

Non-Parametric Kernel Learning with robust pairwise constraints

Changyou Chen · Junping Zhang ·
Xuefang He · Zhi-Hua Zhou

Received: 20 November 2010 / Accepted: 1 September 2011 / Published online: 17 September 2011
© Springer-Verlag 2011

Abstract For existing kernel learning based semi-supervised clustering algorithms, it is generally difficult to scale well with large scale datasets and robust pairwise constraints. In this paper, we propose a new *Non-Parametric Kernel Learning* (NPKL) framework to deal with these problems. We generalize the graph embedding framework into kernel learning, by reforming it as a semi-definitive programming (SDP) problem, smoothing and avoiding over-smoothing the functional Hilbert space with Laplacian regularization. We propose two algorithms to solve this problem. One is a straightforward algorithm using SDP to solve the original kernel learning problem, dented as *TRAnsductive Graph Embedding Kernel* (TRAGEK) learning; the other is to relax the SDP problem and solve it with a constrained gradient descent algorithm. To accelerate the

learning speed, we further divide the data into groups and used the sub-kernels of these groups to approximate the whole kernel matrix. This algorithm is denoted as *Efficient Non-PArametric Kernel Learning* (ENPAKL). The advantages of the proposed NPKL framework are (1) supervised information in the form of pairwise constraints can be easily incorporated; (2) it is robust to the number of pairwise constraints, i.e., the number of constraints does not affect the running time too much; (3) ENPAKL is efficient to some extent compared to some related kernel learning algorithms since it is a constraint gradient descent based algorithm. Experiments for clustering based on the learned kernels show that the proposed framework scales well with the size of datasets and the number of pairwise constraints. Further experiments for image segmentation indicate the potential advantages of the proposed algorithms over the traditional k -means and N -cut clustering algorithms for image segmentation in term of segmentation accuracy.

This work was done when C. Chen was at Fudan University, Shanghai, China.

C. Chen (✉)
Research School of Information Sciences and Engineering,
The Australian National University, Canberra, Australia
e-mail: cchangyou@gmail.com

J. Zhang
Shanghai Key Laboratory of Intelligent Information Processing
and School of Computer Science, Fudan University,
Shanghai, China
e-mail: jpzhang@fudan.edu.cn

X. He
School of Software and Information Engineering,
Beijing University of Aeronautics & Astronautics, Beihai, China
e-mail: ahexuefang@gmail.com

Z.-H. Zhou
National Key Laboratory for Novel Software Technology,
Nanjing University, Nanjing, China
e-mail: zhouzh@nju.edu.cn

Keywords Kernel learning · Semi-definitive programming · Graph embedding · Pairwise constraint · Semi-supervised learning

1 Introduction

Semi-supervised clustering based on kernel learning is a popular research topic in machine learning since one can incorporate the information of a limited number of labeled data or a set of pairwise constraints into the kernel learning framework [1]. The reason is that for clustering, the pairwise constraints provide useful information about which data pairs are in the same category and which ones are not. To learn such kinds of kernel matrices, Kulis et al. [2]

proposed to construct a graph based kernel matrix which unifies the vector-based and graph-based semi-supervised clustering. A further refinement on learning kernel matrices for clustering was investigated by Li et al. [3]. In their approach, data are implicitly projected onto a feature space which is a unit hyperball, subjected to a collection of pairwise constraints. However, the above clustering algorithms via kernel matrices either can not scale well with the increasing number of pairwise constraints and the amount of data, or lacks theoretical guarantee for the positive semi-definite property of the kernel matrices. In another aspect, Yueng et al. [4] proposed an efficient kernel learning algorithm through low rank matrix approximation. However, in their algorithm, the form of kernel matrix is assumed to be linear combination of several base kernel matrices. Note that this might reduce the dimension of the hypothesis kernel space, we call such kinds of algorithms parametric kernel learning. In addition, Cortes et al. [5] proposed a kernel learning algorithm by taking the non-linear combinations of kernels, which is a generalization of the linear combination case but still lies in the framework of parametric kernel learning. Addressing these two limitations is the major purpose of this paper.

On the other hand, we note that many algorithms based on the graph embedding framework often achieve an enhanced discriminant ability by utilizing the marginal information, e.g., making the dissimilarity data points near the margin as far as possible and meanwhile compacting the points in the same class [6, 7]. It is therefore worthwhile to generalize the graph embedding framework into kernel learning.¹

Based on the aforementioned goals, in this paper we propose a new scalable kernel learning framework NPCL (*Non-Parametric Kernel Learning with robust pairwise constraints*), and apply it for semi-supervised clustering. First, we generalize the graph embedding framework on a feature space which is assumed to be a possibly infinite subspace of the l^2 Hilbert space with unit norm, which is similar to [3]. Then the unknown feature projection function ϕ is implicitly learned by transforming the criterion of the graph embedding (i.e., maximizing the sum of distances of between-class data pairs while minimizing that of within-class data pairs) into an SDP problem. To get smoother solution of the predictive function, smoothing technique using some kind of Laplacian regularizer is introduced. By this, ideally, data from the same class would be projected into the same location in the feature

space. Meanwhile, the distances between the locations of different classes should be as large as possible, as illustrated in Fig. 1. We propose two algorithms to solve this problem. One is to optimize the objective function directly by solving an SDP problem, which we call *TRAnsductive Graph Embedding Kernel* (TRAGEK) learning, since the SDP problem is derived from a transductive graph embedding formulation. In TRAGEK, it is not necessary to explicitly specify which pair of data points should lie close, and the running time is much less sensitive to the number of pairwise constraints than Li et al.'s work [3]. However, the SDP problem in TRAGEK limits the application of the proposed algorithm to large scale datasets. To alleviate this problem, we propose to solve the SDP problem via a constrained gradient descent algorithm that iteratively projects the unconstrained solutions to the cone formed by the constraints. Furthermore, we divide the whole dataset into groups of sub-data sets, and the corresponding sub-kernels are learned for these sub-data separately. Finally, the global kernel matrix is obtained through the combination of these sub-kernels. In this way, not only is the positive semi-definite property of the kernel matrix well preserved, but also the computational complexity scales at most linearly with the size of the dataset, which is very efficient. We call this algorithm *Efficient Non-Parametric Kernel Learning* (ENPAKL).

The remaining of this paper is organized as follows. Section 2 reviews some related work for clustering using kernel learning. Section 3 formulates our problem and presents the TRAGEK algorithm. Section 4 elaborates the efficient algorithm ENPAKL for our problem. And experiment results are shown in Sect. 5. Finally, Sect. 6 concludes the paper.

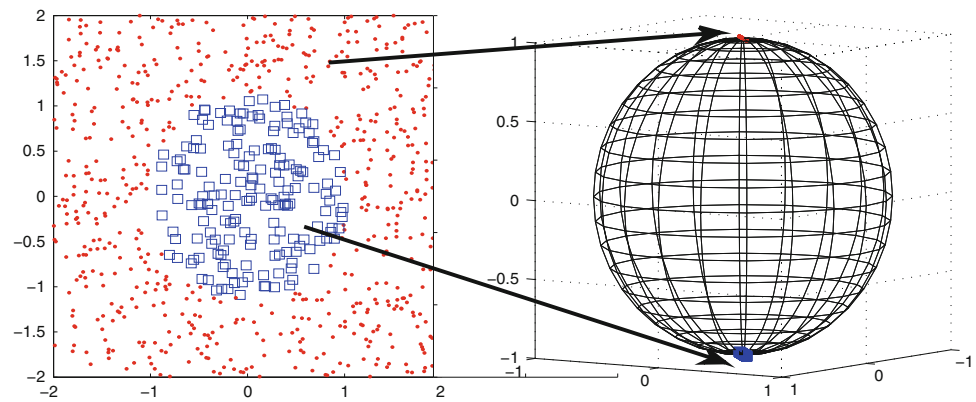
2 Clustering and kernel learning

Learning with kernels [8] is a popular research topic in machine learning. In this section, however, rather than reviewing the theoretical aspects of kernel learning algorithms such as [9, 10], and the online kernel learning such as [11], we put our emphasize on the applications of kernel learning algorithms. More specifically, we discuss the work of kernel learning for semi-supervised clustering. To a certain extent, kernel learning for clustering can be viewed as metric learning, because the kernel matrix can be regarded as some specific distance between data points equipped with a specific metric.

