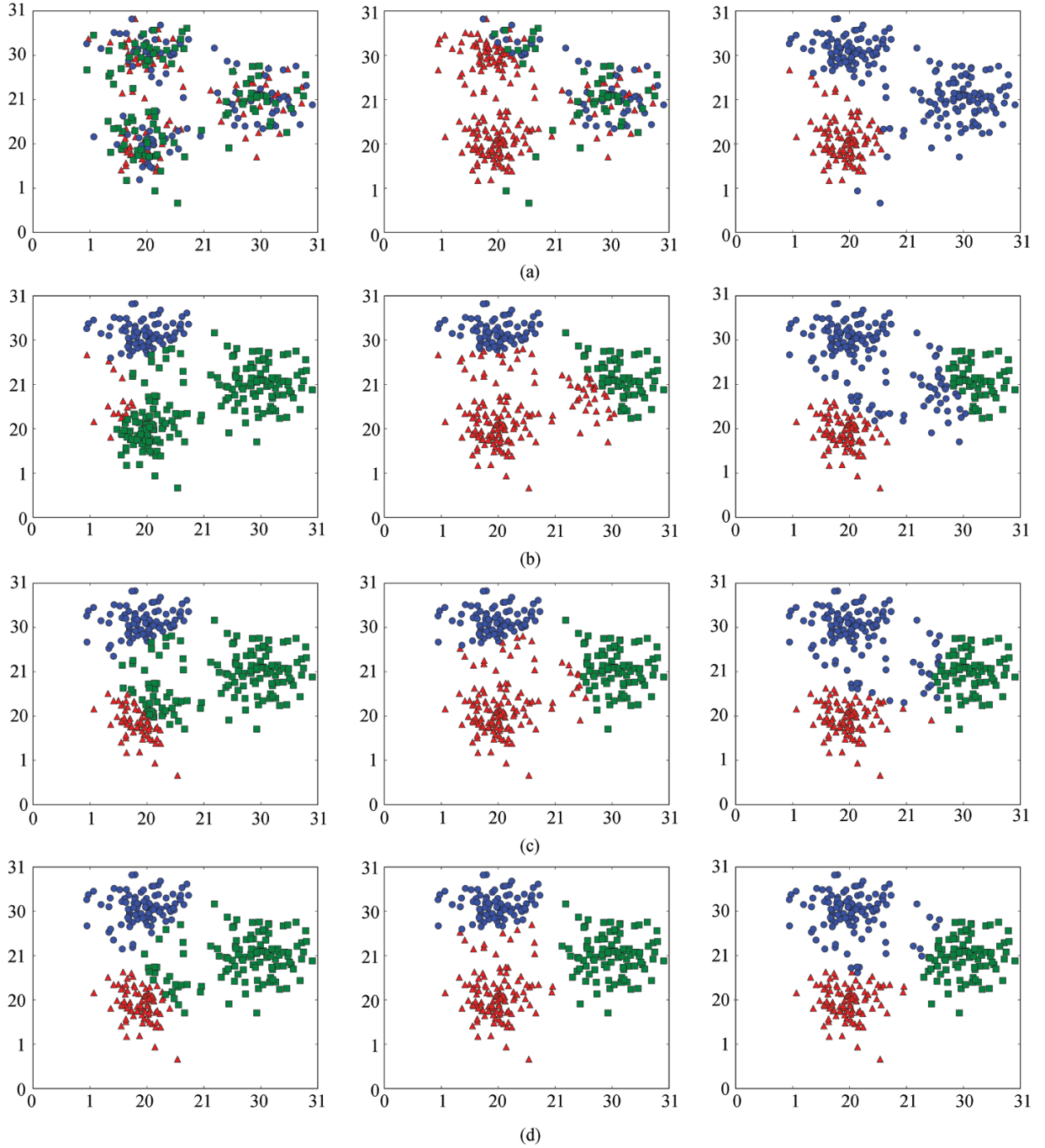


Fig.2. Demonstration of the shaking meta-algorithm on three Gaussian clusters. The first image contains the initial random assignments of the cluster labels. The next three images correspond to the first round of the outermost loop during which each class claims a large number of data points at a time. During the following rounds, the classes keep claiming data points but the number of points claimed decreases exponentially with the index of the outermost loop.

## 4 Experiments

In the experiments we evaluate the accuracy of the proposed unsupervised multi-class regularized least-squares algorithm (UMC-RLS) on several real world and synthetic datasets with various baseline methods.

### 4.1 Methods

As baseline methods we employ the cutting plane multi-class MMC method (CPMMC)<sup>[14]</sup>,  $K$ -means clustering (more specifically,  $k$ -means++<sup>[22]</sup>), Gaussian mixture models fitted using expectation maximization



with full covariance structure (GMM)<sup>[23]</sup>, and spectral clustering using a 10-nearest neighbors-graph as similarity graph (SC)<sup>[24]</sup>. Additionally, to demonstrate the importance of the shaking heuristic, we also provide results for simpler variants of the proposed method, where the combinatorial search over the cluster assignments is based either on pure stochastic search (Stoc-RLS) or on the direction of the steepest descent without the shaking heuristic (Steep-RLS). As discussed above, the stochastic variant extends our previous implementation for the binary case<sup>[13]</sup> to the multi-class setting.

