



## Indefinite kernel spectral learning

Siamak Mehrkanoon<sup>a</sup>, Xiaolin Huang<sup>b,\*</sup>, Johan A.K. Suykens<sup>a</sup>

<sup>a</sup> Department of Electrical Engineering (ESAT-STADIUS), KU Leuven, Kasteelpark Arenberg 10, Leuven B-3001, Belgium

<sup>b</sup> Institute of Image Processing and Pattern Recognition, and the MOE Key Laboratory of System Control and Information Processing, Shanghai Jiao Tong University, Shanghai 200240, PR China

### ARTICLE INFO

#### Article history:

Received 21 February 2017

Revised 18 October 2017

Accepted 14 January 2018

Available online 3 February 2018

#### Keywords:

Semi-supervised learning

Scalable models

Indefinite kernels

Kernel spectral clustering

Low embedding dimension

### ABSTRACT

The use of indefinite kernels has attracted many research interests in recent years due to their flexibility. They do not possess the usual restrictions of being positive definite as in the traditional study of kernel methods. This paper introduces the indefinite unsupervised and semi-supervised learning in the framework of least squares support vector machines (LS-SVM). The analysis is provided for both unsupervised and semi-supervised models, i.e., Kernel Spectral Clustering (KSC) and Multi-Class Semi-Supervised Kernel Spectral Clustering (MSS-KSC). In indefinite KSC models one solves an eigenvalue problem whereas indefinite MSS-KSC finds the solution by solving a linear system of equations. For the proposed indefinite models, we give the feature space interpretation, which is theoretically important, especially for the scalability using Nyström approximation. Experimental results on several real-life datasets are given to illustrate the efficiency of the proposed indefinite kernel spectral learning.

© 2018 Elsevier Ltd. All rights reserved.

### 1. Introduction

Kernel-based learning models have shown great success in various application domains [1–3]. Traditionally, kernel learning is restricted to positive semi-definite (PSD) kernels as the properties of Reproducing Kernel Hilbert Spaces (RKHS) are well explored. However, many positive semi-definite kernels such as the sigmoid kernel [4] remain positive semi-definite only when their associated parameters are within a certain range, otherwise they become non-positive definite [5]. Moreover, the positive definite kernels are limited in some problems due to the need of non-Euclidean distances [6,7]. For instance in protein similarity analysis, the protein sequence similarity measures require learning with a non-PSD similarity matrix [8].

The need of using indefinite kernels in machine learning methods attracted many research interests on indefinite learning in both theory and algorithm. Theoretical discussions are mainly on Reproducing Kernel Kreĭn Spaces (RKKS, [9,10]), which is different to the RKHS for PSD kernels. In algorithm design, a lot of attempts have been made to cope with indefinite kernels by regularizing the non-positive definite kernels to make them positive semi-definite [11–14]. It is also possible to directly use an indefinite kernel in e.g., support vector machine (SVM) [4]. Though an indefinite ker-

nel makes the problem non-convex, it is still possible to get a local optimum as suggested by Lin and Lin [15]. One important issue is that kernel trick is no longer valid when an indefinite kernel is applied in SVM and one needs new feature space interpretations to explain the effectiveness of SVM with indefinite kernels. The interpretation is usually about a pseudo-Euclidean (pE) space, which is a product of two Euclidean vector spaces, as analyzed in [10,16]. Notice that “indefinite kernels” literally covers asymmetric ones and complex ones. But this paper restricts “indefinite kernel” to the kernels that correspond to real symmetric indefinite matrices, which is consistent to the existing literature on indefinite kernel.

Indefinite kernels are also applicable to the least squares support vector machines [17]. In LS-SVM, one solves a linear system of equations in the dual and the optimization problem itself has no additional requirement on the positiveness of the kernel. In other words, even if an indefinite kernel is used in the dual formulation of LS-SVM, it is still convex and easy to solve, which is different from indefinite kernel learning with SVM. However, like in SVM, using an indefinite kernel in LS-SVM loses the traditional interpretation of the feature space and a new formulation has been recently discussed in [18].

Motivated by the success of indefinite learning for some supervised learning tasks, we in this paper introduce indefinite similarities to unsupervised as well as semi-supervised models that can learn from both labeled and unlabeled data instances. There have been already many efficient semi-supervised models, such as

\* Corresponding author.

E-mail addresses: [siamak.mehrkanoon@esat.kuleuven.be](mailto:siamak.mehrkanoon@esat.kuleuven.be) (S. Mehrkanoon), [xiaolinhuang@sjtu.edu.cn](mailto:xiaolinhuang@sjtu.edu.cn) (X. Huang), [johan.suykens@esat.kuleuven.be](mailto:johan.suykens@esat.kuleuven.be) (J.A.K. Suykens).

Laplacian support vector machine [19], which assumes that neighboring point pairs with a large weight edge are most likely within the same cluster. However, to the best of our knowledge, there is no work that extends unsupervised/semi-supervised learning to indefinite kernels.

Since using indefinite kernels in the framework of LS-SVM does not change the training problem, here we focus on multi-class semi-supervised kernel spectral clustering (MSS-KSC) model proposed by Mehrkanoon et al. [20]. MSS-KSC model and its extensions for analyzing large-scale data, data streams as well as multi-label datasets are discussed in [21–23] respectively. When one of the regularization parameters is set to zero, MSS-KSC becomes the kernel spectral clustering (KSC), which is an unsupervised learning algorithm introduced by Alzate and Suykens [24]. It is a special case of MSS-KSC. Due to the link to LS-SVM, it can be expected and also will be shown here that MSS-KSC with indefinite similarities are still easy to solve. However, the kernel trick is no longer valid and we have to find corresponding feature space interpretations. The purpose of this paper is to introduce indefinite kernels for semi-supervised learning as well as unsupervised learning as a special case. Specifically, we propose indefinite kernels in MSS-KSC and KSC models. Subsequently, we derive their feature space interpretation. Besides of theoretical interests, the interpretation allows us to develop algorithms based on Nyström approximation for large-scale problems.

The paper is organized as follows. Section 2 briefly reviews the MSS-KSC with PSD kernel. In Section 3, the MSS-KSC with an indefinite kernel is derived and the interpretation of the feature map is provided. As a special case of MSS-KSC, the KSC with an indefinite kernel and its feature interpretation is discussed in Section 4. In Section 5, we discuss the scalability of the indefinite KSC/MSS-KSC model on large-scale problems. The experimental results are given in Section 6 to confirm the validity and applicability of the proposed model on several real life small and large-scale datasets. Section 7 ends the paper with a brief conclusion.

## 2. MSS-KSC with PSD kernel

Consider training data

$$D = \underbrace{\{x_1, \dots, x_{n_{UL}}\}}_{(D_U)} \cup \underbrace{\{x_{n_{UL}+1}, \dots, x_n\}}_{(D_L)}, \quad (1)$$

where  $\{x_i\}_{i=1}^n \in \mathbb{R}^d$ . The first  $n_{UL}$  points do not have labels whereas the last  $n_L = n - n_{UL}$  points have been labeled. Assume that there are  $Q$  classes ( $Q \leq N_c$ ), then the label indicator matrix  $Y \in \mathbb{R}^{n_L \times Q}$  is defined as follows:

$$Y_{ij} = \begin{cases} +1 & \text{if the } i\text{th point belongs to the } j\text{th class,} \\ -1 & \text{otherwise.} \end{cases} \quad (2)$$

The primal formulation of multi-class semi-supervised KSC (MSS-KSC) described by Mehrkanoon et al. [20] is given as follows:

$$\min_{w^{(\ell)}, b^{(\ell)}, e^{(\ell)}} \frac{1}{2} \sum_{\ell=1}^Q w^{(\ell)T} w^{(\ell)} - \frac{\gamma_1}{2} \sum_{\ell=1}^Q e^{(\ell)T} V e^{(\ell)} + \frac{\gamma_2}{2} \sum_{\ell=1}^Q (e^{(\ell)} - c^{(\ell)})^T \tilde{A} (e^{(\ell)} - c^{(\ell)}) \quad (3)$$

subject to  $e^{(\ell)} = \Phi w^{(\ell)} + b^{(\ell)} \mathbf{1}_n$ ,  $\ell = 1, \dots, Q$ ,

where  $c^{(\ell)}$  is the  $\ell$ th column of the matrix  $C$  defined as

$$C = [c^{(1)}, \dots, c^{(Q)}]_{n \times Q} = \begin{bmatrix} \mathbf{0}_{n_{UL} \times Q} \\ Y \end{bmatrix}_{n \times Q}. \quad (4)$$

Here

$$\Phi = [\varphi(x_1), \dots, \varphi(x_n)]^T \in \mathbb{R}^{n \times h}$$

where  $\varphi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^h$  is the feature map and  $h$  is the dimension of the feature space which can be infinite dimensional.  $\mathbf{0}_{n_{UL} \times Q}$  is a zero matrix of size  $n_{UL} \times Q$ ,  $Y$  is defined previously, and the right hand of (4) is a matrix consisting of  $\mathbf{0}_{n_{UL} \times Q}$  and  $Y$ . The matrix  $\tilde{A}$  is defined as follows:

$$\tilde{A} = \begin{bmatrix} \mathbf{0}_{n_{UL} \times n_{UL}} & \mathbf{0}_{n_{UL} \times n_L} \\ \mathbf{0}_{n_L \times n_{UL}} & I_{n_L \times n_L} \end{bmatrix},$$

where  $I_{n_L \times n_L}$  is the identity matrix of size  $n_L \times n_L$ .  $V$  is the inverse of the degree matrix defined as follows:

$$V = D^{-1} = \text{diag}\left(\frac{1}{d_1}, \dots, \frac{1}{d_n}\right),$$

where  $d_i = \sum_{j=1}^n K(x_i, x_j)$  is the degree of the  $i$ th data point.