For traditional clustering algorithms, the most frequently used ones include k-means and fuzzy k-means [12]. Recently, Yang et al. [13] proposed an improved fuzzy k-means algorithm to assign a value of 1 to data pairs with a defined cluster score. Moreover, Trappey et al. [14]

¹ Note that although we want to learn a kernel matrix from the aspect of graph embedding, it has little relationship with some algorithms using graph embedding framework such as marginal factor analysis (MFA) [6]. The reason is that such kinds of algorithms aim at supervised learning for classification, thus there is no need to compare the proposed algorithm with them.

Fig. 1 With a feature projection function ϕ , data from two classes are projected from two-dimensional data space into three-dimensional feature space in which each point corresponds to one class



presented a fuzzy ontology schema for hierarchically clustering of documents, which can solve the inconsistent and ineffective problem encountered by the traditional keyword-based methods. For the traditional k-means clustering algorithm, Xiong et al. [15] proposed to use the coefficient of variation (CV) as some criterion to analysis the performance of k-means algorithm under skewed data distribution.

To utilize label information to enhance the cluster performance, semi-supervised clustering approaches were proposed, which can be roughly categorized into constraint- and metric-based ones. The former utilizes either labeled data or pairwise constraints to improve the performance of clustering [16, 17], while the latter learns a more rational metric to fit the constraints by utilizing the provided label information [18, 19], thus the semi-supervised kernel learning is closely related to the clustering algorithms.

More specifically, Wagstaff et al. [20] modified the k-means algorithm by considering pairwise similarities and dissimilarities, known as the constrained k-means. To boost the performance of constrained k-means, Hong et al. [21] refined the assignment order of this algorithm by ranking all instances in the dataset according to their clustering uncertainty. Based on the hidden Markov random field (HMRF), Basu et al. [16] used the provided pairwise constraints in the objective function for semi-supervised clustering, while Lu and Leen [22] proposed a clustering algorithm with the must-link and cannot-link constraints using the Gaussian mixture model (GMM). Xing et al. [19] learned a distance matrix for clustering by explicitly minimizing the distances between similar samples and maximizing those between dissimilar ones. Bar-Hillel et al. [23] proposed a simpler but more efficient algorithm using the relevant component analysis (RCA). Furthermore, they proposed another metric learning algorithm [24] that learns a non-parametric distance function, but without the guarantee that the function is actually a metric.

Although kernel methods have been studied for decades, not much work has focused on learning a non-parametric kernel using only the training data. In contrary, most of the work focuses on learning the kernel from some predefined base kernels [25]. Recently, kernel learning for clustering has attracted more and more attentions because one can easily incorporate some useful information into the kernel learning framework [2, 26]. Earlier kernel learning algorithms mainly focus on linear or non-linear combination of some base kernels [27]. For example, Yeung et al. [4] proposed a scalable kernel learning algorithm in which some low rank kernel matrices obtained by the eigenvectors of the initial kernel matrix are used as base kernels, and a collection of optimal weights of the base kernels need to be learned. Cortes [5] proposed to learn the kernel matrix by defining some non-linear terms of the basis kernels and use a projection-based gradient descent algorithm to learn the weights. One disadvantage of this algorithm is that only the must-link constraint information is incorporated. To employ both must-link and cannot-link constraint information, Hoi et al. [28] proposed a kernel learning algorithm by formulating it into a semi-definite programming (SDP) problem. To the best of our knowledge, this is the first non-parametric kernel learning algorithm that does not need to explicitly take the base kernels into consideration. Latter, an efficient algorithm for solving the SDP problem mentioned above is proposed by Zhuang et al. [29] by introducing an extra low rank constraint on the objective function. Furthermore, assuming that the feature space is a unit hyperball in which must-link data pairs are constrained to be one point and cannot-link pairs should be orthogonal with each other, Li et al. [3] proposed another SDP based algorithm for kernel learning. Note that one problem for the SDP related algorithms introduced above is the computational complexity, how to avoid this complexity is one of the major concerns of this paper.

3 TRAGEK: transductive non-parametric kernel learning

Our non-parametric kernel learning problem assumes that the feature space of the data is a subspace of the l^2 Hilbert space with unit norm. More specifically, given a data point x , there exists a projection ϕ from the data space to the feature space endowed with a unit norm, i.e., $\|\phi(x)\| = 1$. In this way, data are mapped to the surface of a unit hyperball in which data of different classes can be separated easily. To find such a mapping function ϕ , we tried to generalize the graph embedding [6] into the kernel space and transformed the problem into the semi-definite programming with pairwise constraints. We further introduced some regularization terms such as the Laplacian regularizer to smooth the prediction function. Finally, the kernel k-means clustering algorithm is performed in the learned kernel matrix for clustering.

3.1 Transductive kernel learning with graph embedding

One goal of the graph embedding is to learn a low-dimensional discriminant subspace in which the sum of within-class distances is minimized meanwhile that of between-class distances is maximized [6]. This idea is formulated in Eq. 1:

$$F_g = \min_Y \frac{\text{tr}\{YLY^T\}}{\text{tr}\{YLY^T\}}, \quad (1)$$

where $Y = (y_1, y_2, \dots, y_M)$, $y_i \in \mathbb{R}^m$ is the data representation in the feature space, M is the number of samples. $L = D - W$ and $L^p = D^p - W^p$ are two graph Laplacian matrices with $D(i, i) = \sum_j W(i, j)$ and $D^p(i, i) = \sum_j W^p(i, j)$ representing two diagonal matrices. Two similarity matrices W and W^p represent the within-class relationship and between-class relationship of data points, respectively, which are defined according to the pairwise constraints as:

$$\begin{aligned} W(i, j) &= \begin{cases} 1, & i \text{ and } j \text{ belong to the same class} \\ 0, & \text{otherwise.} \end{cases} \\ W^p(i, j) &= \begin{cases} 1, & i \text{ and } j \text{ belong to different classes} \\ 0, & \text{otherwise.} \end{cases} \end{aligned} \quad (2)$$

For better generalization, the two quadratic forms in Eq. 1 are lifted into linear forms in the kernel space as:

$$F_1 = \text{tr}\{YLY^T\} = \text{tr}\{LY^TY\} = \langle l, K_v \rangle, \quad (3)$$

where $l = \text{vec}(L)$ is the vectorization of the matrix L , and K_v ² is the vectorization of the kernel matrix K defined by $K(i, j) = \langle y_i, y_j \rangle$.

² We use such kind of convention without declaration below.

In contrast to the graph embedding, the goal of TRAGEK is to cluster data in the feature space defined on the unit norm subspace of the l^2 Hilbert space. It is easy to get the inner product of two vector in the feature space using the re-producing property of the re-producing Hilbert space:

$$K(\phi(x), \phi(y)) = k_{\phi(y)}(\phi(x)) = \langle \phi(x), \phi(y) \rangle, \quad (4)$$

where $k_{\phi(y)}(\cdot) \in \mathcal{H}$, x and y are two data points in the data space, ϕ is the mapping function to be learned. Furthermore, in our transductive inference, the matrices W and W^p are defined using the pairwise constraints as in Eq. 2. Based on this, the objective function of TRAGEK is to kernelize the criterion of graph embedding [6] as:

$$\begin{aligned} F_t &= \min_K \frac{\langle l, K_v \rangle}{\langle l^p, K_v \rangle} \\ \text{s.t. } & l^T K_v > 0, \quad l^{pT} K_v > 0, \\ & K(i, i) = 1, \quad \text{for all } i, \\ & K(i, j) \leq 1, \quad \text{for all } i, j. \end{aligned} \quad (5)$$

The optimization function means that the learned kernel K should cluster within-class data as close as possible and between-class data as far as possible in the feature space. The first constraint is required according to the positive semi-definite property of the Laplacian matrix L , and the second and third constraints stem from the assumption that data in the feature space should lie on the surface of a hyperball with unit radius.

3.2 Conic optimization programming relaxation

To transform this problem into a convex optimization problem which has a unique optimal solution and polynomial time computational complexity, we first prove Theorem 1 as follows.