UMC-RLS, Steep-RLS, Stoc-RLS, and CPMMC are based on `Python` and the `Numpy` and `Scipy` libraries. Further, the CPMMC implementation uses the `CVXOPT` optimization library for solving the quadratic programs arising during training the method. The method was originally formulated only for the linear case, but as suggested by Zhao *et al.*<sup>[14]</sup>, the method can be straightforwardly kernelized by running it on a feature representation generated by eigen decomposing the kernel matrix (see, e.g., [25] for a detailed discussion). The implementation, which is based on the primal formulation of the optimization problem, does not scale to large problem sizes, in effect restricting our comparison to problems with few hundred training examples, and few clusters at most. Zhao *et al.*<sup>[14]</sup> noted that one can derive a dual version of the method with much better scalability, but neither technical details nor a corresponding implementation are provided. The *K*-means, GMM and SC implementations are from the `scikit-learn` library available at <http://scikit-learn.org>.

## 4.2 Datasets

We perform experiments on eight tasks that represent a wide variety of application domains (see Table 1). The Iris, Letter, and USPS datasets are standard benchmarks from the UCI repository. From Letter, we choose the first 4 classes in the data. We split the USPS into two tasks, USPS 1~4 and USPS 5~8,

which both contain four classes from the original data. We use the full COIL image recognition dataset<sup>[26]</sup>, as well as two subsets, COIL 1~4 and COIL 5~8. Moons is a well-known artificial benchmark dataset with a non-linear structure. The Letter and USPS datasets are sub-sampled so that they have at most 500 examples, in order to allow for running the experiments with CPMMC. For the full COIL dataset we do not present results for CPMMC, as the implementation does not scale to the considered number of patterns and clusters.

## 4.3 Clustering Performance

We estimate the clustering performance of the compared methods using the adjusted rand index (ARI)<sup>[27]</sup>. After parameter selection, each method is run 10 times, with mean ARI of the repetitions being used to represent the final performance. All kernel methods employ a Gaussian kernel in our experiments.

Similarly to the setups of [14, 28], we choose the regularization parameter  $\lambda$  and kernel width  $\sigma$  for UMC-RLS, Steep-RLS, Stoc-RLS, and CPMMC using grid search. In preliminary experiments we noticed that the methods tended to favor small regularization parameter values, therefore  $\lambda$  is chosen from grid  $\{2^{-10}, 2^{-9}, \dots, 2^{-1}\}$ . Kernel width  $\sigma$  is chosen from the grid  $\{0.1\sigma_0, 0.2\sigma_0, \dots, \sigma_0\}$ , where  $\sigma_0$  is the maximum distance between two data points in the dataset. For each tested parameter, the performance is computed as the mean over 10 repetitions of clustering with different random initializations. Following [14, 28], we set the error tolerance parameters  $\alpha$  and  $\epsilon$  for CPMMC both to 0.01. The parameter  $l$  of CPMMC is set to 10, based on preliminary experiments.

Table 1 presents the mean ARI with standard deviations for the considered methods and data, with the results from best performing methods highlighted in each row. On seven of the eight considered datasets, UMC-RLS either outperforms all the other methods, or shares the place of best performing method with one other baseline method. On the USPS 5~8 data, UMC-

**Table 1.** Comparison of the Clustering Methods  
(Mean ARI and the One Standard Deviation Based on 10 Repetitions are Provided)

Dataset	Size	Classes	SC	GMM	<i>K</i> -Means	CPMMC	Stoc-RLS	Steep-RLS	UMC-RLS
Coil	1 440	20	0.74 ± 0.01	0.55 ± 0.01	0.58 ± 0.02	–	0.31 ± 0.10	0.43 ± 0.07	<b>0.84 ± 0.06</b>
Coil 1~4	288	4	0.70 ± 0.00	0.50 ± 0.00	0.49 ± 0.00	0.47 ± 0.08	0.29 ± 0.11	0.29 ± 0.09	<b>1.00 ± 0.00</b>
Coil 5~8	288	4	0.74 ± 0.06	0.61 ± 0.03	0.60 ± 0.04	0.62 ± 0.05	0.35 ± 0.06	0.37 ± 0.10	<b>0.97 ± 0.09</b>
Iris	150	3	0.43 ± 0.00	<b>0.96 ± 0.00</b>	0.70 ± 0.09	0.69 ± 0.08	0.45 ± 0.25	0.56 ± 0.13	<b>0.96 ± 0.00</b>
Letter	500	4	0.38 ± 0.00	0.45 ± 0.00	0.43 ± 0.00	0.39 ± 0.10	0.19 ± 0.12	0.20 ± 0.12	<b>0.46 ± 0.09</b>
Moons	500	2	0.99 ± 0.00	0.78 ± 0.00	0.65 ± 0.00	0.69 ± 0.45	0.21 ± 0.35	0.07 ± 0.05	<b>1.00 ± 0.00</b>
USPS 1~4	500	4	0.66 ± 0.00	0.48 ± 0.00	0.61 ± 0.01	0.59 ± 0.19	0.38 ± 0.12	0.38 ± 0.08	<b>0.85 ± 0.02</b>
USPS 5~8	500	4	<b>0.93 ± 0.00</b>	0.64 ± 0.00	0.67 ± 0.00	0.57 ± 0.11	0.26 ± 0.12	0.34 ± 0.14	0.91 ± 0.02