As stated in [20], the object function in the formulation (3), contains three terms. The first two terms together with the set of constraints correspond to a weighted kernel PCA formulation in the least squares support vector machine framework given in [24] which is shown to be suitable for clustering and is referred to as kernel spectral clustering (KSC) algorithm. The last regularization term in (3) aims at minimizing the squared distance between the projections of the labeled data and their corresponding labels. This term enforces the projections of the labeled data points to be as close as possible to the true labels. Therefore by incorporating the labeled information, the pure clustering KSC model is guided so that it respects the provided labels by not misclassifying them. In this way, one could learn from both labeled and unlabeled instances. In addition thanks to introduced model selection scheme in [20], the MSS-KSC model is also equipped with the out-of-sample extension property to predict the labels of unseen instances.

It should be noted that, ignoring the last regularization term, or equivalently setting  $\gamma_2 = 0$  and  $Q = N_c - 1$ , reduces the MSS-KSC formulation to kernel spectral clustering (KSC) described in [24]. Therefore, KSC formulation in the primal can be covered as a special case of MSS-KSC formulation. As illustrated by Mehrkanoon et al. [20], given  $Q$  labels the approach is not restricted to finding just  $Q$  classes and instead is able to discover up to  $2^Q$  hidden clusters. In addition, it uses a low embedding dimension to reveal the existing number of clusters which is important when one deals with large number of clusters.

When the feature map  $\varphi$  in (3) is not explicitly known, in the context of PSD kernel, one may use the kernel trick and solve the problem in the dual. Elimination of the primal variables  $w^{(\ell)}$ ,  $e^{(\ell)}$  and making use of Mercer's Theorem result in the following linear system in the dual [20]:

$$\gamma_2 \left( I_n - \frac{R \mathbf{1}_n \mathbf{1}_n^T}{\mathbf{1}_n^T R \mathbf{1}_n} \right) c^{(\ell)} = \alpha^{(\ell)} - R \left( I_n - \frac{\mathbf{1}_n \mathbf{1}_n^T R}{\mathbf{1}_n^T R \mathbf{1}_n} \right) \Omega \alpha^{(\ell)}, \quad (5)$$

where  $R = \gamma_1 V - \gamma_2 \tilde{A}$ . In (5), there are two coefficients, namely  $\gamma_1$  and  $\gamma_2$ , which reflect the emphasis on unlabeled and labeled samples, respectively, as shown in (3). Besides, there could be one or multiple parameters in the kernel. All of these parameters could be tuned by cross-validation.

## 3. MSS-KSC with indefinite kernel

Traditionally, the kernel used in MSS-KSC is restricted to be positive semi-definite. When the kernel in (5) is indefinite, one still requires to solve a linear system of equations. However, the feature space has different interpretations compared to definite kernels. In what follows we establish and analyze the feature space interpretations for MSS-KSC.

**Theorem 3.1.** Suppose that for a symmetric but indefinite kernel matrix  $K$ , the solution of the linear system (5) is denoted by  $[\alpha^*, b^*]^T$ .

Then there exist two feature mappings  $\varphi_1$  and  $\varphi_2$ , which correspond to the matrices  $\Phi_1$  and  $\Phi_2$ , respectively, such that

$$w_1^{(\ell)} = \sum_{i=1}^n \alpha_{*,i}^{(\ell)} \varphi_1(x_i), \ell = 1, \dots, Q, \quad (6)$$

and

$$w_2^{(\ell)} = \sum_{i=1}^n \alpha_{*,i}^{(\ell)} \varphi_2(x_i), \ell = 1, \dots, Q, \quad (7)$$

which is a stationary point of the following primal problem:

$$\begin{aligned} \min_{w_1^{(\ell)}, w_2^{(\ell)}, b_*^{(\ell)}, e^{(\ell)}} & \frac{1}{2} \sum_{\ell=1}^Q w_1^{(\ell)T} w_1^{(\ell)} - \frac{1}{2} \sum_{\ell=1}^Q w_2^{(\ell)T} w_2^{(\ell)} \\ & + \frac{\gamma_2}{2} \sum_{\ell=1}^Q (e^{(\ell)} - c^{(\ell)})^T \tilde{A} (e^{(\ell)} - c^{(\ell)}) - \frac{\gamma_1}{2} \sum_{\ell=1}^Q e^{(\ell)T} V e^{(\ell)} \\ \text{subject to} & e^{(\ell)} = \Phi_1 w_1^{(\ell)} + \Phi_2 w_2^{(\ell)} + b_*^{(\ell)} \mathbf{1}_n, \ell = 1, \dots, Q. \end{aligned} \quad (8)$$

Then, the dual problem of (8) is given in (5), with the kernel matrix  $\Omega$  defined as follows,

$$\Omega_{i,j} = K_1(x_i, x_j) - K_2(x_i, x_j), \quad (9)$$

where,  $K_1(x_i, x_j)$  and  $K_2(x_i, x_j)$  are two PSD kernels.

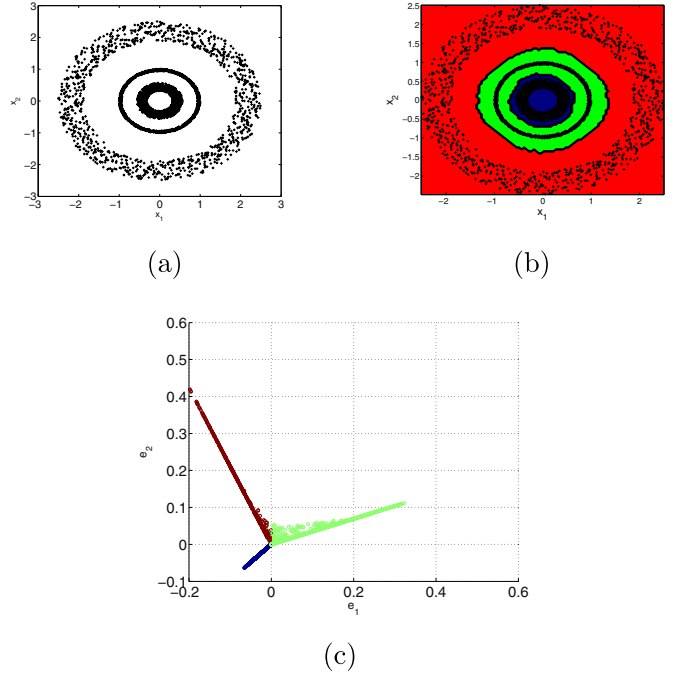
**Proof.** The Lagrangian of the constrained optimization problem (8) becomes

$$\begin{aligned} \mathcal{L}(w_1^{(\ell)}, w_2^{(\ell)}, b_*^{(\ell)}, e^{(\ell)}, \alpha_*^{(\ell)}) &= \frac{1}{2} \sum_{\ell=1}^Q w_1^{(\ell)T} w_1^{(\ell)} - \sum_{\ell=1}^Q w_2^{(\ell)T} w_2^{(\ell)} \\ & - \frac{\gamma_1}{2} \sum_{\ell=1}^Q e^{(\ell)T} V e^{(\ell)} \\ & + \frac{\gamma_2}{2} \sum_{\ell=1}^Q (e^{(\ell)} - c^{(\ell)})^T \tilde{A} (e^{(\ell)} - c^{(\ell)}) \\ & + \sum_{\ell=1}^Q \alpha_*^{(\ell)T} \left( e^{(\ell)} - \Phi_1 w_1^{(\ell)} \right. \\ & \left. - \Phi_2 w_2^{(\ell)} - b_*^{(\ell)} \mathbf{1}_n \right), \end{aligned}$$

where  $\alpha_*^{(\ell)}$  is the vector of Lagrange multipliers. Then the KKT optimality conditions are as follows,

