

SimpleMKKM: Simple Multiple Kernel K-means

Xinwang Liu, Li Liu, Jian Xiong, En Zhu, Junwei Han, Meng Wang, Dinggang Shen, and Wen Gao

Abstract—We propose a simple yet effective multiple kernel clustering algorithm, termed simple multiple kernel k-means (SimpleMKKM). It extends the widely used supervised kernel alignment criterion to multi-kernel clustering. Our criterion is given by an intractable minimization-maximization problem in the kernel coefficient and clustering partition matrix. To optimize it, we re-formulate the problem as a smooth minimization one, which can be solved efficiently using a reduced gradient descent algorithm. We theoretically analyze the performance of SimpleMKKM in terms of its clustering generalization error. Furthermore, we develop comprehensive experiments to study the proposed SimpleMKKM from the perspective of clustering accuracy, advantage on the formulation and optimization, evolution of the learned consensus clustering matrix, clustering with number of sample, clustering with number of base kernels, the learned kernel weight analysis, the running time and convergence. As indicated, our algorithm delivers its effectiveness by significantly and consistently outperforming state of the art multiple kernel clustering alternatives. Our work provides a more effective approach to fuse multi-view data for clustering, which could trigger novel research on multiple kernel clustering. Our codes and data are publicly available at <https://xinwangliu.github.io/>.

Index Terms—multiple kernel clustering, multiple view learning, kernel alignment maximization

1 INTRODUCTION

IN Multi-view clustering (MVC) [1], we aim to combine a set of pre-specified kernel matrices to improve clustering performance. These kernel matrices could encode heterogeneous sources or views of the data [2], [3], [4]. One popular method, multiple kernel k-means (MKKM) [5], has been studied intensively and used in various applications [2], [6], [7], [8], [9], [10], [11]. The approach is attractive also from a theoretical perspective, as it unifies the search of the optimal base kernel coefficient and the clustering partition matrix into a single objective function, which is usually solved by using two-step alternating optimization on the coefficients and clustering partition matrix.

Several variants of MKKM have been developed to further improve the clustering performance [2], [6], [12], [13], [14]. Notably, [6] substantially increases the expressiveness of MKKM by allowing for a locally adaptive kernel mixtures, which can better capture sample-specific characteristics of data. [12] proposes an extension that optimizes

a localized kernel alignment criterion. It aligns the local density of the samples given by the k -nearest neighbours with an ideal similarity matrix. This alignment helps to keep neighbouring sample pairs together, which avoids unreliable similarity evaluation. Such an alignment helps the clustering algorithm to focus on neighboring sample pairs, in that they shall stay together. This avoids unreliable similarity evaluation for farther sample pairs. Observing that existing MKKM algorithms do not sufficiently consider the correlation among these kernels, [13] employs matrix regularization to reduce the redundancy and enhance the diversity of the selected kernels. Most of existing MKKM algorithms assume that the optimal kernel is a linear combination of a group of base kernels. This assumption is challenged in [14], who proposes an optimal neighborhood kernel clustering (ONKC) algorithm to enhance the representability of the optimal kernel and strengthen the negotiation between kernel learning and clustering. More recently, MKKM algorithms have been extended to handle missing views [15]. By assuming the optimal kernel is a linear combination of the base kernel matrices, [16] develop a minimization-maximization framework that aims to be robust to adversarial perturbation. More recently, many work has been devoted to extend existing MKKM to handle multiple kernel clustering with incomplete kernels [15], [17], [18], [19]. All these variants potentially improve standard MKKM and achieve promising clustering performance in various applications.

The objective functions of the mentioned methods differ, but they all share one commonality: they learn the kernel coefficient and the clustering partition matrix *jointly*. By this way, the learned kernel coefficient can best serve the clustering, leading to superior clustering performance. However, simultaneously solving for the kernel coefficients *and* the clustering partition is intractable. One commonly adopted remedy is to decouple the optimization of the kernel coefficients and the clustering partition through a block coordinate descent algorithm, which optimizes the two alternately.

- X. Liu and E. Zhu are with College of Computer, National University of Defense Technology, Changsha, 410073, China. E-mail: {xinwangliu,enzhu}@nudt.edu.cn.
- L. Liu is with College of System Engineering, National University of Defense Technology, Changsha, China, and also with the Center for Machine Vision and Signal Analysis, University of Oulu, 90014 Oulu, Finland (E-mail: li.liu@oulu.fi).
- J. Xiong is with School of Business Administration, Southwestern University of Finance and Economics, Chengdu, Sichuan, 611130, China (E-mail: xiongjian2017@swufe.edu.cn).
- J. Han is with School of Automation, Northwestern Polytechnical University, Xian, China (E-mail: junwei.han2010@gmail.com).
- M. Wang is with School of Computer Science and Information Engineering, Hefei University of Technology, Hefei, China (E-mail: eric.mengwang@gmail.com).
- D. Shen is with Department of Radiology and BRIC, University of North Carolina at Chapel Hill, North Carolina 27599, USA, and also with Department of Brain and Cognitive Engineering, Korea University, Seoul 02841, Republic of Korea (E-mail: dgshen@med.unc.edu).
- W. Gao is with School of Electronics Engineering and Computer Science, Peking University, Beijing, China, 100871 (E-mail: wgao@pku.edu.cn).

Manuscript received November 18, 2020.

This means, one block of variables is minimized while the other is kept fixed. However, such alternate optimization algorithms can get trapped into a local optima of the objective function. As a remedy, [12], [13] propose regularization strategies to avoid getting trapped into local minimum. The incorporation of these regularization terms comes at a price: the approach has additional hyper-parameters, which are difficult to select, given the unsupervised nature of clustering tasks.

In this paper, we propose Simple MKKM (SimpleMKKM)—a novel formulation for multiple kernel clustering that addresses the aforementioned shortcomings. Unlike previous approaches, SimpleMKKM optimizes the unsupervised kernel alignment criterion directly. Specifically, it minimizes kernel alignment with respect to the kernel coefficient and maximizes it with respect to the clustering matrix. This minimization-maximization optimization problem cannot readily be solved using existing alternate optimization frameworks. However, we show that this min-max problem actually leads to a more efficient and effective optimization algorithm. Specifically, we reformulate the min-max problem as a minimization problem, whose objective relies on the known optimal solution to kernel k-means. We then prove the differentiability of the optimal value function and calculate its reduced gradient. This leads to a solution using a reduced gradient descent algorithm, without alternating optimization. We show a generalization error bound for our approach, thus theoretically guaranteeing its clustering performance. We conduct comprehensive experiments on eleven benchmark datasets, where we compare SimpleMKKM to eight baseline methods in terms of four common evaluation criteria. We observe that SimpleMKKM consistently outperforms its competitors. Moreover, we conduct extra experimental study from the following aspects: advantage on the formulation and optimization, evolution of the learned consensus clustering matrix, clustering with number of sample, clustering with number of base kernels, the learned kernel weight analysis, the running time and convergence.

We end up this section by summarizing the main contributions of this paper as follows:

- We develop a simple while effective criterion for multiple kernel clustering, which is given by an intractable minimization-maximization. The problem is reformulated as a smooth minimization, which can be solved efficiently using reduced gradient descent.
- We theoretically analyze the performance of SimpleMKKM in terms of its clustering generalization error on test data.
- We conduct comprehensive experiments to validate the effectiveness of the proposed algorithm.

In addition, the proposed SimpleMKKM is parameter-free, making it readily applicable in practice. More importantly, SimpleMKKM can be taken as a strong baseline to trigger new research on multiple kernel clustering.

2 RELATED WORK

In this section, we briefly review the most related, including multiple kernel k-means (MKKM) and robust MKKM

clustering using min-max optimization [16].

2.1 MKKM

Given a group of pre-calculated kernel matrices $\{\mathbf{K}_p\}_{p=1}^m$, MKKM assumes that the optimal kernel matrix \mathbf{K}_γ can be parameterized as $\mathbf{K}_\gamma = \sum_{p=1}^m \gamma_p^2 \mathbf{K}_p$, where $\gamma \in \Delta = \{\gamma \in \mathbb{R}^m \mid \sum_{p=1}^m \gamma_p = 1, \gamma_p \geq 0, \forall p\}$ represents the kernel weights of these base kernel matrices. It jointly learns the kernel weights γ and the clustering partition matrix \mathbf{H} by optimizing Eq. (1).

$$\min_{\gamma \in \Delta} \min_{\mathbf{H}} \text{Tr}(\mathbf{K}_\gamma(\mathbf{I} - \mathbf{H}\mathbf{H}^\top)) \quad (1)$$

$$s.t. \mathbf{H} \in \mathbb{R}^{n \times k}, \mathbf{H}^\top \mathbf{H} = \mathbf{I}_k.$$

In literature, the optimization problem in Eq. (1) is usually be solved by alternatively updating \mathbf{H} and γ : (i) **Optimizing \mathbf{H} given γ** . With the kernel coefficients γ fixed, \mathbf{H} can be obtained by solving a kernel k-means clustering optimization problem; (ii) **Optimizing γ given \mathbf{H}** . With \mathbf{H} fixed, γ can be optimized via solving the following quadratic programming with linear constraints,

$$\min_{\gamma \in \Delta} \sum_{p=1}^m \gamma_p^2 \text{Tr}(\mathbf{K}_p(\mathbf{I}_n - \mathbf{H}\mathbf{H}^\top)), \quad (2)$$

which has a closed-form solution.

As noted in [2], [6], using a convex combination of kernels $\sum_{p=1}^m \gamma_p \mathbf{K}_p$ to replace $\sum_{p=1}^m \gamma_p^2 \mathbf{K}_p$ is not a viable option, because this could make only one single kernel activate and all the others assigned with zero weight, as seen from Eq. (2). Other recent work using ℓ_2 -norm combinations can be found in [15], [20], [21].