**Table 2.** Comparison of the Clustering Methods (Maximum ARI out of 10 Repetitions is Provided)

Dataset	SC	GMM	$K$ -Means	CPMMC	Stoc-RLS	Steep-RLS	UMC-RLS
Coil	0.75	0.57	0.62	–	0.43	0.58	<b>1.00</b>
Coil 1~4	0.70	0.50	0.49	0.56	0.50	0.43	<b>1.00</b>
Coil 5~8	0.79	0.65	0.65	0.67	0.42	0.51	<b>1.00</b>
Iris	0.43	<b>0.96</b>	0.62	0.76	0.85	0.92	<b>0.96</b>
Letter	0.38	0.45	0.43	0.49	0.37	0.36	<b>0.57</b>
Moons	0.99	0.78	0.65	0.99	0.95	0.14	<b>1.00</b>
USPS 1~4	0.66	0.48	0.62	<b>0.92</b>	0.51	0.52	0.88
USPS 5~8	<b>0.93</b>	0.65	0.68	0.68	0.51	0.57	<b>0.93</b>

RLS is the second best performing method. SC and GMM are the most competitive baselines. SC outperforms the other methods on USPS 5~8 data, and is the only baseline method also able to solve the Moons problem. GMM performs as well as UMC-RLS on Iris, and almost as well on Letter.  $K$ -means and CPMMC are not competitive with UMC-RLS, but can still on some of the datasets outperform either SC or GMM. The mean clustering performances of Stoc-RLS and Steep-RLS are very poor. On most runs CPMMC, Stoc-RLS and Steep-RLS seem to get stuck in bad local minima. Thus, the shaking heuristic implemented by UMC-RLS proves to be crucial in order to achieve stable and good performance.

Next, we compare the methods based on the maximum ARI achieved out of the final ten runs. The experiment allows us to estimate whether the performance of some of the considered methods could be significantly improved by using random restarts based meta-heuristics. The results are presented in Table 2. In this setting, CPMMC becomes competitive with the SC and GMM methods, even outperforming all the methods on USPS1 datasets. Stoc-RLS and Steep-RLS results are also greatly improved compared with mean results. Still, even if we compare the maximum ARI out of 10 repetitions for the other methods (Table 2) with the mean ARI out of 10 repetitions for UMC-RLS (Table 1) UMC-RLS still appears to be the overall best performing method.

To conclude, the experimental results suggest that the proposed UMC-RLS method often achieves high clustering performance on real world problems, and seems to represent the state-of-the-art among multi-class MMC type of methods. Further, the results highlight the importance of the shaking heuristic, as otherwise the combinatorial search will typically not lead to a good clustering solution.

#### 4.4 Runtime

Finally, we also experimentally explore the scalability of the proposed method. All the runtime experiments were performed on a computer with Intel i7-

3770 3.40 GHz processor, 16 GB of memory and 64-bit Ubuntu 13.04 operating system. The CPMMC is not considered in the comparison, since its implementation is highly suboptimal in terms of running time.

In Fig.3 we have plotted the number of steepest descent steps executed while running UMC-RLS for varying sized random subsets of the COIL dataset with 20 clusters. From the plot it can be seen that the number of steps required grows linearly with the size of the dataset, as can be expected from the complexity considerations. Thus the experiment further verifies that the steepest descent search can be executed efficiently for the proposed method.

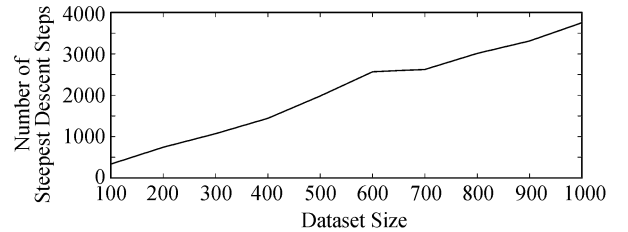


Fig.3. Number of steepest descent steps required by UMC-RLS as a function of the dataset size.

In Fig.4 we have the CPU running time for the different methods on varying sized random subsets of the COIL dataset. It can be observed that when dealing with such a moderately sized dataset the UMC-RLS implementation actually outperforms the baseline methods in terms of computational speed. To evaluate the behaviour of the methods on larger datasets we further performed a scaling experiment on the MNIST handwritten character recognition dataset. The dimensionality of the data is 784, and it contains 10 clusters representing different digits. As can be seen in Fig.5, once the sample size becomes several thousands, the higher computational complexity of UMC-RLS begins to dominate the running time of the method, the  $K$ -means and SC have better scalability. Still, even with 8 000 data points UMC-RLS clustering can be run in a bit over five minutes on a modern desktop computer.