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w_1^{(\ell)}} = 0 \rightarrow w_1^{(\ell)} = \Phi_1^T \alpha_*^{(\ell)}, \ell = 1, \dots, Q, \\ \frac{\partial \mathcal{L}}{\partial w_2^{(\ell)}} = 0 \rightarrow w_2^{(\ell)} = -\Phi_2^T \alpha_*^{(\ell)}, \ell = 1, \dots, Q, \\ \frac{\partial \mathcal{L}}{\partial b_*^{(\ell)}} = 0 \rightarrow \mathbf{1}_n^T \alpha_*^{(\ell)} = 0, \ell = 1, \dots, Q, \\ \frac{\partial \mathcal{L}}{\partial e^{(\ell)}} = 0 \rightarrow \alpha_*^{(\ell)} = (\gamma_1 V - \gamma_2 \tilde{A}) e^{(\ell)} + \gamma_2 c^{(\ell)}, \ell = 1, \dots, Q, \\ \frac{\partial \mathcal{L}}{\partial \alpha_*^{(\ell)}} = 0 \rightarrow e^{(\ell)} = \Phi_1 w_1^{(\ell)} + \Phi_2 w_2^{(\ell)} + b_*^{(\ell)} \mathbf{1}_n, \ell = 1, \dots, Q. \end{cases} \quad (10)$$

Elimination of the primal variables  $w_1^{(\ell)}, w_2^{(\ell)}, e^{(\ell)}$  and making use of the kernel trick ( $\Omega_1 = \Phi_1^T \Phi_1$  and  $\Omega_2 = \Phi_2^T \Phi_2$ ) lead to the linear system of equations in the dual defined in (5) with the indefinite kernel matrix defined in (9). With  $\alpha_*$  obtained from (5), the weight vectors  $w_1^{(\ell)}$  and  $w_2^{(\ell)}$  defined in (6) and (7), satisfy the first-order optimality condition of (8).  $\square$



**Fig. 1.** Illustrating the performance of KSC model with an indefinite kernel (TL1 kernel) on synthetic three concentric example. (a) Original data. (b) The predicted memberships obtained using indefinite KSC model with  $\mu = 0.4$ . (c) The line structure of the score variables,  $e$ , indicating the good generalization performance of indefinite KSC model with  $\mu = 0.4$ .

One can show that from the third KKT optimality condition, the bias term is determined by

$$b_*^{(\ell)} = (1/\mathbf{1}_n^T R \mathbf{1}_n) (-\mathbf{1}_n^T \gamma_2 c^{(\ell)} - \mathbf{1}_n^T R \Omega \alpha_*^{(\ell)}), \ell = 1, \dots, Q, \quad (11)$$

where  $R$  is defined as in (5). Once the solution vector and the bias term are obtained, one can use the out-of-sample extension property of the model to predict the score variables of the unseen test instances as follows:

$$e_{\text{test}}^{(\ell)} = \Omega \alpha_*^{(\ell)} + b_*^{(\ell)}, \ell = 1, \dots, Q. \quad (12)$$

The above discussion gives the feature space interpretation for indefinite MSS-KSC. The discussion in a pE space is similar to indefinite SVM; see, [10,16,18]. The main difference from learning algorithms for PSD kernels is that the indefinite learning is to minimize a pseudo-distance. The readers are referred to Fig. 1 in [16], which gives a clear geometric explanation for the distance in a pE space.

In practice, the performance of the MSS-KSC model depends on the choice of the parameters. In this aspect, there is no difference between a PSD kernel and an indefinite kernel. Therefore the following model selection scheme introduced in [20] for MSS-KSC can be employed:

$$\max_{\gamma_1, \gamma_2, \mu} \eta \text{Sil}(\gamma_1, \gamma_2, \mu) + (1 - \eta) \text{Acc}(\gamma_1, \gamma_2, \mu). \quad (13)$$

It is a combination of Silhouette (Sil) and classification accuracy (Acc).  $\eta \in [0, 1]$  is a user-defined parameter that controls the trade off between the importance given to unlabeled and labeled instances. The MSS-KSC algorithm with an indefinite kernel is summarized in Algorithm 1. One can note that the main difference with respect to Algorithm 1 discussed in [20] is at the level employing the indefinite kernel and all the other steps remain unchanged.

**Algorithm 1** Indefinite kernel in multi-class semi-supervised classification model.

- 1: **Input:** Training data set  $\mathcal{D}$ , labels  $Z$ , tuning parameters  $\{\gamma_i\}_{i=1}^2$ , kernel parameter  $\mu$ , test set  $\mathcal{D}^{\text{test}} = \{x_i^{\text{test}}\}_{i=1}^{N_{\text{test}}}$  and codebook  $CB = \{c_q\}_{q=1}^Q$
- 2: **Output:** Class membership of test data  $\mathcal{D}^{\text{test}}$
- 3: Construct the indefinite kernel matrix  $\Omega$  (see (9)).
- 4: Solve the dual linear system (5) with the indefinite kernel matrix  $\Omega$  to obtain  $\{\alpha^\ell\}_{\ell=1}^Q$  and compute the bias term  $\{b_*^\ell\}_{\ell=1}^Q$  using (11).
- 5: Estimate the test data projections  $\{e_{\text{test}}^{(\ell)}\}_{\ell=1}^Q$  using (12).
- 6: Binarize the test projections and form the encoding matrix  $[\text{sign}(e_{\text{test}}^{(1)}), \dots, \text{sign}(e_{\text{test}}^{(Q)})]_{N_{\text{test}} \times Q}$  for the test points (Here  $e_{\text{test}}^{(\ell)} = [e_{\text{test},1}^{(\ell)}, \dots, e_{\text{test},N_{\text{test}}}^{(\ell)}]^T$ ).
- 7: For each  $i$ , assign  $x_i^{\text{test}}$  to class  $q^*$ , where  $q^* = \underset{q}{\text{argmin}} d_H(e_{\text{test},i}^{(\ell)}, c_q)$  and  $d_H(\cdot, \cdot)$  is the Hamming distance.

#### 4. KSC with indefinite kernels - as a special case

As a special case of MSS-KSC formulation (8), when  $\gamma_2 = 0$  and  $Q = N_c - 1$ , we obtain (17), i.e., the KSC model given by Alzate and Suykens [24]. This dual problem itself does not require the positiveness of  $\Omega$ . Thus, an indefinite kernel is applicable here and one still solves an eigenvalue problem. However, the kernel trick, which is the key to build primal-dual relationship for definite kernels, cannot be used for indefinite kernels, which follows that different feature space interpretations are needed. In this section, we establish and analyze the feature space interpretations, similar to the discussion for indefinite MSS-KSC.

**Theorem 4.1.** Suppose that the solution of the eigenvalue problem (17), in the dual, for a symmetric but indefinite kernel matrix  $K$  is denoted by  $[\alpha_*, b_*]^T$ . Then there exist two feature mappings  $\varphi_1$  and  $\varphi_2$ , such that

$$w_1^{(\ell)} = \sum_{i=1}^n \alpha_{*,i}^{(\ell)} \varphi_1(x_i), \ell = 1, \dots, N_c - 1, \quad (14)$$

and

$$w_2^{(\ell)} = \sum_{i=1}^n \alpha_{*,i}^{(\ell)} \varphi_2(x_i), \ell = 1, \dots, N_c - 1, \quad (15)$$

which is a stationary point of the following primal problem:

$$\min_{w_1^{(\ell)}, w_2^{(\ell)}, b_*^{(\ell)}, e^{(\ell)}} \frac{1}{2} \sum_{\ell=1}^{N_c-1} w_1^{(\ell)T} w_1^{(\ell)} - \frac{1}{2} \sum_{\ell=1}^{N_c-1} w_2^{(\ell)T} w_2^{(\ell)} - \frac{\gamma_1}{2} \sum_{\ell=1}^{N_c-1} e^{(\ell)T} V e^{(\ell)} \quad (16)$$

$$\text{subject to } e^{(\ell)} = \Phi_1 w_1^{(\ell)} + \Phi_2 w_2^{(\ell)} + b_*^{(\ell)} \mathbf{1}_n, \ell = 1, \dots, N_c - 1.$$

Then, the dual problem of Haasdonk (16) is given as:

$$VP_v \Omega \alpha^{(\ell)} = \lambda \alpha^{(\ell)}, \quad (17)$$

where  $\lambda = n/\gamma_\ell$ ,  $\alpha^{(\ell)}$  are the Lagrange multipliers and  $P_v$  is the weighted centering matrix:

$$P_v = I_n - \frac{1}{\mathbf{1}_n^T V \mathbf{1}_n} \mathbf{1}_n \mathbf{1}_n^T V.$$

Here  $I_n$  is the  $n \times n$  identity matrix and the kernel matrix  $\Omega$  is defined as follows,

$$\Omega_{i,j} = K_1(x_i, x_j) - K_2(x_i, x_j), \quad (18)$$

where,  $K_1(x_i, x_j)$  and  $K_2(x_i, x_j)$  are two PSD kernels.