2.2 Robust MKKM Using Min-Max Optimization

Recently, [16] proposed a MKKM clustering method with the aim to be robust against adversarial perturbation. To achieve this goal, the authors use a $\min_{\mathbf{H}}\text{-max}_{\gamma}$ formulation that combines views so as to achieve high within-cluster variance in the combined space \mathbf{W}_γ and then updates clusters by minimizing such variance. Its optimization problem is,

$$\min_{\mathbf{H}} \max_{\gamma \in \Theta} \text{Tr}(\mathbf{W}_\gamma(\mathbf{I} - \mathbf{H}\mathbf{H}^\top)) \quad (3)$$

$$s.t. \mathbf{H} \in \mathbb{R}^{n \times k}, \mathbf{H}^\top \mathbf{H} = \mathbf{I}_k,$$

where $\Theta = \{\gamma \in \mathbb{R}^m \mid \sum_{p=1}^m \gamma_p^2 \leq 1, \gamma_p \geq 0, \forall p\}$ and $\mathbf{W}_\gamma = \sum_{p=1}^m \gamma_p \mathbf{K}_p$.

Note that in contrast to Eq. (1), the above approach adopts an ℓ_2 -norm constraint on the kernel weights to avoid sparse solutions. It is observed that using an ℓ_2 -norm constraint can obtain non-sparse kernel coefficients, which is helpful to better utilize the complementary information in the data. Similar to MKKM, the problem in Eq. (3) can be solved by following the same alternate optimization framework.

Although the objective functions of MKKM and its variants may vary, they share a common alternate optimization routine. The aforementioned alternate framework could cause the optimization w.r.t γ to produce high redundant or overly sparse solutions [13]. This in turn would make the multiple kernel matrices less utilized, and adversely affects

the clustering performance. A direct remedy is to incorporate some regularization on γ to help its optimization [12], [13]. However, the incorporation of regularization may introduce extra hyper-parameters. How to determine those in unsupervised learning tasks such as clustering is difficult. In the following, we introduce our simple MKKM objective, and design a novel optimization procedure for it that avoids these issues.

3 SIMPLEMKKM: SIMPLE MKKM

In this section, we first give the proposed SimpleMKKM kernel alignment-based objective. We then reformulate it as the minimization of an optimal value function, and prove its differentiability. After that, we develop a reduced gradient descent algorithm to solve it efficiently and effectively.

3.1 SimpleMKKM Formulation

Kernel alignment criterion has been widely used for kernel tuning in supervised learning due to its simplicity and effectiveness [22], [23]. Our new formulation is based on unsupervised multiple kernel alignment criterion, inspired by existing supervised kernel learning. One can optimize this criterion by maximizing over both γ and \mathbf{H} . Though theoretically elegant, we empirically observe that such $\max_{\gamma} \max_{\mathbf{H}}$ formulation does not achieve promising clustering performance, which is different from supervised kernel learning. We conjecture this is caused by the over-fitted optimization between γ and \mathbf{H} . On the other hand, from the optimization perspective of MKKM in Eq. (1), $\text{Tr}(\mathbf{K}_{\gamma}(\mathbf{I} - \mathbf{H}\mathbf{H}^{\top}))$ should be minimized. This objective can be decomposed into two terms, $\text{Tr}(\mathbf{K}_{\gamma})$ and $-\text{Tr}(\mathbf{K}_{\gamma}\mathbf{H}\mathbf{H}^{\top})$. The first term can be regarded as regularization on γ , which should be optimized via minimizing γ . The other one is the opposite of kernel alignment, which should be minimized via maximizing \mathbf{H} . By taking both regularisation and partitioning into account, our SimpleMKKM proposes to optimize the kernel alignment criterion by minimizing γ and maximizing \mathbf{H} as:

$$\min_{\gamma \in \Delta} \max_{\mathbf{H}} \text{Tr}(\mathbf{K}_{\gamma}\mathbf{H}\mathbf{H}^{\top}) \quad (4)$$

$$\text{s.t. } \mathbf{H} \in \mathbb{R}^{n \times k}, \mathbf{H}^{\top}\mathbf{H} = \mathbf{I}_k,$$

where $\Delta = \{\gamma \in \mathbb{R}^m \mid \sum_{p=1}^m \gamma_p = 1, \gamma_p \geq 0, \forall p\}$ and $\mathbf{K}_{\gamma} = \sum_{p=1}^m \gamma_p^2 \mathbf{K}_p$.

Though simple, the SimpleMKKM formulation in Eq. (4) has the following merits: (1) It is the first MKKM objective that, strictly coincides with the kernel alignment criterion via $\text{Tr}(\mathbf{K}_{\gamma}\mathbf{H}\mathbf{H}^{\top})$ to tune kernel weights. In contrast, MKKM and its all variants adopt $\text{Tr}(\mathbf{K}_{\gamma}(\mathbf{I} - \mathbf{H}\mathbf{H}^{\top}))$ as the criterion by extending the objective of classic kernel k-means to multiple kernels. It is worth noting that the kernel alignment criterion is more general and can be used for any kernel tuning tasks. As a result, it can be used for multiple kernel clustering. (2) According to [16], regularisation by min-max optimization of γ and \mathbf{H} generates more robust clusters by avoiding overfitting to noisy views or datapoints. (3) As we shall see next, while our formulation looks intractable, it actually leads to a more efficient and effective optimisation algorithm than the standard alternating strategies used for MKKM. Furthermore, unlike alternatives [12], [13] relying on regularisation by penalizing γ , SimpleMKKM introduces

no additional parameters beyond the number of clusters to form.

Our new formulation in Eq. (4) cannot be readily solved by the widely adopted alternate optimization strategy, as done in MKKM and its variants. In the following, we design an efficient and effective reduced gradient descent algorithm. Firstly, we equivalently rewrite the optimization in Eq. (4) as,

$$\min_{\gamma \in \Delta} \mathcal{J}(\gamma), \quad (5)$$

with

$$\mathcal{J}(\gamma) = \left\{ \max_{\mathbf{H}} \text{Tr}(\mathbf{K}_{\gamma}\mathbf{H}\mathbf{H}^{\top}) \text{ s.t. } \mathbf{H}^{\top}\mathbf{H} = \mathbf{I}_k \right\}. \quad (6)$$

In this way, the min-max optimization is transformed to a minimization one, where its objective is a kernel k-means optimal value function. In the following, we first prove the differentiability of $\mathcal{J}(\gamma)$, and apply the reduced gradient descent algorithm to decrease Eq. (5).

3.2 The Calculation of Reduced Gradient

In the literature, several works discuss the existence and computation of derivatives of optimal value functions $\mathcal{J}(\gamma)$ [24], [25], [26]. The most appropriate reference for our case is Theorem 4.1 in [24], which has already been utilized to tune the hyper-parameters of SVM [25] and optimize the kernel weights in multiple kernel learning [26]. The following Theorem 1 shows that $\mathcal{J}(\gamma)$ in Eq. (5) is differentiable.

Theorem 1. $\mathcal{J}(\gamma)$ in Eq. (6) is differentiable. Further, $\frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_p} = 2\gamma_p \text{Tr}(\mathbf{K}_p \mathbf{H}^* \mathbf{H}^{*\top})$, where $\mathbf{H}^* = \{\arg \max_{\mathbf{H}} \text{Tr}(\mathbf{K}_{\gamma} \mathbf{H}\mathbf{H}^{\top}) \text{ s.t. } \mathbf{H}^{\top}\mathbf{H} = \mathbf{I}_k\}$.

Proof. For any given $\gamma \in \Delta$, the maximum of optimization problem $\max_{\mathbf{H}} \text{Tr}(\mathbf{K}_{\gamma} \mathbf{H}\mathbf{H}^{\top}) \text{ s.t. } \mathbf{H}^{\top}\mathbf{H} = \mathbf{I}_k$ is unique, with $\tilde{\mathbf{H}}^* \in \{\tilde{\mathbf{H}}^* | \tilde{\mathbf{H}}^* = \mathbf{H}^* \mathbf{U}, \mathbf{U} \mathbf{U}^{\top} = \mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_k\}$ the corresponding maximizer. According to Theorem 4.1 in [24], $\mathcal{J}(\gamma)$ in Eq. (6) is differentiable, and $\frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_p} = 2\gamma_p \text{Tr}(\mathbf{K}_p \tilde{\mathbf{H}}^* (\tilde{\mathbf{H}}^*)^{\top}) = 2\gamma_p \text{Tr}(\mathbf{K}_p \mathbf{H}^* \mathbf{H}^{*\top})$. \square

3.3 The Optimization Algorithm

We propose to solve the optimization in Eq. (5) with reduced gradient descent algorithms. We firstly calculate the gradient of $\mathcal{J}(\gamma)$ according to Theorem 1, and then update γ with a descent direction by which the equality and non-negativity constraints on γ can be guaranteed.

To fulfill this goal, we firstly handle the equality constraint by computing the reduced gradient by following [26]. Let γ_u be a non-zero component of γ and $\nabla \mathcal{J}(\gamma)$ denote the reduced gradient of $\mathcal{J}(\gamma)$. The p -th ($1 \leq p \leq m$) element of $\nabla \mathcal{J}(\gamma)$ is

$$[\nabla \mathcal{J}(\gamma)]_p = \frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_p} - \frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_u} \quad \forall p \neq u, \quad (7)$$

and

$$[\nabla \mathcal{J}(\gamma)]_u = \sum_{p=1, p \neq u}^m \left(\frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_u} - \frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_p} \right) \quad (8)$$

Following the suggestion in [26], we choose u to be the index of the largest component of vector γ which is considered to provide better numerical stability.