Theorem 1 *The nonlinear optimization problem defined in Eq. 5 can be relaxed to a second order cone programming as:*

$$\begin{aligned} \min & : \{t\} \\ \text{s.t. } & t + \langle l^p, K_v \rangle \geq \left\| \begin{matrix} \langle l, K_v \rangle \\ t - \langle l^p, K_v \rangle \end{matrix} \right\|, \\ & l^T K_v > 0, \quad l^{pT} K_v > 0, \\ & K(i, i) = 1, \quad \text{for all } i, \\ & K(i, j) \leq 1, \quad \text{for all } i, j. \end{aligned} \quad (6)$$

Proof we first relax the F_t in Eq. 5 as:

$$F'_t = \min_f \frac{(\langle l, K_v \rangle)^2}{4 \cdot \langle l^p, K_v \rangle} \quad (7)$$

This does not bring a significant reduction to the optimization formula since both $\langle l, K_v \rangle$ and $(\langle l, K_v \rangle)^2$ are

monotonously increasing when $l^T K_v > 0$. Now let's introduce an extra variable t such that $t \geq \frac{\langle l, K_v \rangle^2}{4 \cdot \langle l^p, K_v \rangle}$. As a consequence, the optimization equation in Eq. 7 is equivalent to:

$$\begin{aligned} \min : & \{t\} \\ \text{s.t. } & t \geq \frac{\langle l, K_v \rangle^2}{4 \cdot \langle l^p, K_v \rangle}. \end{aligned} \tag{8}$$

Furthermore, the constraint in Eq. 8 can be further transformed as:

$$\begin{aligned} 4t \cdot \langle l^p, K_v \rangle & \geq \langle l, K_v \rangle^2 \\ \Leftrightarrow (t + \langle l^p, K_v \rangle)^2 & \geq \langle l, K_v \rangle^2 + (t - \langle l^p, K_v \rangle)^2, \end{aligned}$$

which is a second order cone constraint:

$$t + \langle l^p, K_v \rangle \geq \left\| \begin{matrix} \langle l, K_v \rangle \\ t - \langle l^p, K_v \rangle \end{matrix} \right\|. \tag{9}$$

Combining Eqs. 7, 8 and 9 we get the conclusion in Theorem 1.

To further simplify Eq. 5, we here introduce Theorem 2.

Theorem 2 *If K is a positive semi-definite matrix of size $n \times n$, then*

$$K(i, j) \leq \max\{K(k, k)\}, \quad 1 \leq i, j, k \leq n. \tag{10}$$

Proof Since matrix K is positive semi-definite, we can decompose K as:

$$K = (k_1, k_2, \dots, k_n)^T (k_1, k_2, \dots, k_n), \tag{11}$$

where k_i is a vector of arbitrary dimension. So we have:

$$\begin{aligned} K(i, j) & = \langle k_i, k_j \rangle \\ & \leq \sqrt{\langle k_i, k_i \rangle \cdot \langle k_j, k_j \rangle} \\ & \leq \max\{K(i, i), K(j, j)\} \\ & \leq \max\{K(k, k)\}, 1 \leq k \leq n. \end{aligned} \tag{12}$$

By Theorem 2, we can introduce an extra positive semi-definite constraint to replace the last n^2 inequality constraints in Eq. 6, resulting in:

$$\begin{aligned} \min : & \{t\} \\ \text{s.t. } & t + \langle l^p, K_v \rangle \geq \left\| \begin{matrix} \langle l, K_v \rangle \\ t - \langle l^p, K_v \rangle \end{matrix} \right\|, \\ & K \succeq 0, \\ & K(i, i) = 1, \quad \text{for all } i, \\ & l^T K_v > 0, \quad l^p T K_v > 0. \end{aligned} \tag{13}$$

3.3 Smoothness controlling

While the number of constraints has been greatly reduced, it is necessary to refine Eq. 13 since the current prediction

function is not smooth enough in the Hilbert space. Note that if all the elements of K in Eq. 13 were the same, it would lead to $\langle l, K \rangle = 0$ and $\langle l^p, K \rangle = 0$. Consequently, the optimization problem would be infeasible since there were no interior points in the second order cone. We regard this case as over-smoothness. Therefore, we propose two strategies to smooth the predicted function and avoid over-smoothness.

Note that the ‘‘smoothness’’ of the manifold, which is measured by the Laplace–Beltrami operator on Riemannian manifolds, can be substituted by a discrete analogue operator defined as the graph Laplacian on the graph [30]. Therefore, we can employ Laplacian regularizer, denoted as S , to smooth the prediction function in its Hilbert space. Similar to Li et al.’s work [3], we here introduce a refined regularizer by incorporating a global normalized graph Laplacian³ into our optimization framework, which is defined as:

$$\begin{aligned} S & = \sum_{i,j=1}^n W'(i, j) \left\| \frac{\phi(x_i)}{\sqrt{D'(i, i)}} - \frac{\phi(x_j)}{\sqrt{D'(j, j)}} \right\|^2 \\ & = \text{tr}\{(I - (D')^{-\frac{1}{2}} W' (D')^{-\frac{1}{2}}) K\} \\ & = \langle \bar{l}, K_v \rangle \end{aligned} \tag{14}$$

where K, K_v are defined as the same as the previous ones, I is the identity matrix, W' is defined as:

$$W'(i, j) = \begin{cases} e^{-\frac{\|k_i - k_j\|^2}{2\sigma^2}}, & i \neq j \\ 0, & i = j, \end{cases} \tag{15}$$

where σ is a scale factor, $(I - (D')^{-\frac{1}{2}} W' (D')^{-\frac{1}{2}})$ is a normalized Laplacian matrix corresponding to W' , and $\bar{l} = \text{vec}(I - (D')^{-\frac{1}{2}} W' (D')^{-\frac{1}{2}})$ is the vectorization of the Laplacian matrix. Obviously, this formulation enforces smoothness over all the data globally. Finally, adding this term as an regularizer of the optimization problem in Eq. 13 results in the following optimization problem:

$$\begin{aligned} \min : & \{t + \lambda S\} \\ \text{s.t. } & t + \langle l^p, K_v \rangle \geq \left\| \begin{matrix} \langle l, K_v \rangle \\ t - \langle l^p, K_v \rangle \end{matrix} \right\|, \\ & K \succeq 0, \\ & K(i, i) = 1, \quad \text{for all } i \\ & \langle l, K_v \rangle > 0, \quad \langle l^p, K_v \rangle > 0, \end{aligned} \tag{16}$$

where λ is a parameter controlling the degree of smoothness on the predicted function.

As we stated before, over-smoothness is likely to happen. In addition, dropping the last two constraints

³ More satisfactory results might be attained if employing more sophisticated regularizers.

$l^T K_v > 0, l^{p^T} > 0$ is prone to numerical problem.⁴ We solve this by adding two extra terms into the optimization framework.

Remember that our goal is to make within-class samples as close as possible, and between-class samples as far as possible, this can be formulated in the following two formulas:

$$f_1 = \min \operatorname{tr}\{YLY^T\} \equiv \min\{-\langle w, K_v \rangle\}, \tag{17}$$

$$f_2 = \max \operatorname{tr}\{YL^p Y^T\} \equiv \min\{\langle w^p, K_v \rangle\}, \tag{18}$$

where w, w^p are the vectorization forms of the two similarity matrices W and W^p defined in Eq. 2. Denote these two regularizers as $S_1 = -\langle w, K_v \rangle$ and $S_2 = \langle w^p, K_v \rangle$, we can incorporate them into the objective function to get the final optimization formula as.

$$\begin{aligned} \min : & \{t + \lambda S + \lambda_1 S_1 + \lambda_2 S_2\} \\ \text{s.t. } & t + \langle l^p, K_v \rangle \geq \left\| \begin{matrix} \langle l, K_v \rangle \\ t - \langle l^p, K_v \rangle \end{matrix} \right\|, \\ & K \succeq 0, \\ & K(i, i) = 1, \quad \text{for all } i, \end{aligned} \tag{19}$$

where the two parameters λ_1 and λ_2 control the weights of the two graphs. As the two regularization terms are introduced, the two terms $\langle l, K_v \rangle$ and $\langle l^p, K_v \rangle$ will be enforced larger than zero. Therefore, we can drop the last two constraints in Eq. 16 in practice without causing any problem. Furthermore, a remarkable advantage of Eq. 19 is that it is a conic optimization programming (also SDP problem) which can be solved using the popular conic optimization software such as SeDuMi, which is of polynomial time complexity and has a theoretically proven $O(\sqrt{n} \log(\frac{1}{\epsilon}))$ worst-case iteration bound [31]. TRAGEK is illustrated in Algorithm 1.

Algorithm 1 The TRAGEK Algorithm

- 1: **Input:** data $X = (x_1, x_2, \dots, x_M)$, must-link constrains \mathcal{M} , cannot-link constrains \mathcal{C} , cluster number k .
- 2: **Output:** k clusters.
- 3: Construct the two Laplacian matrices in Eq. (5) as in graph embedding framework [30] using the provided must-link constraints \mathcal{M} and cannot-link constraints \mathcal{C} , respectively.
- 4: Construct the global smoothness Laplacian matrix defined in Eq. (14).
- 5: Solve the conic optimization programming defined in Eq. (19).
- 6: Run the kernel k-means algorithm on the learnt kernel matrix K to form k clusters.