**Proof.** It follows the proof of indefinite MSS-KSC model described in (3).  $\square$

From the link between KSC and LS-SVM, the above theorem also could be regarded as a weighted and multi-class extension of the result obtained by Huang et al. [18]. To give an intuitive idea that using indefinite kernels in KSC is possible, we show a simple example that applies the truncated  $\ell_1$  distance (TL1) kernel [25], which is indefinite and takes the following formulation,

$$K(s, t) = \max\{\mu - \|s - t\|_1, 0\}. \quad (19)$$

For this problem, one can observe that KSC with an indefinite kernel indeed can successfully cluster the points, as shown in Fig. 1. Here the Silhouette index is used for model selection (see [26] for overview of the internal clustering quality metrics).

Theorem 4.1 and Theorem 4.2 are both based on the positive decomposition of an indefinite kernel matrix  $\Omega$ : since it is a symmetric and real matrix, we can surely find two PSD matrices  $K_1$  and  $K_2$  such that

$$\Omega_{ij} = K_{1ij} - K_{2ij}.$$

For example,  $K_1$  and  $K_2$  can be constructed from the positive and negative eigenvalues of  $\Omega$ . This decomposition indicates that a PSD kernel is a special case of indefinite kernel with  $K_{2ij} = 0$ . Therefore, the use of indefinite kernel in spectral learning provides flexibility to improve the performance of PSD learning, if the kernel, which could be indefinite or definite, is suitably designed.

#### 5. Scalability

Kernel based models have shown to be successful in many machine learning tasks. However, unfortunately, many of them scale poorly with the training data size due to the need for storing and computing the kernel matrix which is usually dense.

In the context of kernel based semi-supervised learning with PSD kernels, attempts have been made to make the kernel based models scalable, see [21,27,28]. Mehrkanoon, et al. [21] introduced the Fixed-Size MSS-KSC (FS-MSS-KSC) model for classification of large-scale partially labeled instances. FS-MSS-KSC uses an explicit feature map approximated by the Nyström method [17,29] and solves the optimization problem in the primal. The finite dimensional approximation of the feature map is obtained by numerically solving a Fredholm integral equation using the Nyström discretization method which results in an eigenvalue decomposition of the kernel matrix  $\Omega$ ; see [29].

The  $i$ th component of the  $n$ -dimensional feature map  $\hat{\varphi} : \mathbb{R}^d \rightarrow \mathbb{R}^n$ , for any point  $x \in \mathbb{R}^d$ , can be obtained as follows:

$$\hat{\varphi}_i(x) = \frac{1}{\sqrt{\lambda_i^{(s)}}} \sum_{k=1}^n u_{ki} K(x_k, x), \quad (20)$$

where  $\lambda_i^{(s)}$  and  $u_i$  are eigenvalues and eigenvectors of the kernel matrix  $\Omega_{n \times n}$ . Furthermore, the  $k$ th element of the  $i$ th eigenvector is denoted by  $u_{ki}$ . In practice when  $n$  is large, we work with a subsample (prototype vectors) of size  $m \ll n$  of which the elements are selected using an entropy based criterion. In this case, the  $m$ -dimensional feature map  $\hat{\varphi} : \mathbb{R}^d \rightarrow \mathbb{R}^m$  can be approximated as follows:

$$\hat{\varphi}(x) = [\hat{\varphi}_1(x), \dots, \hat{\varphi}_m(x)]^T, \quad (21)$$

where

$$\hat{\varphi}_i(x) = \frac{1}{\sqrt{\lambda_i^{(s)}}} \sum_{k=1}^m u_{ki} K(x_k, x), i = 1, \dots, m. \quad (22)$$

Here,  $\lambda_i^{(s)}$  and  $u_i$  are the eigenvalues and eigenvectors of the constructed kernel matrix  $\Omega_{m \times m}$  with the selected prototype vectors.

When an indefinite kernel is used, the matrix  $K$  has both positive and negative eigenvalues. Thus, according to the previous feature interpretations, one can then construct two approximations for the feature maps  $\Phi_1$  and  $\Phi_2$  based on positive and negative eigenvalues, respectively. Here we give the following lemma to explain the approximation for indefinite MSS-KSC and a similar result is valid for indefinite KSC as well.

**Lemma 5.1.** *Given the  $m$ -dimensional approximation to the feature map, i.e.  $\hat{\Phi}_1 = [\hat{\varphi}(x_1), \dots, \hat{\varphi}(x_n)]^T \in \mathbb{R}^{n \times m_1}$ , and  $\hat{\Phi}_2 = [\hat{\varphi}(x_1), \dots, \hat{\varphi}(x_n)]^T \in \mathbb{R}^{n \times m_2}$ , and regularization constants  $\gamma_1, \gamma_2, \in \mathbb{R}^+$ , the solution to (8) is obtained by solving the following linear system of equations in the primal:*

$$\begin{bmatrix} \hat{\Phi}_1^T R \hat{\Phi}_1 + I_{m_1} & \hat{\Phi}_1^T R \hat{\Phi}_2 & \hat{\Phi}_1^T R 1_n \\ \hat{\Phi}_2^T R \hat{\Phi}_1 & \hat{\Phi}_2^T R \hat{\Phi}_2 - I_{m_2} & \hat{\Phi}_2^T R 1_n \\ 1_n^T R \hat{\Phi}_1 & 1_n^T R \hat{\Phi}_2 & 1_n^T R 1_n \end{bmatrix} \begin{bmatrix} w_1^{(\ell)} \\ w_2^{(\ell)} \\ b^{(\ell)} \end{bmatrix} = \gamma_2 \begin{bmatrix} \hat{\Phi}_1^T c^{(\ell)} \\ \hat{\Phi}_2^T c^{(\ell)} \\ 1_n^T c^{(\ell)} \end{bmatrix}, \quad \ell = 1, \dots, Q, \quad (23)$$

where  $R = \gamma_2 A - \gamma_1 V$  is a diagonal matrix,  $V$  and  $R$  are given previously.  $I_{m_1}$  and  $I_{m_2}$  are the identity matrix of size  $m_1 \times m_1$  and  $m_2 \times m_2$  respectively.

**Proof.** Substituting the explicit feature maps  $\hat{\Phi}_1$  and  $\hat{\Phi}_2$  into formulation (8), one can rewrite it as an unconstrained optimization problem. Subsequently setting the derivative of the cost function with respect to the primal variables  $w_1^{(\ell)}$ ,  $w_2^{(\ell)}$  and  $b^{(\ell)}$  to zero results in the linear system (23).  $\square$

The score variables evaluated at the test set  $\mathcal{D}^{\text{test}} = \{x_i\}_{i=1}^{n_{\text{test}}}$  become:

$$e_{\text{test}}^{(\ell)} = \hat{\Phi}_1^{\text{test}} w_1^{(\ell)} + \hat{\Phi}_2^{\text{test}} w_2^{(\ell)} + b^{(\ell)} 1_{n_{\text{test}}} \quad \ell = 1, \dots, Q, \quad (24)$$

where  $\hat{\Phi}_1^{\text{test}} = [\hat{\varphi}(x_1), \dots, \hat{\varphi}(x_{n_{\text{test}}})]^T \in \mathbb{R}^{n_{\text{test}} \times m_1}$  and  $\hat{\Phi}_2^{\text{test}} = [\hat{\varphi}(x_1), \dots, \hat{\varphi}(x_{n_{\text{test}}})]^T \in \mathbb{R}^{n_{\text{test}} \times m_2}$ . The decoding scheme consists of comparing the binarized score variables for test data with the codebook  $\mathcal{CB}$  and selecting the nearest codeword in terms of Hamming distance.