We then take the positivity constraints on γ into consideration in the descent direction. Note that $-\nabla \mathcal{J}(\gamma)$ is a descent direction since our aim is to minimize $\mathcal{J}(\gamma)$. However, directly using this direction would violate the positivity constraints in the case that if there is an index p such that $\gamma_p = 0$ and $[\nabla \mathcal{J}(\gamma)]_p > 0$. In such case, the descent direction for that component should be set to 0. This gives the descent direction for updating γ as

$$d_p = \begin{cases} 0 & \text{if } \gamma_p = 0 \text{ and } [\nabla \mathcal{J}(\gamma)]_p > 0 \\ -[\nabla \mathcal{J}(\gamma)]_p & \text{if } \gamma_p > 0 \text{ and } p \neq u \\ -[\nabla \mathcal{J}(\gamma)]_u & \text{if } p = u. \end{cases} \quad (9)$$

After a descent direction $\mathbf{d} = [d_1, \dots, d_m]^\top$ is computed by Eq. (9), γ can be calculated via the updating scheme $\gamma \leftarrow \gamma + \alpha \mathbf{d}$, where α is the optimal step size. It can be selected by a one-dimensional line search strategy such as Armijo's rule. The whole algorithm procedure solving the optimization problem in Eq. (4) is outlined in Algorithm 1.

Algorithm 1 SimpleMKKM

```

1: Input:  $\{\mathbf{K}_p\}_{p=1}^m$ ,  $k$ ,  $t = 1$ .
2: Initialize  $\gamma^{(1)} = \mathbf{1}/m$ ,  $\text{flag} = 1$ .
3: while flag do
4:   compute  $\mathbf{H}$  by solving a kernel k-means with
      $\mathbf{K}_{\gamma^{(t)}} = \sum_{p=1}^m (\gamma_p^{(t)})^2 \mathbf{K}_p$ .
5:   compute  $\frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_p}$  ( $p = 1, \dots, m$ ) and the descent di-
     rection  $\mathbf{d}^{(t)}$  in Eq. (9).
6:   update  $\gamma^{(t+1)} \leftarrow \gamma^{(t)} + \alpha \mathbf{d}^{(t)}$ .
7:   if  $\max |\gamma^{(t)} - \gamma^{(t-1)}| \leq 1e - 4$  then
8:     flag=0.
9:   end if
10:   $t \leftarrow t + 1$ .
11: end while
```

3.4 Computational Complexity and Convergence

We discuss the computational complexity of SimpleMKKM. From Algorithm 1, at each iteration, SimpleMKKM needs to solve a kernel k-means problem, calculate the reduced gradient, and search optimal step size. Therefore, its computational complexity at each iteration is $O(n^3 + m * n^3 + m * n_0)$, where n_0 is the maximal number of operations required to find the optimal step size. As observed, SimpleMKKM does not significantly increase the computational complexity of existing MKKM algorithms, as also validated by the experimental results in Figure 6.

We then briefly discuss the convergence of SimpleMKKM. Note that Eq. (6) is a traditional kernel k-means which has a global optimum. Under this condition, the gradient computation in Theorem 1 is exact, and our algorithm performs reduced gradient descent on a continuously differentiable function $\mathcal{J}(\gamma)$ defined on the simplex $\{\gamma \in \mathbb{R}^m | \sum_{p=1}^m \gamma_p = 1, \gamma_p \geq 0, \forall p\}$, which does converge to the minimum of $\mathcal{J}(\gamma)$ [26]. The quick convergence of SimpleMKKM is validated by the experimental results in Figure 5.

We conclude this section by discussing the differences with MKKM-MM [16]. Though both works share a min-max

(max-min) framework, their differences can be summarized from the following three aspects: (1) The objectives are different. SimpleMKKM adopts the unsupervised kernel alignment criterion while MKKM-MM inherits the objective of MKKM, which can be clearly seen from Eq. (3) and Eq. (4). Further, MKKM-MM applies the ℓ_2 -norm constraints on γ to avoid sparse solutions. However, although using the ℓ_1 -norm constraint, our SimpleMKKM still obtains non-sparse solution, as shown by the results in Figure 4. (2) More importantly, the optimization strategies are totally different. MKKM-MM follows the widely used alternating optimization paradigm to solve Eq. (3). In contrast, we, for the first time, reformulate the MKKM as a minimization problem, and develop a reduced gradient descent algorithm to efficiently solve it. (3) The clustering performance is different. We empirically compare their clustering performance, and observe that SimpleMKKM consistently and significantly outperforms MKKM-MM on all 11 benchmark datasets, as shown in Table 1.

4 THE GENERALIZATION ANALYSIS

Generalization error for k-means clustering has been studied by fixing the centroids obtained in the training process and computing their generalization to testing data [27], [28]. In this section, we study how the centroids obtained by the proposed SimpleMKKM generalizes onto test data by deriving its generalization bound.

We now define the error of SimpleMKKM. Let $\hat{\mathbf{C}} = [\hat{\mathbf{C}}_1, \dots, \hat{\mathbf{C}}_k]$ be the learned matrix composed of the k centroids and $\hat{\gamma}$ the learned kernel weights by the proposed SimpleMKKM, where $\hat{\mathbf{C}}_v = \frac{1}{|\hat{\mathbf{C}}_v|} \sum_{j \in \hat{\mathbf{C}}_v} \phi_{\hat{\gamma}}(\mathbf{x}_j)$, $1 \leq v \leq k$. By defining $\Theta = \{\mathbf{e}_1, \dots, \mathbf{e}_k\}$, effective SimpleMKKM clustering should make the following error small

$$1 - \mathbb{E}_{\mathbf{x}} \left[\max_{\mathbf{y} \in \Theta} \langle \phi_{\hat{\gamma}}(\mathbf{x}), \hat{\mathbf{C}}\mathbf{y} \rangle_{\mathcal{H}^k} \right], \quad (10)$$

where $\phi_{\hat{\gamma}}(\mathbf{x}) = [\hat{\gamma}_1 \phi_1^\top(\mathbf{x}), \dots, \hat{\gamma}_m \phi_m^\top(\mathbf{x})]^\top$ is the learned feature map associated with the kernel function $K_{\hat{\gamma}}(\cdot, \cdot)$ and $\mathbf{e}_1, \dots, \mathbf{e}_k$ form the orthogonal bases of \mathbb{R}^k . Intuitively, it says the expected alignment between test points and their closest centroid should be high. We show how the proposed algorithm achieves this goal.

Let us define a function class first:

$$\mathcal{F} = \left\{ f : \mathbf{x} \mapsto 1 - \max_{\mathbf{y} \in \Theta} \langle \phi_{\hat{\gamma}}(\mathbf{x}), \mathbf{C}\mathbf{y} \rangle_{\mathcal{H}^k} \mid \gamma^\top \mathbf{1}_m = 1, \right. \\ \left. \gamma_p \geq 0, \mathbf{C} \in \mathcal{H}^k, |K_p(\mathbf{x}, \tilde{\mathbf{x}})| \leq b, \forall p, \forall \mathbf{x} \in \mathcal{X} \right\}, \quad (11)$$

where \mathcal{H}^k stands for the multiple kernel Hilbert space.

Theorem 2. For any $\delta > 0$, with probability at least $1 - \delta$, the following holds for all $f \in \mathcal{F}$:

$$\mathbb{E}[f(\mathbf{x})] \leq \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) + \frac{\sqrt{\pi/2}bk}{\sqrt{n}} + (1+b) \sqrt{\frac{\log 1/\delta}{2n}}. \quad (12)$$

The detailed proof is provided in the appendix due to conciseness and readability.

According to Theorem 2, for any learned $\hat{\gamma}$ and $\hat{\mathbf{C}}$, to achieve a small

$$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] = 1 - \mathbb{E}_{\mathbf{x}} \left[\max_{\mathbf{y} \in \Theta} \langle \phi_{\hat{\gamma}}(\mathbf{x}), \hat{\mathbf{C}}\mathbf{y} \rangle_{\mathcal{H}^k} \right], \quad (13)$$

the corresponding $\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$ needs to be as small as possible. Assume that γ and \mathbf{C} are obtained by minimizing $\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$ and that \mathbf{H} is constrained to be orthogonal, we have

$$\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) \leq 1 - \frac{1}{n} \text{Tr}(\mathbf{K}_\gamma \mathbf{H} \mathbf{H}^\top) \quad (14)$$

because the proposed algorithm poses a constraint $\mathbf{H}^\top \mathbf{H} = \mathbf{I}_k$ which will make the corresponding centroids non-optimal for minimizing $\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$. This means that $1 - \frac{1}{n} \text{Tr}(\mathbf{K}_\gamma \mathbf{H} \mathbf{H}^\top)$ is an upper bound of $\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$. To minimize the upper bound, we may have to maximize over γ and \mathbf{H} , leading to $\max_\gamma \max_{\mathbf{H}} \text{Tr}(\mathbf{K}_\gamma \mathbf{H} \mathbf{H}^\top)$. However, it is intractable to find a good solution to γ and \mathbf{H} under this criterion, and it is prone to over-fitted solutions [16]. Instead, we take one of its lower bounds, $\min_\gamma \max_{\mathbf{H}} \text{Tr}(\mathbf{K}_\gamma \mathbf{H} \mathbf{H}^\top)$ as the the objective of SimpleMKKM in Eq. (4). This analysis verifies the good generalization ability of the proposed SimpleMKKM.

5 EXPERIMENTAL RESULTS

In this section, we conduct a comprehensive experimental study to evaluate the proposed SimpleMKKM in terms of clustering performance, the learned kernel weights, the running time, and convergence.