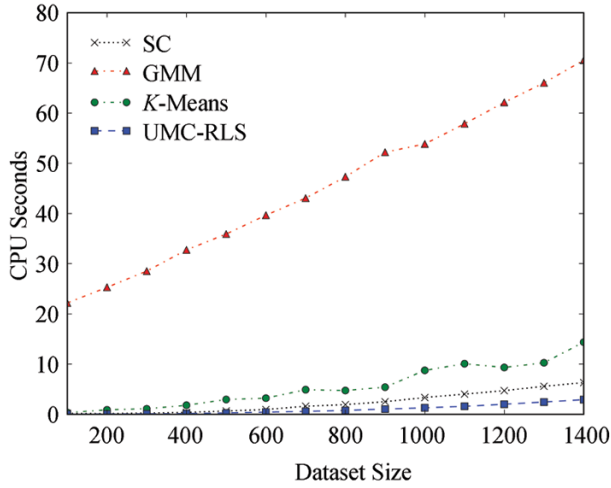


Fig.4. CPU runtime on the COIL dataset.

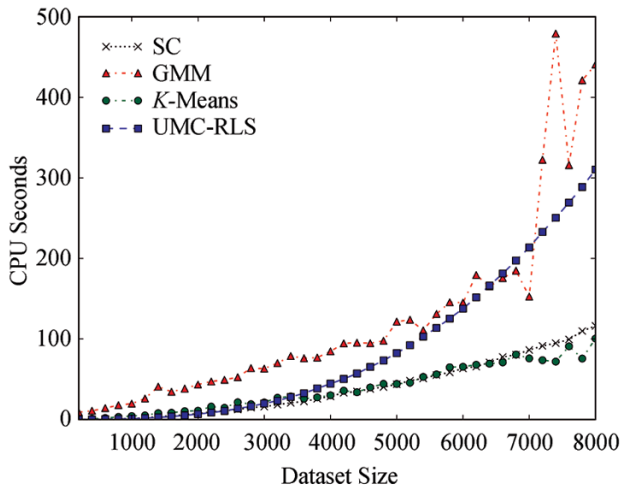


Fig.5. CPU runtime on MNIST dataset.

The computational bottleneck for UMC-RLS is still the computation of the eigen decomposition of the kernel matrix needed during initialization. Here, kernel matrix approximation techniques could be of great benefit, in case one needs to scale the method to very large datasets.

## 5 Discussion and Future work

As shown by the experimental results, the proposed shaking heuristic provides considerably better results than the simple greedy approach only relying on the steepest descent directions. Still, the heuristic is the first non-trivial one we tested, and hence it is of very ad-hoc nature. We expect that far better heuristics can be produced if we can encode prior knowledge about the classification problem into it, as is often possible in practical problems. Heuristics could also be designed for other variations of unsupervised classification, such as learning with partial class memberships<sup>[29]</sup>.

The main computational bottleneck of the proposed algorithm is computing the eigen decomposition of the kernel matrix, whose time complexity is cubic with respect to dataset size in the worst case. For large datasets, a standard practice in kernel-based learning is to employ sparse kernel matrix approximation techniques, such as the well-known Nyström method<sup>[30]</sup>, which will decrease the complexity of performing the decomposition to linear time, usually without considerably harming the classification performance. Another bottleneck is due to the linear cost of a single steepest descent step, which causes the overall time complexity to become quadratic if the amount of steps also grows linearly with respect to dataset size, as is usually the case with the shaking heuristic. To remedy this, we intend to develop such variations of the method and the heuristic that, instead of doing global steepest descent steps with linear cost, would employ the steps on small local subsets of the data at the time so that the step costs would scale with the subset sizes rather than the overall dataset size. This would also reduce the memory size of the cache matrices required by the method (see Appendix), since the caches would have to be constructed for the subsets only and reconstructed when the subset would be changed. If the size of the subsets can be fixed to a small constant, the overall computational and memory complexities of the method will be linear with the overall dataset size. This requires considerably more sophisticated heuristics that also take care of changing the local subsets when needed.

The performance of kernel-based learning methods depends considerably on the hyper-parameters values, such as those of the regularization and kernel parameters. However, tuning the values properly is very challenging in unsupervised learning. There exist several methods for measuring the cluster validity that do not depend on the class-labels of the data points but work in completely unsupervised fashion<sup>[31-32]</sup>. In the future, we also intend to investigate the potential of these methods for setting the hyper-parameter values. The problem is not so severe in tasks like image segmentation, where the method is used in an interactive fashion so that the user can tune the parameter values by hand in order to get results that are satisfactory enough.