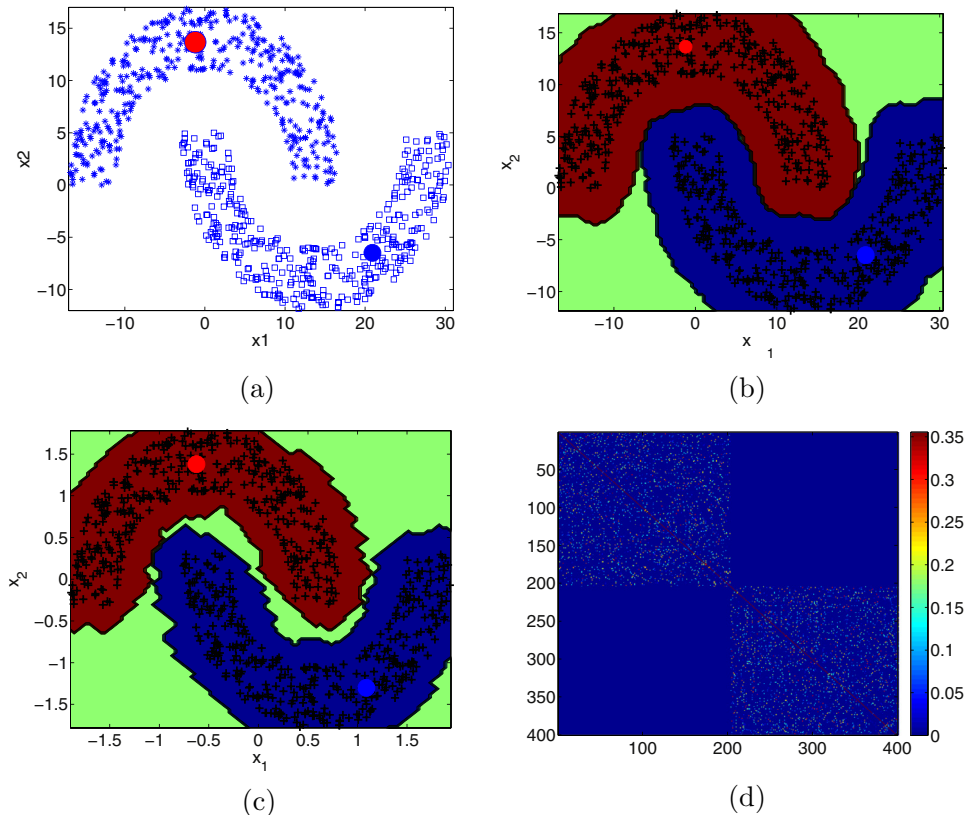
## 6. Numerical experiments

In this section, experimental results on a synthetic as well as several real-life datasets from the UCI machine learning repository [30] are given. We also show the applicability of the proposed indefinite method on a simple image segmentation task. Furthermore, the performance of the model for classification of partially labeled large-scale datasets using indefinite kernels will be studied in this section.

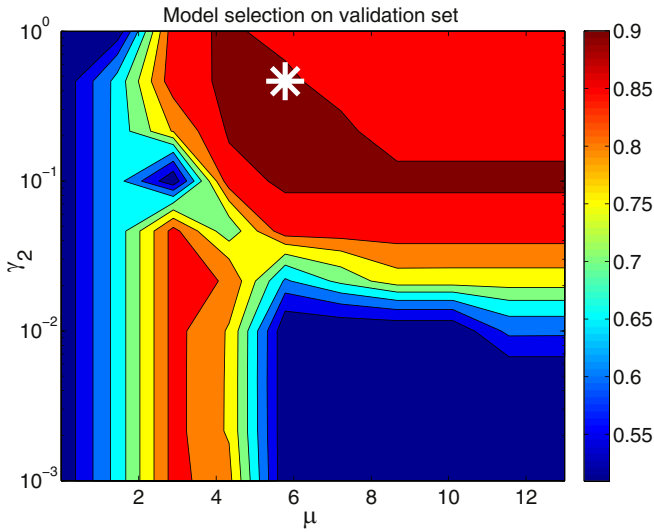
The performance of kernel learning relies on the choice of kernel. In this paper, we consider two indefinite kernels in KSC/MSS-KSC. One is the TL1 kernel (19) and the other is the tanh kernel with parameters  $c, d$ :

$$K(s, t) = \tanh(cs^T t + d). \quad (25)$$

Notice that when  $c > 0$ , the tanh kernel is conditionally positive definite; otherwise, it is indefinite. In the following experiments,  $c$  is selected from both positive and negative values, and hence the tanh kernel is regarded as an indefinite kernel in this paper. The performance of these indefinite kernels will be compared with the RBF kernel, which is the most popular PSD kernel and takes the



**Fig. 2.** Illustrating the performance of MSS-KSC model on synthetic single labeled example. (a) Original labeled and unlabeled points. (b) The predicted memberships obtained using MSS-KSC model with the RBF kernel. (c) The predicted memberships obtained using MSS-KSC model with an indefinite kernel. (d) The associated similarity matrix indicating the cluster structure in the data.



**Fig. 3.** Illustrating the sensitivity of the MSS-KSC model with respect to its parameters,  $\gamma_2$  and  $\mu$  in the case of the TL1 kernel for the Wine dataset.

following formulation:

$$K(s, t) = \exp(-\|s - t\|_2^2 / \sigma^2). \quad (26)$$

### 6.1. Semi-supervised classification

First, Two-moons dataset, a 2-dimensional synthetic problem, is considered to visualize the performance of indefinite kernels in a semi-supervised setting. The results obtained via the RBF kernel and the TL1 kernel are shown in Fig. 2, from which it can be seen that the two classes have been successfully classified by both the PSD and the non-PSD kernel. One may notice that the decision boundaries obtained by the TL1 kernel is not as smooth as the RBF

kernel. It is due to the piecewise linearity of the TL1 kernel and could be different if other non-PSD kernels are used.

Next, we conduct experiments on real-life datasets from UCI repository [30]. Here, 60% of the whole data (at random) is used as test set and the remaining 40% as training set. We randomly select part of the training data as the labeled and the remaining ones as the unlabeled training data. The ratio of the labeled training data points that is used in our experiments is denoted as follows:

$$\text{ratio}_{\text{label}} = \frac{\# \text{ labeled training data points}}{\# \text{ training data points}}.$$

The considered ratios for forming a labeled training set are one-fourth, one-third and half of the whole training dataset. To reduce the randomness of the experiment, we repeat this process 10 times. At each run, 10-fold cross validation is performed for model selection. The parameters to tune are the regularization constants  $\gamma_1$ ,  $\gamma_2$  and kernel parameters. In our experiments, we set  $\gamma_1 = 1$  and then find reasonable values for  $\gamma_2$ ,  $\mu$  in the range  $[10^{-3}, 10^0]$  and  $[0, d]$ , respectively. For the RBF kernel, and  $\sigma \in \{10^{-4}, 10^{-3}, \dots, 10^4\}$ . For tanh kernel, the candidate sets are  $c \in \{-0.5 - 0.2, -0.1, 0, 0.1, 0.2, 0.5\}$  and  $d \in \{2^{-10}, 2^{-7}, \dots, 2^3\}$ . The cross-validation performance on the Wine dataset for the TL1 kernel is shown in Fig. 3, from which and other experiments, we empirically observed that the TL1 kernel enjoys good stability on its kernel parameters. This makes its performance for a pre-given value, e.g.,  $\mu = 0.7d$ , satisfactory in many tested examples.

The average accuracy on the test dataset over 10 trials are reported in Table 1, where the details of the datasets are provided as well. From the results, one can observe that the performance of the unsupervised KSC model with an indefinite kernel is generally comparable to that with the RBF kernel. For most problems, the TL1 kernel with a pre-given  $\mu$  outputs good results. Moreover, there are indeed some problems, like Monk3 and Ionosphere, for which indefinite kernel learning can improve the performance significantly.

**Table 1**

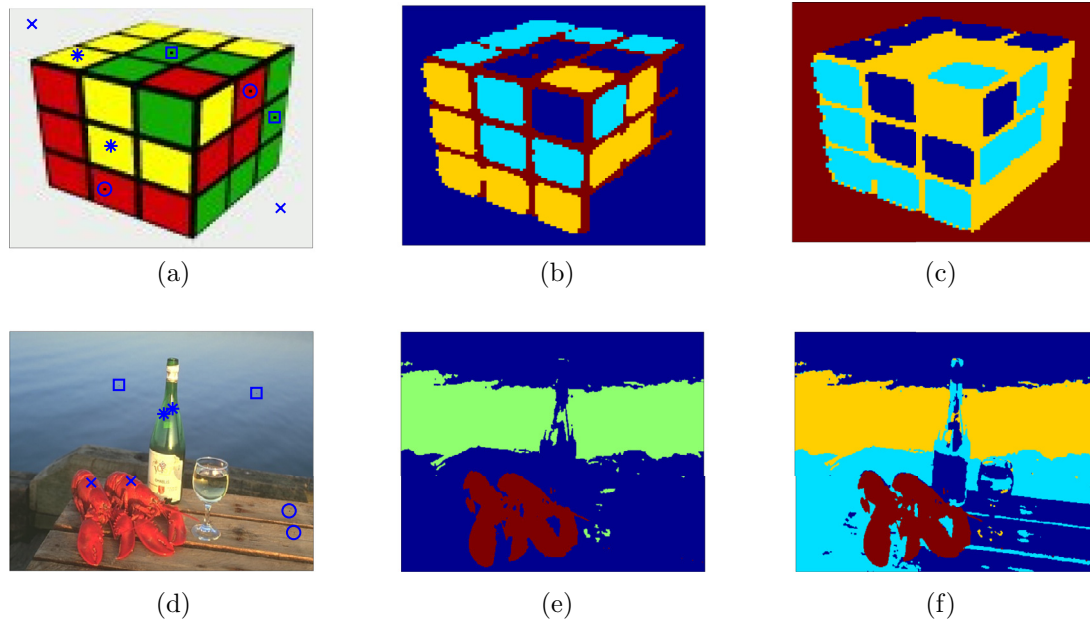
The average accuracy and the standard deviation of the LapSVMp [19] and MSS-KSC on the test set using PSD and indefinite kernels.