5.1 Experimental Settings

A number of standard MKKM benchmark datasets are adopted to evaluate SimpleMKKM, including *Flo17*¹, *Flo102*², *PFold*³, *CCV*⁴, *Digit*⁵, *Cal*⁶. Meanwhile, six sub-datasets, i.e. *Cal-5*, *Cal-10*, *Cal-15*, *Cal-20*, *Cal-25* and *Cal-30*, are constructed via selecting the first 5, 10, 15, 20, 25 and 30 samples from each class respectively from the *Caltech102* data. Their details are shown in Table 2. It can be observed that the number of samples, kernels and categories of these datasets shows considerable variation, providing a good platform to compare the performance of different clustering algorithms.

TABLE 2: Specification of our 11 benchmark datasets.

Dataset	Number of		
	Samples	Kernels	Clusters
Flo17	1360	7	17
Flo102	8189	4	102
PFold	694	12	27
CCV	6773	3	20
Digit	2000	3	10
Cal-5	510	48	102
Cal-10	1020	48	102
Cal-15	1530	48	102
Cal-20	2040	48	102
Cal-25	2550	48	102
Cal-30	3060	48	102

For all data sets, the number of clusters k is assumed known and is set as the true number of classes. The widely

used clustering accuracy (ACC), normalized mutual information (NMI), purity and rand index are applied to evaluate the clustering performance.

For all algorithms, we repeat each experiment 50 times with random initialization to reduce the effect of randomness caused by k-means, and report the means and variation. We next thoroughly study SimpleMKKM in terms of: clustering performance, ablation study on the formulation and optimization, evolution of the learned \mathbf{H} , clustering with number of samples, clustering with number of base kernels, the learned kernel weights, running time and algorithm convergence. Along with SimpleMKKM, we ran another eight comparative algorithms in recent MKC literature, including

- **Average kernel k-means (Avg-KKM)**. The consensus kernel is the uniformly combined base kernels, which is taken as the input of kernel k-means.
- **Multiple kernel k-means (MKKM)** [5]. The base kernels are linearly combined into the consensus kernel. In addition, the combination weights are optimized along with clustering.
- **Localized multiple kernel k-means (LMKKM)** [6]. The base kernels are combined with sample-adaptive weights.
- **Optimal neighborhood kernel clustering (ONKC)** [29]. The consensus kernel is chosen from the neighbor of linearly combined base kernels.
- **Multiple kernel k-means with matrix-induced regularization (MKKM-MiR)** [13]. The optimal combination weights are learned by introducing a matrix-induced regularization term to reduce the redundancy among the base kernels.
- **Multiple kernel clustering with local alignment maximization (LKAM)** [12]. The similarity of a sample to its k -nearest neighbors, instead of all samples, is aligned with the ideal similarity matrix.
- **Multi-view clustering via late fusion alignment maximization (LF-MVC)** [30]. Base partitions are firstly calculated using each single view and then optimally integrated into a consensus partition.
- **MKKM-MM** [16]. It proposes a $\min_{\mathbf{H}} \max_{\gamma}$ formulation that combines views in a way to reveal high within-cluster variance in the combined kernel space and then updates clusters by minimizing such variance.

The implementations of the above algorithms are publicly available in corresponding papers, and we directly adopt them without revision in our experiments. Among all the compared algorithms, ONKC [29], MKKM-MiR [13], LKAM [12] and LF-MVC [30] have hyper-parameters to be tuned. Note that the issue of hyper-parameter tuning in clustering tasks is still an open problem. By following the same way in literature, we reuse their released codes and tune the hyper-parameters by grid search to produce the best possible results on each dataset. By this way, the reported results of these algorithms with hyper-parameters would be over-estimated. As a result, the hyper-parameter tuning would prohibit these multiple kernel (view) clustering algorithm from practical applications. It is therefore desired that a

1. www.robots.ox.ac.uk/~vgg/data/flowers/17/
2. www.robots.ox.ac.uk/~vgg/data/flowers/102/
3. mkl.ucsd.edu/dataset/protein-fold-prediction
4. www.ee.columbia.edu/ln/dvmm/CCV/
5. <http://ss.sysu.edu.cn/py/>
6. www.vision.caltech.edu/Image_Datasets/Caltech101/

TABLE 1: Empirical evaluation and comparison of SimpleMKKM with eight baseline methods on five benchmark datasets in terms of clustering accuracy (ACC), normulaized mutual information (NMI), Purity and Rand Index. Boldface means no statistical difference from the best one.

DATASETS	AVG-KKM	MKKM [5]	LMKKM [6]	ONKC [29]	MKKM-MiR [13]	LKAM [12]	LF-MVC [30]	MKKM-MM [16]	SIMPLEMKKM PROPOSED
ACC									
FLO17	51.0± 1.3	43.6± 1.7	42.7± 1.5	43.4± 2.1	58.0± 1.2	48.9± 0.9	57.2± 1.3	51.0±1.3	59.1± 1.2
FLO102	27.1± 0.8	22.4± 0.5	-	39.2± 0.9	39.1± 1.3	40.4± 1.0	29.0± 1.0	27.1±0.8	42.5± 0.8
PFOLD	29.0± 1.6	27.0± 1.1	22.4± 0.7	35.3± 1.3	34.3± 1.7	33.8± 1.7	31.6± 1.7	29.0±1.6	34.7± 1.9
CCV	19.6± 0.6	18.0± 0.5	18.6± 0.1	22.1± 0.6	20.9± 0.9	18.9± 0.3	23.1± 0.9	19.6±0.6	22.2± 0.7
DIGIT	88.8± 0.1	47.3± 0.7	47.3± 0.7	89.5± 0.1	87.4± 0.1	95.0± 0.1	89.1± 0.1	88.8±0.7	90.3± 0.1
AVG.	43.1	31.7	-	45.9	47.9	47.4	46.0	43.1	49.8
NMI									
FLO17	49.6± 0.8	44.3± 1.3	43.8± 1.0	43.1± 1.3	56.2± 0.6	48.2± 0.6	54.6± 0.9	49.7±0.8	57.5± 0.8
FLO102	46.0± 0.5	42.7± 0.2	-	55.7± 0.4	55.9± 0.6	55.8± 0.3	47.5± 0.3	46.0±0.5	58.6± 0.5
PFOLD	40.3± 1.2	38.0± 0.6	34.7± 0.6	44.0± 0.8	43.1± 1.0	43.6± 1.0	41.8± 0.9	40.3±1.3	44.4± 1.1
CCV	16.8± 0.4	15.1± 0.5	14.4± 0.1	18.4± 0.3	17.9± 0.4	16.8± 0.2	19.3± 0.3	16.8±0.4	18.2± 0.3
DIGIT	80.8± 0.2	48.8± 0.7	48.7± 0.7	81.7± 0.1	79.6± 0.1	89.4± 0.1	81.1± 0.2	80.8±0.2	83.3± 0.1
AVG.	46.7	37.8	-	48.6	50.5	50.8	48.9	46.7	52.4
PURITY									
FLO17	52.0± 1.0	45.1± 1.4	44.5± 1.4	45.2± 1.9	59.4± 0.9	50.1± 0.6	58.1± 1.4	52.0±1.0	60.5± 1.4
FLO102	32.3± 0.6	27.8± 0.4	-	45.1± 0.9	45.2± 1.0	46.7± 0.6	34.5± 0.5	32.3±0.6	48.6± 0.7
PFOLD	37.4± 1.7	33.7± 1.1	31.2± 1.0	41.9± 1.0	41.2± 1.4	41.6± 1.3	38.9± 1.5	37.4±1.7	41.8± 1.5
CCV	23.8± 0.5	22.2± 0.5	22.0± 0.1	24.3± 0.5	23.4± 0.7	22.2± 0.3	26.1± 0.5	23.8±0.5	25.3± 0.5
DIGIT	88.8± 0.1	50.1± 0.7	50.1± 0.7	89.5± 0.1	87.4± 0.1	95.0± 0.1	89.1± 0.1	88.8± 0.1	90.3± 0.1
AVG.	46.9	35.8	-	49.2	51.3	51.1	49.3	46.9	53.3
RAND INDEX									
FLO17	32.3± 1.0	26.4± 1.3	26.0± 1.1	24.3± 1.6	39.6± 0.8	30.2± 0.8	38.6± 1.0	32.3±1.3	41.3± 1.1
FLO102	15.5± 0.5	12.1± 0.4	-	24.5± 0.6	24.9± 1.0	26.3± 0.6	17.2± 0.8	15.5±0.5	28.5± 0.8
PFOLD	14.4± 1.8	12.1± 0.7	7.8± 0.4	17.6± 1.3	17.4± 1.6	17.3± 1.7	16.2± 1.7	14.4±1.8	17.6± 1.9
CCV	6.6± 0.2	5.8± 0.2	5.6± 0.1	7.5± 0.3	7.0± 0.4	6.2± 0.1	8.4± 0.5	6.6±0.2	7.5± 0.2
DIGIT	77.5± 0.2	31.4± 0.6	31.3± 0.6	78.7± 0.1	75.4± 0.1	89.2± 0.1	78.2± 0.2	77.5±0.2	80.3± 0.1
AVG.	29.3	17.6	-	30.5	32.9	33.8	31.7	29.3	35.0

clustering algorithm is parameter-free, as the proposed SimpleMKKM does.