4 ENPAKL: an efficient kernel learning algorithm

Although TRAGEK is a convex optimization problem which means that there exists a global optimal solution, the

⁴ Note that for two positive semi-definite matrices A and B , $\operatorname{tr}\{AB\} \geq 0$ holds, but not always > 0 . We thus can not drop the last two constraints directly. Otherwise it is easy to run into numerical problem. Because the constraint still holds if $\langle l^p, K_v \rangle$ is equal to a small enough positive constant, but this is far from the goal that $\langle l^p, K_v \rangle$ should be as large as possible.

computational cost is often relatively high and it often results in unstable solutions for large datasets, furthermore, it is sensitive to the parameter settings such as the choices of $\lambda, \lambda_1, \lambda_2$, etc. In this section, we propose to resolve these problems by two strategies. Firstly, we propose a constrained gradient descent based algorithm to make the learning procedure much stable and efficient. Secondly, we reduce the large-scale kernel learning problem into sub-kernel learning and combine these sub-kernels to approximate the global kernel matrix, this trick makes the computational complexity rely linearly on the number of data points. Experimental results in Sect. 5 show that the proposed strategies can approximate the true kernels well.

4.1 Constrained gradient descent

Note that the original optimization problem of Eq. 19 is not efficient enough, by taking the advantage of the iterative projection algorithm [19], we propose a constrained gradient descent based algorithm for training. The algorithm iteratively projects the solution obtained by gradient descent to the cones formed by the constraints.

Specifically, we want to avoid the SDP formulation above to reduce the computational complexity, thus we reformulate the original kernel learning problem of Eq. 19 in the following form:

$$\begin{aligned} F' &= \min_K \frac{\langle l, K_v \rangle}{\langle l^p, K_v \rangle} \\ \text{s.t. } & K \succeq 0, \\ & K(i, i) = 1, \quad \text{for all } i. \end{aligned} \tag{20}$$

Taking the logarithm of F' , and using the Laplacian smoother S in Eq. 19 as a regularization term, the objective function of ENPAKL is:

$$\begin{aligned} F &= \min_K \{\log(\langle l, K_v \rangle) - \log(\langle l^p, K_v \rangle) + \lambda \log(\langle \bar{l}, K_v \rangle)\} \\ \text{s.t. } & K \succeq 0, \\ & K(i, i) = 1, \quad \text{for all } i. \end{aligned} \tag{21}$$

It is straightforward to derive the gradient of F in Eq. 21:

$$\frac{\partial}{\partial K_v} F = \frac{l}{\langle l, K_v \rangle} - \frac{l^p}{\langle l^p, K_v \rangle} + \lambda \frac{\bar{l}}{\langle \bar{l}, K_v \rangle}. \tag{22}$$

Thus, we can update the kernel matrix K_v by constrained gradient descent as:

$$\begin{aligned} K_v^t &= K_v^{t-1} - \omega \cdot \frac{\partial}{\partial K_v^{t-1}} F \\ \text{s.t. } & K^t \succeq 0, \\ & K^t(i, i) = 1, \quad \text{for all } i, \end{aligned} \tag{23}$$

where ω is the step size for the current update, and t is the iteration index. Actually, we can regard Eq. 23 as the optimization problem on manifolds [32], however, instead of solving this problem directly on manifolds, we use the projection method by iteratively projecting the updated values into the cones formed by the constraints until converged. The constrained gradient descent algorithm is described in Algorithm 2.

Algorithm 2 Constraint Gradient Descent Algorithm

```

1: Input: Initial kernel matrix  $K_v^0$ , learning rate  $\omega$ .
2: Output: Kernel matrix  $K$ .
3: Calculate for the first update:  $K_v^1 = K_v^0 - \omega \cdot \frac{\partial}{\partial K_v} F$ .
4: Set  $t = 1$ .
5: while Not Converged do
6:   while Not Found rational  $K''$  satisfying the constraints do
7:     Solve:
           
$$K' = \arg \min_K \|K - K^t\|^2$$

           s. t.  $K \succeq 0$ ,
           
$$K'' = \arg \min_{K'} \|K - K'\|^2,$$

           s. t.  $K(i, i) = 1$ , for all  $i$ .
7:   
$$K_v^t = K_v^{t-1} - \omega \cdot \frac{\partial}{\partial K_v^{t-1}} F.$$

8:   end while
9:    $t = t + 1$ .
10:   $K_v^t = K_v^t - \omega \cdot \frac{\partial}{\partial K_v^t} F$ .
11: end while

```

Note that in Algorithm 2, the solution of Eqs. 24 and 25 can be calculated based on the following theorems:

Theorem 3 *The solution to Eq. 24 is to set all the negative eigenvalues of K_v to 0, that is, $K = V \max(D, 0) V^T$, where V, D are eigenvectors and eigenvalues of K_v .*

Theorem 4 *The solution to Eq. 25 is equal to K' , except that all the diagonal elements of K' are set to 1.*

Similar problem and proof for Theorem 3 can be found in [33]. Here we prove Theorem 4.

Proof of Theorem 4 Unfolding the norm and omitting the terms independent of K' , we can rewrite Eq. 25 as:

$$K'' = \arg \min_{K'} \{ KK^T - 2KK'^T \} \tag{26}$$

s.t. $tr\{KE_i\} = 1$.

where E_i is a matrix of the same size with K , and with all elements being 0 except the i th element of the diagonal being 1. Then the corresponding Lagrangian function with the corresponding Lagrangian multipliers λ_i 's is:

$$g(K, \lambda) = tr \left\{ KK^T - 2KK'^T - \sum_i \lambda_i (tr\{KE_i\} - 1) \right\}. \tag{27}$$

Taking the derivative of $g(K, \lambda)$ with respect to K , we have:

$$\begin{aligned} \frac{\partial}{\partial K} g(K, \lambda) &= 2K - 2K' - \sum_i \lambda_i E_i \\ &= 2K - 2K' - diag(\lambda). \end{aligned} \tag{28}$$

Setting Eq. 28 to zero leads to:

$$K = K' + \frac{1}{2} diag(\lambda). \tag{29}$$

Remember the constraint is $diag(K) = 1$, where 1 is a vector with all elements being 1, and note the form of the solution of K in Eq. 28, we can get the solution K by setting the diagonal elements of K' to 1. This completes the proof.

4.2 Learning the global kernel from sub-kernels

In this section, the second strategy to improve the efficiency of the proposed algorithm is presented. First note that it is not scalable and efficient enough for large datasets using Algorithm 2 directly, since we need to perform an eigen-decomposition for the current kernel matrix to solve Eq. 24, which is time consuming when the number of data points is large. Therefore, we propose to approximate the global kernel matrix using local kernel matrices (or sub-kernel matrices) formed by a subset of data points.

Suppose we start with a small subset of data (namely, m data points) denoted as $D = \{x_1, x_2, \dots, x_m\}$, and the corresponding sub-kernel matrix K_D has been learned using the constrained gradient descent algorithm described in Algorithm 2. The idea is to approximate the other elements of the global kernel matrix using this sub-kernel matrix. Note that because data of the same class in the feature space \mathcal{H} is assume to be flat (they are clustered into one point ideally in the feature space), it is reasonable to approximate all other data points $\phi(x_i)$ using the linear combination of this subset of data $\phi(D)$, that is: $\phi(x_i) = \sum_j w_{ij} \phi(x_j)$, where w_{ij} are the weights to be learned. There are two situations for $x_i \notin D$:

1. If x_i has at least one link constraint with some points x_j in D , according to our assumption, this means in the feature space, $\phi(x_i) = \phi(x_j)$, $x_j \in D$. Taking all such points into consideration, we relax $\phi(x_i)$ to be the linear combination of other points in the feature space, then we get $\phi(x_i) = \sum_j \xi \phi(x_j)$, where $w_{ij} = \xi$ is equal for all x_j .
2. If x_i has no link constraints with the points in D , then we approximate $\phi(x_i)$ using the weighted combination of $\phi(D)$ in the feature space. We assume these weights should be approximately the same with those learned by minimizing the reconstruction error in the original data space. This makes our approximation different from the one proposed by Yueng et al. [4]. While the objective function for w_{ij} is similar to local linear

embedding (LLE) [34], the definition of neighborhood is different. The objective function for w_{ij} is:

$$E = \min_{w_{ij}} \sum_i \left\| x_i - \sum_{j \in \mathcal{N}(i)} w_{ij} x_j \right\|^2, \tag{30}$$

where $\mathcal{N}(i)$ is defined to be the k nearest data points of x_i except for those having cannot link constraints with x_i .