Dataset	$d$	$Q$	Ratio <sub>label</sub>	$D^{\text{train}}_{\text{Labeled}} / D^{\text{train}}_{\text{Unlabeled}} / D^{\text{test}}$	MSS-KSC method				
					RBF kernel $\sigma$ is tuned	TL1 kernel $\mu$ is tuned	TL1 kernel $\mu = 0.7d$	Tanh-kernel $c, d$ is tuned	LapSVMp
Iris	4	3	1/4	15/45/90	$0.85 \pm 0.09$	$0.88 \pm 0.07$	$0.86 \pm 0.09$	$0.65 \pm 0.11$	$0.70 \pm 0.12$
			1/3	20/40/90	$0.87 \pm 0.07$	$0.88 \pm 0.09$	$0.86 \pm 0.03$	$0.71 \pm 0.07$	$0.76 \pm 0.11$
			1/2	30/30/90	$0.92 \pm 0.03$	$0.90 \pm 0.08$	$0.88 \pm 0.09$	$0.77 \pm 0.10$	$0.83 \pm 0.10$
Wine	13	3	1/4	18/54/106	$0.89 \pm 0.07$	$0.90 \pm 0.08$	$0.89 \pm 0.03$	$0.59 \pm 0.12$	$0.73 \pm 0.11$
			1/3	24/48/106	$0.92 \pm 0.01$	$0.93 \pm 0.01$	$0.92 \pm 0.03$	$0.75 \pm 0.11$	$0.84 \pm 0.09$
			1/2	36/36/106	$0.94 \pm 0.01$	$0.95 \pm 0.02$	$0.93 \pm 0.03$	$0.84 \pm 0.12$	$0.90 \pm 0.10$
Zoo	16	7	1/4	11/30/60	$0.89 \pm 0.05$	$0.84 \pm 0.10$	$0.75 \pm 0.17$	$0.60 \pm 0.10$	$0.78 \pm 0.08$
			1/3	14/27/60	$0.89 \pm 0.04$	$0.90 \pm 0.04$	$0.80 \pm 0.10$	$0.66 \pm 0.09$	$0.82 \pm 0.11$
			1/2	21/20/60	$0.90 \pm 0.04$	$0.89 \pm 0.04$	$0.83 \pm 0.17$	$0.72 \pm 0.12$	$0.85 \pm 0.10$
Seeds	7	3	1/4	21/63/126	$0.87 \pm 0.05$	$0.88 \pm 0.03$	$0.85 \pm 0.09$	$0.62 \pm 0.10$	$0.80 \pm 0.10$
			1/3	28/56/126	$0.88 \pm 0.09$	$0.86 \pm 0.09$	$0.85 \pm 0.04$	$0.70 \pm 0.12$	$0.83 \pm 0.11$
			1/2	42/42/126	$0.90 \pm 0.01$	$0.88 \pm 0.02$	$0.88 \pm 0.02$	$0.79 \pm 0.11$	$0.87 \pm 0.09$
Monk1	6	2	1/4	56/167/333	$0.63 \pm 0.04$	$0.66 \pm 0.03$	$0.63 \pm 0.03$	$0.59 \pm 0.09$	$0.60 \pm 0.10$
			1/3	75/148/333	$0.67 \pm 0.03$	$0.69 \pm 0.03$	$0.64 \pm 0.03$	$0.60 \pm 0.03$	$0.65 \pm 0.11$
			1/2	112/111/333	$0.68 \pm 0.07$	$0.70 \pm 0.08$	$0.70 \pm 0.03$	$0.63 \pm 0.07$	$0.69 \pm 0.08$
Monk2	6	2	1/4	61/180/360	$0.63 \pm 0.08$	$0.61 \pm 0.06$	$0.54 \pm 0.03$	$0.57 \pm 0.02$	$0.58 \pm 0.11$
			1/3	81/160/360	$0.64 \pm 0.06$	$0.62 \pm 0.05$	$0.55 \pm 0.03$	$0.61 \pm 0.06$	$0.63 \pm 0.10$
			1/2	121/120/360	$0.71 \pm 0.04$	$0.65 \pm 0.06$	$0.58 \pm 0.02$	$0.63 \pm 0.03$	$0.66 \pm 0.11$
Monk3	6	2	1/4	56/166/332	$0.74 \pm 0.03$	$0.81 \pm 0.03$	$0.81 \pm 0.02$	$0.68 \pm 0.10$	$0.77 \pm 0.08$
			1/3	74/148/332	$0.79 \pm 0.02$	$0.85 \pm 0.03$	$0.83 \pm 0.04$	$0.74 \pm 0.02$	$0.80 \pm 0.09$
			1/2	111/111/332	$0.81 \pm 0.02$	$0.87 \pm 0.03$	$0.87 \pm 0.02$	$0.77 \pm 0.04$	$0.84 \pm 0.10$
Pima	8	2	1/4	77/231/460	$0.70 \pm 0.01$	$0.70 \pm 0.03$	$0.70 \pm 0.03$	$0.62 \pm 0.14$	$0.70 \pm 0.08$
			1/3	74/148/460	$0.71 \pm 0.02$	$0.72 \pm 0.03$	$0.71 \pm 0.01$	$0.69 \pm 0.02$	$0.71 \pm 0.10$
			1/2	154/154/460	$0.72 \pm 0.02$	$0.72 \pm 0.02$	$0.72 \pm 0.02$	$0.70 \pm 0.05$	$0.72 \pm 0.06$
Ionosphere	33	2	1/4	36/105/210	$0.77 \pm 0.05$	$0.81 \pm 0.08$	$0.75 \pm 0.07$	$0.69 \pm 0.04$	$0.77 \pm 0.09$
			1/3	47/94/210	$0.83 \pm 0.06$	$0.88 \pm 0.03$	$0.77 \pm 0.07$	$0.71 \pm 0.05$	$0.83 \pm 0.08$
			1/2	71/70/210	$0.86 \pm 0.07$	$0.88 \pm 0.03$	$0.79 \pm 0.05$	$0.73 \pm 0.03$	$0.86 \pm 0.09$

**Table 2**

Comparison of the KSC model with PSD and indefinite kernel, K-means and landmark-based spectral clustering algorithm using two internal clustering quality metrics, i.e. Silhouette and DB index, on some real datasets.

Dataset	$n$	$d$	$N_c$	Silhouette index			DB index		
				RBF	TL1	K-means	RBF	TL1	K-means
Wine	178	13	3	0.44	0.46	<u>0.50</u>	1.41	<u>1.06</u>	1.22
Thyroid	215	3	2	0.68	<u>0.81</u>	0.75	0.52	<u>0.43</u>	0.97
Breast	699	9	2	0.75	0.75	0.75	<u>0.77</u>	0.86	0.76
Glass	214	9	7	0.81	<u>0.84</u>	0.63	1.20	<u>1.09</u>	0.64
Iris	150	4	3	0.77	0.77	0.64	0.73	<u>0.59</u>	0.70



**Fig. 4.** Illustrating the performance of MSS-KSC model with an indefinite kernel (TL1) on image segmentation. (a,d) The labeled images. (b,e) The segmentations obtained by unsupervised KSC model with the TL1 kernel. (c,f) The segmentation obtained by semi-supervised MSS-KSC model with the TL1 kernel.

## 6.2. Clustering

The experimental results on several real world clustering datasets<sup>1</sup> using KSC model with the RBF and the TL1 kernel are reported in Table 2. The cluster memberships of these datasets are not known beforehand, therefore the clustering results can be evaluated by internal clustering quality metrics such as the widely used silhouette index (Sil-index) and the Davies Bouldin index (DB-index) [26]. Larger values of Sil-index imply better clustering quality. While, the lower the value of DB-index means that the clustering quality is better. In Table 2, the best indices are underlined where one can observe the good performance of the TL1 kernel. Notice that simply from these experiments, we cannot conclude indefinite kernel is better or worse than the definite ones. But the results indicate that for some problems, it is worth to consider the proposed indefinite unsupervised learning methods, which may further improve the performance from the traditional PSD kernel learning methods.

## 6.3. Image segmentation

Here we show the application of the proposed indefinite Kernel on unsupervised and semi-supervised image segmentation. Following the lines of Mehrkanoon et al. [22], for each image, a local color histogram with a  $5 \times 5$  local window around each pixel

is computed using minimum variance color quantization of eight levels. A subset of 500 unlabeled pixels together with some labeled pixels are used for training and the whole image for test. The original and labeled images together with segmentation results are shown in Fig. 4. One can qualitatively observe that thanks to the provided labeled pixels, the semi-supervised model performs better than completely unsupervised model on the test images.