5.2 Experimental Results

5.2.1 Clustering Performance

Table 1 presents the ACC, NMI and purity comparison of the above algorithms. From this table, we have the following observations:

- The proposed SimpleMKKM consistently and significantly outperforms MKKM. For example, it exceeds MKKM by 12.7%, 16%, 6.1%, 3.1%, 34.6%, 4.4%, 7.2%, 8.9%, 10.1%, 10.6% and 11.7% in terms of ACC on all benchmark datasets. These results demonstrate the efficacy of its min-max formulation and associated optimization algorithm.
- MKKM-MM [16] is the first try in literature to improve MKKM via minimization-maximization. As observed, it does improve the MKKM. However the improvement over MKKM is marginal on all datasets. Meanwhile, the proposed SimpleMKKM significantly outperforms MKKM-MM. This once again demonstrates the advantage of our formulation and the associated optimization strategy.
- Our SimpleMKKM achieves comparable or slightly better performance than MKKM-MiR [13], ONKC

[29], and LF-MVC [30], all of which are considered the state of the art in multi-kernel clustering. Note that all of these algorithms have several hyper-parameters to tune due to the incorporation of regularization on the kernel weight γ . Though demonstrating promising clustering performance, these algorithms need to take a lot of effort to determine the best hyper-parameters in practical applications. And parameter tuning may be impossible in real applications where there is no ground truth clustering to optimize. In contrast, our SimpleMKKM is parameter-free.

In summary, SimpleMKKM demonstrates superior clustering performance over the alternatives on all datasets and has no hyper-parameter to be tuned. We expect that the simplicity and efficacy of SimpleMKKM will make it a good option to be considered for practical clustering applications. Note that some results of LMKKM [6] are not reported due to out-of-memory errors, which are caused by its cubic computational and memory complexity.

5.2.2 Ablation Study on the Formulation and Optimization

In order to show the advantage of the proposed formulation and optimization algorithm, we conduct an ablation study on all benchmark datasets to compare the alternatives MKKM-R and SimpleMKKM-C. MKKM-R denotes

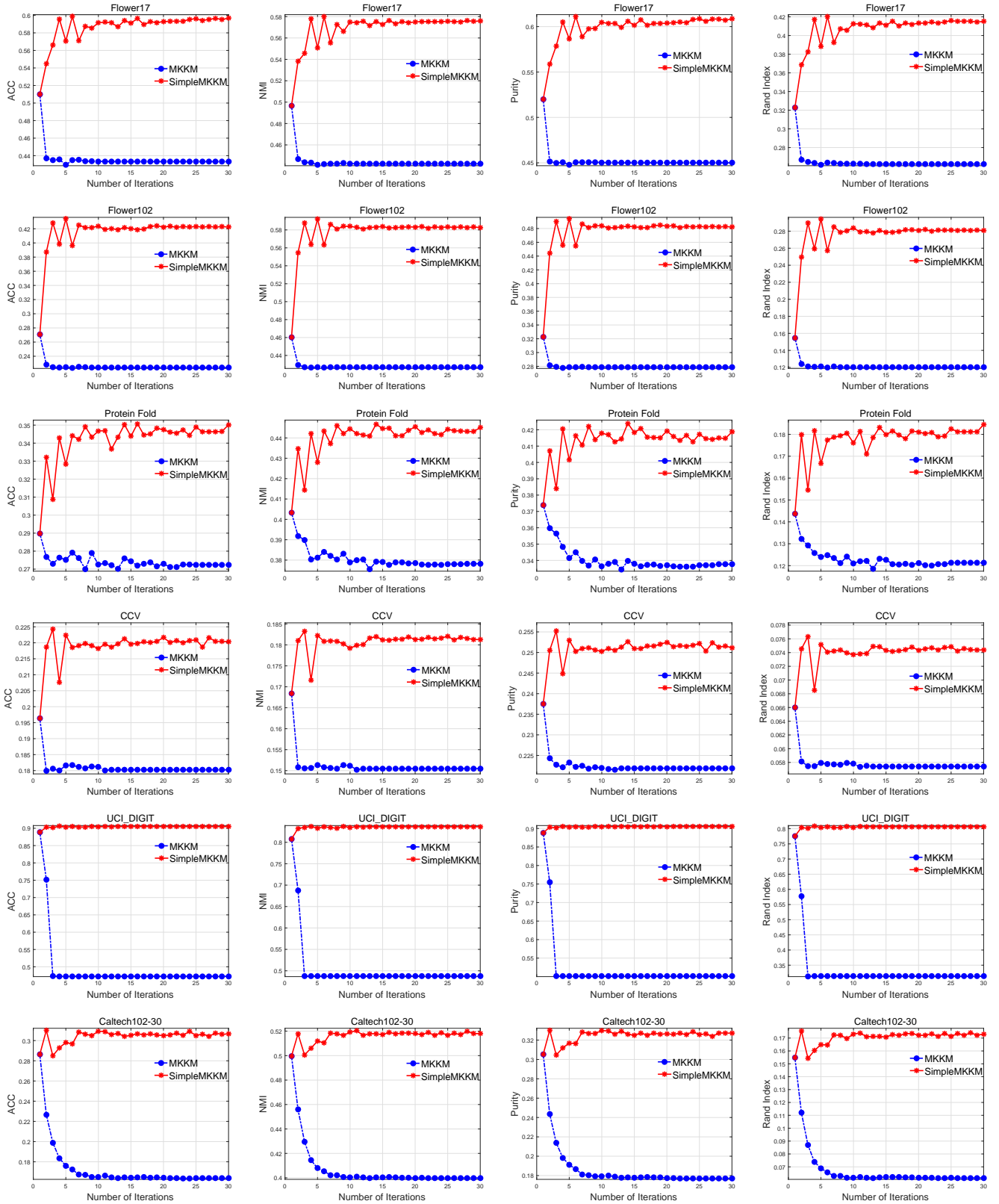


Fig. 1: The clustering comparison of the learned \mathbf{H} by MKKM and the proposed SimpleMKKM with iterations.

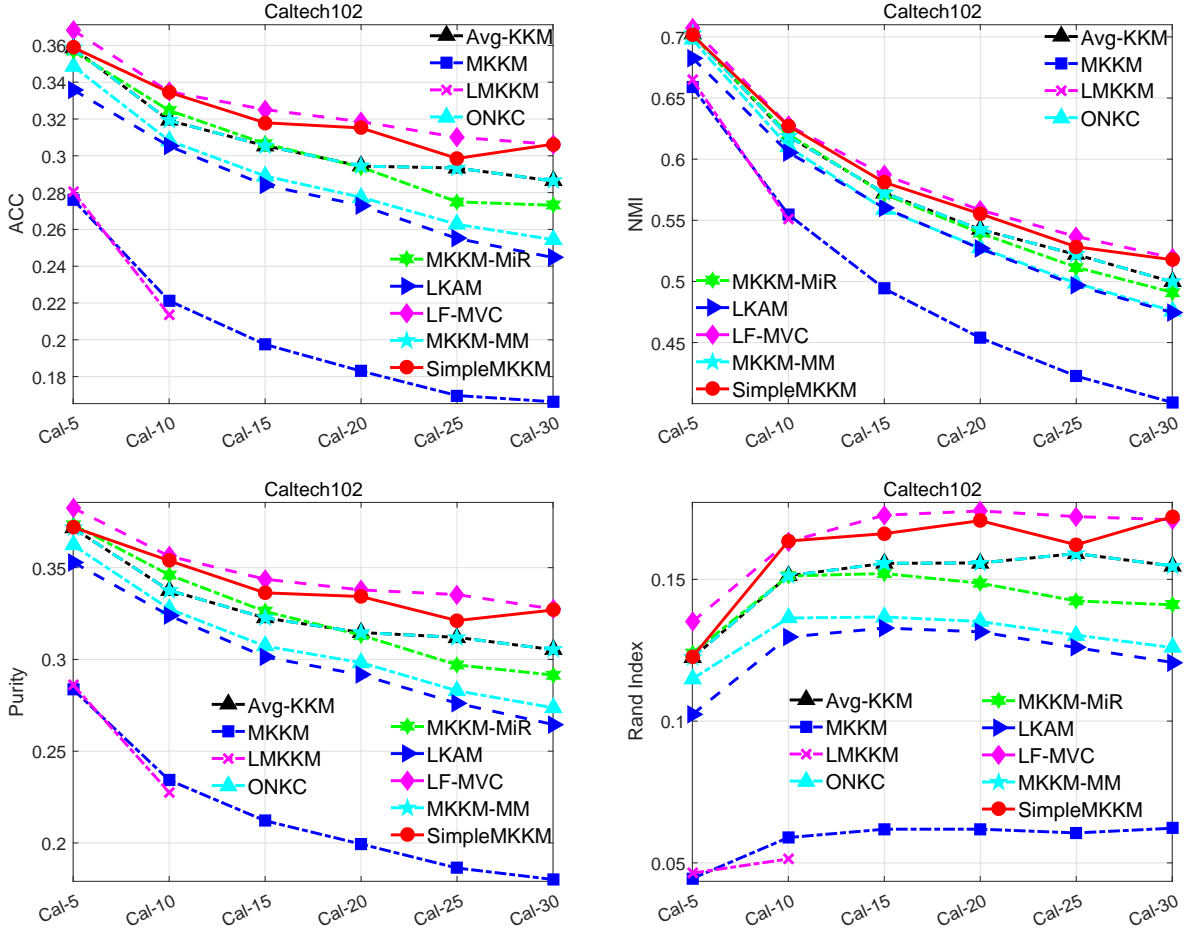


Fig. 2: The clustering performance of the aforementioned algorithms with the variation of number of samples on Caltech102.

optimizing the objective of existing MKKM in Eq. (1) with reduced gradient descent, while SimpleMKKM-C denotes optimizing the criterion in Eq. (4) with coordinate descent optimization (see Section 3.1 for discussion). Note that SimpleMKKM-C has the same objective as SimpleMKKM, but it uses the widely adopted alternate optimization to solve it in place of our newly derived reduced gradient algorithm.

From the results reported in Table 3, we clearly observe that: (1) Our SimpleMKKM and SimpleMKKM-C formulations have significant advantages over MKKM and MKKM-MiR, demonstrating the value of our novel min-max objective; (2) It is also observed that our SimpleMKKM outperforms SimpleMKKM-C, which confirms that our new gradient-based optimization algorithm is also much better than the widely used alternate optimization. This ablation study well demonstrates that both our novel formulation and new optimization attribute to the improvement of clustering performance.