To sum up, the weights w_{ij} 's are defined as:

$$w_{ij} = \begin{cases} 1 & x_i, x_j \in D \text{ and } i = j, \\ \frac{1}{T} & x_i \notin D, x_j \in D \text{ has a linked constraint,} \\ \text{Eq. (30)} & x_i \notin D, x_j \in D \text{ is } x_i\text{'s neighboring point,} \\ & \text{but has no cannot linked constraint,} \\ 0 & \text{otherwise,} \end{cases}$$

where T is the number of link constraints for $x_i \notin D$ and $x_j \in D$. Thus, the whole dataset in the feature space can now be written in the matrix form as:

$$\phi(X) = \phi(D)W^T, \tag{31}$$

where X is the whole dataset, $W = (w_1^T, w_2^T, \dots, w_n^T)^T$, $w_i = (w_{i1}, w_{i2}, \dots, w_{iM})^T$. Then, the whole kernel matrix can be approximated using Eq. 31 as:⁵

$$K_X = \phi(X)^T \phi(X) = WK_D W^T. \tag{32}$$

From Eq. 32 it can be seen that the kernel matrix K_X of the whole dataset can be approximated by a sub-kernel matrix K_D . However, given arbitrary must-link and cannot-link constraints, only one sub-kernel matrix might not approximate the whole kernel matrix well because all the pairwise constraints might not be included in one sub-dataset. To solve this problem, we propose a sub-data set picking schema that scales at most linearly with the size of the dataset to partition the whole dataset into several sub-data sets, then we use the corresponding sub-kernels to approximate the whole kernel. It can be proved that the computational complexity for this strategy is at most $O(n)$ times larger than that of using only one sub-kernel.

In this schema, we use the number of constraints (degrees of nodes in the graph) in the sub-data as the measure of the prior information this sub-data contains. The larger degree of one data point, the more prior

⁵ There are two points to be declared here. One is that it is easy to prove that K_X in Eq. 32 is a positive semi-definite matrix if K_D is positive semi-definite, this property makes K_X of the whole dataset still be a kernel matrix, which does not violate our objective. The second point is that for unknown points, in order to constrain the feature space be a hyperball, we need to normalize the weights calculated in Eq. 30 by dividing the weights by a normalized scalar $w_i^T K_D w_i$, that is, $w_i = \frac{w_i}{w_i^T K_D w_i}$.

information it has, and thus the higher probability the data point should be used to learn the sub-kernel. This sub-data set picking schema is described in Algorithm 3.

Algorithm 3 Sub-data sets picking schema

- 1: Sort the data points by the degrees of themselves.
- 2: Choose the first B data points that have the largest degrees as basic landmarks.
- 3: For the rest of the data points with degrees larger than 0, each time choose R data points in descending order of degrees, then combine them with the basic landmarks to get one sub-data set, loop until all data points have been chosen.

We can see from Algorithm 3 that the number of sub-data sets scales at most linearly with the number of the whole data points, thus is very efficient. Suppose at last we divide the whole dataset into L sub-data sets, and for each of such sub-data set a sub-kernel is learned by some kernel learning algorithms such as Algorithm 2, also we denote K_1, K_2, \dots, K_L as the approximated kernel matrices calculating using Eq. 32, then the final kernel matrix for the whole dataset is set to be:

$$K = \sum_i \alpha_i K_i, \quad \text{s.t. } \sum_i \alpha_i = 1, \tag{33}$$

where α_i is the weight for the i th kernel. In the experiments, we set α_i proportional to the total degree of their data points.

Note that this algorithm is efficient because it solves the original SDP problem of TRAGEK using a constraint gradient descend based algorithm, we will compare these two algorithms with respect to their efficiency and accuracy for clustering in the experiments.

5 Experiments

To test the proposed kernel learning algorithms TRAGEK and ENPAKL, we employed them for clustering and also used ENPAKL for image segmentation. We carried out the evaluations on two simulated datasets and ten datasets from the UCI machine learning repository [35]. The details of these datasets are tabulated in Table 1, where the first nine datasets have been often used in evaluating the performance of semi-supervised clustering algorithms [3, 28]. We compared ENPAKL with the PCP algorithm [3] and the SSKK algorithm [2] as well as the traditional k-means algorithm. We also investigated the influence of the number of pairwise constraints to the clustering performance for ENPAKL. To measure the clustering performance, we adopted the metric defined in [28]:

$$acc = \sum_{i>j} \frac{2 \cdot \mathcal{I}(\mathcal{I}(c_i, c_j), \mathcal{I}(\hat{c}_i, \hat{c}_j))}{n(n-1)}, \tag{34}$$

where $\mathcal{I}(a, b)$ is an indicator function returning 1 if $a = b$, and 0 otherwise, c denotes the true cluster

membership and \hat{c} denotes predicted cluster membership, and n is the number of samples. Without loss of generality, moreover, we set parameters λ , λ_1 and λ_2 in Eq. 19 to 1, and the scale factor σ in Eq. 15 to the average pairwise distance of the data set.

Other than this, we also applied the proposed ENPAKL together with the k -means and N -cut algorithms [36] on the MSRC Object Category Image Database (v2) [37] for image segmentation. Details are described in Sect. 5.5.

5.1 An illustrative experiment

To show that the proposed algorithms can propagate label information through the datasets, so that data in different classes can be separated as far as possible and those in the same class are clustered as close as possible, we run TRAGEK on the two synthetic datasets as used in [28], i.e., the chessboard and double-spiral datasets in Fig. 2. For better illustration, we rearranged the order of the data such that the first part of the data matrix belongs to one class and the last part to the other. It can be seen in Fig. 2 that the two classes are well separated. More specifically, we observed that the elements in the learned kernel K corresponding to the same class tend to 1 (black), while those corresponding to different classes tend to -1 (white). This means that the data from different classes are projected onto the opposite points on the hyperball of the feature space.

5.2 On small datasets

To compare the proposed kernel learning algorithms with some related algorithms for clustering, we tested them on nine small-scale data sets in Table 1 ranging from glass to wine. For ENPAKL, we set $B = 20, R = 10$ in Algorithm 3. Note that pairwise constraints are required for TRAGEK, ENPAKL, PCP, and SSKK, so we randomly generated k must-link constraints in each class and k cannot-link constraints between each two classes, where k ranges from 10 to 100 with an interval of 10. We thus have a total of $\frac{c(c+3c)k}{2}$ pairwise constraints for each experiment with a dataset of c classes. For each k , we randomly generated 20 different pairwise constraints, resulting in 20 different realizations of the pairwise constraints. The reported results were the average of the 20 different realizations together with 10 repetitions in the kernel k -means clustering step in Algorithm 1 for each pairwise constraint realization. The results are illustrated in Fig. 3.

We can observe from Fig. 3 that:

- TRAGEK outperforms the other three algorithms in all datasets except for the Iris and Wine datasets when the number of pairwise constraints is <30 .

Table 1 Twelve datasets used in our experiments

Data set	# Classes	Dimension	# Samples
Chessboard	2	2	100
Double-spiral	2	3	100
Glass	7	9	214
Heart	2	13	270
Iris	3	4	150
Protein	6	20	116
Sonar	2	60	208
Soybean	4	35	47
Wine	3	13	178
Wisconsin	2	30	720
Digital04	5	6	1000
Waveform	3	21	1800

The first two databases are artificial datasets, the rest ones are from the UCI machine learning repository