Yet another research direction not yet covered in this paper is the semi-supervised extension of the proposed method. Given a small set of data points with known class labels and a large set with unknown labels, the proposed ideas can be easily modified so that the known class labels are kept fixed, while the labels of the other points are switched with the meta-heuristic just like in the fully unsupervised case. The problem, of course, becomes considerably easier, since the classes can in

most practical cases be assumed to be concentrated around the points with known labels and the method does not have to be initialized with a completely random labeling or with a simpler clustering algorithm. Moreover, instead of the shaking proposed here for the completely unsupervised learning, new meta-heuristics can be specially tailored for the semi-supervised case, for example, such that keep claiming unlabeled points to the classes one by one from the vicinity of the labeled points until all points have been labeled. With the help of the efficient computational shortcuts proposed in this paper for searching the steepest label switching directions, it becomes possible to develop and test several different types of new heuristics also for the semi-supervised problems.

## 6 Conclusions

In this work we proposed a multi-class extension of the binary maximum margin principle. Our framework is based on the least-squares variant of the original problem formulation, which has been experimentally proven to be a valuable candidate for such clustering settings, e.g., [12] or [13]. So far, only little work has been done related to the interesting extension of these schemes to the multi-class case though. This is the focus of the work at hand dealing with the unsupervised extension of the corresponding one-vs-all multi-class setting.

The key contributions provided in this work are 1) a carefully designed steepest descent strategy, and 2) its extremely efficient implementation. The former contribution is based on a new shaking strategy that effectively avoids getting trapped in bad local optima during early stages of the optimization process. The latter contribution is based on a series of non-trivial matrix-based update steps that take care of the intermediate optimization tasks induced by the global shaking steepest descent framework. The experimental evaluation takes into account a variety of real world datasets, and the clustering accuracy of our approach was compared with the ones of state-of-the-art methods. The results demonstrate the superior performance of the proposed framework and, hence, indicate that the unsupervised regularized least-squares approach is a promising clustering variant, given that one addresses the induced combinatorial optimization task appropriately.

**Acknowledgments** We would like to thank the anonymous reviewers for their comments.

## References

- [1] Hastie T, Tibshirani R, Friedman J. The Elements of Statistical Learning: Data Mining, Inference, and Prediction (2nd edition). New York, NY, USA: Springer, 2009.

- [2] Bao T, Cao H, Chen E, Tian J, Xiong H. An unsupervised approach to modeling personalized contexts of mobile users. *Knowledge and Information Systems*, 2012, 31(2): 345-370.
- [3] Jain A, Dubes R. Algorithms for Clustering Data. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1988.
- [4] Schölkopf B, Smola A. Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. Cambridge, MA, USA: MIT Press, 2001.
- [5] Steinwart I, Christmann A. Support Vector Machines. New York, NY, USA: Springer-Verlag, 2008.
- [6] Xu L, Neufeld J, Larson B, Schuurmans D. Maximum margin clustering. In *Advances in Neural Information Processing Systems 17*, Saul L, Weiss Y, Bottou L (eds.), MIT Press, 2005, pp.1537-1544.
- [7] Pahikkala T, Airola A, Gieseke F, Kramer O. Unsupervised multi-class regularized least-squares classification. In *Proc. the 12th IEEE International Conference on Data Mining (ICDM 2012)*, Dec. 2012, pp.585-594.
- [8] Boyd S, Vandenberghe L. Convex Optimization. New York, NY, USA: Cambridge University Press, 2004.
- [9] Valizadegan H, Jin R. Generalized maximum margin clustering and unsupervised kernel learning. In *Advances in Neural Information Processing Systems 19*, Schölkopf B, Platt J, Hoffman T (eds.), MIT Press, 2007, pp.1417-1424.
- [10] Zhao B, Wang F, Zhang C. Efficient maximum margin clustering via cutting plane algorithm. In *Proc. the SIAM International Conference on Data Mining*, Apr. 2008, pp.751-762.
- [11] Li Y, Tsang I, Kwok J, Zhou Z. Tighter and convex maximum margin clustering. In *Proc. the 12th International Conference on Artificial Intelligence and Statistics*, Apr. 2009, pp.344-351.
- [12] Zhang K, Tsang I, Kwok J. Maximum margin clustering made practical. In *Proc. the 24th International Conference on Machine Learning*, June 2007, pp.1119-1126.
- [13] Gieseke F, Pahikkala T, Kramer O. Fast evolutionary maximum margin clustering. In *Proc. the 26th International Conference on Machine Learning*, June 2009, pp.361-368.
- [14] Zhao B, Wang F, Zhang C. Efficient multiclass maximum margin clustering. In *Proc. the 25th International Conference on Machine Learning*, July 2008, pp.1248-1255.
- [15] Xu L, Schuurmans D. Unsupervised and semi-supervised multi-class support vector machines. In *Proc. the 20th National Conference on Artificial Intelligence*, July 2005, pp.904-910.
- [16] Rifkin R, Klautau A. In defense of one-vs-all classification. *Journal of Machine Learning Research*, 2004, 5(Jan): 101-141.
- [17] Rifkin R, Yeo G, Poggio T. Regularized least-squares classification. In *Advances in Learning Theory: Methods, Models and Applications*, Suykens J, Horvath G, Basu S, Micchelli C, Vandewalle J (eds.), Amsterdam, The Netherlands: IOS Press, 2003, pp.131-154.
- [18] Kimeldorf G, Wahba G. Some results on Tchebycheffian spline functions. *Journal of Mathematical Analysis and Applications*, 1971, 33(1): 82-95.
- [19] Girosi F, Jones M, Poggio T. Regularization theory and neural networks architectures. *Neural Computation*, 1995, 7(2): 219-269.
- [20] Russell S, Norvig P. Artificial Intelligence: A Modern Approach (3rd edition). Upper Saddle River, NJ, USA: Prentice Hall Press, 2009.
- [21] Kirkpatrick S, Gelatt C, Vecchi M. Optimization by simulated annealing. *Science*, 1983, 220(4598): 671-680.
- [22] Arthur D, Vassilvitskii S. *k*-means++: The advantages of careful seeding. In *Proc. the 18th Annual ACM-SIAM Symposium on Discrete Algorithms*, Jan. 2007, pp.1027-1035.