5.2.3 Evolution of the Learned \mathbf{H}

To compare the clustering performance of the proposed SimpleMKKM and existing MKKM with iterations, we take \mathbf{H} at each iteration to calculate ACC, NMI, purity and rand index, and report them in Figure 1. As observed, the start points of both SimpleMKKM and MKKM on all

sub-figures are the same. This is because both algorithms are initialized with the unified weights, which generates the same \mathbf{H} , learning to the same clustering performance. The clustering performance of the proposed SimpleMKKM firstly is increased with iterations, and then kept stable, which sufficiently demonstrates the effectiveness of our algorithm. In contrast, the clustering performance of MKKM is decreased with iterations on all sub-figures, implying that existing MKKM is inferior to average kernel k-means. This states that the widely used MKKM may not be a good choice to fuse multiple base kernels. Comparable, our proposed SimpleMKKM significantly outperforms average kernel k-means on all sub-figures, considerably showing the effectiveness and necessity of the learning procedure.

5.2.4 Clustering Performance with Number of Samples

In this subsection, we conduct an experiment to compare the clustering performance of the proposed SimpleMKKM with the variation of number of samples on Caltech102. In specific, we evaluate their clustering performance on Cal-5, Cal-10, Cal-15, Cal-20, Cal-25 and Cal-30, which are constructed by selecting the first 5, 10, 15, 20, 25 and 30 samples from each class respectively from the Caltech102 data.

The ACC, NMI, purity and rand index of these algorithms with the variation of number of samples are plotted

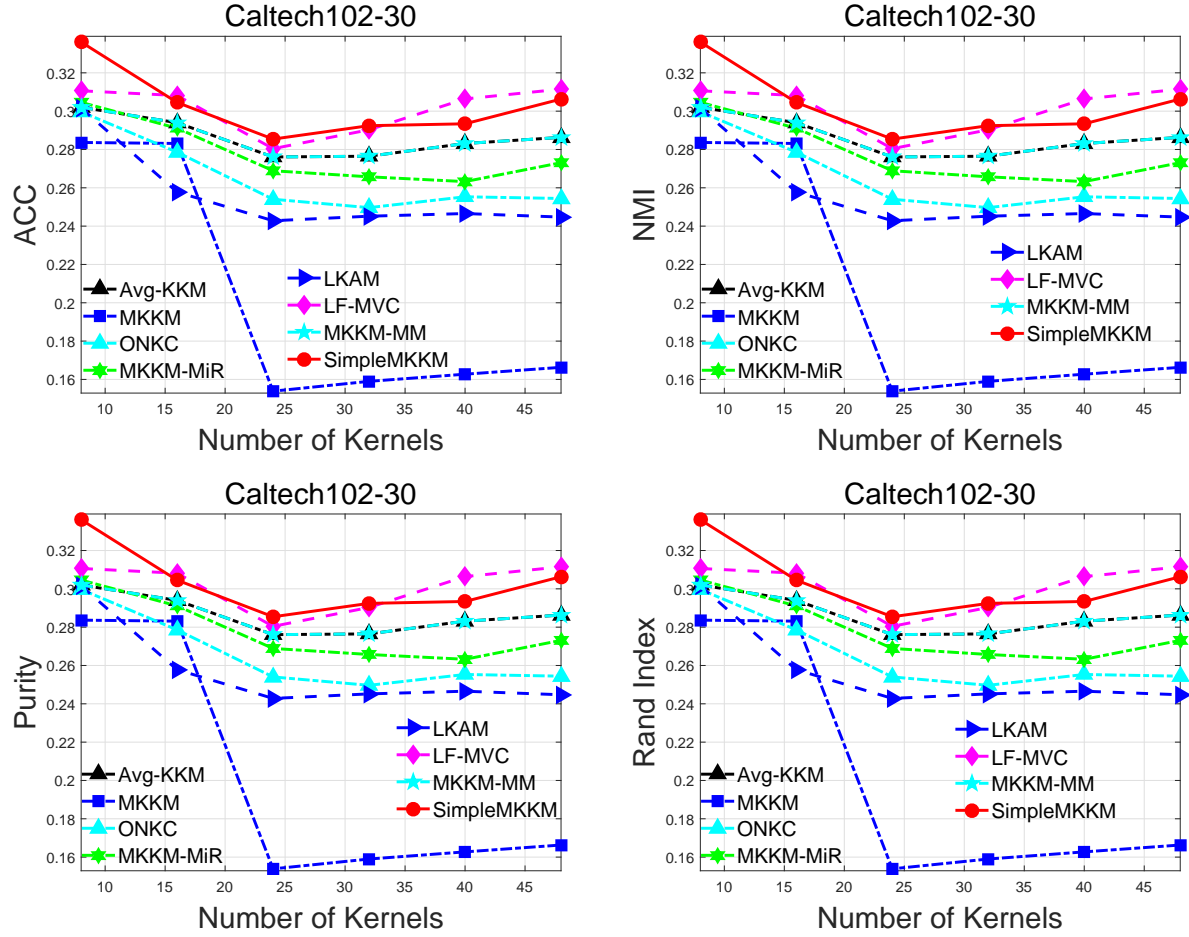


Fig. 3: The clustering performance of the aforementioned algorithms with the variation of number of base kernels on Caltech102.

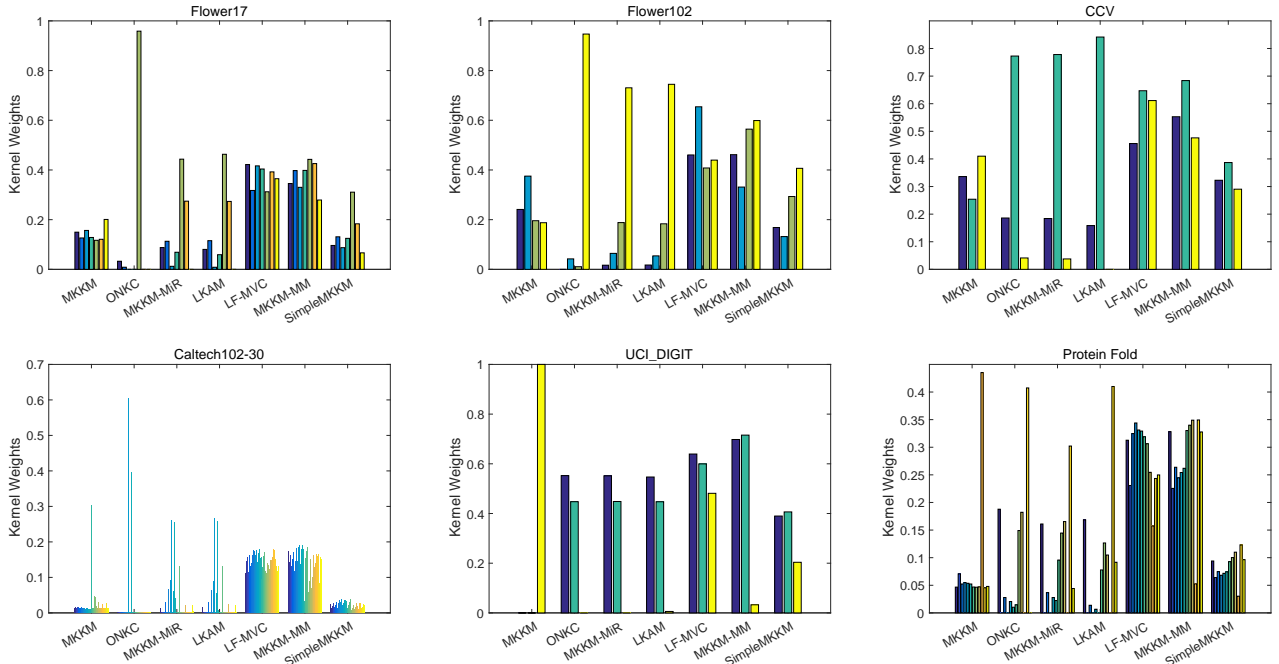


Fig. 4: The kernel weights learned by different algorithms. SimpleMKKM maintains reduced sparsity compared to several competitors. Other datasets omitted due to space limit.

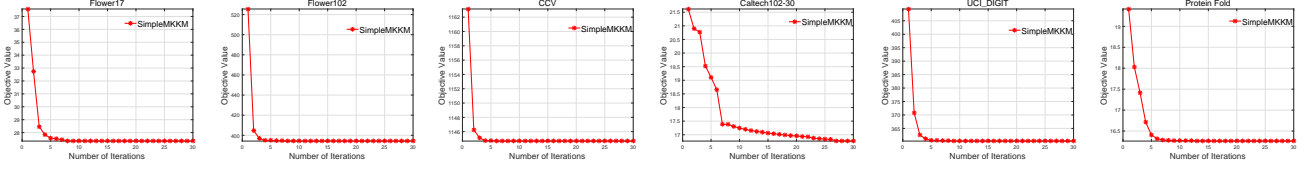


Fig. 5: The objective of SimpleMKKM decreases with iterations. The curves for other datasets are omitted due to space limit.

in Figure 2. As observed, the proposed SimpleMKKM considerably improves the clustering performance of existing MKKM and its variants. Taking the results in sub-figure 2a for example, SimpleMKKM outperforms MKKM by 8.3%, 11.3%, 12.1%, 13.2%, 12.9% and 14% with different number of samples for each cluster, respectively. It exceeds the newly developed MKKM variant, i.e., MKKM-MM [16], by 0.1%, 1.5%, 1.3%, 2.1%, 0.5% and 2%, respectively. We also observe that SimpleMKKM achieves comparable clustering performance with LF-MVC, which is considered as the state-of-the-art in existing multi-view clustering [30]. In sum, the proposed SimpleMKKM achieves the best clustering performance among MKKM based clustering algorithms, and is comparable with the strongest baseline among multi-view clustering in terms of ACC, NMI, purity and rand index.