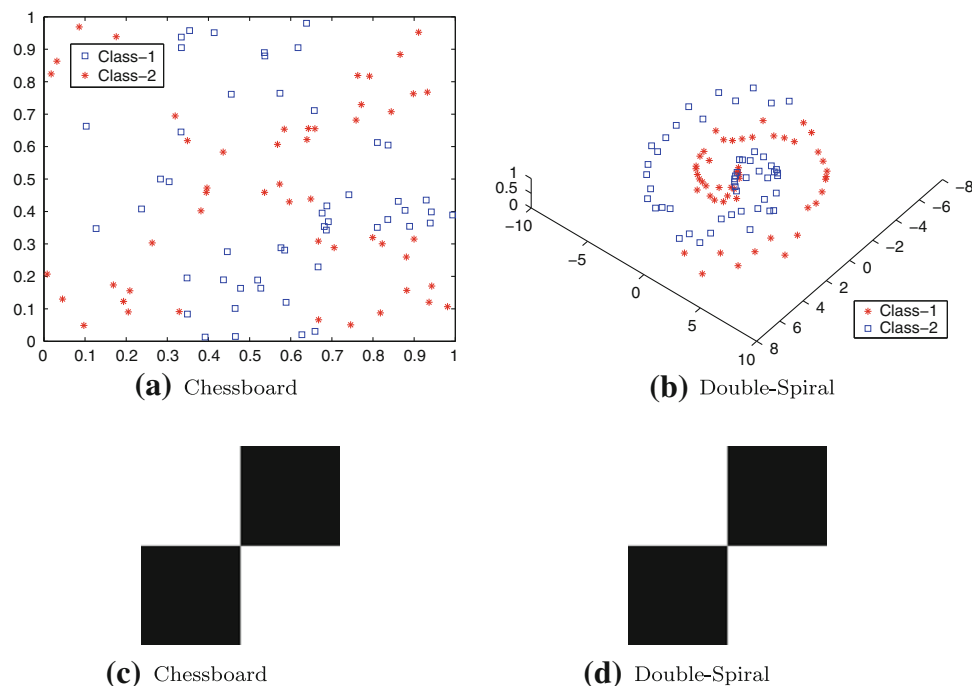
- PCP is worse than TRAGEK in clustering accuracy in most cases and it runs into numerical problems when the number of pairwise constraints is large or when some noisy constraints (constraints that are wrongly labeled) are added. Furthermore, we observed in the experiments that the running time of TRAGEK fluctuated little when varying the number of pairwise constraints, which can be seen in Sect. 5.4.
- ENPAKL approximates amazingly well to the original kernel learning problem TRAGEK, sometimes even gets better performance. Another merit of ENPAKL is that it is much faster than TRAGEK and PCP. We will give some examples below.

5.3 On larger datasets

Note that the datasets used in Sect. 5.2 are small, though often used in evaluating semi-supervised clustering algorithms [3, 28]. TRAGEK and other algorithms such as PCP can not scale well with large datasets and robust pairwise constraints. To evaluate the scalability of the proposed ENPAKL algorithm, we performed the experiments on three larger datasets described in the last three rows of Table 1. The `optdigits` dataset is a subset of a large digital dataset, and it contains the digits from 0 to 4 with each class containing 200 instances (Digital04); `Waveform` also comes from a large dataset and it has 1,800 instances.

The intrinsic disadvantage of PCP prevents it from being applied on such kind of large data with robust constraints. In order to enable it to work, we reduced the number of constraints by sampling. Specifically, the number of sampled constraints was set to the final number of constraints after the reduction in TRAGEK. We repeated the experiments for ten times with random sampling for the PCP algorithm, and

Fig. 2 Clustering results on chessboard and double-spiral datasets. **c, d** *Black color* means the corresponding values in the kernel matrix is 1, and *white color* means -1



picked up the best result and reported as the performance of PCP, denoted as PCP_{sample}. In this experiment, we varied the number of pairwise constraints relative to the number of total data samples, and the parameters in Algorithm 3 were set to $B = 100$, $R = 100$. Also ENPAKL¹ is a variant of ENPAKL by replacing the base sub-kernel learning algorithm in Algorithm 2 with PCP. The results for these algorithms are shown in Fig. 4.

It is found in the experiments that:

- ENPAKL is a little faster than ENPAKL¹, meanwhile both of them are much faster than PCP and SSKK.
- TRAGEK is apparently superior to PCP and SSKK, where these two algorithms even fail to compete with the traditional k-means algorithm.
- The performances of ENPAKL and ENPAKL¹ are competitive, and represent the best algorithms in terms of effectiveness and efficiency.

5.4 Running time

This section shows the running time of several related algorithms and we claim that the running time of TRAGEK is not sensitive to the number of pairwise constraints. To test this, we perform experiments on the Heart dataset and the Chessboard dataset with increasing number of pairwise constraints from 10 to 100 with an interval of 10. The results are shown in Fig. 5a and b. We can see from the figures that as the number of pairwise constraints increases, the running time of TRAGEK varies little, whereas that of PCP increases dramatically.

Next we examined the efficiency of ENPAKL. We used the Wisconsin dataset and recorded the corresponding running time. Note that PCP is too time consuming when the constraints are large, thus we do not show its running time here. We compared ENPAKL with TRAGEK, the results are shown in Fig. 5c. Obviously, ENPAKL is much more efficient than TRAGEK in term of computational complexity⁶.

5.5 Image segmentation

In this section we applied ENPAKL for image segmentation by doing clustering on images. We tested our algorithm on the MSRC Object Category Image Database (v2) [37], which contains 791 images of size approximately 320×210 , and includes different scenes such as grasses, forests, streets, etc. In this experiment, we do not care about what feature we used. Instead, we want to test the effectiveness and robustness of the proposed algorithm against other popular clustering algorithms such as k -means, N -cut, and etc.. As a result we simply used the histogram features in the experiments (richer features for image segmentation would be our future work). Specifically, we divided each image into 5×5 patches, and extracted the color histograms of each patch as its features, and finally used these features to do the segmentation. We set the number of clusters to the ground truth, for ENPAKL, we randomly generated 50 must-link and cannot-link

⁶ This experiment was run on an Intel Core 2 Duo CPU T6400 2.00 GHZ with 2 GB of DDR2 memory