- [23] Dempster A, Laird N, Rubin D. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 1977, 39(1): 1-38.
- [24] Shi J, Malik J. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2000, 22(8): 888-905.
- [25] Schölkopf B, Mika S, Burges C, Knirsch P, Müller K R, Rätsch G, Smola A. Input space versus feature space in kernel-based methods. *IEEE Transactions on Neural Networks*, 1999, 10(5): 1000-1017.
- [26] Nene S, Nayar S, Murase H. Columbia object image library (COIL-100). Technical Report, CUCS-006-96, Department of Computer Science, Columbia University, 1996.
- [27] Hubert L, Arabie P. Comparing partitions. *Journal of Classification*, 1985, 2(1): 193-218.
- [28] Wang F, Zhao B, Zhang C. Linear time maximum margin clustering. *IEEE Transactions on Neural Networks*, 2010, 21(2): 319-332.
- [29] Waegeman W, Verwaeren J, Slabbinck B, De Baets B. Supervised learning algorithms for multi-class classification problems with partial class memberships. *Fuzzy Sets and Systems*, 2011, 184(1): 106-125.
- [30] Williams C, Seeger M. Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems 13*, Leen T, Dietterich T, Tresp V (eds.), MIT Press, 2001, pp.682-688.
- [31] Halkidi M, Batistakis Y, Vazirgiannis M. On clustering validation techniques. *Journal of Intelligent Information Systems*, 2001, 17(2/3): 107-145.
- [32] Zhao Q. Cluster validity in clustering methods [Ph.D. Thesis]. University of Eastern Finland, 2012.



**Tapio Pahikkala** currently acts as a professor of intelligent systems in University of Turku, Finland. He received his Ph.D. degree in computer science from University of Turku in 2008 and held an Academy of Finland postdoctoral position during 2010~2012. His research focuses on machine learning, pattern recognition, algorithmics, and computational intelligence.

He has authored more than 80 peer-reviewed scientific publications and served in program committees of numerous scientific conferences. He is a member of both ACM and IEEE, and has served as a member in several committees of IEEE societies.



**Antti Airola** received his M.Sc. degree in software engineering and the D.Sc. (Tech.) degree in information and communication technology from the University of Turku, Finland, in 2006 and 2011 respectively. He is currently a university teacher at the University of Turku. His research interests include both basic and applied research in machine learning,

with focus on regularized kernel methods.



**Fabian Gieseke** received his Diploma degrees in mathematics and computer science from the University of Münster, Germany, and his Ph.D. degree in computer science from the Carl von Ossietzky University of Oldenburg, Germany. He is currently working as a postdoctoral researcher at the University of Copenhagen, Denmark. His research

interests include support vector machines and their extensions to semi- and unsupervised learning settings as well as large-scale applications in astronomy.



**Oliver Kramer** is a junior professor for computational intelligence at the University of Oldenburg in Germany. His main research interests are data mining, evolutionary optimization, and the application of computational intelligence techniques to renewable energy systems. He received a Ph.D. degree in computer science from the University of

Paderborn, Germany, in 2008. After a postdoc stay at the TU Dortmund, Germany, from 2007 to 2009, and the International Computer Science Institute in Berkeley (USA) in 2010, he became a junior professor at the Bauhaus University Weimar, Germany. Since August 2011 he is affiliated to the Department of Computing Science at the Carl von Ossietzky University of Oldenburg.

## Appendix

In Section 3 we gave an intuitive description of the proposed search algorithm, and presented a claim about its overall computational complexity. Here, we show in detail how the claimed complexity can be achieved. The consideration can be divided into the following three fundamental parts:

1) *Initialization of Cache Memories*. Before starting the actual search, certain cache memories have to be initialized that the subsequent parts will take advantage of them.

2) *Computation of the Steepest Descent Directions*. The proposed search algorithm computes these directions before deciding, for which data point the cluster label should be switched next.

3) *Update of the Caches*. Whenever a cluster label of a data point is switched, the cache memories have to be updated in order to maintain the ability to compute the steepest descent directions efficiently.

We go through these phases one by one before summarizing them in the proof of Theorem 1.