## 6.4. Large scale datasets

Here we show the possibility of applying the TL1 kernel in the context of semi-supervised learning on large-scale datasets. The size of the real-life data, on which the experiments were conducted, ranges from medium to large and covering both binary and

**Table 3**  
Dataset statistics.

Dataset	# points	# attributes	# classes
Adult	48,842	14	2
IJCNN	141,691	22	3
Cod-RNA	331,152	8	2
Coverttype	581,012	54	3
SUSY	5,000,000	18	2
Sensorless	58,509	48	11
letter	20,000	16	26
Satimage	6435	36	6
texture	5500	40	11
USPS	9298	256	10

<sup>1</sup> <http://cs.joensuu.fi/sipu/datasets/> (accessed: 2015-12-29).

**Table 4**

Comparing the average test accuracy, standard deviation and computation time of the FS-MSS-KSC model [21] with the RBF kernel and the TL1 kernel on real-life datasets over 10 simulation runs.

Dataset	$p$	Ratio <sub>label</sub>	$\mathcal{D}_{tr}^L$	$\mathcal{D}_{tr}^U$	$\mathcal{D}_{test}$	Test accuracy		Computation time (in seconds)	
						RBF	TL1	RBF	TL1
USPS	2	1/3	1000	2000	1859	<u>0.86 ± 0.002</u>	<u>0.86 ± 0.002</u>	<u>0.02</u>	0.16
		1/3	2000	4000	1859	<u>0.88 ± 0.003</u>	<u>0.89 ± 0.002</u>	<u>0.02</u>	0.81
Texture	3	1/4	500	1500	1100	<u>0.85 ± 0.002</u>	<u>0.87 ± 0.002</u>	<u>0.01</u>	0.02
		1/4	1000	3000	1100	<u>0.89 ± 0.004</u>	<u>0.91 ± 0.001</u>	<u>0.02</u>	0.05
Satimage	3	1/4	500	1500	1287	<u>0.83 ± 0.003</u>	<u>0.85 ± 0.003</u>	<u>0.01</u>	0.02
		1/4	1000	3000	1287	<u>0.85 ± 0.001</u>	<u>0.86 ± 0.002</u>	<u>0.02</u>	0.05
Adult	3	1/4	4000	12,000	9768	<u>0.844 ± 0.003</u>	<u>0.847 ± 0.006</u>	<u>0.08</u>	0.20
		1/4	8000	24,000	9768	<u>0.846 ± 0.003</u>	<u>0.852 ± 0.005</u>	<u>0.22</u>	0.34
Letter	3	1/4	2000	6000	4000	<u>0.65 ± 0.002</u>	<u>0.68 ± 0.003</u>	<u>0.05</u>	0.12
		1/4	4000	12,000	4000	<u>0.69 ± 0.004</u>	<u>0.71 ± 0.002</u>	<u>0.12</u>	0.25
Sensorless	3	1/4	4000	12,000	11,701	<u>0.92 ± 0.002</u>	<u>0.93 ± 0.002</u>	<u>0.24</u>	1.46
		1/4	8000	24,000	11,701	<u>0.94 ± 0.001</u>	<u>0.96 ± 0.001</u>	<u>0.54</u>	3.21
IJCNN	5	1/6	4000	20,000	28,338	<u>0.935 ± 0.004</u>	<u>0.933 ± 0.001</u>	<u>0.51</u>	1.70
		1/6	16,000	80,000	28,338	<u>0.956 ± 0.002</u>	<u>0.953 ± 0.001</u>	<u>2.53</u>	6.01
Cod-RNA	5	1/6	8000	40,000	66,230	<u>0.959 ± 0.001</u>	<u>0.952 ± 0.001</u>	<u>0.92</u>	1.57
		1/6	32,000	160,000	66,230	<u>0.962 ± 0.0005</u>	<u>0.958 ± 0.001</u>	<u>6.63</u>	8.27
Coverttype	5	1/6	8000	40,000	116,202	<u>0.732 ± 0.001</u>	<u>0.740 ± 0.003</u>	<u>1.80</u>	8.01
		1/6	64,000	320,000	116,202	<u>0.781 ± 0.001</u>	<u>0.772 ± 0.002</u>	<u>12.20</u>	25.7
SUSY	2	1/3	500,000	1,000,000	1,000,000	<u>0.771 ± 0.001</u>	<u>0.771 ± 0.001</u>	<u>4.91</u>	15.98
		1/3	1,000,000	2,000,000	1,000,000	<u>0.783 ± 0.001</u>	<u>0.787 ± 0.001</u>	<u>10.01</u>	34.70

multi-class classification. The classification of these datasets is performed using different number of training labeled and unlabeled data instances. In our experiments, for all the datasets, 20% of the whole data (at random) is used for test, and the training set is constructed from the remaining 80% of the data. In order to have a realistic setting, the number of unlabeled training points are considered to be  $p$  times more than that of labeled training points, where, in our experiments, depending on the size of the dataset,  $p$  ranges from 2 to 5. Descriptions of the considered datasets can be found in Table 3.

The average results of the proposed MSS-KSC model with the TL1 kernel together with that of Fixed-size MSS-KSC [21] are tabulated in Table 4. From Table 4, one can observe that the proposed MSS-KSC algorithm with an indefinite kernel has been successfully applied on large scab data and its accuracy is comparable to that of the RBF kernel. This is an interesting point as in many applications one need to address the scalability of the models when using indefinite kernel. It should be mentioned that as expected, the computational time of MSS-KSC with the RBF kernel is faster than that of MSS-KSC with the TL1 kernel. This can be explained by the fact that in the RBF kernel, one feature map is constructed where as in the TL1 kernel one needs to calculate two feature maps.

## 7. Conclusions

Motivated by success of indefinite kernels in supervised learning, we in this paper proposed to use indefinite kernels in the semi-supervised learning framework. Specifically, we studied the indefinite KSC and MSS-KSC models. For both models the optimization problems remain easy to solve if indefinite kernels are used. The interpretations of the feature map in the case of indefinite kernels are provided. Based on these interpretations, Nyström approximation can be used for the scalability of indefinite KSC and MSS-KSC. The proposed indefinite learning methods are evaluated on real datasets in comparison with the existing methods with the RBF kernel. One can observe that for some datasets, the indefinite kernel shows its superiority, which implies that there are some semi-supervised tasks requiring indefinite learning methods. For example, when some (dis)similarity induces to indefinite kernels, it is better to directly use those indefinite kernel rather than to find approximate PSD ones. Furthermore, if an indefinite kernel is

suitably selected or designed, the indefinite learning performance could be very promising.

## Acknowledgments

The authors are grateful to the anonymous reviewer for insightful comments.

The research leading to these results received funding from the [European Research Council](#) under the European Union's Seventh Framework Programme (FP7/2007–2013)/ERC AdG A-DATADRIVE-B (290923). This letter reflects only our views: The EU is not responsible for any use that may be made of the information in it. The research leading to these results received funds from the following sources: Research Council KUL: GOA/10/09 MaNet, CoE PFV/10/002 (OPTEC), BIL12/11T; PhD/Postdoc grants; Flemish Government: FWO: PhD /Postdoc grants, projects: G.0377.12 (Structured systems), G.088114N (Tensor based data similarity); IWT: PhD/Postdoc grants, projects: SBO POM (100031); iMinds Medical Information Technologies SBO 2014; Belgian Federal Science Policy Office: IUAP P7/19 (DYSCO, Dynamical systems, control and optimization, 2012–2017). Siamak Mehrkanoon was supported by a Postdoctoral Fellowship of the Research Foundation-Flanders (FWO). Xiaolin Huang is supported by [National Natural Science Foundation of China](#) (no. 61603248). Johan Suykens is a full professor at KU Leuven, Belgium.

## References

- [1] D. Wang, X. Zhang, M. Fan, X. Ye, Hierarchical mixing linear support vector machines for nonlinear classification, *Pattern Recognit.* 59 (2016) 255–267.
- [2] Y. Li, X. Tian, M. Song, D. Tao, Multi-task proximal support vector machine, *Pattern Recognit.* 48 (2015) 3249–3257.
- [3] J. Richarz, S. Vajda, R. Grzeszick, G.A. Fink, Semi-supervised learning for character recognition in historical archive documents, *Pattern Recognit.* 47 (2014) 1011–1020.
- [4] V. Vapnik, *Statistical Learning Theory*, Wiley, 1998.
- [5] Q. Wu, Regularization networks with indefinite kernels, *J. Approx. Theory* 166 (2013) 1–18.
- [6] E. Pekalska, B. Haasdonk, Kernel discriminant analysis for positive definite and indefinite kernels, *IEEE Trans. Pattern Anal. Mach. Intell.* 31 (2009) 1017–1032.
- [7] F.M. Schlei, P. Tino, Indefinite core vector machine, *Pattern Recognit.* 71 (2017) 187–195.
- [8] Y. Chen, M.R. Gupta, B. Recht, Learning kernels from indefinite similarities, in: *Proceedings of the 26th International Conference on Machine Learning*, 2009, pp. 145–152.