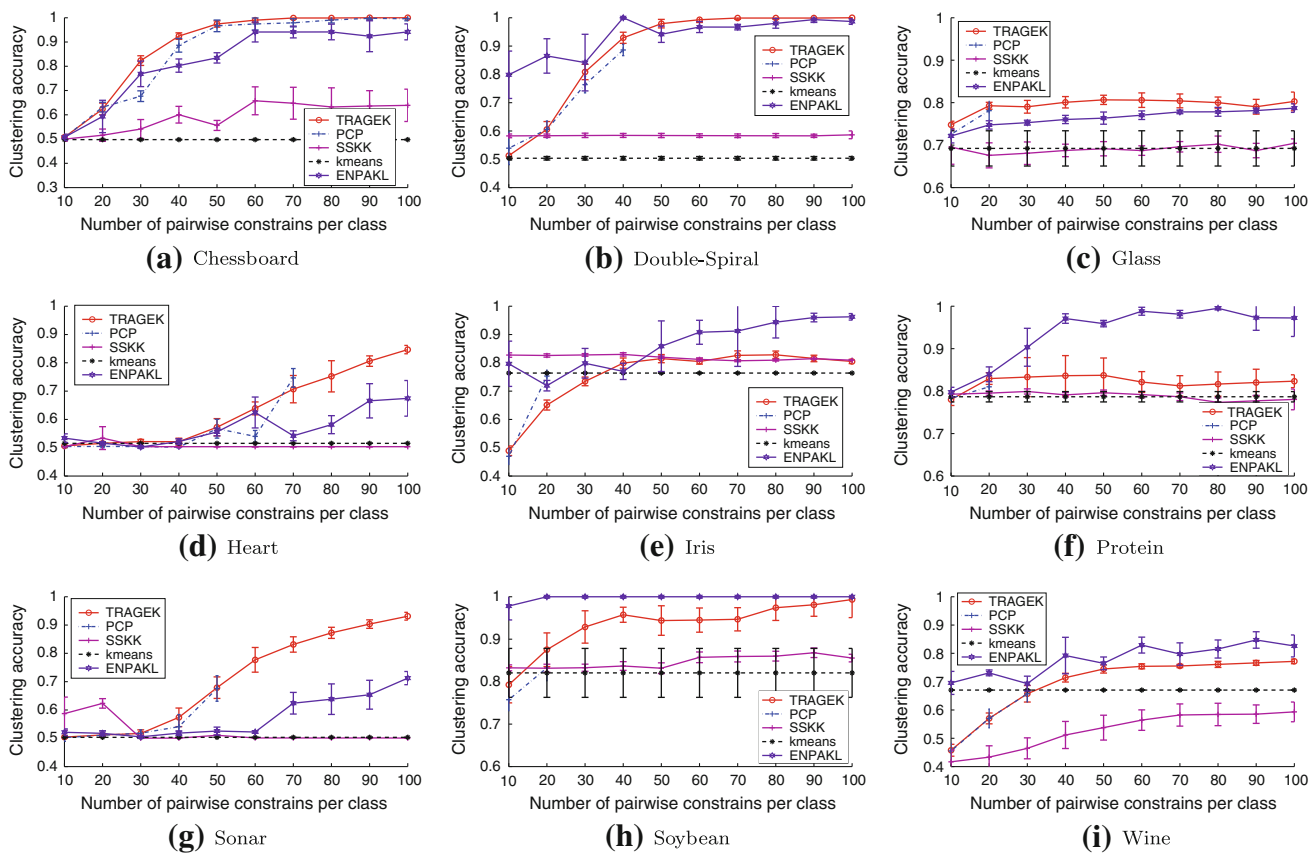


Fig. 3 Clustering performance on nine small UCI datasets

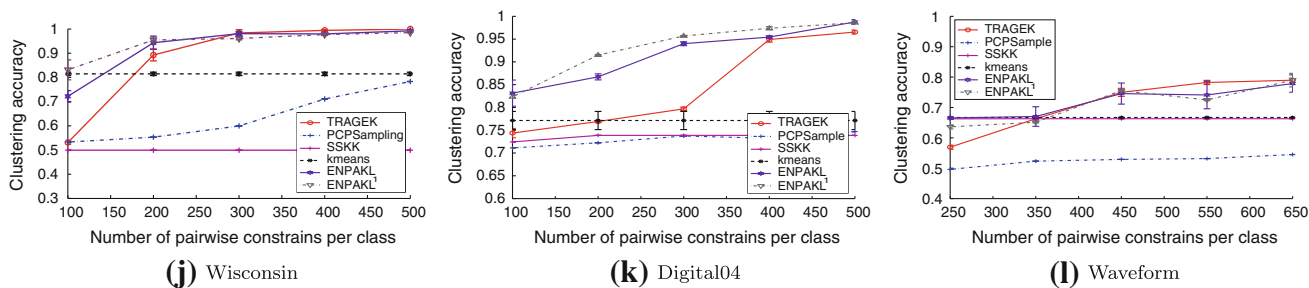


Fig. 4 Clustering performance on three large UCI datasets

constraints for each cluster in the images. For simplicity, we compared the proposed ENPAKL algorithm with the k -means and N -cut algorithms⁷ which are popularly used in image segmentation, and also because the above experiments have shown the superior of the k -means algorithm over PCP and SSKK. Some examples of the images and their segmentation results are shown in Fig. 6. From these results we can see the superior of ENPAKL over the k -means and N -cut algorithms in term of segmentation accuracy, though it runs much slower, which is a typical

problem for kernel based algorithms.⁸ We used Eq. 34 as the segmentation accuracy criterion, and the corresponding accuracies are also shown in the figure. We see from the figure that ENPAKL performs best while N -cut and k -means are comparable. Also note that for some images, the segmentations learned by ENPAKL are very close to the Ground Truth, while those learned by the k -means and the

⁷ We used an efficient implementation of the N -cut algorithm in [38]

⁸ The k -means algorithm takes about 1 s for one image, the N -cut algorithm takes about 2 s, while ENPAKL needs about 5 min, and PCP cannot run in this experiment because the corresponding data is too large. How to accelerate the speed of the proposed algorithm further is our future work.

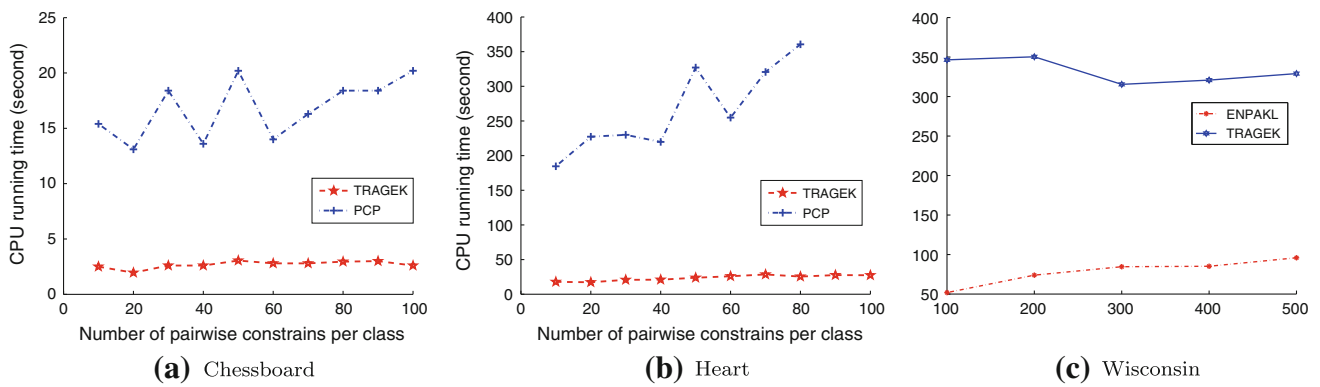


Fig. 5 Running time comparison. The x -axis is the number of constraints, the y -axis represents the running time in seconds

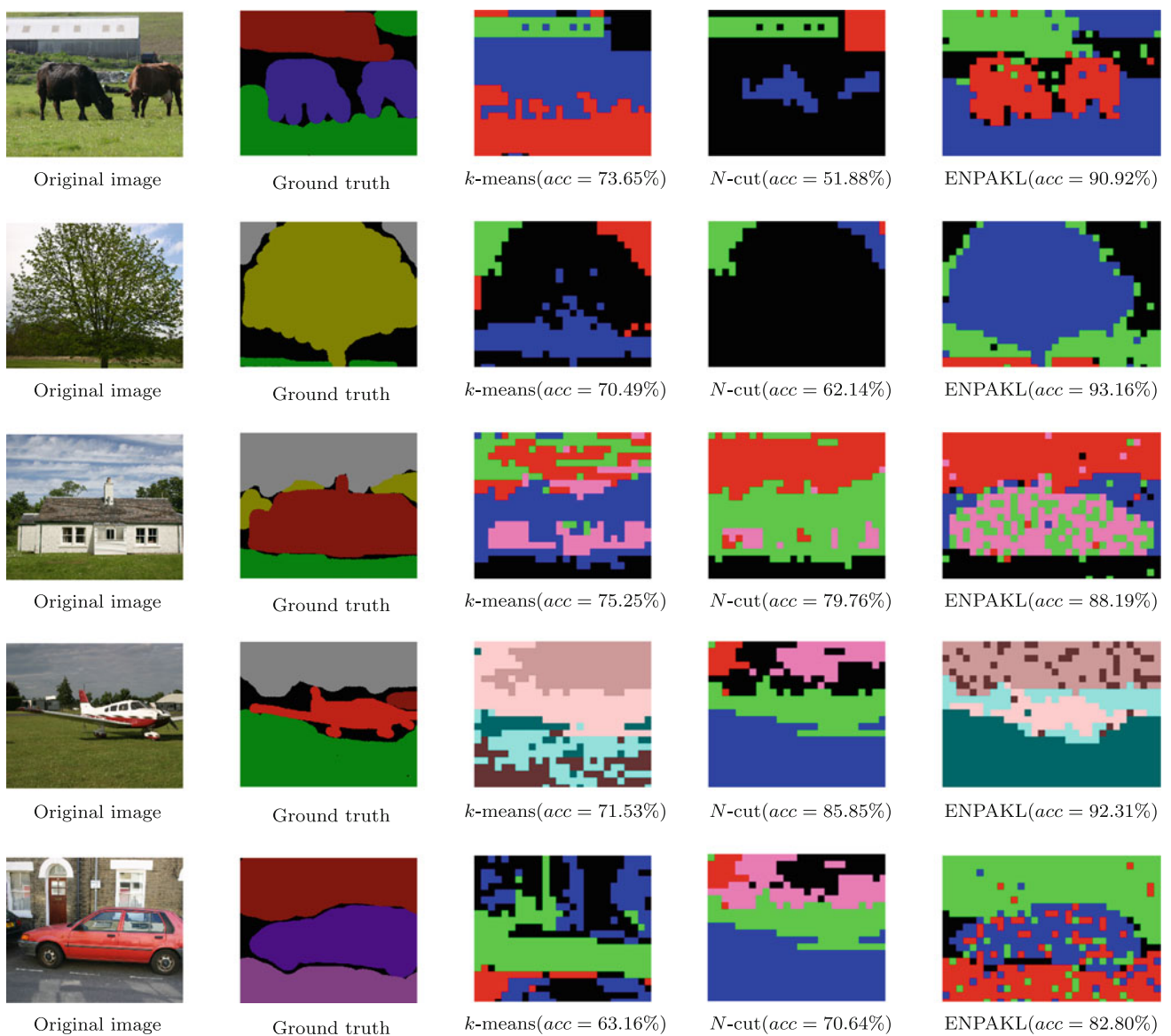


Fig. 6 Image segmentation using ENPAKL, k -means and N -cut. Here acc means segmentation accuracy evaluating using Eq. 34

N -cut are much worse, this indicates that supervisory information could help image segmentation a lot, and it is encouraged to use such kind of information to boost the segmentation accuracy. We believe better segmentation results can be obtained by choosing the constraints carefully, by using other kinds of features such as the sift features [39] and rich textual features [40], and also by taking the spatial information into consideration.