### A.1 Initialization of Cache Memories

First, we reformulate the objective function of the regularized least-squares framework. Let  $\mathbf{K} = \mathbf{V}\mathbf{A}\mathbf{V}^T$  be the eigen decomposition of the kernel matrix, let  $\tilde{\mathbf{A}} = (\mathbf{A} + \lambda\mathbf{I})^{-1}$ , and let  $F(\mathbf{y})$  be defined as in (7). We can prove the following:

**Lemma A1.**

$$F(\mathbf{y}) = 1 - \mathbf{y}^T \mathbf{V} \mathbf{A} \tilde{\mathbf{A}} \mathbf{V}^T \mathbf{y}.$$

*Proof.* Using standard linear algebra techniques, we obtain the following decomposition

$$\begin{aligned} F(\mathbf{y}) &= (\mathbf{y} - \mathbf{K}\mathbf{G}\mathbf{y})^T (\mathbf{y} - \mathbf{K}\mathbf{G}\mathbf{y}) + \lambda \mathbf{y}^T \mathbf{G} \mathbf{K} \mathbf{G} \mathbf{y} \\ &= \mathbf{y}^T \mathbf{V} \left( \mathbf{I} - 2\mathbf{A}\tilde{\mathbf{A}} + \mathbf{A}^2\tilde{\mathbf{A}}^2 + \lambda\mathbf{A}\tilde{\mathbf{A}}^2 \right) \mathbf{V}^T \mathbf{y} \\ &= \mathbf{y}^T \mathbf{V} \left( \mathbf{I} + (-2\mathbf{I} + \tilde{\mathbf{A}}\mathbf{A} + \lambda\tilde{\mathbf{A}})\mathbf{A}\tilde{\mathbf{A}} \right) \mathbf{V}^T \mathbf{y} \\ &= \mathbf{y}^T \mathbf{V} \left( \mathbf{I} + (-2\mathbf{I} + \tilde{\mathbf{A}}(\mathbf{A} + \lambda\mathbf{I}))\mathbf{A}\tilde{\mathbf{A}} \right) \mathbf{V}^T \mathbf{y} \\ &= \mathbf{y}^T \mathbf{V} \left( \mathbf{I} + (-2\mathbf{I} + \mathbf{I})\mathbf{A}\tilde{\mathbf{A}} \right) \mathbf{V}^T \mathbf{y} \\ &= \mathbf{y}^T \mathbf{V} \left( \mathbf{I} - \mathbf{A}\tilde{\mathbf{A}} \right) \mathbf{V}^T \mathbf{y} \\ &= 1 - \mathbf{y}^T \mathbf{V} \mathbf{A} \tilde{\mathbf{A}} \mathbf{V}^T \mathbf{y}. \end{aligned} \quad (\text{A1})$$

□

Given the kernel matrix containing all pairwise kernel evaluations between the training data points, the computation of the compact decomposition, in which only the eigen vectors corresponding to the nonzero eigenvalues are computed, requires  $O(r^2n)$  time, where  $r$  is the rank of the kernel matrix. In the compact decomposition, the matrix  $\mathbf{V} \in \mathbb{R}^{n \times r}$  contains the  $r$  eigenvectors and  $\mathbf{A} \in \mathbb{R}^{r \times r}$  is a diagonal matrix containing the  $r$  nonzero eigenvalues. It is also worth pointing out that the eigen decomposition of the kernel matrix is often used to turn the kernel-based clustering setting into a linear one (as is done in the multi-class MMC experiments by, e.g., Zhao et al.<sup>[14]</sup>); it therefore forms a common computational bottleneck for the kernel-based competitors considered in our experimental evaluation.

**Assumption A1.** Assume that we are given the compact eigen decomposition of the kernel matrix  $\mathbf{K} = \mathbf{V}\mathbf{A}\mathbf{V}^T$ , where  $\mathbf{V} \in \mathbb{R}^{n \times r}$ ,  $\mathbf{A} \in \mathbb{R}^{r \times r}$ , and  $r$  is the rank of the kernel matrix, as well as an initial vector of class labels  $\mathbf{c} \in \mathbb{C}^n$ . In the initialization phase, we prepare the following cache memories which are updated whenever the vector of class labels is changed:

- the  $n \times n$ -matrix

$$\mathbf{R} = \mathbf{V} \mathbf{A} \tilde{\mathbf{A}} \mathbf{V}^T, \quad (\text{A2})$$

- the vectors  $\mathbf{R}p_c(\mathbf{c}), \forall c \in \mathbb{C}$ ,
- as well as the values  $Q(\mathbf{c})$  and  $F(p_h(\mathbf{c})), \forall h \in \mathbb{C}$ .

The computational complexity of the initialization phase is dominated by the computation of  $\mathbf{R}$ , which can be done in  $O(rn^2)$  time given the compact decomposition of the kernel matrix of rank  $r$ .

### A.2 Computation of the Steepest Descent Directions

The next lemma concerns the efficient computation of  $S(\mathbf{c}, j, d)$ , given that certain intermediate results have already been computed and cached. Its proof also encompasses implementation details of the search algorithms that take advantage of the computational shortcuts.