- [9] C.S. Ong, X. Mary, S. Canu, A.J. Smola, Learning with non-positive kernels, in: Proceedings of the 21st International Conference on Machine Learning, 2004, pp. 639–646.
- [10] G. Loosli, S. Canu, C.S. Ong, Learning SVM in Kreĭn spaces, *IEEE Trans. Pattern Anal. Mach. Intell.* 38 (2016) 1204–1216.
- [11] E. Pekalska, P. Paclik, R.P.W. Duin, A generalized kernel approach to dissimilarity-based classification, *J. Mach. Learn. Res.* 2 (2002) 175–211.
- [12] R. Luss, A. d'Aspremont, Support vector machine classification with indefinite kernels, in: *Advances in Neural Information Processing Systems*, 2008, pp. 953–960.
- [13] J. Chen, J. Ye, Training SVM with indefinite kernels, in: *Proceedings of the 25th International Conference on Machine Learning*, 2008, pp. 136–143.
- [14] Y. Ying, C. Campbell, M. Girolami, Analysis of SVM with indefinite kernels, *Adv. Neural Inf. Process. Syst.* 22 (2009) 2205–2213.
- [15] H.-T. Lin, C.-J. Lin, A study on sigmoid kernels for SVM and the training of non-PSD kernels by SMO-type methods, 2003. Internal report. <https://www.csie.ntu.edu.tw/~cjlin/papers/tanh.pdf>.
- [16] B. Haasdonk, Feature space interpretation of SVMs with indefinite kernels, *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (2005) 482–492.
- [17] J.A.K. Suykens, T.V. Gestel, J. De Brabanter, B. De Moor, J. Vandewalle, *Least Squares Support Vector Machines*, World Scientific Pub. Co, Singapore, 2002.
- [18] X. Huang, A. Maier, J. Hornegger, J.A.K. Suykens, Indefinite kernels in least squares support vector machine and kernel principal component analysis, *Appl. Comput. Harmon. Anal.* 43 (2017) 162–172.
- [19] M. Belkin, P. Niyogi, V. Sindhwani, Manifold regularization: a geometric framework for learning from labeled and unlabeled examples, *J. Mach. Learn. Res.* 7 (2006) 2399–2434.
- [20] S. Mehrkanoon, C. Alzate, R. Mall, R. Langone, J.A.K. Suykens, Multiclass semisupervised learning based upon kernel spectral clustering, *IEEE Trans. Neural Netw. Learn. Syst.* 26 (2015) 720–733.
- [21] S. Mehrkanoon, J.A.K. Suykens, Large scale semi-supervised learning using KSC based model, in: *proceedings of the 2014 International Joint Conference on Neural Networks*, 2014, pp. 4152–4159.
- [22] S. Mehrkanoon, O.M. Agudelo, J.A.K. Suykens, Incremental multi-class semi-supervised clustering regularized by kAlman filtering, *Neural Netw.* 71 (2015) 88–104.
- [23] S. Mehrkanoon, J.A.K. Suykens, Multi-label semi-supervised learning using regularized kernel spectral clustering, in: *proceedings of the 2016 International Joint Conference on Neural Networks*, 2016, pp. 4009–4016.
- [24] C. Alzate, J.A.K. Suykens, Multiway spectral clustering with out-of-sample extensions through weighted kernel PCA, *IEEE Trans. Pattern Anal. Mach. Intell.* 32 (2010) 335–347.
- [25] X. Huang, J.A.K. Suykens, S. Wang, A. Maier, J. Hornegger, Classification with truncated  $\ell_1$  distance kernel, *IEEE Trans. Neural Netw. Learn. Syst.* doi:10.1109/TNNLS.2017.2668610.
- [26] J.C. Bezdek, N.R. Pal, Some new indexes of cluster validity, *IEEE Trans. Syst. Man Cybern. Part B Cybern.* 28 (1998) 301–315.
- [27] G.S. Mann, A. McCallum, Simple, robust, scalable semi-supervised learning via expectation regularization, *proceedings of the 24th International Conference on Machine Learning* (2007) 593–600.
- [28] W. Liu, J. He, S.-F. Chang, Large graph construction for scalable semi-supervised learning, *proceedings of the 27th International Conference on Machine Learning* (2010) 679–686.
- [29] C. Williams, M. Seeger, Using the Nyström method to speed up kernel machines, in: *Advances in Neural Information Processing Systems*, 2001, pp. 682–688.
- [30] A. Asuncion, D.J. Newman, UCI machine learning repository, 2007. <http://archive.ics.uci.edu/ml/index.php>.



**Siamak Mehrkanoon** received the B.Sc. degree in pure mathematics and the M.Sc. degree in applied mathematics from the Iran University of Science and Technology, Tehran, Iran, in 2005 and 2007, respectively. He is holder of Ph.D. degrees in Numerical Analysis and Machine Learning from Universiti Putra Malaysia, Seri Kembangan, Malaysia, and KU Leuven, Belgium, in 2011 and 2015, respectively. He was a Visiting Researcher with the Department of Automation, Tsinghua University, Beijing, China, in 2014, a Postdoctoral Research Fellow with the University of Waterloo, Waterloo, ON, Canada, from 2015 to 2016, and a visiting postdoctoral researcher with the Cognitive Systems Laboratory, University of Tübingen, Tübingen, Germany, in 2016. He is currently an FWO Postdoctoral Research Fellow with the STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics, KU Leuven.

His current research interests include deep learning, neural networks, kernel-based models, unsupervised and semi-supervised learning, pattern recognition, numerical algorithms, and optimization. Dr. Mehrkanoon received several fellowships for supporting his scientific studies including Postdoctoral Mandate (PDM) Fellowship from KU Leuven and Postdoctoral Fellowship of the Research Foundation-Flanders (FWO).



**Xiaolin Huang** received the B.S. degree in control science and engineering, and the B.S. degree in applied mathematics from Xi'an Jiaotong University, Xi'an, China in 2006. In 2012, he received the Ph.D. degree in control science and engineering from Tsinghua University, Beijing, China. From 2012 to 2015, he worked as a postdoctoral researcher in ESAT-STADIUS, KU Leuven, Leuven, Belgium. After that he was selected as an Alexander von Humboldt Fellow and working in Pattern Recognition Lab, the Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany, where he was appointed as a group head. From 2016, he has been an Associate Professor at Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, Shanghai, China. In 2017, he has been awarded as "1000-Talent" (Young Program). His current research areas include machine learning, optimization, and their applications on medical image processing.



**Johan A.K. Suykens** was born in Willebroek Belgium, on May 18, 1966. He received the M.S. degree in Electro-Mechanical Engineering and the Ph.D. Degree in Applied Sciences from the Katholieke Universiteit Leuven, in 1989 and 1995, respectively.

In 1996 he has been a Visiting Postdoctoral Researcher at the University of California, Berkeley. He has been a Postdoctoral Researcher with the Fund for Scientific Research FWO Flanders and is currently a Professor (Hoogleraar) with KU Leuven. He is author of the books *Artificial Neural Networks for Modelling and Control of Non-linear Systems* (Kluwer Academic Publishers) and *Least Squares Support Vector Machines* (World Scientific), co-author of the book *Cellular Neural Networks, Multi-Scroll Chaos and Synchronization* (World Scientific) and editor of the books *Nonlinear Modeling: Advanced Black-Box Techniques* (Kluwer Academic Publishers) and *Advances in Learning Theory: Methods, Models and Applications* (IOS Press).

Prof. Suykens received an IEEE Signal Processing Society 1999 Best Paper (Senior) Award and several best paper awards at international conferences. He was a recipient of the International Neural Networks Society 2000 Young Investigator Award for significant contributions in the field of neural networks. He has been awarded an ERC Advanced Grant 2011 and has been elevated IEEE Fellow 2015 for developing least squares support vector machine. In 1998, he organized an International Workshop on Nonlinear Modeling with Timeseries Prediction Competition. He served as an Associate Editor of the IEEE Transactions on Circuits and Systems from 1997 to 1999 and 2004 to 2007, and the IEEE Transactions on Neural Networks from 1998 to 2009. He served as a Director and an Organizer of the NATO Advanced Study Institute on Learning Theory and Practice, Leuven, in 2002, a Program Co-Chair of the International Joint Conference on Neural Networks in 2004 and the International Symposium on Nonlinear Theory and its Applications in 2005, an Organizer of the International Symposium on Synchronization in Complex Networks in 2007, a Co-Organizer of the Conference on Neural Information Processing Systems Workshop on Tensors, Kernels and Machine Learning in 2010, and the Chair of the International Workshop on Advances in Regularization, Optimization, Kernel methods and Support vector machines in 2013.