6 Conclusion

In this paper, we proposed a non-parametric kernel learning framework. It generalizes the graph embedding framework [6] into kernel space and is reformed as a conic optimization programming. A global Laplacian regularizer is used to smooth the functional space. Two algorithms are proposed for the corresponding kernel learning problem, one is to solve the original optimization problem through semi-definite programming. The other is to relax the SDP problem and solve with a constrained gradient descent based algorithm. To further reduce the computational complexity, the whole data is proposed to be divided into groups, and sub-kernels for these groups are learned separately, then the global kernel is constructed by combining these sub-kernels. Experiments are performed on nine datasets for clustering and one image dataset for image segmentation. Experimental results show that the proposed ENPAKL algorithm is superior to the recently developed algorithms [2, 3] in terms of computational effectiveness and clustering accuracy, and often achieves better image segmentation.

We will study the parameters setting problem in the future. For example, the regularizer S in Eq. 19 may be replaced by a more sophisticated regularizer such as the s -weighted Laplacian operator [41]. The algorithms should also be evaluated with different settings of B and R in Algorithm 3, the k in the graph construction, etc. Furthermore, we can incorporate ENPAKL into other kernel methods such as kernelization of some dimensional reduction algorithms. In addition, we will apply the proposed algorithm to more real applications, and explore more efficient algorithms for this problem since the current methods is not fast enough for large scale datasets.

Acknowledgments This work was supported in part by the NFSC (No. 60975044, 61073097) and 973 Program (No. 2010CB327900).

References

- Chapelle O, Schölkopf B, Zien A (2006) Semi-supervised learning. MIT Press, Cambridge
- Kulis B, Basu S, Dhillon I, Mooney R (2005) Semi-supervised graph clustering: a kernel approach. In: Proceedings of the 22nd international conference on machine learning, pp 457–464
- Li Z, Liu J, Tang X (2008) Pairwise constraint propagation by semidefinite programming for semi-supervised classification. In: Proceedings of the 25th international conference on machine learning, pp 576–583
- Yeung D-Y, Chang H, Dai G (2007) A scalable kernel-based algorithm for semi-supervised metric learning. In: Proceedings of the 20th international joint conference on artificial intelligence, pp 1138–1143
- Cortes C, Mohri M, Rostamizadeh A (2009) Learning non-linear combinations of kernels. In: Advances in neural information processing systems, vol 21
- Yan SC, Xu D, Zhang BY, Zhang H-J, Yang Q, Lin S (2007) Graph embedding and extensions: a general framework for dimensionality reduction. IEEE Trans Pattern Anal Mach Intell 29(1):40–51
- Yang J, Yan SC, Fu Y, Li XL, Huang TS (2008) Non-negative graph embedding. In: Proceedings of IEEE computer society conference on computer vision and pattern recognition
- Schölkopf B, Smola AJ (2001) Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT Press, Cambridge
- Cortes C, Mohri M, Rostamizadeh A (2010) Two-stage learning kernel algorithms. In: Proceedings of the 27th international conference on machine learning
- Cortes C, Mohri M, Rostamizadeh A (2010) Generalization bounds for learning kernels. In: Proceedings of the 27th international conference on machine learning
- Jin R, Hoi SCH, Yang T (2010) Online multiple kernel learning: algorithms and mistake bounds. In: Proceedings of the 21st international conference algorithmic learning theory, pp 390–404
- Baraldi A, Blonda P (1999) A survey of fuzzy clustering algorithms for pattern recognition—part II. IEEE Trans Syst Man Cybern Part B 29(6):786–801
- Yang MS, Wu KL, Hsieh JN, Yu J (2008) Alpha-cut implemented fuzzy clustering algorithms and switching regressions. IEEE Trans Syst Man Cybern Part B 38(3):904–915
- Trappey AJC, Trappey CV, Hsu F-C, Hsiao DW (2009) A fuzzy ontological knowledge document clustering methodology. IEEE Trans Syst Man Cybern Part B 39(3):123–131
- Xiong H, Wu J, Chen J (2009) k-Means clustering versus validation measures: a data-distribution perspective. IEEE Trans Syst Man Cybern Part B 39(2):318–331
- Basu S, Bilenko M, Mooney R (2004) A probabilistic framework for semi-supervised clustering. In: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, pp 59–68
- Bilenko M, Basu S, Mooney R (2004) Integrating constraints and metric learning in semi-supervised clustering. In: Proceedings of the 21st international conference on machine learning, pp 81–89
- Kamvar SD, Klein D, Manning C (2003) Spectral learning. In: Proceedings of the 18th international joint conference on artificial intelligence, pp 561–566
- Xing EP, Ng AY, Jordan MI, Russell S (2003) Distance metric learning, with application to clustering with side-information. In: Advances in neural information processing systems, vol 15
- Wagstaff K, Cardie C, Rogers S, Schroedl S (2001) Constrained k-means clustering with background knowledge. In: Proceedings of the 18th international conference on machine learning, pp 798–803
- Hong Y, Kwong S (2009) Learning assignment order of instances for the constrained k-means clustering algorithm. IEEE Trans Syst Man Cybern Part B 39(2):568–574
- Lu Z, Leen TK (2005) Semi-supervised learning with penalized probabilistic clustering. In: Advances in neural information processing systems, vol 17, pp 849–856

23. Bar-Hillel A, Hertz T, Shental N, Weinshall D (2005) Learning a Mahalanobis metric from equivalence constraints. *J Mach Learn Res* 6:937–965
24. Hertz T, Bar-Hillel A, Weinshall D (2004) Boosting margin based distance function for clustering. In: *Proceedings of the 21st international conference on machine learning*, pp 393–400
25. Xu Z, Dai M, Meng D (2009) Fast and efficient strategies for model selection of Gaussian support vector machine. *IEEE Trans Syst Man Cybern Part B* 39(5):1292–1307
26. Dhillon I, Guan Y, Kulis B (2004) Kernel k-means, spectral clustering and normalized cuts. In: *Proceedings of the tenth ACM SIGKDD international conference on knowledge discovery and data mining*, pp 551–556
27. Bousquet O, Herrmann D (2003) On the complexity of learning the kernel matrix. In: *Advances in neural information processing systems*, vol 15, pp 399–406
28. Hoi SCH, Jin R, Lyu MR (2007) Learning nonparametric kernel matrices from pairwise constraints. In: *Proceedings of the 24th international conference on machine learning*, pp 361–368
29. Zhuang J, Tsang IW, Hoi SCH (2009) SimpleNPKL: simple nonparametric kernel learning. In: *Proceedings of the 26th international conference on machine learning*, pp 1273–1280
30. Zhou DY, Huang J, Schölkopf B (2005) Learning from labeled and unlabeled data on a directed graph. In: *Proceedings of the 22nd international conference on machine learning*, pp 1036–1043
31. Sturm JF (1999) Using SeDuMi 1.02, a matlab toolbox for optimization over symmetric cones. *Optim Methods Softw* 11(2):625–653
32. Adler RL, Dedieu JP, Margulies JY, Martens M, Shub M (2002) Newton’s method on Riemannian manifolds and a geometric model for the human spine. *IMA J Numer Anal* 22(3):359–390
33. Golub GH, Loan CFV (1996) *Matrix computation*. Johns Hopkins University Press, Baltimore
34. Roweis S, Saul L (2000) Nonlinear dimensionality reduction by locally linear embedding. *Science* 290(5500):2323–2326
35. Asuncion A, Newman DJ (2007) UCI machine learning repository. University of California, Irvine, School of Information and Computer Sciences. <http://www.ics.uci.edu/~mllearn/MLRepository.html>
36. Shi J, Malik J (2000) Normalized cuts and image segmentation. *IEEE Trans Pattern Anal Mach Intell* 22(8):888–905
37. Criminisi A MSRC Category Image Database (v2). MSRC. http://research.microsoft.com/en-us/um/people/antcrim/data_objrec/msrc_objcategoimagedatabase_v2.zip
38. Shi J MATLAB Normalized Cuts Segmentation Code. <http://www.cis.upenn.edu/jshi/software/>
39. Lowe DG (2004) Distinctive image features from scale-invariant keypoints. *Int J Comput Vis* 60(2):91–110
40. Tan B, Zhang J, Wang L (2011) Semi-supervised elastic net for pedestrian counting. *Pattern Recognit* 44(10–11):2297–2304
41. Duchenne O, Audibert J-Y, Keriven R, Ponce J, Segonne F (2008) Segmentation by transduction. In: *Proceedings of IEEE computer society conference on computer vision and pattern recognition*