5.2.5 Clustering with Variation of Base Kernels

To explore the ability of the proposed SimpleMKKM in dealing with different number of base kernels, we design an experiment on Caltech102 by selecting the first

8, 16, 24, 32, 40 and 48 base kernels. The clustering performance in terms of ACC, NMI, purity and rand index of the aforementioned algorithms with different number of base kernels are shown in Figure 3. As observed, we conclude that: i) The proposed SimpleMKKM demonstrates the overall best clustering performance among all compared ones in terms of ACC, NMI, purity and rand index. ii) With the increase of number of base kernels, the clustering performance of MKKM is dramatically decreased. In contrast, the clustering performance of SimpleMKKM is relatively stable with different number of base kernels, demonstrating its advantages in handling large number of base kernels. iii) The results in Figure 3 show that more base kernels is not necessarily helpful for improving clustering performance. In some applications, larger number of base kernels may result in worse clustering performance. This motivates us to automatically select a subset from a group of pre-specified base kernels and optimally combined the selected subset for multiple kernel clustering. This strategy could further significantly improve the clustering performance, which will be explored in our future work.

5.2.6 Kernel Weight Analysis

We next investigate the kernel weights learned by the compared algorithms. The results are plotted in Figure 4. We can see that the kernel weights learned by MKKM are extremely sparse on some datasets such as UCI-Digital, which is caused by the alternate optimization. This sparsity insufficiently exploits the multiple kernel matrices and explains the weak performance of MKKM. For example, the clustering accuracy of MKKM on UCI-Digital is only 47.2%. However, despite the ℓ_1 -norm constraint on γ , the kernel weights learned by our SimpleMKKM are all non-sparse on all datasets, which contributes to its superior clustering performance. This non-sparsity of the learned kernel weights is attributed to our new reduced gradient descent algorithm, which in turn is derived based on our new min-max kernel alignment objective.

5.2.7 Runtime and Convergence

We also report the running time of the compared algorithms in Figure 6. As observed, in addition to significantly improving performance, SimpleMKKM does not considerably increase the running time compared with MKKM and its variants. The objective of SimpleMKKM with iterations is reported in Figure 5. From these figures, we observe that the objective is monotonically decreased and the algorithm usually converges in less than ten iterations on all datasets. This corroborates our earlier theoretical analysis of the nature of our proposed objective and efficient optimisation algorithm.

TABLE 3: Empirical comparison of SimpleMKKM with MKKM, MKKM-R and SimpleMKKM-C on all benchmark datasets.

Dataset	MKKM [5]	MKKM-R	SimpleMKKM-C	SimpleMKKM
ACC				
Flo17	43.6± 1.2	43.7± 1.4	54.2± 1.8	59.1± 1.2
Flo102	22.4± 0.5	22.4± 0.5	41.8± 1.2	42.5± 0.8
PFold	27.0± 1.1	26.6± 1.1	29.0± 1.4	34.7± 1.9
CCV	18.0± 0.5	17.9± 0.6	22.1± 0.7	22.2± 0.7
Digit	47.3± 0.7	47.3± 0.7	90.4± 0.9	90.3± 0.6
Cal-30	16.6± 0.4	16.7± 0.4	30.4± 1.1	30.6± 0.9
NMI				
Flo17	44.3± 1.3	44.3± 1.1	54.3± 1.4	57.5± 0.8
Flo102	42.7± 0.2	42.6± 0.2	58.0± 0.5	58.6± 0.5
PFold	38.0± 0.6	37.5± 0.8	38.4± 0.8	44.4± 1.1
CCV	15.1± 0.5	14.8± 0.4	18.2± 0.3	18.2± 0.3
Digit	48.8± 0.7	48.7± 0.7	83.5± 0.2	83.3± 0.1
Cal-30	40.1± 0.3	40.2± 0.3	51.8± 0.6	51.8± 0.5
Purity				
Flo17	45.1± 1.4	44.9± 1.4	55.1± 1.8	60.5± 1.4
Flo102	27.8± 0.4	27.8± 0.4	47.9± 0.8	48.6± 0.7
PFold	33.7± 1.1	33.1± 0.9	35.7± 1.0	41.8± 1.4
CCV	22.2± 0.5	22.3± 0.4	25.2± 0.5	25.3± 0.5
Digit	50.1± 0.7	50.1± 0.7	90.4± 0.9	90.3± 0.6
Cal-30	18.0± 0.5	18.1± 0.4	32.5± 1.0	32.7± 0.8
Rand Index				
Flo17	45.1± 1.4	44.9± 1.4	55.1± 1.8	60.5± 1.4
Flo102	27.8± 0.4	27.8± 0.4	47.9± 0.8	48.6± 0.7
PFold	33.7± 1.1	33.1± 0.9	35.7± 1.0	41.8± 1.4
CCV	22.2± 0.5	22.3± 0.4	25.2± 0.5	25.3± 0.5
Digit	50.1± 0.7	50.1± 0.7	90.4± 0.9	90.3± 0.6
Cal-30	18.0± 0.5	18.1± 0.4	32.5± 1.0	32.7± 0.8

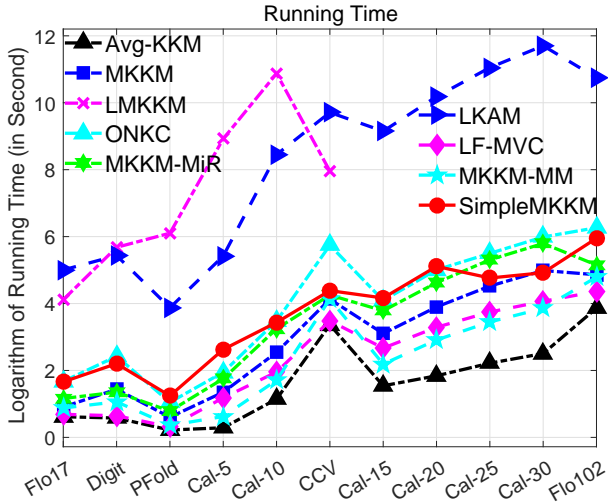


Fig. 6: Running time of different algorithms on 11 benchmark datasets (in second). The experiments are conducted on a PC with Intel(R) Core(TM)-i7-5820 3.3 GHz CPU and 32G RAM in MATLAB environment. SimpleMKKM is comparably fast to alternatives while providing superior performance and requiring no hyper-parameter tuning.

6 CONCLUSION

In this paper, we have extended the widely used supervised kernel alignment criterion to clustering, and introduce a novel clustering objective of by minimizing alignment for γ and maximizing it for H . We show that this novel objective can be transformed into a minimization problem which is differentiable and amenable to a solution by reduced gradient descent. This makes SimpleMKKM unique among MKC alternatives, in not requiring a local-minimum prone alternating coordinate descent strategy.

We derive a generalization bound for our approach using global Rademacher complexity analysis. Comprehensive experiments demonstrate the effectiveness of SimpleMKKM. We expect that the simplicity, lack of hyper-parameters, and efficacy of SimpleMKKM will make it a go-to solution for practical multi-kernel clustering applications in future. Future work may aim to extend SimpleMKKM to handle incomplete kernels, study further applications, and derive convergence rates using local Rademacher complexity analysis [31], [32]. In addition, we plan to automatically select a subset from a large number of base kernels, and optimally combined them for multiple kernel clustering.

ACKNOWLEDGEMENTS

We thank Prof. Timothy M. Hospedales and Prof. Marius Kloft for their constructive comments on revising the manuscript. This work was supported by the Natural Science Foundation of China (project no. 61773392 and 61922088).

REFERENCES

[1] B. Zhao, J. T. Kwok, and C. Zhang, "Multiple kernel clustering," in *SDM*, 2009, pp. 638–649.

[2] S. Yu, L.-C. Tranchevent, X. Liu, W. Glänzel, J. A. K. Suykens, B. D. Moor, and Y. Moreau, "Optimized data fusion for kernel k-means clustering," *IEEE TPAMI*, vol. 34, no. 5, pp. 1031–1039, 2012.

[3] C. Zhang, Y. Cui, Z. Han, J. T. Zhou, H. Fu, and Q. Hu, "Deep partial multi-view learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.

[4] S. Sun, W. Dong, and Q. Liu, "Multi-view representation learning with deep gaussian processes," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.

[5] H. Huang, Y. Chuang, and C. Chen, "Multiple kernel fuzzy clustering," *IEEE Trans. Fuzzy Systems*, vol. 20, no. 1, pp. 120–134, 2012.

[6] M. Gönen and A. A. Margolin, "Localized data fusion for kernel k-means clustering with application to cancer biology," in *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8–13 2014, Montreal, Quebec, Canada*, 2014, pp. 1305–1313.

[7] X. Peng, Z. Huang, J. Lv, H. Zhu, and J. T. Zhou, "COMIC: multi-view clustering without parameter selection," in *Proceedings of the 36th International Conference on Machine Learning, ICML, 2019*, pp. 5092–5101.

[8] A. Kumar and H. Daumé, "A co-training approach for multi-view spectral clustering," in *ICML*, 2011, pp. 393–400.

[9] A. Kumar, P. Rai, and H. Daumé, "Co-regularized multi-view spectral clustering," in *NIPS*, 2011, pp. 1413–1421.

[10] Z. Zhang, L. Liu, F. Shen, H. T. Shen, and L. Shao, "Binary multi-view clustering," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 7, pp. 1774–1782, 2019.

[11] X. Li, H. Zhang, R. Wang, and F. Nie, "Multi-view clustering: A scalable and parameter-free bipartite graph fusion method," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.