**Lemma A2.** Assume that we have cache memories given in Assumption A1 available. Then, the value of  $S(\mathbf{c}, j, d)$  can be computed in a constant time.

*Proof.* Let  $\mathbf{y} \in \{-1, 1\}^n$  and let us denote  $\hat{\mathbf{y}} = \mathbf{y} - 2y_j \mathbf{e}^j$ , that is,  $\mathbf{y}$  and  $\hat{\mathbf{y}}$  are two  $\pm 1$ -valued vectors differing from each other only by the  $j$ -th entry. Moreover, we denote  $\mathbf{t} = \mathbf{R}\mathbf{y}$ . Then, continuing from (A1), if we already know  $F(\mathbf{y})$ ,  $\mathbf{R}$ , and  $\mathbf{t}$ , the value of  $F$  for  $\hat{\mathbf{y}}$  can be computed from

$$\begin{aligned} F(\hat{\mathbf{y}}) &= n - \hat{\mathbf{y}}^T \mathbf{R} \hat{\mathbf{y}} \\ &= n - (\mathbf{y} - 2y_j \mathbf{e}^j)^T \mathbf{R} (\mathbf{y} - 2y_j \mathbf{e}^j) \\ &= n - \mathbf{y}^T \mathbf{R} \mathbf{y} + 4y_j \mathbf{y}^T \mathbf{R} \mathbf{e}^j - 4y_j^2 \mathbf{e}^{jT} \mathbf{R} \mathbf{e}^j \\ &= F(\mathbf{y}) + 4y_j \mathbf{e}^{jT} \mathbf{R} \mathbf{y} - 4\mathbf{e}^{jT} \mathbf{R} \mathbf{e}^j \\ &= F(\mathbf{y}) + 4y_j t_j - 4R_{j,j}. \end{aligned}$$

Moreover, we observe that we can define a simple formula for calculating the difference in the objective values if a single entry of the  $\pm 1$ -valued vector is flipped as follows:

$$D(\mathbf{y}, j) = F(\hat{\mathbf{y}}) - F(\mathbf{y}) = 4y_j t_j - 4R_{j,j}. \quad (\text{A3})$$

Putting together (12), (13), and (A3),  $S(\mathbf{c}, j, d)$  can be formally written as

$$S(\mathbf{c}, j, d) = Q(\mathbf{c}) + D(p_{c_j}(\mathbf{c}), j) + D(p_d(\mathbf{c}), j).$$

The formula contains the objective value adjustments of both the old cluster  $c_j$ , and the new one  $d$  of the aggregate objective function (13) that allows the cluster assignments of a single data point to have only one positive entry. Thus, given the assumptions about cached intermediate results, we can calculate the objective value change caused by switching a single entry of the cluster label vector  $\mathbf{c}$  in  $O(1)$  time. □

### A.3 Updates of the Caches

After the steepest descent direction is found and the cluster label of the corresponding patterns is switched,

part of the cache memories given in Assumption A1 have to be updated accordingly in order to maintain the ability to perform fast searches. As shown in the following lemma, the update operation does now slow down the computation of the steepest descent directions.

**Lemma A3.** *The cache memories given in Assumption A1 can be updated in  $O(n)$  time after the cluster label of a single pattern is switched.*

*Proof.* Given that  $\mathbf{R}p_{c_j}(\mathbf{c})$  is stored in memory, the vector  $\mathbf{R}(p_{c_j}(\mathbf{c}) - 2\mathbf{e}^j)$  can be obtained from

$$\mathbf{R}(p_{c_j}(\mathbf{c}) - 2\mathbf{e}^j) = \mathbf{R}p_{c_j}(\mathbf{c}) - 2(\mathbf{R}_j)^T$$

in  $O(n)$  time. The vector  $\mathbf{R}(p_d(\mathbf{c}) + 2\mathbf{e}^j)$  can be computed analogously. The values  $Q(\mathbf{c})$ ,  $F(p_{c_j}(\mathbf{c}))$ , and  $F(p_d(\mathbf{c}))$  are obtained in a constant time as implied by the proof of Lemma A2. The matrix  $\mathbf{R}$  does not depend on  $\mathbf{c}$  and, thus, does not have to be updated.  $\square$

#### A.4 Runtime Proof

Putting everything together, we arrive to the proof of Theorem 1.

*Proof of Theorem 1.* As shown in Lemma A2, the evaluation of  $S(\mathbf{c}, j, d)$  can be performed in constant time by taking advantage of the cache memories defined in Assumption A1. Since there are  $|\mathcal{C}|$  clusters and  $n$  data points, finding the steepest descent direction requires  $O(|\mathcal{C}|n)$  time. Lemma A3 in turn shows that the cache memories can be updated according to the steepest descent direction in  $O(n)$ , but this is dominated by the time required for finding the next steepest descent direction. Altogether, Algorithm 3 performs  $O(s|\mathcal{C}|n^2)$  calls for the function  $S(\mathbf{c}, j, d)$ , each with a constant time complexity, and the initialization of the caches requires  $O(n^2r)$  time.  $\square$