[12] M. Li, X. Liu, L. Wang, Y. Dou, J. Yin, and E. Zhu, "Multiple kernel clustering with local kernel alignment maximization," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9–15 July 2016*, 2016, pp. 1704–1710.

[13] X. Liu, Y. Dou, J. Yin, L. Wang, and E. Zhu, "Multiple kernel k-means clustering with matrix-induced regularization," in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12–17, 2016, Phoenix, Arizona, USA, 2016*, pp. 1888–1894.

[14] X. Liu, M. Li, L. Wang, Y. Dou, J. Yin, and E. Zhu, "Multiple kernel k-means with incomplete kernels," in *AAAI*, 2017, pp. 2259–2265.

[15] X. Liu, X. Zhu, M. Li, L. Wang, E. Zhu, T. Liu, M. Kloft, D. Shen, J. Yin, and W. Gao, "Multiple kernel k-means with incomplete kernels," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 5, pp. 1191–1204, May 2020.

[16] S. Bang, Y. Yu, and W. Wu, "Robust multiple kernel k-means clustering using min-max optimization," 2018.

[17] X. Liu, L. Wang, X. Zhu, M. Li, E. Zhu, T. Liu, L. Liu, Y. Dou, and J. Yin, "Absent multiple kernel learning algorithms," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 6, pp. 1303–1316, 2020.

[18] X. Liu, X. Zhu, M. Li, L. Wang, C. Tang, J. Yin, D. Shen, H. Wang, and W. Gao, "Late fusion incomplete multi-view clustering," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 10, pp. 2410–2423, Oct 2019.

[19] X. Liu, M. Li, C. Tang, J. Xia, J. Xiong, L. Liu, M. Kloft, and E. Zhu, "Efficient and effective regularized incomplete multi-view clustering," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–13, 2020.

[20] M. Kloft, U. Brefeld, S. Sonnenburg, and A. Zien, " l_p -norm multiple kernel learning," *JMLR*, vol. 12, pp. 953–997, 2011.

[21] C. Cortes, M. Mohri, and A. Rostamizadeh, "L2 regularization for learning kernels," in *UAI*, 2009, pp. 109–116.

[22] —, "Algorithms for learning kernels based on centered alignment," *JMLR*, vol. 13, pp. 795–828, 2012.

[23] N. Cristianini, J. Shawe-Taylor, A. Elisseeff, and J. S. Kandola, "On kernel-target alignment," in *Advances in Neural Information Processing Systems 14*, 2002.

[24] J. F. Bonnans and A. Shapiro, "Optimization problems with perturbations: A guided tour," *SIAM Review*, vol. 40, no. 2, pp. 228–264, 1998.

[25] O. Chapelle, V. Vapnik, O. Bousquet, and S. Mukherjee, "Choosing multiple parameters for support vector machines," *Machine Learning*, vol. 46, no. 1, pp. 131–159, Jan 2002.

- [26] A. Rakotomamonjy, F. R. Bach, S. Canu, and Y. Grandvalet, "Simplemkl," *JMLR*, vol. 9, pp. 2491–2521, 2008.
- [27] A. Maurer and M. Pontil, " k -dimensional coding schemes in Hilbert spaces," *IEEE Transactions on Information Theory*, vol. 56, no. 11, pp. 5839–5846, 2010.
- [28] T. Liu, D. Tao, and D. Xu, "Dimensionality-dependent generalization bounds for k -dimensional coding schemes," *Neural computation*, vol. 28, no. 10, pp. 2213–2249, 2016.
- [29] X. Liu, S. Zhou, Y. Wang, M. Li, Y. Dou, E. Zhu, and J. Yin, "Optimal neighborhood kernel clustering with multiple kernels," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4–9, 2017, San Francisco, California, USA, 2017*, pp. 2266–2272.
- [30] S. Wang, X. Liu, E. Zhu, C. Tang, J. Liu, J. Hu, J. Xia, and J. Yin, "Multi-view clustering via late fusion alignment maximization," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10–16, 2019*, 2019, pp. 3778–3784.
- [31] M. Kloft and G. Blanchard, "On the convergence rate of l_p -norm multiple kernel learning," *J. Mach. Learn. Res.*, vol. 13, pp. 2465–2502, 2012.
- [32] C. Cortes, M. Kloft, and M. Mohri, "Learning kernels using local rademacher complexity," in *Advances in Neural Information Processing Systems*, 2013, pp. 2760–2768.



Xinwang Liu received his PhD degree from National University of Defense Technology (NUDT), China. He is now Professor at School of Computer, NUDT. His current research interests include kernel learning and unsupervised feature learning. Dr. Liu has published 60+ peer-reviewed papers, including those in highly regarded journals and conferences such as IEEE T-PAMI, IEEE T-KDE, IEEE T-IP, IEEE T-NNLS, IEEE T-MM, IEEE T-IFS, NeurIPS, CVPR, ICCV, AAAI, IJCAI, etc. More information can be found at <https://xinwangliu.github.io/>.



Li Liu received the BSc degree in communication engineering, the MSc degree in photogrammetry and remote sensing and the Ph.D. degree in information and communication engineering from the National University of Defense Technology (NUDT), China, in 2003, 2005 and 2012, respectively. She joined the faculty at NUDT in 2012, where she is currently an Associate Professor with the College of System Engineering. She was a cochair of seven International Workshops at CVPR, ICCV, and ECCV. She is going to lecture a tutorial at CVPR'19. She was a guest editor of special issues for IEEE TPAMI and IJCV. Her current research interests include facial behavior analysis, texture analysis, image classification, object detection and recognition. Her papers have currently over 1800 citations in Google Scholar. She currently serves as Associate Editor of the Visual Computer Journal.



Jian Xiong received the B.S. degree in engineering, and the M.S. and Ph.D. degrees in management from National University of Defense Technology, Changsha, China, in 2005, 2007, and 2012, respectively. He is an Associate Professor with the School of Business Administration, Southwestern University of Finance and Economics. His research interests include data mining, multiobjective evolutionary optimization, multiobjective decision making, project planning, and scheduling.



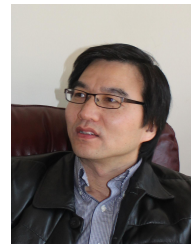
En Zhu received his PhD degree from National University of Defense Technology (NUDT), China. He is now Professor at School of Computer Science, NUDT, China. His main research interests are pattern recognition, image processing, machine vision and machine learning. Dr. Zhu has published 60+ peer-reviewed papers, including IEEE T-CSVT, IEEE T-NNLS, PR, AAAI, IJCAI, etc. He was awarded China National Excellence Doctoral Dissertation.



Junwei Han (Senior Member, IEEE) received the B.S. and Ph.D degrees from Northwestern Polytechnical University, Xian, China, in 1999 and 2003, respectively. He is currently a Professor in Northwestern Polytechnical University. His research interests include computer vision and pattern recognition. He is an Associate Editor of IEEE Trans. on Neural Networks and Learning Systems, IEEE Trans. on Multimedia, and IEEE Trans. on Circuits and Systems for Video Technology.



Meng Wang is a professor at the Hefei University of Technology, China. He received his B.E. degree and Ph.D. degree in the Special Class for the Gifted Young and the Department of Electronic Engineering and Information Science from the University of Science and Technology of China (USTC), Hefei, China, in 2003 and 2008, respectively. His current research interests include multimedia content analysis, computer vision, and pattern recognition. He has authored more than 200 book chapters, journal and conference papers in these areas. He is the recipient of the ACM SIGMM Rising Star Award 2014. He is an associate editor of IEEE Transactions on Knowledge and Data Engineering (IEEE TKDE), IEEE Transactions on Circuits and Systems for Video Technology (IEEE TCSVT), IEEE Transactions on Neural Networks and Learning Systems (IEEE TNNLS), and IEEE Transactions on Multimedia (IEEE TMM).



Dinggang Shen is Jeffrey Houpt Distinguished Investigator, and a Professor of Radiology, Biomedical Research Imaging Center (BRIC), Computer Science, and Biomedical Engineering in the University of North Carolina at Chapel Hill (UNC-CH). He is currently directing the Center for Image Analysis and Informatics, the Image Display, Enhancement, and Analysis (IDEA) Lab in the Department of Radiology, and also the medical image analysis core in the BRIC. He was a tenure-track assistant professor in the University of Pennsylvania (UPenn), and a faculty member in the Johns Hopkins University. Dr. Shens research interests include medical image analysis, computer vision, and pattern recognition. He serves as an editorial board member for eight international journals. He has also served in the Board of Directors, the Medical Image Computing and Computer Assisted Intervention (MICCAI) Society, in 2012–2015, and will be General Chair for MICCAI 2019. He is Fellow of IEEE, AIMBE and IAPR.



Wen Gao received his PhD degree from University of Tokyo, Japan. He is now Boya Chair Professor and the Director of Faculty of Information and Engineering Sciences at Peking University, and the founding director of National Engineering Lab. for Video Technology (NELVT) at Peking University. Prof. Gao works in the areas of multimedia and computer vision, topics including video coding, video analysis, multimedia retrieval, face recognition, multimodal interfaces, and virtual reality. He published seven books, over 220 papers in refereed journals, and over 600 papers in selected international conferences. He served or serves on the editorial board for several journals, such as IEEE T-IP, IEEE T-CSVT, IEEE T-MM, IEEE T-AMD. He chaired a number of prestigious international conferences on multimedia and video signal processing, such as IEEE ICME 2007, ACM Multimedia 2009, IEEE ISCAS 2013, and also served on the advisory and technical committees of numerous professional organizations. Prof. Gao has been featured by IEEE Spectrum in June 2005 as one of the "Ten To Watch" among China's leading technologists. He is a fellow of IEEE, a fellow of ACM, and a member of Chinese Academy of Engineering.